# The mechanics of GF

Krasimir Angelov

University of Gothenburg

August 22, 2013

- Well known grammar formalism (Seki at al., 1991)

- Natural extension of CFG that produces tuples of strings instead of simple strings

- It is trivial to implement classical context-sensitive languages - $\{a^n b^n c^n | n \geq 0\}$:

The parser uses a language which is a subset of GF.

- The linearization types are flat tuples of strings:

$$\textbf{lincat } C = Str * Str * \ldots * Str;$$

- The linearizations are simple concatenations:

$$\textbf{lin } f \; x \; y = < x.p1, x.p2 ++ y.p3 >;$$

- No operations are allowed
- No variants are allowed
- No parameters and tables
- No pattern matching
- No gluing is allowed (i.e. $++$ but not $+$)

**cat** $N, S$

**fun** $z : N$

$\quad s : N \rightarrow N$

$\quad c : N \rightarrow S$

**lincat** $N = Str * Str * Str$

$\quad S = Str$

**lin** $z \quad = <\text{""}, \text{""}, \text{""}>$

$\quad s\ x = <\text{"a"} + x.p1, \text{"b"} + x.p2, \text{"c"} + x.p3>$

$\quad c\ x = x.p1 + x.p2 + x.p3$

- Operations elimination

- Variants elimination

- Parameter types elimination

- Linearization rules transformations

- Common subexpressions optimization

# Operations elimination

The operations are **NONRECURSIVE** functions. They are evaluated at compile time. *(macroses)*

## GF

```
oper mkN noun = case noun of {
    _ + "s" ⇒ < noun, noun + "es" >;
    _       ⇒ < noun, noun + "s" >
    };
  lin apple_N = mkN "apple";
      plus_N = mkN "plus";
```

## GF Core

```
lin apple_N = < "apple", "apples" >;
    plus_N = < "plus", "pluses" >;
```

*Note: the pattern matching in mkN was eliminated*

The variants are just expanded:

**GF**

$$\textbf{lin } girl\_N = mkN \text{ (}"tjej" \mid "flicka" \text{)};$$

**GF Core**

$$\textbf{lin } girl\_N_1 = mkN \text{ }"tjej";$$
$$girl\_N_2 = mkN \text{ }"flicka";$$

# Parameter Types Elimination

$$\textbf{lincat } NP = \{s : Case \Rightarrow Str; g : Gender; n : Number; p : Person\}$$

$$\textbf{param } Case = Nom|Acc|Dat;$$

$$Gender = Masc|Fem|Neutr;$$

$$Number = Sg|Pl;$$

$$Person = P1|P2|P3;$$

A value of type $Case \Rightarrow Str$ looks like:

$$\textbf{table } \{Nom \Rightarrow s_1; \ Acc \Rightarrow s_2; \ Dat \Rightarrow s_3\}$$

We could replace it with the tuple:

$$< s_1, s_2, s_3 >$$

Then in general type like $A \Rightarrow Str$ is equivalent to:

$$\underbrace{Str * Str * \ldots * Str}_{n \text{ times}}$$

where $n$ is the number of values in the parameter type $A$.

# Parameter Fields Elimination

## GF

$$\text{lincat } NP = \{s : \ldots; g : Gender; n : Number; p : Person\}$$

## GF Core

$$\text{lincat } NP_1 = Str * Str * Str; \quad - Masc; Sg, P1$$
$$NP_2 = Str * Str * Str; \quad - Masc; Sg, P2$$
$$NP_3 = Str * Str * Str; \quad - Masc; Sg, P3$$
$$NP_4 = Str * Str * Str; \quad - Masc; Pl, P1$$
$$\vdots$$
$$NP_{18} = Str * Str * Str; \quad - Neutr; Pl, P3$$

# Linearization Rules Transformation

## GF

$$\textbf{fun } AdjCN : AP \rightarrow CN \rightarrow CN;$$
$$\textbf{lin } AdjCN \ ap \ cn = \{$$
$$s = ap.s!cn.g \mathbin{+\!\!+} cn.s;$$
$$g = cn.g$$
$$\};$$

## GF Core

$$\textbf{fun } AdjCN_1 : AP \rightarrow CN_1 \rightarrow CN_1; \qquad -Masc$$
$$\textbf{lin } AdjCN_1 \ ap \ cn = \ <ap.p1 \mathbin{+\!\!+} cn.p1>$$

$$\textbf{fun } AdjCN_2 : AP \rightarrow CN_2 \rightarrow CN_2; \qquad -Fem$$
$$\textbf{lin } AdjCN_2 \ ap \ cn = \ <ap.p2 \mathbin{+\!\!+} cn.p1>$$

$$\textbf{fun } AdjCN_3 : AP \rightarrow CN_3 \rightarrow CN_3; \qquad -Neutr$$
$$\textbf{lin } AdjCN_3 \ ap \ cn = \ <ap.p3 \mathbin{+\!\!+} cn.p1>$$

# No pattern matching

## Allowed

$$\textbf{oper } mkN \; noun = \textbf{case } noun \textbf{ of } \{$$
$$\_ + "\text{s}" \Rightarrow \; < noun, noun + "\text{es}" >;$$
$$\_ \qquad \Rightarrow \; < noun, noun + "\text{s}" >$$
$$\};$$

## Not Allowed

$$\textbf{lin } DetCN \; det \; cn = \textbf{case } det.s \textbf{ of } \{$$
$$"" \Rightarrow \ldots$$
$$\_ \Rightarrow \ldots$$
$$\}$$

*Hint: use parameter which says whether the string is empty*

# No gluing

**lin** *DetCN det cn* = **case** *det.spec* **of** {

  . . .

    *Indefinite* ⇒ **case** *cn.g* **of** {*Utr* ⇒ "en"; *Neutr* ⇒ "ett"} ++ *cn.s*

  }

**lin** *DetCN det cn* = **case** *det.spec* **of** {

    *Definite* ⇒ *cn.s* + **case** *cn.g* **of** {*Utr* ⇒ "en"; *Neutr* ⇒ "et"};

  . . .

  }

*Hint: for agglutinative languages (Turkish, Finnish, Estonian, Hungarian, ...) use custom lexer*

## Agglutinatination

- Some languages have pottentially infinite set of words:

  Turkish:
  anlamiyorum = anla(root) -mi(negation) -yor(continuous) -um(first person)
  I don't understand

- The grammar could be based on roots and suffixes instead of on words:

  "anla" ++ "&+" ++ "mi" ++ "&+" ++ "yor" ++ "&+" ++ "um"

- The lexer/unlexer are responsible to produce the real words

- GF $\Rightarrow$ (GF Core $\equiv$ PMCFG)

- Linearization is overload resolution

- Parsing is search