



Canolfan ESRC Centre
dros Ymchwil i Ddwyieithrwydd for Research
on Bilingualism

Mining corpora using the Bangor Autoglosser

Kevin Donnelly

ESRC Centre for Research on Bilingualism
Bangor, Wales



Cyngor Cyllido Addysg
Uwch Cymru
Higher Education Funding
Council for Wales

hefcw



Outline

- Briefly **introduce** the autoglosser
- Give examples of **using its output** to mine the ESRC corpora in various ways
- Seek **input** re other possible uses

The Bangor Autoglosser

Bangor corpora

	<i>Chats</i>	<i>Hours</i>	<i>Words</i>	<i>Date</i>
Welsh-English (Siarad)	69	40	456k	2009
Welsh-Spanish (Patagonia)	32	20	161k	2011
Spanish-English (Miami)	31	20	126k	2011
	132	80	743k	

All available under the GPL.

- Part-of-speech tags
- **walks:** `walk.V.3S.PRES`
- Allow non-fluent speakers to parse the conversation
- Labour-intensive, can be inconsistent
- Difficult to add further tags later
- **Solution:** use a computer to generate the glosses!

Sample file

- **Demo:** import and autogloss a file containing two utterances

Benefits

- Handles multiple languages simultaneously
- Handles broken text
- Highly configurable
- Reasonably fast: **1000 words a minute**
- Accuracy **98%+**
- First tagger for Welsh - other uses possible
- Open license (**GPL**) - reproducible science

Mining the corpora

Data analysis

- Text now in a database
- **Demo:** data grid interface
- Easy basic queries (eg) how many words per speaker, or per language tag
- **Demo:** browse columns
- **Demo:** query for words not in the dictionary
- Effective and flexible typo correction

Presentation of output

- Typeset files using \LaTeX gloss alignment packages
- **ExPex** (John Frampton), **Xscribe** (Natalie Weber)
- **Demo**: sample file pdf
- Easier to read and review the text
- Easy to generate indexes and concordances
- **Demo**: pages from two summaries

Sequence selection

- Surface and glossed tiers in transcription file are linked at **utterance** level
- Database allows linking at **word** level
- Words can be selected based on the gloss
- Sub-select based on position (before or after the target word)
- Filter by gloss if necessary

Sequences

- Det+N+Adj (Spanish/English)
- **the fair estúpido**
- **el tremendo tip**
- About 26 sequences out of 126,000 words
- Non-sandwich (non-nested) possessive sequences (Welsh/English)
- **cath fi** instead of **fy nghath [i]**
- 340 1S instances – only 5 are sandwich-type
- **Handout**

- Can be done more rigorously
- Order the utterances
- Select the utterance at the 25%, 50% and 75% points
- Sub-select ...
 - ▶ 5 utterances before and after
 - ▶ by each speaker
 - ▶ containing a finite verb

Clause delineation

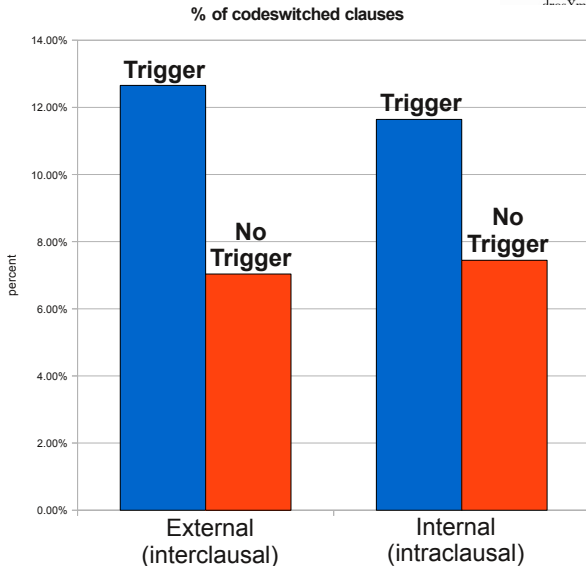
- Ideal? – build dependency trees
- Surprisingly good proxy can be achieved by using
 - ▶ verbs in Welsh
 - ▶ conjunctions in Spanish
 - ▶ subject pronouns in English
- Welsh splitter has error rates of 2-5%
- **Demo**: clause-split the simple file

Data aggregation

- Outputs can be combined with metadata
- Sociolinguistic analysis
- Are particular structures correlated with speakers' background?
- **Demo:** spreadsheet

- Are codeswitches more frequent in the vicinity of triggers (cognates)?
 - ▶ names, proper nouns
 - ▶ identical words: **car**
 - ▶ related words: **perswadio/persuade**
- Split into clauses
- Group into speech-turns
- Mark triggers and count them
- Count codeswitches between and within clauses

Results



Timespend

- Original 2,500-word Dutch text took 50+ hours to analyse
- 456,000-word Siarad corpus analysed in 4 hours

Running text

- Input written formal text
(instead of transcribed colloquial dialogue)
- Historical Corpus of Welsh (420,000 words)
- *Berr Hanes* (James Groniosaw) – 1779

Results

- 1,200-word sample passage, glossed in 2 minutes
- 57% of words correctly glossed
- 3% wrongly glossed
- 40% had two or more glosses offered
- Grammar rules need to be extended to cope with this type of text.

Comments?
Suggestions?

<http://bangortalk.org.uk>