

Linearisation-only PGF

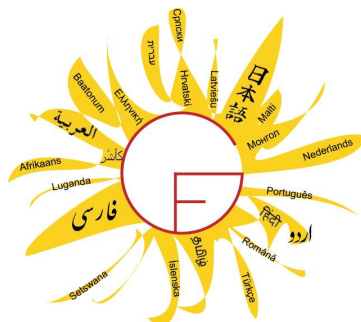
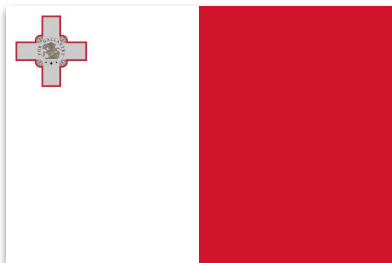
GF Summer School 2021

John J. Camilleri

digital  grammars
Language technology to rely on.

About me

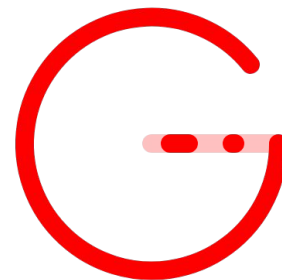
- From Malta
- Attended first GF summer school in Gothenburg, 2009
- Ph.D. in Computer Science at Chalmers & University of Gothenburg
- CTO at Digital Grammars



GF Summer School
Gothenburg, Sweden
August 17-28, 2009

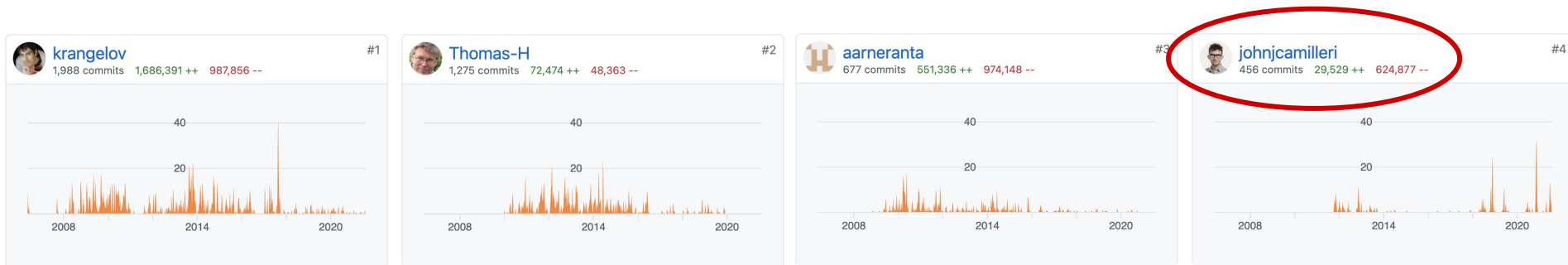


**UNIVERSITY OF
GOTHENBURG**





Me and GF



Source: <https://github.com/GrammaticalFramework/gf-core/graphs/contributors>

- Worked with GF on/off for 12 years
- Some grammars, but mostly integration and internal stuff
- At DG we use GF in the "real world"
- What I've learnt is that...

parsing is hard

Problems with parsing

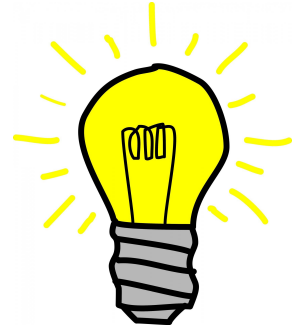
- Ambiguity is impossible to avoid
 - probabilities?
 - statistical models?
 - exceptions?
- Fitting grammar to a corpus is unsatisfying work
- Corpus is often wrong, inconsistent, or otherwise problematic
- Coverage comes at cost of overgeneration

generation is easy

What's nice about generation

- It's not parsing
- There's no ambiguity (if you ignore variants)
- Data is usually more structured & consistent
- If customer data already exists, half the work is already done
- If not, designing abstract data more *is* a satisfying task
- Much better control of the process
- Generation is a "simpler" operation

Motivation behind LPGF



1. GF grammars enable both parsing and generation via PGF.
2. Compiling to PGF can be slow & memory-intensive.
3. Often we don't want to parse at all.
4. Let's compile to something smaller which only supports linearisation: LPGF.
5. Anticipated benefits:
 - a. Compiling GF to LPGF will be less resource-demanding than to PGF
 - b. An LPGF file will be smaller than a PGF file
 - c. The speed of linearisation will be faster and require less memory
 - d. Not needing to support parsing can open up other possible features
6. The theory already exists!

Published: 15 December 2009

PGF: A Portable Run-time Format for Type-theoretical Grammars

[Krasimir Angelov](#) , [Björn Bringert](#) & [Aarne Ranta](#)

Journal of Logic, Language and Information **19**, 201–228 (2010) | [Cite this article](#)

104 Accesses | **11** Citations | [Metrics](#)

Abstract

Portable Grammar Format (PGF) is a core language for type-theoretical grammars. It is the target language to which grammars written in the high-level formalism Grammatical Framework (GF) are compiled. Low level and simple, PGF is easy to reason about, so that its

A bit of history

Paper (2009)	Today (2021)
PGF (section 2)	LPGF
PMCFCG (section 3)	PGF

- Initially:
 - PGF for linearisation
 - PMCFCG for parsing
 - GF compiled to both
- Later discovered that:
 - PGF can also be used for parsing (but very inefficient)
 - PMCFCG can also be used for linearisation (but has size issues)
- Currently:
 - GF is compiled only to PMCFCG
 - Makes compiler and runtime simpler (with some tradeoffs in performance and features)
- Proposal:
 - Revive linearisation-only PGF
 - Pure generation is a common use case
 - Reclaim performance tradeoffs & add extra features

The situation has been clear for a long time:

- *PMCFG-only is a simpler solution to implement*
- *but it is too costly for some grammars,
and unnecessarily so if only linearization is needed*
- *some tasks that one could in principle perform with GF
are therefore simply not possible in practice*

Aarne Ranta, personal communication (2021)

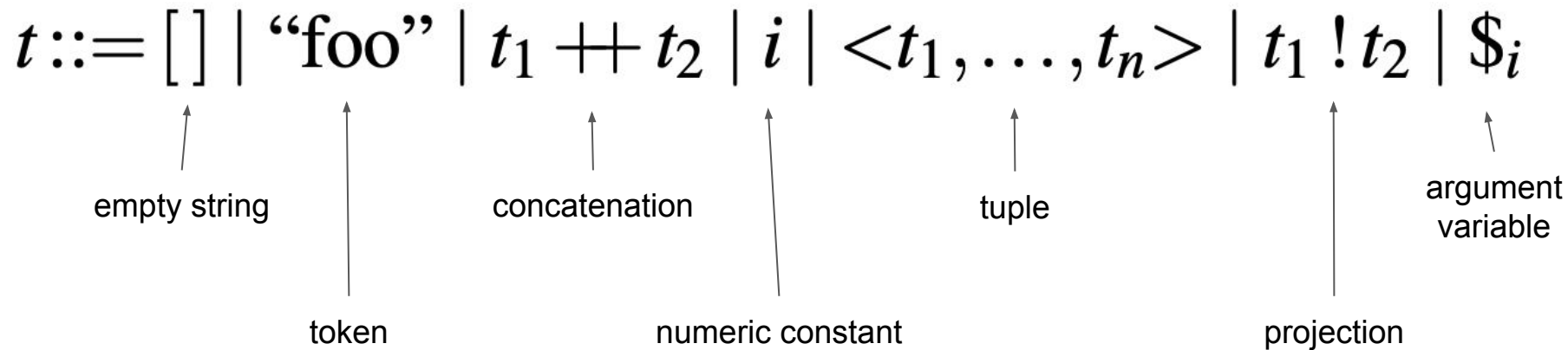


I'm sold!



down to the nitty-gritty...

LPGF: syntax



LPGF: type system

Strings:

$$[] : \text{Str} \quad \text{“foo”} : \text{Str} \quad \frac{s, t : \text{Str}}{s ++ t : \text{Str}}$$

Bounded integers:

$$i : \text{Int}_i \quad \frac{i : \text{Int}_m}{i : \text{Int}_n} \quad m < n$$

Tuples:

$$\frac{t_1 : T_1 \dots t_n : T_n}{\langle t_1, \dots, t_n \rangle : T_1 * \dots * T_n}$$

Projections:

$$\frac{t : T^n \quad u : \text{Int}_n}{t ! u : T} \quad \frac{t : T_1 * \dots * T_n}{t ! i : T_i} \quad i = 1, \dots, n$$

Argument variables:

$$\frac{}{T_1, \dots, T_n \vdash \$i : T_i} \quad i = 1, \dots, n$$

LPGF: operational semantics \Downarrow

Strings:

$$[] \Downarrow [] \quad \text{“foo”} \Downarrow \text{“foo”} \quad \frac{s \Downarrow v \quad t \Downarrow w}{s ++ t \Downarrow v ++ w}$$

Bounded integers:

$$i \Downarrow i$$

Tuples:

$$\frac{t_1 \Downarrow v_1 \dots t_n \Downarrow v_n}{\langle t_1, \dots, t_n \rangle \Downarrow \langle v_1, \dots, v_n \rangle}$$

Projections:

$$\frac{t \Downarrow \langle v_1, \dots, v_n \rangle \quad u \Downarrow i}{t!u \Downarrow v_i} \quad i = 1, \dots, n$$

Argument variable:

$$v_1, \dots, v_n \vdash \$i \Downarrow v_i \quad (i = 1, \dots, n)$$

LPGF: linearisation operation \mapsto

$$\frac{a_1 \mapsto t_1 \quad \dots \quad a_n \mapsto t_n \quad t_1, \dots, t_n \vdash t \Downarrow v}{f a_1 \dots a_n \mapsto v} \mathbf{lin} f = t$$

Implementation: starting point

- LPGF datatype is implemented as:
`LPGF` (src/runtime/haskell/LPGF.hs)
- GF parser gives you a term of type:
`SourceGrammar` (src/compiler/GF/Grammar/Grammar.hs)
- Top-level PGF compile function has type:
`mkCanon2pgf :: ... SourceGrammar -> IO PGF`
- Top-level LPGF compile function has type:
`mkCanon2lpgf :: ... SourceGrammar -> IO LPGF`

At some point we supported both formats and used them for different purposes. The problem was that the compilation to linearization-only format was full of bugs. After I spent many hours of fixing the compiler, I just dropped it and started using the parsing-only format for everything.



Krasimir Angelov, personal communication (2021)

At some point we supported both formats and used them for different purposes. The problem was that the compilation to linearization-only format was full of bugs. After I spent many hours of fixing the compiler, I just dropped it and started using the parsing-only format for everything.



I remember there were some problems, but some of them were due to the normalization of expressions not working properly. Now that Thomas [Hallgren] fixed it long ago (normalization by evaluation), the task should be easier.

GF Canonical format

src/compiler/GF/Grammar/Canonical.hs

- A subset of GF
- What's left after high-level constructions such as functors and ops have been eliminated by partial evaluation.
- Intended as a common intermediate representation to simplify export to other formats.

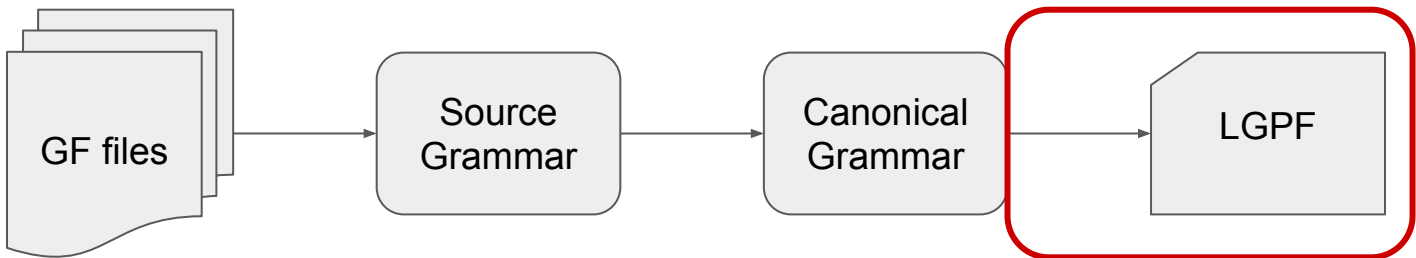
Using this should allow me to avoid the previous issues! 🙌

Hint: use `gf --output-format=canonical_gf` to see for yourselves

Using canonical format

```
grammar2canonical :: Options  
  -> ModuleName  
  -> GF.Grammar.Grammar.Grammar -- SourceGrammar  
  -> GF.Grammar.Canonical.Grammar
```

src/compiler/GF/Compile/Grammar/GrammarToCanonical.hs



The background is a vibrant blue digital tunnel. It features numerous bright blue and white light streaks that radiate from the center, creating a sense of depth and motion. Scattered throughout the scene are various binary digits (0s and 1s) and small, glowing square shapes, some of which appear to be floating or moving along the light paths. The overall effect is one of high-speed data flow or a journey through a digital space.

weeks pass...

LPGF: Linearisation-only PGF format #103

Edit

Open with ▾

johnjcamilleri wants to merge 112 commits into `master` from `lpgf` 📄

💬 Conversation 0

🔗 Commits 112

📄 Checks 15

📄 Files changed 217

+311,269 -35 🟢🟢🟢🟢



johnjcamilleri commented on 17 Mar

Member



Introduction

Recently I've been working on resurrecting on old idea, which is adding support for a PGF file format which *only* supports linearisation, since this is actually quite a common use for GF. The motivations are:

1. Faster & less memory-intensive compilation
2. Smaller binary files
3. Faster linearisation at runtime
4. New features impossible with parsing, e.g. dynamic lexicon.

The format itself is described in section 2 of the paper:

"PGF: A Portable Run-Time Format for Type-Theoretical Grammars"

Reviewers



Suggestions



Thomas-H

Request

Assignees



No one—assign yourself

Labels



None yet

Projects



None yet

Implementation: summary

1. New output format for GF compilation:

```
gf --make --output-format=lpgf FoodsEng.gf
```

2. Binary format for grammar distribution: `Foods.lpgf`
3. Haskell runtime library: `import LPGF` (next slide)
4. Testsuite for checking correctness w.r.t. treebank (`testsuite/lpgf`)
5. Benchmark for comparing LPGF with PGF/PGF2

Using LPGF Haskell runtime

```
1  import LPGF
2  import qualified Data.Map as M
3
4  main :: IO ()
5  main = do
6      ..lpgf <- readLPGF "Foods.lpgf"
7      ..let
8          ...Just eng = readLanguage "FoodsEng"
9          ...Just concr = M.lookup eng (concretes lngf)
10         ...Just tree = readExpr "Pred (This (Mod Italian Cheese)) (Very Expensive)"
11         ...str = linearizeConcrete concr tree
12     ..putStrLn str
```

Two notable issues

1. Variants
2. Missing linearisation functions

Variants

```
lin car_N = mkN "Auto" "Autos" neuter  
          | mkN "Wagen" "Wagen" masculine ;
```

Source: LexiconGer.gf (RGL)

- Forgivable in parsing, problematic in linearisation
- Can cause combinatorial explosions during compilation
- Somewhat debated (see issues [#14](#), [#36](#), [#37](#))
- Uncontrollable variation is **not in the spirit of NLG**
- My solution: don't implement them (always pick first variant)



Missing linearisation functions

Phrasebook.gf

```
fun Indian : Citizenship ;
```

PhrasebookSpa.gf

```
-- lin Indian = ... TODO
```

Expression

```
PQuestion (HowFar (ThePlace (CitRestaurant Indian)))
```

PGF

```
¿a qué distancia el restaurante [Indian] está?
```

LPGF

```
¿a qué distancia el [Indian] está?
```

benchmarks & profiling

is it really as fast as we hoped?

Benchmark results

Using Phrasebook grammar in all languages (assuming *.gfo files precompiled)

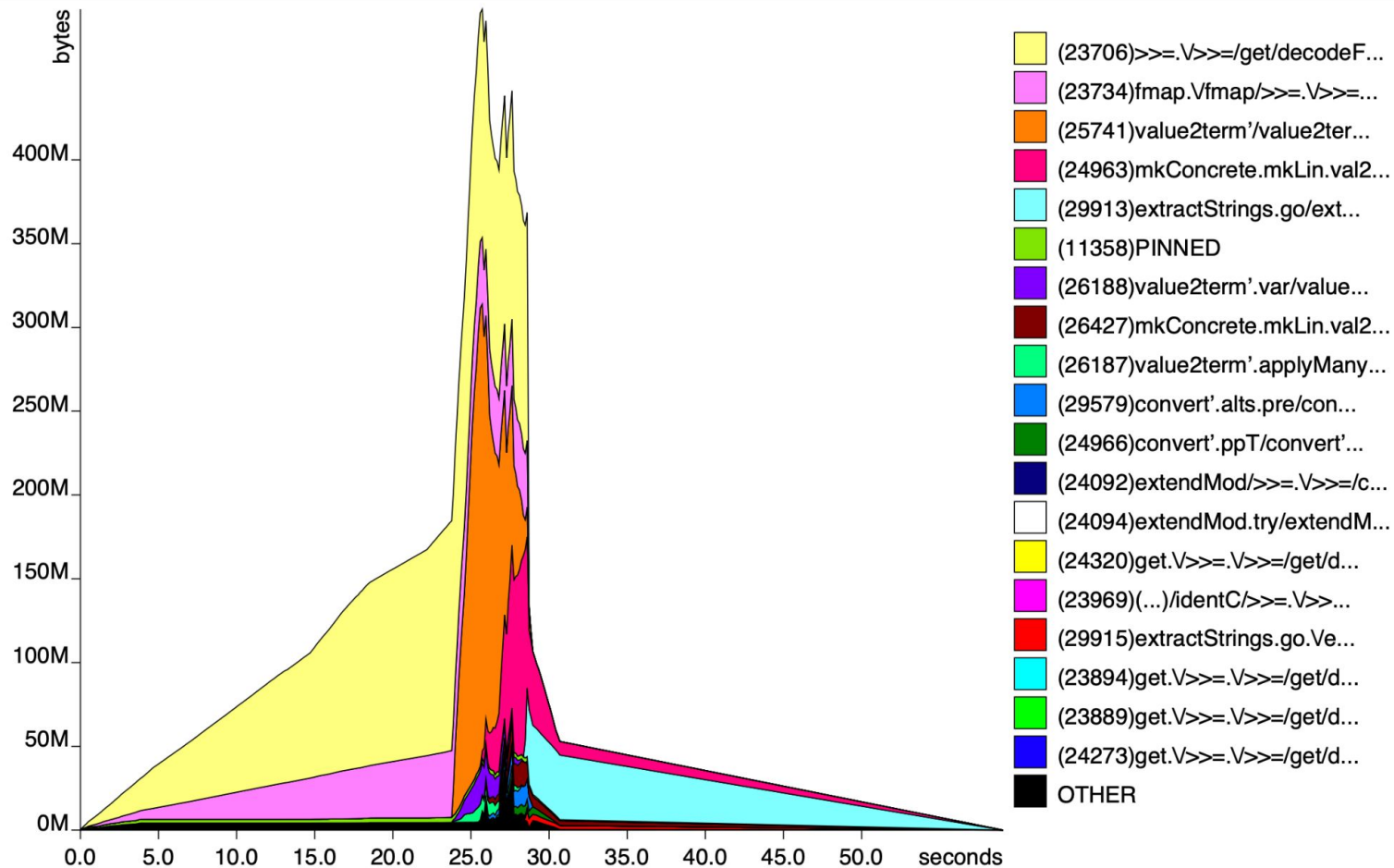
Compilation	Time	Memory	File size
PGF	33.47s	6.42 GB	39.75 MB
LPGF	53.60s	9.84 GB	13.12 MB

Linearisation (10,000 trees)	Time	Memory
PGF	5.67s	2650.00 MB
PGF2 (C runtime)	50.66s	109.05 MB
LPGF	3.49s	718.27 MB

lpgf-bench compile lpgf testsuite/lpgf/phrasebook/PhrasebookFre.gf +RTS -p -T -p

4,966,763,858 bytes x seconds

Wed Aug 4 12:04 2021



lpgf-bench +RTS -p -T -p -h -RTS compile lpgf testsuite/lpgf/phrasebook/PhrasebookFre.gf

total time = 22.68 secs (22681 ticks @ 1000 us, 1 processor)

total alloc = 11,568,701,440 bytes (excludes profiling overheads)

COST CENTRE	MODULE	SRC	%time	%alloc
>>=.\ return fmap.\ value2term' mkConcrete.mkLin.val2lin >>= >> convert'.ppT extractStrings.go.\.str mconcatMap extractStrings.go tableTypes.tabtys value2term'.v2txs mkConcrete.mkLin.val2lin.grps compare convert'.alts.pre put unzipR.(...) fmap convert'.ppP mkConcrete.mkLin.val2lin.\ zipAssign ident2utf8 convert'.ppP.fields unzipR	Data.Binary.Get Data.Binary.Get Data.Binary.Get GF.Compile.Compute.Concrete GF.Compile.GrammarToLPGF GF.Data.ErrM Data.Binary.Put GF.Compile.GrammarToCanonical GF.Compile.GrammarToLPGF GF.Grammar.Macros GF.Compile.GrammarToLPGF GF.Compile.GrammarToCanonical GF.Compile.Compute.Concrete GF.Compile.GrammarToLPGF GF.Infra.Ident GF.Compile.GrammarToCanonical Data.Binary GF.Grammar.Macros GF.Data.ErrM GF.Compile.GrammarToCanonical GF.Compile.GrammarToLPGF GF.Grammar.Macros GF.Infra.Ident GF.Compile.GrammarToCanonical GF.Grammar.Macros	src/runtime/haskell/Data/Binary/Get.hs:(129,28)-(130,54) src/runtime/haskell/Data/Binary/Get.hs:126:5-34 src/runtime/haskell/Data/Binary/Get.hs:(117,27)-(118,48) src/compiler/GF/Compile/Compute/Concrete.hs:(504,1)-(553,30) src/compiler/GF/Compile/GrammarToLPGF.hs:(123,9)-(274,71) src/compiler/GF/Data/ErrM.hs:(38,3)-(39,21) src/runtime/haskell/Data/Binary/Put.hs:(97,5)-(100,35) src/compiler/GF/Compile/GrammarToCanonical.hs:(178,5)-(203,47) src/compiler/GF/Compile/GrammarToLPGF.hs:345:15-38 src/compiler/GF/Grammar/Macros.hs:504:1-30 src/compiler/GF/Compile/GrammarToLPGF.hs:(334,5)-(363,19) src/compiler/GF/Compile/GrammarToCanonical.hs:(128,5)-(132,31) src/compiler/GF/Compile/Compute/Concrete.hs:536:5-32 src/compiler/GF/Compile/GrammarToLPGF.hs:180:15-44 src/compiler/GF/Infra/Ident.hs:58:17-19 src/compiler/GF/Compile/GrammarToCanonical.hs:(243,9)-(247,55) src/runtime/haskell/Data/Binary.hs:318:5-22 src/compiler/GF/Grammar/Macros.hs:170:35-51 src/compiler/GF/Data/ErrM.hs:(53,3)-(54,24) src/compiler/GF/Compile/GrammarToCanonical.hs:(218,5)-(233,49) src/compiler/GF/Compile/GrammarToLPGF.hs:(242,47)-(244,43) src/compiler/GF/Grammar/Macros.hs:182:1-51 src/compiler/GF/Infra/Ident.hs:(83,1)-(88,16) src/compiler/GF/Compile/GrammarToCanonical.hs:232:9-57 src/compiler/GF/Grammar/Macros.hs:170:1-51	71.7 1.9 1.7 1.4 0.9 0.9 0.9 0.7 0.6 0.5 0.5 0.5 0.5 0.3 0.3 0.3 0.2 0.2 0.2 0.2 0.2 0.2 0.2 0.2 0.1	9.3 3.8 5.3 11.7 2.4 3.7 5.2 1.1 3.9 3.1 2.1 3.4 6.1 1.4 3.7 1.1 1.1 1.4 1.0 1.5 1.7 1.3 1.2 1.2 1.1

Optimisations

In place

- Pruning impossible subterms
- Raw strings stored exactly once in binary file
- Make `GF.Grammar.Canonical.Id` a type synonym for `GF.Infra.Ident.RawIdent`

Tested without success


- Use `String` instead of `Data.Text` ([lpgf-string](#))
- Memoisation of subterms ([lpgf-memo](#))

Untested ?

- Playing with Haskell laziness
- Optimising `extractStrings` optimisation

conclusions

Status summary

- It works*!
- LPGF files are [often] smaller than PGFs 😄
- Runtime linearisation in LPGF is faster than both PGF and PGF2 🎉
- Compiling to LPGF is at least as slow/memory-consuming as PGF, often significantly worse 😞
- Extra features... 

* except for variants and missing functions

Roadmap

Immediate priorities

- As fast and small as possible
 - compilation
 - runtime (is Haskell a problem here?)
- More features to runtime API
 - type-checking of user-generated expressions

Probably also want

- GF shell support
- Runtime in other languages (native or via bindings)
 - JavaScript?
 - Python?
 - Java?
 - C?

Extra features

1. Dynamic lexicon

- No recompiling when lexicon changes
- Linearise function in runtime takes extra `Lexicon` argument
- But: access to smart paradigms?
- Syncing lexicon with grammar

2. No need to know all tokens at compile time

- it could allow things such as runtime gluing
- conversion of integers to numeral expressions, which is often expected by users
- lin-only extensions of the GF source language

3. Multi-PGF format

- some concrete syntaxes are full PMCFG
- some concrete syntaxes are linearisation-only PGF
- parsing with a linearisation-only language is an error
- all other cases work seamlessly

text:ual

<https://textual.ai/>

- Presented at GFSS 2018, South Africa
- Pure NLG for product descriptions in e-commerce
- 25 languages
- Lexicon is huge and constantly updated
- Recompiling to PGF on-the-fly is impractical
- Home-grown dynamic lexicon solution using placeholders

Learn more & contribute

- Questions to the community
 - Who else is this valuable for?
 - Which features do you want to see/prioritise?
- Draft pull request
 - <https://github.com/GrammaticalFramework/gf-core/pull/103>
 - Should we merge it?
- LPGF Dev README
 - <https://github.com/GrammaticalFramework/gf-core/blob/lpgf/testsuite/lpgf/README.md>
 - How to run tests/benchmarks
 - Details from various investigations
- Help with optimisations

Credits

- Aarne, for reviving this idea
- Krasimir, who knows everything about PGF
- Andreas "Anka" Källberg, for helping with Haskell profiling

