

دانشگاه آزاد اسلامی واحد هشتروند

کامپایلر ها

اصول، ابزار ها، روش ها

مؤلف: آلفردوی. آهو - راوی سدی - جفری دی. المن

مدرس: مهندس باقری نیاء

تهیه کننده: سدی همتی

ترجمه:

عبارت است از تبدیل یک ساختار از یک زبان (زبان مبدا) به زبان دیگر (مقصد).

مترجم:

برنامه ای است که عمل ترجمه از زبان مبدا به زبان مقصد را انجام می دهد.

انواع مترجم:

□ **کامپایلر:** اگر زبان مبدا یکی از زبان های سطح بالا (پاسکال، فرترن PL/1) و زبان مقصد زبان ماشین باشد به چنین مترجمی کامپایلر گویند.

□ **مفسر:** مفسر نیز یک زبان سطح بالا را دستور به دستور ترجمه کرده و اجرا می کند و زبان ماشین تولید نمی کند.
نکته ۱: فرق کامپایلر با مفسر در این است که کامپایلر کد قابل اجرا تولید می کند ولی برای اجرای یک برنامه توسط مفسر هر بار می بایست ترجمه و اجرا صورت گیرد..

نکته ۲: حجم مفسر از حجم کامپایلر کمتر است ولی سرعت کامپایلر بیشتر است
نکته ۳: از لحاظ مصرف حافظه بهتر است از مفسر استفاده شود، چون مفسر فضای کمتری اشغال می کند و از لحاظ امنیت کامپایلر بهتر است زیرا کل برنامه را یک جا ترجمه می کند..

□ **اسمبلر:** ورودی اسمبلر زبان اسمبلی (کد های یادمان) می باشد که آن را به زبان ماشین تبدیل می کند. اسمبلر های اولیه هر دستور زبان اسمبلی را تبدیل به یک دستور زبان ماشین می کردند که جهت سرعت بخشیدن به عمل ترجمه ، اسمبلر هائی ابداع گردید که قادرند یک دستور اسمبلی (دستور کلان) را تبدیل به چندین دستور زبان ماشین کنند.

□ **پیش پردازنده:** یک زبان سطح بالا را تبدیل به یک زبان سطح بالای دیگری می کند که اغلب زبان مبدا شکل توسعه یافته زبان است مثل Turbo C، و زبان مقصد شکل استاندارد زبان است مثل ANSI C

مراحل کامپایل در شش مرحله صورت می گیرد.

۱- تحلیل لغوی (Lexical analysis)

۲- تحلیل نحوی (syntax analysis)

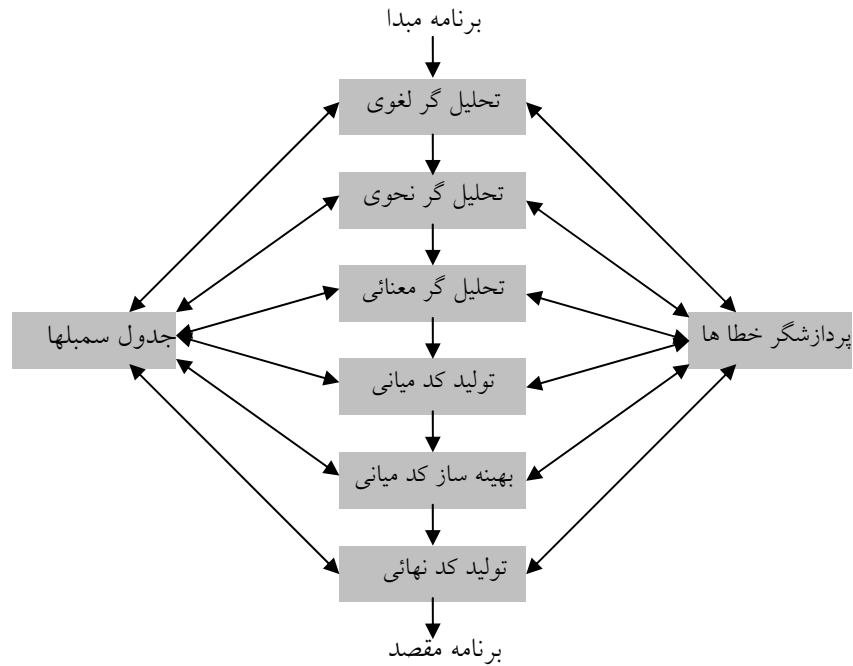
۳- تحلیل معنایی (semantic analysis)

۴- تولید کد میانی (intermediate code generation)

۵- بهینه سازی کد (code optimization)

۶- تولید کد نهائی (code generation)

این شش بخش به همراه بخش های "پردازش خطا" و "جدول سمبل ها" یک کامپایلر را تشکیل می دهند ارتباط بین این بخش ها در زیر نشان داده شده است.



تحلیلگر لغوی:

در یک کامپایلر به تحلیل خطی (Linear Analysis)، تحلیل لغوی (lexical Analysis)، تحلیل پوی (scanning Analysis) گفته می شود. کاراکترها را از چپ به راست خوانده و با رسیدن به یک جدا کننده آنها را در گروه های منطقی به هم مرتبط به نام Token قرار می دهد و Token ها که خروجی تحلیل گر لغوی می باشند در جدول سمبلها با فرمت خاصی ذخیره می شوند و همچنین به عنوان ورودی تحلیل گر نحوی استفاده می شوند به عبارتی تحلیل گر لغوی واسط بین برنامه مبدا و تحلیلگر نحوی است.

انواع مختلف نشانه ها (Tokenها) عبارتند از

کلمات کلیدی (keywords)

عملگرها (operators)

ثابتها (literals)

شناسه ها (identifiers) که به اسامی متغیرها، توابع، رویه ها و به طور کلی اسامی که کاربر انتخاب می کند گفته می شود.

$\left. \begin{array}{l} \text{space} \\ \text{tab} \\ \text{ } \\ \vdots \end{array} \right\} \text{ جداکننده ها (delimiters)}$

ساختار Token برای $a=b+2$

Token type	Token value
id	a

نکته:

تحلیلگر لغوی commentها (توضیحات) را رد می کند و فاصله های خالی (white space) را که کاراکترهای شناسه ها را که از هم جدا می کند را حذف می کند و اگر خطائی رخ بدهد آن را گزارش می دهد. مثلا با دیدن عبارت "6a" تحلیل گر لغوی خطائی را گزارش می کند.

تحلیلگر نحوی:

تحلیل سلسله مراتبی (Hierarchical analysis)، تجزیه (parsing analysis) یا نحوی (syntax analysis) نامیده می شود. تحلیلگر نحوی با داشتن گرامر مربوط به زبان و Token های دریافتی از تحلیلگر لغوی برای یک دستور درخت اشتقاق آن را ساخته و اگر دستور مربوط به زبان نباشد خطا گزارش می دهد وگرنه درخت اشتقاق را به عنوان خروجی به تحلیلگر معنایی می دهد. مثلا در زبان C اگر دستوری به شکل $for(i=0; i++)$ داشته باشیم خطای نحوی رخ میدهد.

تحلیلگر معنایی:

مهمترین وظیفه تحلیلگر معنایی کنترل نوع یا (Type checking) است. اگر تبدیل نوع مجاز باشد عمل تبدیل نوع را انجام می دهد وگرنه خطای تبدیل نوع را گزارش میدهد
تبدیل نوع $float + int \leftarrow$
خطا $float + char \leftarrow$

تولید کننده کد میانی:

در این بخش برنامه ورودی به یک زبان میانی تبدیل می شود. میتوان این زبان میانی را به برنامه ای برای یک ماشین انتزاعی تشبیه کرد. این زبان میانی حداقل بایستی خواص ذیل را دارا باشد.

سهولت تولید کد: زبان ماشین منطقی تا حد امکان باید ساده در نظر گرفته شود تا تولید کد برای آن آسان باشد.

سهولت ترجمه به زبان مقصد: زبان ماشین منطقی باید به گونه ای باشد که تبدیل آن به زبان ماشین آسان باشد.

بهینه سازی کد میانی:

در این بخش سعی می شود تا کد میانی به صورتی در آید که سریع تر اجرا شود مثلا در کد $a = func(b) + func(b)$ به جای دو بار فراخوانی تابع $func()$ این تابع با پارامتر b یک بار فراخوانی شده و مقدار برگشتی آن در یک متغیر موقت ریخته می شود و سپس از مقدار متغیر موقت استفاده می شود.

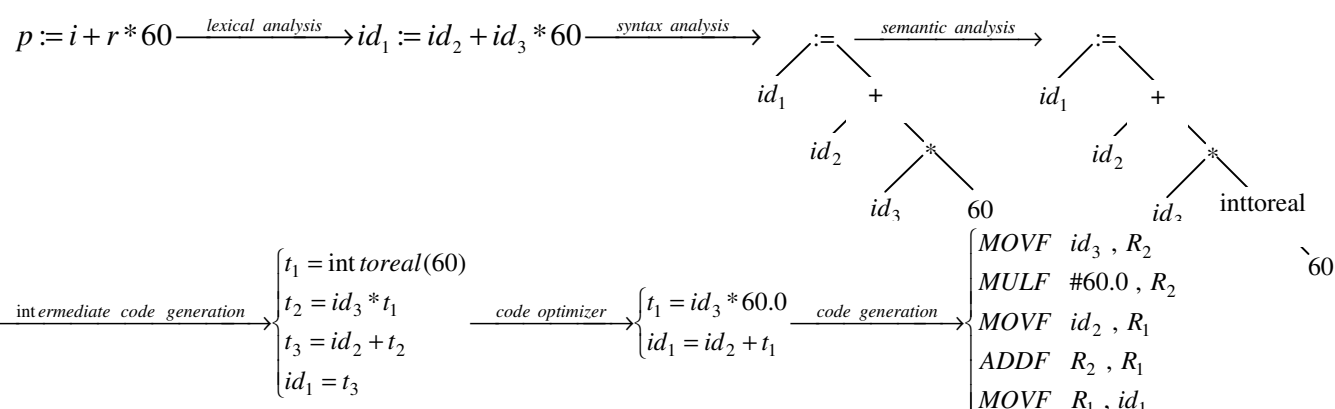
دلیل این که عمل بهینه سازی همراه با کد میانی انجام نمی شود چیست؟

جواب- عمل تولید کد میانی همراه با درخت انجام می شود و کار کردن با درخت مشکل است ولی بهینه ساز کد، یک فایل را به عنوان ورودی می گیرد.

تولید کننده کد نهایی:

در این مرحله کد اسمبلی قابل اجرا بر روی ماشین ایجاد می شود. در واقع کد تهیه شده میانی که در آن از ثبات استفاده نشده می بایست به کدی تبدیل شود که در آن از ثبات استفاده می شود.

به عنوان مثال شکل زیر مراحل کامپایلر $p = i + r * 60$ و تاثیر هر مرحله از کامپایلر بر روی این دستور را نشان می دهد (فرض شده است که i, p, r همگی از نوع $real$ هستند)



نکات تکمیلی جلسه اول

- بخش front-end یک کامپایلر برای تمام ماشین ها یکسان است ولی back-end برای ماشین های مختلف متفاوت است
- اگر کاراکتر غیر مجاز در متن برنامه نوشته شود تحلیل لغوی (scanner) آن را تشخیص میدهد (علوم کامپیوتر-۷۹)
- در فرایند کامپایلر یک برنامه توکن ها در مرحله آنالیز لغوی (Lexical) شناسائی میشوند. (مهندسی کامپیوتر آزاد-۷۱)
- در مرحله تحلیل لغوی کلیه شناسه های موجود در برنامه وارد جدول علائم (symbol Table) میشود. (مولف-راهیان ارشد)
- back-end بخشی از مرحله بهینه سازی و مرحله تولید کد نهائی که وابسته به ماشین است، می باشد. (مولف-راهیان ارشد)
- تفاوت بین compiler و preprocessor در این است که کامپایلر تحلیل لغوی و معنایی را روی برنامه انجام میدهد، در حالی که preprocessor این کار را انجام نمیدهد (مولف-راهیان ارشد)
- دلیل استفاده از front-end و back-end در کامپایلر ها، تدارک ترکیب چند front-end و back-end در یک خانواده کامپایلر است (مولف-راهیان ارشد)
- در صورتی که شناسه قبلا اعلان نشده باشد، یک خطای معنایی رخ میدهد. (نوپردازان-مولف)
- کشف خطاهای مربوط به ساختار تک تک لغات وظیفه تحلیل گر لغوی است (کتاب-پیام نور)
- در دستور `var 7temp := integer;` که به زبان پاسکال می باشد، تحلیل گر لغوی لغت `7temp` را به عنوان خطا گزارش میدهد زیرا در پاسکال یک شناسه نباید با عدد شروع شود.
- در زبان پاسکال عبارت، `A B :=` به دلیل عدم رعایت ترتیب صحیح انتساب دارای خطای نحوی می باشد.
- در زبان پاسکال عبارت، `if(a = b then` به دلیل عدم توازن پرانتزها دارای خطای نحوی می باشد.
- تحلیل گر معنایی معنی دار بودن عباراتی که از نظر نحوی درست بوده اند را مورد بررسی قرار میدهد.

دو مزیت عمده کد میانی عبارت است از:

- کامپایلر را می توان مستقل از ماشین نوشت لذا با تغییر ماشین ها و تکنولوژی آنها کافی است تنها بخش مولد کد تغییر کند و نه دیگر بخش های کامپایلر.
- یک بهینه ساز کد مستقل از ماشین را می توان برای سریع تر کردن کد میانی استفاده کرد، در نتیجه بدون توجه به نوع ماشین، کد تا حد ممکن بهینه خواهد بود.

به طور کلی کدهای میانی شامل انواع زیر می باشد :

- ۱- کد شبه اسمبلی
 - ۲- درخت خلاصه نحوی
 - ۳- جملات سه آدرسه
- با یک پیمایش پسوندی به سادگی می توان از درخت خلاصه نحوی کد اسمبلی تولید کرد.

وظایف تحلیل گر نحوی:

- ۱- بررسی صحت و درستی ترتیب لغات برنامه (بررسی ساختار برنامه مبدا)
- ۲- بررسی جریان کنترل (Flow of control checks)، مثلاً عبارت Break در زبان C موجب خروج از نزدیکترین حلقه (while، for، switch) می گردد که در صورت نبودن چنین حلقه ای خطا رخ میدهد.
- ۳- کنترل منحصر به فرد بودن، مثلاً نام یک Label می بایست دقیقاً یک بار در برنامه ظاهر شود.
- ۴- چک کردن نام های مرتبط، مثلاً در زبان Ada یک حلقه دارای یک نام در ابتدا و انتهای ساختار خود است که می بایست یکسان باشند.
- ۵- چک کردن ساختار های تودرتو

برخی از مواردی که توسط تحلیل گر معنایی انجام میشود عبارتند از:

بررسی هماهنگی پارامترها: فراخوانی یک روال باید با تعریف روال هماهنگی داشته باشد. مثلاً در زبان C تعداد پارمتر های ارسالی، نوع آنها و ترتیب آنها باید در تعریف تابع هماهنگی داشته باشد.

بررسی و کنترل نوع: در این قسمت نوع عملوند های یک عملگر مورد بررسی قرار می گیرد. مثلاً در $a[1.5] := 12$ خطای معنایی وجود دارد زیرا اندیس یک آرایه (a) نبایستی یک عدد اعشاری باشد.

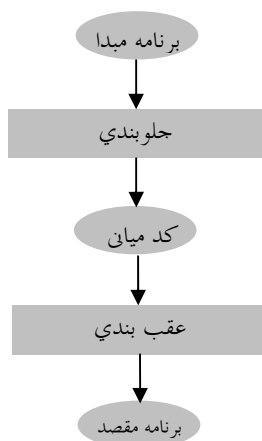
تبدیل نوع: عملیات جمع برای اعداد صحیح و اعشاری متفاوت است برای انجام عمل جمع در عبارت $c = a + b$ در صورتی که c و b از نوع float و a از نوع int باشد، تحلیلگر معنایی عمل ارتقاء نوع را انجام میدهد. بدین معنا که موقتا متغیر a به نوع float تبدیل میشود.

تعریف دوباره متغیر: تحلیلگر معنایی بررسی می کند تا هیچ متغیری دو بار تعریف نشده باشد

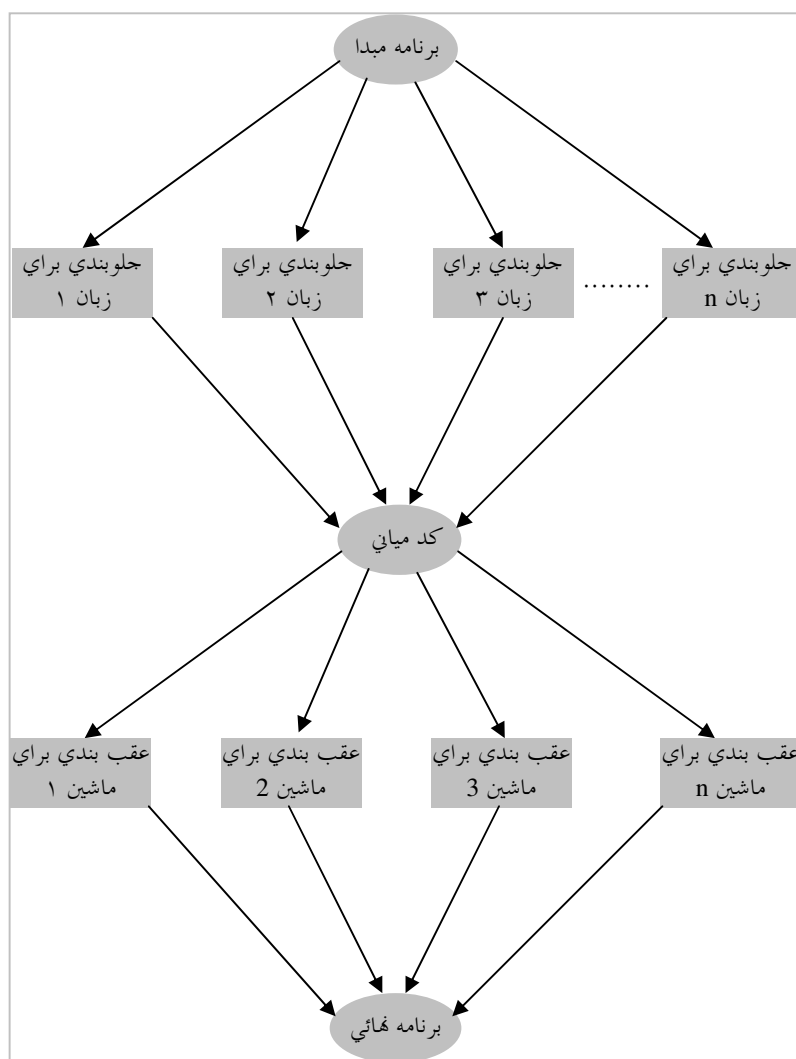
جلسه دوم

جلوبندی (Front End) و عقب بندی (Back End) کامپایلر

می خواهیم برنامه ای را به زبان C روی انواع مختلف کامپیوتر (مانند Vax, Mainfram, IBM) اجرا کنیم در این صورت برای هر نوع کامپیوتر، باید یک کامپایلر جداگانه بسازیم. اگر n تعداد زبان های برنامه سازی (مانند C و پاسکال و...) و K تعداد انواع مختلف کامپیوترها باشد در این صورت به nk کامپایلر نیاز است. ایجاد این تعداد کامپایلر بسیار زمانبر و پرهزینه است. برای حل این مشکل از تقسیم کامپایلر به جلو بندی و عقب بندی استفاده می کنیم. به این شکل که یک زبان میانی در نظر می گیریم در ابتدا برنامه مبدا را به این زبان میانی ترجمه کرده سپس از زبان میانی به زبان مقصد ترجمه می کنیم. بخشی از کامپایلر که وظیفه ترجمه برنامه مبدا به برنامه زبان میانی را بر عهده دارد و وابسته به زبان ماشین نیست را جلو بندی و بخشی از کامپایلر که وظیفه ترجمه برنامه از زبان میانی را به زبان مقصد را بر عهده دارد و وابسته به زبان ماشین است را عقب بندی کامپایلر می گوئیم (شکل زیر).



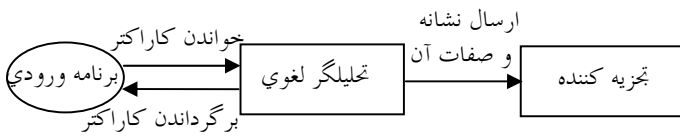
با استفاده از این روش برای n زبان مبدا و k کامپیوتر مختلف به n جلو بندی و k عقب بندی نیاز است که در مجموع $n+k$ برنامه (کامپایلر) احتیاج است (شکل زیر این موضوع را نشان می دهد).



تقسیم بندی کامپایلر به عقب بندی و جلو بندی چه مزایایی دارد؟

- سادگی طراحی
- استقلال جلو بندی از زبان مقصد
- استقلال عقب بندی از زبان مبدا
- کاهش پیچیدگی
- افزایش قابلیت استفاده مجدد
- افزایش سرعت تولید کامپایلر برای سخت افزار جدید و زبان های جدید.

تحلیگر لغوی (scanner): واسط بین کامپایلر و برنامه مبدا می باشد و مهمترین وظیفه آن خواندن برنامه ورودی به صورت کاراکتر به کاراکتر و تشخیص نشانه ها یا token ها میباشد (شکل زیر نحوه قرار گرفتن تحلیگر لغوی بین ورودی و تجزیه کننده را نشان می دهد).



Scanner میتواند به صورت هم روال یا زیرروال با parser پیاده سازی شود، Scanner به محض درخواست token بعدی از parser، توکن بعدی را به پارسر می دهد.

وظایف دیگر scanner

- حذف فضا های خالی و توضیحات (comment)
- گزارش خطا ها

از آنجا که scanner برنامه را به صورت کاراکتر به کاراکتر می خواند، می تواند تعداد کاراکتر های هر خط و به تبع آن شماره خطوطی را که توکن ها در آن قرار دارد را نیز مشخص کند. بنابر این هنگام اعلام خطا شماره خطی را که خطا رخ داده است را نیز گزارش می کند. به چه دلایلی بهتر است که Scanner و parser به صورت مجزا پیاده سازی شوند؟

- سادگی
- کارایی بالاتر
- قابلیت حمل

سادگی:

پارسری که دربردارنده قواعد مربوط به حذف فضا های خالی و توضیحات می باشد پیچیده تر از پارسری است که فرض می کند این اعمال را scanner انجام می دهد.

کارایی بالاتر:

خواندن کاراکتر به کاراکتر یک عمل ورودی است قرار دادن این اعمال در Scanner و در نظر گرفتن تمهیداتی برای بالا بردن سرعت در scanner سبب افزایش کارایی می شود.

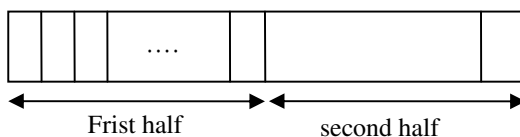
قابلیت حمل:

اعمال قرارداد های مربوط به کاراکتر های ویژه (مثلا علامت \wedge در پاسکال) در scanner و همین طور اعمال مدیریت ورودی وابسته به هر زبان در scanner قابلیت حمل تحلیگر را بالا می برد.

□ وقت گیرترین کار کامپایلر همان کاری است که scanner انجام می دهد پس باید سرعت را بالا ببریم، باید از بافر استفاده شود چون بافر باعث بالا رفتن سرعت می شود.

افزایش سرعت:

تحلیگر لغوی یا scanner تنها فازی از کامپایلر است که عمل ورودی را انجام می دهد بنابراین به دلیل این که به ازای هر کاراکتر می بایست یک عمل ورودی انجام داد استفاده از روشی جهت افزایش سرعت مفید می باشد این روش استفاده از بافر می باشد.



هر یک از انواع بافر ها دارای دو نوع اشاره گر مکی باشند
Forward: با خواندن هر کاراکتر یک واحد جلو می رود.

Lexeme_beginning: ابتدای توکنی را که می خواهیم مشخص کنیم نشان می دهد

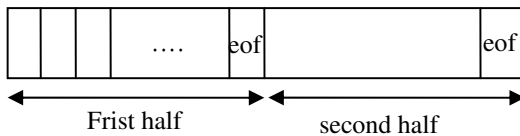
کد مربوط به پیشروی اشاره گر پیش رو (forward) به شکل زیر خواهد بود.

```
if forward at the end of frist half then begine
reload second half;
forward = forward + 1
end
else if forward at the end of second half then begine
reload frist half;
move forward to beginning of frist half
end
else forward = forward + 1
```

به ازای خواندن هر کاراکتر به دو تا مقایسه نیاز داریم می خواهیم تعداد مقایسه ها را کمتر کنیم (یک مقایسه) راه حل آن است که از نگهبان استفاده کنیم

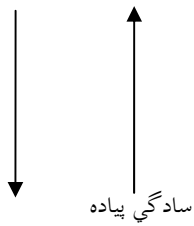
کاراکتر نگهبان: eof

کاراکتر نگهبان انتهای هر نیمه بافر قرار می گیرد و کاراکتر ویژه ای است که جزء متن نیست.



کد مربوط به پیشروی اشاره گر پیش رو (forward) همراه با کاراکتر نگهبان به شکل زیر خواهد بود.

```
forward = forward + 1;
if forward ↑ = eof then begine
if forward at the end of frist half then begine
reload second half;
forward = forward + 1
end
else if forward at the end of second half then begine
reload frist half;
move forward to beginning of frist half
end
else /* eof within a buffer signifying end of input */
end
```



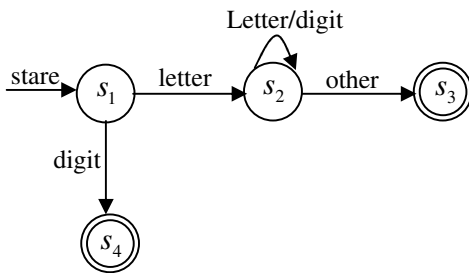
۲- پیاده سازی به کمک زبان برنامه سازی سیستم و استفاده از امکانات مدیریت ورودی در آن زبان

۳- پیاده سازی به کمک زبان اسمبلی و مدیریت مستقیم ورودی

□ به ابزار هائی که خودشان کامپایلر تولید می کنند کامپایلر کامپایلرها می گویند مثل **lex** در سیستم عامل **unix**

□ الگوهای توکن ها را به وسیله عبارت منظم نمایش می دهند و هر عبارت منظمی قابل تبدیل به **dfa** است.

شکل زیر نمودار تغییر حالت برای شناسه ها را نشان می دهد، قاعده تشخیص شناسه عبارت است از دنباله ای از حروف و ارقام که اولین نماد آن حرف باشد.



حال **table** ها را مشخص می کنیم.

جدول **char-class**

class	letter	letter	digit	other
char	a-z	A-Z	0-9	other

جدول **state**

state	letter	digit	other
s_1	s_2	s_4	s_4
s_2	s_2	s_2	s_3
s_3	-	-	-
s_4	-	-	-

< token type , token value >

هر **token** از دو مولفه تشکیل می شود یکی نوع **token** و دیگری مقدار **token**

بعد از یافتن هر توکن می بایست توکن در جدول نمادها ذخیره شود.

در جدول نمادها نوع و نام و شماره اولین سطر ذخیره می شود.

جلسه سوم

دیاگرام تشخیص عملگرهای رابطه ای

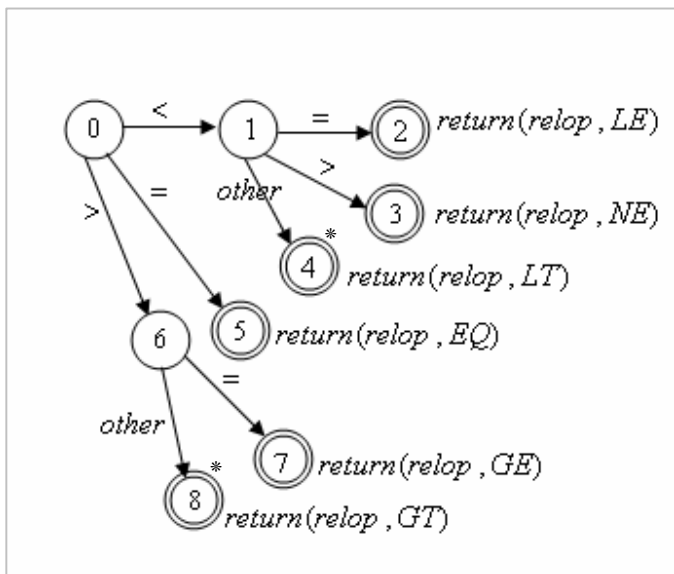
عملگرهای رابطه ای $>$, $>=$, $=$, $<$, $<=$

L : lese

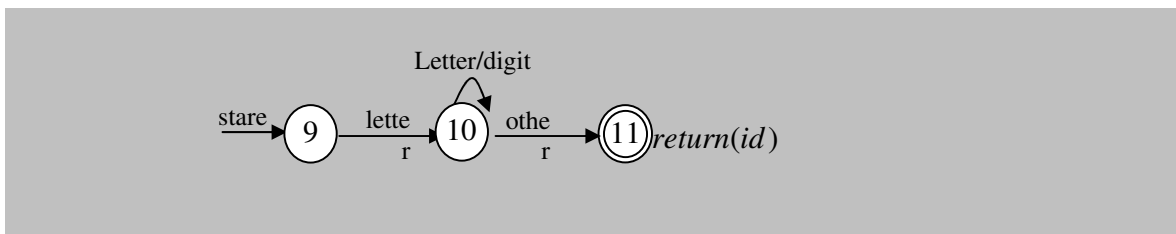
E : Equal

G : grather

کاراکتر خوانده شده بایستی به بافر برگردانده شود.



دیاگرام تشخیص شناسه (identifire).



دیاگرام تشخیص اعداد

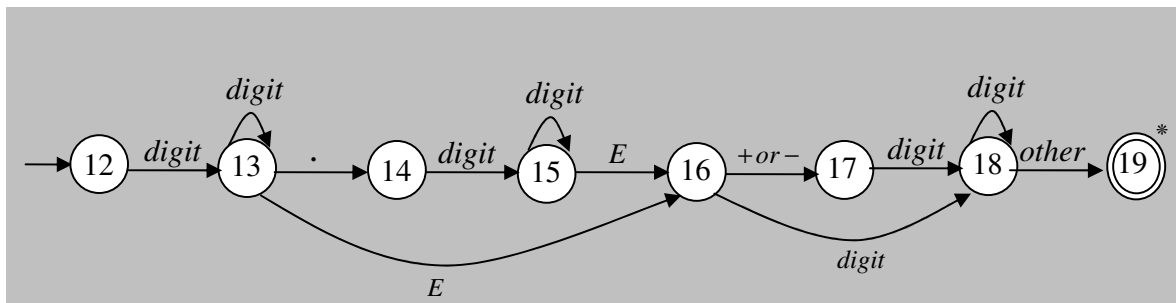
تشخیص اعداد به صورت نماد علمی

$digit \cdot digit E (+/-) digit$ 12.34E+1

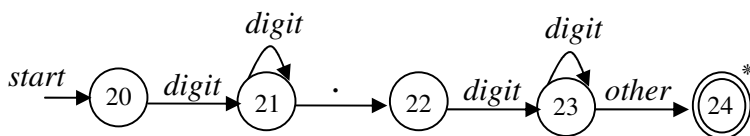
$digit \cdot digit E digit$ 12.34 E1

$digit E (+/-) digit$ 12E+1

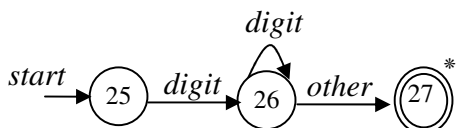
$digit E digit$ 12E1



تشخیص اعداد به صورت اعشاری
digit.digi



تشخیص اعداد به صورت صحیح



برنامه به زبان C برای پیدانمودن حالت شروع بعدی

```
int state=0 , start=0;
int fail()
{
switch (state) {
case 0 :    start = 9;break;
case 9 :    start = 12;break;
case 12 :   start = 20;break;
case 20 :   start = 25;break;
case 25 :   recover():break;
default :   /* compiler error */
}
return start;
}
```

Default

رشته خوانده شده با الگوی هیچ کدام از توکن ها مطابقت نمی کند در این صورت خطا گزارش می شود ولی کامپایلر نباید متوقف شود بلکه اسکن بقیه رشته ها را ادامه می دهد.

برنامه به زبان C برای تحلیلگر لغوی در صفحه بعد

```

while(1){
    switch(state){
case 0: c = nextchar();if (isspace(c)){lexeme_beginnin ++ , state = 0;}
        else if (c == '<') state = 1;
        else if (c == '=') state = 5;
        else if (c == '>') state = 6;
        else state = fail();
        break;
case 1 c = nextchar();
        if (c == '=') state = 2;
        else if (c == '>') state = 3;
        else state = 4;
        break;
case 2 : tokentype = relop; tokenvalue = LE; break;
case 3 : tokentype = relop; tokenvalue = NE; break;
case 4 : ungetch(c); tokentype = relop; tokenvalue = L; break;
:
case 9 : c = nextchar();
        if (isletter(c)) state = 10;
        else state = fail();
        break;
case 10: c = nextchar();
        if (isletter(c) || isdigit(c)) state = 10;
        else state = 11; break;
case 11: ungetch(c); install(token value);
        return(gettoken(token value));
        break;
case 12: c = nextchar();
        if (isdigit(c)) state = 13;
        else state = fail();
        break;
case 13: c = nextchar();
        if (isdigit(c)) state = 13;
        else if (c == '.') state = 14;
        else if (c == 'E') state = 16;
        else state = fail();
        break;
case 14: c = nextchar();
        if (isdigit(c)) state = 15
        else state = fail();
        break;
case 15: c = nextchar();
        if (isdigit(c)) state = 15
        else if (c == 'E') state = 16;
        else state = fail();
        break;

```

توضیحات: تابع `install()` به جدول سمبل ها نگاه می کند اگر `Token` خوانده شده کلمه کلیدی باشد صفر را بر می گرداند اگر کلمه کلیدی نباشد و شناسه باشد دو حالت داریم اگر این شناسه در جدول سمبلها موجود باشد اشاره گر به مدخلی که این توکن در آن قرار دارد برگردانده میشود و اگر در جدول توکن ها موجود نباشد در جدول سمبل ها قرار داده میشود و اشاره گر به مدخلی که قرار داده میشود بر گردانده می شود(شکل زیر جدول سمبل).

جدول سمبل ها

If	k
wa	id
a2	id

اگر برای تشخیص کلمات کلیدی از دیاگرام انتقال حالت استفاده شود، تعداد حالات زیاد شده و از این رو `state` گردانی مشکل خواهد بود و سرعت پایین میاید ، پس بهتر است کلمات کلیدی در آغاز کار در جدول سمبلها یا جدول نشانه ها قرار داده شوند

نام تابع شناسه است

تابع `gettoken()` اگر توکن خوانده شده کلمه کلیدی باشد نشانه متناظر آن برگردانده میشود و اگر شناسه باشد `id` را برمیگرداند.

9 , 10 , 11 برای تشخیص کلمات کلیدی است.

بهتر است که توکن هائی که بیشتر تکرار می شوند، در `state` های آغازین قرار گیرند، یعنی در `case` های اول که باعث می شود سرعت بالا رود.

```
cace 16: c = nextchar();
    if (isdigit(c)) state = 18;
    else if (c == '+' || c == '-') state = 17;
    else state = fail();
    break;
cace 17: c = nextchar();
    if (isdigit(c)) state = 18;
    else state = fail();
    break;
cace 18: c = nextchar();
    if (isdigit(c)) state = 18;
    state = 19;
    break;
cace 19: ungetch(c); tokentype = 'num'; tokenvalue = ;
    if (isdigit(c)) state = 18;
    break;
cace 20: c = nextchar();
    if (isdigit(c)) state = 21;
    else state = fail();
    break;
cace 21: c = nextchar();
    if (isdigit(c)) state = 21;
    else if (c == '.') state = 22;
    else state = fail();
    break;
cace 22: c = nextchar();
    if (isdigit(c)) state = 23;
    else state = fail();
    break;
cace 23: c = nextchar();
    if (isdigit(c)) state = 23;
    else state = 24;
    break;
cace 24: c = ungetch(c); tokentype = 'num'; tokenvalue = ;
cace 25: c = nextchar();
    if (isdigit(c)) state = 26;
    else state = fail();
    break;
cace 26: c = nextchar();
    if (isdigit(c)) state = 26;
    else state = 27;
    break;
cace 27: ungetch(c); tokentype = 'num'; tokenvalue = ;
}
}
```

جلسه چهارم تشخیص خطا

Scanner خطاهای زیادی را نمی تواند تشخیص دهد، زیرا دیدگاه محلی نسبت به برنامه دارد به عبارتی برنامه مبدا را به صورت کاراکتر به کاراکتر می خواند به عنوان مثال در عبارت زیری **Scanner** نمی تواند تشخیص دهد که **if** به صورت غلط نوشته شده است (**fi**) و آن را به عنوان یک شناسه معتبر به **parser** تحویل می دهد و کنترل این قبیل خطاها را به **parser** و مولفه های بعد آن محول می سازد. $f_i(x) = f(x)$

→ به عنوان شناسه معتبر می شناسد

□ اگر رشته ای از کاراکترها با الگوی هیچ کدام از توکنها مطابقت نداشته باشد در این صورت خطای لغوی رخ می دهد و **Scanner** می بایست این قابلیت را داشته باشد که از این خطا و خطاهای دیگر گذر کرده و عمل **Scan** (تحلیل لغوی) را تا پایان فایل ادامه دهد.

□ در روش **panic mode** اسکنر از رشته ورودی حذف می کند تا جایی که ادامه کاراکترها با الگوی یکی از **token**ها تطابقت داشته باشد.

مکانیزم های گذر از خطا

panic mode -

- اضافه کردن یک کاراکتر جدید مثل \oplus که دو نقطه را اضافه کردیم
- تعویض یا جایگزینی دو کاراکتر مجاور مثل $\langle \rangle \rightarrow < >$
- تغییر یک کاراکتر

تحلیل نحوی

ساختار هر زبان با قواعد آن مشخص می شوند ساختار زبان های برنامه سازی با گرامرهای مستقل از متن پیاده سازی میشود.

$G(S, NT, T, P)$

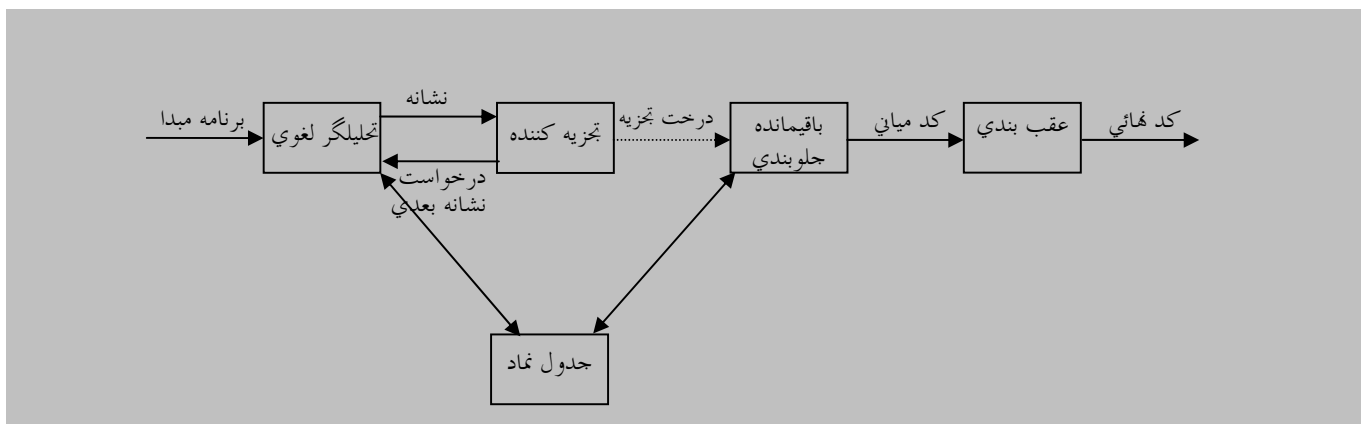
S : start symbol

NT : None Terminal Set

T : Terminal Set

P : Production Rule

شکل زیر موقعیت تجزیه کننده در مدل کامپایلر را نشان می دهد.



- وظیفه اسکنر تشخیص توکن ها می باشد.
- وظیفه اصلی **parser** تشخیص این که رشته توکن های تولید شده توسط **scanner** ، آیا توسط قواعد زبان قابل تولید هست یا نه

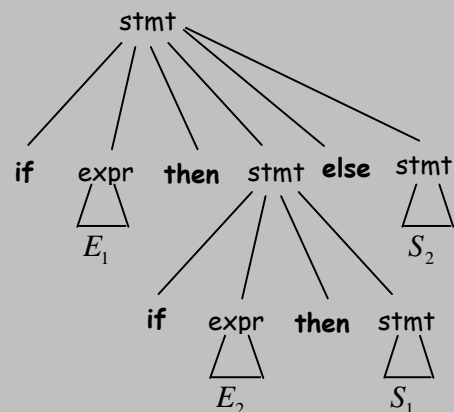
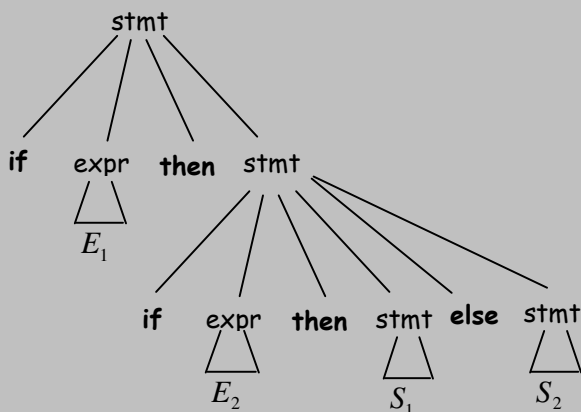
انواع پارسر ها

- پارسرهای غیر پیشگو کننده، باید قابلیت **back tracking** داشته باشد.
- پارسرهای پیشگو کننده ۱- پارسر های بالا به پایین ۲- پارسر های پایین به بالا
- در غیر پیشگو کننده با گرفتن رشته ای از توکن ها با سعی و خطا سعی می کنند که تشخیص دهند آیا این رشته قابل تولید توسط گرامر زبان هست یا نه که این پارسر ها باید خاصیت **back tracking** داشته باشند
- در پارسرهای بالا به پایین از جمله **start symbol** شروع کرده و با اشتقاق سعی در تولید جمله مورد نظر را دارند به عبارتی جمله مورد نظر برگ های درخت اشتقاق می باشد.
- در پارسر های پایین به بالا از جمله ورودی(رشته توکن ها) شروع کرده و سعی در رسیدن به **start symbol** را دارند به عبارتی به جای عبارت سمت چپ، سمت راست را می گذارند.
- در پیشگو کننده گرامر ها بایستی فاقد ابهام باشند.

ابهام:

گرامری دارای ابهام است اگر بتوان برای حداقل یک جمله از زبان تولید شده توسط آن گرامر دو درخت اشتقاق چپ یا دو درخت اشتقاق راست مختلف پیدا کرد به عنوان مثال. گرامر زیر را در نظر بگیرید ملاحظه می شود که برای تولید جمله S_2 **else** S_1 **then** E_2 **if** E_1 **then** چپ یا دو درخت اشتقاق چپ وجود دارد و گرامر مبهم است. و می بایست رفع ابهام شود.

```
stmt → if expr then stmt
      | if expr then stmt else stmt
      | other
```



- اگر زبانی داشته باشیم که نتوانیم گرامر مبهمی برای آن پیدا کنیم زبان ذاتا غیر مبهم است.
رفع ابهام گرامر صفحه قبل به شکل زیر خواهد بود

```

stmt → m_stmt | un_stmt
m_stmt → if expr then m_stmt | other
un_stmt → if expr then m_stmt else un_stmt

```

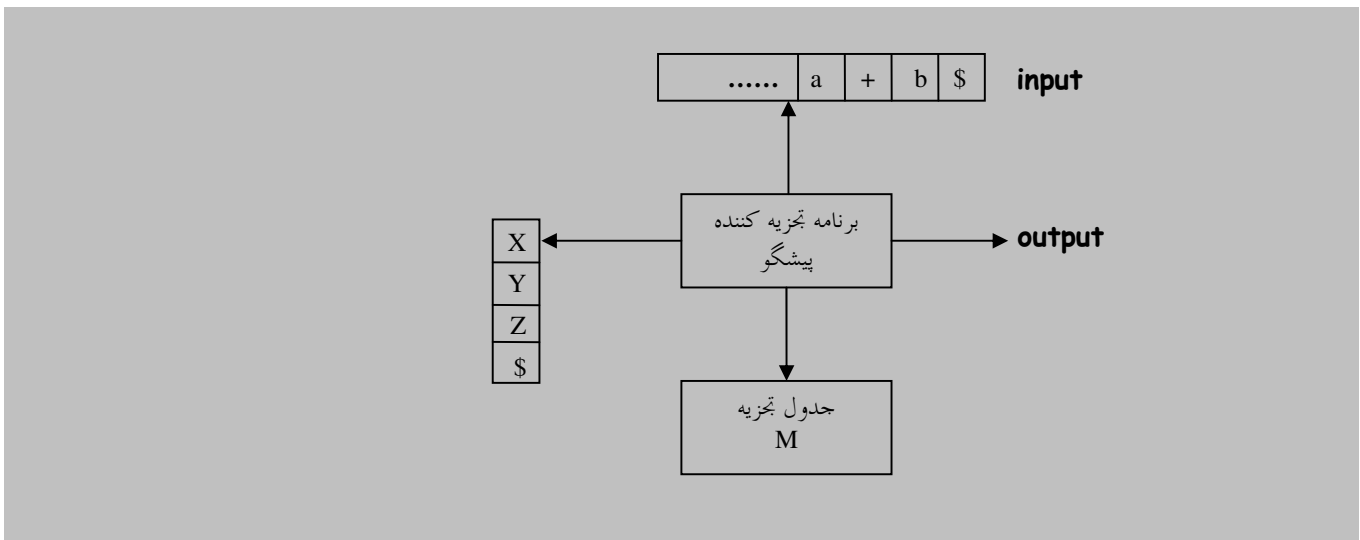
- در اسکنر ورودی کاراکتر است در پارسر ورودی یک توکن است

تجزیه گر پیشگویی پذیر غیر بازگشتی (بالا به پایین)

$LL(1)$ در این روش "L" اول به این معنی است که رشته ورودی از سمت چپ خوانده میشود و "L" دوم به این معنی است که پارسر از سمت چپ ترین اشتقاق استفاده می کند.

ورودی رشته ای از توکن هاست که پایان آن را با \$ نمایش می دهیم

دارای یک پشته است که در ابتدای پشته \$ و **startsymbol** قرار می گیرد، همچنین دارای یک جدول تجزیه می باشد، ساختار کلی این نوع پارسر به فرم زیر است.



مجموعه آغازین $first(\alpha)$

برای سمبل α مجموعه آغازین یا $first$ آن مجموعه ترمینال هائی است که شکل های جمله ای مشتق شده از α با آن شروع می شوند
قواعد ساخت مجموعه آغازین x

- اگر α پایانه باشد $first(\alpha) = \{\alpha\}$

- اگر α بتواند \mathcal{E} را تولید کند آنگاه \mathcal{E} را به مجموعه $first(x)$ اضافه کنید.

- اگر $X \rightarrow Y_1 Y_2 \dots Y_k$ یک غیر پایانه و $first(Y_1), first(Y_2), first(Y_3) \dots first(Y_n)$ مجموعه های $first(Y_i)$ شامل \mathcal{E} باشند یعنی همه Y_i ها بتوانند تهی را تولید کنند در نتیجه X نیز می تواند \mathcal{E} را تولید کند که در این صورت \mathcal{E} به مجموعه $first(X)$ اضافه می شود.

- اگر $X \rightarrow Y_1 Y_2 \dots Y_k$ یک غیر پایانه و $first(Y_1)$ (به جز \mathcal{E}) به مجموعه $first(X)$ اضافه میشود.

- اگر $X \rightarrow Y_1 Y_2 \dots Y_k$ یک غیر پایانه و $first(Y_1)$ (به جز \mathcal{E}) در مجموعه $first(X)$ باشد $first(Y_2)$ (به جز \mathcal{E}) نیز به $first(X)$ اضافه میشود.

- با توجه به قاعده تولید $X \rightarrow Y_1 Y_2 \dots Y_i \dots Y_k$ اگر $first(Y_1), first(Y_2), \dots, first(Y_{i-1})$ (به جز \mathcal{E}) باشد مجموعه $first(Y_i)$ (به جز \mathcal{E}) به مجموعه $first(X)$ اضافه میشود.

مثال: گرامر زیر را در نظر بگیرید.

$$\begin{aligned} \text{first}(T') &= \{ *, \epsilon \} \\ \text{first}(F) &= \{ (, id \} \\ \text{first}(E') &= \{ +, \epsilon \} \\ \text{first}(T) &= \{ (, id \} \\ \text{first}(E) &= \{ (, id \} \end{aligned}$$

پاسخ

$$\begin{aligned} 1) \quad E &\rightarrow TE' \\ 2) \quad E' &\rightarrow +TE' \mid \epsilon \\ 3) \quad T &\rightarrow FT' \\ 4) \quad T' &\rightarrow *FT' \mid \epsilon \\ 5) \quad F &\rightarrow (E) \mid id \end{aligned}$$

مثال. رشته a را توسط گرامر زیر تولید کنید، آیا a عضو این زبان هست یا نه

$$S \rightarrow aB \mid b$$

$$B \rightarrow S \mid \epsilon$$

پاسخ. $S \rightarrow aB \rightarrow a$ و $\text{first}(B) = \{\epsilon, a, b\}$ ، $\text{first}(s) = \{b, a\}$

مثال. در گرامر زیر $\text{first}()$ ها را مشخص کنید.

$$S \rightarrow BAa \mid b$$

$$B \rightarrow aF \mid \epsilon$$

$$A \rightarrow a \mid b \mid aF$$

$$F \rightarrow a \mid c$$

پاسخ

$$\text{first}(B) = \{a, \epsilon\}$$

$$\text{first}(A) = \{a, b\}$$

$$\text{first}(f) = \{a, c\}$$

$$\text{first}(s) = \{a, b\}$$

مثال. در گرامر زیر $\text{first}()$ ها را مشخص کنید.

$$S \rightarrow BAa \mid b \mid Ba$$

$$B \rightarrow aF \mid \epsilon$$

$$A \rightarrow a \mid b \mid aF \mid \epsilon$$

$$F \rightarrow a \mid c$$

پاسخ

$$\text{first}(B) = \{a, \epsilon\}$$

$$\text{first}(A) = \{a, b, \epsilon\}$$

$$\text{first}(f) = \{a, c\}$$

$$\text{first}(s) = \{a, b\}$$

مجموعه Follow(x) (مجموعه دنباله (x))

Follow یک None Terminal (X):

مجموعه ترمینال هائی است که در شکل های جمله ای بلافاصله سمت راست x می آیند (اولین سمت راست آن)

$$\alpha \rightarrow Bx\gamma \rightarrow follow(X) = \{a\}$$

قواعد بدست آوردن مجموعه follow(α)

- اگر α start symbol باشد آنگاه \$ عضو $follow(\alpha)$ می باشد.
- اگر مولدی به صورت $B \rightarrow A\alpha$ داشته باشیم، آنگاه هر چیزی در $first(\alpha)$ به جز ϵ به مجموعه $follow(A)$ اضافه می شود.
- اگر مولدی به صورت $B \rightarrow \alpha A$ داشته باشیم، آنگاه هر چیزی در $follow(B)$ به مجموعه $follow(A)$ اضافه می شود.
- اگر مولدی به صورت $A \rightarrow \alpha B\beta$ وجود داشته باشد، آنگاه هر چیزی در $first(\beta)$ به جز ϵ به مجموعه $follow(B)$ اضافه می شود.
- اگر مولدی به صورت $A \rightarrow \alpha B\beta$ که $first(\beta)$ حاوی ϵ باشد آنگاه هر چیزی در مجموعه $follow(A)$ به $follow(B)$ اضافه می شود.

مثال. FIRST() و Follow() هر کدام را حساب کنید.

$$\begin{aligned} follow(E) &= \{ \$,) \} \\ follow(E') &= \{ \$,) \} \\ follow(T) &= \{ + , , \$,) \} \\ follow(T') &= \{ + , , \$,) \} \\ follow(F) &= \{ * , + , \$,) \} \end{aligned}$$

$$\begin{aligned} first(T') &= \{ * , \epsilon \} \\ first(F) &= \{ (, id \} \\ first(E') &= \{ + , \epsilon \} \\ first(T) &= \{ (, id \} \\ first(E) &= \{ (, id \} \end{aligned}$$

پاسخ

- 1) $E \rightarrow TE'$
- 2) $E' \rightarrow +TE' \mid \epsilon$
- 3) $T \rightarrow FT'$
- 4) $T' \rightarrow *FT' \mid \epsilon$
- 5) $F \rightarrow (E) \mid id$

مثال. FIRST() و Follow() هر کدام را حساب کنید.

$$\begin{aligned} follow(S) &= \{ \$, b , d \} \\ follow(A) &= \{ \$, b , d \} \\ follow(B) &= \{ a , b , d , \$ \} \\ follow(D) &= \{ a , b , d , \$ \} \end{aligned}$$

$$\begin{aligned} first(S) &= \{ a , b , d \} \\ first(A) &= \{ a , b , d \} \\ first(B) &= \{ b , d , \epsilon \} \\ first(D) &= \{ d , \epsilon \} \end{aligned}$$

پاسخ

- 1) $S \rightarrow ABD$
- 2) $A \rightarrow a \mid BSB$
- 3) $B \rightarrow b \mid D$
- 4) $D \rightarrow d \mid \epsilon$

مثال. FIRST() و Follow() هر کدام را حساب کنید.

$$\begin{aligned} follow(S) &= \{ b , \$ \} \\ follow(A) &= \{ a \} \\ follow(B) &= \{ a , b , c , d \} \\ follow(C) &= \{ a , d \} \end{aligned}$$

$$\begin{aligned} first(S) &= \{ \epsilon , b , c , d \} \\ first(A) &= \{ \epsilon , b , c , d \} \\ first(B) &= \{ b , \epsilon \} \\ first(C) &= \{ c , b , \epsilon \} \end{aligned}$$

پاسخ

- 1) $S \rightarrow BCD \mid \epsilon$
- 2) $A \rightarrow AaSb \mid SbC \mid \epsilon$
- 3) $B \rightarrow b \mid \epsilon$
- 4) $C \rightarrow c \mid B$

گرامر زیر که در آن S, E, F سمبل های غیر ترمینال و $+, -, *$ سمبل های ترمینال هستند را در نظر بگیرید مجموعه $First(E)$ و $Follow(E)$ را بیابید.

$$\begin{array}{l} S \rightarrow +E - \\ E \rightarrow F | * \\ F \rightarrow +E* | \lambda \end{array} \quad \begin{array}{l} first(E) = \{+, *\} \\ follow(E) = \{-, *\} \end{array} \quad \text{پاسخ}$$

□ ترمینال ها **follow** ندارند.

گرامر بالا به پایین پیشگوکننده:

□ مبهم نباشد

□ بازگشتی چپ نداشته باشد.

□ فاکتورگیری از چپ بر روی آن اعمال شده باشد.

بازگشتی چپ (left Rotation)

$$A \rightarrow A\alpha \quad \text{۱- بازگشتی چپ آشکار (بدیهی)}$$

$$A \rightarrow B\alpha \rightarrow Ad\alpha$$

این بازگشتی چپ برای A ضمنی است

$$A \rightarrow B\alpha | c$$

$$B \rightarrow Ad | Bc | d$$

۲- بازگشتی چپ ضمنی

حذف بازگشتی چپ بدیهی:

چنانکه قانونی به شکل $A \rightarrow A\alpha_1 | A\alpha_2 | \dots | A\alpha_m | \beta_1 | \beta_2 | \dots | \beta_n$ داشته باشیم به طوریکه β_i ها با A شروع نشوند و هیچکدام از α_i ها نبایستی ϵ باشند می توان مشکل بازگشتی چپ را با جایگزینی دو قانون زیر حل کرد

$$A \rightarrow \beta_1 A' | \beta_2 A' | \dots | \beta_n A'$$

$$A' \rightarrow \alpha_1 A' | \alpha_2 A' | \dots | \alpha_m A' | \epsilon$$

مثال. بازگشتی گرامر زیر را حذف کنید.

$$E \rightarrow EAc | Ea | d | Bc | Ed$$

$$E \rightarrow dE' | BcE'$$

$$E' \rightarrow AcE' | aE' | dE' | \epsilon \quad \text{پاسخ}$$

حذف بازگشتی چپ ضمنی:

□ گرامر G قاعده افسیلون نداشته باشد.

□ دورهای $A \xrightarrow{+} A$ در گرامر G نباشند

الگوریتم

for $i = 1$ to n do

for $j = 1$ to $(i - 1)$ do

- به جای هر قانون به شکل کلی $A_i \rightarrow A_j \beta$ به طوریکه $\alpha_1 | \alpha_2 | \alpha_3 | \dots | \alpha_k$ قرار بده

$$A_i \rightarrow \alpha_1 \beta | \alpha_2 \beta | \dots | \alpha_k \beta$$

- حال اگر A_i دارای بازگشتی چپ آشکار است، آن را حذف کن.

مثال. بازگشتی چپ را از گرامر زیر حذف کنید.

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Sd$$

$$A \rightarrow Ac \mid Aad \mid bd$$

$$S \rightarrow Aa \mid b$$

برای حذف بازگشتی ضمنی چپ داریم

$$S \rightarrow Aa \mid b$$

$$A \rightarrow bdA'$$

$$A' \rightarrow cA' \mid adA' \mid \varepsilon$$

حال با حذف بازگشتی چپ آشکار داریم

فاکتور چپ.

یک گرامر دارای فاکتور چپ است اگر در سمت چپ حداقل دو تا از تولیدات آن یک عنصر مشترک باشد. مثلاً $A \rightarrow a\beta_1 \mid a\beta_2$ دارای فاکتور چپ a می باشد به روشنی می توان دید که از این نوع گرامر ها نمی توان در پارسر های بالا به پایین استفاده کرد چون در هنگام پارس کردن نمی توان تعیین کرد که کدام یک از این قوانین تولید ما را به جواب می رساند. برای رفع مشکل بالا می توان از فاکتور گرفت.

□ در حالت کلی اگر داشته باشیم $A \rightarrow a\beta_1 \mid a\beta_2 \mid \dots \mid a\beta_m \mid \delta_1 \mid \delta_2 \mid \dots \mid \delta_n$ آنگاه با فاکتور گیری از a می توان فاکتور چپ را به صورت زیر حذف نمود

$$A \rightarrow \alpha A' \mid \delta_1 \mid \delta_2 \mid \dots \mid \delta_n$$

$$A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_m$$

مثال گرامر زیر را در نظر بگیرید

$$S \rightarrow iEtS$$

$$S \rightarrow iEtSeS \mid a$$

$$E \rightarrow b$$

$$S \rightarrow iEtSS' \mid a$$

$$S' \rightarrow eS \mid \varepsilon$$

$$E \rightarrow b$$

پس از فاکتور گیری داریم

$$A \rightarrow abcE \mid abcdf \mid abcdT \mid abcdmT'$$

مثال گرامر زیر را در نظر بگیرید

$$A \rightarrow abcB'$$

$$B' \rightarrow E \mid df \mid dT \mid dmT'$$

فاکتور گیری مرحله اول

$$A \rightarrow abcB'$$

$$B' \rightarrow E \mid dZ'$$

$$Z' \rightarrow f \mid T \mid mT'$$

فاکتور گیری مرحله دوم

ساخت جدول تجزیه (parse Table)

گرامر ورودی جدول تجزیه بایستی خواص زیر را داشته باشد.

- گرامر مبهم نباشد
- بازگشتی چپ نداشته باشد
- فاکتور گیری از چپ روی آن اعمال شده باشد.

جدول تجزیه:

- ماتریس دو بعدی است
- سطر ها غیر پایانه ها را نشان می دهد
- ستون ها پایانه ها را نشان می دهد.

برای تشکیل جدول پارسینگ به صورت زیر عمل می کنیم

۱- برای هر قانون به شکل $A \rightarrow \alpha$ مراحل زیر را تکرار کنید

(a) به ازای $a \in first(\alpha)$ شماره قاعده $A \rightarrow \alpha$ را در محل $M[A, a]$ قرار می دهیم.

(b) اگر $\epsilon \in first(\alpha)$ آنگاه به ازای هر $B \in follow(A)$ شماره قاعده $A \rightarrow \alpha$ را در محل $M[A, B]$ قرار می دهیم.

(c) اگر $\epsilon \in first(\alpha)$ و $\$ \in follow(A)$ آنگاه در محل $M[A, \$]$ شماره قاعده $A \rightarrow \epsilon$ را قرار می دهیم.

مثال.

$follow(E) = \{ \$,) \}$	$first(T') = \{ *, \epsilon \}$	1) $E \rightarrow TE'$
$follow(E') = \{ \$,) \}$	$first(F) = \{ (, id \}$	2-3) $E' \rightarrow +TE' \epsilon$
$follow(T) = \{ +, , \$,) \}$	$first(E') = \{ +, \epsilon \}$	4) $T \rightarrow FT'$
$follow(T') = \{ +, , \$,) \}$	$first(T) = \{ (, id \}$	5-6) $T' \rightarrow *FT' \epsilon$
$follow(F) = \{ *, +, \$,) \}$	$first(E) = \{ (, id \}$	7-8) $F \rightarrow (E) id$

	+	*	()	id	\$
E	E	E	1	E	1	E
E'	2	E	E	3	E	3
T	E	E	4	E	4	E
T'	6	5	E	6	E	6
F	E	E	7	E	8	E

مثال. برای گرامر زیر جدول پارسینگ تشکیل دهید.

(1,2,3) $S \rightarrow A; | for(A; C; A)S | B$

(4,5) $A \rightarrow V = E | \epsilon$

(6,7) $C \rightarrow E | \epsilon$

(8) $E \rightarrow V$

(9) $V \rightarrow idX$

(10,11) $X \rightarrow 0V | \epsilon$

(12) $B \rightarrow \{L\}$

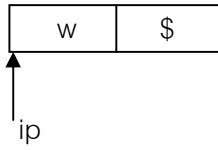
(13,14) $L \rightarrow SL\epsilon | \epsilon$

	;	for	()	=	id	0	{	}	\$
S	1	2	E	E	E	1	E	3	E	E
A	5	E	E	5	E	4	E	E	E	E
C	7	E	E	E	E	6	E	E	E	E
E		E	E	E	E	8	E	E	E	E
V		E	E	E	E	9	E	E	E	E
X	11	E	E	11	11	E	10	E	E	E
B		E	E	E	E	E	12	E	E	E
L	13	13	E	E	E	13	E	13	14	E

	S	A	C	E	V	X	B	L
$first()$	$id ; for \{$	$id \epsilon$	$id \epsilon$	id	id	0ϵ	$\{$	$id ; for \{ \epsilon$
$follow()$	$\$ id ; for \{$	$);$	$;$	$);$	$=);$	$=);$	$\$ id ; for \{$	$\}$

الگوریتم تجزیه پیشگو کننده: LL(1)

- به انتهای پشته ورودی \$ اضافه کنید.
- Ip را طوری مقدار دهی کنید که به ابتدای رشته w اشاره کند.
- به پشته \$ و startsy را اضافه کنید (push)
- رشته ورودی توکن ها هستند
- در پشته، ترمینال، non ترمینال و \$ قرار دارد.



با استفاده از جدول تجزیه و همچنین یک پشته به راحتی می توان عمل تجزیه بالا به پایین را انجام داد. روش کار به این صورت است در ابتدا به انتهای رشته ورودی علامت \$ که نشان دهنده خاتمه فایل است افزوده می گردد از سوی دیگر علامت سرترم گرامر نیز به همراه \$ در داخل پشته قرار می گیرد. عمل تجزیه زمانی با موفقیت به پایان می رسد که در ورودی \$ و در پشته تنها علامت \$ باقی مانده باشد.

پس از قرار دادن سرترم گرامر و \$ در پشته از ورودی یک علامت خوانده می شود، اگر فرض کنیم که علامت موجود در ورودی a و عنصر بالایی پشته X باشد حالات زیر ممکن است اتفاق بیفتد

۱- اگر $a = X = \$$ باشد عمل تجزیه با موفقیت به پایان رسیده است

۲- اگر $a = X \neq \$$ باشد، X از بالای پشته حذف می شود و عنصر بعدی از ورودی خوانده می شود

۳- اگر X یک ترمینال باشد و $a \neq x$ خطای نحوی رخ میدهد

۴- اگر X یک ترم میانی (غیر ترمینال) باشد برنامه به محل $M[X, a]$ در جدول مراجعه می کند در این جا دو حالت امکان دارد الف) قاعده ای به صورت $X \rightarrow UVW$ در خانه $M[X, a]$ وجود دارد در این حالت عنصر X از پشته pop شده و UVW به درون پشته push می شود به طوریکه عنصر U در بالای پشته قرار می گیرد. ب) خانه $M[X, a]$ خالی است در این حالت یک خطای نحوی رخ میدهد.

توضیحات بالا را به صورت الگوریتم زیر نیز می توان نوشت.

```

set ip to point to the first symbol of W$;
repeat
  let X be the top stack symbol and α the symbol pointed to by ip;
  if X is a terminal or $ then
    if X = α then
      pop X from the stack and Advance ip;
    else error();
  else /* X is nonterminal */
    if  $M[X, a] = X \rightarrow Y_1 Y_2 \dots Y_k$  then begin
      pop X from the stack;
      push  $Y_k, Y_{k-1}, \dots, Y_1$  on to the stack, with  $Y_1$  on top;
      output the production  $X \rightarrow Y_1 Y_2 \dots Y_k$ 
    end
  else error
Until X = $ /* stack is empty */
    
```

عمل تجزیه برای جمله $id + id * id$ در زیر تشریح شده است (با توجه به جدول تجزیه صفحه قبل).

پشته	ورودی	قانون استفاده شده
\$E	$id + id * id$	
$E'T$	$id + id * id$	$E \rightarrow TE'$
$E'T'F$	$id + id * id$	$T \rightarrow FT'$
$E'T'id$	$id + id * id$	$F \rightarrow id$
$E'T'$	$+ id * id$	
E'	$+ id * id$	$T' \rightarrow \epsilon$
$E'T +$	$+ id * id$	$E' \rightarrow +TE'$
$E'T$	$id * id$	
$E'T'F$	$id * id$	$T \rightarrow FT'$
$E'T'id$	$id * id$	$F \rightarrow id$
$E'T'$	$* id$	
$E'T'F *$	$* id$	$T' \rightarrow *FT'$
$E'T'F$	id	
$E'T'id$	id	$F \rightarrow id$
$E'T'$	\$	
E'	\$	$T' \rightarrow \epsilon$
\$	\$	$E' \rightarrow \epsilon$

گذر از خطا

خطا

۱- بالای پشته پایانه باشد و با ورودی فعلی تطابق نداشته باشد

۲- بالای پشته غیر پایانه ولی محل $M[X, token]$ خالی (null) باشد.

دو روش برای اصلاح خطای نحوی در روش $LL(1)$ وجود دارد.

- *panic Mode*

- *pharse level*

در روش *panic Mode* اگر پارسر با مراجعه به یک خانه خالی جدول تجزیه یک خطای نحوی بیابد، آنقدر از ورودی حذف می کند تا به یکی از اعضای مجموعه ای موسوم به مجموعه *synchornizing* برسد در روش *panic Mode* با ازای هر غیر پایانه در گرامر یک مجموعه *synchornizing* در نظر گرفته میشود. کارائی روش *panic Mode* نیز بستگی به انتخاب مناسب مجموعه *synchornizing* دارد. این

مجموعه باید به گونه ای انتخاب شود که عمل تجزیه بتواند بدون حذف قسمت زیادی از ورودی به کار خود ادامه دهد. یک انتخاب مناسب در نظر گرفتن مجموعه *follow* هر غیر پایانه ای به عنوان مجموعه *synchronizing* آن غیر پایانه است. با این وجود در نظر گرفتن مجموعه *follow* تنها برای *synchronizing* کافی نیست. برای این که حذف کمتری در برنامه صورت بگیرد می توان نماد های بیشتری را به این مجموعه افزود، مثلاً می توان مجموعه *first* غیر پایانه ها را نیز به مجموعه *synchronizing* آنها افزود.

توسعه الگوریتم تجزیه جهت اعمال گذر از خطا

```

tos = 0;
stack[tos++] = $;
stack[tos++] = start symbol;
token = get token();
do{ // X is top of the stack
if X is a terminal or X = $;
    if X = token {
        pop X;
        token = get token();}
    else
        pop X;
else // X is a non-terminal;
if M[X,token]=NoT NULL && X ≠ synch;
{
    pop X;
    push Yk,Yk-1,...,Y1
}
else if X = synch;
    if X is not sign N.T in stack;
        pop X;
else
    token = get token();
else
    token = get token();
} while(Token ≠ $);

```

1. $S \rightarrow iEtS S'$
2. $S \rightarrow a$
3. $S' \rightarrow eS$
4. $S' \rightarrow \epsilon$
5. $E \rightarrow b$

مثال. آیا گرامر زیر LL(1) هست؟

حل. ابتدا First() و Follow() ها را پیدا می کنیم.

$$\text{First}(S) = \{i, a\}$$

$$\text{First}(S') = \{e, \epsilon\}$$

$$\text{First}(E) = \{b\}$$

$$\text{Follow}(S) = \{e, \$\}$$

$$\text{Follow}(S') = \{e, \$\}$$

$$\text{Follow}(E) = \{t\}$$

	a	b	e	i	t	\$
S	2			1		
S'			3,4			4
E		5				

جدول تجزیه

با توجه به جدول تجزیه گرامر فوق LL(1) نیست.
اگر گرامری در ورودی یا در مدخلی از جدول تجزیه بیش از یک مقدار داشته باشد آن گرامر LL(1) نیست.

بعضی گرامر ها را با

- رفع ابهام
- حذف بازگشتی چپ
- فاکتورگیری از چپ

می توان به LL(1) تبدیل کرد

تشخیص LL(1) بودن بدون استفاده از جدول

گرامری LL(1) است که برای هر قاعده آن که به فرم $A \rightarrow \alpha \mid \beta$ باشد شرایط زیر برقرار باشد.

$$1- \text{First}(\beta) \cap \text{First}(\alpha) = \emptyset$$

2- حداکثر یکی از α و β رشته ϵ را تولید کند.

3- اگر $\alpha \xrightarrow{*} \epsilon$ (به عبارتی $\epsilon \in \text{First}(\alpha)$) باشد $\text{First}(\beta) \cap \text{Follow}(A) = \emptyset$

نکته: اگر $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ آنگاه بایستی $\text{First}(\alpha_1) \cap \text{First}(\alpha_2) \cap \dots \cap \text{First}(\alpha_n) = \emptyset$

مثال. آیا گرامر زیر LL(1) هست.

$$1,2. S \rightarrow A \mid B$$

$$3,4. A \rightarrow Ab \mid f$$

$$5. B \rightarrow cdA \mid ceA$$

$$\text{First}(cdA) \cap \text{First}(ceA) = c \neq \emptyset$$

پارسرهای پایین به بالا ، انتقال - کاهش

از رشته ورودی شروع کرده و با جایگزینی قواعد سعی در تولید Start Symbol را دارند.

S → aABe

A → Abc \ b

B → d

مثال تجزیه پایین به بالا - انتقال کاهش

دنباله ورودی abcde

S → aABe → aAde → aAbcde → abcde اشتقاق راست

abcde انتقال

aAbcd کاهش

aAde انتقال

aABe کاهش

S

پارسرهای پایین بالا به جهت عکس اشتقاق راست عمل می کنند.

عبارت (Phars): بخشی از یک فرم جمله ای است که از یک غیر پایانه بوجود آمده باشد به عنوان نمونه در مثال زیر، β یک عبارت

محسوب می شود. $S \Rightarrow^* \alpha A \gamma \Rightarrow^+ \alpha \beta \gamma$

عبارت ساده (Simple Phars): عبارتی است که در یک قدم بوجود آمده باشد در بسط زیر β یک عبارت ساده است

$S \Rightarrow^* \alpha A \gamma \Rightarrow \alpha \beta \gamma$

دستگیره (Handle): عبارت ساده ای است که در جهت عکس یک اشتقاق راست در نظر گرفته می شود سمت راست Handle هیچ غیر ترمینالی وجود ندارد.

تعریف دیگر: زیر رشته ای منطبق بر سمت راست یک قانون می باشد که ایجاد کننده یک کاهش به غیر پایانه سمت چپ آن قانون می باشد.

تعریف دقیق تر: اگر ترتیب کاهش رشته ورودی به نماد شروع معکوس سمت راست ترین اشتقاق (اشتقاق راست) باشد دنباله ای را که در هر

مرحله کاهش می یابد، **دستگیره** می نامیم.

کاهش	دستگیره
bccdef	c
bBcdef	Bcd
bBef	e
bBCf	bBCf
S	

مثال. دستگیره ها را در کاهش bccdef در با توجه به گرامر زیر مشخص کنید.

S

bBCf

bBef

bBcdef

bccdef

$S \rightarrow bBCf$

حل. ابتدا bccdef را بوسیله سمت راست ترین اشتقاق تولید می کنیم $B \rightarrow Bcd \mid c$

$C \rightarrow e$

پیاده سازی پارسرهای پایین به بالا با استفاده از یک پشته:

در این روش از یک پشته و یک بافر ورودی جهت نگهداری رشته ورودی استفاده می شود. در شروع عمل پارسینگ روی پشته و انتهای رشته ورودی \$ را قرا می دهیم.

اعمال مورد استفاده در این روش:

انتقال (shift): در انتقال، سمبول ها (Token های) ورودی بالای پشته قرار داده می شوند تا زمانی که یک **handel** در بالای پشته ظاهر شود. به عبارتی عمل انتقال تا وقتی که علایم روی پشته با سمت راست یکی از قوانین گرامر منطبق شود ادامه می یابد.

کاهش (Reduce): مجموعه عناصر روی پشته که منطبق با سمت راست یکی از قوانین است را حذف و به جای آن یک غیر ترمینال که در سمت چپ قانون مذکور است را قرار می دهد.

قبول ورودی (Accept): پارسر پایان موفقیت امیز تجزیه را اعلام می کند.

تشخیص خطا (Error): پارسر یک خطای نحوی تشخیص داده و رویه خطا پرداز را فرا می خواند. در زیر مراحل پارسینگ رشته $id+id*id$ بر اساس گرامر زیر دیده می شود.

پشته	ورودی	عمل
\$	$id+id*id\$$	انتقال id
$\$id$	$+id*id\$$	کاهش $E \rightarrow id$
$\$E$	$+id*id\$$	انتقال $+$
$\$E+$	$id*id\$$	انتقال id
$\$E+id$	$*id\$$	کاهش $E \rightarrow id$
$\$E+E$	$*id\$$	انتقال $*$
$\$E+E*$	$id\$$	انتقال id
$\$E+E* id$	$\$$	کاهش $E \rightarrow id$
$\$E+E*E$	$\$$	کاهش $E \rightarrow E * E$
$\$E+E$	$\$$	کاهش $E \rightarrow E + E$
$\$E$	$\$$	Accept

$$(1) E \rightarrow E + E$$

$$(2) E \rightarrow E * E$$

$$(3) E \rightarrow (E)$$

$$(4) E \rightarrow id$$

در روش پارسینگ انتقال-کاهش دو نوع مشکل وجود دارد.

۱- تصمیم گیری در این مورد که کدام مجموعه از عناصر روی پشته با سمت راست قوانین منطبق است (مشکل در انتخاب دستگیره)

۲- انتخاب قانونی که کاهش تحت آن انجام شود. این مشکل وقتی پیش می آید، که سمت راست بیش از یک قانون با مجموعه ای از عناصر روی پشته منطبق باشد. چنین حالتی را تداخل کاهش-کاهش گویند

انواع تداخل در پارسینگ انتقال- کاهش

۱- **تداخل انتقال-کاهش (Shift-Reduce conflict):** زمانی روی می دهد که پارسر نتواند تصمیم بگیرد که عمل انتقال را انجام بدهد یا عمل کاهش را.

۲- **تداخل کاهش-کاهش (Reduce-Reduce conflict):** اگر بیش از یک قاعده برای کاهش در یک لحظه قابل استفاده باشد این نوع مشکل رخ داده است.

روش پارسینگ LR

روش پارسینگ LR جزء دسته الگوریتم های پارسینگ پایین به بالا می باشد، الگوریتم های LR خود به سه دسته کوچکتر تقسیم می شوند. که تفاوتشان در ساخت جدول تجزیه شان می باشد. در حالی که روش پارسینگ برای آنها مشابه است

- (Simple LR)SLR

- (Look Ahead LR)LALR

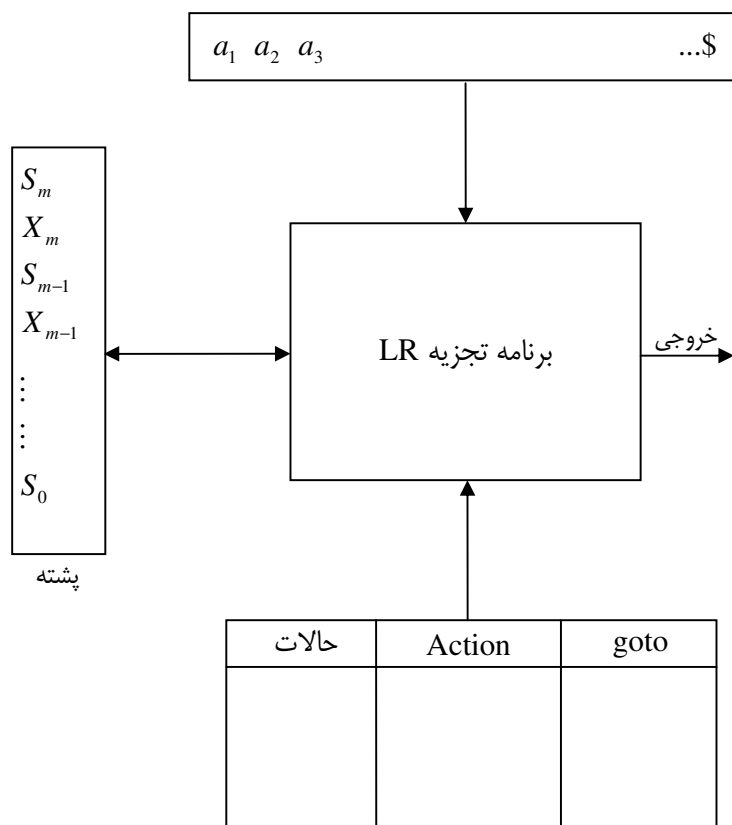
- (Canonical)CLR

- برای اکثریت قریب به اتفاق زبان های برنامه نویسی می توان از این روش پارسینگ استفاده نمود پارسر های LR دارای 4 جزء اصلی هستند
- **بافر ورودی**، رشته ورودی در این قسمت قرار می گیرد و به انتهای رشته ورودی \$ اضافه می شود LR با دیدن \$ پایان رشته ورودی را تشخیص می دهد
- **پشته**: محتوای پشته به صورت $S_0 X_1 S_1 \dots X_m S_m$ نشان داده می شود. S_i نشان دهنده حالات و X_i پایانه ها و غیر پایانه های گرامر است، نمادهای ورودی به پشته منتقل می شوند تا یک دستگیره در پشته یافت شود، پس از کشف دستگیره کاهش انجام میگیرد
- **برنامه پارسینگ**: این برنامه قسمت اصلی تجزیه کننده است برنامه تجزیه کننده بر اساس نماد جاری رشته ورودی و حالات بالای پشته و محتوای جدول تجزیه مرحله بعدی را تعیین می کند
- **جدول پارسینگ**: این بخش شامل دو قسمت **action** و **goto** است. قسمت **action** عملی را که باید انجام شود و بخش **goto** حالت بعدی را مشخص می کند

مزیت: سرعت اعلان خطا های نحوی از این پارسرها بالاست

عیب: جهت ساخت جدول تجزیه این پارسرها حجم محاسبات دستی خیلی زیاد است.

در شکل زیر اجزاء مختلف پارسر LR و رابطه آنها با هم نشان داده شده است



یک قلم (1-item) برای LR(0) (قلم به طور خلاصه) از گرامر G مولدی از G با یک نقطه (یا علامت خاص دیگر) در مکانی در سمت راست آن است بنابراین مولد $A \rightarrow XYZ$ چهار قلم را ایجاد می کند.

$A \rightarrow .XYZ$

$A \rightarrow X.YZ$

$A \rightarrow XY.Z$

$A \rightarrow XYZ.$

مولد $A \rightarrow \epsilon$ تنها یک قلم به صورت $A \rightarrow .$ را تولید می نماید، یک قلم می تواند با یک زوج عدد صحیح نشان داده شود، اولین عدد شماره مولد و دومین عدد موقعیت نقطه را مشخص می نماید، در نتیجه قلم مشخص می نماید که چه مقداری از مولد، در یک نقطه از فرایند تجزیه دیده شده است برای مثال اولین قلم فوق نشان میدهد که انتظار دیدن رشته مشتق پذیر از XYZ را در ورودی داریم. دومین قلم نشان می دهد که اخیراً در ورودی رشته مشتق پذیر از X دیده شده است و این که در ادامه انتظار دیدن یک رشته مشتق پذیر از YZ را داریم.

عمل Closure :

اگر I مجموعه ای از اقلام برای گرامر G باشد، آنگاه Closure(I) مجموعه اقلام ایجاد شده از I با استفاده از دو قانون زیر است
 ۱- در ابتدا هر قلم موجود در I به Closure(I) اضافه می شود.

- اگر $B \rightarrow \alpha$ در Closure(I) موجود باشد و $B \rightarrow \gamma$ یک مولد باشد آنگاه قلم $B \rightarrow \gamma$ اگر قبلاً در I نباشد به آن اضافه می گردد. این قانون تا زمانی به کار گرفته می شود که هیچ قلم جدید دیگری نتواند به Closure(I) اضافه شود.

مثال. باتوجه به گرامر زیر اگر I مجموعه ای یک قلم به صورت $\{[E' \rightarrow .E]\}$ باشد آنگاه Closure(I) شامل چه اقلام هائی خواهد شد؟

$E' \rightarrow .E$
$E \rightarrow .E + T$
$E \rightarrow .T$
$T \rightarrow .T * F$
$T \rightarrow .F$
$F \rightarrow .(E)$
$F \rightarrow .id$

پاسخ.

$E' \rightarrow E$
 $E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid id$

ساخت جدول تجزیه پارسر SLR

ابتدا می بایست دیاگرام حالت این پارسر رسم شود جهت رسم دیاگرام حالت یک پارسر به صورت زیر عمل می کنیم

- اگر S علامت شروع گرامر باشد ابتدا قاعده ای به فرم $S' \rightarrow S$ ($S' \rightarrow S\$$) که در آن S' یک غیر پایانه جدید است به گرامر اضافه می کنیم، گرامر حاصل را گرامر افزوده (augmented grammar) گویند

- حالت جدیدی به نام S_0 ایجاد می کنیم و ایتیم $S_0 \rightarrow .S_0$ در S_0 قرار داده و سپس Closure این ایتیم را به S_0 اضافه می کنیم

- اگر به طور کلی در حالت S_i ایتیم هائی به فرم $\left\{ \begin{array}{l} \{A_1 \rightarrow \alpha_1 \cdot x \beta_1 \\ \vdots \\ \{A_n \rightarrow \alpha_n \cdot x \beta_n \end{array} \right.$ داشته باشیم (که در آن X می تواند پایانه و یا غیر پایانه باشد) در

در این صورت حالت جدیدی به نام S_j ایجاد کرده، S_i را توسط لبه ای با برچسب X به S_j منتقل می کنیم و ایتیم فوق را با این تغییر که در همه علامت . به بعد از علامت X منتقل شده است در وضعیت جدید قرار می دهیم و سپس بستر این ایتیم ها را محاسبه و در S_j قرار می دهیم. چنانکه در دیاگرام حالتی مانند S_k وجود داشته باشد که دقیقا مطابق S_j باشد در این صورت S_j ایجاد نشده و در عوض S_i توسط لبه ای با برچسب X به S_k متصل می گردد، این قدم را آنقدر تکرار می کنیم تا دیگر حالت جدیدی به دیاگرام اضافه نشود.

مثال. دیاگرام حالت گرامر زیر را رسم کنید.

$$E \rightarrow E + T$$

$$E \rightarrow T$$

$$T \rightarrow id$$

$$T \rightarrow (E)$$

در قدم اول یک Start Symbol مانند $S \rightarrow ES$ را به گرامر اضافه می کنیم و حالت جدید S_0 را ایجاد کرده و مراحل گفته شده را پی

می گیریم

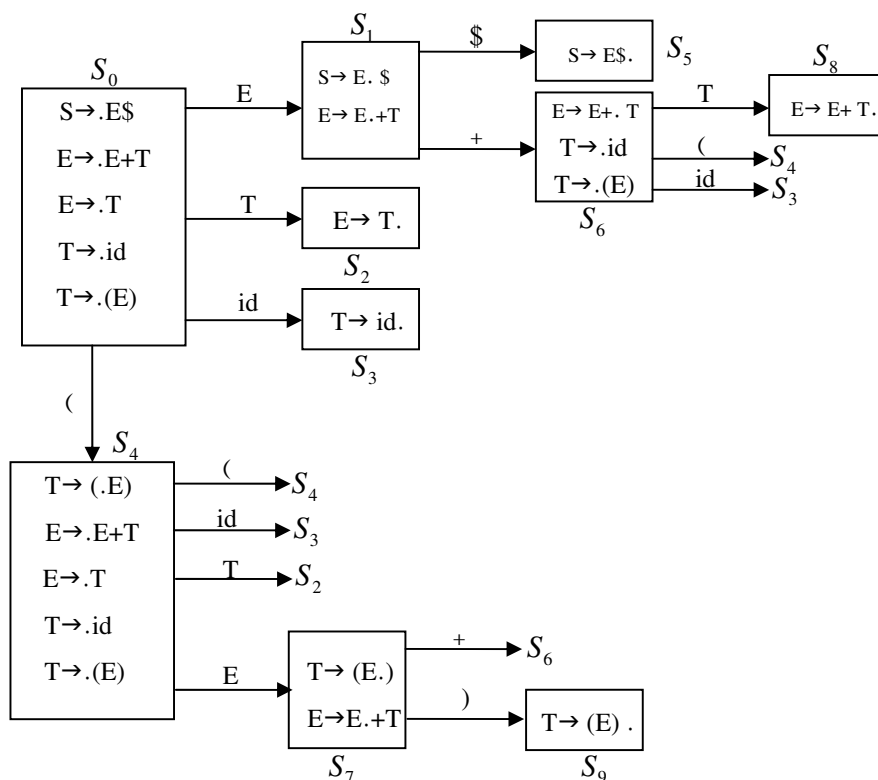
$$S \rightarrow E$$

$$E \rightarrow E + T$$

$$E \rightarrow T$$

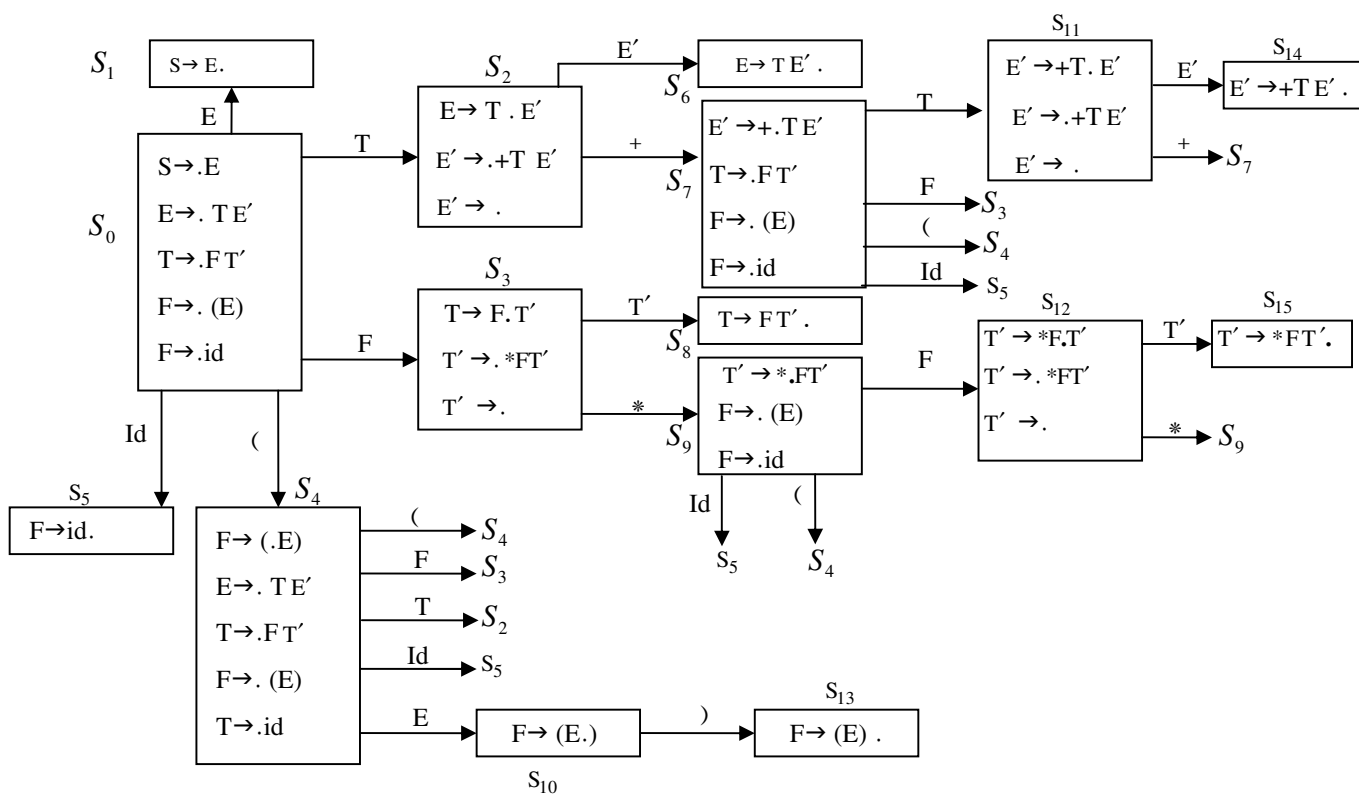
$$T \rightarrow id$$

$$T \rightarrow (E)$$



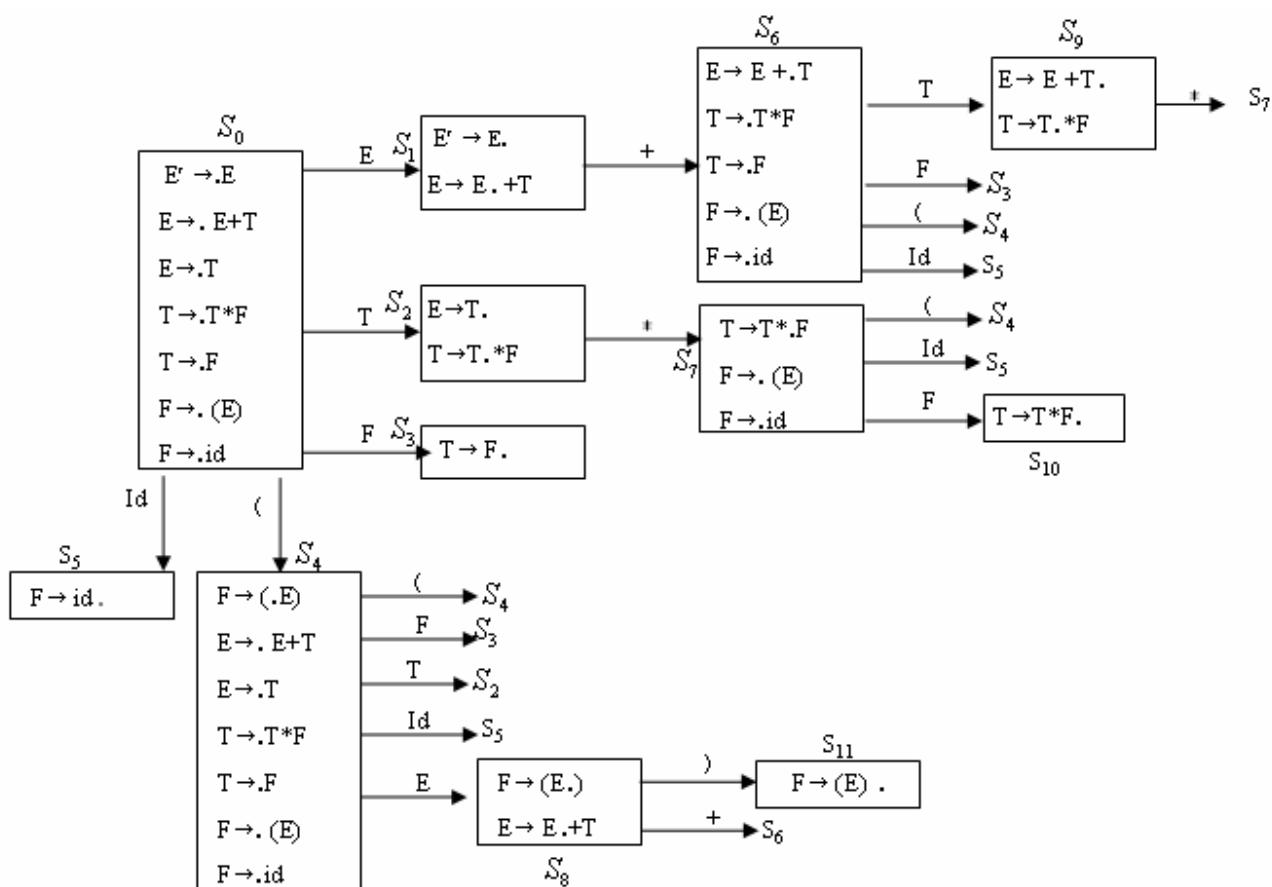
مثال. دیاگرام حالت گرامر زیر را رسم کنید.

- 1) $E \rightarrow TE'$
- 2) $E' \rightarrow +TE'$
- 3) $E' \rightarrow \epsilon$
- 4) $T \rightarrow FT'$
- 5) $T' \rightarrow *FT'$
- 6) $T' \rightarrow \epsilon$
- 7) $F \rightarrow (E)$
- 6) $F \rightarrow id$



مثال. دیاگرام حالت گرامر زیر را رسم کنید.

- | |
|---|
| $E' \rightarrow E$
1) $E \rightarrow E + T$
2) $E \rightarrow T$
3) $T \rightarrow T * F$
4) $T \rightarrow F$
5) $F \rightarrow (E)$
6) $F \rightarrow id$ |
|---|



در جدول تجزیه SLR :

- سطرها State ها را نشان می دهند.
- ستون های قسمت اول (action)، ترمینال ها $+$ را نشان می دهند
- ستون های قسمت دوم (goto) غیرترمینال ها را نشان می دهند.

برای تکمیل قسمت **action** و **goto** به صورت زیر عمل می کنیم

- اگر a یک ترمینال باشد و در حالت S_i با دیدن a به خانه S_j برویم در $action[i, a]$ قرار می دهیم S_j
- در حالت S_i به ازای تمامی دستورات $A \rightarrow \alpha$ ، به ازای تمام اعضای مجموعه $Follow(A)$ (مثلا b) در $action[i, b]$ قرار می دهیم R_n به طوریکه n شماره دستور $A \rightarrow \alpha$ است.
- اگر در حالت S_i دستور $S' \rightarrow S$ قرار داشته باشد در $action[i, \$]$ فرمان $accept$ قرار می گیرد.
- برای تکمیل بخش $goto$ کافی است چنانکه از حالت S_i با یک غیر ترمینال (مثلا A) به حالت S_j می رویم در $goto[i, A]$ قرار می دهیم j

مثال. جدول تجزیه دیاگرام صفحه قبل را رسم کنید.

$Follow(E) = \{+,), \$\}$

$Follow(T) = \{*, +,), \$\}$

$Follow(F) = \{*, +,), \$\}$

	Action					goto			
	id	+	*	()	\$	E	T	F
S_0	S_5			S_4			1	2	3
S_1		S_6				Accept			
S_2		R_2	S_7		R_2	R_2			
S_3		R_4	R_4		R_4	R_4			
S_4	S_5			S_4			8	2	3
S_5		R_6	R_6		R_6	R_6			
S_6	S_5			S_4				9	3
S_7	S_5			S_4					10
S_8		S_6			S_{11}				
S_9		R_1	S_7		R_1	R_1			
S_{10}		R_3	R_3		R_3	R_3			
S_{11}		R_5	R_5		R_5	R_5			

جهت تجزیه یک رشته با استفاده از جدول SLR مراحل زیر را طی می کنیم.

ابتدا علامت $\$$ را به انتهای رشته ورودی اضافه می کنیم و حالت اولیه به پشته اضافه می گردد.

۱- **انتقال:** اگر $action[i, a] = S_n$ ، برنامه تجزیه کننده ابتدا نماد ورودی a و سپس n را به بالای پشته منتقل میکند.

۲- **کاهش:** اگر $action[i, a] = R_n$ باشد یک دستگیره یافت شده است که باشد کاهش یابد اگر n شماره قاعده تولید $A \rightarrow \beta$ باشد، برنامه تجزیه کننده دنباله β را از بالای پشته حذف می کند و پس از حذف β اگر حالت بالای پشته m باشد برنامه تجزیه کننده ابتدا A و سپس، عدد موجود در $goto[m, A]$ را به بالای پشته اضافه می کند، نکته این است که برای حذف β از بالای پشته به تعداد $2|\beta|$ از بالای پشته حذف می گردد زیرا بین هر نماد در پشته یک شماره حالت نیز جود دارد.

۳- **پذیرش:** اگر $action[i, a] = accept$ باشد تجزیه رشته ورودی با موفقیت انجام شده است

۴- **خطا:** اگر $action[i, a] = error$ باشد تجزیه رشته ورودی با عدم موفقیت روبرو شده است

مثال. با استفاده از جدول تجزیه زیر رشته $id*id+id$ را تجزیه کنید

- $E' \rightarrow E$
 1) $E \rightarrow E + T$
 2) $E \rightarrow T$
 3) $T \rightarrow T * F$
 4) $T \rightarrow F$
 5) $F \rightarrow (E)$
 6) $F \rightarrow id$

مراحل تجزیه رشته $id*id+id$

جدول تجزیه SLR

پشته	رشته ورودی	خروجی
0	$id*id+id\$$	Shift5
0id5	$*id+id\$$	$F \rightarrow id$
0F3	$*id+id\$$	$T \rightarrow F$
0T2	$*id+id\$$	Shift7
0T2*7	$id+id\$$	Shift5
0T2*7id5	$+id\$$	$F \rightarrow id$
0T2*7F10	$+id\$$	$T \rightarrow T * F$
0T2	$+id\$$	$E \rightarrow T$
0E1	$+id\$$	Shift6
0E1+6	$id\$$	Shift5
0E1+6id5	$\$$	$F \rightarrow id$
0E1+6F3	$\$$	$T \rightarrow F$
0E1+6T9	$\$$	$E \rightarrow E + T$
0E1	$\$$	Accept

	Action						goto		
	id	+	*	()	\$	E	T	F
S_0	S_5			S_4			1	2	3
S_1		S_6				Accept			
S_2		R_2	S_7		R_2	R_2			
S_3		R_4	R_4		R_4	R_4			
S_4	S_5			S_4			8	2	3
S_5		R_6	R_6		R_6	R_6			
S_6	S_5			S_4				9	3
S_7	S_5			S_4					10
S_8		S_6			S_{11}				
S_9		R_1	S_7		R_1	R_1			
S_{10}		R_3	R_3		R_3	R_3			
S_{11}		R_5	R_5		R_5	R_5			

```

while(!EOF){
IF(action[ $S_i, x$ ] =  $S_j$ )
    Push(x);
    Push( $S_j$ );
    x = next Token();
else IF(action[ $S_i, x$ ] =  $R : A \rightarrow \alpha$ )
    For(i = 0 to  $2|\alpha|$ )
        pop();
        Push(A);
        Push(goto( $S_i, A$ ))
else IF[action(s, x) = Accept] Then
    return
else
    error();
} // end
    
```

با توجه به توضیحات صفحه قبل می توان الگوریتم تجزیه یک رشته را به این شکل نوشت.
 نکته: اگر گرامری SLR باشد آنگاه LALR و CLR نیز هست.

مزیت استفاده از گرامر های مبهم در تجزیه LR :

با استفاده از ابهام می توان جداول کوچکتری برای تجزیه پایین به بالا ایجاد کرد(مزیت). استفاده از گرامر های مبهم در بسیاری از موارد موجب ایجاد تداخل هائی در جدول تجزیه خواهد شد برای این که عمل تجزیه پایین به بالا امکان پذیر باشد می بایست این تداخل ها را از بین برد . بنابراین با مبهم کردن گرامر گرچه تعداد قوانین و در نتیجه حالات جدول کاهش می یابد اما مشکل تداخل ایجاد می شود.

$$E \rightarrow E + T \mid T$$

$$T \rightarrow T * F \mid F$$

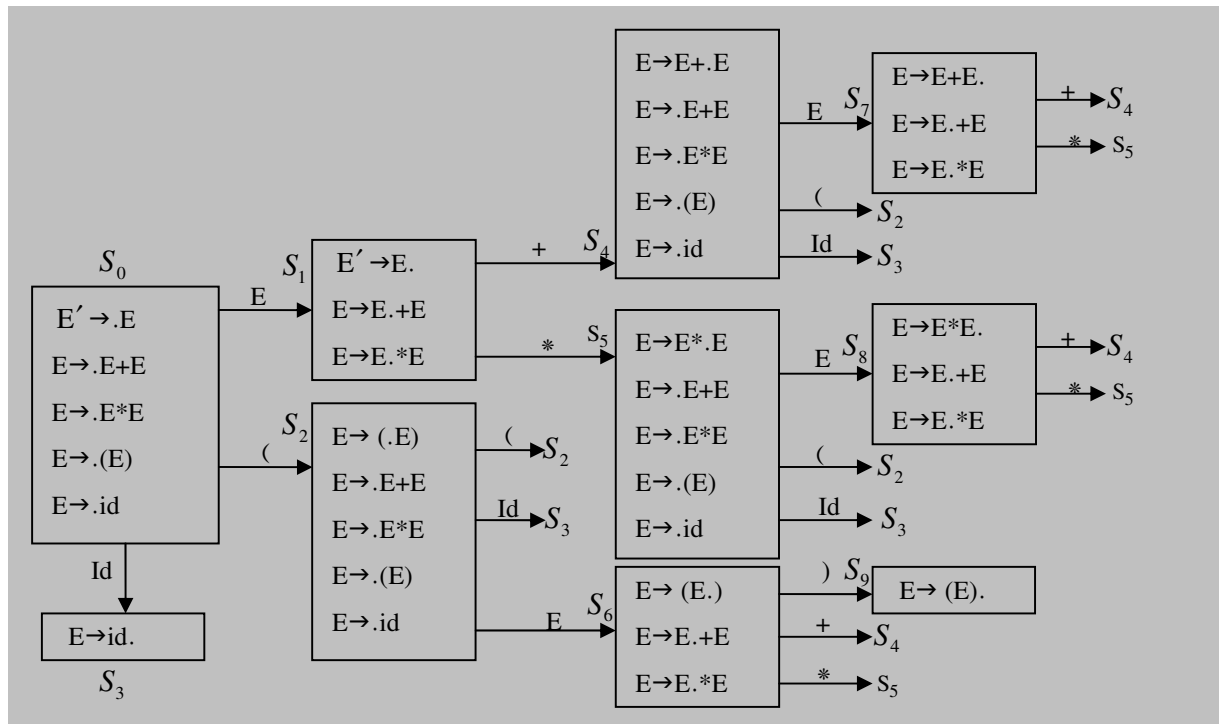
$$F \rightarrow id \mid (E)$$

برای روشن شدن موضوع گرامر را در نظر بگیرید، این گرامر (اگر دیاگرام آن رسم) به یازده حالت (State) نیاز دارد حال می خواهیم

دیاگرام مبهم شده این گرامر را که در زیر آمده است رسم کنیم و همچنین جدول تجزیه را ایجاد کرده و به بررسی تداخل ها و تعداد حالت(که کاهش یافته) بپردازیم.

- $E \rightarrow E + E$
- $E \rightarrow E * E$
- $E \rightarrow (E)$
- $E \rightarrow id$

رسم دیاگرام:



حال با توجه به دیاگرام رسم شده در صفحه بعد جدول تجزیه را تشکیل می دهیم.

جدول تجزیه دیاگرام صفحه قبل.

$$\text{Follow}(E) = \{\$, +, *, \}$$

state	action					goto	
	id	+	*	()	\$	E
0	S ₃			S ₂			1
1	S ₄	S ₅				acc	
2	S ₃			S ₂			6
3		R ₄	R ₄		R ₄	R ₄	
4	S ₃			S ₂			7
5	S ₃			S ₂			8
6		S ₄	S ₅		S ₉		
7		S ₄ , R ₁	S ₅ , R ₁		R ₁	R ₁	
8		S ₄ , R ₂	S ₅ , R ₂		R ₂	R ₂	
9		R ₃	R ₃		R ₃	R ₃	

ملاحظه می شود که در چهار خانه تداخل shift/Reduce رخ داده است. جهت از دست ندادن سرعت از گرامر مبهم استفاده می شود و جهت حل تداخل از الویت عملگر ها استفاده می شود.

تجزیه CLR :

آیتم LR(1) : یک آیتم LR(1) یک زوج مرتب متشکل از یک آیتم LR(0) و یک مجموعه پایانه به نام مجموعه پیش بینی (Lookahead) است، و معمولا به صورت $[A \rightarrow \alpha \cdot, LA]$ نمایش داده می شود رابطه زیر در مورد مجموعه پیش بینی و مجموعه Follow غیر پایانه سمت راست این آیتم ها وجود دارد. $LA \subseteq \text{Follow}(A)$

توضیح این که جهت پر کردن جدول تجزیه CLR ابتدا می بایست دیاگرام حالت آن را رسم کنیم دیاگرام حالت CLR مانند SLR است و آیتم شروع $\langle S' \rightarrow \cdot S, \{\$ \} \rangle$ می باشد

الگوریتم محاسبه تابع بستار که در مورد آیتم های LR(0) توضیح داده شده بسیار شبیه الگوریتم یافتن بستار آیتم های LR(1) است با این تفاوت که در این جا باید به صورت زیر برای آیتم های غیر هسته ای جدیدی که به مجموعه بستار اضافه می شوند، مجموعه پیش بینی تعیین نمود.

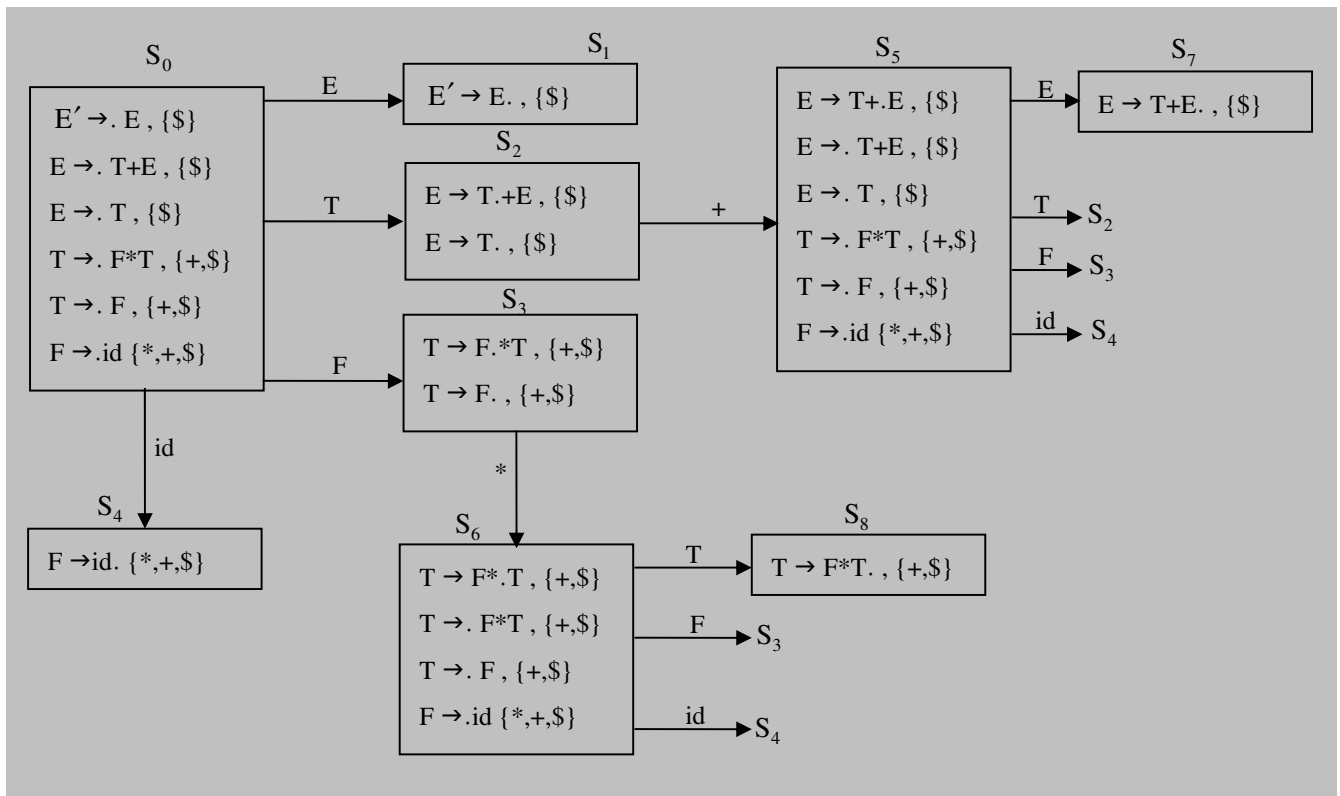
$$[A \rightarrow \alpha \cdot B \beta, \{a\}]$$

$$[B \rightarrow \cdot \delta, \{b\}] \mid b \in \text{First}(\beta a)$$

مثال. دیاگرام حالت CLR گرامر زیر را رسم کنید.

$E' \rightarrow E$
1) $E \rightarrow T + E$
2) $E \rightarrow T$
3) $T \rightarrow F * T$
4) $T \rightarrow F$
5) $F \rightarrow \text{id}$

صفحه بعد دیاگرام رسم شده است.



نحوه تکمیل جدول CLR(1) :

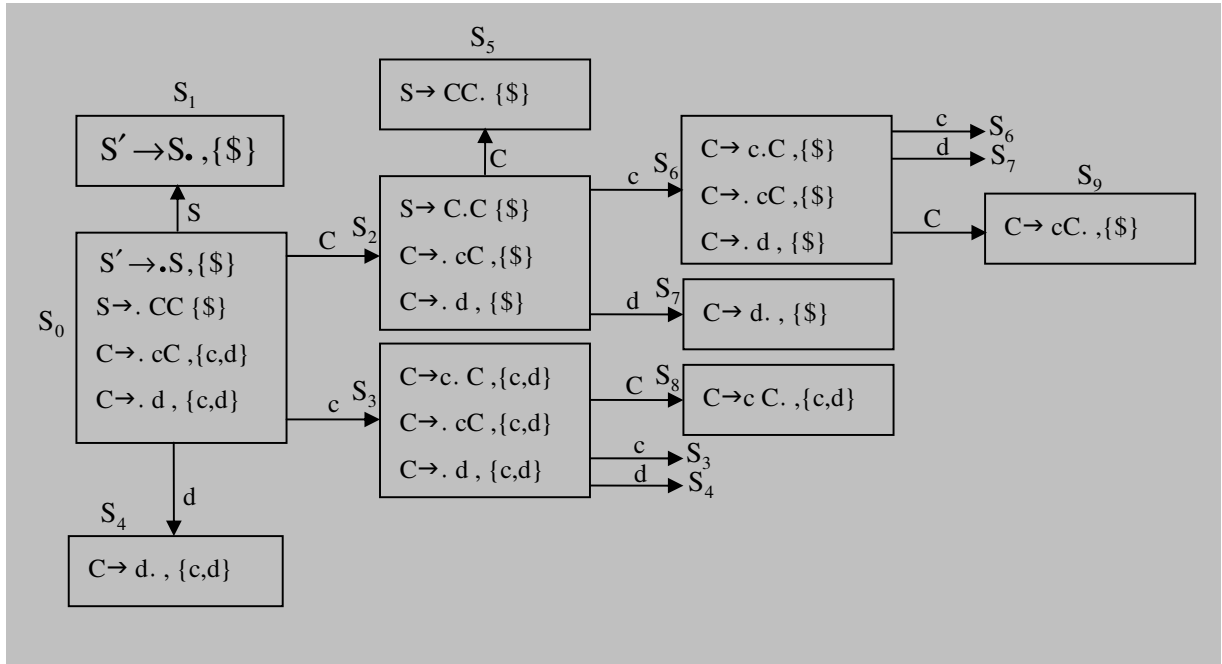
روش تهیه جدول CLR(1) بسیار شبیه روشی است که در رابطه با SLR(1) بیان گردید، با این تفاوت که اگر در این جا در وضعیت I آیتمی به فرم $[A \rightarrow \alpha., \{b\}]$ داشته باشیم، در خانه های $action[i, b]$ دستور Reduce $A \rightarrow \alpha$ قرار می دهیم، به عبارت دیگر در این جا به جای استفاده از مجموعه Follow ها از مجموعه پیش بینی شده استفاده می کنیم.

مثال. جدول تجزیه را برای دیاگرام مثال قبلی تکمیل کنید.

state	action				goto		
	id	+	*	\$	E	T	F
S ₀	S ₄				1	2	3
S ₁				ac			
S ₂		S ₅		R ₂			
S ₃		R ₄	S ₆	R ₄			
S ₄		R ₅	R ₅	R ₅			
S ₅	S ₄				7	2	3
S ₆	S ₄					8	3
S ₇				R ₁			
S ₈		R ₃		R ₃			

مثال. دیاگرام حالت CLR گرامر زیر را رسم کنید.

- $S' \rightarrow S$
 1) $S \rightarrow CC$
 2) $C \rightarrow cC$
 3) $C \rightarrow d$



جدول تجزیه دیاگرام رسم شده.

state	action			goto	
	c	d	\$	S	C
S ₀	S ₃	S ₄		1	2
S ₁				ac	
S ₂	S ₆	S ₇			5
S ₃	S ₃	S ₄			8
S ₄	R ₃	R ₃			
S ₅			R ₁		
S ₆	S ₆	S ₇			9
S ₇			R ₃		
S ₈	R ₂	R ₂			
S ₉			R ₂		

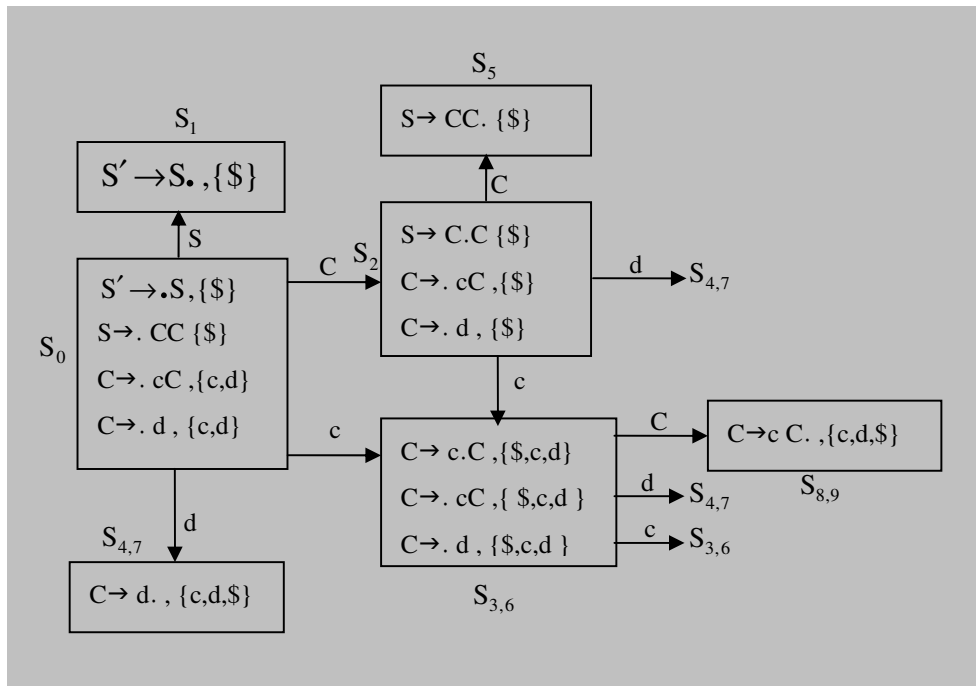
رسم دیاگرام و جدول تجزیه LALR :

برای رسم دیاگرام LALR یک گرامر باید، ابتدا دیاگرام CLR گرامر را رسم کرد. سپس از روی دیاگرام CLR با انجام دو عمل زیر دیاگرام LALR بدست می آید.

ابتدا در دیاگرام CLR به دنبال وضعیت هائی می گردیم که بخش آیتیم های $LR(0)$ آنها (که به آن اصطلاحاً هسته یا Core یک وضعیت می گویند) یکسان باشد. سپس در دیاگرام وضعیت هائی را که هسته یکسانی دارند ادغام می کنیم و مجموعه پیش بینی های این وضعیت ها را نظیر به نظیر با هم ادغام می کنیم.

مثال. با توجه به دیاگرام CLR صفحه 39، دیاگرام LALR مربوطه را رسم کنید.

ملاحظه می شود که در دیاگرام مربوطه S_3 با S_6 ، S_4 با S_7 و S_8 با S_9 قابل ادغام هستند (هسته یکسان دارند)

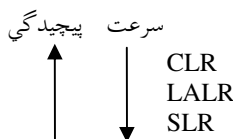


اگر گرامری CLR باشد و پس از ادغام وضعیت های مزبور در جدول تجزیه LALR گرامر نیز تداخلی بوجود نیاید می توان نتیجه گرفت که گرامر LALR است. حال اگر گرامری CLR بوده و LALR نباشد، پس از ادغام در جدول تجزیه LALR گرامر تداخل نوع کاهش/کاهش بروز خواهد کرد، هیچگاه در اثر ادغام در جدول LALR تداخل نوع انتقال-کاهش بوجود نخواهد آمد مگر آنکه گرامر CLR هم نبوده باشد که در این صورت قطعاً LALR هم نیست، به عبارتی در چنین حالتی، تداخل قطعاً قبل از ادغام نیز وجود داشته است. طریقه تهیه جدول تجزیه LALR از روی دیاگرام ادغام شده عیناً همانند روش تهیه جدول CLR است. به عنوان یک نمونه جدول تجزیه LALR مثال قبل به صورت زیر خواهد بود.

state	action			goto	
	c	d	\$	S	C
0	$S_{3,6}$	$S_{4,7}$		1	2
1			acc		
2	$S_{3,6}$	$S_{4,7}$			5
36	$S_{3,6}$	$S_{4,7}$			89
47	R_3	R_3	R_3		
5			R_1		
89	R_2	R_2	R_2		

تعداد وضعیت های جداول تجزیه SLR و LALR هر گرامری دقیقا یکسان است لیکن تعداد وضعیت های جدول تجزیه CLR گرامرها به مراتب بزرگتر از جداول SLR و LALR است. تنها عیب جزئی روش LALR نسبت به CLR این است که خطاهای نحوی را ممکن است قدری دیرتر کشف کند. البته هیچ یک از این سه روش پس از رسیدن به یک توکن غلط آن را به داخل انباره انتقال نخواهد داد. در CLR پس از رسیدن به یک توکن خطا، دیگر حتی عمل کاهش نیز انجام نخواهد شد، لیکن در مورد SLR و LALR ممکن است خطا پس از اجرای چند عمل کاهش اضافی کشف شود. با این حساب CLR سریعتر است.

اگر گرامر CLR باشد آیا می توان نتیجه گرفت که گرامر LALR است؟
 خیر زیرا در جدول تجزیه LALR ممکن است تداخل r-r داشته باشیم.
 اگر گرامری SLR باشد آن گرامر CLR و هم LALR است
 مقایسه سرعت و پیچیدگی



تحلیل معنایی Type checking

به صورت پویا : در زمان اجرا عمل تحلیل معنایی انجام میشود. کامپایلر باید علامت خاصی را در لابه لای کد اجرا قرار دهد به عنوان مثال فرض کنید آرایه ۱۰ تایی تعریف کرده اید در برنامه طبق زیر داریم

```
Int A[10];      A[i]=5
```

در اینجا مقدار i در زمان اجرا مشخص میشود در این جا کامپایلر علامتی را قرار میدهد تا مشخص کند که مقدار i در حین اجرا تعیین میشود

به صورت ایستا:

۱- یکسان بودن گونه (Type checking): اجازه نداریم آرایه را با متغیر جمع کنیم. عملگرها باید عملوند یکسان داشته باشند.

۲- یکسان بودن نام (واحد بودن نام): نمیتوان یک نام را دو جا تعریف کرد

۳- کنترل جریان: مثلا در swich در case باید Break داشته باشیم وگرنه خطا رخ میدهد

۴- بررسی ساختارهای تودرتو: مثلا در زبان Ada هر حلقه یک نام مخصوص برای خود دارد

```
For1(.....)
  For2(.....)
  End for2
End for1
```

مثال:

اگر a یک متغیر از نوع آرایه یک بعدی با اندیس های ۱ تا ۱۰ و I یک عدد صحیح باشد خطای خارج از محدوده (sub scrip out of vong) که در دستورات I=11 ; B:=A[i] وجود دارد در حالت کلی در چه زمانی و توسط چه برنامه ای قابل کشف است؟

۱- زمان اجرا توسط سیستم عامل

۲- زمان کامپایل توسط تجزیه کننده دستوری

۳- زمان کامپایل توسط تجزیه کننده معنایی

۴- زمان اجرا توسط خود برنامه حاوی دستورات فوق

جواب: گزینه 4 درست است، اگر B := A[11] بود آنگاه خطا در مرحله کامپایل توسط تحلیل گر مفهومی کشف می شد

نکته: پیش پردازنده ها در زمان کامپایل مشخص میشوند.
سوال: کدام یک از گزینه ها در خصوص گرامر نوع مقابل صحیح است؟

- 1) $S \rightarrow aACb$
 2-3) $A \rightarrow b \mid \lambda$
 3-4) $C \rightarrow cC \mid \lambda$

- ۱) گرامر LL(1) نیست، SLR(1) هم نیست
 ۲) گرامر LL(1) نیست ولی SLR(1) هست
 ۳) گرامر LL(1) و SLR(1) است
 ۴) گرامر LL(1) هست ولی SLR(1) نیست
 جواب. گزینه 1 درست است

	S	A	C
First	{a}	{b, ε}	{c, ε}
Follow	{\$}	{c,b}	{b}

جدول پارسینگ LL(1) به شکل زیر خواهد شد.

	a	b	c	\$
S	1			
A		2/3	3	
C		4	4	

پس گرامر بالا LL(1) نیست ، اگر دیاگرام و جدول تجزیه این گرامر را رسم کنیم ملاحظه می شود که ، SLR(1) هم نیست

مثال: گرامر زیر برای یک پارسر بالا به پایین پیشگو کننده مناسب نیست مشکل را مشخص کنید و گرامر را طوری بنویسید که برای بالا به پایین پیشگو کننده مشخص کنید و بعد برای گرامر جدید جدول تجزیه LL(1) را بسازید.

- $L \rightarrow Ra \mid Qba$
 $R \rightarrow aba \mid Caba \mid Rbc$
 $Q \rightarrow bbc \mid bc$

اگر گرامری بازگشتی چپ داشته باشد آیا LL(1) هست یا نه؟

خیر چون برای LL(1) بودن باید مبهم نباشد و بازگشتی چپ و.. نداشته باشد

نکته اگر گرامری مبهم نباشد و فاکتورگیری چپ روی آن اعمال شود بازگشتی راست دارد LL(1) است و اگر گرامری بازگشتی راست داشته باشد

ممکن است LL(1) باشد یا نباشد

نکته: اگر گرامر CLR نباشد حتما LALR نیز نیست