# Language Technology
# for Language Communities:
# An Overview based on Basque Experience 2020

**Kepa Sarasola**, Olatz Perez-de-Viñaspre. Iñaki Alegria
Ixa group. HiTZ Center
University of the Basque Country

*Happily year after year we feel that*
   *the recovery processes of Welsh and Basque*
      *go hand in hand,*

*They are quite parallel processes.*

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

*Technology makes it easier    :-))*

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

- **32 years** working on Language Technology

- **Basque**-centred research group but also other languages

- **Multidisciplinary**: computer scientists, linguistics...

- **Text**-based resources and apps (speech with Aholab group)
- **3 levels: r**esources, basic tools, applications (with Elhuyar)
- **Local**                                  and        **Global**
- Basque **community**     and      International research **community**
- **Collaboration**: Basque academy, Elhuyar Foundation, publications
- **Alternative forums**: Basque Summer University, NGOs...

ixa

eman ta zabal zazu

Universidad
del País Vasco

Euskal Herriko
Unibertsitatea

HiTZ
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

ixa taldea

People
Members

Products

**Education
Masters**

Master Ofiziala
Official Master'sDegree

Hizkuntzaren
Azterketa eta
Prozesamendua (HAP)

Language
Analysis and
Processing (LAP)

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

Bangor 2020

Universidad del País Vasco
Euskal He... Unibertsita...

**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

# Successful applications since 2018

- Machine translation
- **Use of Basque in Health services**


- Digital humanities
- Speech synthesis
- Speech recognition
- Conversational interfaces, chatbots
- ...

ixa

eman ta zabal zazu

Universidad
del País Vasco

Euskal Herriko
Unibertsitatea

HiTZ
Hizkuntza Teknologiako Zentroa
Basque **Center for Language Technology**

# Successful applications since 2018

- Machine translation
- Use of Basque in Health services
  - **Olatz Perez de Viñaspre**
    *Basque, NLP and Clinical domain*



Olatz Perez de
Viñaspre

- Digital humanities
- Speech synthesis
- Speech recognition
- Conversational interfaces, chatbots
- ...

**ixa**

eman ta zabal zazu
Universidad
del País Vasco
Euskal Herriko
Unibertsitatea

**HiTZ**
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

# Basque

- Old language (pre-Indo-European)
- ~ 800.000 speakers (25-30% of people)
- Standardization in 1968 but rich dialects
- New political rule a bit later
- Basque schools: Ikastolak
- Declining on France side (non official)
- No monolingual speakers

ixa

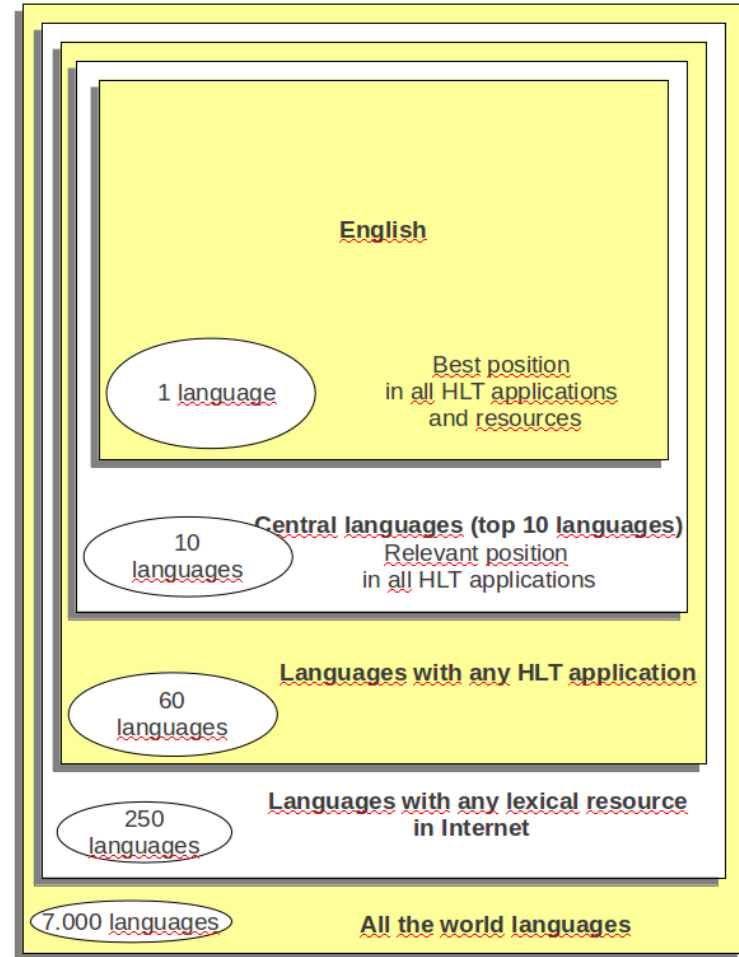eman ta zabal zazu

Universidad del País Vasco  Euskal Herriko Unibertsitatea

HiTZ
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

# Foundations

Collecting texts...

- Standardization (1968)
- (Digital) Contents (school books, small dictionaries…)
  - Readers: School (Ikastolak) → University
- Open/Free software   /   open contents
- Wikimedia / **Wikipedia**
- Digital community
- Need of incremental design and development of language foundations, tools, and applications

***Research & development***

**End-user applications**
**Language tools**
*Basic & applied research*

**Linguistic foundations**
**Linguistic resources**

Universidad
del País Vasco

Euskal Herriko
Unibertsitatea

HiTZ
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

# Under-resourced languages (levels)

- *Basic LAnguage Resource Kit* (BLARK) *Krauwer (2003)*

- Typology based on (digital) resources

- Associations:

  Special Interest Group:
  Under-resourced Languages (SIGUL)



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

mercator
European Research Centre on Multilingualism and Language Learning

SOAS University of London — 100 Years
World Languages Institute

Foundation for Endangered Languages

SaLT MIL
Speech And Language Technology for Minority Languages

EUROPEAN LANGUAGE RESOURCES ASSOCIATION
ELRA elDa

ixa

eman ta zabal zazu
Universidad del País Vasco
Euskal Herriko Unibertsitatea

HiTZ
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

# Basic resources

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

ixa

eman ta zabal zazu
Universidad
del País Vasco
Euskal Herriko
Unibertsitatea

HiTZ
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

# Corpora (digital texts)

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

ixa

eman ta zabal zazu

Universidad del País Vasco
Euskal Herria Unibertsitatea

HiTZ
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

# 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

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
       1.000
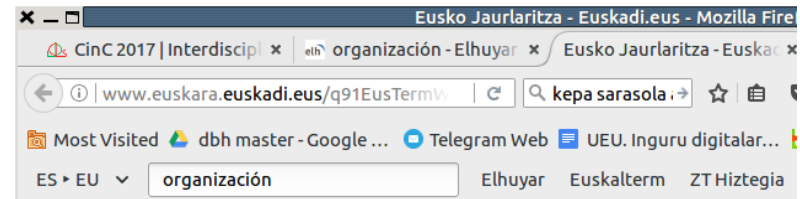2004
2003
2002 — Hasiera

# Basic tools and applications

- On-line dictionary      → Games

- Morphology              → Spelling corrector

- Lemmatizer/POS_tagger → Search engine

- Machine translator or Normalizator

# **Dictionaries** (mono- or bilingual)

- Basic tool for students, journalists and writers
- Historical evolution: Paper cards → Word Proc. → XML/TEI
- XML/TEI → **Multimedia**: DVD, Web/phone, paper
    - Unique maintenance → 3 products
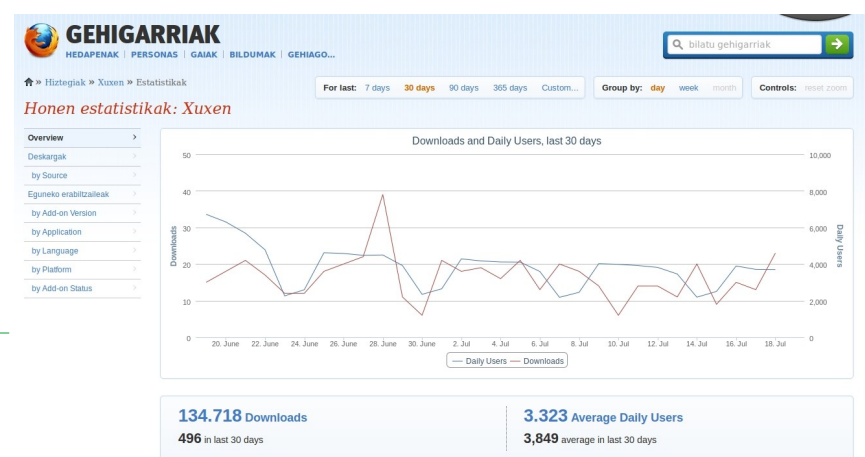- Integration: *Euskalbar* (browser)
- Some projects:

    Garabide NGO (Nahuatl) and Cuba (*DBE*).
    Scrable in Basque

Universidad
del País Vasco

Euskal Herriko
Unibertsitatea

HiTZ
Hizkuntza Teknolo
Basque **Center** for Language Technology

# Morphology / Spelling corrector

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

Universidad del País Vasco    Euskal Herriko Unibertsitatea    | Basque Center for Language Technology

# 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

**HITZ**
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

# Successful applications since 2018

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

eman ta zabal zazu

Universidad
del País Vasco

Euskal Herriko
Unibertsitatea

HiTZ
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

# Machine Translation

- Not perfect, it needs postedition,
  but the quality is very high
- Something incredible five years ago.
- From Basque to Spanish but also to French, English…
  … in both senses
- We have 4 free translation-services via web for Basque
- We are in a new world
  where the use of translation can be enormous
- We need human translators and with their help translation
  could be extended to many new fields and situations
- New horizons for under-resourced languages

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

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

(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).

# Spanish Plan for Language technology

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

## Spanish Plan for Artificial Intelligence

- 2021-2027
- Now being designed by the Spanish Government.

## Basque Plan for Language technology

- 2021-2025  (?)
- Now being designed by the Basque Government.

ixa

eman ta zabal zazu
Universidad          Euskal Herriko
del País Vasco      Unibertsitatea

HiTZ
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology

# Language Technology
# for Language Communities:
# An Overview based on Basque Experience 2020

**Kepa Sarasola**, Olatz Perez-de-Viñaspre. Iñaki Alegria
Ixa group. HiTZ Center
University of the Basque Country

Universidad del País Vasco
Euskal Herriko Unibertsitatea

**HiTZ**
Hizkuntza Teknologiako Zentroa
Basque Center for Language Technology