

The Apertium machine translation platform: five years on

Mikel L. Forcada^{1,2,3}, Francis M. Tyers²,
Gema Ramírez-Sánchez³

¹Centre for Next Generation Localisation, Dublin City University, Dublin 9 (Ireland)

²Departament de Llenguatges i Sistemes Informàtics, Universitat d'Alacant,
E-03071 Alacant (Spain)

³Prompsit Language Engineering, S.L., St. Francesc, 74, 1-L, E-03195 l'Altet
(Spain)

FreeRBMT'09, Alacant, Nov. 2, 2009

Contents

- 1 History
- 2 The Apertium platform
- 3 The Apertium philosophy
- 4 Technology
- 5 The Apertium community
- 6 Research
- 7 Business
- 8 Recent developments
- 9 Lots of work ahead
- 10 Funding

The inception/1

- April 2004: MLF's letter to main HLT groups in Spain: "let's lobby government agencies to fund building of FOS MT system for the languages of Spain (es, gl, ca, eu)".
- July 2004: Spanish Ministry of Industry, Tourism and Commerce launches call to fund consortia to develop linguistic technology for the languages of Spain.

The inception/2

- August 2004, consortium to build FOS MT systems forms:
 - Universities: EHU, UA, UPC, UVigo
 - Companies: Eleka, Elhuyar, Imaxin Software
- Project gets funded to develop two systems:
 - Apertium ($es \leftrightarrow ca$, $es \leftrightarrow gl$)
 - Matxin ($es \rightarrow eu$)
- Project later renamed *Opentrad*
 - Name sometimes erroneously used to refer to the systems Apertium and Matxin.
 - Now used by two companies as a trademark to offer services on Apertium.

Technology/1

- Apertium not built from scratch.
- Complete FOS re-specification, rewriting and extension of closed-source systems built by Transducens at the UA:
 - **interNOSTRUM** (`interNOSTRUM.com`, `es↔ca`)
 - **Tradutor Universia** (`tradutor.universia.net`, `es↔pt`)
- Linguistic data for `es↔ca` and `es↔gl` built combining in-house resources with existing FOS data (e.g., in Freeling).

Technology/2

- Apertium 1.0: designed to treat with closely-related language pairs ($es \leftrightarrow ca$, $es \leftrightarrow pt$, etc.)
- Apertium 2.0: extended to treat less-related languages ($fr \leftrightarrow ca$, then $en \leftrightarrow ca$, etc.)
- Apertium 3.0: Unicode compliance to deal with any written language in the world

A conservative design? /1

Most of the design of Apertium is rather “conservative”:

- **No “rocket science”**: tested and established techniques and technologies: finite state transducers, finite-state pattern matching, hidden Markov models.
- **High-school linguistics**: representation based on well-known and widely-accepted linguistic concepts (morphology, parts of speech and just a little bit of syntax).

A conservative design? /2

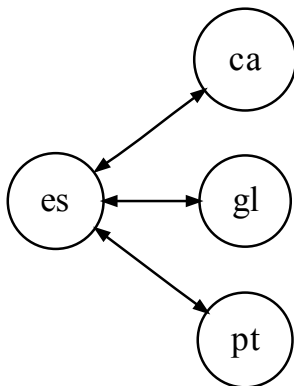
- **Good-old 70's Unix style:** modularity achieved “the Unix way”:
 - little programs “that do one thing and do it well” (McIlroy 1978)
 - “simple parts that are connected by clean interfaces” (Raymond 2004)
 - text, pipes & filtersfor easy diagnosis, extension, to build *frankensteins*, etc.

Development of language pairs as a driving force for innovation/1

On top of the original two language pairs, $es \leftrightarrow ca$ and $es \leftrightarrow gl$, pairs built outside the initial consortium:

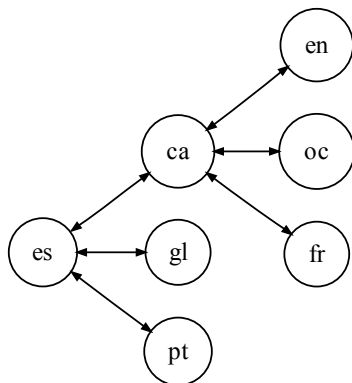
- as research projects by students: $fr \leftrightarrow ca$, $pt \leftrightarrow ca$, $cy \rightarrow en$, etc.
- with other research groups: $en \leftrightarrow ca$, $en \leftrightarrow es$, $ro \rightarrow es$, etc.
- within the Apertium community: $en \rightarrow eo$, $nn \leftrightarrow nb$, $sv \leftrightarrow da$, $br \rightarrow fr$, etc.
- by companies involved in development: $fr \leftrightarrow es$, $pt \leftrightarrow gl$, $oc \leftrightarrow ca$, etc.

Development of language pairs as a driving force for innovation/2



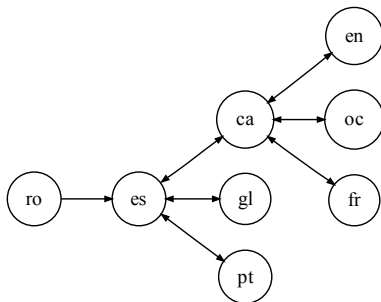
2005

Development of language pairs as a driving force for innovation/3



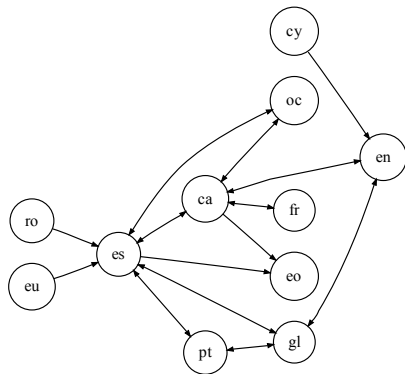
2006

Development of language pairs as a driving force for innovation/4



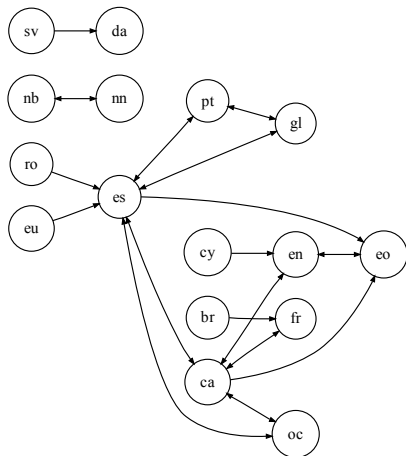
2007

Development of language pairs as a driving force for innovation/5



2008

Development of language pairs as a driving force for innovation/6



2009

Development of language pairs as a driving force for innovation/7

Language-pair development has also motivated changes in the platform:

- three-stage structural transfer was introduced to deal with $en \leftrightarrow ca$
- multi-stage (> 3) structural transfer for $eo \rightarrow en$
- integration of VISL constraint grammar, motivated by
 - FOS grammars for no (nn, nb) and the Sámi languages
 - their utility to deal with the morphology of Celtic languages.

Build on top of word-for-word translation/1

To generate translations which are

- reasonably intelligible and
- easy to correct (*postedit*)

between related languages such as *es-ca*, *es-pt*, *nn-nb*, *ga-gd*, one can just augment *word for word* translation with

- robust lexical processing (including multi-word units)
- lexical categorial disambiguation (part-of-speech tagging)
- local structural processing based on simple and well-formulated rules for frequent structural transformations (reordering, agreement)

Build on top of word-for-word translation /2

For harder, not so related, language pairs:

- It should be possible to build on that simple model.
- It should be possible to generalize its concepts so that linguistic complexity is kept as low as possible.

Clear and effective separation of translation engine and language-pair data/1

- It should be possible to generate the whole system from linguistic data (monolingual and bilingual dictionaries, grammar rules) specified in a declarative way.
- This information, i.e.,
 - (language-independent) rules to treat text formats
 - specification of the part-of-speech tagger
 - morphological and bilingual dictionaries and dictionaries of orthographical transformation rules
 - structural transfer rules

should be provided in an interoperable format ⇒ XML.

Clear and effective separation of translation engine and language-pair data/2

- It should be possible to have a single generic (language-independent) engine reading language-pair data (“separation of algorithms and data”).
- Language-pair data should be preprocessed so that the system is fast ($>10,000$ words per second) and compact; for example, lexical transformations are performed by minimized finite-state transducers (FSTs).

Apertium as free/open-source software /1

Reasons for the development of Apertium as free/open-source software:

- To give everyone free, unlimited access to the best possible machine-translation technologies.
- To establish a modular, documented, open platform for shallow-transfer machine translation and other human language processing tasks.
- To favour the interchange and reuse of existing linguistic data.
- To make integration with other free/open-source technologies easier.

Apertium as free/open-source software /2

More reasons for the development of Apertium as free/open-source software:

- To benefit from collaborative development
 - of the machine translation engine
 - of language-pair data for currently existing or new language pairs

from industries, academia and independent developers.

- To help shift MT business from the obsolescent *licence-centered* model to a *service-centered* model.
- To radically guarantee the *reproducibility* of machine translation and natural language processing research.
- Because public research investments must be made available to the public.

Reasons for the use of copyleft

What is *copyleft*?

- Obviously a play on the word *copyright*.
- Copyleft, when added to a free license, means that modifications have to be distributed with the same (copylefted) license.

Apertium chose *copylefted* free/open-source licences from the very beginning:

- To enable communities of programmers to build a machine translation *commons*, a shared body of FOS machine translation software and data
- while allowing for commercial uses (more later).

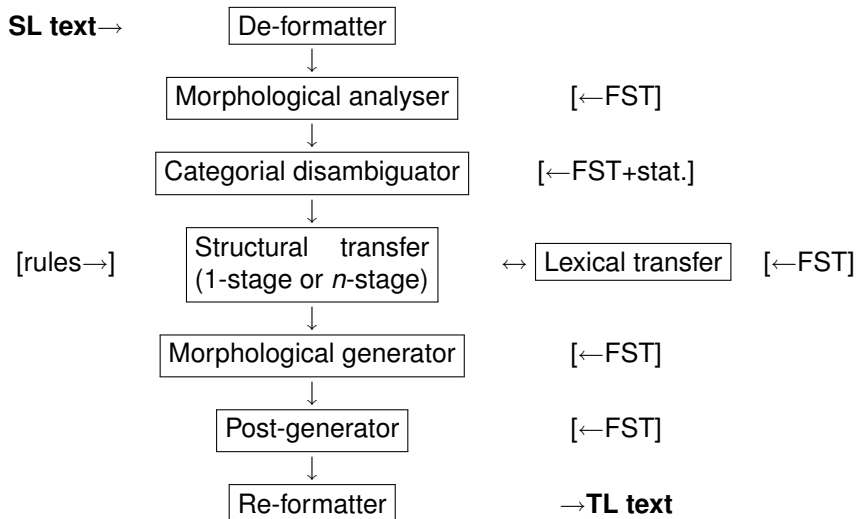
The license chosen was the GNU General Public License (GPL)

The Apertium platform

Apertium is a free/open-source machine translation platform (<http://www.apertium.org>) providing:

- 1 A free/open-source modular shallow-transfer machine translation **engine** with:
 - text format management
 - finite-state lexical processing
 - statistical lexical disambiguation
 - shallow transfer based on finite-state pattern matching
- 2 Free/open-source **linguistic data** in well-specified XML formats for a variety of language pairs
- 3 Free/open-source tools: **compilers** to turn linguistic data into a fast and compact form used by the engine and software to learn disambiguation or structural transfer rules.

Architecture/1



Architecture/2

XML linguistic data compiled for speed:

- Lexical information (SL and TL morphological dictionaries, SL–TL bilingual dictionaries, post-generation rules) → finite-state transducers (FST).
- Patterns identifying the left-hand side of structural transfer rules → finite-state pattern matchers
- Disambiguation rules and probabilities obtained from text corpora → hidden Markov models (HMM)
- etc.

The Apertium community/1

Not the ideal community development situation, but close. In addition to the original (funded) developers, a community (instigated by Francis Tyers) formed around the platform.

- More than 100 developers in `sourceforge.net/projects/apertium/`, many outside the original group (thank you all!)
- Code updated very frequently: hundreds of monthly SVN commits
- A collectively-maintained *wiki* shows the current development and tips for people building new language pairs or code.

The Apertium community/2

- Externally developed tools and code:
 - a graphical user interface `apertium-tolk`, and the related diagnostic tool `apertium-view` and `apertium-view`
 - plugins for OpenOffice.org, the Pidgin (previously Gaim) messaging program, for the Wordpress content management system, the Virtaal translation software, the Jubler film-subtitling application, etc.
 - A standalone film subtitling application (`apertium-subtitles`)
 - Dictionaries adapted to mobile phones and handhelds (`tinylex`)
 - Windows ports.
- Many people gather and interact in the `#apertium` IRC channel (at `freenode.net`).
- Stable packages ported to Debian GNU/Linux (and therefore to Ubuntu and gNewSense).

Research/1

- Apertium is also a MT research platform.
- **New code** (`apertium-tagger-training-tools`, `apertium-transfer-tools`) or **language-pair data** have often been released simultaneously to research publications.
- The research undertaken has even produced a PhD thesis (Felipe Sánchez-Martínez 2008) and four master's theses (Gema Ramírez-Sánchez, Carme Armentano-Oller, Francis M. Tyers, Ángel Seoane).
- A survey of published research may be found in the paper.
- Apertium has also been used to obtain resources for other MT systems.

Research/2

Access to FOS software like Apertium

- guarantees the reproducibility of all of the above experiments
- “lowers the bar for entry to your project for new colleagues” (Pedersen 2008)

Research/3

Together with other FOS machine translation software, such as

- the Giza++ statistical aligner,
- the Moses statistical MT engine,
- the IRSTLM language-model toolkit,
- the Cunei example-based MT platform,
- the Anymalign aligner,
- the Matxin MT system for Basque, and
- the OpenLogos MT system,

Apertium contributes to the reproducibility and the advancement of MT research and experiments.

Business/1

Companies in the initial consortium offer services based on Apertium:

- Eleka Ingeniaritza Linguistikoa
- imaxin|Software

Prompsit Language Engineering, started in 2006:

- works almost exclusively on Apertium
- currently one of the main developers of the platform

Business/2

Traditional model: Companies holding copyright of an Apertium component can also sell non-free, closed-source software based on it.

New model: Any company can sell services around the Apertium platform (allowed by the GPL licence):

- installing and supporting translation servers
- maintaining and extending language-pair data for a particular application
- integrating Apertium in multilingual documentation management systems

Business/3

Advantages of the FOS business model:

- From technological dependency (*vendor lock-in*) to technological partnership.
- Customers may prefer the social image of a company developing FOSS software like Apertium.
- The FOSS community improves your business by improving Apertium.

Vulnerability:

- Any company can offer the same services on Apertium (however, companies deeply involved in Apertium development have a definite competitive edge over those who are not that involved).

Business/4

Two commercial success stories:

- imaxin|software and the Universidade de Vigo integrated Apertium in the online edition of one of the most known Galician newspapers originally published only in *es*, *La Voz de Galicia*; (<http://lavozdegalicia.es/>), so that it has now a *gl* edition.
- The company Taller Digital (owned by the Universitat d'Alacant) hired Prompsit Language Engineering to build the Catalan Government's official MT systems for *es*↔*oc* and *ca*↔*oc* on Apertium:
<http://traductor.gencat.cat/>

The 2009 Google Summer of Code

Apertium was selected to participate as a mentoring organisation in the 2009 Google Summer of Code. Successful projects:

- two new language pairs: $nn \leftrightarrow nb$ and $sv \leftrightarrow da$
- a morphological analyser for bn
- an improved part-of-speech tagger
- a web-service infrastructure
- porting of the lexical component to Java
- hybridising Apertium with other systems

Ongoing work

- Universidá d'Uviéu: $es \leftrightarrow ast$
- University of Reykjavík: $is \leftrightarrow en$
- Universitat d'Alacant and Prompsit: $es \leftrightarrow it$
- University of Tromsø: $sme \leftrightarrow smj$

Lots of work ahead: known limitations

- No successful, general-purpose lexical selection for polysemic words
- No deep (parse-tree-based) structural transfer, needed for syntactically divergent language pairs
- Current lexical processing not adequate for agglutinative languages or languages with non-catenative morphology.
- The representation of morphological inflection is still too low-level.
- No support to segment long compound words (de: *Kontaktlinse**verträglichkeitstest*)
- Apertium is a *transfer* system: generating a new pair involves the creation of explicit bilingual resources.
`apertium-dixtools` helps build pair *A–B* from *A–C* and *C–B*, but task is far from trivial.

Lots of work ahead: an inconsistency

An inconsistency in the way Apertium is developed:

- Apertium is FOS software, but
- its development is hosted in Sourceforge, which includes non-free closed-source software

We are considering:

- Migration to a completely FOS software development platform (Gna!, Savannah, etc.), or
- self-hosting.

No decision has been made yet.

Funding

Apertium has been funded by

- The Ministry of Industry, Tourism and Commerce of Spain (also, the Ministries of Education and Science and of Science and Technology of Spain)
- The Secretariat for Technology and the Information Society of the Government of Catalonia
- The Ministry of Foreign Affairs of Romania
- The Universitat d'Alacant
- The Ofis ar Brezhoneg (Breton Language Board)
- Google (Google Summer of Code 2009) scholarships
- Companies: Prompsit Language Engineering, ABC Enciklopedioj, Eleka Ingeniartiza Linguistikoa, imaxin|software, etc.

License

This work may be distributed under the terms of

- the Creative Commons Attribution–Share Alike license:

`http:`

`//creativecommons.org/licenses/by-sa/3.0/`

- the GNU GPL v. 3.0 License:

`http://www.gnu.org/licenses/gpl.html`

Dual license! E-mail me to get the sources: `mlf@ua.es`