

توضیحات پیاده‌سازی فازهای پروژه درس کامپایلر

۱-۱- مفاهیم کامپایلری

قبل از این که به بررسی نرم‌افزار لکس و یک پردازیم، ابتدا تعدادی از مفاهیم کامپایلری را به طور اختصار مورد مطالعه قرار می‌دهیم.

الفبا: مجموعه‌ای متناهی از عناصر را الفبا می‌گویند. مانند الفبای زبان فارسی که از ۳۲ حرف تشکیل شده است. از خصوصیات الفبا، ساده و تجزیه‌ناپذیر بودن آن است.

رشته: به مجموعه‌ای متناهی از الفبا، رشته گفته می‌شود. مانند کلمات در زبان فارسی که ترکیبی از الفبای فارسی هستند. البته در زبان محاوره‌ای و یا نوشتاری روزمره و کتاب‌ها، رشته‌ها یا همان کلمات باید دارای معنا و مفهوم خاصی باشند که بیان‌کننده‌ی منظور شخص گوینده یا نویسنده است ولی در تعریف بالا، کلمات می‌توانند بدون معنا و یا حتی از نظر املائی غلط باشند؛ مانند اگل، ضضع، خاهر و... در مثال‌های ذکرشده، دو مورد اول از نظر معنا غلط و مورد سوم از نظر املائی غلط هستند ولی با این حال، یک رشته به حساب می‌آیند.

مجموعه: مجموعه‌ای که در این جا مورد نظر است، شبیه مجموعه‌ها در ریاضی است که از دو آکولاد باز و بسته "{}"، که بین این دو آکولاد عناصر الفبا قرار می‌گیرند، تشکیل شده است. مانند:

$\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$ یا $\{a, b, c\}$

همان‌طور که در ریاضیات بر روی مجموعه‌ها، عملیات اجتماع، اشتراک و ضرب^۱ تعریف می‌شود، برای این مجموعه‌ها نیز چنین اعمالی وجود دارد که شبیه عملیات مشابه در ریاضیات است.

مجموعه‌های باقاعده: مجموعه‌هایی را باقاعده گویند که دارای شرایط زیر باشد:

- مجموعه‌های تک عضوی و هم‌چنین مجموعه‌های تهی از دسته‌ی مجموعه‌های باقاعده هستند.

$\{\}$ ، $\{b\}$

- اگر دو مجموعه به‌طور جداگانه هر یک مجموعه‌ای باقاعده باشند، در این صورت اجتماع، اشتراک و ضرب این دو مجموعه نیز مجموعه‌های باقاعده به شمار می‌روند.
- مجموعه‌ای باقاعده است که بتوان آن را از مجموعه‌های تک عضوی با استفاده‌ی مکرر از عملگرهای اتصال، اتحاد و ستاره تولید نمود.

عبارات باقاعده: همان مجموعه باقاعده است که فاقد علامت مجموعه "{}" باشد و تمام قوانین مربوط به مجموعه‌های باقاعده در مورد عبارات باقاعده صدق می‌کند.

زبان: مجموعه‌ی رشته‌هایی از الفبا است که از ساختار خاصی پیروی می‌کنند. برای تعریف یک زبان، مطلوب آن است که:

- الفبای آن تعریف شود.
- ابزاری جهت تعریف ساختار زبان تولید نمود.
- ابزاری جهت تولید رشته‌های زبان و جهت ارزیابی تعلق یک رشته به زبان تولید نمود. (رشته‌ای متعلق به یک زبان است که از ساختار آن زبان پیروی کند و از الفبای آن ساخته شده باشد).

^۱ در بعضی از کتاب‌ها به عملیات ضرب مجموعه‌ها، اتصال یا الحاق مجموعه‌ها نیز گفته می‌شود.

گرامر: برای بیان ساختار یک زبان و تعریف عمومی و کلی یک زبان، از مفهومی به نام گرامر استفاده می‌شود.

برای درک بیشتر مفاهیم بالا، به مثال شکل ۱ توجه کنید:

{کتاب، علی، شیشه، شکست، را} = پایانه

{<فاعل>، <فعل>، <مفعول>، <اسم>، <جمله>} = ناپایانه

<جمله> = نماد شروع

<فاعل> <مفعول> را <فعل> → <جمله> : قواعد

شکست → <فعل>

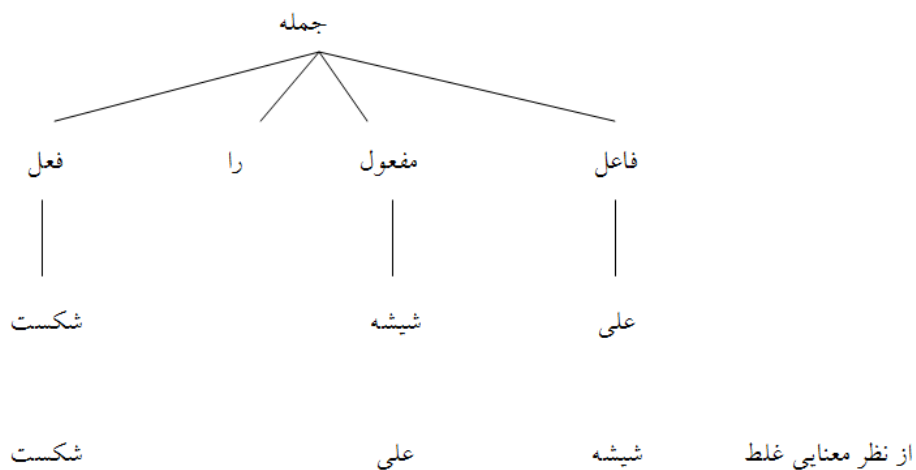
<اسم> → <فاعل>

<اسم> → <مفعول>

شیشه → <اسم>

کتاب → <اسم>

علی → <اسم>



شکل-۱ ساختار یک جمله در زبان فارسی

گرامر می‌تواند جملاتی را تولید نماید که از لحاظ ساختار، کلیه‌ی مشخصات زبان را داشته باشد اما تضمینی بر تعلق رشته به زبان بر اساس معنای آن وجود ندارد. به عبارت دیگر، گرامر می‌تواند روی ساختار زبان صحبت کند ولی در مورد معنای آن نمی‌تواند بحث کند.

۲-۱- بررسی نرم افزار لکس و یک

نرم افزار لکس و یک از جمله ابزارهایی هستند که برای تسهیل روند ساخت کامپایلر ایجاد شده اند. این دو نرم افزار برای تحلیل لغوی و نحوی زبان ها به کار می روند. به این صورت که به ترتیب عبارات باقاعده و گرامر را از کاربر گرفته و زیربرنامه های تحلیل لغوی و نحوی را در اختیار قرار می دهند. لکس برای تهیه برنامه هایی که ورودی ها را تجزیه و تحلیل می کنند، به کار می رود. به عنوان مثال، اگر شما می خواهید یک برنامه ی ماشین حساب ساده بنویسید که وقتی عبارت $۸۵ + ۱۷۵$ را به آن دادید، جواب محاسبه را به شما برگرداند، می توانید از لکس برای تشریح رشته حرفی ورودی استفاده کنید. عبارت بالا از سه کلمه "۱۷۵"، "+" و "۸۵" تشکیل شده است؛ بنابراین لکس، برنامه ای برای تحلیل ورودی و تبدیل کلمات به نشانه^۲ - که مشخص کننده نوع کلمه است - را تولید می کند.

مزیت استفاده از لکس در این است که می توان با دقت و کارایی بیشتری، جزئیات فنی را که در تحلیل و ترجمه باید به کار روند، مورد استفاده قرار داد. کد تهیه شده توسط لکس، شناسایی داده های ورودی و مقایسه ی آن ها با چیزی که شما می خواهید را برآورده می کند. در یک، قوانینی را که باید بر نشانه ها اعمال شوند، به صورت گرامر و به عنوان ورودی داده شده و ابزار یک، به طور خودکار زیربرنامه ای تولید می کند که بتواند این قوانین را اعمال نماید. به بیان ساده تر، در ورودی یک، طرز استفاده از نشانه ها مشخص می شوند.

این دو برنامه می توانند با یکدیگر و یا به طور جداگانه مورد استفاده قرار بگیرند. برنامه ی لکس می تواند وظیفه ی تحلیل لغوی ورودی را بر عهده داشته باشد و برنامه ای برای استفاده از این تحلیل نوشته شود، یا با نوشتن برنامه ای کار تحلیل لغوی انجام شود و سپس با استفاده از یک، برنامه ای برای به کار بردن آن در اختیار قرار بگیرد. گرچه بهتر است این دو برنامه با یکدیگر مورد استفاده قرار بگیرند، زیرا در این صورت کارایی برنامه بالاتر خواهد بود.

قابل ذکر است که کامپایلر از دو بخش تحلیل گر لغوی و تحلیل گر نحوی بهره می برد. از لکس برای نوشتن تحلیل گر لغوی و از یک برای نوشتن تحلیل گر نحوی استفاده می شود.

۳-۱- خروجی های دو برنامه

این دو برنامه می توانند برنامه هایی به زبان های پاسکال، سی++^۳، سی و جاوا تهیه کنند که در ساختن کامپایلرها استفاده می شوند. همچنین می توان به این برنامه ها که به صورت خودکار تولید شده اند، کدهای دیگری اضافه شود که نیاز خاصی مورد نظر است. البته با ذکر این نکته که نباید کد اضافه شده، اختلالی در کدهای دیگر برنامه ایجاد کند.

خروجی لکس

هدف ابزار لکس، تولید کد برای زیربرنامه ای به نام yylex است که این زیربرنامه بدون آرگومان، و خروجی آن یک عدد صحیح است. اگر عدد صحیح خروجی صفر باشد، نشانه ی پایان پرونده، و در غیر این صورت، کدی که معرف یک نشانه مربوط به متن ورودی است، برگردانده می شود. متغیرها و زیربرنامه های دیگر آن عبارتند از:

```
char *ytext;
int yyleng;
int yylineno;
FILE *yyin;
```

^۲ Token

^۳ C++

```
FILE *yyout;
int yylex(void);
void yy_reset(void);
int yygetc(void);
void yyerror(char*, ...);
int yywrap(void);
```

:yytext

متغیری که حروفی که یک نشانه را تشکیل می‌دهند، در خود نگه می‌دارد.

:yyleng

متغیری که طول نشانه را نگه می‌دارد.

:yyin

مکانی که نشانه‌ها باید از آن‌جا خوانده شوند را در خود نگه می‌دارد.

:yyout

مکانی که خروجی‌های لکس باید در آن قرار بگیرند.

:yyreset()

این زیربرنامه برای مقداردهی اولیه به متغیرها به کار می‌رود. بعد از هر بار شناسایی، این زیربرنامه صدا زده می‌شود.

:yygetc()

زیربرنامه‌ای که برای خواندن ورودی مورد استفاده قرار می‌گیرد.

:yyerror()

زیربرنامه استاندارد پیغام خطا که به وسیله‌ی یک نیز مورد استفاده قرار می‌گیرد.

:yywrap

این زیربرنامه زمانی فراخوانی می‌شود که در خواندن ورودی به انتهای پرونده رسیده شد.

:yylval

یک متغیر عمومی است که می‌توان از آن برای انجام کارهای مورد نظر استفاده کرد، برای نمونه، در مثال بالا می‌توان از آن برای قراردادن مقدار واقعی نشانه‌ی عددی که لکس پیدا می‌کند، استفاده کرد.

خروجی یک

هدف یک تولید کد برای زیربرنامه `yyparse()` است. این زیربرنامه از `yylex()` برای به‌دست‌آوردن نشانه استفاده می‌کند و همچنین در مواقعی که نیاز به مقدار واقعی نشانه‌ها داشته باشد، از `yylval` استفاده می‌کند. در سی و پاسکال این زیربرنامه بدون هیچ آرگومانی صدا زده می‌شود. خروجی این زیربرنامه، اگر ورودی پردازش شده درست باشد، صفر است و مقدار آن اگر ورودی پردازش شده، نادرست باشد، یک است. لازم به ذکر است که می‌توان این مقادیر را تغییر داد.

۴-۱- تعریف نشانه‌ها

اگر به قسمت‌های قبل توجه کرده باشید، `yylex()` کدی را که معرف نشانه‌ی پیدا شده است را برمی‌گرداند و `yyparse()` بر روی آن عمل می‌کند. بدیهی است که این دو زیربرنامه باید بر مقدارهای متفاوت نشانه‌ها با یکدیگر هماهنگ باشند.

نشانه‌ها را می‌توان در برنامه‌ی یک به صورت زیر تعریف کرد:

<i>%token</i>	<i>token_name</i>
---------------	-------------------

این نشانه‌ها در برنامه‌ی تولیدشده به زبان سی به صورت متغیرهای ثابت که با `#define` شروع می‌شوند، و در پاسکال با استفاده از `CONST`، نمایش داده می‌شوند. مقدار این نشانه‌ها از ۲۵۷ شروع می‌شود. زیرا از اعداد ۱ تا ۲۵۶ برای کدهای آسکی حرف‌ها استفاده می‌شود. وقتی یک اجرا می‌شود، یک پرونده که شامل زیربرنامه تحلیل‌گر نحوی است و همچنین یک سرآیند^۴ پرونده که شامل تعریف نشانه‌ها می‌باشد، ایجاد می‌کند.

جدول ۱ - برنامه‌های تجزیه‌گر و سرآیند آن در زبان‌های مختلف

زبان	سرآیند پرونده	پرونده برنامه تجزیه‌گر
سی	ytab.h	ytab.c
سی++	ytab.hpp	ytab.cpp
پاسکال	ytab.hp	ytab.pas

برای مثال، در یک برنامه‌ی ماشین‌حساب ساده که ورودی‌های آن به صورت زیر هست:

789 + 45, 9045 - 234

دو نوع نشانه وجود دارد:

- عملوند عدد صحیح
- عملگر ریاضی

اگر `yylex` یک عملگر شبیه + یا - پیدا کند، یک کد مبنی بر پیدا کردن آن برمی‌گرداند و اگر یک عملوند عدد صحیح پیدا کند، باید مقدار آن را در `yyval` قرار دهد و یک کد مبنی بر پیدا کردن یک عدد صحیح برگرداند. در این صورت برای اطمینان از این که `yylex` و `yyparse` با یکدیگر در کد تعریف‌شده موافقت دارند، باید در یک به صورت زیر تعریف شوند:

<i>%token</i>	<i>INTEGER_</i>
---------------	-----------------

لازم به ذکر است نمی‌توان از کلمات رزرو شده در زبان‌های برنامه‌نویسی برای تعریف نشانه‌ها استفاده کرد. برای مثال، در پاسکال، کلمه‌ی رزرو شده‌ی `INTEGER` وجود دارد؛ بنابراین در تعریف نشانه، علامت `_` در آخر آن اضافه شده است. اگر بخواهیم برنامه ماشین‌حساب را بسط داده و متغیرها را از اعداد تشخیص دهیم، می‌توانیم به صورت زیر عمل کنیم:

^۴ Heading

%token	INTEGER_
%token	VARIABLE

۱-۵- ساختار برنامه لکس

ساختار کلی برنامه‌ی ورودی لکس به صورت زیر است:

definitions (تعاریف)

%%

translation (مترجم)

%%

functions (زیربرنامه‌ها)

قسمت‌های تعریف و زیربرنامه‌ها، قسمت‌های اختیاری هستند و می‌توانید آن‌ها را به کار نبرید. در قسمت زیربرنامه‌ها می‌توان کدهای موردنیاز خود را نوشت. این قسمت به انتهای پرونده لکس تولیدشده اضافه می‌شود. در قسمت مترجم، قوانین عباراتی که لکس باید آن‌ها را مقایسه کند، قرار می‌گیرد. در قسمت تعاریف، متغیرها و ساختارهایی قرار می‌گیرد که برای زبان موردنظر شما لازم است. برای مثال، مقدار دادن به yyin در پاسکال.

عبارت‌های حرفی: ساده‌ترین راه برای تشریح عبارت‌های حرفی، لیستی از حروف است که به‌وسیله‌ی "" مشخص و محصورشده باشند. مانند "علی" یا "حسن" یا "۱۲۳۴۵۶۷۸۹" که رشته‌های حرفی نامیده می‌شوند. ابزار لکس، کدهای کنترلی داخل رشته را نیز تشخیص می‌دهد و می‌توان آن‌ها را به‌صورت زیر مشخص کرد:

\n خط جدید

\b یک فاصله به عقب^۵

\t هشت فاصله^۶

... ...

این کدها در عبارت‌هایی که حرف‌های کنترلی و غیرقابل چاپ وجود دارند، می‌توانند مورد استفاده قرار بگیرند. همچنین می‌توان با استفاده از حرف‌های ^ و \$، برای بیان شروع و پایان خط استفاده کرد. به عنوان مثال، "ke" ^ نشان‌دهنده‌ی شروع یک عبارت باقاعده و "\$" ^ abcdefg نشان‌دهنده‌ی شروع و پایان یک خط است.

کلاس‌های حرفی: یک کلاس حرفی به وسیله‌ی تعدادی حرف که در داخل یک جفت براکت [] قرار گرفته‌اند، مشخص می‌شود. مانند [1-9]. این عبارت، بیان‌کننده‌ی کلاس حرفی اعداد است. همچنین برای حروف الفبا و چیزهای دیگر که نیاز است، می‌توان کلاس حرفی تشکیل داد؛ مانند کلاس حرفی حروف الفبای انگلیسی^۷ [A-Za-z]. در کلاس‌های حرفی هر بار فقط یک حرف انتخاب می‌شود. یک عبارت باقاعده که با یک علامت * همراه باشد، می‌تواند صفر بار یا بیشتر تکرار شود.

^۵ Backspace

^۶ Tab

^۷ "-" به مفهوم تا است، یعنی هرچه حرف اسکی بین این دو حرف وجود دارند.

عبارت $[0-9][0-9]^*$ بیان‌کننده‌ی آن است که رشته‌ی مورد نظر، با یک عدد شروع شده و با تعدادی صفر یا بیشتر ادامه پیدا می‌کند. به عبارت دیگر عبارت بالا، الگویی از حرف‌ها برای بیان اعداد است. برای مثال، $[A-Z][a-z]^*$ بیان‌کننده‌ی رشته‌ای است که اول آن، حرف بزرگی است و به دنبال آن تعدادی صفر یا بیشتر حروف کوچک تکرار می‌شوند. بدین ترتیب، برای تعریف متغیرها در زبان سی باید به صورت زیر تعریف کرد:

$[A-Za-z_][A-Za-z_0-9]^*$

این عبارت نشان‌دهنده‌ی یک حرف یا حرف _ است که به دنبال آن، حرف یا عدد یا حرف _ تکرار می‌شود. حرف (عملگر) * نشان‌دهنده‌ی تکرار صفر یا بیشتر است و علامت + بیان‌گر تکرار تعداد یک و یا بیشتر است. برای مثال، $[0-9]^+$ نشان‌دهنده‌ی تکرار یک یا بیشتر از حرف‌های عددی است که معادل با عبارت $[0-9][0-9]^*$ است. همچنین، روش دیگری نیز وجود دارد که می‌توان تعداد تکرارها را نیز مشخص کرد. برای مثال عبارت $\{0-9\}3$ بیانگر این است که عبارت موردنظر باید سه بار تکرار شود.

علاوه بر این، می‌توان برای عبارتی، یک محدوده‌ی تکرار قائل شد. این کار به صورت زیر انجام می‌گیرد:

$\{0-9\}1,9$

عبارت بالا نشان‌دهنده‌ی یک عدد یک تا نه رقمی است. این روش را زمانی به کار می‌برند که استفاده‌کننده نمی‌خواهد عددش از محدوده‌ای بزرگ‌تر شود.

$[A-Za-z_][A-Za-z_0-9]\{0,31\}$

عبارت بالا نشان‌دهنده‌ی تعریف متغیرها در زبان سی است (۱ تا ۳۲ حرف).

عبارت‌های انتخابی: اگر در آخر یک عبارت باقاعده، یک علامت سؤال (?) قرار بگیرد، نشان‌دهنده‌ی اختیاری بودن آن عبارت باقاعده است. برای مثال، $A?$ نشان‌دهنده‌ی صفر یا یک اتفاق از A است؛ یعنی A می‌تواند باشد یا نباشد. اگر در بین دو عبارت باقاعده علامت | یا پایپ قرار داده شود، نشان‌دهنده‌ی این است که یکی از این عبارت‌های باقاعده می‌تواند به کار رود ولی هر دو به کار نمی‌روند^۸؛ مانند $[A-Z][a-z]$ که نشان‌دهنده‌ی حروف بزرگ یا کوچک است.

گروه‌بندی: از حرف () می‌توان برای گروه‌بندی عبارات باقی‌مانده استفاده کرد. برای مثال:

("high"|"medium"|"low")

که نشان‌دهنده‌ی هر سه رشته‌ی low, medium, high است.

توجه: علامت " عمل گروه‌بندی را انجام نمی‌دهد و یک اشتباهی که ممکن است رخ دهد، به صورت زیر است:

"is"?

این الگو نشان‌دهنده‌ی حرف i است که با s اختیاری دنبال می‌شود. برای این که کل رشته را اختیاری کنیم، باید از گروه زیر استفاده کنیم:

("is")?

تشریح قسمت‌های مختلف برنامه لکس

۱-۵-۱- قسمت تعاریف:

این قسمت در ابتدای برنامه‌ی لکس است و برای خوانایی بیشتر برنامه‌ی لکس و همچنین تعریف متغیرهای موردنیاز به کار می‌رود. در این قسمت می‌توان به عبارات باقاعده یک نام اختصاص داد. آخرین خط این قسمت شامل % است که علامت پایان قسمت تعاریف می‌باشد.

یک تعریف لکس شامل انتساب یک الگوی حرفی به یک نام است که قالب کلی آن به صورت زیر است:

^۸ مانند عملگر «یا» در مجموعه‌ها

<i>name</i>	<i>regular_expression</i>
-------------	---------------------------

regular_expression الگویی که نام قسمت name را به خود می‌گیرد، تشریح می‌کند. مانند:

<i>digit</i>	<i>[0-9]</i>
<i>lower</i>	<i>[a-z]</i>
<i>upper</i>	<i>[A-Z]</i>

این سه تعریف به الگوهای مختلف نام اختصاص می‌دهد.
یک تعریف در قسمت تعاریف می‌تواند از تعاریف دیگر با قرار دادن آن نام تعاریف در داخل {} دوباره استفاده کند.
مانند:

<i>letter</i>	<i>{lower}{upper}</i>
---------------	-----------------------

مثال قبل بیان می‌کند که letter شامل الگوی upper و lower که قبلاً تعریف شده‌اند، می‌باشد.

<i>variable</i>	<i>{letter}({letter}/{digit})*</i>
-----------------	------------------------------------

مثال بالا نیز نشان‌دهنده‌ی یک حرف است که به وسیله‌ی صفر یا بیشتر حرف یا رقم دنبال می‌شود.

۱-۵-۲- مترجم:

بعد از %/ که علامت پایان تعاریف است، قسمت مترجم شروع می‌شود. این قسمت توصیف‌کننده‌ی نشانه‌های ورودی و اعمالی که باید بر روی آن‌ها صورت گیرد، است. این قسمت با علامت %/ پایان می‌پذیرد. قالب کلی این قسمت به صورت زیر است:

<i>token-pattern</i>	<i>{action}</i>
----------------------	-----------------

token-pattern یک عبارت باقاعده می‌باشد و ممکن است شامل تعاریفی باشد که پیش از این، در قسمت تعریف معرفی شده است.

action یک سری دستورات به زبان موردنظر است که باید بر روی الگوی موردنظر صورت گیرد و می‌تواند در چندین سطر باشد.

توجه: اگر شما دستورات پاسکال را در قسمت action قرار می‌دهید، بهتر است از (* و *) به جای {} استفاده کنید.
برای مثال:

<i>[/t r]</i>	<i>(* ignored *)</i>
---------------	----------------------

مثال برای قسمت مترجم:

یک ماشین حساب ساده را در نظر بگیرید. قانون قسمت مترجم به صورت زیر خواهد بود:

<i>[0-9]+ {</i>
<i> yylval=atoi(yytext);</i>
<i> return INTEGER;</i>
<i>}</i>

عبارت باقاعده جمله‌ای از اعداد را مقایسه می‌کند و نشانه‌ی پیداشده در داخل yytext قرار می‌گیرد، شما باید مقدار واقعی نشانه پیداشده را در yyval قرار دهید. بعد از این، تبدیل در قسمت action نوع نشانه را برمی‌گرداند. برای مثال در پاسکال:

```
[-/*+]{
    getval := INTEGER (yytext[1]);
    goto yylexReturn; }
```

عبارت بالا شامل چهار حرف عملگر است. اگر yylex یکی از این عملگرها را پیدا کند، در قسمت action اولین حرف موجود در yytext را برمی‌گرداند.

برای پاسکال، با استفاده از زیربرنامه INTEGER() به یک مقدار عددی تبدیل می‌شود. در پاسکال با استفاده از goto yylexReturn از yylex خارج می‌شویم. اگر مقداری برای برگرداندن وجود داشته باشد، باید قبل از goto مقدار آن را در yyval قرار دهیم.

۱-۵-۳- اظهارها^۹

اظهارها می‌توانند با استفاده از ساختار

```
%{      %}
```

به قسمت تعاریف یا مترجم اضافه شوند.

اظهارها جهت تعریف متغیرها مورد استفاده واقع می‌شوند. در قسمت مترجم‌ها، ممکن است اظهارها در ابتدا قرار بگیرند. در این صورت، متغیرهای تعریف‌شده در آن، به صورت محلی در زیربرنامه yylex() به کار گرفته می‌شوند. مانند:

```
%%
%{
    int wordcount=0;
}%
[^\\t\\n] + {wordcount++}
[\\t]+
[\\n]{
    printf("%d\\n",wordcount);
    wordcount=0;
}
```

این برنامه برای شمارش تعداد کلمات در هر خط به کار می‌رود. اگر yylex() تعداد یک یا بیشتر از حرف‌هایی را که حرف فاصله خالی، یا حرف هشت فاصله و یا حرف خط جدید نیستند را پیدا کند، به متغیر wordcount یک عدد اضافه می‌کند و در جملاتی که حرف فاصله خالی یا حرف هشت فاصله وجود دارد، کاری انجام نمی‌دهد و وقتی به حرف خط جدید رسید، تعداد کلماتی که پیدا کرده است را چاپ می‌کند و متغیر wordcount را صفر می‌کند. اظهارها را در بخش تعاریف نیز می‌توان به کار برد. در این صورت، این متغیرها برای زیربرنامه yylex() متغیر خارجی به حساب می‌آیند. مانند:

^۹ Declarations

```
%{
    int characters = 0;
    int words = 0;
    int lines = 0;
}%
%%
\n { ++lines;
    ++ characters;
}
[/t] characters+=yylength;
[^t\n]+{
    ++words;
characters+=yylength;
}
```

توجه: اگر سمت action تنها یک سطر باشد، می‌توان علامت {} را حذف کرد.

در مثال بالا، بخش تعاریف فقط شامل اظهار است و بعد از %/ مترجم می‌آید. در مترجم اول اگر n را پیدا کند، به متغیرهای تعداد خط و متغیر تعداد حرف‌ها یک عدد اضافه می‌کند. در مترجم دوم، اگر فاصله یا t پیدا کند، به متغیر تعداد حرف‌ها به اندازه‌ی yyleng اضافه می‌کند و در مترجم سوم، اگر دنباله‌ای از حرف‌ها را پیدا کند، به متغیر تعداد کلمه یک عدد و به متغیر تعداد حرف به اندازه‌ی yyleng اضافه می‌کند.

از این متغیرها می‌توان در زیربرنامه main() به صورت زیر استفاده کرد:

```
main()
{
    extern int characters, words, lines;
    (void)yylex();
    printf("%d characters, ", characters);
    printf(" %d words, ", words);
    printf(" %d lines \n ", lines);
}
```

در این زیربرنامه ابتدا زیربرنامه yylex() برای تحلیل ورودی فراخوانی می‌شود. سپس جواب‌ها به زیربرنامه اصلی برگردانده می‌شوند و در زیربرنامه اصلی نتایج چاپ می‌شوند. با توجه به زبانی که لکس در آن کد تولید می‌نماید (به عنوان مثال زبان سی)، زبان لکس خاص آن زبان (سی) ایجاد شده است. برای آشنایی بیشتر با زبان‌های مختلف لکس، برنامه‌ای که عملیات ماشین حساب را انجام می‌دهد، به زبان‌های لکس خاص زبان‌های مختلف برنامه‌نویسی در ادامه آورده شده است.

۱-۶- برنامه لکس برای ماشین حساب

برای زبان سی (شکل ۲-۷)

```
%{
#include "ytab.h"
extern int yylval;
}%

%%

[0-9]+{
    yylval = atoi(yytext);
    return INTEGER;
}

[-+/*\n] return *yytext;

[\\t]+;
```

شکل ۲-۷ برنامه ماشین حساب به زبان لکس مختص زبان سی

برای زبان سی++ (شکل ۳-۷)

```
%{
#include "ytab.hpp"
int yylval;
}%

%%

[0-9] + {
    yylval = atoi((char *) yytext);
    return INTEGER;
}

[-+/*\n] return *yytext;

[\\t]+;
```

شکل ۳-۷ برنامه ماشین حساب به زبان لکس مختص زبان سی++

```
%{
unit dc1;
interface

function yylex: integer;
procedure yysetup;
{$I dc1.hp}
var yyval: YYSTYPE;
var code: integer;
var valbuffer: string [40]

implementation

function yywrap: integer;
begin
  yywrap:=1
end;
%}

%%

[0-9]+{
  (* do not handle octal or hex constants (yet)*)
  valbuffer:=Copy(yytext,1,yytext[0]);
  val(valbuffer,yyval,code);
  retval:=CONSTANT;
  goto yylexReturn
}

\n {
  retval:=ord('^');
  goto yylexReturn
}

[-+/*] {
  retval:=Integer(yytext[0]);
  goto yylexReturn
}

[|t|r]+;
%%
procedure yysetup;
begin
```

```

assign(yyin,");
reset(yyin);
assign(yyout,");
rewrite(yyout)
end;

```

شکل ۷-۴ برنامه ماشین حساب به زبان لکس مختص زبان پاسکال

در زیر برنامه yysetup، متغیرهای yyin و yyout مقدار اولیه شده‌اند. در برنامه‌های بالا، زیر برنامه‌هایی برای تشخیص تمام نشانه‌هایی که در یک ماشین حساب وجود دارد، نوشته شده است. در ادامه، برنامه‌ی یک برای این ماشین حساب نوشته شده است.

۷-۱- گرامر یک

معمولاً برنامه‌ی یک را گرامر می‌گویند. یک برای تولید زیر برنامه تحلیل گر نحوی کامپایلرهای زبان‌های برنامه‌نویسی به وجود آمده است. در پرونده ورودی یک ساختار این زبان‌ها تشریح می‌شود و یک برنامه‌ی تحلیل گر نحوی مورد نیاز را در اختیار کاربر قرار می‌دهد. قالب کلی دستورات یک به صورت زیر است:

DECLARATION اظهار

%%

RULES; قوانین

%%

FUNCTIONS زیر برنامه‌ها

تشریح قسمت‌های مختلف ساختار کلی یک

۷-۱-۱- قسمت اظهار

در این قسمت نشانه‌هایی که گرامر مورد نظر را می‌سازند، تعریف می‌شوند. تعریف نشانه‌ها به صورت زیر است:

%token name

name نام نشانه‌ها است. مانند:

%token variable

مثال بالا تعریف نشانه یک متغیر در زبان می‌باشد.

پرونده تعریف نشانه‌ها:

خروجی یک شامل دو پرونده است که در یک پرونده برنامه تحلیل گر نحوی است و یک سرآیند پرونده که شامل تعریف نشانه‌ها است. این سرآیند پرونده باید به برنامه‌ی لکس متصل شود تا لکس بتواند تعریف دقیقی از نشانه‌ها را

داشته باشد. برای انجام این کار، می‌توان به‌طور مستقیم از این سرآیند پرونده در برنامه‌ی لکس در قسمت تعاریف به صورت زیر استفاده کرد:

```
%{
#include "ytab.h"
%}
```

۱-۷-۱-۲- قوانین اولویت

در قسمت اظهار از یک، می‌توانند قوانین اولویت و انجمن‌پذیری نیز قرار بگیرند. برای درک قوانین اولویت به این مثال توجه کنید. در قوانین ریاضی، ضرب و تقسیم از جمع و تفریق اولویت بیشتری دارند؛ مگر این‌که از پرانتز برای تغییر این اولویت استفاده شود. هم‌چنین، این دو نسبت به یکدیگر اولویت یکسانی دارند ولی هردوی آن‌ها از جمع و تفریق اولویت بیشتری دارند.

عبارات زیر را در زبان سی در نظر بگیرید:

$$A = 2 - 2 - 1;$$

$$A = B - 7 * 2;$$

برای محاسبه‌ی عبارت اول، محاسبه از چپ به راست انجام می‌شود:

$$A = ((2 - 2) - 1);$$

ولی در عبارت دوم، ابتدا عمل ضرب انجام می‌شود و سپس عملیات تفریق؛ یعنی محاسبه از راست به چپ انجام می‌شود:

$$A = (B - (7 * 2));$$

چون ضرب اولویت بیشتری نسبت به تفریق دارد، ضرب ابتدا انجام می‌شود و نتیجه از B کم می‌شود و سپس نتیجه به A اختصاص داده می‌شود. این کار به وسیله‌ی قوانین اولویت در یک صورت می‌گیرد. عملیاتی که از راست به چپ انجام می‌گیرند را انجمن‌پذیری از راست و عملیاتی که از چپ به راست انجام می‌گیرند را انجمن‌پذیری از چپ می‌نامند. برای مثال عبارت $a/b/c$ ممکن است به دو صورت مختلف پردازش شود.

اگر عملیات تقسیم، انجمن‌پذیری از راست داشته باشد، عملیات به‌صورت زیر انجام می‌شود:

$$a / (b / c)$$

و اگر عملیات تقسیم، انجمن‌پذیری از چپ داشته باشد، عملیات به‌صورت زیر انجام می‌شود:

$$(a / b) / c$$

می‌توان انجمن‌پذیری و اولویت عملیات را در قسمت اظهار به‌صورت زیر مشخص کرد:

```
%left  operator
%right operator
```

در یک خط می‌توانند چند عملگر بیابند که همه‌ی آن‌ها دارای اولویت یکسانی می‌شوند. مانند:

```
%left  '+', '-'
```

```
%left  '*, '/'
```

left و right به ترتیب نشان‌دهنده‌ی انجمن‌پذیری از چپ و انجمن‌پذیری از راست هستند. عملگرهایی که در یک سطر تعریف می‌شوند، از عملگرهایی که در سطرها‌ی بالا تعریف می‌شوند، اولویت بیشتری دارند. این عبارت‌ها بیان‌کننده‌ی این مطلب هستند که * و / دارای اولویت بیشتری نسبت به + و - می‌باشند.

۱-۷-۳- تعاریف زبان

در قسمت اظهار از گرامر یک، می‌توان تعاریف موردنظر خود را قرارداد. همان‌طور که در لکس این کار را با استفاده از $\{\% \}$ انجام می‌دادیم.

۱-۷-۴- قوانین

یک قانون، یک ساختار گرامری درست را تشریح می‌کند. این ساختار گرامری برای تشخیص نشانه‌های ورودی یا ساختار گرامری و یا ترکیبی از آن دو به کار می‌رود. برای درک این مطلب به قانون ماشین حساب توجه نمایید:

expression: INTEGER
expression: expression '+' expression
expression: expression '-' expression
expression: expression '/' expression
expression: expression '' expression;*

این قوانین ساختار گرامری را تشریح می‌کنند که *expression* نامیده می‌شود. ساده‌ترین عبارت یک نشانه *integer* است که با *INTEGER* شناسایی می‌شود. ممکن است عبارات پیچیده‌تری را با اضافه کردن جمع و ضرب و تقسیم و تفریق ایجاد کنید. یک قانون مانند:

expression: expression '+' expression

دارای سه جزء است؛ یک *expression* و یک نشانه "+" و *expression* دیگر. اگر یک برنامه از گرامر قبلی برای تشخیص عبارت $3 + 2 + 1$ استفاده کند، اول از همه، عدد ۱ وارد می‌شود که یک نشانه *INTEGER* است. بنابراین می‌تواند به یک *expression* تبدیل شود. در نتیجه ورودی به صورت زیر می‌شود:

expression + 2 + 3

البته ۲ نیز یک *INTEGER* است. بنابراین یک *expression* است و به شکل زیر تبدیل می‌شود:

expression + expression + 3

برنامه، قسمت اول ورودی را به صورت عبارت درستی از *expression* تشخیص داده است. بنابراین به صورت زیر تبدیل می‌شود:

expression + 3

به روش مشابه عبارت بالا نیز به *expression* تبدیل می‌شود.

۱-۷-۱-۵- عمل^{۱۰}

در قسمت قانون گرامر یک، فقط ساختار گرامر تشریح نمی‌شود؛ بلکه بیان می‌شود وقتی ساختاری تشخیص داده شد چه کاری انجام شود. به عبارت دیگر، می‌توان قانون و عمل را با هم به کار برد. قالب کلی یک قانون به صورت زیر است:

name: definition{action};

name نام ساختار تعریف شده است.

definition تعریف ساختار گرامری با جملاتی از نشانه‌ها و یا ناپایانه‌ها است.

action شامل صفر یا بیشتر دستورالعمل است که وقتی ورودی به صورتی باشد که در definition تعریف شده است، انجام می‌شود.

در یک متغیر \$1 برای مقدار جزء اول، متغیر \$2 برای مقدار جزء دوم، \$3 برای مقدار جزء سوم و \$\$ برای نمایش مقدار ساختار تعریف شده به کار می‌رود. برای محاسبه جمع، جزء اول و جزء سوم با هم جمع می‌شوند.

expression : expression '+' experssion { \$\$ = \$1 + \$3;};

برای محاسبه‌ی تمام عبارت می‌توان به صورت زیر آن را بیان کرد:

expression: expression '-' expression { \$\$ = \$1 - \$3;};

expression: expression '' expression { \$\$ = \$1 * \$3;};*

expression: expression '/' expression { \$\$ = \$1 / \$3;};

expression: INTEGER_ { \$\$ = \$1;};

قانون آخر بیان کننده‌ی این مطلب است که اگر عبارتی دارای یک نشانه INTEGER_ بود مقدار عبارت فقط مقدار نشانه می‌باشد.

اگر عمل مشخصی در قانون انجام نگیرد، عمل تعریف شده برابر با { \$\$ = \$ 1; } خواهد بود. این جمله بیان گر این است که مقدار پیش فرض در یک ساختار، مقدار اولین جزء است.

مثال بالا برای زبان سی ++ و سی بود. در زیر مثالی برای پاسکال آورده شده است:

expression: expression '-' expression { \$\$ = \$1 - \$3;};

expression: expression '' expression { \$\$ = \$1 * \$3;};*

expression: expression '/' expression { \$\$ = \$1 div \$3;};

expression: INTEGER_ { \$\$ = \$1;};

۱-۷-۱-۶- فشرده سازی قوانین

اگر یک قانون دارای شکل های مختلفی در یک ساختار باشد، می‌توان آن را به شکل زیر ترکیب کرد:

name:

definition1 {action}

/definition2 {action}

```

/definon3 {action}
/definition4 {action}
/...
;

```

توجه کنید که نقطه ویرگول^{۱۱} انتهای قانون را مشخص می کند و هر قانون دارای action مربوط به خود است. با این تعریف مثال بالا را می توان به صورت زیر نوشت:

```

expression:
    INTEGER_      { $$ = $1; }
    / expression '-' expression { $$ = $1 - $3; }
    / expression '+' expression { $$ = $1 + $3; }
    / expression '*' expression { $$ = $1 * $3; }
    / expression '/' expression { $$ = $1 / $3; }

```

۱-۷-۱-۷- نماد آغازگر^{۱۲}

اولین ساختار گرامری تعریف شده در قسمت قوانین، باید شامل ساختار کلی باشد. برای مثال، اگر ورودی یک، گرامر یک زبان باشد، قانون اول باید شامل یک برنامه ی کامل باشد. نام این قانون اول را نماد آغازگر می نامند. yyparse() ورودی ها را در یک توصیف مناسب جمع آوری می کند. اگر گرامر یک زبان برنامه نویسی را تشریح کند و نماد آغازگر، یک برنامه کامل را توضیح دهد، yyparse() بعد از انجام عملیات متوقف می شود. نماد آغازگر باید طوری تعریف شود که همه ی ورودی های معتبر را شامل شود. در مثال ماشین حساب، باید یک برنامه با قانون زیر تعریف شود:

```

Program :
/*null*/
/program expression '\n'
;

```

این دو قانون تعریف برنامه را بیان می کنند که می تواند شامل یک عبارت باشد که یک حرف خط جدید بعد از آن می آید (سطری که باید محاسبه شود)، و به دنبال آن، چنین خط های بیشتری می تواند بیاید یا اینکه می تواند شامل چیزی نباشد. تعریف «هیچی» در اولین قانون آمده است.

۱-۷-۱-۸- تعریف کد استفاده کننده^{۱۳}

در قسمت قوانین می توان متغیرهایی که در داخل yyparse() مورد استفاده قرار می گیرند را قرار داد. این کار به وسیله ی % { } صورت می گیرد. این متغیرها به صورت محلی در زیر برنامه yyparse مورد استفاده قرار می گیرند.

۱-۷-۱-۹- زیر برنامه ها:

این قسمت در گرامر یک اختیاری است و اغلب استفاده نمی شود. در این قسمت زیر برنامه های کاربر قرار می گیرند که عیناً در انتهای زیر برنامه ی تولید شده توسط یک قرار می گیرد.

^{۱۱} semicolon

^{۱۲}Start Symbol

^{۱۳}User code declaration

۸-۱- برنامه‌ی یک برای ماشین حساب

همانند لکس، برای هر زبان برنامه‌نویسی متداول مانند سی، زبان یک مختص آن زبان برنامه‌نویسی (سی) ایجاد شده است. در ادامه، برنامه مربوط به یک ماشین حساب را در زبان‌های مختلف یک بررسی می‌کنیم.

برنامه برای زبان سی (شکل ۵-۷)

```
%token INTEGER
%left '+' '-'
%left '*' '/'

%%

program:
    program expression '\n' = {printf("%d\n", $2);}
    / /*NULL*/
;

expression :
    INTEGER
    / expression '+' expression = { $$ = $1 + $3; }
    / expression '-' expression = { $$ = $1 - $3; }
    / expression '*' expression = { $$ = $1 * $3; }
    / expression '/' expression = { $$ = $1 / $3; }
;
```

شکل ۵-۷ برنامه ماشین حساب به زبان یک مختص زبان سی

برنامه برای زبان سی++ (شکل ۶-۷)

```
%{
#include "dc11.hpp"
#include "dc1.hpp"
#ifdef _ZTC_
#include <stream.hpp>
#else
#include <iostream.h>
#endif
}%

%token INTEGER
%left '+' '-'
```

```

%left '*' '/'

%%

program:
    program expression '\n' {cout << $2 << "\n";}
/      /*NULL*/
;

expression:
    INTEGER
/ expression '+' expression { $$ = $1 + $3; }
/ expression '-' expression { $$ = $1 - $3; }
/ expression '*' expression { $$ = $1 * $3; }
/ expression '/' expression { $$ = $1 / $3; }
;
%%
main()
{
yy_scan scan;
yy_parse parse;
return parse.yyparse(&scan);
}

```

شکل ۶-۷ برنامه ماشین حساب به زبان یک مختص زبان سی++

برنامه برای زبان پاسکال (شکل ۷-۷)

```

%token CONSTANT
%left '+' '-'
%left '*' '/'
%{
program dc1;
uses dc1l;

procedure yyerror (msg: string);
begin
    writeln(msg);halt(1)
end;
%}

%%

program:
    program expression '\n' { writeln($2);}

```

```
/  /*NULL*/  
;  
  
expression:  
    CONSTANT  
/  expression '+' expression  { $$:= $1 + $3; }  
/  expression '-' expression  { $$:= $1 - $3; }  
/  expression '*' expression  { $$:= $1 * $3; }  
/  expression '/' expression  { $$:= $1 div $3; }  
;  
%%  
begin  
yysetup;  
{ $IFDEF YYDEBUG }  
yydebug:= 1;  
{ $ENDIF }  
halt(yyvsparse)  
end.
```

شکل ۷-۷ برنامه ماشین حساب به زبان یک مختص زبان پاسکال