# Grammar Engineering Tools

John Camilleri, Ramona Enache, Thomas Hallgren, Aarne Ranta

MOLTO Final Presentation 2013



## Grammars as Software

Key to high-quality translation: control over details, debugging

As opposed to: holistic systems, more data, parameter tuning

Similar to: compilers (translators of computer languages)

- expected to translate correctly
- pipeline: parsing + semantic analysis + generation
- semantics encoded in **abstract syntax**

# Compilation via abstract syntax



## Translation via abstract syntax



# **Translation example**



Catalan: Guernica està pintat sobre llenç per Pablo Picasso en 1937. Dutch: Guernica werd in 1937 door Pablo Picasso op canvas geschilderd. English: Guernica was painted on canvas by Pablo Picasso in 1937. Finnish: Guernican maalasi Pablo Picasso kankaalle vuonna 1937. French: Guernica a été peint sur canvas par Pablo Picasso en 1937.

# Multilingual grammar in GF

Declarative program defining the translation relation among any number n of languages

- Abstract: fun Painted : Painting -> Painter -> Fact
- English: lin Painted x y = x ++ "painted" ++ y
- Finnish: lin Painted x y = x ++ "maalasi" ++ y
- French: lin Painted x y = x ++ "a peint" ++ y

But isn't this too simple-minded?

# The complexity of concrete syntax

French: agreement, clitics, ... (*il a peint X* vs. *j'ai peint X* vs. *il les a peintes* ...)

lin

```
Painted x y = x.s ! Nom ++ case y.isPron of {
  True => y.s ! Acc ++ avoir_V ! x.agr ++ peindre_V ! PastPart y.agr ;
  False => avoir_V ! x.agr ++ peindre_V ! PastPart MascSg ++ y.s ! Acc
  }
```

```
avoir_V = table ["avoir","ai","as","a","avons",...]
```

Moreover: tenses, negation, question forms, ...

# The complexity of multilingual systems

Two dimensions: semantic components X languages. For example:

module	Bulgarian	Catalan	Dutch	English	
Answer	AnwerBul	AnswerCat	AnswerDut	AnswerEng	•••
Query	QueryBul	QueryCat	QueryDut	QueryEng	•••
Text	TextBul	TextCat	TextDut	TextEng	•••
Lexicon	LexiconBul	LexiconCat	LexiconDut	LexiconEng	•••
Data	DataBul	DataCat	DataDut	DataEng	

Museum Library (WP8): (1 + 15) \* 5 = 80 modules

Mathematics Library (WP6): (1 + 15) \* 16 + 27 = 676 modules

# Mastering the complexity

**Programming language**: GF - functions, types, modules

**Compiler**: type checking, optimizations

**Library**: low-lever linguistic details

**Development environment**: projects consistency, code navigation

**Documentation**: tutorials, reference manuals, best practices

**Training**: tutorial events, on-line courses

**Community**: helping each other

# The GF programming language

First created at Xerox Research in 1998

For **CS people**: a special-purpose functional language for grammars (like YACC, but more powerful)

For **MT people**: a formalism for synchronous grammar (like TAG, but more powerful)

For **language theory people**: a front-end to PMCFG (Parallel Multiple Context-Free Grammars)

New things during MOLTO:

• probabilistic GF grammars

# The GF compiler

From high-level GF to low-level PGF (Portable Grammar Format)

Separate compilation of modules

Code generation to different formats (e.g. Nuance, XFST/Lexc, Giza)

#### New things during MOLTO:

- the PGF format
- optimized compilation
- run-time bindings from C, C++, Java, Python
- compilation as cloud service

# The GF Resource Grammar Library

Complete morphology engine + comprehensive syntax + lexicon

Afrikaans	Bulgarian	Catalan	Chinese	Danish	Dutch	English
Finnish	French	German	Greek	Hindi	Italian	Japanese
Latvian	Maltese	Nepali	Norwegian	Persian	Polish	Punjabi
Romanian	Russian	Sindhi	Spanish	Swedish	Thai	Urdu

#### New during MOLTO:

- 13 new languages (built outside MOLTO): 9 Asian, 3 EU
- big lexicon resources (10-100k lemmas) for 11 languages

#### The library **API**

#### Cl - declarative clause, with all tenses

Function	Туре	Example	
genericCl	<u>VP</u> -> <u>Cl</u>	one sleeps	
mkCl	<u>NP</u> -> <u>V</u> -> <u>Cl</u>	she sleeps	
mkCl	<u>NP</u> -> <u>V2</u> -> <u>NP</u> -> <u>Cl</u>	she loves him	
mkCl	<u>NP</u> -> <u>V3</u> -> <u>NP</u> -> <u>NP</u> -> <u>Cl</u>	she sends it to him	
mkCl	<u>NP</u> -> <u>VV</u> -> <u>VP</u> -> <u>Cl</u>	she wants to sleep	
mkCl	<u>NP</u> -> <u>VS</u> -> <u>S</u> -> <u>C1</u>	she say • API: mkUtt (mkCl sh	e NP want VV (mkVP sleep V))
mkCl	<u>NP</u> -> <u>VQ</u> -> <u>QS</u> -> <u>Cl</u>	she wol • Afr: sy wil te slaap	· · · · · · · · ( · · · · · ·
mkCl	<u>NP</u> -> <u>VA</u> -> <u>A</u> -> <u>Cl</u>	• Bul: тя иска да спи • Cat: ella vol dormir	
mkCl	<u>NP</u> -> <u>VA</u> -> <u>AP</u> -> <u>Cl</u>	she bec • Dan: hun vil sove	
mkCl	<u>NP</u> -> <u>V2A</u> -> <u>NP</u> -> <u>A</u> -> <u>Cl</u>	Dut: ze wil slapen     Fng: she wants to sle	en l
mkCl	<u>NP</u> -> <u>V2A</u> -> <u>NP</u> -> <u>AP</u> -> <u>Cl</u>	she pai • Fin: hän tahtoo nukk	иа
mkCl	<u>NP</u> -> <u>V2S</u> -> <u>NP</u> -> <u>S</u> -> <u>Cl</u>	• Fre: elle veut dormir she ans	
mkCl	<u>NP</u> -> <u>V2Q</u> -> <u>NP</u> -> <u>QS</u> -> <u>C1</u>	she ask • Hin: वह सोना चाहती	है
mkCl	<u>NP</u> -> <u>V2V</u> -> <u>NP</u> -> <u>VP</u> -> <u>Cl</u>	she beg • Ita: lei vuole dormire	
mkCl	<u>NP</u> -> <u>VPSlash</u> -> <u>NP</u> -> <u>Cl</u>	Jpn: 彼女は寝たが     she beg	っている
mkCl	<u>NP</u> -> <u>A</u> -> <u>Cl</u>	she is a • Nep: उनी सुत्न चाहन्दि	बन्
mkCl	<u>NP</u> -> <u>A</u> -> <u>NP</u> -> <u>Cl</u>	she is a • Nor: hun vil sove	
mkCl	<u>NP -&gt; A2 -&gt; NP -&gt; Cl</u>	او می خواهد بخوابد :Pes • she is n	1
mkCl	<u>NP</u> -> <u>AP</u> -> <u>Cl</u>	<ul> <li>Pnb: او سوبا چاندی /ے she is v</li> <li>Pol: ona chce spać</li> </ul>	
mkCl	$NP \rightarrow NP \rightarrow Cl$	she is t • Ron: ea vrea sã doar	mă
mkCl	<u>NP -&gt; N -&gt; Cl</u>	• Rus: <i>она хочет cnar</i>	пь 9 сф
mkCl	<u>NP -&gt; CN -&gt; Cl</u>	she is a • Spa: ella quiere dorm	nir
mkCl	NP -> Adv -> Cl	she is h	
mkCl	NP -> VP -> Cl	• Tha: หลอนอยากนอน she alw • Urd: • مسونا جامتي هم	ทลบ
mkCl	<u>N</u> -> Cl	there is a nouse	

### The painted predicate with RGL

One-liners in every language - grammar writer can ignore details

lin Painted x y = mkS pastTense (mkCl x paint\_V2 y)

lin Painted x y = mkS pastTense (mkCl x maalata\_V2 y)

lin Painted x y = mkS perfectTense (mkCl x peindre\_V2 y)

# **GF** development environments

GF shell: support for interactive compilation and testing

IDE (Integrated Development Environment) - an Eclipse plug-in

Cloud-based grammar editor: on-line grammar development

New during MOLTO:

- the Eclipse IDE
- the cloud-based grammar editor

# **GF** documentation

http://www.grammaticalframework.org/

100+ articles on GF

#### New during MOLTO:

- 30+ articles
- Best practices
- The GF book: Aarne Ranta, *Grammatical Framework: Programming with Multilingual Grammars*, CSLI Publications, Stanford, 2011.
- Chinese translation of the book by Yan Tian, Shanghai, 2013.

CSLI Studies in Computational Linguistics

**GRAMMATICAL FRAMEWORK** is a programming language designed for writing grammars, which has the capability of addressing several languages in parallel. This thorough introduction demonstrates how to write grammars in Grammatical Framework and use them in applications such as tourist phrasebooks, spoken dialogue systems, and natural language interfaces. The examples and exercises presented here address several languages, and the readers are shown how to look at their own languages from the computational perspective.

Since the book requires no previous knowledge of linguistics, it can be an effective and useful resource for computer scientists and programmers, while introducing linguists to a novel approach to multilingual grammars inspired by the theory of programming languages.

Aarne Ranta is professor of computer science at the University of Gothenburg, Sweden. He is the acting coordinator of the European Union research project MOLTO (Multilingual On-Line Translation), which develops techniques for highguality translation among fifteen languages. Grammatical Framework
Programming with Multilingual Grammars

Aarne Ranta









# GF training events

Tutorials in large conferences: LREC-2010, CADE-2011, ICFP-2012

GF Summer Schools: 2009 Gothenburg, **2011 Barcelona**, 2013 Frauenchiemsee (Bavaria)

• 2-week event with 30 participants from 15 countries

# **GF** community

117 members in gf-dev mailing list

Around 50 resource grammar developers

Coverage of world's languages: http://www.postcrashgames.com/gf\_world/

Developers in most of these countries



# What is possible

Size of an average application: 15 languages, 200 functions

Size of the biggest application: 5 languages, 56k functions

Effort for building an average grammar: days for the first language, hours for the next ones

Skills required:

- to get a project started: domain expertise, some days of GF training
- to add a language: practical language skills, some hours of GF training

## Bootstrapping a grammar

To get started: design abstract syntax to fit an ontology

The first language: concrete syntax using RGL API and parsing examples

Later languages: change the words, and perhaps a few syntax functions

Extend vocabulary: extract words from other sources (wordnet, Wikipedia, Wiktionary)

#### Example: abstract syntax for CRM ontology

```
abstract QueryPainting = {
  cat
    Painting ; Query ;
  fun
    QPainter : Painting -> Query ; -- who painted x
    QYear : Painting -> Query ; -- when was x painted
    QMuseum : Painting -> Query ; -- where is x displayed
    QColour : Painting -> Query ; -- what colours does x have
    QSize : Painting -> Query ; -- what is the size of x
    QMaterial : Painting -> Query ; -- what material is x painted on
```

# Example: concrete syntax for English

```
concrete QueryPaintingEng of QueryPainting =
 open LexiconPaintingEng, SyntaxEng, ParadigmsEng in {
 lincat
   Painting = NP ; Query = QS ;
 lin
   QPainter t = mkQS pastTense (mkQCl who_IP paint_V2 t);
   QYear t = mkQS pastTense (mkQCl when_IAdv (mkCl t (passiveVP paint_V2))
   QMuseum t = mkQS (mkQCl where_IAdv (mkCl t displayed_VP))
   QColour t = mkQS (mkQCl whatPl_IP (mkNP thePl_Det (mkCN (mkN2 colour_N))
   QMaterial t = mkQS (mkQCl whatSg_IP (mkNP the_Det (mkCN (mkN2 material_)
   QSize t = mkQS (mkQCl whatSg_IP (mkNP the_Det (mkCN (mkN2 size_N) t)))
```

# Example: concrete syntax for German

```
concrete QueryPaintingGer of QueryPainting =
 open LexiconPaintingGer, SyntaxGer, ParadigmsGer in {
 lincat
   Painting = NP ; Query = QS ;
 lin
   QPainter t = mkQS pastTense (mkQCl who_IP malen_V2 t) ;
   QYear t = mkQS pastTense (mkQCl when_IAdv (mkCl t (passiveVP malen_V2))
   QMuseum t = mkQS (mkQCl where_IAdv (mkCl t ausgestellt_VP))
   QColour t = mkQS (mkQCl whatPl_IP (mkNP thePl_Det (mkCN (mkN2 farbe_N) ·
   QMaterial t = mkQS (mkQCl whatSg_IP (mkNP the_Det (mkCN (mkN2 material_)
   QSize t = mkQS (mkQCl whatSg_IP (mkNP the_Det (mkCN (mkN2 groesse_N) t))
```

# The smartest solution: functor

```
incomplete concrete QueryPaintingI of QueryPainting =
  open LexiconPainting, Syntax in {
    lincat
    Painting = NP ; Query = QS ;
    lin
        QPainter t = mkQS pastTense (mkQCl who_IP paint_V2 t) ;
        QYear t = mkQS pastTense (mkQCl when_IAdv (mkCl t (passiveVP paint_V2))
        QMuseum t = mkQS (mkQCl where_IAdv (mkCl t displayed_VP))
        QColour t = mkQS (mkQCl whatPl_IP (mkNP thePl_Det (mkCN (mkN2 colour_N)
        QMaterial t = mkQS (mkQCl whatSg_IP (mkNP the_Det (mkCN (mkN2 material_I
        QSize t = mkQS (mkQCl whatSg_IP (mkNP the_Det (mkCN (mkN2 size_N) t)))
```

sharing all code but the lexicon (works for 90% of rules)

# Example-based grammar writing

Extract translation rule by parsing an example

Abstract syntax	Like She He	first grammarian
English example	she likes him	first grammarian
German translation	er gefällt ihr	ORACLE
resource tree	mkCl he_Pron gefallen_V2 she_Pron	GF parser
concrete syntax rule	Like x y = mkCl y gefallen_V2 x	variables renamed

ORACLE = native speaker or statistical sentence alignment

Methodology with some tool support

# The MOLTO heritage

More languages in RGL: reason to build more applications

Applications: reason to support more languages in RGL

Tool of choice for controlled language implementation

Community growth, enterprise awareness

Next step: scaling up to open-domain translation (first experiments in MOLTO)

#### Demo: eclipse-film.m4v



Grammar cloning, library browsing, regression testing

# Publications related to MOLTO grammar tools

K. Angelov and A. Ranta. Implementing Controlled Languages in GF. N. Fuchs (ed.), *CNL-2009 Controlled Natural Languages*, LNCS/LNAI 5972, 2010.

J. Camilleri. An IDE for the Grammatical Framework. *Third International Workshop* on Free/Open-Source Rule-Based Machine Translation (FreeRBMT 2012).

G. Détrez and A. Ranta. Smart Paradigms and the Predictability and Complexity of Inflectional Morphology. EACL (European Association for Computational Linguistics), Avignon, April 2012.

R. Enache, A. Ranta, and K. Angelov. An Open-Source Computational Grammar of Romanian. A. Gelbukh (ed.), *CiCLING-2010*, LNCS 6008, 2010.

A. Ranta. Example-Based Grammar Writing. In S. Larsson and L. Borin (eds), *From Quantification to Conversation. Festschrift for Robin Cooper on the Occasion of his 65th Birthday.* College Publications, London, 2012.

A. Ranta. Machine Translation and Type Theory. In P. Dybjer, S. Lindström, E. Palmgren, and G. Sundholm (eds), *Epistemology versus Ontology. Essays on the Philosophy and Foundations of Mathematics in Honour of Per Martin-Löf.* Springer, Heidelberg, 2012. pp. 281-312.

A. Ranta, R. Enache, and G. Détrez. Controlled Language for Everyday Use: the MOLTO Phrasebook. In N. Fuchs and M. Rosner (eds), *Controlled Natural Language 2010*, Springer LNCS/LNAI, vol. 7175, 2012. pp. 115-136.

A. Ranta, *Grammatical Framework: Programming with Multilingual Grammars*, CSLI Publications, Stanford, 2011.

A. Ranta, K. Angelov, and T. Hallgren. Tools for multilingual grammar-based translation on the web. *Proceedings of the ACL 2010 System Demonstrations*, ACM Digital Library, 2010.

S. Virk, M. Humayoun, and A. Ranta. An Open-Source Punjabi Resource Grammar. Proceedings of RANLP-2011, Recent Advances in Natural Language Processing, Hissar, Bulgaria, 12-14 September, 2011. pp. 70-76.

S. Virk, M. Humayoun, and A. Ranta. An Open Source Urdu Resource Grammar. *Proceedings of the 8th Workshop on Asian Language Resources (Coling 2010 work-shop)*, 2010.

S. Virk. *Computational Linguistics Resources for Indo-Iranian Languages*, PhD Thesis, University of Gothenburg, 2013.