

به نام خدا

اصول طراحی کامپایلر - شرح پروژه دانشگاه اصفهان

مقدمه

در این پروژه قصد داریم برای یک زبان برنامه‌نویسی به نام PL یک کامپایلر طراحی و پیاده‌سازی کنیم. زبان PL یک زبان دستوری شبیه به زبان C است. یک کامپایلر از تحلیل‌گر لغوی، تحلیل‌گر نحوی، تحلیل‌گر معنایی، تولیدکننده و بهینه‌ساز کد میانی، و تولیدکننده و بهینه‌ساز کد اسمبلی تشکیل شده است. یک کامپایلر کامل، یک برنامه نوشته شده در یک زبان برنامه‌نویسی را از یک فایل متنی دریافت کرده، معادل کد اسمبلی آن را تولید می‌کند. کامپایلر طراحی شده در این پروژه یک کامپایلر کامل نیست، بدین معنا که این کامپایلر کد اسمبلی تولید نمی‌کند. کامپایلر طراحی و پیاده‌سازی شده در این پروژه، یک برنامه به زبان برنامه‌نویسی PL را دریافت کرده، توسط تحلیل‌گر لغوی توکن‌های آن را استخراج می‌کند، و توکن‌های استخراج شده را به تحلیل‌گر نحوی ارسال می‌کند. تحلیل‌گر نحوی، بر اساس گرامر زبان، و یک یا دو الگوریتم تجزیه (مانند تجزیه بالا به پایین و پایین به بالا)، یک درخت نحوی تولید کرده و در صورت وجود خطا، پیام خطا صادر می‌کند. در صورتی که در برنامه ورودی خطایی وجود نداشت، کامپایلر برنامه را پس از تحلیل معنایی اجرا می‌کند.

پیاده‌سازی کامپایلر توسط هر زبانی (مانند سی++، جاوا، و پایتون) می‌تواند انجام شود. تهیه یک گزارش جامع برای توضیح نحوه طراحی هر یک از قسمت‌های کامپایلر الزامی است. در طراحی و پیاده‌سازی کامپایلر دانشجویان می‌توانند از هر منبعی استفاده کنند، اما باید منبع را حتما ذکر کنند. در صورت استفاده از موتورهای هوش مصنوعی، لازم است منبعی که موتور هوش مصنوعی استفاده کرده است ذکر شود. برای مثال بینگ، منابعی را که از آن اطلاعات استخراج می‌کند ذکر می‌کند.

۱ تحلیل‌گر لغوی

در این قسمت از پروژه می‌خواهیم یک تحلیل‌گر لغوی برای زبان PL طراحی و پیاده‌سازی کنیم. تحلیل‌گر لغوی یک فایل حاوی یک برنامه را دریافت می‌کند. در صورتی که برنامه حاوی توکن‌های معتبر در زبان PL باشد، دنباله‌ای از توکن‌های تشخیص داده شده چاپ می‌شوند، و در غیر اینصورت پیام خطا چاپ می‌شود. واژه‌های زبان حساس به حروف کوچک و بزرگ¹ هستند، بنابراین X و x دو کلمه متفاوت‌اند.

۱.۱ کلمات کلیدی

کلمات کلیدی زبان با حروف کوچک نوشته می‌شوند که در زیر ذکر شده‌اند.

```
bool break char continue else false
for if int print return true
```

۲.۱ شناسه‌ها

یک شناسه نامی برای یک موجودیت در یک زبان برنامه‌نویسی است. دو موجودیت در این زبان برنامه‌نویسی عبارتند از متغیر و تابع. یک متغیر موجودیتی است که یک یا دنباله‌ای از خانه‌های حافظه را اشغال می‌کند و دارای یک نام است. یک تابع موجودیتی است که تعدادی ورودی دریافت می‌کند، محاسباتی را انجام می‌دهد و در برخی موارد یک مقدار بازمی‌گرداند. یک تابع با یک نام مشخص می‌شود. شناسه‌ها با یک حرف یا علامت زیرخط² آغاز می‌شوند و می‌توانند حاوی ارقام، حروف و علامت‌های زیرخط باشند. شناسه‌ها نمی‌توانند برابر با هیچ‌یک از کلمات کلیدی باشند.

۳.۱ علامت‌های نشانه‌گذاری

علامت‌های نشانه‌گذاری در این زبان به شرح زیر هستند.

¹ case-sensitive
² underline

- علامت‌های آکولاد باز { و آکولاد بسته } در تعریف بلوک‌ها استفاده می‌شوند.
- علامت‌های پرانتز باز (و پرانتز بسته) در تعریف توابع، فراخوانی توابع، و عبارت‌های محاسباتی استفاده می‌شوند.
- علامت‌های گروه باز [و گروه بسته] برای تعریف آرایه‌های استفاده می‌شوند.
- علامت ویرگول ، برای جدا کردن ورودی‌های تابع از یکدیگر در تعریف و فراخوانی تابع استفاده می‌شود.
- علامت نقطه‌ویرگول ; در پایان تعریف متغیرها، دستورات محاسبه‌ای و فراخوانی توابع، و همچنین در تعریف حلقه‌ها استفاده می‌شوند.

۴.۱ توضیحات

توضیحات³ با دو علامت اسلش⁴ یا خط اریب // آغاز می‌شوند و با کاراکتر پایان خط معادل کد اسکی^۵ یا \n پایان می‌یابند. تحلیل‌گر لغوی توضیحات را به تحلیل‌گر نحوی ارسال نمی‌کند، اما پس از تحلیل لغات لیست همهٔ توکن‌ها را چاپ می‌کند.

۵.۱ مقادیر عددی

مقادیر عددی می‌توانند در مبنای ده (دهدهی یا دسیمال)^۵ یا در مبنای شانزده (شانزده‌شانزدهی یا هگزادسیمال)^۶ باشند. یک عدد دهدهی می‌تواند مثبت یا منفی باشد که در صورت منفی بودن با علامت - آغاز می‌شوند. اعداد هگزادسیمال با دو کاراکتر 0x آغاز می‌شوند.

۶.۱ کاراکترها و رشته‌های ثابت

یک کاراکتر ثابت، یکی از حروف الفبای اسکی^۷ است. یک کاراکتر ثابت بین دو علامت آپوستروف ' ' قرار می‌گیرد. برای نشان دادن علامت آپوستروف در یک کاراکتر ثابت از '\'' و برای نشان دادن کاراکتر خط اریب وارون (بک‌اسلش)^۸ از '\\\'' استفاده می‌شود. یک رشتهٔ ثابت را با دنباله‌ای از کاراکترها که در بین دو علامت نقل قول " " قرار گرفته‌اند، نشان می‌دهیم. برای نشان دادن علامت نقل قول در یک رشته ثابت از "\" استفاده می‌شود.

۷.۱ عملگرها

در یک عبارت می‌توان از عملگرهای حسابی^۹ مانند + ، - ، * ، / برای جمع، تفریق، ضرب، و تقسیم، استفاده کرد. برای به دست آوردن باقیمانده از عملگر % استفاده می‌شود. عملگرهای یگانی + و - برای تعیین مثبت و منفی بودن اعداد به کار می‌روند. عملگرهای یگانی + و - بالاترین اولویت را دارند و پس از آنها * ، / ، % هم اولویت بوده و در درجهٔ دوم اولویت قرار دارند و در نهایت + و - اولویت سوم قرار می‌گیرند.

عملگرهای رابطه‌ای^{۱۰} > ، < ، <= ، >= و != برای مقایسه دو مقدار به کار می‌روند. عملگر > مقدار درست را باز می‌گرداند اگر عملوند اول از عملوند دوم بزرگتر باشد. به همین ترتیب عملگرهای بعدی مقدار درست را باز می‌گردانند اگر عملوند اول آنها از عملوند دوم بزرگتر یا مساوی، کوچکتر، کوچکتر یا مساوی، مساوی، نامساوی باشد. عملگرهای رابطه‌ای نسبت به عملگرهای حسابی اولویت کمتری دارند. عملگرهای منطقی^{۱۱} && و && و فصل || و نقیض ! نیز در عبارات منطقی به کار می‌روند. یک عبارت منطقی عبارتی است که از متغیرهای منطقی و عملگرهای منطقی تشکیل شده و مقدار آن درست یا نادرست است. اولویت عملگرهای منطقی از عملگرهای حسابی و رابطه‌ای کمتر است. عملگر انتساب = برای مقداردهی یک متغیر به کار می‌رود و اولویت آن از عملگرهای حسابی، مقایسه‌ای، و منطقی کمتر است.

³ comments

⁴ slash

⁵ decimal

⁶ hexadecimal

⁷ American Standard Code for Information Interchange (ASCII)

⁸ backslash

⁹ arithmetic operators

¹⁰ relational operators

¹¹ logical operator

۸.۱ فاصله‌های خالی

توکن‌ها توسط یک فاصله خالی^{۱۲} و یا ترکیبی از فاصله‌های خالی از یکدیگر جدا می‌شوند. یک فاصله خالی شامل کاراکتر فاصله با کد اسکی ۳۲، کاراکتر خط جدید با کد اسکی ۱۰، و کاراکتر ستون جدید با کد اسکی ۹ می‌شود.

۹.۱ لیست توکن‌ها

در جدول زیر توکن‌های زبان PL به همراه نام توکن‌ها ذکر شده‌اند. خروجی تحلیل‌گر لغوی دنباله‌ای از نام توکن‌های یک برنامه ورودی است.

	Lexeme	Token		Lexeme	Token
۱	bool	T_Bool	۲۲	!=	T_ROp_NE
۲	break	T_Break	۲۳	==	T_ROp_E
۳	char	T_Char	۲۴	&&	T_LOp_AND
۴	continue	T_Continue	۲۵		T_LOp_OR
۵	else	T_Else	۲۶	!	T_LOp_NOT
۶	false	T_False	۲۷	=	T_Assign
۷	for	T_For	۲۸	(T_LP
۸	if	T_If	۲۹)	T_RP
۹	int	T_Int	۳۰	{	T_LC
۱۰	print	T_Print	۳۱	}	T_RC
۱۱	return	T_Return	۳۲	[T_LB
۱۲	true	T_True	۳۳]	T_RB
۱۳	+	T_AOp_PL	۳۴	;	T_Semicolon
۱۴	-	T_AOp_MN	۳۵	,	T_Comma
۱۵	*	T_AOp_ML	۳۶	variable or function names	T_Id
۱۶	/	T_AOp_DV	۳۷	decimal integers	T_Decimal
۱۷	%	T_AOp_RM	۳۸	hexadecimal integers	T_Hexadecimal
۱۸	<	T_ROp_L	۳۹	constant strings "[string]"	T_String
۱۹	>	T_ROp_G	۴۰	constant characters '[character]'	T_Character
۲۰	<=	T_ROp_LE	۴۱	// [string] \n	T_Comment
۲۱	>=	T_ROp_GE	۴۲	whitespace (newline, tab, and space characters)	T_Whitespace

¹² whitespace

۲ تحلیل‌گر نحوی

در این قسمت از پروژه می‌خواهیم یک تحلیل‌گر نحوی برای زبان PL طراحی و پیاده‌سازی کنیم. تحلیل‌گر نحوی دنباله‌ای از توکن‌ها را، که توسط تحلیل‌گر لغوی از یک فایل حاوی برنامه‌ای به زبان PL استخراج شده‌اند، دریافت می‌کند. در صورتی که برنامه فاقد خطای نحوی باشد، درخت نحوی را در خروجی چاپ می‌کند، در غیر اینصورت خطاهای برنامه را گزارش می‌کند. ساده‌ترین روش برای پیاده‌سازی تحلیل‌گر نحوی استفاده از تجزیه‌کننده پیش‌بینی‌کننده است. در صورتی که ورودی با خطا روبرو شد، کامپایلر باید با استفاده از روش‌هایی بازیابی خطا انجام دهد.

۱.۲ متغیرها و نوع‌های داده‌ای

سه نوع داده اصلی در این زبان وجود دارند. نوع داده‌های صحیح، کاراکتر، و منطقی که با `int` و `char` و `bool` نمایش داده می‌شوند. همچنین می‌توان توسط عملگر براکت باز و بسته آرایه تعریف کرد. یک آرایه متغیری است که دنباله‌ای محدود از خانه‌های حافظه را توسط یک نام معین توصیف می‌کند. اندازه آرایه با یک عدد صحیح ثابت تعیین می‌شود. برای تعریف یک متغیر، نوع متغیر و نام آن مشخص می‌شوند. متغیرهایی که هم‌نوع هستند می‌توانند توسط عملگر کاما از یکدیگر جدا شوند. همچنین متغیرها می‌توانند در هنگام تعریف با استفاده از عملگر انتساب مقداردهی شوند. در زیر چند نمونه تعریف متغیر نشان داده شده است.

```
۱ int x;
۲ bool b;
۳ int array[20];
۴ bool booleans[10];
۵ int i=5,j,k=10;
۶ bool t=true, f=false;
۷ char c1, c2='p';
۸ char string[50], s[12] = "PL Compiler";
```

۲.۲ دستورات شرطی

یک دستور شرطی برای بررسی یک شرط و اجرای پاره‌ای از دستورات در صورت برقراری شرط استفاده می‌شود. در یک دستور شرطی از کلمات کلیدی `if` و `else` استفاده می‌شود. در صورتی که شرط متعلق به `if` برقرار باشد، دستورات متعلق به آن که در بلوک متعلق به `if` با استفاده از علامت‌های آکولاد باز و بسته مشخص شده‌اند، اجرا می‌شوند. در صورتی که شرط برقرار نباشد، دستورات متعلق به بلوک `if` اجرا نمی‌شوند. در صورتی که بخواهیم هنگامی که شرط برقرار نیست، دستوراتی را اجرا کنیم می‌توانیم از کلیدواژه `else` و بلوک متعلق به آن استفاده کنیم. امکان تعریف `if` و `else` های تودرتو نیز وجود دارد، بدین معنی که شرطی بررسی شده، در صورت برقراری شرط دستوراتی اجرا می‌شوند و در صورتی که شرط برقرار نباشد و شرط دیگری برقرار باشد، پاره‌ای دیگر از دستورات اجرا می‌شوند. در زیر چند نمونه دستورات شرطی تعریف شده‌اند.

```
۱ if (x==1) {
۲     y = 6;
۳     z = 7;
۴ }
۵ if (y<=2 && z>7 || w!=9) {
۶     x=4+y;
۷ } else {
۸     x=3*z;
۹ }
۱۰ if (x>=5 || z<7 && !b) {
۱۱     b = true;
۱۲ } else if (x==3) {
۱۳     array[4] = 43;
۱۴ } else {
۱۵     c = 'q';
```

۳.۲ حلقه‌ها

یک حلقه برای تکرار تعدادی دستورات استفاده می‌شود. در یک حلقه از کلمه کلیدی `for` استفاده می‌شود. حلقه شامل سه بخش مقداردهی اولیه، شرط حلقه، و گام حلقه می‌شود. در مقداردهی اولیه یک متغیر مقداردهی می‌شود، در شرط حلقه یک شرط بررسی شده، و در صورت برقراری شرط، حلقه دوباره تکرار می‌شود، و در گام شرط مقدار متغیر حلقه تغییر می‌کند. این سه قسمت با عملگر نقطه ویرگول از یکدیگر جدا می‌شوند. درون حلقه می‌توان از دستورات `break` و `continue` استفاده کرد. اگر از دستور `break` در یک حلقه استفاده شود، تکرار حلقه متوقف می‌شود و کنترل برنامه از حلقه خارج می‌شود و از اگر دستور `continue` استفاده شود، دستورات بعد از آن اجرا نمی‌شوند و کنترل برنامه به ابتدای حلقه باز می‌گردد. در زیر چند نمونه حلقه تعریف شده‌اند.

```

۱ for (int i = 0 ; i<10; i=i+1) {
۲     array[i] = i;
۳ }
۴ for (x = 100 ; x>0; x=x-2) {
۵     sum = sum+x;
۶ }
۷ for (j = 0 ; ; j=j+1) {
۸     sum=sum+j;
۹     if (j>10) {
۱۰         break;
۱۱     }
۱۲ }
۱۳ for (int i = 0 ; i<100; i=i+1) {
۱۴     if (i % 7 == 0) {
۱۵         continue;
۱۶     }
۱۷     s = s+i;
۱۸ }
```

۴.۲ توابع

یک تابع با تعیین نوع مقدار بازگردانده شده توسط تابع، نام تابع، و پارامترهای ورودی تابع و نوع آنها تعریف می‌شود. یک تابع شامل دسته‌ای از دستورات است که در بلوک تابع بین دو علامت آکولاد باز و بسته نوشته می‌شوند. یک تابع نمی‌تواند یک آرایه بازگرداند. دستور `return` یک مقدار را از تابع باز می‌گرداند. در زیر چند نمونه تابع تعریف شده‌اند.

```

۱ int multiply(int x, int y) {
۲     int z = x*y;
۳     return z;
۴ }
۵ int multiply(int x, bool b, char c) {
۶     int r;
۷     if (b && c == 'y') {
۸         r = x*10;
۹     } else {
۱۰         r = x*5;
۱۱     }
```

```
۱۲     return r;
۱۳ }
```

۵.۲ دستور چاپ

از دستور چاپ برای چاپ مقادیر عددی، کاراکترها و رشته‌ها استفاده می‌شود. این دستور یک رشته، و صفر یا تعدادی متغیر را دریافت کرده، مقادیر متغیرها را در درون رشته چاپ می‌کند. رشته دریافت شده، رشته قالب‌بندی نامیده می‌شود. در درون رشته قالب‌بندی تعدادی تعیین‌کننده نوع استفاده شده‌اند که برای چاپ اعداد و کاراکترها استفاده می‌شوند. تعیین‌کننده‌های نوع می‌توانند %d برای اعداد صحیح یا %c برای کاراکترها یا %s برای آرایه‌ای رشته‌ها (آرایه‌ای از کاراکترها) باشند. در زیر چند نمونه دستور چاپ نشان داده شده است.

```
۱ print("PL Programming Language");
۲ print("%d", i*j);
۳ print("The value of x is %d and the value of c1 is %c .", x, c1);
۴ print("%s\n", string);
```

۳ تحلیل‌گر معنایی

در این قسمت از پروژه می‌خواهیم یک تحلیل‌گر معنایی برای زبان PL طراحی و پیاده‌سازی کنیم. تحلیل‌گر معنایی درخت تجزیه را که توسط تحلیل‌گر نحوی تولید شده دریافت می‌کند و با پیمایش درخت تجزیه خطاهای معنایی را گزارش می‌کند. یک تحلیل‌گر معنایی همچنین می‌تواند از درخت تجزیه¹³ درخت نحوی¹⁴ تولید کند. به عنوان بخش امتیازی پروژه می‌توانید از درخت تجزیه درخت نحوی تولید کرده و درخت نحوی تولید شده را اجرا کنید. قوانین زیر باید در تحلیل‌گر معنایی رعایت شوند.

- یک شناسه باید قبل از استفاده تعریف شده باشد. یک شناسه می‌تواند نام یک متغیر یا یک تابع باشد.
- هر شناسه دارای یک حوزه تعریف¹⁵ است. یک حوزه تعریف توسط یک بلوک مشخص می‌شود و یک بلوک با آکولاد باز آغاز می‌شود و با آکولاد بسته به اتمام می‌رسد. حوزه‌های تعریف می‌توانند تو در تو باشند و در یک حوزه تعریف می‌توان یک حوزه تعریف دیگر داشت. بنابراین حوزه‌های تعریف دارای یک ساختار سلسله‌مراتبی¹⁶ هستند. دو متغیر همنام می‌توانند در دو حوزه تعریف وجود داشته باشند، اما در یک حوزه تعریف نمی‌توان دو متغیر همنام تعریف کرد. همچنین دو تابع نمی‌توانند همنام باشند.
- نوع عملوندهای عملگرهای محاسباتی `int` و نوع عملوندهای عملگرهای منطقی باید `bool` باشد.
- اندیس یک آرایه باید مقداری بزرگتر از صفر داشته باشند و نوع آن الزاماً `int` باشد.
- نوع عبارتی که در شرط دستور `if` قرار می‌گیرد باید `bool` باشد.
- برنامه دارای یک تابع به نام `main` است که پارامتر ورودی ندارد و نوع خروجی آن `int` است.
- در فراخوانی توابع باید نوع و تعداد آرگومان‌ها با نوع و تعداد پارامترها در تعریف تابع همخوانی داشته باشند.
- در عبارت‌های انتساب مقدار، نوع سمت چپ و نوع سمت راست عملگر انتساب باید یکسان باشند. در صورتی که در سمت راست عملگر انتساب فراخوانی تابع وجود داشته باشد، نوع تابع و نوع متغیر سمت چپ عملگر انتساب باید یکسان باشند.
- نوع تابع و نوع عبارتی که بعد از کلمه کلیدی `return` در تابع قرار گرفته باید یکسان باشند.

¹³ parse tree

¹⁴ syntax tree

¹⁵ scope

¹⁶ hierarchical structure