

NLP tools for historical lexicography

Kevin Scannell
Cadhan Aonair
20 March 2025

What is historical lexicography?

banter, *n.*

Text size: A A

View as: [Outline](#) | [Full entry](#)

Quotations: [Show all](#) | [Hide all](#) Keywords: [On](#) | [Off](#)

Pronunciation: Brit. ▶ /'bʌntə/, U.S. ▶ /'bæn(t)ər/

Frequency (in current use): ●●●●●●●●

Etymology: ... [\(Show More\)](#)

1. Wanton nonsense talked in ridicule of a subject or person; *hence*, humorous ridicule generally; (now usually) good-humoured raillery, pleasantry.

[Thesaurus »](#)

1702 *Eng. Theophrastus* 232 The ordinary reasons of War and Peace, are very little better than Banter and Paradox.

1705 S. WHATELY in W. S. Perry *Hist. Coll. Amer. Colonial Church: Virginia* (1870) I. 172 I know no better way of answering bombast, than by banter.

1710 SWIFT *Tale of Tub* (ed. 5) Apol. sig. A8 Peter's Banter (as he calls it in his Alsatia Phrase) upon Transubstantiation.

1843 DICKENS *Martin Chuzzlewit* (1844) xxiv. 298 She took it for banter, and giggled excessively.

1880 L. STEPHEN *Alexander Pope* v. 113 Gay..had an illimitable flow of good-tempered banter.

Source: <https://oed.hertford.ox.ac.uk/oed-editions/oed-online/>

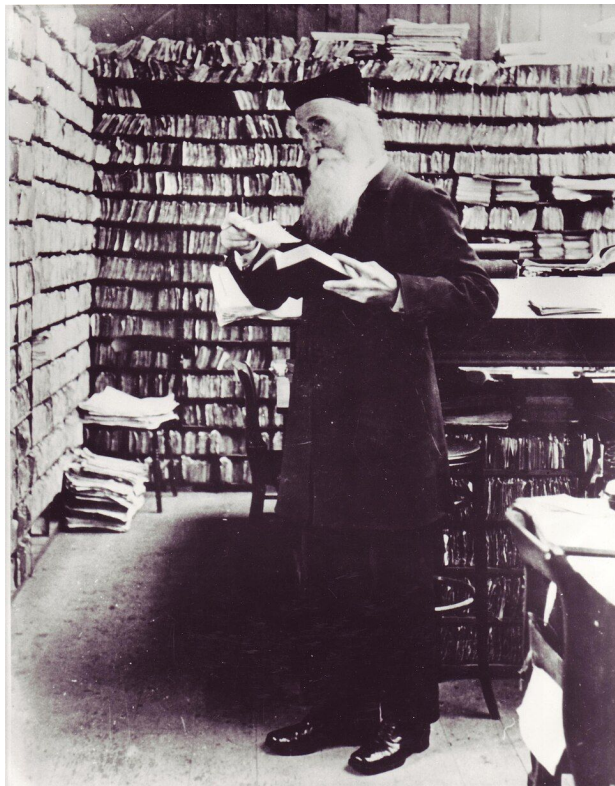
Irish language lexicography

- Virtually all dictionaries are bilingual English-Irish or Irish-English
- New English-Irish dictionary published 2020, print and online (previous 1957)
- New Irish-English dictionary hopefully by 2030 (current 1977)
- **No** full-scale monolingual dictionary
- **No** etymological dictionary
- **No** historical dictionary

Foclóir Stairiúil na Gaeilge (Historical Dictionary of Irish)

- Project of the Royal Irish Academy in Dublin
- Partial funding from the Department of the Gaeltacht in Ireland
- Intended to cover the period 1600–2000 (“modern Irish”)
- Project has existed since the 1970s

How are dictionaries compiled?



CONCORDANCE

English Web 2021 (enTenTen21)

simple banter • 120,672
1.96 per million tokens • 0.0002%

Sample 200 • 200
less than 0.01 • 3.2e-7%

Sort word x

Left context KWIC Right context ↑

1	<input type="checkbox"/>	thesun.ie	en them, explaining there's been "flirty	banter	" between them.</s></s>Speaking to Or
2	<input type="checkbox"/>	newsgroove.co.u...	uded it was "friendly and good-natured	banter	".</s></s>Ballance said he felt he had "r
3	<input type="checkbox"/>	instantprint.co...	leagues, here is our top list of "banned	banter	" in the workplace:</s></s>Hints and tip:
4	<input type="checkbox"/>	redhotpawn.com	a tear to my eye.</s></s>We must say '	banter	' not 'battle'. 😊</s></s>Originally posted
5	<input type="checkbox"/>	worksmart.org.u...	nent is sometimes dismissed as being '	banter	' or just a joke.</s></s>In fact it can havi
6	<input type="checkbox"/>	wikipedia.org	was also at this time that the 'Stand-up	Banter	' performances began, forerunners of to
7	<input type="checkbox"/>	repeatfanzine.c...	</s></s>After some joyfully puerile crowd	banter	(haha- pianist sounds like penis), a firm
8	<input type="checkbox"/>	would-buy-again...	stage, they had great energy and a fun	banter	(Stone making fun of Andres' accent of
9	<input type="checkbox"/>	vdare.com	ers.</s></s>Most were likely just office	banter	, of the kind that was common thirty or f
10	<input type="checkbox"/>	skuds.org	ere were all the usual little in-jokes and	banter	, which I am no longer part of and I felt I
11	<input type="checkbox"/>	all-art.org	s handled with great flexibility; the light,	bantering	, somewhat ironic tone—later to become
12	<input type="checkbox"/>	thewalrus.ca	:</s></s>Still, the focus is Vic and Flo's	banter	, which reveals an easy rapport and lurk
13	<input type="checkbox"/>	gutenberg.ca	d his quaking voice into tones of gentle	banter	, forced himself to smile, to tweak her cl
14	<input type="checkbox"/>	kent.ac.uk	Jer children depend on proficiency with	banter	, which in turn frequently involves verba
15	<input type="checkbox"/>	casting-call.us	o About Nothing - They bicker and they	banter	, they mock and charm, they are Beatric
16	<input type="checkbox"/>	splicetoday.com	people and blogs.</s></s>There was no	banter	, no antagonizing the crowd.</s></s>Co

Irish corpora: corpas.ria.ie and corpas.ie

Tiomna Nuadh ár dTighearna agus ár Slánuightheora Íosa Críost (1602)

anaithnid

Search Term	Page	Line Number	Context
amharus	Mt28:17	3713	amharus ag cuid díobh air.
amharus	Lc4:23	6732	Agus adubhairt sé ríu, adéirthaoi rium gan amharus an seanfhocalsa,
amharus	En10:24	11661	adubhradar ris; gá fad choingéubhas tú sinn an amharus? Innis duinn
amharas	Gn2:12	13282	Agus do ghabh adhuáthmhuireachd agus amharas iád uile, ag rádh, an fear ré
amharus	Gn10:20	14434	Uime sin éirídh agus imthigh síos, agus imthigh ríu, gan amharus ar
amharus	Gn11:12	14580	Agus adubhairt an spiorad ríom dhul ríu, gan amharus ar biot do bheith
amharus	Gn21:22	15978	Ar anadhbharsoin créd doghénam? ní fuil amharus gurab éigean don
amharus	Gn21:39	16047	Agus adubhairt Pól ris, gan amharus air as lúdaighe mhísi, ó
amhras	Rm7:12	17537	Ar anadhbhasin atá an reachd gan amhras náomhtha, agus anaithe
amharas	Rm8:17	17640	amharas do Dhía, agus comhoidhreacha do Chríod: má fhuilngmíd
amhrus	Rm14:23	18217	Achd an tí ar a mbí amhrus, atá sé damanta dá nithidh sé, ar ó
amharus	1Cr7:14	18917	gan chreideamh arna náomhadh san bhfear: no no dobháidh gan amharus ar
amharus	2Cr1:8	19961	brughadh tar mhogh sinn ós díonn ar neirt, ionnus gu raibhe amharus
amhras	Gl3:29	21148	Agus más lé Chríod sibh, as sibh siól Abrahám gan amhras, agus is
amharus	1Tm2:8	23187	tógbháil suás lámh náomhtha, gan fheirg, gan amharus.
amharus	1Tm3:16	23260	Agus gan amharus as mór seicréid na diaghchad: dorinneadh Día
amhrus	Sm1:6	25153	bioth: oír an tí ar a mbí amhrus, is cosmhúil hé ré tuinn fhairge

Iomarbhágh na bhFileadh I (1604)

anaithnid

Search Term	Page	Line Number	Context
amhras	54	1495	is an amhras i dtéid sibh
n-	88	2749	bhar n-amharas orm dá mbeath
amharus	150	4494	m' amharus noch a n-iongnadh

Irish Bardic Poetry (1609)

anaithnid

Corpas Náisiúnta na Gaeilge

Tuilleadh ▾

Gaeilge

English

amhras

▾ 🔍

Abairtí

Cmhchordacht

☒
☐
Abairtí gonta chun tosaigh

☐
☒
Rogha randamach

Líon torthaí: 18,624

▾

↻ 5

↻ 4

↻ 3

↻ 2

↻ 1

⏮

⏪

1 / 466

⏩

⏭

1 ↻

2 ↻

3 ↻

4 ↻

5 ↻

ALT

in chuid eile den oileán. </s> <s> Gan

ALT

e pobal na scoile go léir. </s> <s> Gan

ALT

í fhiú a dturas? </s> <s> Ní raibh aon

ALT

y breathnú amach an fhuinneog agus

ALT

ch sásúil a bheadh ann. </s> <s> Gan

ALT

háidí mar thairgeoir ile. </s> <s> Gan

ALT

a na mbuachaillí agus na gcailíní san

ALT

a dtrioblóidí a bhí ann agus daoine in

ALT

/s> <s> Cuidíonn na beartais seo gan

ALT

a ar chúrsaí filíochta an lae inniu gan

LEA

infort. </s> <s> Cúis amháin go raibh

DOC

l anuas sa taobh seo tíre. </s> <s> Níl

LEA

nnasaí aici. </s> <s> Bhí cuid mhaith

ALT

íomh Columba sa Ghaill. </s> <s> Níl

amhras

amhras

amhras

amhras

amhras

amhras

amhras

amhras

amhras

amhras

amhras

amhras

amhras

amhras

amhras

amhras

amhras

amhras

amhras

amhras

amhras

tá an chríochdheighilt tar éis éagsúlachtaí a

ní bheadh an scoil faoi bhláth mar atá murac

air ná gurbh fhiú. </s> <s> "Tá sé tábhachtac

uirthi. </s> <s> Bhí an chuid ba mheasa de O

beifear ag plé na ceiste seo arís. </s> <s> Ní f

ní hionann na téarmaí creidmheasa a tugadh

orm. </s> <s> 'Úna Nic Fhlannchadha, an ea

ar a chéile. </s> <s> Ré dhorchá go deimhin.

ach ní réiteach iomlán ar an fhadhb iad. </s>

. </s> <s> Scriobh nóta gairid ar gach ceann

ar an bhfeadóg faoi airm a dháileadh ar a ch

ar bith ann ach go mbeidh siad chun tosaigh

ann faoin athrú mór a bhí á bheartú. </s> <s>

ar éinne ach gur scríobh sé féin na leabhair a

NLP tools: “One-click” dictionaries?

- Start with annotated corpus; minimally with headwords + POS, ideally parsed
- Produce headword list (frequency, language ID, noise filtering)
- For a given headword, we want to identify all senses and subsenses
- Fundamentally this is a (hierarchical) clustering problem
- Long-standing and challenging problem in NLP: “word sense induction”
- Sketch Engine implements clustering on contextual sentence embeddings
- GDEX algorithm for selecting example sentences (Kilgarriff et al 2008)
- Lots of post-editing needed! Splitting/lumping, crafting definitions etc.

Lexicographers ❤️ recall

- If we think of corpus query software as a search engine, recall is critical!
- Especially for historical dictionaries: aim to include each sense of each word
- Better to sift through false positives than potentially miss a word/sense/citation



character as noun 11,999,796x



modifiers of "character"
main ...
the main character
<ul style="list-style-type: none">concentrated in: games ?concentrated in: arts ?concentrated in: culture & entertainment ?
show more (1)
fictional ...
a fictional character
<ul style="list-style-type: none">concentrated in: arts ?concentrated in: culture & entertainment ?concentrated in: reference/encyclopedia ?
show more (1)
female ...
female characters
<ul style="list-style-type: none">concentrated in: arts ?concentrated in: culture & entertainment ?concentrated in: games ?
favorite ...
favorite characters
<ul style="list-style-type: none">concentrated in: culture & entertainment ?concentrated in: home & family & children ?concentrated in: games ?
cartoon ...
cartoon characters
<ul style="list-style-type: none">concentrated in: culture & entertainment ?concentrated in: home & family & children ?
playable ...
playable characters
<ul style="list-style-type: none">concentrated in: culture & entertainment ?

nouns modified by "character"
trait ...
character traits
<ul style="list-style-type: none">concentrated in: culture & entertainment ?concentrated in: education ?concentrated in: games ?
show more (2)
assassination ...
character assassination
<ul style="list-style-type: none">concentrated in: multi-topic ?concentrated in: politics & government ?concentrated in: nature & environment ?
show more (2)
arc ...
character arc
<ul style="list-style-type: none">concentrated in: culture & entertainment ?concentrated in: games ?concentrated in: arts ?
show more (1)
development ...
character development
<ul style="list-style-type: none">concentrated in: arts ?concentrated in: games ?concentrated in: culture & entertainment ?
flaw ...
character flaws
<ul style="list-style-type: none">concentrated in: religion ?concentrated in: arts ?concentrated in: culture & entertainment ?

verbs with "character" as object
portray ...
portrayed the character
<ul style="list-style-type: none">concentrated in: culture & entertainment ?concentrated in: reference/encyclopedia ?
play ...
character played
<ul style="list-style-type: none">concentrated in: culture & entertainment ?concentrated in: games ?concentrated in: arts ?
show more (1)
feature ...
featuring characters
<ul style="list-style-type: none">concentrated in: culture & entertainment ?concentrated in: games ?concentrated in: reference/encyclopedia ?
introduce ...
characters introduced
<ul style="list-style-type: none">concentrated in: culture & entertainment ?concentrated in: arts ?concentrated in: games ?
show more (1)
recur ...
a recurring character
<ul style="list-style-type: none">concentrated in: games ?concentrated in: culture & entertainment ?concentrated in: discussion ?
show more (1)
name ...

Word Sketches use Dependency Parsing

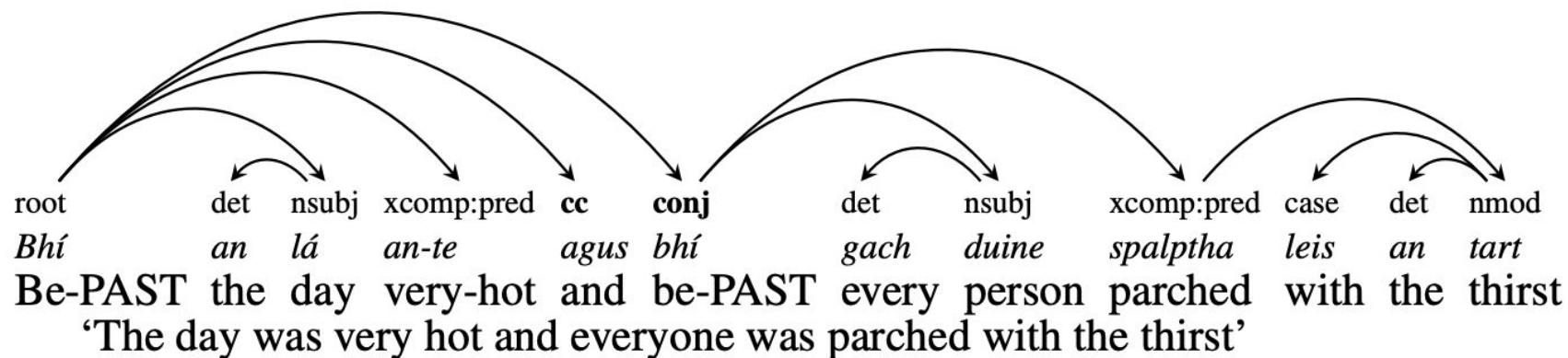


Figure 4 in *Lynn, Teresa and Foster, Jennifer (2016) Universal dependencies for Irish. In: Second Celtic Language Technology Workshop. (CLTW 2016), 4 July 2016, Paris, France.*

Standard Irish

- Official standard introduced in the 1940's and 1950's
- Significant simplifications to spelling and grammar
- Taggers and parsers we have were developed for the standard language
- These perform poorly for pre-standard texts
- Major challenge for the historical dictionary
- <http://corpas.ria.ie/> has 3000 texts published between 1600 and 1926

Standardization

- I developed a tool for standardizing Irish texts, c. 2007
- Shallow statistical MT approach; *does no annotation of pre-standard text*

Cuirfimid ^{anseo} ann so beagán do breugaib na Nua-Ghall a scríobh
Cuirfeam síos ann so beagán do breugaib na Nua-Ghall do scríobh ar Éirinn ar
Lorç Chambrens; agus doigean tosach ar bhréagnú Chambrens féin, mar a n-abair
so raibh cíoscáin as an rí Artúr ar Éirinn, agus surab é am fa' r ceangail an cíos
orthu gCathair ab aois don Tiarna chéad naoi déag
orra i gCathair Leon, an tan fá haois do'n Tigearna cúis céad agus naoi déag, mar
a chuireann ina chroiníc sa den leabhar
cuireas Campion 'na croiníc i san dara caibidil do'n dara leabhar, mar a n-abair

Gaeilge Gaedhilge Gaeilige Gaelige Gailge Gaoidhilge Gaoidheilge Gaelge
Gaidhlige Gaedheilge Gaoidhelge Gailege Gaielge Gaodhailge Gaedilge
Gaeidhilge Gaedhilige Gaoidilge Gaeilgele Gaedhlige Gaédhilge Gaoilge
Gaeillge Gaeilga Gaidhilge Gaelilge Gaodheilge Gaeilge Gaedhilghe
Gadhilge Gaheilge Gaellge Gaoilaige Gaodhilge Gaedhilgé Gaeilege Gaeilge
Gailige Gaeilgé Gaeghilge Gaedhailge Gaoidhlige Gaelgie Gaeiloge Gaeilgle
Gaeilghe Gaeilge Gaeidhlge Gaeidheilge Gaeeilge Gaoilige Gaóilge
Gaoilaga Gaoigheilge Gaoidhlge Gaoidelge Gaoideilge Gaodhéilge Gaieilge
Gaeulge Gaeuilge Gaeolge Gaoidheilge Gaeilgi Gaeilgee Gaeílge Gaeidlge
Gaeidilge Gaeidhelge Gaehilge Gaeeilgee Gaedhlge Gaedhiilge Gaedhelga
Gaédhailge Gaedgilge Gadehilge Gaddhilge Gaoghailge Gaileige Gaidhlhige
Gaidhlge Gaeliage Gaelga Gaéilge Gaedilghe Gaedhulge Gaedhealg Gaedheilg
Gaédheilg Gaedhilg Gaedhilig Gaeilg Gaoidhealg Geadhilge Geailge



Caoimhín Ó Scanail

@kscanne

...

174 litriú ar fhocal amháin, in ord minicíochta 😊

```
kps@borel:~/gaeilge/crubadan/twitter$ egrep -i -o 'ch?[oó](mh?)?gh?[a-z]*r[a-z]*  
d[a-z]*s*' sonrai/ga-tweets.txt | tolow | sed 's/^ch/c/' | sort | uniq -c | sort  
-r -n | sed 's/^ *([0-9]) * //' | tr "\n" " " | fmt  
comghairdeas comghairdeachas comgháirdeachas comghairdeas  
comghairdeachas comghairdeas comgháirdeas comghairdeachas  
comgháirdeachas comghardeas cómhghairdeas comgháirdeachas  
comghairdeaghas cómhgháirdeas comghairdeachais comgháirdeachais  
comghairdis comgheardas comghairdeagas comgháirdeas comghairdneas  
comghairdes cómhgháirdeachas comgháirdis comghairdeagas  
comghairdegas comgháirdeas comghairdeachais comghdeardas comgháirdeas  
comghardachas comgháirdeachas comgháirdeagas comghairdigeas  
comgharidgeas comghairdeacheas comghairdachas comghairdeas  
comghairdeas comghirdeas coghairdeachas comgheardais comghardeachas  
comgairdeas comgairdeachas comghairdeachs cómhghairdeachas comgardaghas  
comghraideas comghiardeas comghairdras comgháirdeachas comghairdeas  
comhgáirdeas comghairdres comghairdegas comghairdech as comghairdeahas  
comghairdas comghairdneas coghairdeas comghsirdeachas comghrds  
comgháirdeachas comghairedeas comghairedeachas comghairdreas  
comghairdeaches comgardech as comgardeachas comgárdachas  
cómhghairdeas comghghairdeachas comghardeas cóghairdeas comgrades  
comghairdeas comghardieas comghardeachais comghairdwas comghairdeas  
comgháirdeachas cómhgháirdeachais comghards comghárdeas comghardeas  
comghardeachais comghairds comghairdgeas comghairdegas comghairdeagais  
comghairdeacheas cómhgháirdeachas cómhghairdeachas cómhghairdeachais  
comghairdeacas comghairdais comgháirdeachas comghardas cómhgháirdeachas  
cógáirdeachas comgradeas comghairdeachas comghairdeachas comghsirdeas  
comghéirdeachas comgheardeachas comghdeardás comghdairdeas  
comghhardis comghárdeas comghárdachas comghaordeas comghairgdeas  
comgháirdeachas comghairdis cómhghairdeachas comghairdeachas  
comghairdess comghairdegas comghairdeas comgháirdech as  
comghairdeass comghairdeasannaagus comghairdeás comghairdeais  
comghairdeaghas comgháirdeacheas comghairdeachás comgháirdeáchas  
comghairdeáchas comgháirdeachas comghairdeacahas comgháirdas  
comghghairdeas comghgardeas comgheardeas comgheardais comgdairdeas  
comghardgas comgardeaghas comghardas comghardaicheas comghairedeachas  
comghairdeis comghairdegeas comghairdecais comgháirdeagas comghairdeachss  
comgháirdeachas comgháirdeachas comghairdeacahas comghairdeas  
comghrdeas comghgháirdeas comghghairdeas comgheardeas comghairgdeas  
comghairedeachas comghairdreas comghairdigeas comghairdegas comghairdecas  
cómhghairdeas comghairdeachais comggairdeas comgardeas comgairdneas  
comgháirdeachas comgairdeachais coghárdas cógardachas cogáirdeachas  
cógairdeachas
```

Traditional processing pipeline

- Run an older text through the standardizer, outputs word-level alignments
- Tag/parse the standardized text using tools for the modern language
- “Project” the annotations back to the original text
 - One-to-many standardization (“naoidheug”): adjust tokenization of source text
 - Many-to-one standardization (“ann so”): DB of 750 most common examples + annotations
- Essentially the pipeline used for the corpas.ria.ie site
- But how well does this work?
- We knew this approach was inherently limited
 - Limited by accuracy of the shallow standardizer, which is very good but will never be perfect
 - More importantly, discards grammatical features which have disappeared in standard Irish

Test corpus of pre-standard Irish texts

- 150 sentences, just under 4000 tokens
- 25 sentences from three 20th c. books, one per major dialect: “Older” corpus
- 25 sentences from three very challenging texts: “Oldest” corpus
 - 1602 Irish New Testament
 - Foras Feasa ar Éirinn (1630s)
 - Cín Lae Amhlaoibh (1820s)
- Manually tagged/parsed following the Universal Dependencies guidelines
- https://github.com/UniversalDependencies/UD_Irish-Cadhan/blob/dev/ga_cadhan-ud-test.conllu

Results (unlabeled attachment)

Model	Standard Irish	Older (1900-1950)	Oldest (1600-1900)
UD	81.8	77.6	61.2
Standardize+Project	81.1	84.8	73.0
UD+Silver training	82.0	84.0	70.6

Observations

- Modern taggers/parsers perform poorly on older texts
- Traditional pipeline using the standardizer gives good results
- But, promising results w/o gold training and w/o using the standardizer directly
- With enough gold training data, we can eliminate the standardizer

Thank you! / Go raibh maith agaibh!

- <https://cadhan.com/>
- <https://github.com/kscanne/>