

Theory of GF and Parsing, Hints for Efficient Grammars

Krasimir Angelov

Chalmers University of Technology

August 21, 2009

1 Introduction

2 GF Core

3 Optimizations

4 Debugging

5 Conclusion

The GF language is:

- domain specific - Grammars first of all
- declarative - What vs How

The GF system is:

- optimizing compiler and interpreter
- not as smart as you

Although the language is declarative, the compiler needs some help to produce efficient grammars

Parsing vs Linearization

There are two major operations that someone could do with a grammar

- Linearization
 - **Efficient** - mapping from tree to string
- Parsing
 - **Search Problem** - find the tree(s) that produce a given string

Speaking about the efficiency I mean efficient parsing and compact grammar

Inside the Black Box

- The GF language is too complex to be used directly for efficient parsing
- The compiler transforms it into simpler **GF Core** language
- The efficiency of the grammar depends on the GF Core
- Small Core = Efficient Grammar

$GF \Rightarrow GF \text{ Core}$

$GF \text{ Core} \equiv PMCFG$

Parallel Multiple Context-Free Grammar (PMCFG)

- Well known grammar formalism (Seki et al., 1991)
- Natural extension of CFG that produces tuples of strings instead of simple strings
- It is trivial to implement classical context-sensitive languages - $\{a^n b^n c^n | n > 0\}$:

fun z = < "", "", "" >

s x = < "a" ++ x.p1, "b" ++ x.p2, "c" ++ x.p3 >

c x = < x.p1 ++ x.p2 ++ x.p3 >

Mildly Context-Sensitive Languages and GF

Joshi Aravind. 1991. Tree Adjoining Grammars: How much context-sensitivity is required to provide reasonable structural descriptions?

- The language must be parsable in polynomial time.
- The language must have constant growth.
- The language should admit limited cross-serial dependencies.

Example language: $\{a^n b^n c^n | n > 0\}$ and all MCFG

More than Mildly Context-Sensitive Languages

The exponential language $\{a^{2^n} \mid n > 0\}$:

fun $z = \langle "a" \rangle$

$s\ x = \langle x.p1 ++ x.p1 \rangle$

is not Mildly Context-Sensitive (**reduplication**).

Four non-MCS Natural Languages

- Mandarin Chinese numeral names
- Mandarin Chinese yes/no questions
- Old Georgian case system
- Lindenmayer system

Mandarin Chinese numeral names

- Concatenation of ten thousands

	wan	10 000
wan wan	yi	100 000 000
wan yi	zhao	1 000 000 000 000

- Composite numbers

wu zhao zhao wu zhao
five trillion trillion five trillion

5 000 000 000 005 000 000 000 000

- Formally

$$\{\text{wu zhao}^{k_1} \text{ wu zhao}^{k_2} \dots \text{wu zhao}^{k_n} \mid k_1 > k_2 > \dots > k_n\}$$

Mandarin Chinese yes/no questions

Zhangsan like play basketball not like play basketball
Zhangsan ai da lanqiu bu ai da lanqiu

Does Zhangsan like to play basketball?

Zhangsan like play basketball not like play volleyball
Zhangsan ai da lanqiu bu ai da paiqiu

Zhangsan likes to play basketball, not to play volleyball

Old Georgian case system

- Example:

tkuenda micemul ars cnob-ad saidumlo-j
to you given is knowing-Adv mistery-Nom

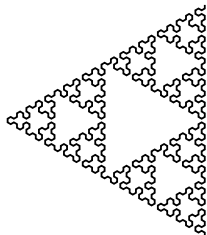
igi sasupevel-isa m-is ymrt-isa-isa-j
Art=Nom kingdom-Gen Art-Gen god-Gen-Gen-Nom

Onto You it is given to know the mistery of the kingdom of God
Mark 4:11

- Formal

N_1 -Nom N_2 -Gen-Nom N_3 -Gen²-Nom ... N_k -Gen ^{$k-1$} -Nom

- Mathematical objects
- The structure of some plants
- The growth of some crystals
- Symmetry in music



$$A \rightarrow a$$

$$B \rightarrow b$$

$$A \rightarrow BRARB$$

$$B \rightarrow ALBLA$$

The reduplication is a norm!

The current status of parsing with GF

Angelov. 2009. Incremental Parsing with PMCFG

Features:

- Very Efficient (**polynomial** - close to **linear**)
- Supports PMCFG for free
- PMCFG allows more compact grammars
- It is incremental !!!

Things to consider:

- the GF \Rightarrow GF Core conversion is often exponential
- The grammar should be carefully written to avoid combinatorial explosion
- In practice careful means linguistically motivated

The Rules of The Game

INITIAL PREDICT

$$\frac{S \rightarrow f[\vec{B}]}{[{}^0_0 S \rightarrow f[\vec{B}]; 1 : \bullet \alpha]} \quad S - \text{start category, } \alpha = \text{rhs}(f, 1)$$

PREDICT

$$\frac{B_d \rightarrow g[\vec{C}] \quad [{}^k_j A \rightarrow f[\vec{B}]; l : \alpha \bullet \langle d; r \rangle \beta]}{[{}^k_k B_d \rightarrow g[\vec{C}]; r : \bullet \gamma]} \quad \gamma = \text{rhs}(g, r)$$

SCAN

$$\frac{[{}^k_j A \rightarrow f[\vec{B}]; l : \alpha \bullet s \beta]}{[{}^{k+1}_j A \rightarrow f[\vec{B}]; l : \alpha s \bullet \beta]} \quad s = w_{k+1}$$

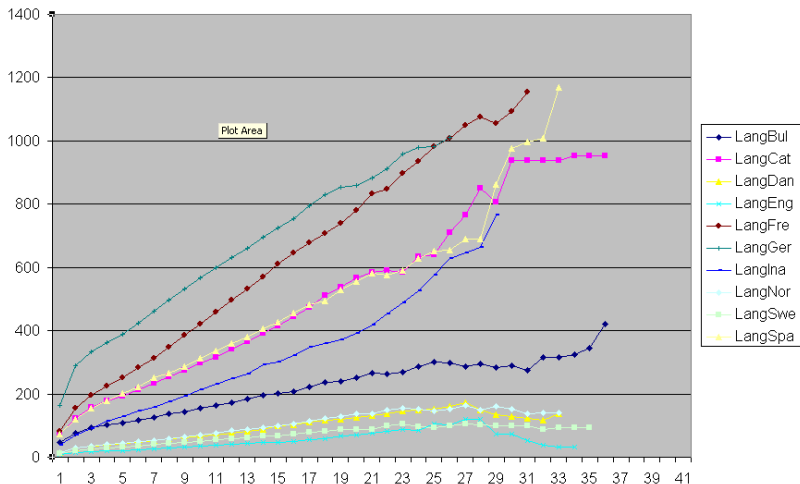
COMPLETE

$$\frac{[{}^k_j A \rightarrow f[\vec{B}]; l : \alpha \bullet]}{N \rightarrow f[\vec{B}] \quad [{}^k_j A; l; N]} \quad N = (A, l, j, k)$$

COMBINE

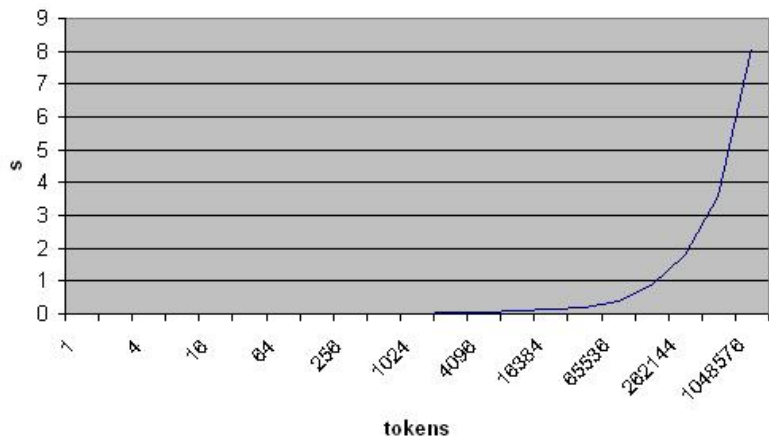
$$\frac{[{}^u_j A \rightarrow f[\vec{B}]; l : \alpha \bullet \langle d; r \rangle \beta] \quad [{}^k_u B_d; r; N]}{[{}^k_j A \rightarrow f[\vec{B}\{d := N\}]; l : \alpha \bullet \langle d; r \rangle \bullet \beta]}$$

Parsing with the resource library



Note: Much faster with application grammars

Parsing of $\{a^{2^n} \mid n > 0\} - O(n \log n)$



1 Introduction

2 GF Core

3 Optimizations

4 Debugging

5 Conclusion

The parser uses language which is simplified version of GF.

- The linearization types are flat tuples of strings:

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

- The linearizations are 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
- 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

Hints for Operations

Since the operations are computed at compile time this doesn't affect the runtime efficiency. However they affect the compilation speed (*slightly*).

Variants elimination

The variants are just expanded:

GF

```
lin girl_N = mkN (variants {"tjej"; "flicka"});
```

GF Core

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

Note: Appropriate for application specific grammars. Should be avoided in resource grammars.

Variants are not always what you want

GF

```
lin Answer pol verb = "I" ++pol ++verb;  
eat = "eat";  
like = "like";  
Pos = "";  
Neg = variants{"do not"; "don't"};  
  
Comp s1 s2 = s1 ++ ";" ++ s2;
```

Comp (Answer Neg like) (Answer Neg eat)

I don't like; I don't eat
I don't like; I do not eat
I do not like; I don't eat
I do not like; I do not eat

Variants are not always what you want

GF

```
lin Answer pol verb = \\style ⇒ "I" ++ (pol!style) ++ verb;  
eat = "eat";  
like = "like";  
Pos = table{ Official ⇒ "" ; Spoken ⇒ "" };  
Neg = table{ Official ⇒ "do not" ; Spoken ⇒ "don't" };  
  
Comp ids1 ids2 = variants{ comp Official ; comp Spoken };  
  
oper comp style = s1!style ++ ";" ++ s2!style;
```

Comp (Answer Neg like) (Answer Neg eat)

I don't like; I don't eat

I do not like; I do not eat

The variants could blow up

When many variants are used in parallel the number of core rules grows exponentially.

GF

```
lin start_word = variants{"open";"start"} ++ variants{"Word";"Microsoft Word"};
```

GF Core

```
lin start_word1 = "open" ++ "Word";  
start_word2 = "open" ++ "Microsoft Word";  
start_word3 = "start" ++ "Word";  
start_word4 = "start" ++ "Microsoft Word";
```

Variants explosion with tables

GF

```
lin close_word = \\tense ⇒ close_V!tense ++ variants{"Word";"Microsoft Word"};  
oper close_V = table Tense {"close";"closed";"have closed";...};
```

Desugared GF

```
lin close_word = table Tense {  
  "close" ++ variants{"Word";"Microsoft Word"};  
  "closed" ++ variants{"Word";"Microsoft Word"};  
  "have closed" ++ variants{"Word";"Microsoft Word"}  
  ...  
};
```

Note: Leads to $2^3 = 8$ possible combinations although there is only one variant in the original code

Variants explosion with tables - CORE

GF Core

```
lin close_word1 = <"close" ++ "Word";  
                "closed" ++ "Word";  
                "have closed" ++ "Word" >;  
close_word2 = <"close" ++ "Microsoft Word";  
              "closed" ++ "Word";  
              "have closed" ++ "Word" >;  
close_word3 = <"close" ++ "Word";  
              "closed" ++ "Microsoft Word";  
              "have closed" ++ "Word" >;  
close_word4 = <"close" ++ "Microsoft Word";  
              "closed" ++ "Microsoft Word";  
              "have closed" ++ "Word" >;  
...
```

Variants explosion with tables - SOLUTION

GF

```
lin close_word = variants{closelt "Word"; closelt "Microsoft Word"};  
oper closelt obj = \\tense ⇒ close_V!tense ++ obj;  
close_V = table Tense {"close"; "closed"; "have closed"; ...};
```

- **Variants should be used with care**
- A variant in the wrong place could lead too many combinations, which is often not what you want
- A combinatorial explosion could kill the compiler with

Out of memory

- Unnecessary combinations will slow down the parser
- Hint: use the command 'l -all'

Parameter Types Elimination

lincat $NP = \{s : Case \Rightarrow Str; g : Gender; n : Number; p : Person\}$
param $Case = Nom|Acc|Dat;$
 $Gender = Masc|Fem|Neutr;$
 $Number = Sg|Pl;$
 $Person = P1|P2|P3;$

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 tuple:

$\langle s_1, s_2, s_3 \rangle$

Then in general 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 .

Parameter Fields Elimination

GF

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

GF Core

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

Note: The number of categories doesn't immediately affect the size of the compiled grammar

Counting Parametric Types

It is important to know how many possible values a given parameter type has because:

- This determines the number of fields in the core:

$$P \Rightarrow Str$$

- This determines the number of categories in the core:

$$\{ \dots ; p : P \}$$

Counting the number of parameter values

Parameter Definition

$$\begin{array}{l} \mathbf{param} P = P_1 Q_{11} Q_{12} \dots Q_{1m_1} \\ \quad | P_2 Q_{21} Q_{22} \dots Q_{2m_2} \\ \quad \dots \\ \quad | P_n Q_{n1} Q_{n2} \dots Q_{nm_n} \end{array}$$

Values Count

$$\begin{array}{l} \mathbb{C}(P) = \mathbb{C}(Q_{11}) * \mathbb{C}(Q_{12}) \dots \mathbb{C}(Q_{1m_1}) \\ \quad + \mathbb{C}(Q_{21}) * \mathbb{C}(Q_{22}) \dots \mathbb{C}(Q_{2m_2}) \\ \quad \dots \\ \quad + \mathbb{C}(Q_{n1}) * \mathbb{C}(Q_{n2}) \dots \mathbb{C}(Q_{nm_n}) \end{array}$$

Counting Parametric Tables and Records

Parametric Records

$$\mathbb{C}(\{q_1 : Q_1; q_2 : Q_2 \dots q_n : Q_n\}) = \mathbb{C}(Q_1) * \mathbb{C}(Q_2) \dots \mathbb{C}(Q_n)$$

Parametric Tables

$$\mathbb{C}(P \Rightarrow Q) = \mathbb{C}(Q)^{\mathbb{C}(P)}$$

Warning: Exponentials should be avoided!!!

- **Keep the lexicon compact:**

$$\text{lincat } N = \{s : NForm \Rightarrow Str; g : DGender\};$$

param *NForm*

= *NF Number Species*

| *NFSgDefNom*

| *NFPICount*

| *NFVocative*

param *NForm*

= *NF Number Species Case*

| *NFPICount*

- Linguistically accurate
- The irregularity is obvious
- Mathematically elegant
- Linguistically overgenerating

Comment

The lexical items are inflection tables. Duplication means overhead for every entry in the lexicon.

- **Keep the syntax elegant:**

$$\text{lincat } CN = \{s : \textit{Number} \Rightarrow \textit{Species} \Rightarrow \textit{Case} \Rightarrow \textit{Str}\};$$

Comment

The syntactic rules are closed set. Compared to the lexicon this is a small set so it is not so important to make them compact. It is much more important to have clear easy to manipulate structure.

The efficiency of the parser is not affected by the number of fields in the linearization types.

Hints for Parameters

- Minimize the number of inherent parameters

$$\text{lincat } N = \{s : NForm \Rightarrow Str; g : DGender\};$$

```
param DGender  
  = DMasc Animacy  
  | DFem  
  | DNeutr
```

```
oper DGender = {g : Gender; a : Animacy}  
param Gender  
  = Masc  
  | Fem  
  | Neutr;  
param Animacy = Animate | Inanimate;
```

- Animacy matters only for Masc
- Animacy given for all genders

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 >
```


Counting Linearization Rules

In general linearization rule like:

fun $f : A_1 \rightarrow A_2 \rightarrow \dots \rightarrow A_n \rightarrow A;$

produces $\mathbb{C}(f)$ rules in the core

$$\mathbb{C}(f) = \mathbb{C}(A_1) * \mathbb{C}(A_2) * \dots * \mathbb{C}(A_n)$$

Comment

The number of rules could be reduced by reducing the number of parameters in the linearization types. The count is also reduced by the optimizations in the compiler.

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"};  
  ...  
}
```

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

Agglutination

- Some languages have potentially infinite set of words:

Turkish:

anlamıyorum = 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

- 1 Introduction
- 2 GF Core
- 3 Optimizations**
- 4 Debugging
- 5 Conclusion

Three main optimizations reduce the exponential size of the grammar:

- Common Subexpressions Optimization
- Common Functions Optimization
- Coercion Rules

Note: the optimizations cannot be expressed in GF Core. PMCFG is needed.

Common Subexpressions Optimization

GF Core

```
lin u x y = < x.p1, x.p2 ++ y.p1 >  
    v x y = < "a", x.p2 ++ y.p1 >
```

PMCFG

$$F_1 := (S_1, S_2) [u]$$
$$F_2 := (S_3, S_2) [v]$$
$$S_1 := \langle 0; 0 \rangle$$
$$S_2 := \langle 0; 1 \rangle \langle 1; 0 \rangle$$
$$S_3 := "a"$$

Common Subexpressions Optimization in the Lexicon

GF Core

lin *good_A* = < "dobřr", "dobra", "dobro", "dobre" >

beautiful_A = < "hubav", "hubava", "hubavo", "hubavo" >

PMCFG

$F_1 := (S_1, S_2, S_3, S_4)$ [*good_A*]

$F_2 := (S_5, S_6, S_7, S_7)$ [*beautiful_A*]

$S_1 :=$ "dobřr" $S_5 :=$ "hubav"

$S_2 :=$ "dobra" $S_6 :=$ "hubava"

$S_3 :=$ "dobro" $S_7 :=$ "hubavo"

$S_4 :=$ "dobre"

Common Functions Optimization

The function symbols in PMCFG could be reused in different productions

PMCFG

$$C_1 \leftarrow F_1[C_2, C_3]$$

$$C_1 \leftarrow F_1[C_4, C_5]$$

$$F_1 := (S1, S2, S3, S4) \quad [u]$$

PMCFG

$$C_1 \leftarrow F_1[C_2, C_{31}, C_{41}, C_5]$$
$$C_1 \leftarrow F_1[C_2, C_{32}, C_{41}, C_5]$$
$$C_1 \leftarrow F_1[C_2, C_{33}, C_{41}, C_5]$$
$$C_1 \leftarrow F_1[C_2, C_{31}, C_{42}, C_5]$$
$$C_1 \leftarrow F_1[C_2, C_{32}, C_{42}, C_5]$$
$$C_1 \leftarrow F_1[C_2, C_{33}, C_{42}, C_5]$$
$$C_1 \leftarrow F_1[C_2, C_{31}, C_{43}, C_5]$$
$$C_1 \leftarrow F_1[C_2, C_{32}, C_{43}, C_5]$$
$$C_1 \leftarrow F_1[C_2, C_{33}, C_{43}, C_5]$$

PMCFG

$$C_1 \leftarrow F_1[C_2, C_3, C_4, C_5]$$

$$C_3 \leftarrow -[C_{31}]$$

$$C_3 \leftarrow -[C_{32}]$$

$$C_3 \leftarrow -[C_{33}]$$

$$C_4 \leftarrow -[C_{41}]$$

$$C_4 \leftarrow -[C_{42}]$$

$$C_4 \leftarrow -[C_{43}]$$

- 1 Introduction
- 2 GF Core
- 3 Optimizations
- 4 Debugging**
- 5 Conclusion

Debugging the Compiler

You can dump the PMCFG representation of the grammar with the following command:

```
c:\gf> gf -make -output-format=pmcfg_pretty LangMy.gf
Reading Lang.pgf...
Refusing to overwrite Lang.pgf
Writing LangEng.pmcfg...
```

This will produce one file with extension `.pmcfg` with four interesting sections:

- productions
- functions
- sequences
- startcats

Thank You and Have Fun !!!