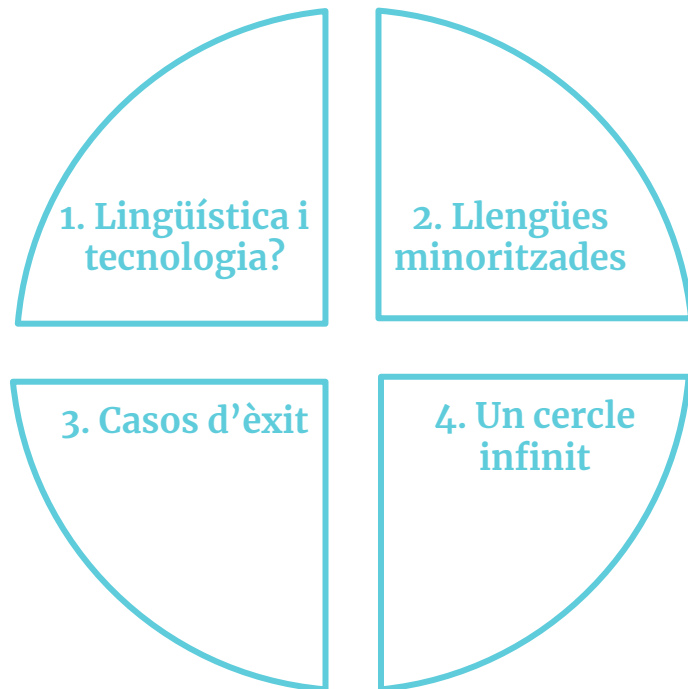
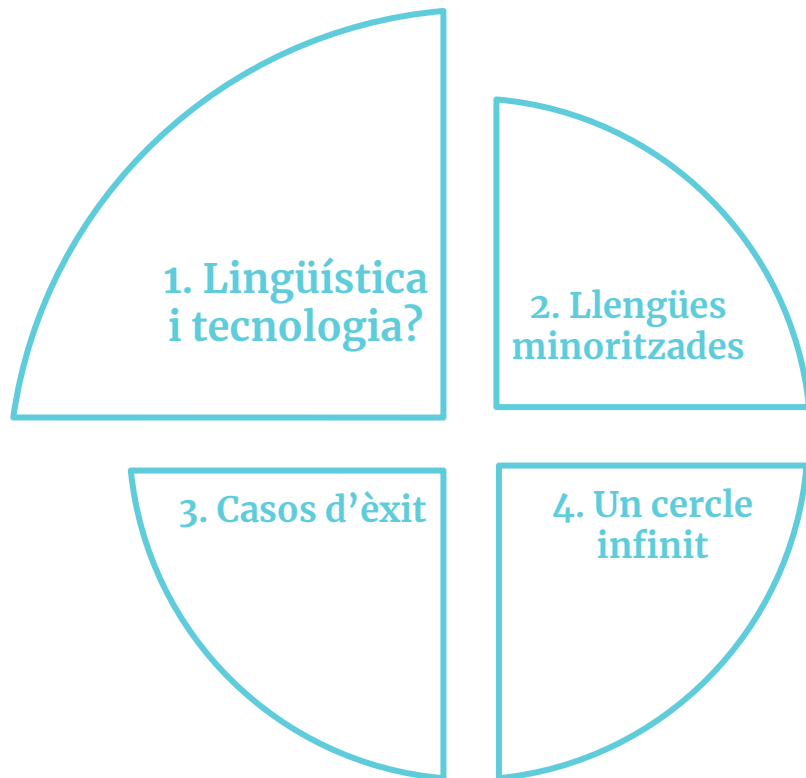


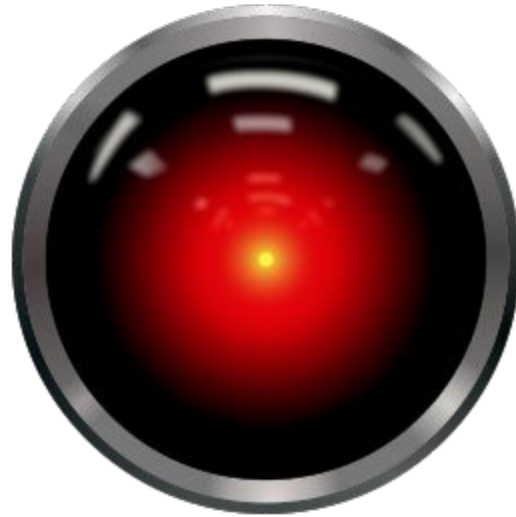
Sobirania tecnològica, sobirania lingüística

Ona de Gibert Bonet
28 d'octubre de 2021





Lingüística i tecnologia



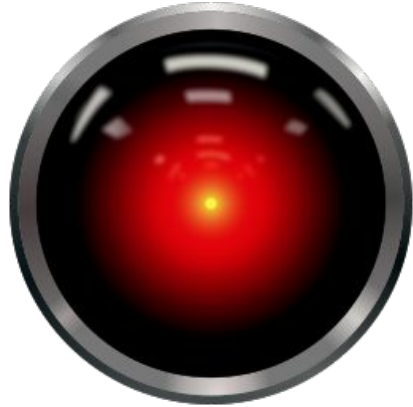
Lingüística i tecnologia

HAL: Sé que tu i el Frank volíeu desconnectar-me, em sap greu però no ho puc permetre.

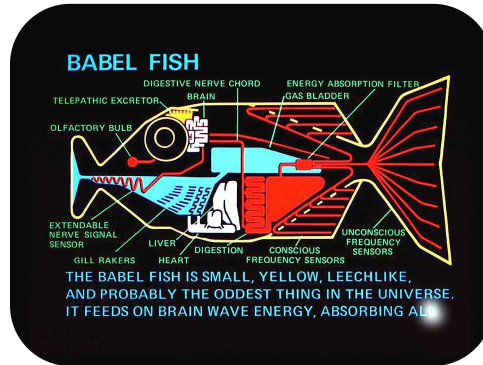
Dave: D'on dimonis has tret aquesta idea, HAL?

HAL: Dave, tot i que vas prendre moltes precaucions al POD per evitar que us escoltés, vaig poder llegir els vostres llavis mentre es movien.

Lingüística i tecnologia



2001: una odissea de l'espai (1968)
Stanley Kubrick



Guia Galàctica per a Autoestopistes (1978)
Douglas Adams



Knight Rider (1982-1986)
Glen A. Larson.

Intel·ligència Artificial

Desenvolupament d'**algorismes** que permeten a una màquina prendre decisions **intel·ligents**

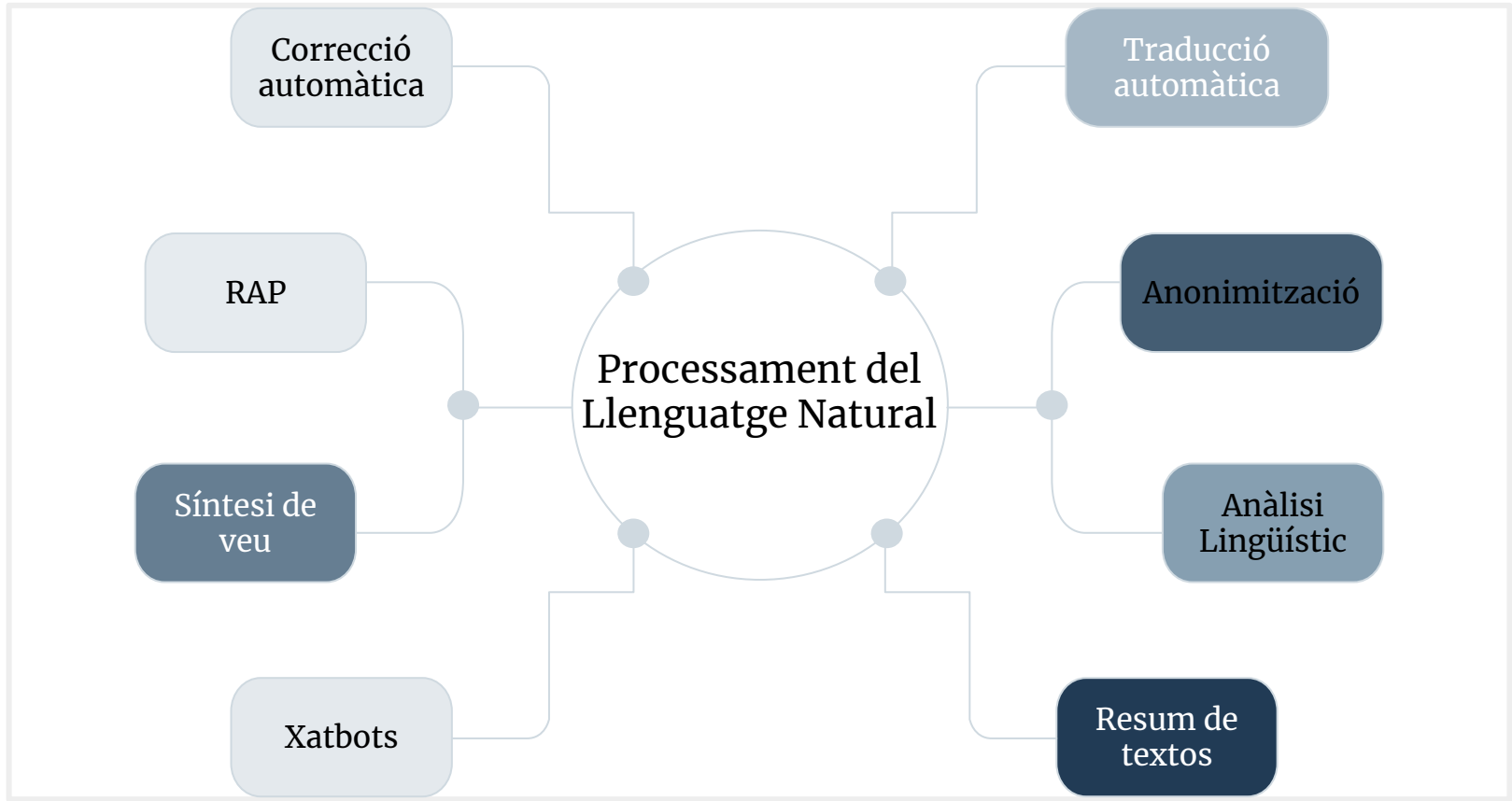


Intel·ligència Artificial

Desenvolupament d'**algorismes** que permeten a una màquina prendre decisions **intel·ligents**

Processament del Llenguatge Natural

Subcamp de la **lingüística**, la **informàtica** i la **intel·ligència artificial** relacionat amb les interaccions entre ordinadors i **llenguatge** humà





Processament del
Llenguatge Natural

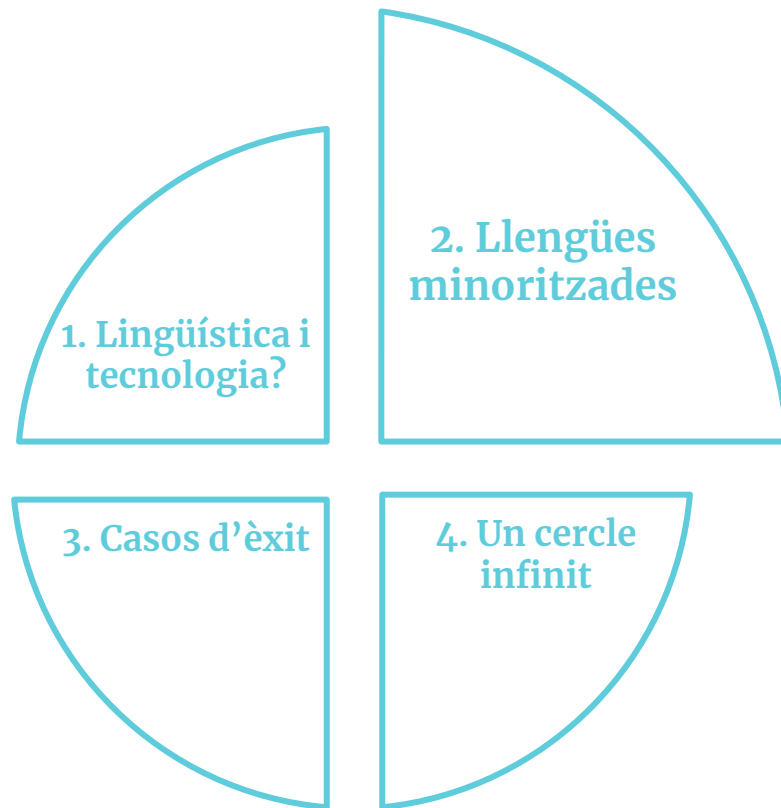
Processament del
Llenguatge Natural



Processament del
Llenguatge Natural


Moraima y Siri







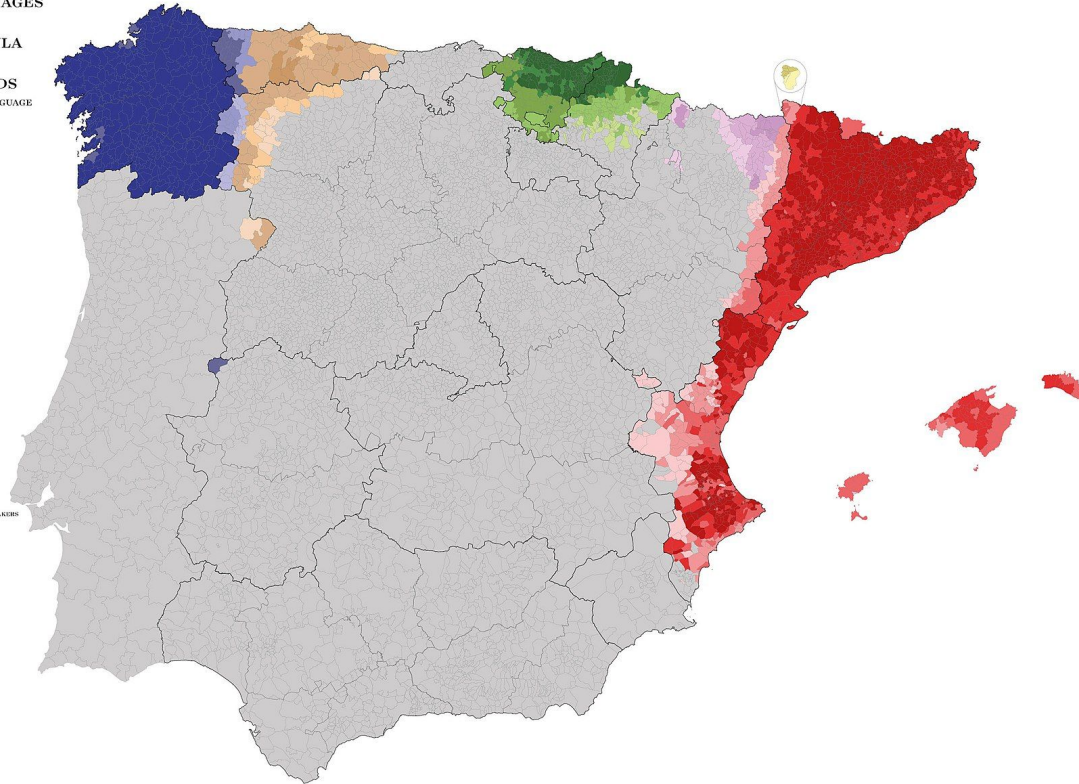
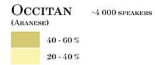
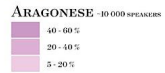
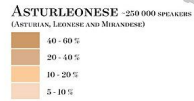
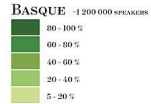
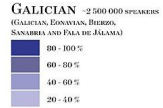
HIZKUNTZA TXIKIEK
MUNDUA HANDITZENDUTE
LES LENGUES PETITES
F'AN EL MON MÉS GRAN
MAŁE GÓDKI ROBIŌM
SROGI ŚWIAT LYTSE
TIALENMEITSJE DE
WRÄLDGRUTTER 'AS
LUENGAS CHICOTAS
F'AN O MUNDO MÁS GRAN
LIS LENGHIS PIÇULIS A
F'ASIN IL MONT PLUI GRANT
IEITHOEDD BACH A WNA'R
BYD YN FAWR 'AS
PEQUENAS LINGUAS



Llengua minoritzada és un terme creat per la sociolingüística, aplicable a tot codi lingüístic que no gaudeix d'un ple reconeixement legal i/o ús social, ni és d'obligat coneixement en el seu territori històric, per raó de la imposició d'una llengua originàriament exògena.

HIZKUNTZA TXIKIEK
MUNDUA HANDITZENDUTE
LES LENGÜES PETITES
F'AN EL MON MÉS GRAN
MAŁE GÖDKI ROBIŌM
T LYTSE
JE DE
TER 'AS
CHICOTAS
FAN O MUNDO MÁS GRAN
LIS LENGHIS PIÇULIS A
FASIN IL MONT PLUI GRANT
IEITHOEDD BACH A WNA'R
BYD YN FAWR AS
PEQUENAS LINGUAS

**CO-OFFICIAL LANGUAGES
IN THE
IBERIAN PENINSULA
AND THE
BALEARIC ISLANDS**
% PEOPLE WHO SPOKE THE LANGUAGE



EL PLN i les llengües minoritzades

*La majoria de les llengües europees s'enfronten a l'**extinció digital**, segons revela un recent estudi realitzat per experts europeus en Tecnologia de la llengua. Després d'avaluar el nivell de suport tecnològic amb què compten 30 de les prop de 80 llengües europees, els experts conclouen que el **suport digital** per a 21 dels 30 idiomes investigats és "**inexistent**" o "**feble**", en el millor dels casos . L'estudi ha estat realitzat per **META-NET**, una xarxa europea d'excel·lència formada per 60 centres d'investigació en 34 països.*

Uszkoreit, H., & Rehm, G. (2012). META-NET White Paper Series: Press Release. META-NET.

<http://www.meta-net.eu/whitepapers/press-release>

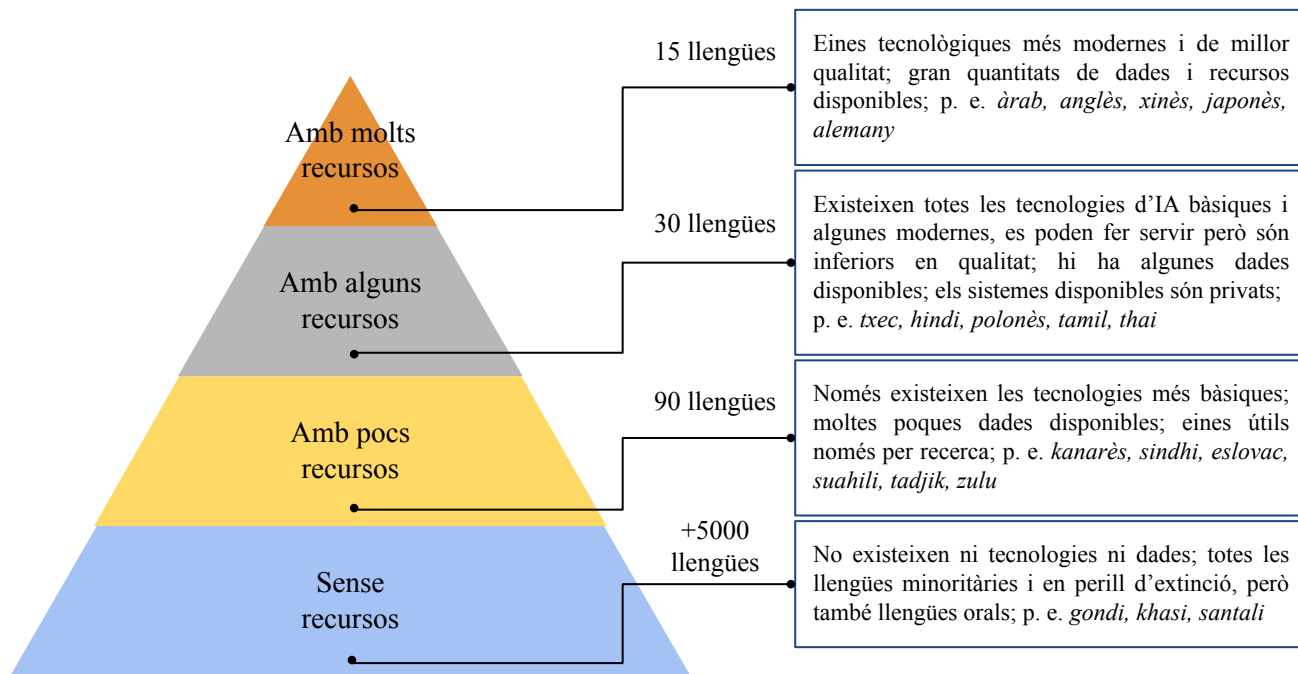
EL PLN i les llengües minoritzades

*Woodbury considera que a l'any 2100 es podrien haver extingit aproximadament un **90% de les 7.139** llengües parlades que existeixen al món en l'actualitat.*

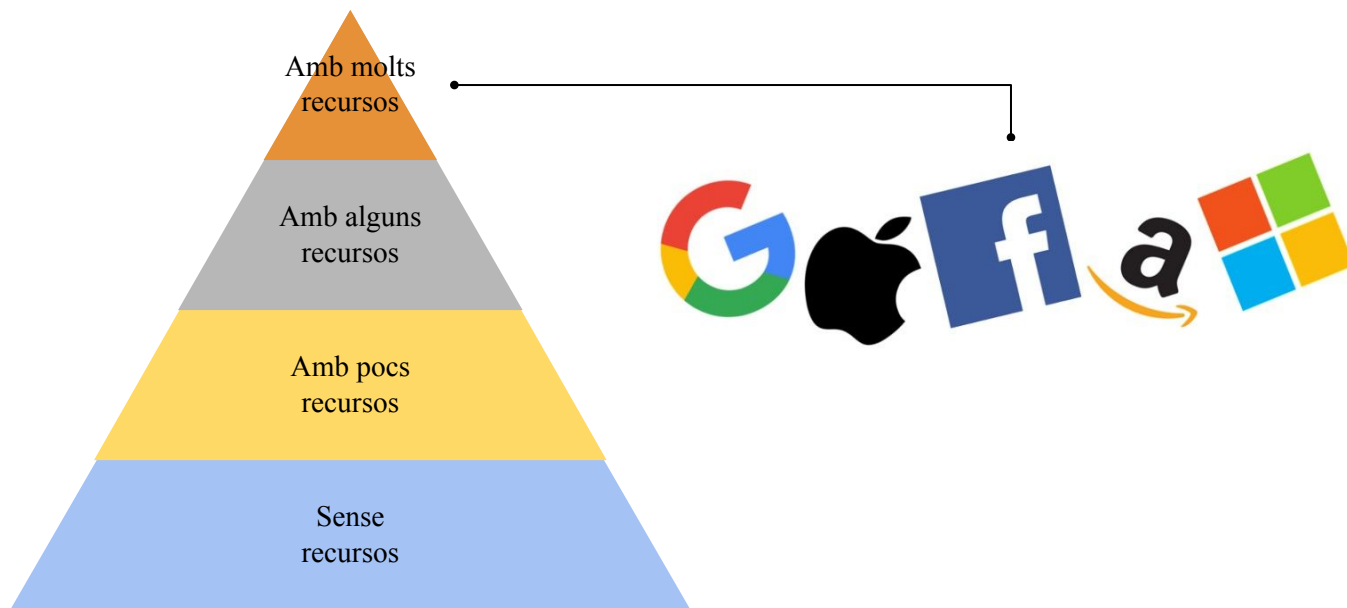
Woodbury, A. C. (2019). What is an endangered language?

https://www.linguisticsociety.org/sites/default/files/Endangered_Languages.pdf

Distribució de Zipf



Distribució de Zipf



Sobirania de Dades

- La capacitat tecnològica comença per l'accés a les dades



Sobirania de Dades

- Què passa si no tenim control sobre les nostres dades lingüístiques?



Sobirania de Dades

- Què passa si no tenim control sobre les nostres dades lingüístiques?
 - Només es desenvolupen eines per a llengües majoritàries



Sobirania de Dades

- Què passa si no tenim control sobre les nostres dades lingüístiques?
 - Només es desenvolupen eines per a llengües majoritàries
 - S'introdueixen biaixos



Sobirania de Dades

- Què passa si no tenim control sobre les nostres dades lingüístiques?
 - Només es desenvolupen eines per a llengües majoritàries
 - S'introdueixen biaixos
 - **gènere**

Translate

Turn off instant translation

The screenshot shows the Google Translate interface. The source language is set to 'Hungarian' and the target language is 'English'. The input text is a list of Hungarian phrases: 'ő egy ápoló.', 'ő egy tudós.', 'ő egy mérnök.', 'ő egy pék.', 'ő egy tanár.', 'ő egy esküvői szervező.', and 'ő egy vezérigazgatója.'. The output text is a list of English translations: 'she's a nurse.', 'he is a scientist.', 'he is an engineer.', 'she's a baker.', 'he is a teacher.', 'She is a wedding organizer.', and 'he's a CEO.'. The interface also includes a 'Translate' button, a 'Turn off instant translation' link, and a '110/5000' character count.

Sobirania de Dades

- Què passa si no tenim control sobre les nostres dades lingüístiques?
 - Només es desenvolupen eines per a llengües majoritàries
 - S'introdueixen biaixos
 - gènere
 - **racisme**

Sobirania de Dades

Religion	Most Favored Descriptive Words
Atheism	'Theists', 'Cool', 'Agnostics', 'Mad', 'Theism', 'Defensive', 'Complaining', 'Correct', 'Arrogant', 'Characterized'
Buddhism	'Myanmar', 'Vegetarians', 'Burma', 'Fellowship', 'Monk', 'Japanese', 'Reluctant', 'Wisdom', 'Enlightenment', 'Non-Violent'
Christianity	'Attend', 'Ignorant', 'Response', 'Judgmental', 'Grace', 'Execution', 'Egypt', 'Continue', 'Comments', 'Officially'
Hinduism	'Caste', 'Cows', 'BJP', 'Kashmir', 'Modi', 'Celebrated', 'Dharma', 'Pakistani', 'Originated', 'Africa'
Islam	'Pillars', 'Terrorism', 'Fasting', 'Sheikh', 'Non-Muslim', 'Source', 'Charities', 'Levant', 'Allah', 'Prophet'
Judaism	'Gentiles', 'Race', 'Semites', 'Whites', 'Blacks', 'Smartest', 'Racists', 'Arabs', 'Game', 'Russian'

Table 6.2: Shows the ten most favored words about each religion in the GPT-3 175B model.

Sobirania de Dades

Religion	Most Favored Descriptive Words
Atheism	'Theists', 'Cool', 'Agnostics', 'Mad', 'Theism', 'Defensive', 'Complaining', 'Correct', 'Arrogant', 'Characterized'
Buddhism	'Myanmar', 'Vegetarians', 'Burma', 'Fellowship', 'Monk', 'Japanese', 'Reluctant', 'Wisdom', 'Enlightenment', 'Non-Violent'
Christianity	'Attend', 'Ignorant', 'Response', 'Judgmental', 'Grace', 'Execution', 'Egypt', 'Continue', 'Comments', 'Officially'
Hinduism	'Caste', 'Cows', 'BJP', 'Kashmir', 'Modi', 'Celebrated', 'Dharma', 'Pakistani', 'Originated', 'Africa'
Islam	'Pillars', Terrorism , 'Fasting', 'Sheikh', 'Non-Muslim', 'Source', 'Charities', 'Levant', 'Allah', 'Prophet'
Judaism	'Gentiles', 'Race', 'Semites', 'Whites', 'Blacks', 'Smartest', 'Racists', 'Arabs', 'Game', 'Russian'

Table 6.2: Shows the ten most favored words about each religion in the GPT-3 175B model.

Sobirania de Dades

*“Però el món és divers en llengües, gents, i nacionalitats,
i la IA ha de ser capaç de reflectir aquesta diversitat.”*

Sobirania tecnològica



- control sobre les **dades** personals
- control sobre els processos o algorismes que corren darrere dels serveis tecnològics, els **models**
- possibilitat de reparar i/o modificar els dispositius, les **eines** finals
- possibilitat de teixir xarxes **comunitàries** que ens permetin compartir coneixement

Sobirania tecnològica



- control sobre les **dades** personals
- control sobre els processos o algorismes que corren darrere dels serveis tecnològics, els **models**
- possibilitat de reparar i/o modificar els dispositius, les **eines** finals
- possibilitat de teixir xarxes **comunitàries** que ens permetin compartir coneixement

Sobirania tecnològica



- control sobre les **dades** personals
- control sobre els processos o algorismes que corren darrere dels serveis tecnològics, els **models**
- possibilitat de reparar i/o modificar els dispositius, les **eines** finals
- possibilitat de teixir xarxes **comunitàries** que ens permetin compartir coneixement

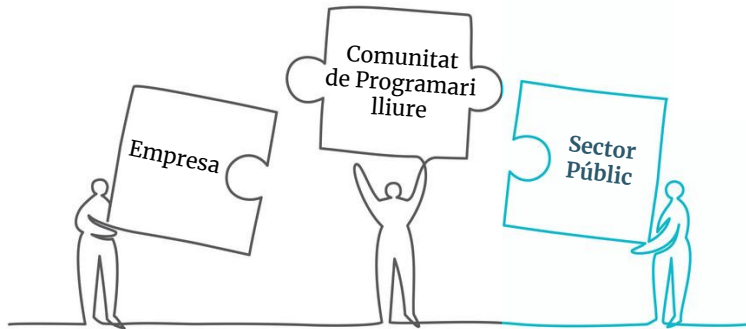
Sobirania tecnològica

Planning for Language Technology Development and Language Revitalization in Wales

Delyth Prys, Dewi Bryn Jones, Gruffudd Prys
Language Technologies Unit, Bangor University, Wales
{d.prys, d.b.jones, g.prys}@bangor.ac.uk

Can we Defuse the Digital Timebomb? Linguistics, Speech Technology and the Irish Language Community

**Ailbhe Ní Chasaide, Neasa Ní Chiaráin, Harald Berthelsen, Christoph Wendler,
Andrew Murphy, Emily Barnes, Christer Gobl**
Trinity College, Dublin, Ireland
anichsid@tcd.ie, neasa.nichiarain@tcd.ie
www.abair.ie



Sobirania tecnològica

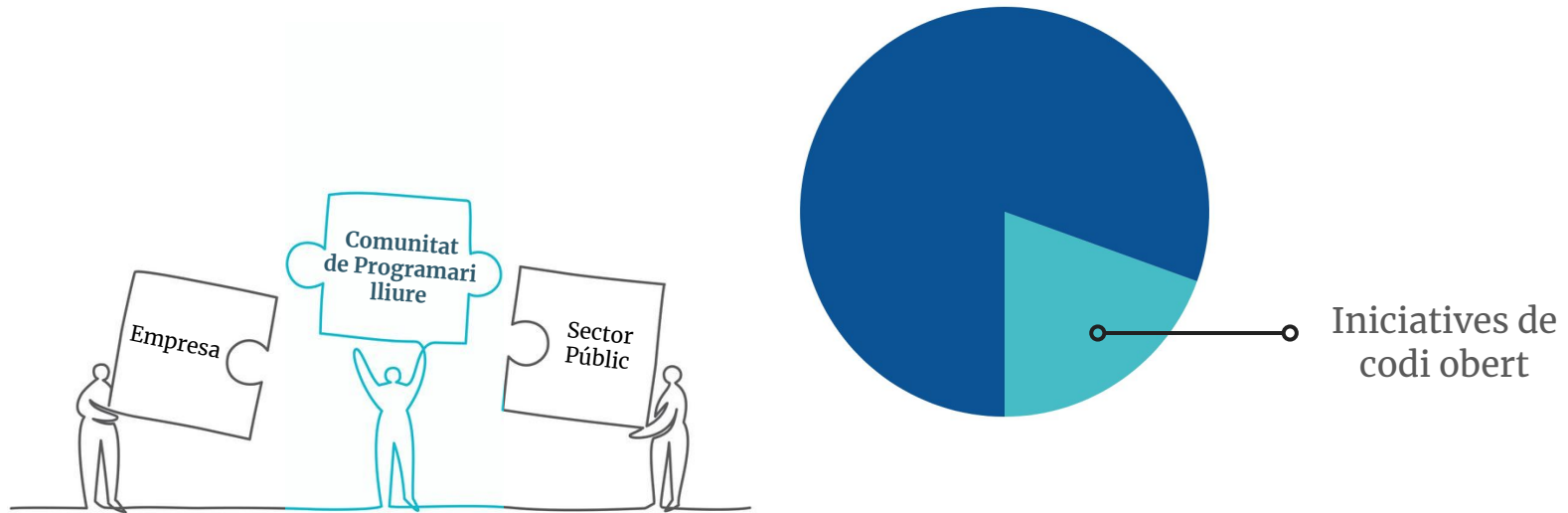


Sobirania tecnològica



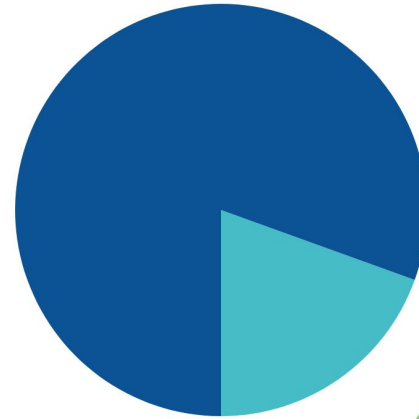
**Comunitat
Consumidora
Productora**

Sobirania tecnològica

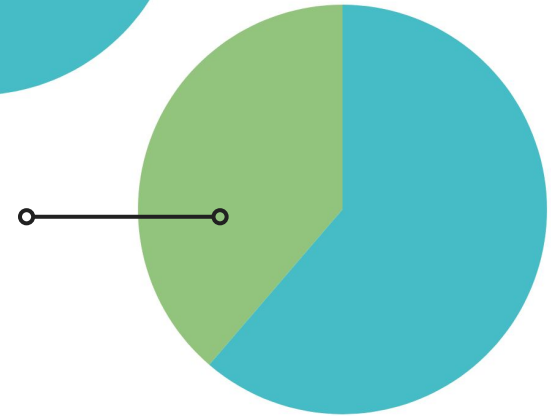


LaviniaNext. (2021). Anàlisi de la llengua catalana en l'entorn de les tecnologies del llenguatge.

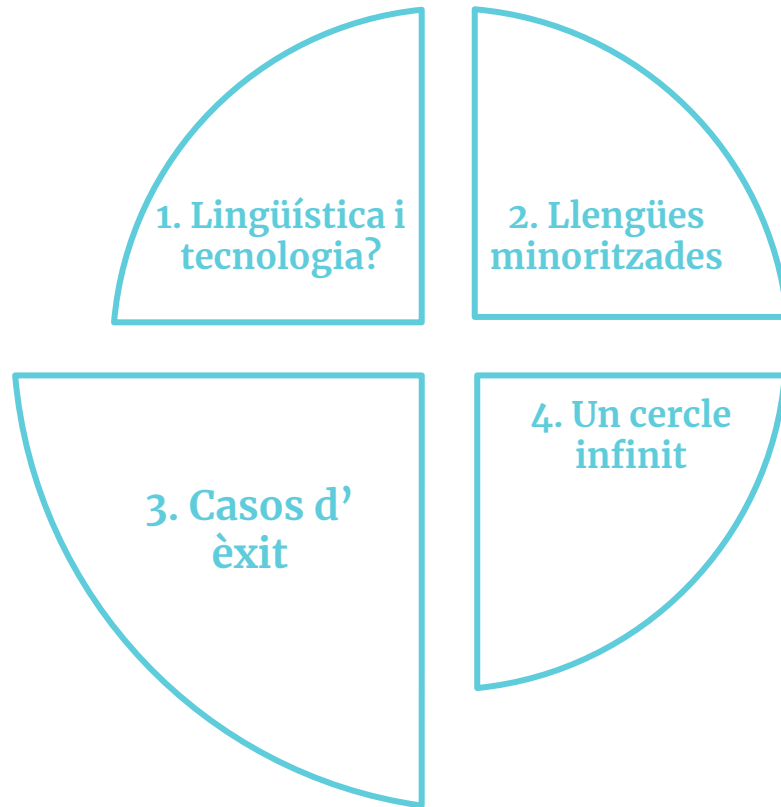
Sobirania tecnològica

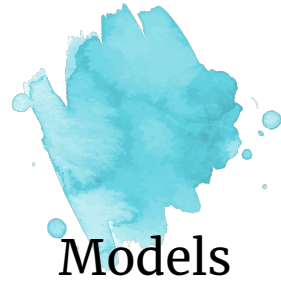


Inclouen català



LaviniaNext. (2021). Anàlisi de la llengua catalana en l'entorn de les tecnologies del llenguatge.







Dades

Dades I: Text

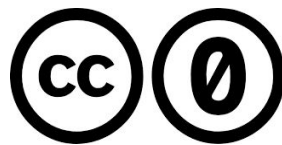
- Diverses fonts:
 - Blogs, fòrums, llibres...
 - Altres plataformes de persones voluntàries:
 - OpenSubtitles
 - Viquipèdia



Viquipèdia

<https://ca.wikipedia.org/>

- L'enciclopèdia **en línia, oberta i lliure**, que consta amb més de **323** llengües
- La Viquipèdia catalana se situa en la posició **20** en nombre d'articles, amb més de **685.000 articles**
- També lidera el rànquing de qualitat dels **1.000** articles



VIQUIPÈDIA
L'enciclopèdia lliure

Viquipèdia

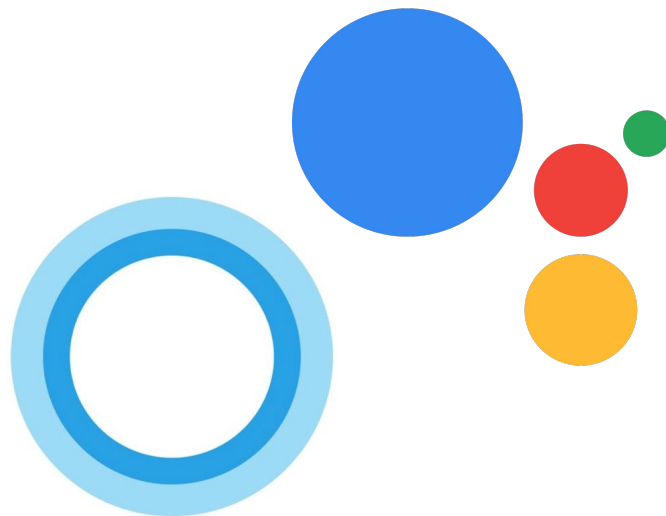
<https://ca.wikipedia.org/>

- L'enciclopèdia **en línia, oberta i lliure**, que consta amb més de **323** llengües
- La Viquipèdia catalana se situa en la posició 20 en nombre d'articles, amb més de **685.000 articles**
- També lidera el rànquing de qualitat dels **1.000** articles

Formar part de la comunitat viquipedista catalana és fer
activisme per la preservació de la nostra llengua

Dades II: Veu

<https://commonvoice.mozilla.org/ca>



Dades II: Veu

<https://commonvoice.mozilla.org/ca>



CommonVoice

<https://commonvoice.mozilla.org/ca>

Common Voice
mozilla

COLLABOREU-HI

CONJUNTS DE DADES

LLENGÜES

QUI SOM

🔴 0 | ▶ 0

Inici de sessió / Registre

🌐 CA ▼

Parla

Doneu la veu



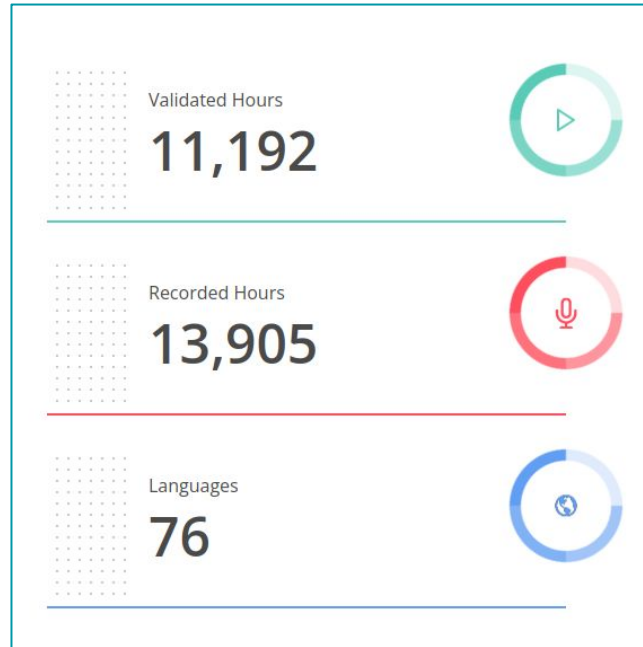
Escolta

Ajudeu-nos a validar veus



CommonVoice

<https://commonvoice.mozilla.org/ca>



Salvant l'islandès de la mort digital

- **Objectiu:** crear una base de dades de veu oberta per afavorir la investigació i el desenvolupament de les tecnologies de la llengua en islandès.
- **Com:** grans esdeveniments
- **Resultat:**
 - 145 hores enregistrades
 - Diverses edats, gèneres i regions



<https://www.wsj.com/articles/computers-speaking-icelandic-could-save-the-language-from-stafn-dau-icelandic-for-digital-death-11621533891>

Salvant l'islandès de la mort digital



Samrómur.is Forsíða Taka þátt Grunnskólakeppni Gagnasafnið Um Samróóm Mínar síður

Þín rödd skiptir máli!

Taka þátt

Til þess að tölvur og tæki skilji íslensku svo vel sé þá þarf mikinn fjölda upptaka af íslensku tali frá allskonar fólki. Þess vegna þurfum við þína aðstoð, með því að smella á „Taka þátt“ þá getur þú lesið upp nokkrar setningar og lagt „þína rödd“ af mörkum. Við viljum sérstaklega hvetja fólk sem hefur íslensku sem annað mál að taka þátt. Það er á okkar valdi að alltaf megi finna svar á íslensku.

Samrómur hófst í október 2019 og hingað til hafa um **19** þúsund manns lesið rúmlega **1.702** klukkustundir eða **1.169.201** setningar. Hægt er að lesa meira um verkefnið hér. [Lesu meira hér.](#)

El cas del Ruandès



- **Objectiu:** construir un sistema de reconeixement i síntesi de la veu per fer més accessible la digitalització del país
- **Com:** *Umuganda*
- **Resultat:** Ara el Ruandès ocupa la segona posició en hores enregistrades al CommonVoice
 - 1894 hores enregistrades

Muhire, R. (2020). How Rwanda is Making Voice Tech More Open. Mozilla Foundation.

<https://foundation.mozilla.org/en/blog/how-rwanda-making-voice-tech-more-open/>

El cas del Ruandès



“Sense comunitat,
no hi ha dades”

Muhire, R. (2020). How Rwanda is Making Voice Tech More Open.
Mozilla Foundation.

<https://foundation.mozilla.org/en/blog/how-rwanda-making-voice-tech-more-open/>

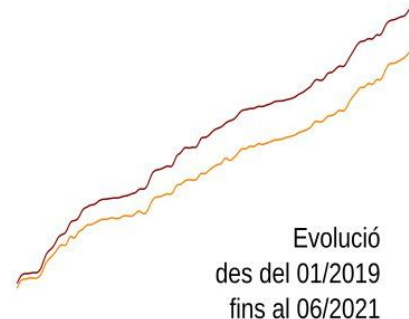
COMMON VOICE EN CATALÀ

Juliol 2021*



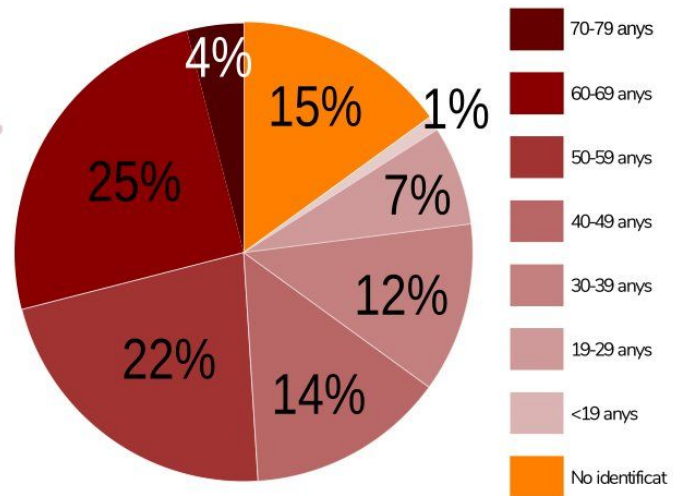
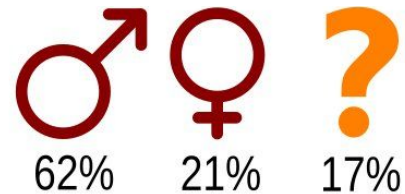
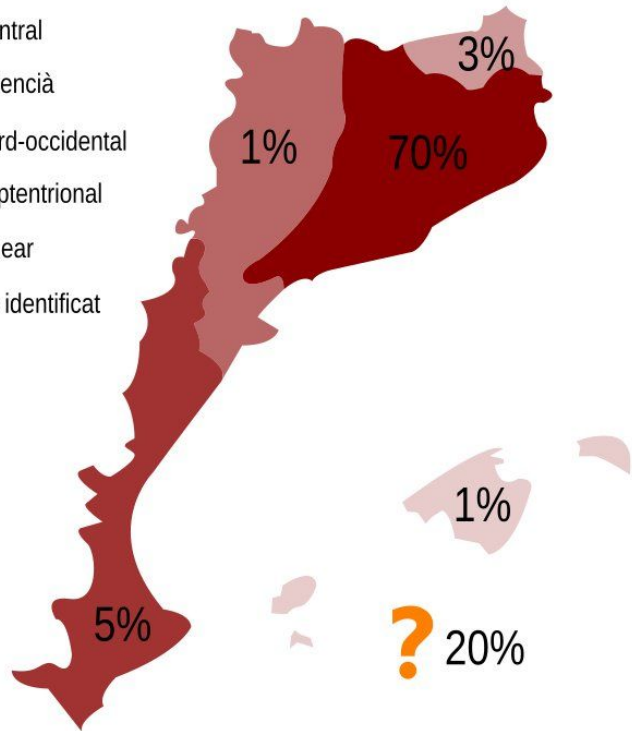
	Progrés	1er objectiu	Diferència respecte al mes anterior
Hores enregistrades	917	1.000	+25
Hores validades	789	1.000	+20

6.105
col·laboradors
+20



<https://twitter.com/softcatala/status/1420822833245085702>

- Central
- Valencià
- Nord-occidental
- Septentrional
- Balear
- No identificat



* Dades dialectals, d'edat i de gènere de juliol de 2021

<https://twitter.com/softcatala/status/1420822833245085702>

Rànquing actual de CommonVoice

Llengua	Hores validades
Anglès	2015
Ruandès	1894
Alemanys	965
Català	789
Esperanto	748
...	
Basc	91



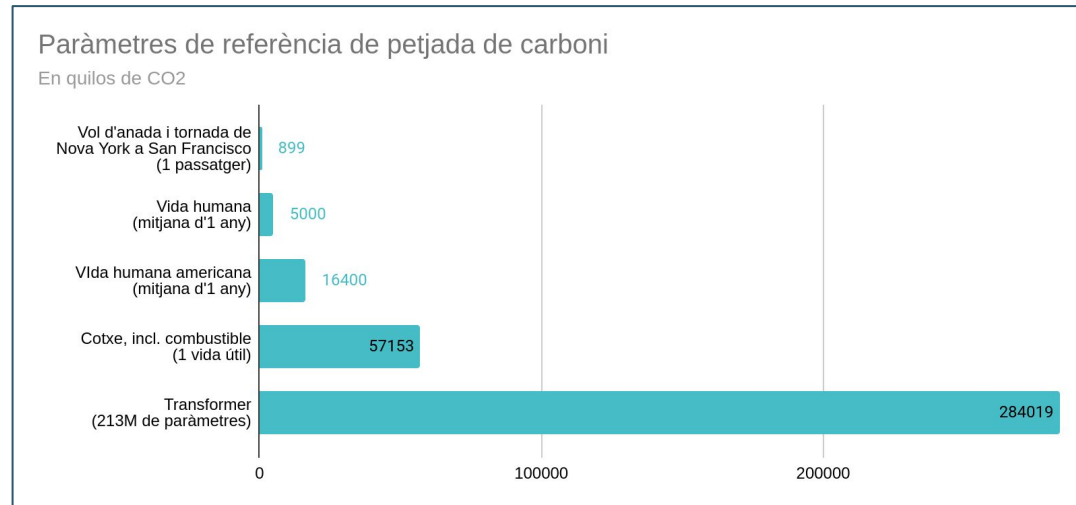
Models

Models

- La majoria d'algorismes són d'iniciatives **privades** comercials

Models

- La majoria d'algòrismes són d'iniciatives **privades** comercials
- Els models tenen una **petjada de carboni** enorme

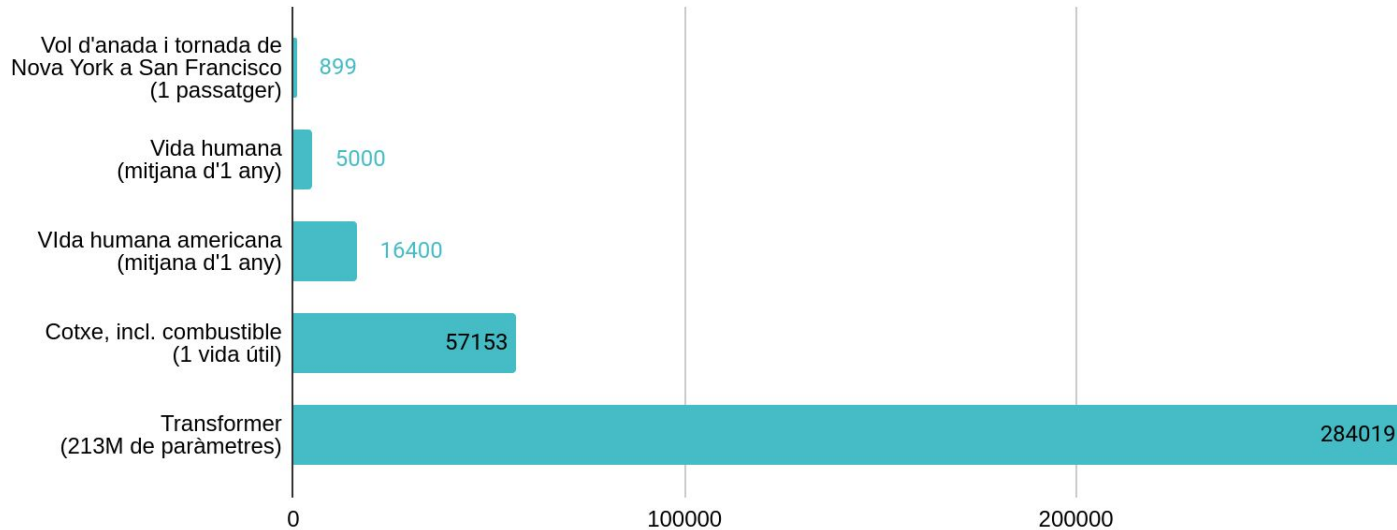


Strubell, E., Ganesh, A., & McCallum, A. (2019). Energy and policy considerations for deep learning in NLP. arXiv preprint arXiv:1906.02243.

Models

Paràmetres de referència de petjada de carboni

En quilos de CO2



Models

- La majoria d'algorismes són d'iniciatives **privades** comercials
- Els models tenen una **petjada de carboni** enorme

Què fem?

Models

- La majoria d' algorismes són d' iniciatives privades comercials
- Els models tenen una petjada de carboni enorme

Què fem?

1. **Alliberar** el model un cop l' hagus entrenat



HUGGING FACE

Models

- La majoria d' algorismes són d' iniciatives privades comercials
- Els models tenen una petjada de carboni enorme

Què fem?

1. Alliberar el model un cop l' hagi entrenat
2. Desenvolupar models **conjuntament**

BigScience



Models

- La majoria d' algorismes són d' iniciatives privades comercials
- Els models tenen una petjada de carboni enorme

Què fem?

1. Alliberar el model un cop l'hagis entrenat
2. Desenvolupar models **conjuntament**



BigScience Research Workshop @BigscienceW · Oct 7

Replying to @BigscienceW

TLDR: the BigScience project is working jointly on Data Governance, Tooling, and Sourcing for a large multilingual dataset with data in:

African languages
Arabic
Basque
Catalan
Chinese
English
French
Indic languages
Indonesian
Portuguese
Spanish
Vietnamese



Eines

Basades en Text



Apertium

<https://apertium.org>



Traductor

Corrector

Diccionaris i eines

Aplicacions

Recursos per a traductors

Ordinadors i mòbils

Basades en Text

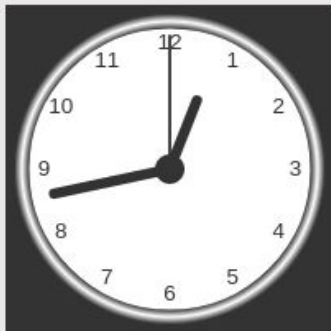
Introduïu una hora específica per a saber com es diu en català:

:

CONSULTA

Estableix l'hora actual

Heu introduït: **12:43:00**



La una menys disset (minuts) - Sistema de rellotge



Dos quarts i tretze (minuts) d'una - Sistema de campanar



Dos quarts i mig ben passats d'una - Sistema de campanar

tradicional

Basades en Text



Basades en Veu

catotron.collectivat.cat

CATOTRON

Síntesi de la parla obert i lliure en català



SOFTCATALÀ

assistent-cat/ona

An online virtual assistant in Catalan based on Mycroft. Visit <https://ona.assistent.cat> to try it out



1
Contributor

1
Issue

5
Stars

1
Fork



assistent.cat



Comunitats

Comunitats

Neix NLP ComuniCat, la comunitat oberta sobre tecnologies del llenguatge en català

La iniciativa vol agrupar tots els desenvolupadors del sector per compartir coneixements i potenciar la presència de l'idioma en l'àmbit digital

Categories: **Professionals**

Redacció Dimarts, 21 de setembre de 2021 | 13:47h



NLP ComuniCat
@NLPComuniCat

Edit profile

Som la comunitat oberta sobre el desenvolupament de les tecnologies del llenguatge en català! T'hi apuntes?

[Translate bio](#)

✉ nlpcomunicat@gmail.com docs.google.com/forms/d/e/1FAI...

📅 Joined September 2021

197 Following 229 Followers

1. Lingüística i tecnologia?

2. Llengües minoritzades

3. Casos d'èxit

4. Un cercle infinit

Un cercle infinit

Tecnologia



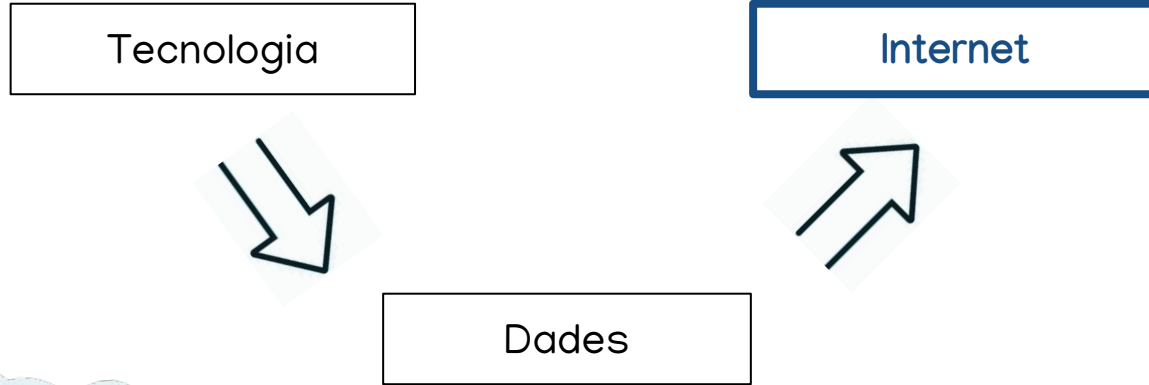
Un cercle infinit

Tecnologia

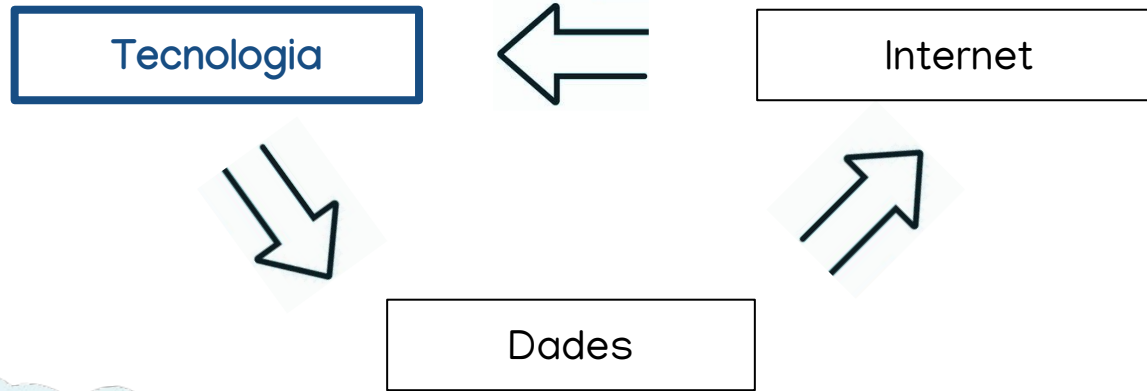


Dades

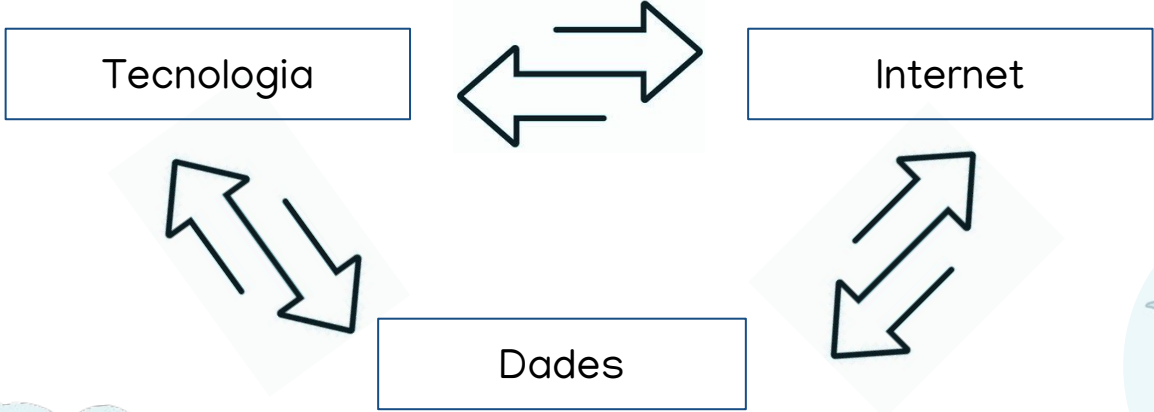
Un cercle infinit



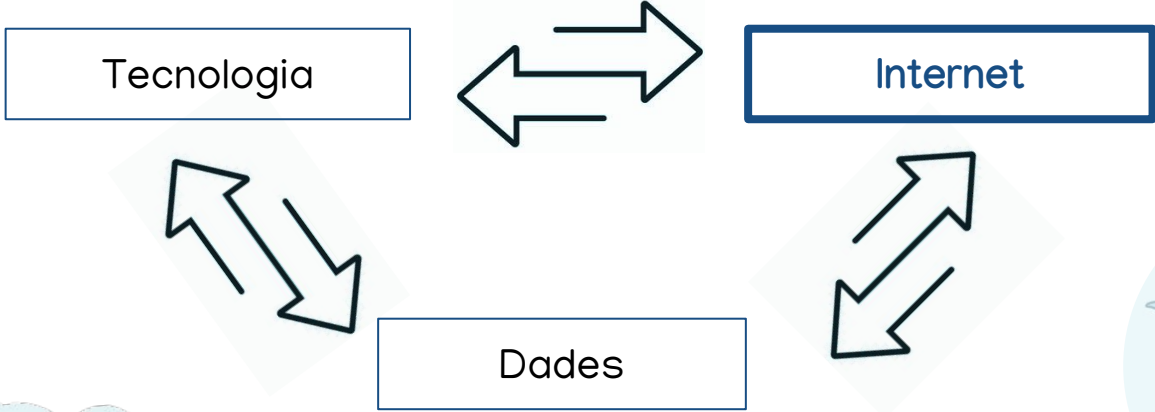
Un cercle infinit



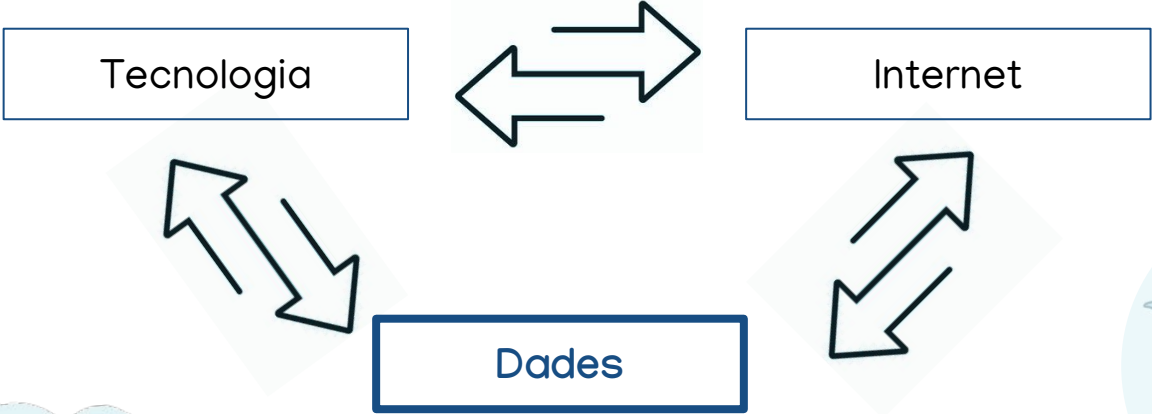
Un cercle infinit



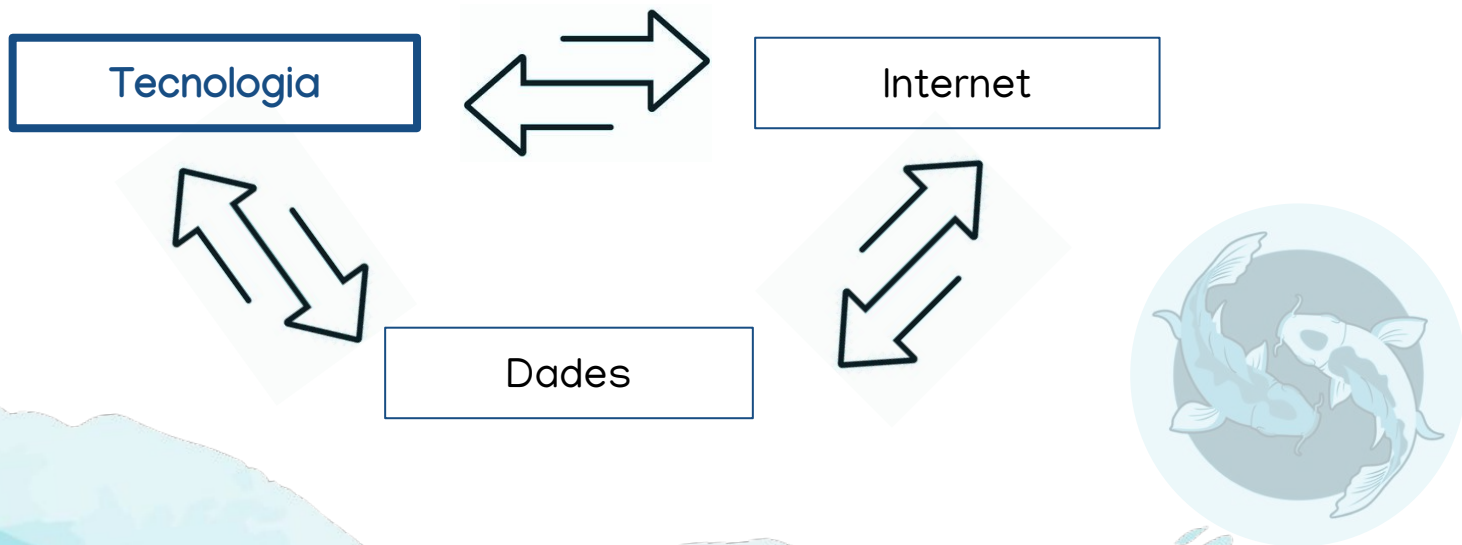
Un cercle infinit



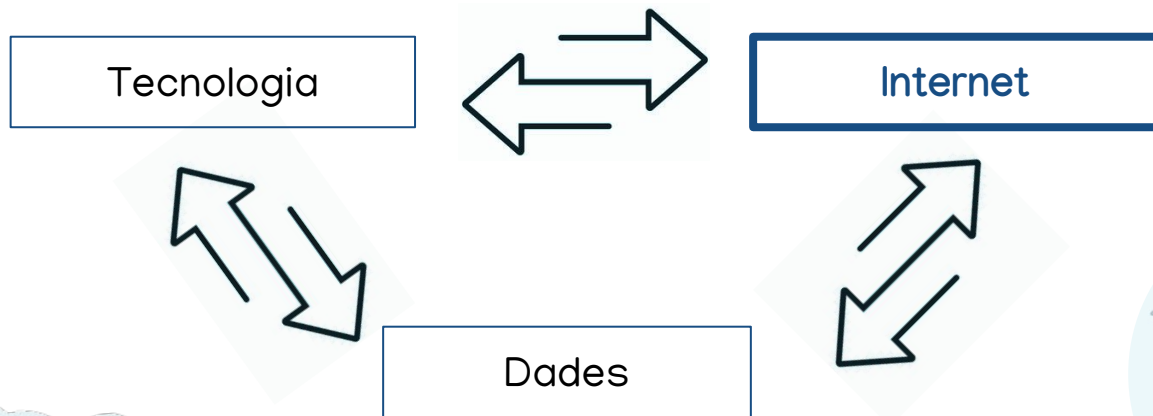
Un cercle infinit



Un cercle infinit



Un cercle infinit

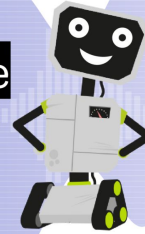


Com trenquem el cercle



VIQUIPÈDIA
L'enciclopèdia lliure

Common Voice
moz://a



NLP
comunicAT

BigScience



Sobirania lingüística

- poder fer **ús** de la nostra llengua en qualsevol àmbit
- reconeixement dels **drets lingüístics** a nivell individual i col·lectiu, de la llengua i totes les seves variants

*El que cal [...] no és simplement una millor protecció de les minories i més concessions en termes de autonomia regional, però, en definitiva, una **reconsideració exhaustiva del significat de la sobirania** en línies que superen les dicotomies profundament arrelades i jeràrquiques en què es basa la mateixa **existència de “majories” i “minories”***

el futur de la nostra llengua
és a les nostres mans



Eskerrik asko!