

apertium-cy: A collaboratively-developed RBMT system for Welsh to English

Francis M. Tyers^{1,2} and Kevin Donnelly³

¹Departament de Llenguatges i Sistemes Informàtics, Universitat d'Alacant, E-03071 Alacant (Spain)

²Prompsit Language Engineering, Avinguda Sant Francesc 74, 1-L. E-03195 L'Altet - Elx (Spain)

³Llanfairpwll, Ynys Môn, LL61 6UX (Wales)

16th January 2009

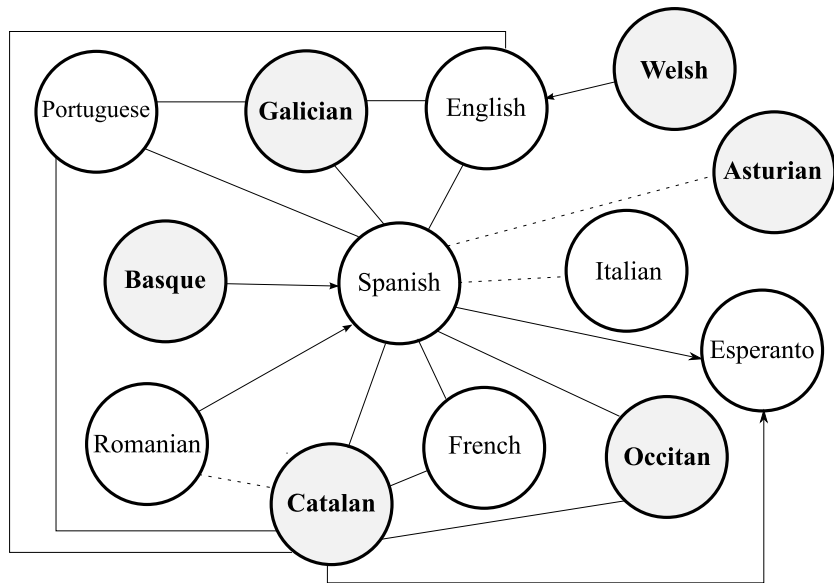
Contents

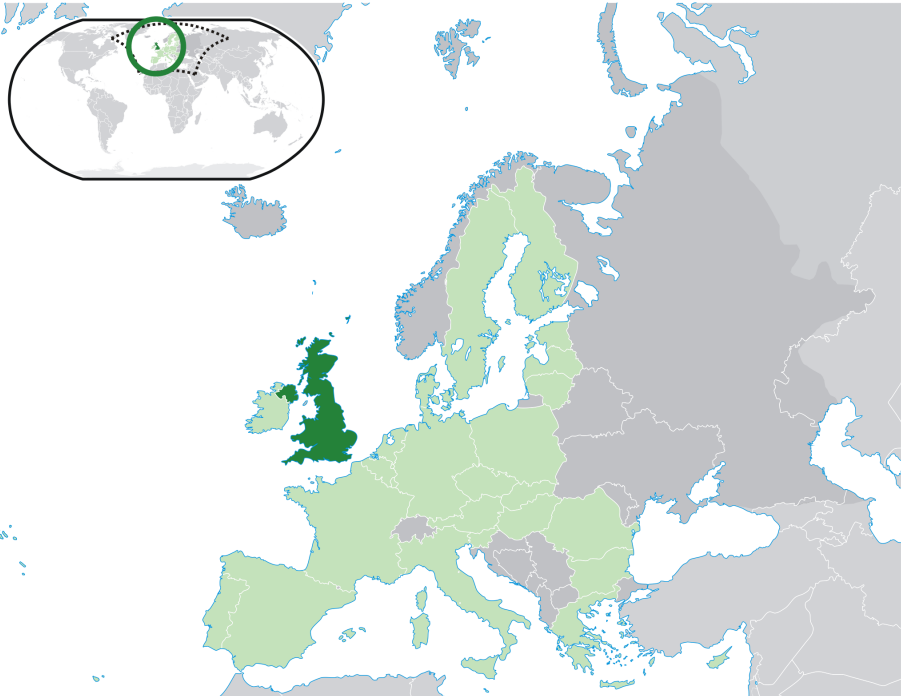
- 1 Introduction
- 2 Development
- 3 Evaluation
- 4 Discussion

So what is Apertium?

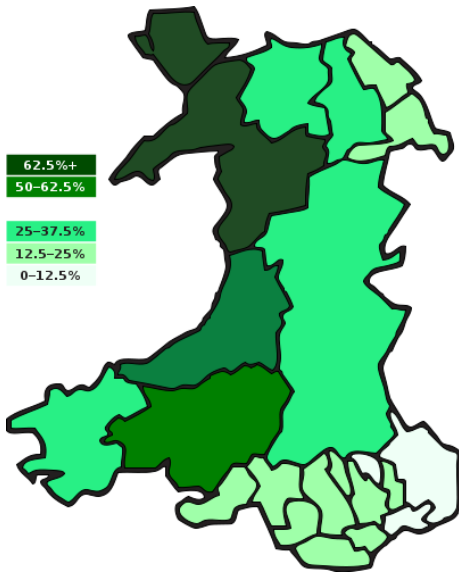
- GPL-licensed platform for machine translation
- Modular – made up of stand-alone programs which communicate using Unix pipes
- Developed by universities, companies and independent developers
- 17 available “stable” language pairs
- More in development
 - And language data in development for many other languages. . . Breton, Icelandic, Hindi, . . .

Language pairs







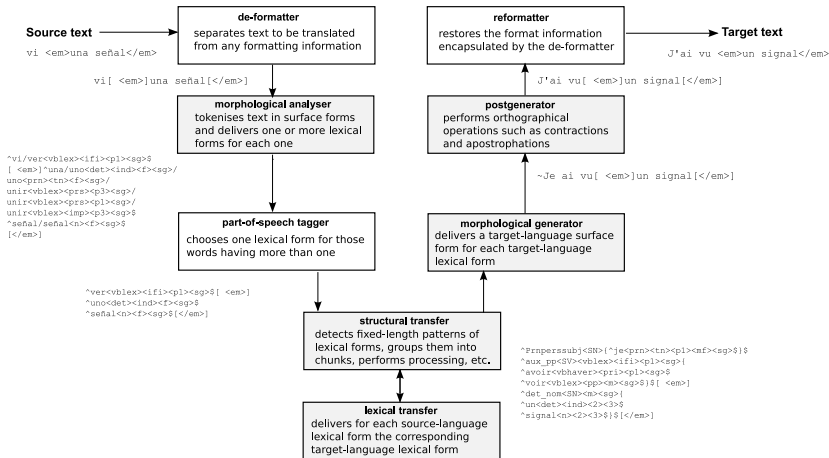




Apertium is a shallow-transfer MT system meaning development consists of:

- Morphological dictionaries (analysis / generation)
- Disambiguation rules and training statistical tagger (including optional target-language training)
- Bilingual dictionary (lexical transfer)
- Shallow syntactic transfer rules
 - Local re-ordering ($\text{nom adj} \rightarrow \text{adj nom}$)
 - Chunking ($\text{adj adj nom} \rightarrow \text{SN}[\text{adj adj nom}]$)
 - Insertions, deletions and substitutions of lexical units and chunks

Translation model



We were able to directly use:

- English morphological dictionary from the `apertium-en-ca` language pair

And the following after conversion:

- English–Welsh bilingual dictionary from Eurfa¹
- Monolingual Welsh dictionary from Eurfa and Konjugator²
- Some rules from the English–Spanish pair

¹<http://www.eurfa.org.uk>

²<http://www.konjugator.org.uk>

Development was performed collaboratively over the internet.

- Wiki and E-mail – For describing transfer and disambiguation rules
- SVN – For version control
- IRC and E-mail – For technical support

(1.3.39) Periphrastic tenses with "bod"

[\[editar\]](#)

Re comments in 1.3.32 above:

For Welsh pattern "VBSEER.pres + subject + [qualifiers] + yn + verb.infin"
 output English "subject + [qualifiers] + is/are + verb + -ing"

mae'r bachgen yn mynd -> The boy is going

Fine - the equivalent of 1.3.8 above, and it seems to be catered for already.

For Welsh pattern "VBSEER.pii + subject + [qualifiers] + yn + verb.infin"
 output English "subject + [qualifiers] + was/were + verb + -ing"

roedd y bachgen yn mynd -> The boy was going

Fine again.

For Welsh pattern "VBSEER.pres + subject + [qualifiers] + wedi + verb.infin"
 output English "subject + [qualifiers] + has/have + verb.pp"

mae'r bachgen wedi mynd -> *The boy is after go - the boy has gone

Done. - Francis Tyers

Done



Proxy: None

Adblock

```
<when>
  <test>
    <and>
      <equal>
        <clip pos="1" part="tipus">
          <lit-tag v="noun"/>
        </equal>
      <equal>
        <clip pos="2" part="tipus"/>
          <lit v="det.def"/>
        </equal>
      <equal>
        <clip pos="3" part="tipus"/>
          <lit-tag v="noun"/>
        </equal>
    </and>
  </test>
  <out>
    <clip pos="2" part="whole"/>
    <b/>
    <clip pos="1" part="whole"/>
    <b/>
    <lit v="of"/>
    <b/>
    <clip pos="2" part="whole"/>
    <b/>
    <clip pos="3" part="whole"/>
  </out>
</when>
```

Transfer rules

“In Welsh, if a noun is followed by the definite article followed by another noun, then output in English definite article, then the first noun, then the preposition ‘of’ then the definite article followed by the second noun.”

	Number	Coverage
Lexicon	10,994	90.1%
Disambiguation	56	-
Chunk	84	-
Inter-chunk	36	-

Table: Statistics from current SVN revision #8140

This page allows you to test the pre-alpha (from SVN) versions of Apertium data. These data will probably not work more often than not. Often the data are woefully incomplete and they are not supported in any way.

Roedd y cwningod yn hapus

Welsh → English

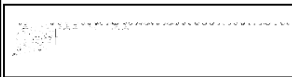
Mark unknown words

Print intermediate representation

Translate

* — Word not found in source.
— Word found but no inflection.
@ — Word not found in target.

^ \$ — Lexical unit start/stop
{ } — Chunk start/stop



Roedd y cwningod yn hapus

```

^Roedd/Bod<vbser><pii><p3><sg>$ ^y/yr<vpart><aff>/yr<det><def><sp>$          lt-proc
^cwningod/cwningen<n><f><pl>$ ^yn/yn<pr>$ ^hapus/hapus<adj>$

^Roedd/Bod<vbser><pii><p3><sg>$ ^y/yr<det><def><sp>$ ^cwningod/cwningen<n>          cg-proc
<f><pl>$ ^yn/yn<pr>$ ^hapus/hapus<adj>$

^Bod<vbser><pii><p3><sg>$ ^yr<det><def><sp>$ ^cwningen<n><f><pl>$          apertium-tagger
^yn<pr>$ ^hapus<adj>$

^Bod<vbser><pii><p3><sg>$ ^yr<det><def><sp>$ ^cwningen<n><f><pl>$          apertium-pretransfer
^yn<pr>$ ^hapus<adj>$

^Verbcj<SV><vbser><past><p3><sg>{^be<vbser><3><4><5>$}$ ^det_nom<SN>          apertium-transfer
<pl>{^the<det><def><sp>$ ^rabbit<n><2>$}$ ^adj<SA>{^happy<adj>
<sint>$}$

^Det_nom<SN><pl>{^the<det><def><sp>$ ^rabbit<n><2>$}$ ^verbcj<SV>          apertium-interchunk
<vbser><past><p3><pl>{^be<vbser><3><4><5>$}$ ^adj<SA>{^happy<adj>
<sint>$}$

^The<det><def><sp>$ ^rabbit<n><pl>$ ^be<vbser><past><p3><pl>$          apertium-postchunk
^happy<adj><sint>$

The rabbits were happy          lt-proc

The rabbits were happy          lt-proc

The rabbits were happy

```


We took two main approaches to evaluation.

- **Quantitative** – To be comparable with other systems, and provide a useful “at a glance” measure of quality.
- **Qualitative** – To give a better idea of where the strengths and weaknesses of the system are.

We also made some tests as to how the output compared with other available systems.

Two corpora were used for quantitative evaluation.

- **Wikipedia** – 318 sentences (5,492 words) selected at random from the Welsh Wikipedia, translated, post-edited then WER, PER, BLEU calculated against post-edited translations.
- **PNAW**³ – 50,000 sentences selected at random from a bilingual parallel corpus translated and then WER, PER, BLEU calculated against reference translations.

³Proceedings of the National Assembly for Wales

	WER	PER	BLEU
Wikipedia true-case	55.78	30.59	30.70
Wikipedia lowercased	53.40	27.22	32.21
PNAW true-case	65.99	35.44	15.12
PNAW lowercased	64.94	34.35	15.68

Table: WER, PER and BLEU metrics for the two corpora

In general terms:

- Short sentences or long sentences made up of sequential parts work reasonably well
- Sentences with marked formations, long multiword units or subordinate clauses often come across “mangled”
 - Pam mae'r bocs yn wag? → *Why the box is empty?
 - Aeth y dyn i ffonio'r heddlu pan welodd y ddamwain. → *The man went to phone the police when the accident saw.
 - Gwasanaeth Tân ac Achub Canolbarth a Gorllewin Cymru → *Fire service and Save Central region and Wales West
- Word-choice is often “unusual”, but rarely “ridiculous”
 - Llai na'r cyfradd chwyddiant → *Smaller than the rate inflation

We also did some brief comparative tests with InterTran⁴ and default Moses configuration trained on the PNAW corpus.⁵

- or *Roedd y Comisiwn yn ymchwilio i'r honiadau bod yr AS wedi methu datgan £103,000 o roddion.*
- in He was the Commission crookedly ymchwiliad I ' group claims be he drives ACE has failed declare he gifts.
- mo The comission to investigate the allegation that the MP has failed to declare £103,000 of roddion.
- ap **The Commission was investigating the allegations that the MP has failed to declare £103,000 gifts.**

Moses and Apertium can both be tested side-by-side online.⁶

⁴<http://www.tranexp.com/>

⁵<http://xixona.dlsi.ua.es/corpora/>

⁶<http://elx.dlsi.ua.es/~fran/welsh/>

Newyddion

Diweddarwyd: Dydd Mawrth, 18 Tachwedd 2008, 12:39 GMT

Gwleidyddiaeth

Fersiwn brintiadwy

Chwaraeon

Cynllun i greu cymuned Aelig

Y Tywydd

Mae ymgyrch wedi dechrau i geisio creu cymuned Aelig yn Ucheldiroedd Yr Alban.

Lleol

Pobl ifanc

Dysgwyr

Dywedodd ymgyrchwyr y byddai cymuned o fwy na 1000 o siaradwyr Gaeleg yn "hwb enfawr" i'r iaith.



Briard y cynllun yw cryfhau'r defnydd o'r iaith Aelig

Mae nifer y siaradwyr wedi lleihau'n fawr yn y blynyddoedd diwethaf, cyfanswm o 58,652 yn ôl Cyfrifiad 2001, a'r mwyafrif yn byw yn yr Ucheldiroedd a'r ynysoedd.

Eisoes mae 30 o bobl wedi ymrwmo i fyw yn y pentref ac mae'r ymgyrchwyr am gynnal cyfarfod yn Inverness ym mis Rhagfyr.

Dim 'cyfle cymdeithasol'

Mae dysgu'r Aelig ar gynydd mewn vsqolion.

HEFYD

- Sianel Aelig newydd yn dechrau 19 Medi 08 | Newyddion
- Gaeleg: Newyddion ar wefan 28 Chwef 08 | Newyddion

CYSYLLTIADAU RHYNGRWYD:

- Alba (Gaeleg / Saesneg)
- Llywodraeth yr Alban (Saesneg)
- Ymddiriedolaeth y BBC yn cymeradwyo'r Gwasanaeth Digidol Gaeleg yn amodol

Dyw'r BBC ddim yn gyfrifol am gynnwys safleoedd rhyngwyr allanol

PRIF STRAEON NEWYDDION

- Dwr: Cyflenwadau yn ôl
- 'Gwledigaeth' darpar arweinydd?
- Cnydau'n 'hwb i gefn gwlad'

[RSS](#) | Beth yw RSS?

BBC Cymru'r Byd

Chwaraeon

Y Tywydd

Dysgu Cymraeg
Learn Welsh

Lleol i Mi

VOCAB

OFF / I FFWRDD

»Tum ON
Troï YMLAEN»What is VOCAB?
Beth yw GEIRFA?

Plan to create Scots Gaelic community

campaign has begun to try create Scots Gaelic community in Scotland Highlands.



the Intention of the plan is strengthen the material of the language Scots Gaelic

campaigners Said community would be he bigger than 1000 speakers Scots Gaelic in "enormous push" to the language.

the number of the speakers have diminished big in the last years, total of 58,652 according to Census 2001, and the majority in life in the Highlands and the islands.

Already 30 peoples have committed oneself to a life in the village and the campaigners are going to hold a meeting in *Inverness in December.

Not 'a social opportunity'

ALSO

- ▶ [new Scots Gaelic Channel beginning](#)
19 Septembers 08 News
- ▶ [a Scots Gaelic: News on a website](#)
28 *Chwef 08 News

INTERNET ASSOCIATIONS:

- ▶ [Alba \(English Scots Gaelic\)](#)
- ▶ [Scotland Government \(English\)](#)
- ▶ [the Trust of the BBC approving *r Scots Gaelic Digital Service conditional](#)

the BBC Are not responsible for positions content an external internet

PRINCIPAL STORIES NEWS

- ▶ [Heap: Supplies according to](#)
- ▶ ['leader' designate Vision?](#)
- ▶ [Crops in 'push to countryside'](#)



Beth yw RSS?

- ▶*Turn ON Tum FORWARD
- ▶*What a lower VOCAB? What



Why not corpus-based MT?

But wouldn't it be quicker to use corpus-based MT?

- No wide-coverage freely available corpus of Welsh–English
- Little chance of finding one – most text is not free
- In this case existing GPL linguistic data was available
- Interested linguist – use available talent

Creating an RBMT system also involves creating useful linguistic tools which can be used by other approaches to MT (e.g. SMT) and other linguistic software.

Newyddion

Diweddarwyd : Dydd Mawrth ,18 Tachwedd 2008 ,12 :39 GMT

Gwleidyddiaeth

Fersiwn brinteadwy

Chwaraeon

Cynllun i greu cymuned Aelig

Y Tywydd

Mae ymgyrch wedi dechrau i geisio creu cymuned Aelig yn Ucheldiroedd Yr Alban .

Lleol

Pobl ifanc

Dysgwyr

BBC Cymru 'r Byd

Chwaraeon

Y Tywydd

Dysgu Cymraeg
Learn Welsh

Lleol i Mi

VOCAB

OFF I FFRWDD

Tum OH

Troï YMLAEN

What is VOCAB

?Beth yw GEIRFA ?

Dyweddodd ymgyrchwyr y byddai cymuned o fwy na 1000 o siaradwyr Gaeleg yn hwb enfawr i 'r iaith .

Mae nifer y siaradwyr wedi lleihau 'n fawr yn y blynyddoedd diwethaf ,cyfanswm o 58,652 yn ôl Cyfrifiad 2001 ,a'r mwyafrif yn byw yn yr Ucheldiroedd .

Eisoes mae 30 o bobl wedi ymrwmo i fyw yn y pentref ac mae 'r ymgyrchwyr am gynnal cyfarfod yn Inverness ym mis Rhagfyr .

Dim cyfle cymdeithasol

Mae dysgu 'r Aelig ar gynnwdd mewn ysgolion .



Briaid y cynllun yw cryfhau 'r defnydd o'r iaith Aelig

HEFYD

- Sianel Aelig newydd yn dechrau 19 Medi 08 Newyddion
- Gaeleg :Newyddion ar wefan 28 Chwef 08 Newyddion

CYSYLLTIADAU RHYNGRWYD :

- Alba (Gaeleg Saesneg)
- Llywodraeth yr Alban (Saesneg)
- Ymddiriedolaeth y BBC yn cymeradwyo'r Gwasanaeth Digidol Gaeleg yn amodol

Dyw 'r BBC ddim yn gyfrifol am gynnwys safleoedd rhyngwyd

PRIF STRAEON NEWYDDION

- Dwr :Cyflenwadau yn ôl
- Gwledigaeth darpar arweinydd ?
- Cnydau 'n hwb i gefn gwlad

[RSS](#) | Beth yw RSS?

For the Apertium platform in general:

- Implementation of a full parser
- Improvement of lexical selection
- Expansion of number of language pairs ... working on es-ast, es-it, br-fr
- Increase ease of contribution

For `apertium-cy` in particular:

- Reverse the direction (`cy`→`en`) – the pitfalls of Scymraeg
- Increase vocabulary coverage
- Other translators with Welsh – e.g. Spanish (Patagonia), Breton (related)

Diolch / Thanks!