

NET LL(1) grammar parser

(Započtový program C#(zima) a Advanced .NET (leto). Pavel Vigilev)

Uživatelska dokumentace.

1.C#-ova část

Aby mohl využít framework, uživatel musí přidat odkaz na assembly frameworku (LL1_Parser.dll). Třída *LL1GrammarParser* (namespace *LL1_Parser*) je hlavní třídou, s kterou uživatel pracuje.

Obsah namespace *LL1_Parser*:

1. Namespace obsahuje definici rozhraní *IToken*, *ITokenWithValue* a *ILexer*.
Uživatelské datové typy tokenů a lexeru musí se dědit z těchto rozhraní.
2. Instance třídy *LL1GrammarParser<T, UToken>* je recursive descent parser pro LL(1) gramatiky.
Generic-parameter *T* definuje typ hodnoty, kterou framework parsuje; *UToken* - typ datové struktury popisující uživatelské tokeny.

Třída *LL1GrammarParser* má následující konstruktor:

LL1GrammarParser(string grammar, ILexer<UToken> lexer, IDictionary<string, UToken> NameTerminalTable, Assembly[] assemblies):

- *grammar* - obsah souboru s gramatikou
- *lexer* - objekt uživatelské třídy lexeru, která tvoří ze vstupního řetězce, posloupnost tokenů typu *UToken*.
- *NameTerminalTable* - tabulka jmen terminalů, aby framework věděl, který symbol je v popisu gramatiky terminal.
- *assemblies* - pole assembly, kde framework hledá typy a metody s kterými pracuje v rámci parsování

Metody *T Parse(string str)* a *bool TryParse(string str, out T result)* mají stejnou semantiku, jakou mají *int.Parse* a *int.TryParse* resp.

Rozhraní *IToken*, *ITokenWithValue* a *ILexer*:

```
interface IToken
{
    bool IsCompatible(IToken other);
}

interface ITokenWithValue<out T> : IToken where T : class
{
    T Value { get; }
}

interface ILexer<T> where T : IToken
{
    IList<T> Tokenize(string str);
}
```

2.Soubor s gramatikou

Gramatika se popisuje nasledujicim zpusobem:

```
nonterminal_name : symbol1 symbol2 ... symboln { expression; expression; ... expression; }
```

kde *symbol* je jmeno terminalu nebo neterminalu.

Nove neterminaly se definuji automaticky. Pokud na prave strane pravidla vyskytuje symbol "symbol" a neexistuje terminal s takovem jmenem, registruje se novy neterminal s timto jmenem. Jmeno terminalu je jednoho z dvou typů:

Vyruzy v '{' a '}'

Složite zavorky naproti každomu pravidlu obsahují neprazdnou posloupnost vyrazu. Každý vyraz podporuje zavolani instančni metody, staticke metody, konstruktoru, využití literalu int, double, string a číslo vysledku parsingu (*\$uint*).

int, *string*, *double* lze využít jako podobne literaly v C#.

\$i odkazuje na rozparsovaný symbol v pravidle, kde vyskytuje.

Zavolani metod mají nasledující syntax. Jako parametr lze využít nějaký jiný libovolný vyraz.

```
// Instance-method invokation
FullMethodName(instance, arg1, arg2, ..., argn);

// static-method invokation
static FullMethodName(arg1, ... argn);

// constructor invokation
new FullTypeName(arg1, ..., argn);
```

FullMethodName (podobně i *FullMethodType*) je úplně jmeno metody včetně všech jmenovych prostoru a tříd, kde se nachází tato metoda.

Příklad:

grammar file:

```
Expr : var { new Lambda.LVar($0); }
Expr : '(' Expr1 ')' { $1; }
Expr1 : Expr Expr { new Lambda.LApplication($0, $1); }
Expr1 : '\' var '.' Expr { new Lambda.LAbstraction(new Lambda.LVar($1), $3); }
```

'(', ')', '\', '.', var jsou terminaly a musí být definovány v *NameTerminalTable* (v příkladu je uvnitř lexeru - *lexer.NameTokenTBL*).

data structure:

```
public abstract class LExprBasic{}

public class LVar :LExprBasic
{
    public string Value { get; }
    public LVar(string val)
    {
        Value = val ?? throw new ArgumentNullException();
    }
    public override string ToString()
    {
        return Value;
    }
}

public class LApplication : LExprBasic
```

```

{
    public LExprBasic What { get; }
    public LExprBasic To { get; }

    public LApplication(LExprBasic w, LExprBasic t)
    {
        What = w ?? throw new ArgumentNullException();
        To = t ?? throw new ArgumentNullException();
    }
    public override string ToString()
    {
        return $"({What.ToString()} {To.ToString()})";
    }
}

public class LAbstraction : LExprBasic
{
    public LVar Variable { get; }
    public LExprBasic Expression { get; }
    public LAbstraction(LVar var, LExprBasic expr)
    {
        Variable = var ?? throw new ArgumentNullException();
        Expression = expr ?? throw new ArgumentNullException();
    }
    public override string ToString()
    {
        return $"(\\{Variable}.{Expression.ToString()})";
    }
}

```

parser-creating:

```

var parser = new LL1GrammarParser<LExprBasic, AbstractLToken>(grammar, lexer, lexer.NameTokenTBL, new
System.Reflection.Assembly[] { typeof(LExprBasic).Assembly });

```

Kde *AbstractLToken* je uživatelská implementace rozhraní *IToken*; *lexer* je instance nějaké třídy *Lexer*, která je implementace rozhraní *ILexer<AbstractLToken>*.

Tedy objekt parser lze využít na parsování lambda-výrazu.

3.Chyby a výjimky

Pokud v průběhu parsování došlo k typové chybě (např. jako argument metody byla předána hodnota nějakého špatného typu), vzniká výjimka **ParsingTypeErrorException**

Pokud nenalezena metoda, která by vyhovovala podmínkám (typy a pořadí parametrů)

ArgumentErrorException

Chyba hledání třídy - **TypeNotFoundException**, chyba hledání metody v třídě - **MethodNotFoundException**.

V případě nalezení dvou úspěšných derivací pro jeden řetězec, vzniká výjimka **AmbiguousGrammarException**.

V případě nalezení levo rekurze - **LeftRecursionGrammarException**.

Výjimky vyšší se dědí z **ParsingErrorException**, která se dědí z **System.Exception**.

Pokud vstupní řetězec nedá se rozparsovat gramatikou, tak vzniká **FormatException** (pro kompatibilitu s *int.Parse*, *uint.Parse*, atd. metodami).