# Language Technology
# for Language Communities:
# An Overview based on Our Experience

Ixa Group (Iñaki Alegria, Kepa Sarasola)

Haskòlinn Í Reykjavik

www.ixa.eus

# Scheme

- To explain **our experience** on Language Technology for Basque

  Perhaps some ideas **can be useful for Icelandic,**
  as Icelandic is a reference for us in the Basque Country


- **Sharing and discussing** these ideas with you


- **Technology** may be useful for language development but…

  A **core work** should be implemented before (or in parallel):
  - **Standardization**
  - **Digital Open Contents** and
  - **Open Source Software**

# Scheme

- Introduction
    - "Ixa" research goup
    - Basque language
    - Basic concepts, basic strategies
    - Languages with scarce resources

- Basic tools and their applications

    Corpus, dictionaries, morphology, lemmatizer

- Language models.
    Good news for less-resourced languages (2020)

- Spanish and Basque plans for Language Technology

- Other associations working to promote the use of Basque

- Discussion

HÁSKÓLINN Í REYKJAVÍK
REYKJAVIK UNIVERSITY

ixa

eman ta zabal zazu

Universidad          Euskal Herriko
del País Vasco     Unibertsitatea

**People**

**Results**

**Master**

Master Ofiziala
Official Master's Degree

Hizkuntzaren
Azterketa eta
Prozesamendua (HAP)

Language
Analysis and
Processing (LAP)

http://ixa.si.ehu.es/master

**Language Technology Applications**

Information Retrieval, Information Extraction and Question Answering
Papers; Projects: Kyoto, paths, Lcloud, opener, skater and Know2; Demo: Ihardetsi (QA system)

Machine Translation
Papers; Project: OpenMT-2, Takardi, qtleap; Demo: Opentrad-Matxin (Spanish to Basque MT system)

Language learning
Papers; Project: Irakazi

**Linguistic processors**

Morphology
Papers; Project: BER2TEK; Demos: Morfeus, Eustagger

Syntax-Morphosyntax
Papers; Project: BER2TEK; Demos: Zatiak (chunker), Maltixa (statistical parser)

Lexicography-Semantics
Papers; Project: Kyoto and Know2; Demos: Know2's demos, Eihera (name entities)

**Linguistic Resources**

Corpus
Papers; Project: Lexikoaren behatokia ; Demos: ZT, Ancora-EPEC , EuSemcor

Dictionaries
Papers; Project: BER2TEK; Demos: EDBL (lexical database), Xuxen (spelling checker)

Ontologies
Papers; Project: Kyoto, Know2 and WNTERM; Demo: Basque Wordnet

# Ixa group (1988) - HiTZ Center (2020)

- **ixa.eus**
- **35 year** working on Language Technology
- **Basque**-centred research group…   but also other languages
- **Multidisciplinary**: computing, linguistics...
- **Text**-based resources and tools (speech in collaboration)
- **3 levels:** resources, basic tools, applications
- **Local**                    and   **Global**

  Basque **community** and   International research **community**
- **Collaboration**: Basque academy, Governments, lexicography centres, publishers, schools...
- **Alternative forums**:
     Basque Summer University,     Wikipedia,    NGOs...

# **Hitz** research center (hitz.eus)

- *Ixa* (text) and *Aholab* (speech) groups works together facing to new challenges (AI) collaboration with government, companies and associations

| Information Extraction and Information Retrieval | Machine Translation | Text Analysis | Speech and audio processing |
|---|---|---|---|
| Human-Computer Interaction | Speech and Language Resources | Medical and Legal domains | Digital humanities and education |

HÁSKÓLINN Í REYKJAVÍK
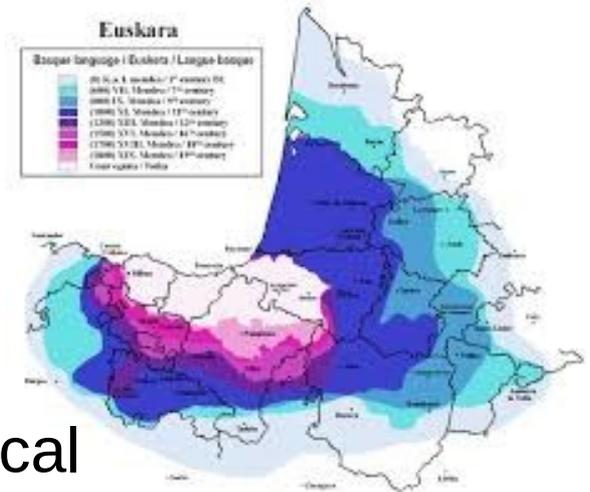REYKJAVIK UNIVERSITY

ixa

eman ta zabal zazu
Universidad del País Vasco   Euskal Herriko Unibertsitatea

# Basque language



Basque Country

Dialects

Historical regression

# Basque language

- Old language (pre-Indo-European)
- Near to 1 million speakers (~35%)
- Revitalization (1960-2023):
  - 1960: Basque schools (immersion)
  - 1968: Standard Basque
  - 1980: New political rule (before it was forbidden)
- Still declining on France side (non official status)
- Newspaper, TV channel, some radio channels and magazines
- No monolingual speakers

# Basque language

- Four grammar-cases in Icelandic

| M1 | Singularra | Plurala | | Singularra | Plurala |
|---|---|---|---|---|---|
| **Nominatiboa** | hest*ur* | hest*ar* | | hest*u*r**inn** | hest*a***nir** |
| **Akusatiboa** | hest | hest*a* | | hest**inn** | hest*a***na** |
| **Datiboa** | hest*i* | hest*um* | | hest*i***num** | hest*u***num** |
| **Genitiboa** | hest*s* | hest*a* | | hest*s***ins** | hest*a***nna** |

- 16 cases in Basque (en.wikipedia)

| Case/Number | Singular | Plural | Undetermined |
|---|---|---|---|
| Absolutive | liburu-a-Ø | liburu-ak | liburu-Ø |
| Ergative | liburu-a-k | liburu-e-k | liburu-k |
| Dative | liburu-a-ri | liburu-e-i | liburu-ri |
| Local genitive | liburu-ko | liburu-e-ta-ko | liburu-tako |
| Possesive genitiveGenitive | liburu-a-ren | liburu-e-n | liburu-ren |
| Comitative (with) | liburu-a-rekin | liburu-e-kin | liburu-rekin |
| Benefactive (for) | liburu-a-rentzat | liburu-e-ntzat | liburu-rentzat |
| Causal (because of) | liburu-a-rengatik | liburu-e-ngatik | liburu-rengatik |
| Instrumental | liburu-a-z | liburu-etaz | liburu-taz |
| Inessive (in, on) | liburu-a-n | liburu-e-ta-n | liburu-tan |
| Ablative (from) | liburu-tik | liburu-e-ta-tik | liburu-tatik |
| Allative (where to: 'to') | liburu-ra | liburu-e-ta-ra | liburu-tara |
| Directive ('towards') | liburu-rantz | liburu-e-ta-rantz | liburu-tarantz |
| Terminative (up to) | liburu-raino | liburu-e-ta-raino | liburu-taraino |
| Prolative | liburu-tzat | | |
| Partitive | liburu-rik | | |

# Basic subjects in Language Technology
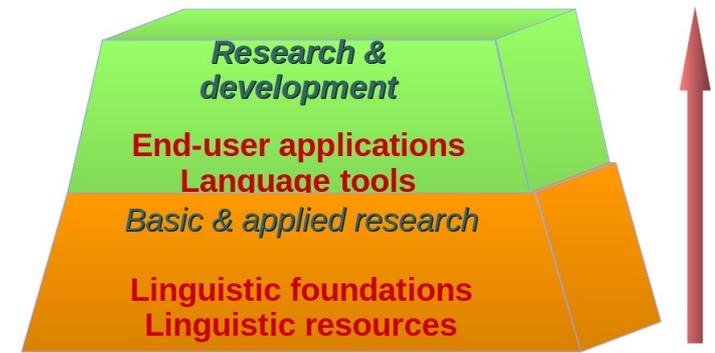
Looking for **texts in Basque**...

- Standardization (1968)
- (Digital) Contents (school books, small dictionaries, magazines…)
    - **Schools** (*Ikastolak*) → **University**
- Open/Free software / open contents
- Wikimedia / Wikipedia
- **Digital community**: websites, blogs, social networking, media...
- Need of **incremental design** and development of technology: **foundations** (resources), **tools**, and **applications**

# Less resourced languages

▪ Typology based on (digital) resources

# Languages with scare resources

- Typology based on (digital) resources
- Associations:



SIGUL:
ISCA Special
Interest Group.
Under-resourced
Languages



**English**

1 language

Best position
in all HLT applications
and resources

**Central languages (top 10 languages)**

10 languages

Relevant position
in all HLT applications

**Languages with any HLT application**

60 languages

**Languages with any lexical resource
in Internet**

250 languages

7.000 languages

**All the world languages**

# Basic resources

- Corpora (digital texts)
- Dictionary (better a digital one)
- Normative grammar (even in paper)

# Basic tools and their applications

- On-line dictionary              → Games

- Morphology                      → Spelling corrector

                                  → Language Learning

- Lemmatizer/POS_tagger → Search engine

- Text normalization
  and Machine translation

HÁSKÓLINN Í REYKJAVÍK
REYKJAVIK UNIVERSITY

ixa

eman ta zabal zazu
Universidad        Euskal Herriko
del País Vasco     Unibertsitatea

# Corpora (digital texts)

- Collecting corpora is not easy ("digital only for paper" → **error**)
- Sources: publishers, schools and **Wikipedia**
  - Alternative way: "*web as a corpus"* techniques and OCR
- Problems: copyrights and formats (pdf, word...)
- Use: data for text mining and for **evaluation**
- An initial (small) digital corpus is a key start point
- Processes: enriching the dictionary, creating/testing the spelling checker, learning (small) language models...





HÁSKÓLINN Í REYKJAVÍK
REYKJAVIK UNIVERSITY

Universidad del País Vasco
Euskal Herriko Unibertsitatea

# **Dictionaries** (mono- or bi-lingual)

- Basic tool for students, journalists, translators and writers
- Historical evolution: Paper cards→MSWord →XML/TEI
- XML/TEI → **Multimedia**: CD/DVD, Web/phone, Paper
  - Unique maintenance → 3 products
- Integration: *Euskalbar* (browser)
- Some projects in collaboration:
  - Small dictionary for Nahuatl, Garabide NGO
    Scholar dictionary in Cuba (*DBE*), CLA institute in Santiago de Cuba
    Scrable and other games in Basque

# Morphology / Spelling corrector

- Computational morphology is compulsory for most of the languages:
  - Dictionary + word-grammar
- The spelling corrector is a key application
  - Basic tool for students, journalists, translators and writers
  - **Key for basque standardization**
- Integration/online: MSoffice, LibreOffice, Mozilla, Android…
- Basic tools:
  - *foma* and *hunspell* (free software)
- Projects: unified Basque, dialectal Basque, Quechua (Univ. Zurich and Cusco), Mapudungun (Chile)

# Lemmatizer / Search engine

- Stemming → Lemmatizer (morphology)→ **POS tagging** (learning)
  
  word    → stem/lemma(root)        → lemma in context
  
  *juego*   → *jugar*(V) / *juego*(N)        → *jugar*(V)

- used for information extraction
  
  + language identifier → Search engine

- A **manually annotated corpus** is needed to create a POS tagger

- Powerful tool for Information Retrieval and Information Extraction

- Some projects for Basque:

Universidad del País Vasco   Euskal Herriko Unibertsitatea

# Successful applications since 2018

- Machine translation
- Use of the local language in Health services
- Digital humanities
- Speech synthesis
- Speech recognition
- Conversational interfaces, chatbots
- ...

Universidad del País Vasco
Euskal Herriko Unibertsitatea

HiTZ
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

# Machine Translation


I translated to Basque 472 Wikipedia articles
(2021-2023)      :-)

- Precondition:   (Parallel) Corpora

- 2006 RBMT (Knowledge based)
  2009 SMT
  2018 NMT

  Collaboration with Elhuyar.eus

- Neural revolution
  - Not perfect, it needs postedition but the quality is very high
  - Something incredible five years ago
  - We have 4 free good translation-services via web for Basque
  - This opens new ways

    where the use of translation could be extended to many situations
  - New horizons for under-resourced languages

- In most of the cases cross-lingual learning is used, but good results are also obtained even only using monolingual corpora (Artetxe et al., 2020)

    → Nice for languages with few parallel resources

HÁSKÓLINN Í REYKJAVÍK
REYKJAVIK UNIVERSITY

# Spanish and Basque plans for Language technology

- Plan for the Advancement of Language Technology
  2015-2020      90 M€

## Spanish Plan for Artificial Intelligence

- 2021-2023
- Spanish Government.
- https://www.lamoncloa.gob.es/lang/en/presidente/news/Paginas/2020/20201202_enia.aspx

## Basque Plan for Language technology

- 2021-2025
- Also the Basque Government.
- https://www.euskadi.eus/gaitu-plan-de-accion-de-las-tecnologias-de-la-lengua-2021-2024/web01-ejeduki/es/

ixa

eman ta zabal zazu
Universidad del País Vasco
Euskal Herriko Unibertsitatea

HiTZ
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

# Good news for less-resourced languages

There have been significant advances, even for less-resourced languages, in several areas:

- lexicon extraction (Artetxe et al., 2019),
- morphology induction (Anastasopoulos&Neubig, 2019)
- POS tagging (Kim et al., 2017),
- machine translation (Artetxe et al., 2017)
- chatbots

(Artetxe et al., 2020)

In most of the cases cross-lingual learning is used, but good results are also obtained even only using monolingual corpora,

→ Nice for languages with few parallel resources

Universidad del País Vasco
Euskal Herriko Unibertsitatea

HiTZ
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

# Good news for less-resourced languages

(Agerri et al., 2020)

Word embeddings and pre-trained language models enabled improvements across most NLP tasks.


Unfortunately they are very expensive to train,

--> small companies and research groups tend to use big models provided by the big companies.


But our mono- & multilingual language BERT models have proven to be very useful in NLP tasks for Basque.

Eventhough they have been created:

▪ with a 500 times smaller corpus than the English one

▪ with a 80 times smaller wikipedia.

ixa

eman ta zabal zazu

Universidad
del País Vasco

Euskal Herriko
Unibertsitatea

HiTZ
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

# Good news for less-resourced languages

The original BERT language model for English was trained in 2018 using Google books corpus with 189 billion words. Almost 500 times bigger than the Basque one (384 millions).

| Source | Text type | Million tokens |
|---|---|---|
| Basque Wikipedia | Encyclopedia | 35M |
| Berria newspaper | News | 81M |
| EiTB Television | News | 28M |
| Argia magazine | News | 16M |
| Local news sites | News | 224.6M |

# Good news for less-resourced languages

(Otegi et al., 2020)

A multilingual language model
pretrained only for English, Spanish and Basque, using:

- The monolingual Basque model
- English Wikipedia (2.5 Gword)
- Spanish Wikipedia (650 Mword)
  (80 and 20 times bigger than the Basque Wikipedia)

Successful to transfer knowledge from English to Basque
in a conversational Question/Answering system

Better than the general Google's official mBERT model
(it covers too many languages, Basque is not well represented).

# Other personal projects

- Basque Summer University

  https://en.wikipedia.org/wiki/Basque_Summer_University

- Association around the Basque Wikipedia

  https://eu.wikipedia.org

- Garabide NGO

  https://www.garabide.eus/english/

- Others:

  Language schools and communities for practice (online also), association for localising open software (*Librezale*) …

  **Repositories** of talks (old people), bibliography in Basque, open literature and critics…

  **.eus** Internet domain

# UEU: Basque Summer University

- The main public university (UPV/EHU) is bilingual
- UEU works only in Basque: community
  - Association of professors, lecturers and students
  - Promoting Basque at the University: teaching, research, dissemination
  - Some official titles in collaboration with official universities
  - Courses, conferences, books, databases…
- Two new interesting projects:
  - **Ikergazte**: pre- and post-doctoral Basque **researchers** together
  - **Open university in Basque**: today is not possible in Basque



JAKINTZA-ARLOA

- [ ] Antropologia (9)
- [ ] Arte eta Letren zientziak (35)
- [ ] Bizi zientziak (32)
- [ ] Filosofia (23)
- [ ] Fisika (41)
- [ ] Geografia (1)
- [ ] Historia (28)
- [ ] Hizkuntzalaritza (68)
- [ ] Informazio eta Komunikazio zientziak (54)
- [ ] Kimika (38)
- [ ] Lurra eta espazioaren zientziak (7)
- [ ] Matematika (45)
- [ ] Medikuntza (20)
- [ ] Nekazaritza zientziak (2)
- [ ] Pedagogia (25)

Aritmetika Lehen Hezkuntzako irakasleentzat

Paperekoa 18 €
Digitala 11.7 €

ESKURATU

Eskola Modernoa

Paperekoa 17 €
Digitala 11.05 €

ESKURATU

Urratuen arrastoan. UEUk begirada literarioa hauspotu zuenekoa

Paperekoa 15 €
Digitala 9.75 €

ESKURATU

IKERGAZTE
NAZIOARTEKO IKERKETA EUSKARAZ

Irudi gehiago

# Wikipedia (open source text)

- No problems with copyrights and formats
- Applications:
    - **language models**
    - **text mining**
    - **Language Technology evaluation**
- Growing and growing
- 2017-2019
  Basque Government Education program
  for creating 1.000 basic articles
  in Basque Wikipedia
  for 12-16 years students
  created by students at the university
  :-))

2020    385.000

2017    250.000

Articles
in Basque
Wikipedia
year
by
year

| Year | Articles |
| --- | --- |
| 2017 | 250.000 |
| 2016 | 219.000 |
| 2015 | 200.000 |
| 2014 | 150.000 |
| 2013 | 130.000 |
| 2012 | 120.000 |
|  | 100.000 |
| 2011 | 60.000 |
| 2010 | 50.000 |
|  | 40.000 |
| 2009 | 30.000 |
|  | 25.000 |
| 2008 | 20.000 |
| 2007 | 10.000 |
|  | 5.000 |
| 2006 | |
| 2005 | |
| 2004 | 1.000 |
| 2003 | |
| 2002 | Hasiera |

ixa

eman ta zabal zazu

Universidad
del País Vasco

Euskal Herriko
Unibertsitatea

HiTZ
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

# Garabide: NGO for language cooperation



- Cooperation with language communities,
  specially Latin America and north of Africa

  *The Revitalization of the Basque Language as a Model for Other Minority Languages in the World*

- https://issuu.com/di-da/docs/garabide-english

- Cooperation in pedagogy, sociolingusitics, technology...

- Master on language management (in Spanish) for leaders

# Basque in the digital era

- The Basque has achieved a good position
  "a successful case"

- Basque (data from a research in 2014 [1])
  - Translated books: 54$^{th}$ among 247
  - Wikipedia: 47$^{th}$ among 134 (34$^{th}$ in 2022)
  - Twitter: 36$^{th}$ among 70

- Globalization and Internet (young people) →
  main languages are more and more prevalent (Netflix...)
  young people learn Basque, have tools, but many do not use them

[1] https://www.pnas.org/doi/pdf/10.1073/pnas.1410931111

# Discussion

- Sustainability / Cost:
- Planning is a need for sustainability
- Fast development could be *expensive*

- Integration on commercial software / open software

  Many popular applications are proprietary and

  the decision to add new languages depends on the company

  → interesting experience in Basque.
  Microsoft become more interested
  after the localization and integration on free software

  → spin-off company for apps

- Standardization/dialect/language/alphabet...

  Each community has to decide...

  But standardization is crucial for text processing!

  A spelling corrector can be a "de facto" proposal

HÁSKÓLINN Í REYKJAVÍK
REYKJAVIK UNIVERSITY

ixa

eman ta zabal zazu
Universidad
del País Vasco
Euskal Herriko
Unibertsitatea

# Discussion

How to face the development of "ChatGPT" for Basque?

▪ Using commercial software?

https://openai.com/customer-stories/government-of-iceland

**English to Icelandic**

**Prompt**
Hvað heitir Donald Duck á íslensku?
What is Donald Duck called in Icelandic?

**GPT-3 completion**
Donald Duck heitir á íslensku Donaldi Kjáni
Donald Duck is called in Icelandic Donaldi the Fool
*Incorrect*

**ChatGPT completion**
Donald Duck heitir sama nafni á íslensku og ensku
Donald Duck has the same name in Icelandic and English
*Incorrect*

**GPT-4 completion**
Donald Duck heitir Andrés Önd á íslensku
Donald Duck is called Andrés Önd in Icelandic
*Correct*

▪ Or creating our open software by our own? Finland

Universidad del País Vasco
Euskal Herriko Unibertsitatea

# Takk fyrir

# Eskerrik asko

# Thank you

*It is always easier to open way
when you have close references*

*Technology makes it easier    :-))*

Spurningar, skoðanir...

Galderak, iritziak...

Questions, opinions...

# Language Technology
# for Language Communities:
# An Overview based on Our Experience

Ixa Group (Iñaki Alegria, Kepa Sarasola)

Haskòlinn Í Reykjavik

www.ixa.eus

HÁSKÓLINN Í REYKJAVÍK
REYKJAVIK UNIVERSITY

ixa

eman ta zabal zazu
Universidad     Euskal Herriko
del País Vasco  Unibertsitatea

# Extra slides (complementary)

- Some Useful applications
  - Aditu, bilingual speech recognition
  - Interprest,
    low cost and portable interpretation services
  - Bidaide. Web service to create and provide multingual contents.

# Aditu, bilingual speech recognition

# Interprest,
## low cost and portable interpretation services

Interpreter



Public

Turin 2021

# Bidaide

Web service to create and provide multingual contents.

In museums, routes or buildings.

Visitors read or listen to explanations on their own mobile

Basque **Center for Language Technology**

# Bidaide