

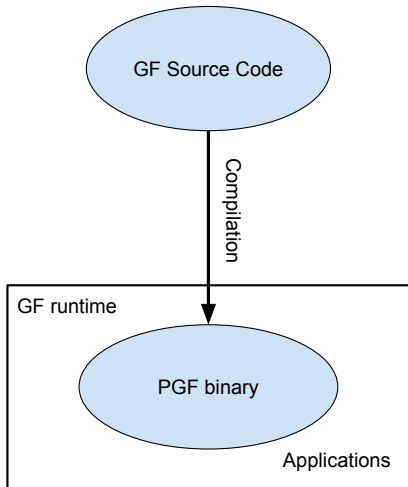
# The mechanics of GF

Krasimir Angelov

University of Gothenburg

August 21, 2017

# Compilation



1 PMCFG<sup>+</sup>

2 Embedded Grammars and Applications

3 Conclusion

The runtime understands only a subset of GF:

- The linearization types are flat tuples of strings:

$$\mathbf{lincat} \ C = \mathit{Str} * \mathit{Str} * \dots * \mathit{Str};$$

- The tuples contain simple concatenations:

$$\mathbf{lin} \ f \ x \ y = \langle x.p1, x.p2 ++ y.p3 \rangle;$$

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

# 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*

# Variants elimination

The variants are just expanded:

GF

```
lin girl_N = mkN ("tjej" | "flicka");
```

GF Core

```
lin girl_N1 = mkN "tjej";  
      girl_N2 = mkN "flicka";
```

# Parameter Types Elimination

**lincat**  $CN = \{s : Case \Rightarrow Str; g : Gender\}$

**param**  $Case = Nom|Acc|Dat;$

$Gender = Masc|Fem|Neutr;$

# Table Types Elimination

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

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

We could replace it with the tuple:

$\langle s_1, s_2, s_3 \rangle$

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

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

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



## GF

**lincat** *CN* = {*s* : ... ; *g* : *Gender*}

## GF Core

**lincat** *NP*<sub>1</sub> = *Str* \* *Str* \* *Str*;    – *Masc*  
*NP*<sub>2</sub> = *Str* \* *Str* \* *Str*;    – *Fem*  
*NP*<sub>3</sub> = *Str* \* *Str* \* *Str*;    – *Neutr*

# Linearization Rules Transformation

## GF

```
fun AdjCN : AP → CN → CN;  
lin AdjCN ap cn = {  
  s = ap.s!cn.g ++ cn.s;  
  g = cn.g  
};
```

## GF Core

```
fun AdjCN1 : AP → CN1 → CN1;      -Masc  
lin AdjCN1 ap cn = < ap.p1 ++ cn.p1 >  
  
fun AdjCN2 : AP → CN2 → CN2;      -Fem  
lin AdjCN2 ap cn = < ap.p2 ++ cn.p1 >  
  
fun AdjCN3 : AP → CN3 → CN3;      -Neutr  
lin AdjCN3 ap cn = < ap.p3 ++ cn.p1 >
```

# No pattern matching

## Allowed

```
oper mkN noun = case noun of {  
  _ + "s" ⇒ < noun, noun + "es" >;  
  _       ⇒ < noun, noun + "s" >  
};
```

## Not Allowed

```
lin DetCN det cn = case det.s of {  
  "" ⇒ ...  
  _  ⇒ ...  
}
```

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

# No gluing

## Allowed

```
lin DetCN det cn = case det.spec of {  
  ...  
  Indefinite ⇒ case cn.g of {Utr ⇒ "en"; Neutr ⇒ "ett"} ++ cn.s  
}
```

## Not Allowed

```
lin DetCN det cn = case det.spec of {  
  Definite ⇒ cn.s + case cn.g of {Utr ⇒ "en"; Neutr ⇒ "et"};  
  ...  
}
```

## Extension: BIND

- Agglutinative languages have potentially infinite set of words:

Turkish:

anlamıyorum = anla(root) -mı(negation) -yor(continuous) -um(first person)  
I don't understand

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

"anla" ++ *BIND* ++ "mı" ++ *BIND* ++ "yor" ++ *BIND* ++ "um"

- The parser uses BIND to do the right segmentation

## Extension: *SOFT\_BIND* and *SOFT\_SPACE*

- *SOFT\_BIND*

|            |   |
|------------|---|
| Expression | "Hi" ++ "John" ++ <i>SOFT_BIND</i> ++ ";" |
| Accepts    | "Hi John," "Hi John␣," , ...              |
| Generates  | "Hi John,"                                |

- *SOFT\_SPACE*

|            |                                    |
|------------|------------------------------------|
| Expression | "10" ++ <i>SOFT_SPACE</i> ++ "000" |
| Accepts    | "10 000" "10000" , ...             |
| Generates  | "10 000"                           |

- Capitalization in German

*CAPIT* + "aktions" + *BIND* + "plan" ⇒ "Aktionsplan"

- Acronyms

*ALL\_CAPIT* + "ibm" ⇒ "IBM"

## Extension: Case-insensitive Parsing

### Enable case-insensitive parsing

```
flags case_sensitive = "off"
```

- All words in the grammar must be in lower-case
- Capitalization is achieved only with *CAPIT* and *ALL\_CAPIT*
- The switching off *case\_sensitive* tells the parser to lower-case the input



## Lexical gaps

```
lin foo_N = table {  
    Sg ⇒ "foo";  
    Pl ⇒ nonExist;  
};
```

### *nonExist*

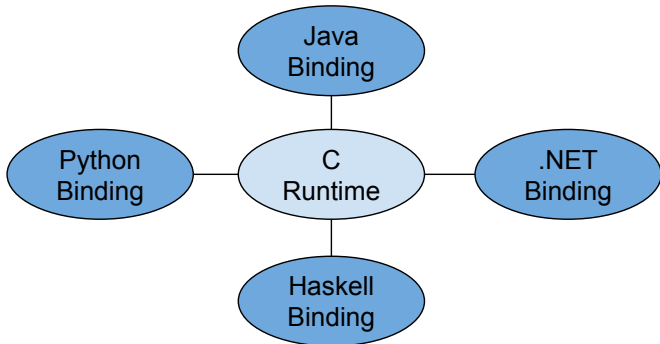
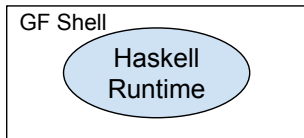
- generates exception in linearization
- blocks parsing with the corresponding function

1 PMCFG<sup>+</sup>

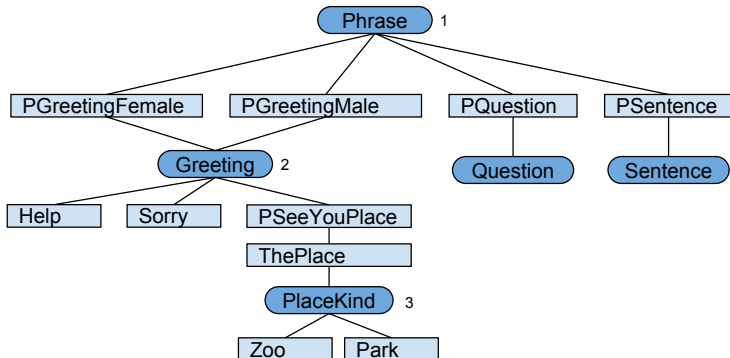
2 Embedded Grammars and Applications

3 Conclusion

# Types of Runtime



# Choice Graph



1 PMCFG<sup>+</sup>

2 Embedded Grammars and Applications

3 Conclusion

- $GF \Rightarrow (GF \text{ Core} \equiv \text{PMCFG})$
- Embedded Grammars:  
Haskell, Python, Java, C, .NET
- User Interfaces