

کامپایلر

مراجع

1. Compilers Principles Techniques & Tools ✓

By Sethi and ulhum (2007) EBook
نسخه 2.0

2. Crafting a Compiler
By Fischer (2007)

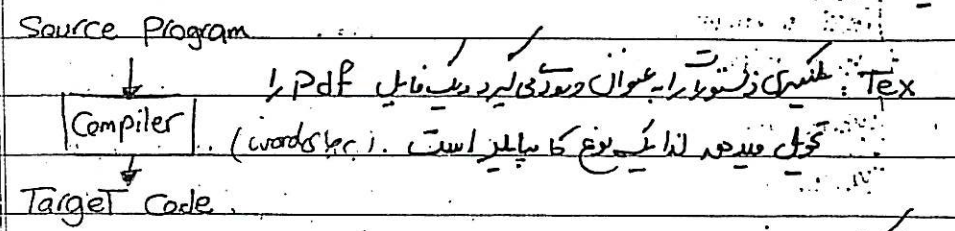
3. The Theory and Practice of Compiler
writing By Thrembley and Soreson (1985)

نظریه : انتهای : تئوری ریاضی ← تئری زیاد
کامپایلر : دانش : کاربرد :
چون برای نوشتن برنامه های کامپیوتری نیاز به کامپایلر داریم

کارهای

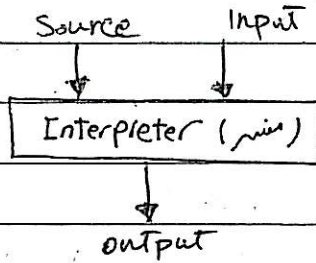
Inter

Source Program به کامپایلر تبدیل می شود Target Code



تفاوت کامپایلر و مترجم:

کامپایلر برنامه را به یک جا اجرایی تبدیل می کند اما مترجم فقط خط به خط را اجرا می کند.



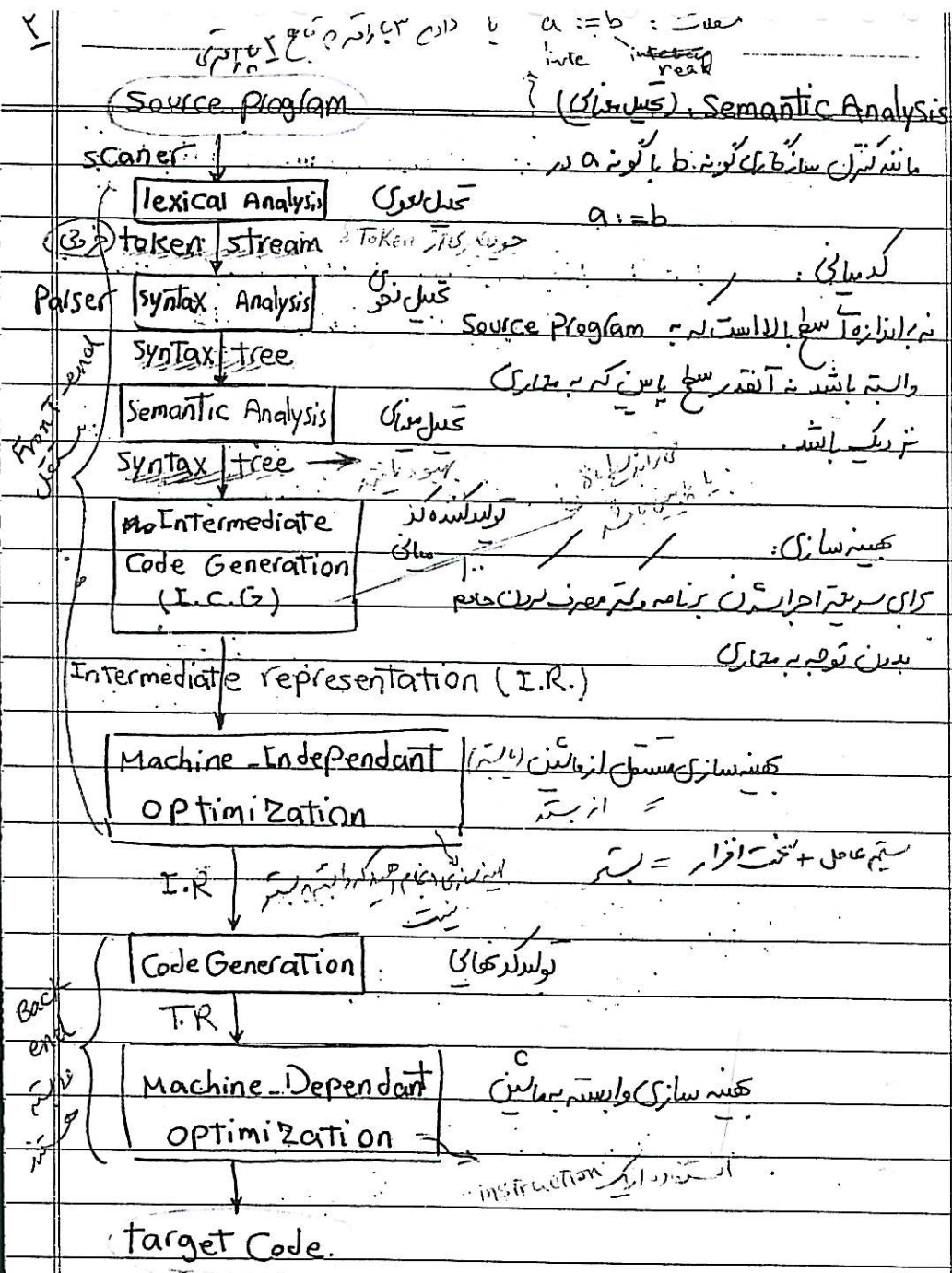
در کامپایلر input نیازیم
و Target مستقل از مترجم
بسیار آسان
مراحل کامپایل:

lexical Analysis: (تحلیل لغوی)

کارشناسان تکلیف لغات برنامه است که مانند عملها، پرانتزها، کاماها و ...
(لغت هر چیزی که به تهاکی لغتی در برنامه ما دارد)

Syntax Analysis: (تحلیل نحوی)

آیین، لغات، قواعد و ترتیب لغات را مورد بررسی قرار می دهد تا مطمئن شویم که
قواعد لغات توسط کامپایلر است. (آیا لغات به ترتیب صحیح آمده اند)
چرا که مشخص کند ترتیب قرار گرفتن آنها؟



از این داخل 5 مرحله اول مستقل از ماشین هستند به آنها Front-end گویند. در 2 مرحله آخر Back-end گویند و وابسته به ماشین هستند.

در این 2 مرحله 2 کامپوزیت داریم که باعث داخل و بیرون در ارتباط هستند

Symbol Table جدول علائم

Identifiers همان متغیرها هستند

error-Handler خطا پرداز

Token: بجای اصطلاح گفته از Token به نشانه ساده می گویند

- Keywords: کلمات طبقه
- Identifiers: نامها
- Operators: عملگرها
- Delimiters: جداکننده ها
- Literals (Constants): ثابتها

مثال: داخل کلمه کامپایل

```
a := b + c * 60
```


attribute

1) اولین مرحله پارسر To Kenize برین رشته های درونی است
 <id, 1> <op, +> <id, 2> <op, +> <id, 3> <op, *> <num, # 60>

جدول نمایی (معمول) Func: array

Name	type	Dim	in	Var	Address
1	a	float	1	var	100
2	b	float	1	"	101
3	c	float	1	"	102
:	:	:	:	:	:

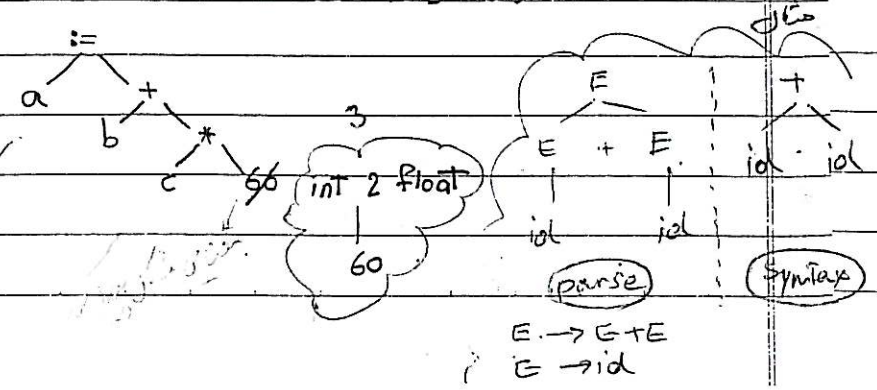
Symbol Table

2) Pars برین رشته ها با استفاده از گرامر الگوریتمی که اگر پارس نشد لذا
 syntax error داریم جواب : Yes / No

Parse Tree : ساختن یک عبارت با گرامر (شامل تمام مراحل میانی و غیر تمینال)

حاصل syntax Tree : نوزدهای مربوط به غیر تمینال ها حذف شده اند و شماره شده اند

ساخته درخت Pars است. syntax Tree:



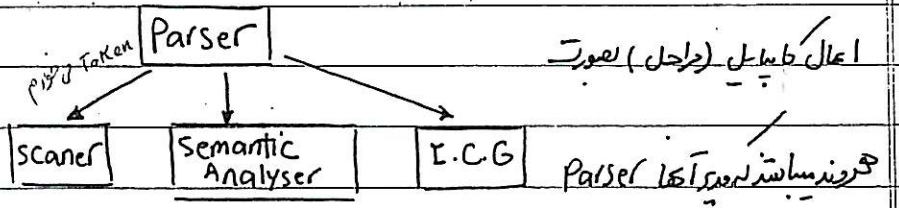
3) تعداد عبارتهای تمینال (افزار شده در تعریف تابع) ریاضیاتی Actual (افزار شده در زمان اجرا) با هم برابرند یا خیر.

در syntax tree چون 60 از نوع int است و از نوع float Semantic A.

60 با 2 از نوع float برشته یک syntax tree بر جای می تابد در وضعیت درونی

Semantic یک syntax tree میگرداند است. Semantic Analyser

مواردی که مربوط به مشکل در برنامه است را می توان تشخیص دهد.



است. مواردی بودن واضح که :
 در یک خط Parser با فعال کرده و آن به Scanner می تولید Token بعد از آن خود

Parser درخت Parse را برشته اش انجام می دهد و بعد تکمیل معنای بلایم داده و بعد تولید

کرده که توسط I.C.G تولید می شود و در آخر نیز دوباره درخواست Token می دهد.

اینکه یک تابع دو پارامتری یک تابع ۳ پارامتری میفرستد : Semantic

درجه صفت که ۷ وجه انجام ندهد. (3 و 4 و 5 از این صفت)

مثال / $500 = 100 + 101 + 500$ (+ , 100 , 101 , 500)
خانه 100 را به خانه 500 بیز $500 = 100$ (= , 100 , 500)
Immediate (= , # 100 , 500)

در برخی کامپایلرها وجود Semantic Analyser برای نسیب خطوط

در بعضی موارد نیاز دارند مثلا زبانهایی بدون آرکیتکچر Optimizer را نیز میتوان حذف

مورد تولید کننده کدهایی را نیز میتوان حذف نمود.

از کدهایی استفاده می شود به دلیل:

1. بهینه تر بودن Optimization به دلیل ساده بودن دستورها

2. Portable شدن برنامه (چون حقوق همکار را نسیب است)

Back-end به 2 عامل وابسته است به همکار داشتن سیستم عامل

حیاتی طراحی کدهایی: Portable Optimization
کدهایی ممکن است بسیار سریع بالا تعریف شود یا در این سطح نزدیک

مهمی باشد trade off است. اگر زبات سطح بالا نوشته شود قابل حمل باشد

کامپایلر Minimum
آرکامپایلر شامل بخش های Scanner, Parser و تولید کدهای

معمولا برای scanner زبانهایی که برای parser زبان های مستقل از من و برای کدهایی حساس به متن کافی است.

نکته: ممکن است یک زبان برنامه نویسی آسانتر ساده باشد به قدرت زبان ها

حساس به متن نیازی نداشته باشد در اینجا Semantic Analyser نیاز

ندارد

تولید کدهایی: از قالب دستور العمل های ساده آرد استفاده می شود

(OP, A1, A2, A3)

کاربرد کدهایی Portability
optimize

مثال کدهایی 1 (int 2, float, # 60, 500)

(* , 102, 500, 501)

(+ , 101, 501, 502) -> سطح پایین

(:= , 2, 502, 500)

2) اگر کدهای Optimizer بهینه

(* , 102, # 60.0, 500)

(+ , 101, 500, 100)

از گاه جمع سرعت کدها بیشتر است

pass =

فرمانها را یک بار میخواند

تولیدهای (3)

```

Mov #60, R1
LD 102, R2(C)
MUL R2, R1 C*60 -> R1
LD 101, R2 (101) -> R2
ADD R2, R1 R2+R1 -> R1
STORE R1, 100 R1 -> 100(a)

```

Peephole Optimization

2 مورد است
 1. جاهایی که متوالاً از صلیبها استفاده کنند. به جای حذف
 در صورتی که سیستمها صلیبها را حذف کنند
 2. در این که استفاده از نرم افزار

برای بهبود کارایی که حاصل زبان ماشین می باشد.

$$a * b + a * c$$

ا در یک سیستم و پارامتر در سیستم دیگر گرفته اند حاصل صلیبها در سیستم a بریزند و نتیجه حاصل است چون دوباره آنها را در load خواهند بود لذا به جای در سیستم b نتیجه می شود آمار سرعت کم می شود

ارتباط Scanner و Symbol Table

اسم Identifier ها در Symbol Table توسط Scanner

ایجاد می شود

ارتباط Parser و Symbol Table

ارتباط Semantic و Symbol Table

type مقیدها اینها می کنند که type خاص Semantic اینها

int a ;

اگر خطا! a, Scanner برای error بریزد

اگر خطا! یک خطا بر آن احتمال میدهد و آنرا به syntax مشکل میزند

نوع Semantic A می دهد به int در type نوشتاری شود و بعد خطا بگردد

a := 2 است می بردی پسند آیا type آنجا یکی است یا نه اگر مشکلی نبود

ادامه دارد می شود.

در عمل نیاز به 2 Pass داریم آنرا اولین بار مقید را می بینیم (مثلاً در a * b) که ادما

است و مقید اظهار نشده (مراجعه کنید به گونه a در گفت گونه اظهار نشده)

گفته: لزومی ندارد اول کار syntax تمام شده باشد تا در Semantic شروع شود

مانند بررسی type در For (is...) قبل از رسم شدن الوداد

ارتباط I.C.G

برای تولید کدی میانی به آدمی مقیدهاست که نیاز به مراجعه به Symbol Table

داریم البته بگردان می آید ICG هم می تواند انجام بدهد.

آیا مشکل بین آنها برود؟
 موازی بودن حاصل کار 2 و 3 : one pass

* در optimize ادعای one pass داریم چون باید چندبار بخوانیم

ارتباط Symbol Table & Optimizer

مشاهده مقادیرهای استفاده شده از حافظه و تعیین آدرس؟

ارتباط Symbol Table & Code Generator

بسته به اینکه سطح کمیاب یا بیش است یا بلا عمل است نیاز به تعیین این

باشه نه ضرب float در float لازم باشه یا int در int سین باید

به ST مراجعه بشود و لیست مقادیر را استخراج کند (ممکن است opcode های

نمی آید برای لازم برای این دو عمل معادله باشد.)

ارتباط Error Handler با مراحل کامپایلر

scanner

lexical error

در زبان c \$ را کاراکتر مجاز نمی آید اگر مشکل کند یا

Identifier ها نباید digit شروع شوند اینجا error کی

لتری هستند Error Handle در مقابل این مسلمات در error کی

لتری، خطاها را Handle می کند و سعی می کند چیزی کامپایلر

متوقف کنیم و معنای خطا را کاربر گزارش دهد

Syntax Error

مثلاً پارنتزها باز نشود و یا باز بسته شدن پارنتزها نامطلوب باشد یا

سعی کالو جا بدهد برای handle کردن خطای کمبود پارنتز بسته می توان

فرض کرد در هر دو خطا گزارش دهد خطاها باید هم گزارش داده شوند و

سعی شود با ما نیز متوقف نشود و خطا handle شوند

Semantic Analyser

warning ها توسط Semantic Analyser گزارش داده می شوند مثلاً

اخطار تخصصی float به int یا استایل مقادیر تعیین نشده باشد type

صورت پیش فرض آن تعیین داده می شود و اگر 2.5 = a: بر روی خطای

type می آید روشن می شود برای این پاراگراف اولی جای که عدد مورد استفاده قرار

بگیرد لیست را تعیین کند که است است اسباه ها در این مرحله باشد و در مراحل

نوعی خطا (Cascade error) error

روشن حکم این است که یک گونه شماره برای این متغیر تعریف نشده در نظر گرفته

در سبب استفاده نوع آن را همین در چون آنها با این خطای undeclared

برنامه نویسی وی را منظم بر اضاها نوع آن نیز می کند.

تولید کن میانی:

با handle کردن خطاها کن میانی تولیدی شود اما اگر خطا وجود داشته

کن Commit می شود در این صورت کن delete می کند تا باطل

ICG با error handle پیام های در مورد گامی نویسنده خاص

برای زنجیره کن میانی یا خطاهای این مبنی داشته

جست در مورد پارسیار و Symbol Table

پارسیار برای جستش در دفتر پارسیار از Symbol Table بهره می برد

محلن است در زمانی برای اضاها توابع آرایه از نحوه مقابل A (E) شماره

ی سویی یا A قبل اضاها شده و نوع آن چیست و توسط Semantic تعیین شده

Pass. یک خواندن که ویری از ابتدا تا آخر

کامپایلرها محلن است یا One-pass یا Multi-pass

واسطی در تولید کردن کامپایلرها One-pass همیشه امداد Optimization بیش از

یک Pass نیاز داریم تا کن میانی را مرتباً جستجو کنند.

تفاوت Interpreter و Compiler ها.

کامپایلر یک بار برنامه را خوانده ترجمه می کند و اجرا می کند اما Interpreter

بارها کند اجرا می کند اما خطا خطا برنامه را می خواند.

زبانها مفسری: Prolog, Pearl, Matlab, Basic

مقایسه: در زبان‌های معین تغییرات خیلی ساده است و راحت می‌توانستی خطا را از آن برداری.

مقایسه: بدون معنی برنامه اجرایی شود در زبان‌های معین معنی است.

syntax error داشته باشد (هرگز یک else نرود) اما در زبان‌ها

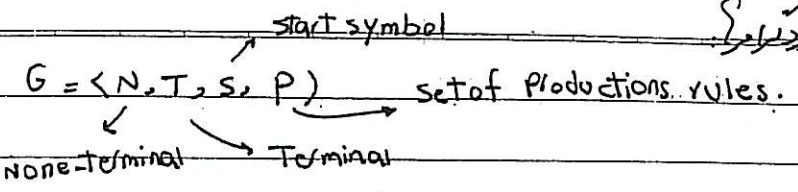
کمپایلری در زمان اجرا syntax گرفته شده

کمپایلر سازگاری در زبان معین مقایسه ندارد.

در زبان‌های کمپایلری نیاز به دانستن source برنامه نداریم اما

در زبان‌های معین source برنامه نیاز است.

نوعی در مورد برنامه



عبارت‌ها $A, B, C \rightarrow D$

نمونه پایانه‌ها $a, b, c \rightarrow$

(عبارت‌هایی که پایانه پیوسته باشند) $x, y, z \rightarrow$ نموده یک عبارت

(طول داشته باشد) $x \rightarrow z$ برای رشته‌ها

نمونه‌های از رشته‌ها $\alpha, \beta, \gamma \rightarrow$

رشته‌های بدون صفر E

عبارت‌ها

$G = \langle N, T, S, P \rangle \equiv \langle \{E\}, \{id, +, *, \wedge, \wedge\}, E, P \rangle$

$P = \{ 1) E \rightarrow E + E$

2) $E \rightarrow E * E$

3) $E \rightarrow (E)$

4) $E \rightarrow id \}$

vocabulary

$V = N \cup T$

انواع لغوها: $\alpha \rightarrow \beta$ (محدود)

1. $\alpha \neq \Lambda$ (Unrestricted) (بدون محدودیت)

2. $|\alpha| \leq |\beta|$ (context-sensitive) (حساس به متن)

کاملاً طول ندارد

انواع لغوها

3. گرامرهای بی‌بافت (Context Free) (مستقل از متن) $\alpha = 1, \alpha \in N$

4. گرامرهای منظم (Regular)

$\beta = 1, \beta \in T$ (مستقل از متن)
 $\beta = 2, \beta = aA$ (منظم از چپ)
 $\beta = 1, \beta \in T$ (مستقل از متن)
 $\beta = 2, \beta = Aa$ (منظم از راست)

BNF (Backus - Naur form)

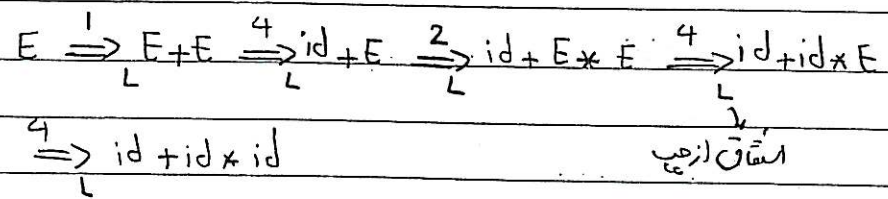
عبارت‌ها را این نوع علامت < > قرار دهد بجای متنی از علامت ::=

استفاده کرده $\langle E \rangle ::= \langle E \rangle + \langle E \rangle$

آی رفع محدودیت تعداد N با قرار دادن حوسبه در $\langle \rangle$ به معنای N

BNF فقط برای گرامرهای مستقل از متن بکار می‌رود

مثال مراحل تولید رشته $id + id * id$

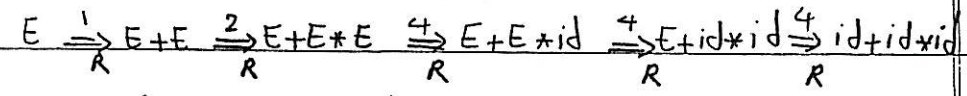


LMD: left - Most Derivation

گزاره تمام مراحل استفاده از چپ ترین غیر پایانی جایگزینی شود

RMD (Right - Most Derivation)

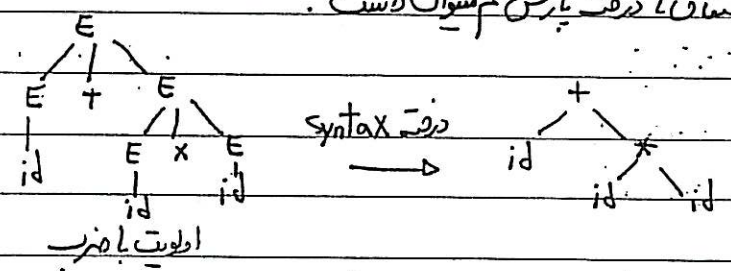
برعکس استفاده از چپ میانه



اگر برای استفاده از روش LMD یا RMD استفاده کنیم اصطلاحاً گفته

Canonical (مابون ضد) است استفاده از قانون مناسب

متناظر با استفاده از درخت پارس هم می‌توان داشت



درخت استفاده از چپ یعنی LMD یا RMD ندارد

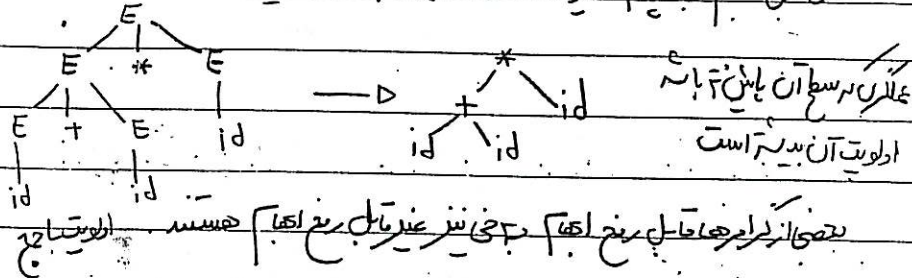
Ambiguous grammer

باید گرامر مستقل از متن می‌گفته می‌شود اگر در زبان آن گزاره حداقل یک

بسته بودن داشته باشند برای آن جمله LMD را بدلیل

RMD را بدلیل تفاوت یا پس داشته باشیم

مثال مثل معجم بود چون میتوان 2 درخت یا پس برای آن کشید



در مثال مثل چون اولویت بین عملوها نبود لذا آنها را برای رعایت

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

اجزای مثال مثل:

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

در اینجا نسبت به جمع یک سطح پایین تر قرار گرفته اند لذا اولویت بدین است

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

اولویت جمع به صورت اول است

مکن است جمع ها هم سطح داشته باشیم اول $a+b+c$ اگر میسر شد $a+b+c$

آی اند اولویت جمع به چه دریم (هیچ کرد) لذا اولویت اول $E \rightarrow E + T$

ابتدا $E \rightarrow T$ و $T \rightarrow E$ چون سرود اولویت جمع است در این سطر

start این تمام بیشتر α رسید
نقطه از ترسیمها

$$L(G) = \{ \alpha \mid \alpha \in T^*, s \Rightarrow^+ \alpha \}$$

SF(G) Step Sentential Form

$$SF(G) = \{ \alpha \mid \alpha \in V^*, s \Rightarrow^* \alpha \} \text{ است}$$

اگر نرم جدا در انتخاب چه درست آمده باشد
در غیر این صورت نقطه چال نرم جدا می گوئیم

رشته های از پایانه و غیر پایانه (می شود از start برای آن رسید $SF(G)$)

$L(G)$ همیشه زیر مجموعه ای از $SF(G)$ است

Reduced grammar

باید 3 خاصیت داشته باشد
1. فاقد قواعدی غیر $A \rightarrow A$ باشد

2. تمام غیر پایانه های آن active باشد

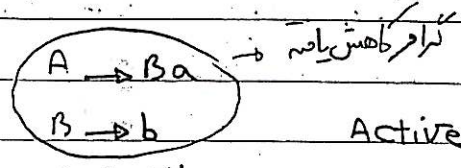
3. اگر نزدیک تمام پایانه ها رسد

$A \Rightarrow^+ \epsilon$ $A \Rightarrow^+ x$

3. تمام غیر خالی‌ها ϵ آن reachable نیست پس

باشند از start برداریم $s \Rightarrow^* \alpha \beta$

اگر خاصیت 2 را نداشته باشد آن برای useless می‌شود



Active = {A, B, D}
 امکان برگشتن به Active
 حسه B, A, D برگشتی است
 اما C غیر برگشتی است

Reachable = {A, B}

حذف شود چون هیچی آن نمی‌برد

تایید هم است اینها Active ها می‌باشند

غیر Active حذف به Reachable می‌شود Unreachable حذف شود

lexical analyser

Scanner

تایید keyword و identifier

Symbol Table نام identifier ها

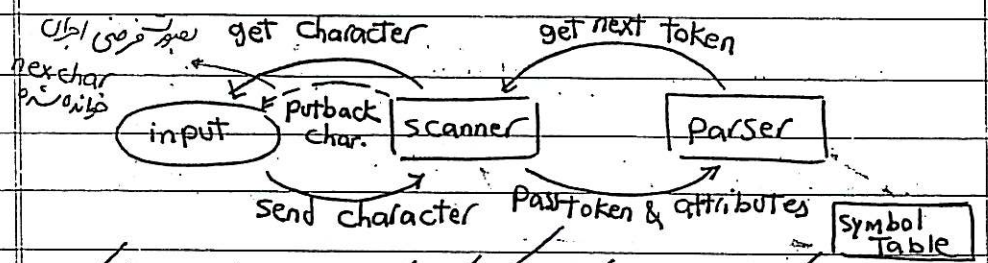
خودترین Comment ها (سنگین و رنگ آبی)

Backup گرفتن از source code

بهترین وقت Compiler optimize کردن Opcode است

زمانی که Compiler می‌خواند Symbol Table است

بهترین وقت optimize کردن است Scanner است



هر Token یک pattern دارد که در آن جستجو می‌شود تا ببینیم خواندن آن

و اگر از الگو پیدا کردیم پس الگو در آن است Pattern علامت‌ها مثل $\langle = \rangle$

بعضی اوقات lexem این pattern است مثل abc12

استخفاف Scanner آن را با attribute آنگاه parser

یازدهمین مدل نظریهٔ اسکینر یا نوع خود عملگر

اگرچه اخیر Scanner را حذف کنیم Parser بجا بیاید Token

کارنامهٔ کویلی می‌گیرد و اگر مثلاً در مثال قبل $E \rightarrow id$ یک Token (Scanner)

لذا برنامه صورت چیزی تر برای آن کارنامه باید نوشته می‌شود طوری که با برنامهٔ قبلی

انجام می‌دهیم در Scanner باید بایست برنامه مستقل از بین انجام دهد و این یک

هزینهٔ اضافی و عدم Optimize بودن را می‌برد در عمل LL(1)

آن نیاز به صرف کردن حافظهٔ بیشتر خواهد بود چون تعداد سطرها بیشتر خواهد شد

(معمولاً است برنامه‌ها) و سبک‌ها (میان‌ها) در نتیجه LL(1) برای حل خواهد شد

در مورد زمان محاسبهٔ مقادیر درست است چون هم زبان‌ها هم متن در هم

زبانها مستقل از متن هرچیز $O(n)$ هستند

کامپایلر

در نتیجه اما کارایی پایین تری آید و اصلاح کردن نیز سخت تر خواهد بود

این Token زبان نه = برسم که سبک می‌شود $ade = x + y$

دوین Token همان = است همین Token و نسبت له + خوانندگی

مثل در زبان فرآیند space اجتناب ناپذیرند چه خوانند چه هستند

assignment $D_{0.5I} = 1.5$

حلقه For $D_{0.5L} = 1.5$

با جدول $D_{0.5I}$ می‌توان Tokenize کرد باید فقط با و برسم لذا سعی اوقات بکن
Tokenize کردن باید زبان بدیشتری بگیرد

در بعضی زبانها مانند PLI، متغیر از keyword ها چون یک نام کار اسماء در مانت

IF then then Else := then; else then := Else;

در اینجا Scanner قادر نیست تفاوت identifier را keyword تشخیص دهد اما

Parser است که تشخیص می‌دهد که identifier است و در هم ^{keyword} Parser با توجه

کل Tokenها تشخیص می‌دهد (برای برنامه تشخیص می‌دهد)

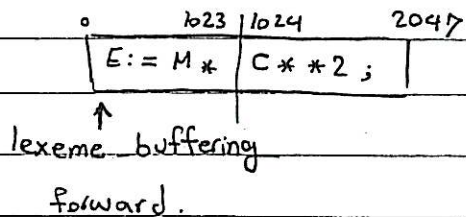
این روش هم برای کامپایلر هم برای برنامه نویسی سخت است

Double buffering
واحد خواندن با همزمان یک است چون از حافظه خوانده می شود.

دوتا pointer متغیر می نامیم

ابتدای یک Token است نشانی نشانی که می خواند lexeme-beginning

کمی که جلوی بردن آن Pattern خارج شود Forward



دلیل استفاده اصلی از double buffering

آنر یک buffer داریم با هم می خواند ابتدای token در حافظه

با reload کردن آن از این می رود و می خواند این Token در

با فریب است باقی می ماند

ابت
حد اکثر نماز Token ، 1024 است

Forward
if Forward = 1023 then load 2nd half of buf
Forward ++

else if Forward == 2047 then load 1st half of buf
Forward := 0

else
Forward ++

باید بداند رساندن این معنی که کار کرده است طرک تمهیدات Sentinel

کرده و وقتی به انتهای بافر رسیدیم کار در نظر برای آن چه می کنیم و اینها کارهای بسیار است

Forward ++

if Forward ↑ = 'eof' then

if Forward == 1023 ...

else if Forward == 2047 ...

else

TERMINATE

Regular Expression (RE) عبارات منظمه

اجزای آنست برای توصیف یک زبان منظم یا گانش یک زبان منظم

رشته‌های طول 2 مشتق از a و b: $(ab)^* (a|b)(ab)^*$

گرامر: اجزای برای تولید رشته‌های زبان

- 1/ \emptyset is a RE; $L(\emptyset) = \{\}$
- 2/ λ is a RE; $L(\lambda) = \{\lambda\}$
- 3/ $a \in \Sigma$ is a RE; $L(a) = \{a\}$
- 4/ if r_1 and r_2 are two REs

$L = \{a, \dots, z\} = [a-z]$
 $D = \{0, \dots, 9\} = [0-9]$

مثال id: $L \cdot (L|D)^*$

مثال عبارت منطقی بنویسید که گزینه‌های زبان زیر باشد؟

$r_1 | r_2$ $L(r_1 | r_2) = L(r_1) \cup L(r_2)$
 $r_1 \cdot r_2$ $L(r_1 \cdot r_2) = L(r_1) \cdot L(r_2)$
 (r_1) $L((r_1)) = L(r_1)$
 r_1^* $L(r_1^*) = (L(r_1))^*$

باز رشته‌های مشتق از a و b که بیشتر از 2 اسیب رسم نداشته باشد

$(b|ab|aab)^* a? a?$
 $\epsilon | a | aa$

تسری عبارت منطقی بنویسید که گزینه زبان زیر باشد؟

are REs too

$a^* = \{\epsilon, a, aa, \dots\}$
 $(ab)^* = \{\epsilon, ab, abab, \dots\}$
 $(a|b|ab) = \{a, ba, b\}$
 $r \cdot \epsilon = \epsilon \cdot r = r$
 $r \cdot (s|t) = rs | rt$
 $(r^*)^* = r^*$ $a? = a | \epsilon$
 $r \cdot r^* = r^+$
 $(a|b)^* = (a^*b^*)^* = (a?b?)^*$

مثال

letter identifier زبان میسرمانند C و پایتون مشتق از letter

digit و underscore ادکویت زیر

underscore با * نباید

underscore لست رسم نباید

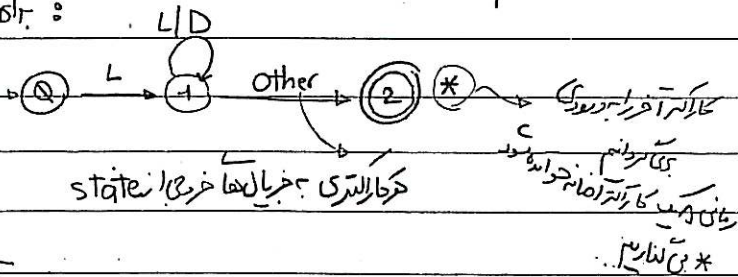
lex Scanner Generator

بزرگ Scanner ها در زبان معنی ندارند لذا Scanner استوانه زبان معنی برای lex توصیف کرد

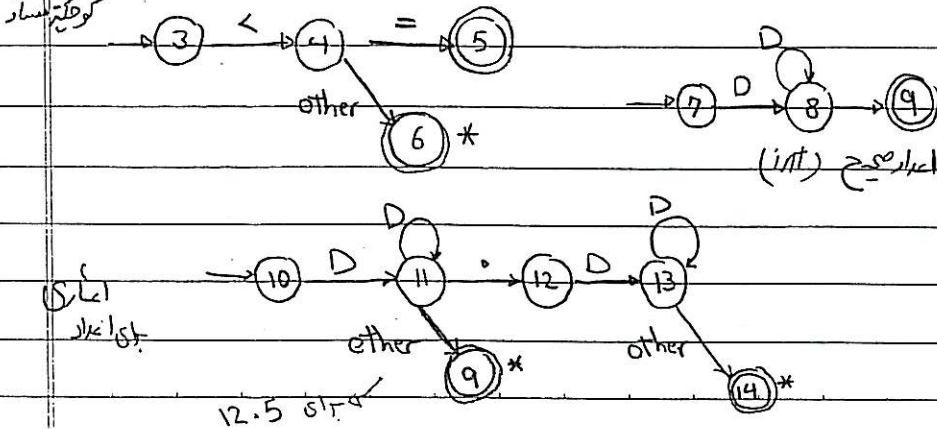
lex: یک زبان معنی را به یک ماشین DFA یا NFA تبدیل می کند

DFA سریعتر از NFA اجرا می شود (O(n)) با معنی که در DFA بسیار ساده است چون تمام حالات بررسی می شوند لذا مطلوب نیست از DFA اکثر استفاده می کنیم

Identifier



Token st



با استفاده از DFA از یک Switch Case استفاده می کنیم

Switch <ch> of



'D':

''':

'<':

end

ساده ترین ماشین با یک Scanner است.

در استفاده از Switch Case ترتیب بررسی Case ها مهم نیست

اگر else if استفاده کنیم آنگاه خطا پس تکرار بیشتری دارد زودتر نوشته می شوند

اگر Scanner ظاهر استخوان ندارد Error handler داریم

روشها وجود با خطاهای لغوی

1 Error Recovery (بسیار بیکی پوشش باز یافتن خطا)

خود را از سر خطا خلاص کنیم

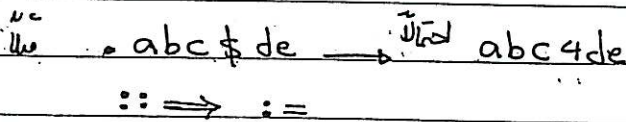
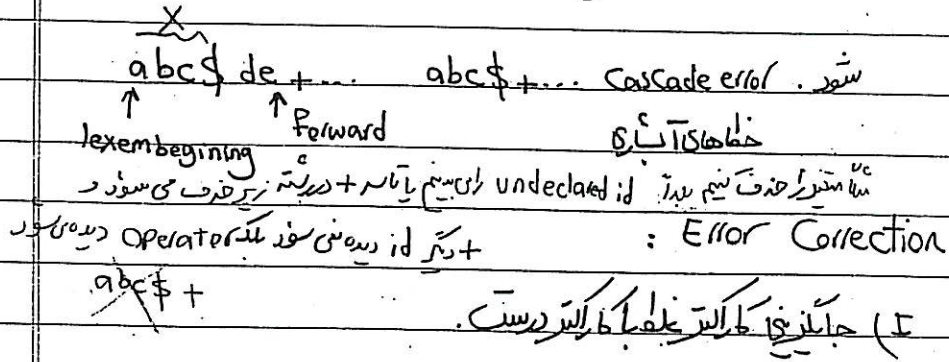
2 Error Correction (تصحیح خطا)

خطا را بصورت صحیح تصحیح کنیم

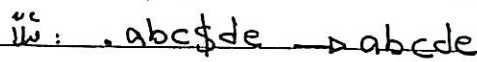
Panic Mode

تجرب زیادی از ورودی را ممکن است در بریزیم (از استای توکن)

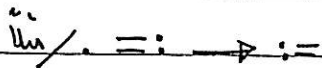
تاکل رفتار خطا در بریزیم (ممکن است بوی cascade errors)



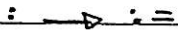
(II) حذف کاراکتر درست.



(III) جای جای کاراکتر درست.



(IV) امان کردن کاراکتر گسسته



Minimal distance

با استفاده از یک distance تعریف کنیم

کدام distance آن 2 است. آنکه با کمترین فاصله را مشخص کنیم.

Scanner وظایف

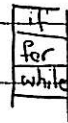
تولید فایل زبان سین identifier keyword ها. یک کاربر

این است زیرا keyword ها ما این را می بینیم با ورودی هر کس

keyword تشخیص داده می شود. اینجا در یک جمله از keyword

است. هم identifierها با این جمله یک کنیم این جمله می تواند

قسمتی از حفظ اثر یا سیمی از Symbol Table این جمله امان



هم

پژتوگراف (get token, install id) و (get token)

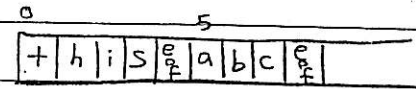
parser فراخوانی می کند به id - install فراخوانی شود اگر id

قبل از این که به شماره برای هر یک از اینها در مان برای Parser بروی کرانیم

اما بعد از (32-1) که برای هر یک از اینها keyword است و شماره

Keyword می بینیم که اینها را می بینیم و Parser می بینیم که

درام keyword است



Name		Symbol Table	
0			
5			
:			

keyword را طول بسیار کمی نسبت به اینها دارند = ذخیره آنها

در st بسیار کم است (از جای حافظه) = در قسمتی از حافظه طالت

کدی را بعد از این که ذخیره می کنیم
و اینها را

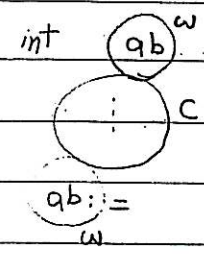
قدرت زیاد است برای Scanner (با اینست)

Syntax Analyzer نتواند از مابقی استفاده کند چون تعدادی از آنها

بازدیده اند یعنی با آن زبان $a^n b^n$ انسان می تواند این را تشخیص دهد.

قدرت زیاد است از این برای Syntax Analyzer می باشد اما برای

Semantic Analyzer می نیست. مثلاً در مورد تعریف \rightarrow بودن مقبول در آن زمان



با اینها به این رجوع کرد \rightarrow WCU

که با هم استفاده کنیم

پارسی (Parsing)

مثال $S \rightarrow cAd$

$A \rightarrow ab|a$

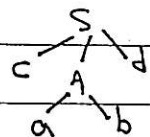
$Cabd$

Cad

مثال Cad را فونمی کنیم و مشاهده می‌کنیم که این عبارت است.

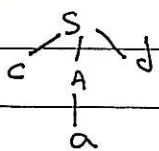
Current token \rightarrow look-ahead (L.A.) token

L.A. = c



L.A. بعد از a است. مقدار a match کرد.

اما از b، حاصل نمی‌شود چون token بعد از آن است. لذا می‌توان ادیت کرد این



خطا فونمی را برطرف نمود. با نیست. Back track کنیم.

هدف: پیدا کردن روش‌های Parse (بین خطاها فونمی)

ما دنبال Parser هایی هستیم که مقدار مراحل آن قابل پیش بینی و محدود باشد

اینه جنبه بار Back tracking داریم می‌کنیم. چه تکرار Back track

(پارسی) عبارت (کلمه یا جمله) است. ما دنبال پارسی هستیم که بتوانیم به عبارت نگاه کنیم

آرگومان مسئله از من و قطعی نباشد. اما نظر اینکه بتوانیم پارسی را به عبارت

ساخته می‌کنیم نیست

ما دنبال کلمه‌ها مسئله از من غیر قطعی نیستیم. چون اگر غیر قطعی باشد از این کلمات پارسی نمی‌توان استفاده کرد. پس نیست. به عبارت دیگری شود

روش‌های Parse:

به دو دسته کلی تقسیم می‌شوند

1. Top-Down (بالا به پایین) \rightarrow LMD

از start شروع می‌کنیم و سعی می‌کنیم آن رشته را ایجاد کنیم

2. Bottom-up \rightarrow RMD (پایین به بالا)

در رشته عمل انتقال می‌کنیم تا به start برسیم

روش‌های Top Down (پایین به بالا) عبارتند از:

Recursive Descent

اجزای آن است که می توانیم
 اجزای آن که توکن به هم وصل می کنند

$$LL(1) \subseteq LR(k)$$

رشته ورودی را از چپ می خوانیم رشته ورودی k توکن k توکن

یا خوانیم هر چه k غیر تر باشد کما و قوی تر است و آن بعد از این

می توان خوانی به سمت راستی پیش قرار داد

در دو حالت یا این یا آن را به هم وصل می کنیم

1. Operator Precedence (OP) تقدم اول

2. Simple Precedence (S.p) تقدم ساده

$$LR(1) \subseteq LR(k)$$

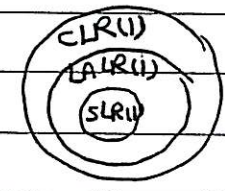
L رشته ورودی از چپ به راست خوانده می شود

R از چپ به راست خوانده می شود

k توکن k توکن می خوانیم

1. SLR(1) : Simple LR برای هر از زیر است
2. LALR(1) : look ahead LR
3. CLR(1) : Canonical LR

توکن است که در این رشته به هم وصل می کنند
 CLR(1) و LALR(1) و SLR(1)

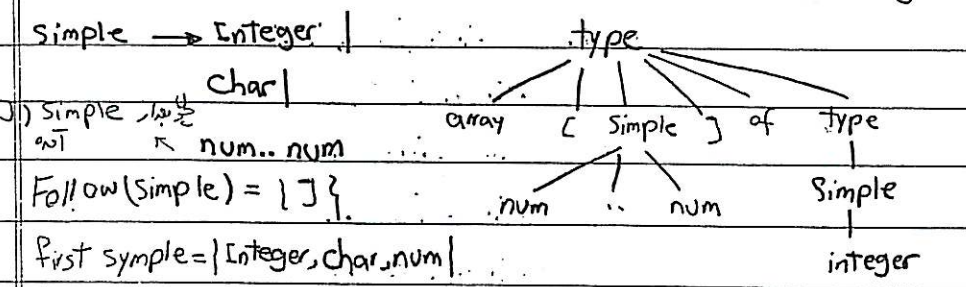


ارزای قدرت

Top-Down Recursive Descent

ابتدا یک مثال کما و قوی شروع می کنیم

array [num..num] of integer
 Simple Integer



First sample = { Integer, char, num }
 First (array) = { array, first (sample) }

match توکن می گیریم که در آن توکن همان توکن LA است یا نه
 اگر نه error می دهیم

بازار هم غیر باید به هم وصل می کنند

```

Proc match (t)
begin
  if t == L.A then L.A1 = Scanner;
else
  error;
end
Proc Type
begin
  if L.A == 'array' then
    begin match ('array');
      match ('[');
      Call Simple;
      match (']');
      match ('of');
      Call Type;
    end;
  else if L.A == '^' then
    begin
      match ('^');
      match ('d');
    end;
  
```

پروسیجر قبسطه در اسر
(بر اهر غیر بلانیه در جوی دارد)

```

else if L.A ∈ {integer, char, num} then
  Call Simple;
else error;
end Type;
Recursive است و معرف نیست  

Recursive function Call دارد  


```

دستور match () کرون برای بلانیه است
کارای این روش به علت مسا کرون تابع پارس است
گرامر باز استی دارد لذا Recursive function Call ندارد
کرون موردی کرا مر سوال شده هدف کرا مر است قرار نیست تبدیل کنیم
نوع کرا مر در رابطه ما مستقل از متن است اگر توانی ناشی از این None Terminal
انده این مرتبه دانسته باشند دو متمم از این نوع رفت
STL → begin STS end
STS → ε | other

```

Proc STL
begin
  if L.A == 'begin' then
    begin
      match ('begin');
      Call STS;
    
```



```

match('end');
end
else error;
end
Follow(STS) = |end|

proc STS
begin
if L.A. == 'other' then
    match('other');
else if L.A. == 'end' then
    return;
else error;
end
    
```

match, ε می خواهد که LL نیست
L.A. == 'end' then
return;

توابع:

$$First(\alpha) = \{ \text{پایه های LL} \} = \{ \alpha \mid \alpha \Rightarrow^* \alpha\beta \}$$

پایه های LL از وقتی که شروع کنیم می توانند در ابتدای رشته واقع شوند

نکته:

اگر تعدادی به شکل $A \rightarrow \alpha \mid \beta$ داشته باشیم شرط اول

$$1) First(\alpha) \cap First(\beta) = \emptyset$$

اگر این شرط منقضی شود از روش Recursive descent نمیتوان استفاده کرد

$$Follow(A) = \{ \text{پایه ها} \} = \{ b \mid S \Rightarrow^* \alpha A \beta \}$$

پایه ها در پدیده غیر LL آنها میتوانند ظاهر شوند اما همیشه در شروع های LL

خوبتر از E جزو Follow می باشد هرگز Follow نمی شود.

Recursive Descent

① RD شرط دوم

$$\text{if } \alpha \Rightarrow^* \epsilon \text{ then } First(\beta) \cap Follow(A) = \emptyset$$

اگر تکراری هر دو باشد ① و ② را باید به اصطلاح آن نامز را در LLU نام

Proc A

```

begin
if L.A ∈ First(α) then use A → α
else if L.A ∈ First(β) then use A → β
* else if L.A ∈ Follow(A) then use A → ε
    
```

* (معمولاً) A نتواند به ε برود اگر نتواند به ε برود این دستور را نمی نویسیم (این خط که را نمی نویسیم)

```

else
error;
end
    
```

همه چیز با کامپایلر LL(1) هستند چون همه چیز با کامپایلر DNA دارند اما هرگز تکراری که LL(1) با هم نمی آید

$E \rightarrow E+T \mid T$ PROC E
 3) $T \rightarrow T * F \mid F$ begin
 5) $F \rightarrow (E) \mid id$ if $L.A \in First(E+T)$ then
 Call E;
 match ('+')
 Call T;
 end
 else if $L.A \in First(T)$ then
 Call T;
 else
 error;
 end E.

$E+T \Rightarrow T+T \Rightarrow F+T \Rightarrow id+T$

گرامر فوق بازگشتی چه طوری E را در T می شناسد؟

گرامر می تواند با جز T شروع می شود و چه طوری در T.

$First(T) \cap First(E) = \{id\} \neq \emptyset$

گرامر نمی تواند با id باشد (LL(1) نیست)

حاصل حذف چپ گردی از گرامر LL(1) است. منظور از چپ گردی این است که در یک گرامر LL(1) هر دو چپ گردی نباید در یک گرامر باشد.

$A \rightarrow \alpha x$
 $B \rightarrow \alpha y$

تبدیل چپ گردی به چپ گردی

$A \rightarrow \alpha_1 \alpha_2 \dots \alpha_n$ چپ گرد

$A \rightarrow \beta_1 \dots \beta_m$ غیر چپ گرد

$\hookrightarrow (\beta_1 \dots \beta_m) (\alpha_1 \dots \alpha_n)^*$

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

$A' \rightarrow \epsilon \mid \alpha_1 A' \mid \dots \mid \alpha_n A'$

چپ گردی از بین رفته

نکته:

چپ گردی گرامر که LL(1) نیست همیشه وجودش وجود دارد

I حذف چپ گردی

II حذف چپ گردی از چپ

مثال گرامر قبل از حذف چپ گردی از چپ

- 1) $E \rightarrow TE'$
- 2,3) $E' \rightarrow \epsilon \mid TE'$
- 4) $T \rightarrow FT'$
- 5,6) $T' \rightarrow \epsilon \mid *FT'$
- 7,8) $F \rightarrow (E) \mid id$

نحوه محاسبه $First(x)$
 $First(x) = \{x\}$ اگر x یک پایانه باشد آنگاه
 (x غایب حالت ندارد است)

2. اگر قاعده ای $x \rightarrow \epsilon$ داشته باشیم آنگاه ϵ را به $First(x)$ اضافه می کنیم

3. اگر قاعده ای $x \rightarrow y_1 y_2 \dots y_k$ داشته باشیم آنگاه:
 $First(x) = \{ \epsilon \} \cup First(y_1)$ اگر $First(y_1) \neq \{ \epsilon \}$
 $First(x) = \{ \epsilon \} \cup First(y_2)$ اگر $First(y_1) = \{ \epsilon \}$
 و همین روند را ادامه می دهیم تا جای y_{k-1} آنگاه $y_{k-1} \Rightarrow \epsilon$ را به $First(x)$ اضافه می کنیم

$A \rightarrow BC$
 $B \rightarrow \epsilon$
 $First(B) = \{ \epsilon \}$
 $\Rightarrow First(A) = \{ C \}$
 $\hookrightarrow First(C) = C$

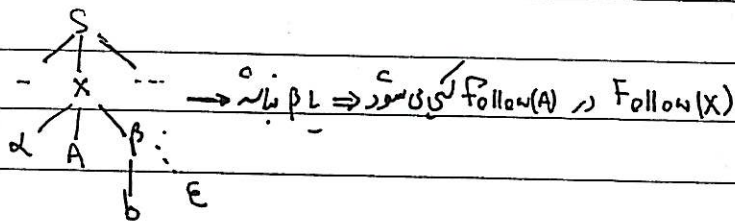
$A \rightarrow Bb | e$ $First(A) = \{ e, b \}$
 $B \rightarrow Ac | b$ $First(B) = \{ e, b \}$

نحوه محاسبه $Follow(A)$
 اگر A حالت شروع باشد آنگاه $\$$ را به $Follow(A)$ اضافه

2. اگر قاعده ای $\alpha A \beta$ داشته باشیم آنگاه β را به $Follow(A)$ اضافه

اگر $Follow(A)$ خالی نباشد (تالی ϵ هرگز جزو $Follow(A)$ نیست)

3. اگر قاعده ای $\alpha A \beta$ داشته باشیم با قاعده ای $x \rightarrow \alpha A \beta$ داشته باشیم آنگاه $\beta \Rightarrow \epsilon$ را به $Follow(x)$ اضافه



مثال (I)

- 1) $E \rightarrow TE'$ $First(E) = \{ id \}$
- 2) $E' \rightarrow e | + TE'$ $First(2) = \{ e \} \cap First(3) = \{ + \} = \emptyset$
- 3) $T \rightarrow FT'$ $First(4) = \{ id \}$
- 5) $T' \rightarrow \epsilon | * FT'$ $First(5) = \{ \epsilon \} \cap First(6) = \{ * \} = \emptyset$
- 7) $F \rightarrow (E) | id$ $First(7) = \{ (\} \cap First(8) = \{ id \} = \emptyset$

Follow(E) = { \$,) }
 این به Follow می‌گویند که در آنجا به کار می‌رود

Follow(E') = { \$,) }
 می‌گویند

قاعده E' می‌تواند 1, 3 است

$$Follow(T) = \underbrace{First(E')}_{\{+\}} + \underbrace{Follow(T)}_{Follow E + Follow E'}$$

$$= \{+, \$,)\}$$

$$Follow(T') = Follow(T)$$

$$Follow(F) = \underbrace{First(T')}_{\{+\}} + \{ \$,) \} = \{+, \$,)\}$$

در 6.6 ظاهره
 ارتباط هم به هم می‌زنند

*
 می‌تواند در 6.6

Follow(T')

Follow(F)

مثلاً قانون دوم LL(1) را می‌توانیم می‌بینیم که

$$\text{First}(P) = \text{First}(+TE') = \{+\}$$

$$Follow(E') = \{ \$,) \} \cap \{+\} = \emptyset$$

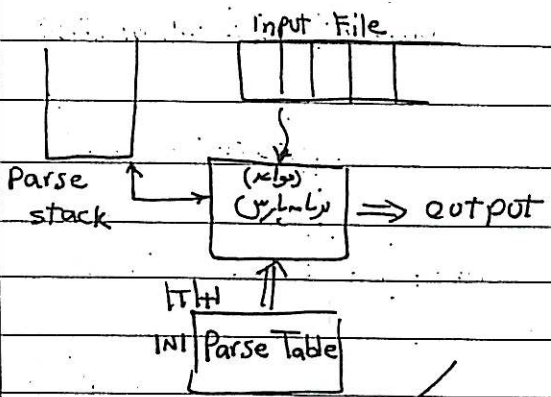
$$5,6 \quad First(*FT') = \{*\}$$

$$Follow(T') = \{+, \$,)\} \cap \{*\} \Rightarrow \emptyset$$

درسته LL(1) است چون هیچ‌کدام از اینها هم
 برقرار است

LL(1) \subset LL(k)

LL(1)



ساختار یک پارسیر LL(1):

در جدول Parse می‌نویسند

انگاه None-Terminal ها

می‌تواند + 1 باشد

هستند

باید در LL(k) استناد بتوان این را در LL(1) کرد.

Parse Table

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

S بران sync در خط میانی

قراری می‌شود

خاسته از Parse Table می‌تواند به سطر و ستون به سطر و ستون

نوعی از عمل Parse برش (LL(1)):

مراحل آغازی:

در ابتدای بار یک علامت \$ به انتهای رشته ورودی اضافه کرده و S را بصورت



برعکس وارد stack می‌شیم.

در هر مرحله از Parse با توجه به علامت بالای stack و توکن جاری (L.A.)
 $(PS(Top) = X)$

1. اگر $x = L.A. = \$$ باشد خاموشی عمل Parse اتمام می‌شود.

2. اگر $x = L.A. \neq \$$ (x یا با توجه به آنجا اسکن می‌شود)

نه توکن جدید را می‌تواند و x از بالا stack حذف می‌شود. (اگر x با y باشد)

استنتاج L.A. تطبیق کند آن‌ها به خطا می‌رود (قوانین)

برعکس طرح stack می‌کنیم

3. اگر علامت بالای stack غیر از آن باشد به همراه توکن جاری با

خطای PT (x, L.A) را چاپ دهد بصورت زیر عمل می‌کنیم.

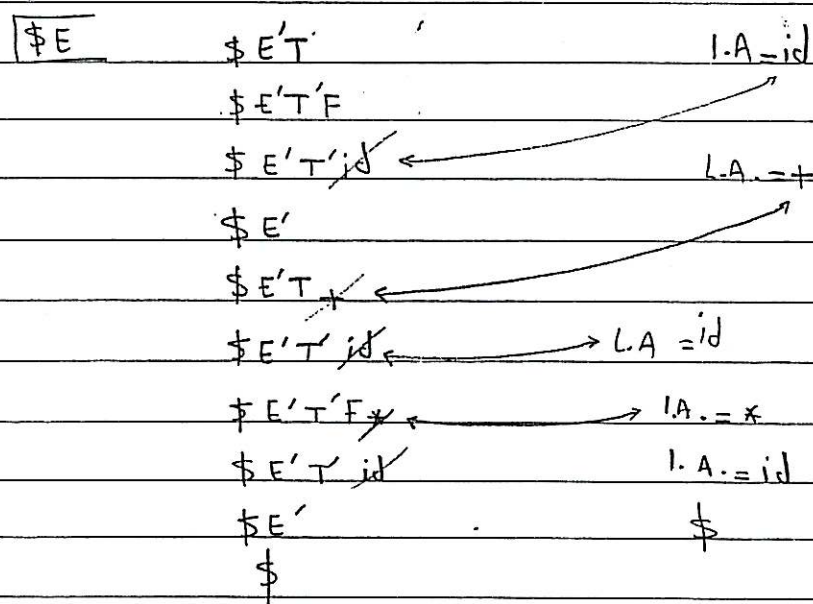
الف) اگر $PT(x, L.A) = x$ باشد آن‌ها از بالا stack حذف کرده

و α بصورت برعکس وارد stack می‌شیم. (اول سمت چپ و بعد غیر از آن که شش می‌باشد)

ب) اگر $PT(x, L.A)$ خطای استنتاج به خطا می‌رود و باید

خطا را چاپ کرده شود.

مثال: $id + id * id$ (با توجه به مثال I)
 ابتدا اسکن \$ در انتهای رشته



Accept

خوبه بپریم چیکه پرسی (1) LL(1) به بازار هر رابطه $A \rightarrow \alpha$ بصورت α عمل می کنیم

1. مقابل غیر پایه A وزیر عناصر (α) First (A) مایه $\alpha \rightarrow PA$

قرار می دهیم

2. اگر α در صورت α باشد به α برود آنگاه مقابل غیر پایه A وزیر عناصر

Follow(A) شماره قواعدی $A \rightarrow \epsilon$ ابتدا می کنیم

تبدیل آنگاه با تمام برون قواعدی $F \rightarrow \epsilon$ در LL(1) می ماند؟ و اگر می ماند

چیکه LL(1) آن را می کشد؟

1.2) $A \rightarrow Bd | Ce$ $first(A \rightarrow Bd) = \{b, d\}$ (مال)

3.4) $B \rightarrow bB | \epsilon$

5.6) $C \rightarrow cC | \epsilon$

$first(C \rightarrow \epsilon) \rightarrow Follow(C) = e$

	b	c	d	e	\$
A	1	2	1	2	
B	3		4		
C		5		6	

Follow(C) = e

اگر LL(1) نبود در خانه ها داخل داریم

چون در هیچ خانه ای تقاطع نداریم \rightarrow LL(1) است.

مدعی اینست که $\$$ خالی است یعنی اینکه هیچ انتظاری نداریم به $none-terminal$ دیده شود یعنی انتظاری نداریم برای $stack$ $none-terminal$ باشد ما $\$$ می بینیم

دخوهی پر خود با خط در روش پارس LL(1):

1. Panic Mode
2. phase level
3. Error productions
4. Global correction

Panic Mode

بیشتر از روشی که ممکن است حذف کند (1) یا پایه $stack$ A و جوابی نداشته باشد و یا (2) با LA NT $stack$ بخانه خالی برسیم شرایط $error$ است.

Synchronizing (مهاجرت کننده):

بازار غیر پایه α تعیین می شوند. تا یک غیر پایه α جسی از رشته ورودی P

تولیدی کنیم اگر همین کار تولید رشته غیر پایه α به خاطر رسیدیم طاکار آن P

ناریدیم بگیریم \rightarrow چه است \rightarrow موق \rightarrow برای $Follow$ آن NT می گیریم

با دور بریزیم