

# 語音詞典

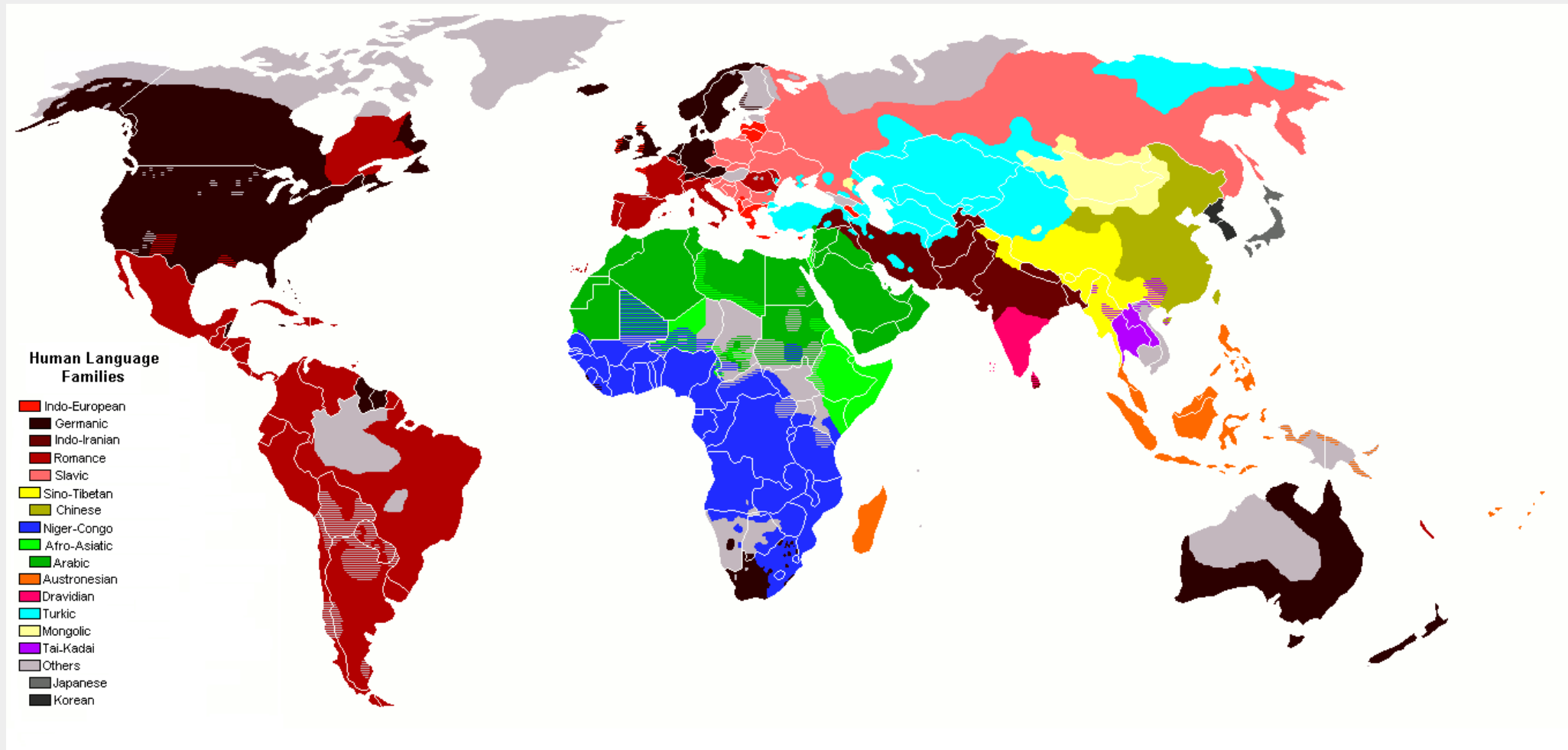
## Recording voices and local languages with Lingualibre

Hugo Lopez

[hugo.lopez@univ-toulouse.fr](mailto:hugo.lopez@univ-toulouse.fr)

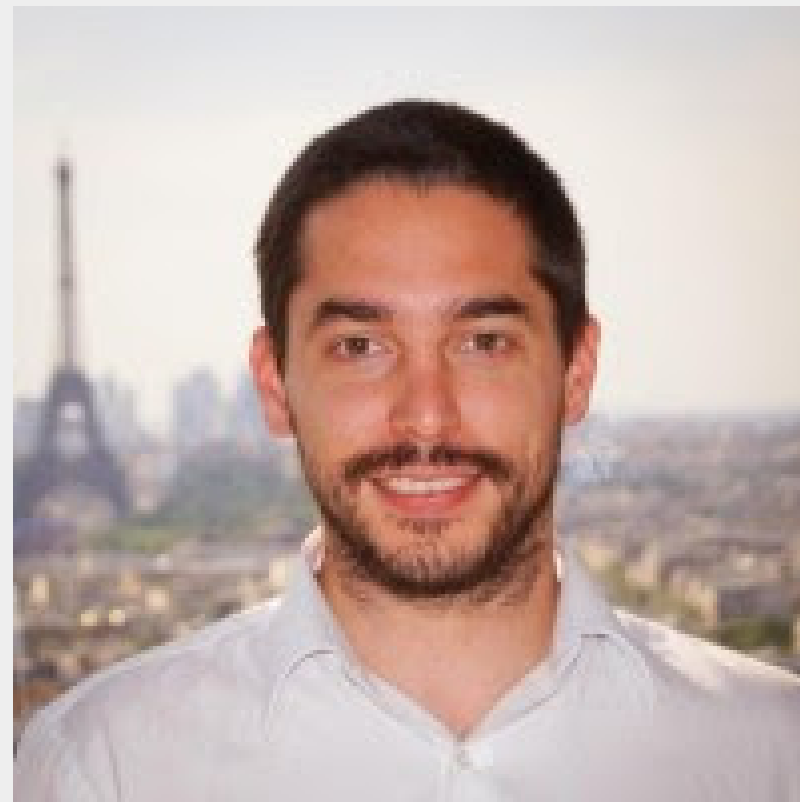
# Outline

- Context
- LinguaLibre & objectives
- Demo of the tool (10mins)
- Current progresses, limits & biases
- Q&A



# Hugo Lopez

- Elearning and language professional
- Open education resources
- Teach new technologies to PhD
- Gascon, not French !



# Attendance

What affinity do you have with :

- Language diversity and minority languages ?
- Web applications ?
- Open source & its community management ?

# LINGUA LIBRE

# What is Lingua Libre

- Wikimedia's open source recording tool
- ...to document languages diversity
- 193 languages, 900,000 records
- For language e-learning services

WIKIPEDIA  
The Free Encyclopedia

Yug

Lingua Libre [edit]

9 languages

Article Talk

Read Edit View history Page Tools

From Wikipedia, the free encyclopedia

**Lingua Libre** is an online collaborative project and tool by the [Wikimédia France](#) [fr] association, which aims to build a [collaborative](#), [multilingual](#), [audiovisual speech corpus](#) under a [free license](#).

**Description** [edit]

Lingua Libre enables the recording of [words](#), [phrases](#) or [sentences](#) of any language, oral ([audio recording](#)) or signed ([video recording](#)).

Words are presented to the speaker in the form of a list, created on the spot, in advance, or by reusing an existing Wikimedia category. The speaker simply reads the word displayed on the screen, and the software moves on to the next word when it detects a silence after the read word.<sup>[1]</sup> This principle, borrowed from the open source software [Shtooka](#) [fr] recorder with the help of its creator, Nicolas Vion, makes it possible to record several hundreds of words per hour. The recordings are then uploaded automatically from the web client to the [Wikimedia Commons](#) media library.

In spring 2021, Lingua Libre was offline due to a fire in Strasbourg,<sup>[2]</sup> but no audio recordings were lost.<sup>[3]</sup>

**Use of the recordings** [edit]


The recordings can be consulted either on Lingua Libre or on [Commons](#). They are mainly used on other Wikimedia projects, for example to illustrate entries on [Wiktionaries](#) or proper nouns in Wikipedia articles.<sup>[1]</sup>

The re-use of the recordings in a language teaching context is envisaged. Language learners can freely download pronunciations and use them on GoldenDict, a popular dictionary software.<sup>[4]</sup> Thus, audio recordings can be used as "Pronunciation Dictionaries" on GoldenDict without needing internet connection.

The recordings are also reused in [Natural Language Processing](#) projects, for example to drive Mozilla's [DeepSpeech](#) speech recognition engines.<sup>[5]</sup>

**Versions** [edit]

**Lingua Libre**



Overview of the website's homepage in December 2020

<b>Type of site</b>	Language recording tool, Online linguistic media library
<b>Available in</b>	Multilingual
<b>Owner</b>	Wikimédia France [fr]
<b>Created by</b>	Wikimedia France and the Wikimedia community
<b>URL</b>	lingualibre.org
<b>Advertising</b>	No
<b>Commercial</b>	No
<b>Registration</b>	Optional, but required for recording
<b>Launched</b>	August 2016; 6 years ago
<b>Current status</b>	Active
<b>Content license</b>	Creative Commons Attribution-ShareAlike 4.0 International (CC BY-SA 4.0)

# Language diversity

- Languages are core and center of our identity, human experience, memories, kindness.
- 7,000+ languages
- 3,500 (half) will die out this century



# Language diversity

- Languages are often segregated, politicized
- Case of France vs Taiwan

# Language diversity

“

The diversity of world languages, their words, expressions, voices, are poorly documented and accessible. Lingualibre.org allows us to record languages vocabulary and audio dictionaries at large scale, in an easy and quick fashion (800 audio/hour).

”

# LINGUA LIBRE (10MINS)

# Oral languages' learning chain

Learners

Speakers

e-services

Data

Lingualibre Studio

**#Section 2**  
Sinogrammes (11) : 中, 国, 日本, 王, 马, 很, 不, 大, 小, 吗.

**2** 中

**1** 中 zhōng  
milieu, moyen ; frapper juste  
口+丨  
● flèche au milieu de la cible  
#S2 · #S2中 · ↗

**3** Ecrit/Audio **4** Audio **5** Effacer

国 guó  
pays, royaume  
口+玉  
● une enceinte protectrice autour une pierre de jade  
#S2 · #S2国 · ↗

Ecrit/Audio Audio Effacer

Latest recordings

仏文学  
Japanese - CKali

英文学  
Japanese - CKali

Record a voice

- Tutorial
- Speaker
- Details
- 4** Studio
- Publish

Studio

▶ accastillage

accostable  
accoster  
affluent  
affouillement  
aï  
aiguilles  
algues  
amateur  
amer  
amphibiotiques  
alluvion  
amérindiens  
amont

Click on below, then read the word aloud.

accastillage

Skip >

1 / 520

Learning

Audio & text

Sharing

# Lingua Libre: audio recording studio


**Record a voice**

- ✓ Tutorial
- ✓ Speaker
- ✓ Details
- 4 Studio**
- 5 Publish

### Studio



phénakistiscope

- phylargie
- phlegmon
- phoniatre
- photocellule
- phratrie
- phylarque
- piaffe
- piaillard
- piédouche
- piéride
- pignole
- pilage
- pilé

Click on  below, then read the word aloud

**phénakistiscope**

Skip to the next word >

  19 / 380

Cancel < Previous Next >

Page: [Lingualibre.org](https://lingualibre.org) Recording Studio

# Lingua Libre: Apps

#Section 2  
Sinogrammes (11) : 中,国,日,本,王,马,很,不,大,小,吗.

**2**




**1** 中 zhōng  
milieu, moyen ; frapper juste  
口+丨  
flèche au milieu de la cible  
#S2 · #S2中 · ↗

**3** **4** **5**

Ecrit/Audio Audio Effacer

**国**



**国 guó**  
pays, royaume  
口+玉  
une enceinte protectrice autour une pierre de jade  
#S2 · #S2国 · ↗

Ecrit/Audio Audio Effacer

## 萌典 [ edit source ]

Page Discussion 汉 漢 不转换 ▾

Tools ▾

From Wikipedia

Open View it!

**萌典**是一部由台灣自由軟體程式設計師唐鳳開發的數位化漢語詞典，是台灣開源社群g0v零時政府的專案之一。作為一部數位化漢語詞典，萌典除了收錄了十六萬筆的中華民國國語詞條之外，還收錄了兩萬筆台灣閩南語、一萬四千筆台灣客家語詞條，以及提供了漢語與英語、法語以及德語的對照。網站作者唐鳳將其以創用CC0協議釋放至公有領域<sup>[2]</sup>。除在線版外，萌典還提供有適用於Windows、macOS和Linux的桌面版，以及使用於Android和iOS的手機版。

**萌典**

g0v  
**萌典**  
g0v.tw



網站類型 數位化漢語詞典

List: Lingua Libre Apps. Future: Moedict.tw ?

# STATED OBJECTIVES

“ Document languages diversity and voices. ”

- Languages
- Accents
- Voices
- Genders

# PROGRESSES, LIMITS & BIASES

After 5 year and 900,000+ recordings, we would like to share past progresses, current analysis and future actions.



# In numbers

## Production

Languages gallery

## Contributors

Lingua Libre map

Log in and record few words

Search by language name

Languages (190)

Languages with over 20k recordings (10)

For activated languages [log in and start recording vocabulary](#), most languages have vocabulary lists ready to record : at Step 3, search "List::your\_ISO/Unix".

<p><b>Langue Française</b></p> <p>210M speakers worldwide</p> <p>Speakers: 405</p> <p>Gender split: 9155 ♂ 21 229 ♀</p> <p>Unique words vs recordings ratio: 174k 267k</p> <p>Recordings gender split: 19k 5k 243k</p> <p>CONTRIBUTE DOWNLOAD</p>	<p><b>ଓଡ଼ିଆ ଭାଷା</b></p> <p>35M speakers worldwide</p> <p>Speakers: 8</p> <p>Gender split: 91 ♂ 7 ♀</p> <p>Unique words vs recordings ratio: 96.4k 122k</p> <p>Recordings gender split: 73 0 122k</p> <p>CONTRIBUTE DOWNLOAD</p>	<p><b>Polszczyzna</b></p> <p>40M speakers worldwide</p> <p>Speakers: 26</p> <p>Gender split: 94 ♂ 1 21 ♀</p> <p>Unique words vs recordings ratio: 91.6k 94.3k</p> <p>Recordings gender split: 13.8k 1 80.5k</p> <p>CONTRIBUTE DOWNLOAD</p>
<p><b>বাংলা</b></p> <p>300M speakers worldwide</p> <p>Speakers: 19</p> <p>Gender split: 92 ♂ 17 ♀</p> <p>Unique words vs recordings ratio: 62k 67.2k</p> <p>Recordings gender split: 368 0 66.8k</p> <p>CONTRIBUTE DOWNLOAD</p>	<p><b>Esperanto</b></p> <p>2M speakers worldwide</p> <p>Speakers: 18</p> <p>Gender split: 90 ♂ 1 17 ♀</p> <p>Unique words vs recordings ratio: 29k 33.8k</p> <p>Recordings gender split: 0 3.9k 29.9k</p> <p>CONTRIBUTE DOWNLOAD</p>	<p><b>English</b></p> <p>750M speakers worldwide</p> <p>Speakers: 109</p> <p>Gender split: 923 ♂ 8 78 ♀</p> <p>Unique words vs recordings ratio: 28.8k 32.8k</p> <p>Recordings gender split: 2.2k 2.7k 27.9k</p> <p>CONTRIBUTE DOWNLOAD</p>



# In numbers

## Production

### Languages gallery

## Reuses

### 2022 review

Languages (190)

Search by language name

Languages with over 20k recordings (10)

*For activated languages [log in and start recording vocabulary](#), most languages have vocabulary lists ready to record : at Step 3, search "List: {your\_ISO}/Unilex".*

<p><b>Langue Française</b></p> <p>210M speakers worldwide</p> <p>Speakers: 405</p> <p>Gender split: 9155 (21) 229 (2)</p> <p>Unique words vs recordings ratio: 174k / 267k</p> <p>Recordings gender split: 19k / 5k / 243k</p> <p>CONTRIBUTE DOWNLOAD</p>	<p><b>ଓଡ଼ିଆ ଭାଷା</b></p> <p>35M speakers worldwide</p> <p>Speakers: 8</p> <p>Gender split: 91 (7) 7 (2)</p> <p>Unique words vs recordings ratio: 96.4k / 122k</p> <p>Recordings gender split: 73 / 0 / 122k</p> <p>CONTRIBUTE DOWNLOAD</p>	<p><b>Polszczyzna</b></p> <p>40M speakers worldwide</p> <p>Speakers: 26</p> <p>Gender split: 94 (1) 21 (2)</p> <p>Unique words vs recordings ratio: 91.6k / 94.3k</p> <p>Recordings gender split: 13.8k / 1 / 80.5k</p> <p>CONTRIBUTE DOWNLOAD</p>
<p><b>বাংলা</b></p> <p>300M speakers worldwide</p> <p>Speakers: 19</p> <p>Gender split: 92 (17) 17 (2)</p> <p>Unique words vs recordings ratio: 62k / 67.2k</p> <p>Recordings gender split: 368 / 0 / 66.8k</p> <p>CONTRIBUTE DOWNLOAD</p>	<p><b>Esperanto</b></p> <p>2M speakers worldwide</p> <p>Speakers: 18</p> <p>Gender split: 90 (1) 17 (2)</p> <p>Unique words vs recordings ratio: 29k / 33.8k</p> <p>Recordings gender split: 0 / 3.9k / 29.9k</p> <p>CONTRIBUTE DOWNLOAD</p>	<p><b>English</b></p> <p>750M speakers worldwide</p> <p>Speakers: 109</p> <p>Gender split: 923 (8) 78 (2)</p> <p>Unique words vs recordings ratio: 28.8k / 32.8k</p> <p>Recordings gender split: 2.2k / 2.7k / 27.9k</p> <p>CONTRIBUTE DOWNLOAD</p>

**Lingua Libre Bot, March 2023 (g)**

Local wiki	First edit	Edit count	%	Groups	Region of most beneficiaries
<b>Wiktionaries</b>					
<a href="#">fr.wiktionary.org</a>	12 June 2018	<a href="#">308,193</a>	54.7%	bot	Europe/France
<a href="#">ku.wiktionary.org</a>	30 November 2021	<a href="#">42,820</a>	7.6%	bot	Asia
<a href="#">oc.wiktionary.org</a>	16 December 2018	<a href="#">20,606</a>	3.7%	bot	Europe/France
<a href="#">shy.wiktionary.org</a>	8 September 2021	<a href="#">1,930</a>	0.34%	bot	Africa
<a href="#">or.wiktionary.org</a>	10 January 2023	<a href="#">249</a>	0.04%	—	Asia/India
All other projects	—	0	0%	—	World
<b>Technical projects</b>					
<a href="#">www.wikidata.org</a>	10 June 2018	<a href="#">62,045</a>		bot	Unclear
<a href="#">meta.wikimedia.org</a>	10 June 2018	<a href="#">6</a>		—	

**Olafbot, March 2023 (g)**

Local wiki	First edit	Edit count	%	Groups	Region of most beneficiaries
<b>Wiktionaries</b>					
<a href="#">pl.wiktionary.org</a>	4 March 2020	<a href="#">189,157</a>	33.6%	bot	Europe/Poland
<b>Technical projects</b>					
<a href="#">lingualibre.org</a>	26 February 2021	<a href="#">5,208</a>		bot	Unclear

# Qualitative

- Per language: large vs minorities
- Per gender
- Per age
- Per per area, income, etc.



# Languages typology & specifics

Large

Medium

Minorities

Resourced

Low resources

Written ?



# Languages typology & specifics

Large

Medium

Minorities

Resourced

Low resources

Written ?



# Languages typology coverage

Demographic <sup>[1]</sup>	World languages		Lili languages		Supported language's profile	Examples	Community's presence
	Number	Ratio	Number	Coverage			
Major (>30M)	30	0.5%	20	66%	Mostly major Western or Indian languages.	FRA, SPA, BEN	<b>Solid:</b> Several productive speakers. Sustained or periodic.
Large (1~30M)	350	5%	100	33%	Mostly Western languages, other notable languages	NLD, AFR, CAT	<b>Emerging:</b> One productive speaker, few not-retained speakers. Fragile.
Marginalized (<1M)	6500	94%	40	<1%	Mostly larger minorities in Western countries.	ATJ, BRE, EUS	<b>Contact point:</b> No productive speaker, one not-retained speaker. Below fragile.

# Key needs

- Mobile **e-dictionaries** for local communities
- For **revitalisation** ! Not documentation.
- Outreach to 6,500 local communities ?

V · T · E		Lingua Libre	[Collapse]
Lingualibre	<b>Repositories</b>	Record Wizard (Recording Studio) · SPARQL2DATA	
	<b>Documentations</b>	{Helps}	
	<b>Technical helps</b>	{Technicals} · Help:SPARQL series	
	<b>Reports</b>	Winter 2021-2022 Public Relations Campaign · Lingua Libre/2022 wishlist#Approach · 2022 Review · Lingua_Libre/Supports/Melody#10. langues régionales et minoritaires_dans les autres pays · Lingua Libre/Supports	
	<b>Referents</b>	Github, Phabricator, Userscripts (Yug, Poslovitch, Pamputt) · Bots (Pamputt, Poslovitch, Lepticed7, Yug) · Onboarding (Yug, WikiLucas, Pamputt) · Events, Outreach (Yug, Adélaïde) · Reports (Yug, Poslovitch) · Funding (Adélaïde)	
Lingua Libre/Signit	<b>Repositories</b>	Record Wizard <a href="#">↗</a> (Recording Studio with video track) · Lingualibre/Signit <a href="#">↗</a> (Firefox Web Extension)	
	<b>Documentations</b>	Minimal video recording tutorial for elegant signed videos <a href="#">↗</a>	
	<b>Reports</b>	2022.09.16 : 2022/Phase_1 light diagnosis · 2022.09.16 : 2022/Phase_2 volunteer coding report and detailed strategy · 2023.04.25 : 2023/Phase_1 freelance coding report and learning patterns · 2023.05-12 : 2023/Phase_2 outreach campaign challenges	
	<b>Referents</b>	<a href="#">Édouard Lopez</a> (concept design) · <a href="#">User:0x010C</a> (creator, developer) · <a href="#">Yug</a> (expansion, coordination)	

# Q&A



# Keep in touch

Role	Contacts
Lead, dev	<a href="#">User:Yug</a>
Wikimedia France	<a href="#">User:Adélaïde Calais WMFr</a>
Home	<a href="http://Lingualibre.org">Lingualibre.org</a>
Code	<a href="https://github.com/lingua-libre">github.com/lingua-libre</a>

# THANK !

## Credits

The following Wikimedia Commons images have been used:

File	Licence	Author
Human_Language_Families.png	CC-BY-SA	JFDP13
WikiLucas00_à_l'Institut_international_pour_la_Francophonie.jpg	CC-BY-SA	WikiLucas00
Daramlagon_and_Maerui-sama_session_on_Bikol_Wiktionary_and_Lingua_Libre_03.jpg	CC-BY-SA	Daramlagon
LinguaLibre_2022_Paris_Surui_training-03.jpg	CC-BY-SA	Yug
Forom_des_langues,_Toulouse,_2023-01.jpg	CC-BY-SA	Yug
Séance_Lingua_Libre_à_Cotonou_en_Mars_2021_-_Photo_17.jpg	CC-BY-SA	Fawaz.tairou
Lingua_Libre_Atikamekw_at_Wikimania_2017_Montreal.jpg	CC-BY-SA	Benoit Rochon