

Syntaktická analýza založená na několika gramatikách

Ondřej Koumar
xkouma02@stud.fit.vutbr.cz

26. března 2024

Motivace

Syntaktická analýza a obecně překladače často zůstávají v praxi poněkud vzdálené od teoretických konceptů. I přes bohatou teorii formálních jazyků, konkrétně například gramatik a gramatických systémů, se v praxi často při syntaktické analýze upřednostňují jednoduché postupy.

Má bakalářská práce se snaží spojit teoretický svět s praktickým využitím. Pro syntaktickou analýzu využívám gramatický systém s několika gramatikami, ale implementace zůstává jednoduchá a srozumitelná pomocí dvou zásobníkových automatů, LL tabulky a precedenční analýzy výrazů. Výsledkem této analýzy bude abstraktní syntaktický strom.

Implementovaný jazyk

Nejdříve se podíváme na jazyk, který je navrhnutými gramatikami popsán. Zatím nemá jméno, nazvěme jej pracovním KOUBP. Je založen na jazyce IFJ22, který je založen na jazyce PHP. Jazyk KOUBP je staticky typovaný jazyk, tedy datový typ všech proměnných a konstant je znám už za překladu.

Konstrukce jazyka

Jazyk KOUBP má čtyři základní datové typy, a to `int`, `float`, `bool`, `string`. Obsahuje konstrukce, které byste očekávali u každého imperativního programovacího jazyka – `while` a `for` cykly, příkazy `if-elseif-else` a `return`, jednoduché aritmetické výrazy a operátory typu sčítání a odčítání, relační operátory, ale například i konkatenci řetězců či logické operátory `&&` a `||`. Kód může být libovolně rozčleněn do funkcí, které mohou a nemusí mít návratový typ.

Oddělovač příkazů

Všechny příkazy musí být ukončeny středníkem, který se používá jako oddělovač příkazů. Výjimkou jsou bloky kódu uzavřené ve složených závorkách. Za takovým blokem středník následovat nemusí, ale uvnitř bloku kódu je žádoucí každý příkaz oddělit.

Blok kódu může být použit samostatně, nicméně typicky jej vidíme za konstrukcemi `while`, `for`, `if`, případně při definici funkce. Příklad bloku kódu lze vidět na obrázku ??.

```
prikaz;  
{  
    prikaz;  
    prikaz;  
}  
prikaz;
```

Obrázek 1: Příklad použití bloku kódu

Středník lze použít nejen jako oddělovač příkazů, může být také chápán jako samostatný příkaz. Jeho význam se dá přirovnat k instrukci `NOP` ve strojovém jazyce/asmbleru.

Gramatický systém

Gramatický systém, který popisuje jazyk KOUBP, je postaven na *CD (cooperating distributed)* gramatických systémech. Tento systém se skládá ze šesti komponent, z nichž každá reprezentuje část navrhovaného jazyka tak, aby pravidla byla systematicky rozdělena. Motivace k tomuto přístupu je podobná členění zdrojových kódů do logických celků. Zároveň je ale důležité, aby se moc nezvýšila náročnost implementace a vyplatilo se rozdělení do více gramatik udělat.

Následují gramatiky/komponenty, které dohromady tvoří gramatický systém jazyka KOUBP. V nadpise vizte ucelenou část jazyka, kterou daná pravidla tvoří.

Struktura programu

```
<program> -> <statement>2 <statementList>1
<program> -> <functionDef>2 <statementList>1

<statementList> -> <statement>2 <statementList>1
<statementList> -> <functionDef>3 <statementList>1
<statementList> -> ε
```

Příkazy a běžné konstrukce

```
<statement> -> if ( <expression>4 ) <codeBlock>5 <if2>2
<statement> -> while ( <expression>4 ) <codeBlock>5
<statement> -> for ( <expression>4 ; <expression>4 ; <expression>4 ) <codeBlock>5
<statement> -> <expression>4 ;
<statement> -> return <returnExp>2 ;
<statement> -> <codeBlock>5
<statement> -> ;

<if2> -> elseif ( <expression>4 ) <codeBlock>5 <if2>2
<if2> -> else <codeBlock>5
<if2> -> ε

<returnExp> -> <expression>4
<returnExp> -> ε
```

Definice funkce

```
<functionDef> -> function funcName ( <params>3 ) : <funcType>3 <codeBlock>5

<params> -> <type>6 variable <params2>3

<params2> -> , <type>6 variable <params2>3
<params2> -> ε

<funcType> -> <type>6
<funcType> -> ε
```

Výrazy

```
<expression> -> <expression>4 [+ - * / . && || == != > < >= <=] <expression>4
<expression> -> [! -] <expression>4
<expression> -> variable
<expression> -> constant
<expression> -> funcName ( <args>4 )
<expression> -> ( <expression>4 )

<args> -> <expression>4 <args2>4
<args> -> ε
```

```
<args2> -> , <expression>4 <args2>4
<args2> -> ε
```

Blok kódu

Všimněme si, že neterminály *codeBlock*₅ a *statement*₂ se mohou navzájem derivovat donekonečna a potenciálně by mohl nastat *deadlock*. Po teoretické stránce se s tímto faktem počítá a gramatiky se nebudou nijak přizpůsobovat, tato problematika bude vyřešena implementačně, zásobníkový automat se s pomocí LL tabulky bude chovat deterministicky.

```
<codeBlock> -> { <statements>5 }
<codeBlock> -> <statement>2

<statements> -> <statement>2 <statements>5
<statements> -> ε
```

Datové typy

```
<type> -> int
<type> -> float
<type> -> string
<type> -> bool
```

Implementace

Implementujeme syntaktický analyzátor v jazyce C++ s dvěma zásobníkovými automaty – jeden pro analýzu s LL tabulkou a druhý pro precedenční analýzu výrazů. Precedenční analýza výrazů bude prováděna obvyklým způsobem, tedy redukcí terminálů a případně pravidel na zásobníku.

Analýza s LL tabulkou bude upravena na používání uspořádaných dvojic (*číslo gramatiky*, *číslo pravidla*) namísto pravidel v položkách tabulky. Číslo gramatiky bude předáno tovární metodě, která vytvoří instanci jedné z gramatik dědicích od abstraktní třídy **Grammar**. Tato třída obsahuje polymorfní metodu **Expand(unsigned ruleNumber)**, která vrací reverzovaný seznam terminálů a neterminálů reprezentující pravou stranu daného pravidla. Reverzované pravidlo se vloží na zásobník a analýza bude pokračovat dalším cyklem.

V každém kroku syntaktické analýzy bude generován uzel abstraktního syntaktického stromu.

Možná rozšíření

Syntaktický analyzátor a jeho datové struktury jsou koncipovány s ohledem na jednoduchou rozšiřitelnost projektu. Aktuálně pracujeme s pevně danou posloupností terminálů/tokenů, kterou analyzátor přijímá jako seznam. Nicméně, máme v plánu budoucí rozšíření projektu o lexikální analyzátor, což umožní psát skutečné programy v jazyce KOUBP.

Týká se to i jednoduchých sémantických kontrol, které jsou v každém kvalitním syntaktickém analyzátoru žádoucí.