

# New Frontiers in Language Technology for Minority Languages

---

Kevin Scannell  
Saint Louis University

# Linguistic Diversity: The Numbers

- About 7100 languages spoken in the world (Ethnologue)
- Almost half are “endangered” (UNESCO)
- 2500-3000 have some online presence
- < 1000 written by the language community
- Wikipedias for about 300 languages
- Tweets in about 275 languages
- Spell checkers for about 180
- Google search interface translated into 150
- Google Translate in about 100

# Language Technology: Examples

- Machine translation
- Speech recognition
- Predictive text
- Search engines
- Dialogue systems
- Spelling/grammar checking
- Text normalization (e.g. modernization)
- Optical character recognition

# “Praistriúchán”

New Irish portmanteau word: “praiseach” = “a mess, a botch job”, “aistriúchán” = “translation”





RTE



News

Pope Francis in Ireland

Ireland

World

Business

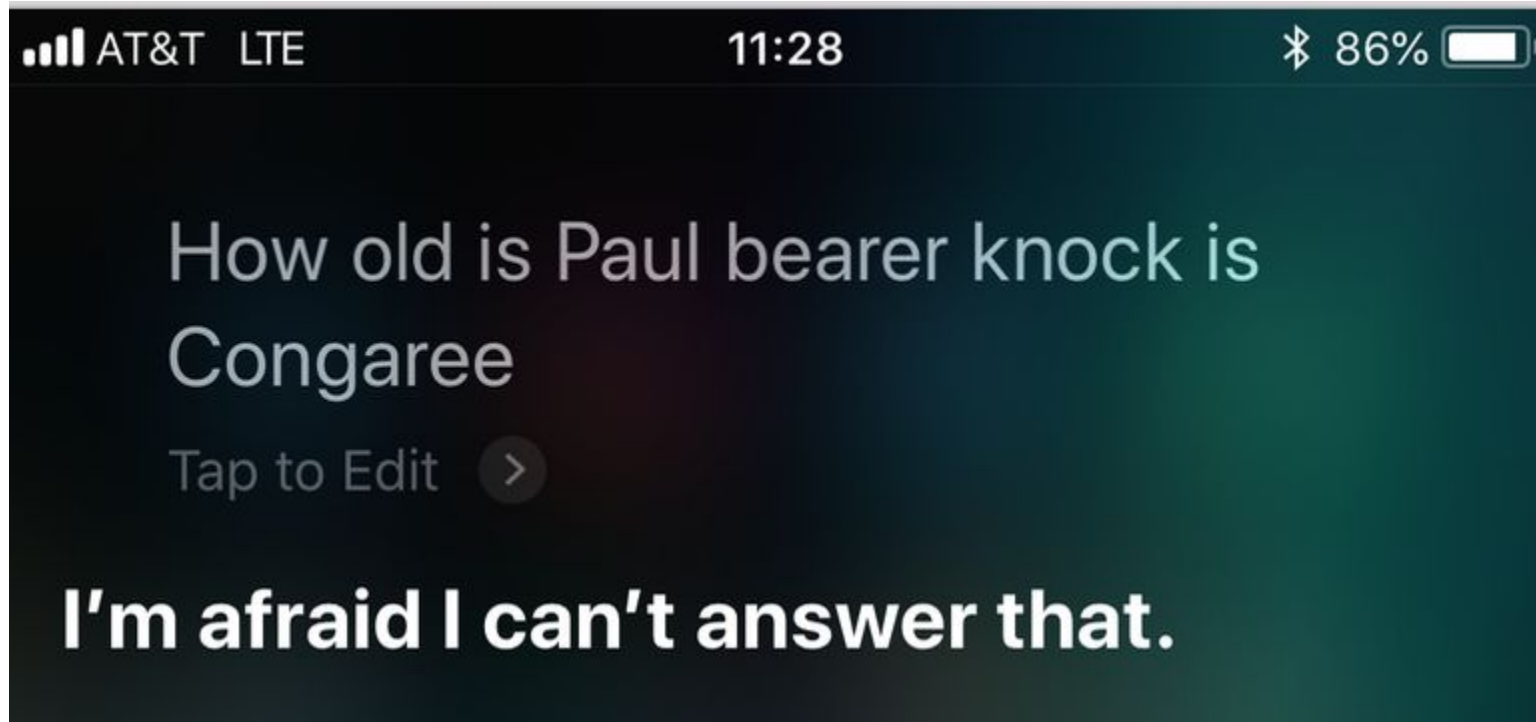
Politics

# Google used to translate English to Irish on 1916 commemorations website

Updated / Thursday, 13 Nov 2014 23:48



## Siri in Irish?

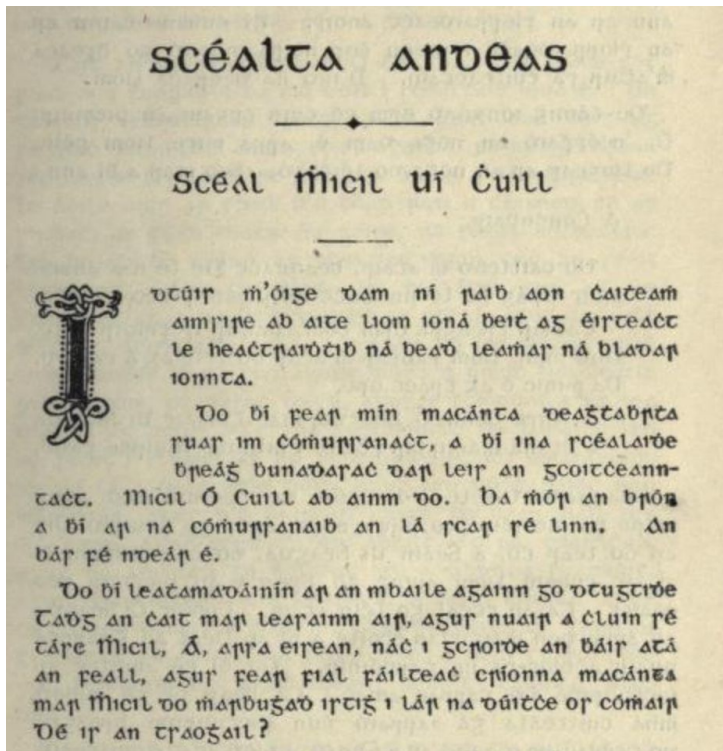


# Google Search

FO/GRAI/ FOSTAI/OCHTA <EMPLOYMENT NOTICES> (ANOIS 2-3,9-10.xi.91)

1. post: Mu/inteoir (Ionadai/) le Mata & Eolai/ocht  
tre/imhse: Samhain go Ma/rta  
fo/n: +353-01-2825872  
seoladh: An Pri/omhoide  
Cola/iste Ra/ithi/n  
Bo/thar Florence  
Bre/  
Co.Chill Mhantain.
2. post: Pri/omhiniu/cho/ir Inmhea/nach  
spriocdha/ta: 15.11.91  
seoladh: An Ceannasai/ Riaracha/in Foirne  
An Post  
Ardoifig an Phoist  
Sra/id Ui/ Chonaill  
Baile A/tha Cliath 1.

# Google Books



SCÉALCA AnXDeAS

scéAi micit m ctiitt

(t^A^S^ T)cúif tfi'óise 'óAni tii fAib Aon óAiteAtfi

^SP^ Aimpife At> xMce tiom lonÁ Deic ^S éifce-Aóc

jl: ^ iieAócf Ai'ócib nÁ beAt) teAtfiAf nÁ t)TAX>At^

jllj lonnCA.

^SS) T) o t)í peAt^ mín mACÁnca ●oeA\$t-At>t\TA

^^ FUAF im óórtiut^fAnAóc, -a bí inA fcéAtAi"óe

t)feÁS GunA'óAf Aó "OAt^ teif An ^coicóeAnii-

CAóc. ITIicit Ó Cuitt aX) -Ainm "oo. t)A tfiót^ ^n bttóu

A ttí A^A nA cómtit^fAnAiO An LÁ yCÁ\< fé tinn. Ati

bÁf f é nT)eÁt\ é.

"Oo t)í leACAmA'oÁinin At^ An mbAite Aj^Ainti 50 ■ociisci'óe  
Ua'ós An ÓAIC m^A^ teAf Ainm -áit^, A^tif ntiAit^ a Otuin fé  
cÁfc Thicit, A, AttfA eifeAn, nÁó 1 scctofóe Art bÁif acá  
An feAtt, A^tif fe^tt fiAt fáitceAó cttionNA mACxinsA  
m^tt fhicit 'oo tfiAttbuSat) ifci\$ 1 tÁf nA "óúitóe of cÓtfiAitt  
"Oé if An cfaO^Ait?



# What's needed?

- **Bigger, better datasets**
- **ML models tailored to the linguistics**
- Technical capacity within communities
- Collaboration with Google, Facebook, Twitter, etc.

# Collecting “everything”

- The Crúbadán project ([crubadan.org](http://crubadan.org)), c. 2000 - present
- Indigenous Tweets ([indigenoustweets.com](http://indigenoustweets.com)), 2011 - present
- RSS feeds (e.g. [chuala.me](http://chuala.me))
- Public Facebook posts
- Feedback loops + crawling
- At least 250 million words of Irish online, before cleaning

# Basic Human Needs



# Digression: Linguistic landscape

- Fully-localized user experience on computer and mobile devices
- Normalization of the language on the computer
- Long but spotty history for Irish!
- Early success via open source software
- Shifting landscape (Thunderbird to GMail, OpenOffice to Google Docs, etc.)
- “Transploitiation”
- Uneven and sometimes uneasy collaboration with big companies
- Not sustainable or scalable
- An Ríomhacadamh ([riomhacadamh.wordpress.com](http://riomhacadamh.wordpress.com))

Bosca Isteach - kscanne@g... x 1140751 - [ga-IE] Fire

https://mail.google.com/mail/u/0/#inbox

Suímh FF IT Gmail Oo KDE Leabhair

Google

Gmail

SCRIOBH

Tábhacht

Bosca Isteach

Tábhachtach

Seolta

Dréachtaí

Gach Rud

Turscar

Bruscar

Ciorcail

Kevin

- Litriú... F7
- Uathsheiceáil an Litriú ↑F7
- Teanga**
- Ionadaitheoir Datha
- Seinnteor Meán
- Íoslighdaigh an Láithreireacht...
- Macraí
- Bainisteoir na nEisinteachtaí...
- Socruithe an Scagaire XML...
- Roghanna UathCheartaithe...
- Saincheap...

- Don Téacs go Léir**
  - Tiontú Hangu/Hanja... ↑F7
  - Aistriúchán na Sinise...
  - Teasáras... ¶F7
  - Fleisciniú
  - Tuilleadh Foclóirí Ar Líne...
- Gaeilge
- Faic (Ná sheiceáil an litriú)
- Úsáid an Teanga Réamhshocrataí
- Tuilleadh...

https://www.facebook.com

Suímh FF IT

Gaeilge

Kevin Scannell

Cuir do Phróifíl in Eagar

CEANÁIN

Fotha Nuachta

Teachtaireachtaí 20+

Ócáidí 6

GRÚPAÍ

Setswana Firefox O...



Gaeilge Amháin

Bailraíocht glactha

Roinn Fógraí

Baill Ócáidí Grianghraif Comhad

Cuir Grianghraif / Fiseán Leis Cuir Ceist

rud éigin...

CHT LE DÉANAÍ

Seanán Ó Coistín

luair

inn aon duine an dearmad faoi deara? Ní fheadar conas mar a seo?

BAILL

Líon na mball - 6.275 (126 nua)

Cuir Daoine Leis an nGrúpa

SPORT RTE NA GAEILGE

Tabhair Culreath Tri Riomphost

CUR SIÓS

See More

CRUTHAIGH GRÚPAÍ NUA

Cuireann grúpaí ar do chumas nithe a roinnt go héasca le chairde le rí 45%

Cruthaigh Grúpa

Comhrá (múchta)

Gan Teideal 1

Sleamhnán

Jótáí Dáileáir

Airíonna

Leaganacha Amach

1

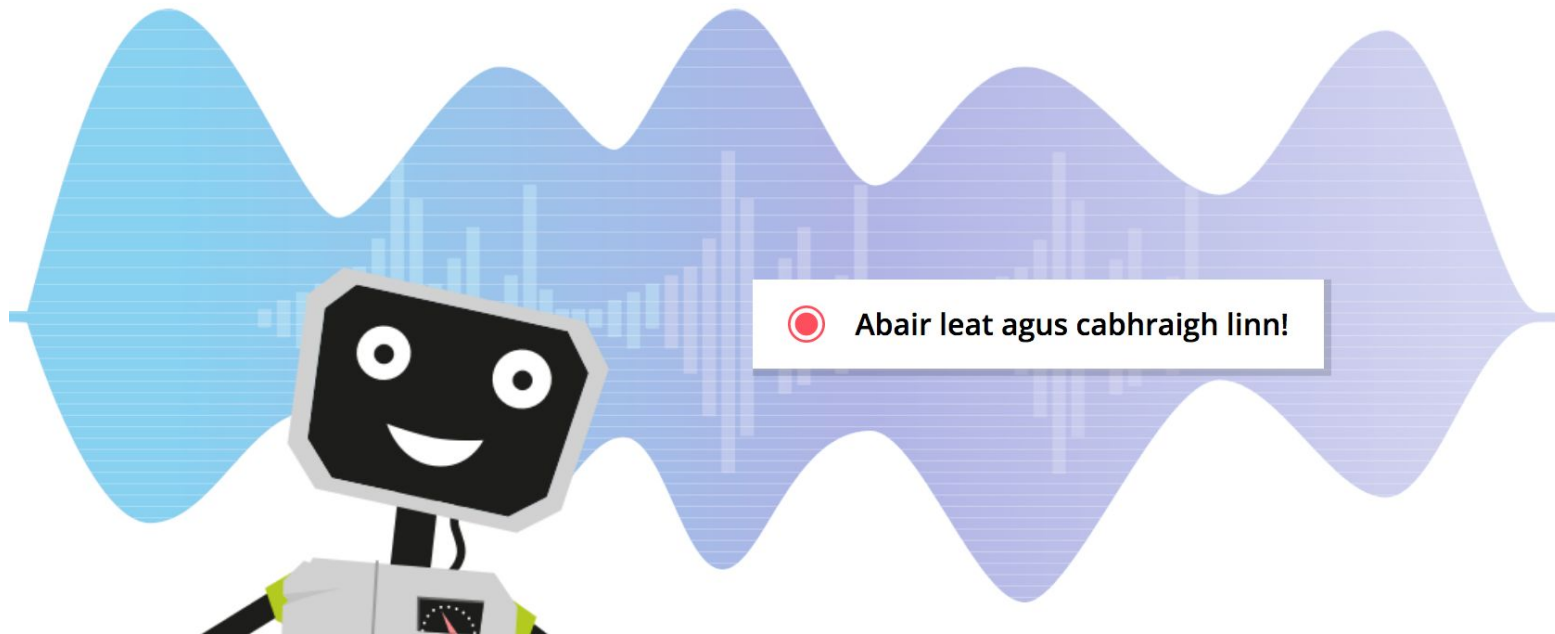
## Common Voice

moz://a

0  
▶ 141



Seo é Common Voice – córas de chuid Mozilla a mhúineann do ríomhairí an chaoi a labhraíonn daoine.



# Language modeling

- Chomsky: “... the notion of ‘probability of a sentence’ is an entirely useless one, under any known interpretation of this term.”
- Let  $S = \text{“this is an entirely useless notion”}$
- $P(S) = P(\text{this} | \wedge) P(\text{is} | \text{this}) P(\text{entirely} | \text{this is}) \dots P(\text{notion} | \text{this is an entirely useless})$
- Usually formulated and computed this way (word prob. given history)
- Humans are pretty good at estimating these, at least
- $P(\text{Friday} | \text{My party is this coming}) > P(\text{Tuesday} | \text{My party is this coming})$
- $P(\text{is} | \text{The man with the glasses}) > P(\text{are} | \text{The man with the glasses})$

# Applications

- Virtually all of the language technologies mentioned above!
- Can be used generatively
- MT, ASR, etc. fundamentally generate text, conditioned on input
- Traditional probabilistic models (“noisy channel”):  $P(T | S) = P(S | T) P(T) / P(S)$
- Baked into many neural models for MT, etc.
- Better language models give better end-to-end performance, generally



# Neural language models

- A flood of recent papers on neural language modeling, big leaps forward
- Almost 100% (and implicitly!) focused on English
- Word “English” doesn’t even appear in Google Brain’s landmark 2016 paper
- Same models applied to Irish show only small gains at best

# Incorporating linguistic knowledge

- Celtic languages have “initial mutations”
- Almost always predictable from previous two words
- *bád seoil* “sailboat”, *mo bhád seoil* “my sailboat”, *ár mbád seoil* “our sailboat”
- N-gram models don’t “see” that these are all the same word
- If most training examples are first type, say, harder to predict collocation
- (Google even gets this wrong in previous image: *tríd an mbóthar*)
- Easy enough to use a *factored language model* to get better results for Irish
- Examples like this abound...

# Some success stories

- Intergaelic MT engines
  - Google Translate-like interface
  - Annotated novels and short stories
  - Social media translations, mouse-over help
- Irish standardization
  - Foclóir Stairiúil na Gaeilge
  - Corpas na Gaeilge 1600-1926
  - New English-Irish dictionary project

FOCLÓIR **AISTRIÚCHÁN**

Tha crois margaid sgoinneil ann bhon 15mh linn. Thog John Cockburn Ormiston mar bhaile modail ann an 1735 mus do lorgadh gual faisg air anns an 19mh linn. Dh'fhàs an cunntas-sluaigh còmhla ris a' ghnìomhachas mèinnearachd. Dùineadh na mèinnean anns na 1960an. Chaidh Ormiston Hall a thogail airson John Cockburn ann an 1745. Cheannaich an Iarla Hopetoun e agus

 Aistrigh »

Tha crois margaid sgoinneil ann bhon 15mh linn. Thog John Cockburn Ormiston mar bhaile

Tá cros margadh iontach ann ón 15ú haois. Thóg John Cockburn Ormiston mar bhaile

modail ann an 1735 mus do lorgadh gual faisg air anns an 19mh linn. Dh'fhàs an

samhail i 1735 sular aimsíodh gual i ngar dó sa 19ú haois. D'fhás an

cunntas-sluaigh còmhla ris a' ghnìomhachas mèinnearachd. Dùineadh na mèinnean anns na

daonáireamh in éineacht leis an ghnó mianadóireachta. Dúnadh na mianaigh sna

1960an. Chaidh Ormiston Hall a thogail airson John Cockburn ann an 1745. Cheannaich an

1960í. Chuaigh Ormiston Hall a thógáil ar son John Cockburn i 1745. Cheannaigh an

Iarla Hopetoun e agus Ormiston fhèin agus leasaicheadh an talla ann an 1772. Sgrìos teine e

Iarla Hopetoun é agus Ormiston féin agus feabhsaíodh an halla in 1772. Sgrìos tine é

anns an Dàrna Cogadh.

sa Dara Cogadh.

## SKEEALAGHT

Brollach

Cabbyl Folliahtagh

Dooinney Va Ceau "Oilskins"

Dooinney Nagh Row Ayn

Shenn Dooinney

Babyr Naight

Buill yn Tarroo-Ushtey

Shenn Ven Veg

Yn Ven Va ny Buitch

Ree Edard ayns y Chashtal

Cabbyl-Ushtey ec Glion Meay

Yn Vible Vooar

Chymsaghey Fuygh

Ny Kirree fo Niaghtey er Slieau

Whuallian

"Jarrood y Theihll, Jarroodit ec y

Theihll"

Juan Beg as yn Dooinney Oie?

## YN VIBLE VOOAR

Tra va mee gobbragh ayns Purt ny hInshy ec oik ny barroosyn, va shenn dooinney cheet dys Purt ny hInshy ny keayrtyr voish Skeristal t'er y raad dys Rhumsaa mysh tree ny kiare meeilaghyn voish Purt ny hInshy. V'eh enmyssit Tommy Kneale as v'eh baghey ayns Skeristal, as v'eh ny ghooinney beg as va symm echey er y Ghaelg as er ny shenn skeealyn. As dooyrt eh dy row bible sy Ghaelg echey ayns y thie as v'ee foddey ny smoo na ny bibleyn cadjin — v'ee ny smoo na daa hrie er e lhurid as foddey ny smoo na ny bibleyn elley, as hug eh cuirrey dooys cheet dys y thie echey dy yeeaghyn urree traa erbee, as va shen goll er daa vlein agh cha row rieu tra aym dy ghoill.

Aghterbee, cha jagh mee rieu lesh shilley er y vible shoh; as eisht cheayll mee dy row eh er ngeddyn baase as va mee boirit mychione shen, as cheayll mee myrgeddin dy row ooilley e stoo-thie tilgit magh ass y thie as aile currit orroo as va ooilley caillt. Aghterbee, va mee shicky dy row yn vible shoh, yn vible neu-chadjin, dy row ee loshtit as stroit as currit mow. Agh cha row mee my vooijer da as cha row cair erbee aym dy vannoo briaht mychione v vible. agh tammvlt nv lurg

## AN BÍOBLA MÓR

Nuair a bhí mé ag obair i bPort na hInse ag oifig na mbusanna, bhí seandúine ag teacht go Port na hInse uaireanta ó Skeristal atá ar an mbealach chuig Rhumsaa timpeall trí nó ceithre mhíle ó Phort na hInse. Bhí sé ainmnithe Tommy Kneale agus bhí sé ina chónaí i Skeristal, agus bhí sé ina dhuine beag agus bhí suim aige sa Mhanainnis agus sna seanscéalta. Agus dúirt sé go raibh bíobla sa Mhanainnis aige sa teach agus bhí sí i bhfad níos mó ná na bíoblaí coitianta — bhí sí níos mó ná dhá thoirigh ar a fhad agus i bhfad níos mó ná na bíoblaí eile, agus thug sé cuireadh domsa teacht go dtí an teach aige le féachaint uirthi am ar bith, agus bhí sin ag dul ar dhá bhliain ach ní raibh riamh am agam le dul.

Ar aon chuma, ní dheachaigh mé riamh le cuairt ar an bhíobla seo; agus ansin chuala mé go raibh sé tar éis bás a fháil agus bhí mé buartha faoi sin, agus chuala mé freisin go raibh uile a thoscáin caite amach as an teach agus tine curtha orthu agus bhí uile caillte. Ar aon chuma, bhí mé cinnte go raibh an bíobla seo, an bíobla neamhghnách, go raibh sé dóite agus scriosta agus scriosta. Ach ní raibh mé i mo mhuintir dó agus ní raibh ceart ar bith agam le déanamh fiosraithe faoin bhíobla.

# Standardization / indexing: [corpas.ria.ie](http://corpas.ria.ie)

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 Gaelge Gaeidhlge Gaeidheilge Gaeilge Gaoilige Gaóilge  
Gaoilaga Gaoigheilge Gaoidhlge Gaoidelge Gaoideilge Gaodhéilge Gaieilge  
Gaeulge Gaeuilge Gaeolge Gaoidheilge Gaeilgi Gaeilgee Gaeílge Gaeidlge  
Gaeidilge Gaeidhelge Gaehilge Gaeilgee Gaedhlge Gaedhiilge Gaedhelga  
Gaédhailge Gaedgilge Gadehilge Gaddhilge Gaoghailge Gaileige Gaidhlige  
Gaidhlge Gaeliage Gaelga Gaéilge Gaedilghe Gaedhulge Gaedhealg Gaedheilg  
Gaédheilg Gaedhilg Gaedhilig Gaeilg Gaoidhealg Geadhilge Geailge