



HAL
open science

Faciliter l'accès des praticiens du Traitement Automatique des Langues à des jeux de données de langues rares : un deuxième point d'étape

Benjamin Galliot, Guillaume Wisniewski, Séverine Guillaume, Guillaume
Jacques, Alexis Michaud

► To cite this version:

Benjamin Galliot, Guillaume Wisniewski, Séverine Guillaume, Guillaume Jacques, Alexis Michaud. Faciliter l'accès des praticiens du Traitement Automatique des Langues à des jeux de données de langues rares : un deuxième point d'étape. Journées Jointes des Groupements de Recherche Linguistique Informatique, Formelle et de Terrain (LIFT) et Traitement Automatique des Langues (TAL), Nov 2022, Marseille, France. hal-03856363

HAL Id: hal-03856363

<https://hal.science/hal-03856363v1>

Submitted on 16 Nov 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - ShareAlike 4.0 International License

Faciliter l'accès des praticiens du Traitement Automatique des Langues à des jeux de données de langues rares : un deuxième point d'étape

Benjamin Galliot¹ Guillaume Wisniewski² Séverine Guillaume¹
Guillaume Jacques³ Alexis Michaud¹

(1) Langues et Civilisations à Tradition Orale (LACITO), CNRS - Sorbonne Nouvelle - INALCO

(2) Laboratoire de Linguistique Formelle (LLF), CNRS - Université de Paris

(3) Centre de Recherches Linguistiques sur l'Asie Orientale (CRLAO), CNRS - École des Hautes Études en Sciences Sociales - Institut National des Langues et Civilisations Orientales

b.g01lyon@gmail.com, Guillaume.Wisniewski@univ-paris-diderot.fr,
severine.guillaume@cnrs.fr, rgyalrongskad@gmail.com,
alexis.michaud@cnrs.fr

RÉSUMÉ

Nous présentons un outil logiciel qui permet d'assembler divers jeux de données de la collection [Pangloss](#) (archive ouverte multimédia de langues rares) en assurant la reproductibilité des expériences menées sur ces données. À titre d'exemple, deux corpus audio transcrits de langues minoritaires de Chine (japhug et na) sont proposés, sous une licence Creative Commons, comme corpus de référence pour des expériences en traitement automatique des langues, et comme exemples d'une chaîne de traitement généralisable à d'autres corpus d'archives ouvertes. L'enjeu global d'une mise à disposition de données de langues rares sous une forme aisément accessible et utilisable est de faciliter le développement et le déploiement d'outils de pointe en traitement automatique des langues naturelles pour tout l'éventail des langues humaines. Cet exposé, qui fait suite à une précédente communication sur le même thème, fait état de nouveautés dont un retour d'expérience concernant un dépôt auprès de Hugging Face.

ABSTRACT

Facilitating NLP specialists' access to language archive materials : an update

We present a software tool to assemble a great range of diverse datasets from the [Pangloss](#) collection (a multimedia open archive of under-documented languages). The tool ensures the reproducibility of experiments conducted on these data. As an example, two transcribed audio corpora of Chinese minority languages (Japhug and Na) are proposed, under a Creative Commons license, as reference corpora for experiments in Natural Language Processing, and as examples of a pipeline that can be generalized to other corpora from open archives. An overarching goal of making language archive data available in an easily accessible and usable form is to facilitate the development and deployment of state-of-the-art natural language processing tools for the full range of human languages. This presentation, which follows a previous paper on the same topic, reports on new developments including feedback on a deposit at Hugging Face Datasets.

MOTS-CLÉS : Corpus de référence, documentation computationnelle des langues, langues rares.

KEYWORDS : Benchmark datasets, Computational Language Documentation, endangered languages.

1 Introduction

Le déploiement d'outils de traitement automatique de la parole comporte des enjeux évidents pour la documentation des langues, à une époque où le déclin de la diversité linguistique s'accélère (parallèlement au déclin de la biodiversité). Inversement, les langues rares présentent à la recherche en informatique tout un éventail de défis, dont l'intérêt est de plus en plus clairement perçu ANASTASOPOULOS et al., 2020. Dans ce contexte, la mise à disposition de corpus de langues rares aisément accessibles, clairement versionnés et faciles d'utilisation paraît une nécessité tout à fait centrale. Cet exposé, qui fait suite à une précédente communication sur le même thème (GALLIOT et al., 2021), présente notre contribution à cette thématique. Suivant l'exemple de la publication du corpus mbochi (bantou) (GODARD et al., 2018), nous avons déposé dans Zenodo et dans Hugging Face Datasets deux corpus audio (avec transcriptions) de langues rares : le japhug et le na, langues minoritaires de Chine. Ces corpus ont été utilisés dans des travaux en reconnaissance automatique de la parole (ADAMS et al., 2018; ADAMS et al., 2021; GUILLAUME, WISNIEWSKI, MACAIRE, et al., 2022; MACAIRE, 2021) et dans des réflexions interdisciplinaires associant TAListes et linguistes (GUILLAUME, WISNIEWSKI, GALLIOT, et al., 2022; MICHAUD et al., 2018; MICHAUD et al., 2020). Les corpus sont disponibles en ligne dans une archive ouverte, la collection Pangloss¹ (MICHAUD et al., 2016), archive ouverte de langues rares hébergée par la plateforme Cocoon², de sorte que l'accès en est ouvert.

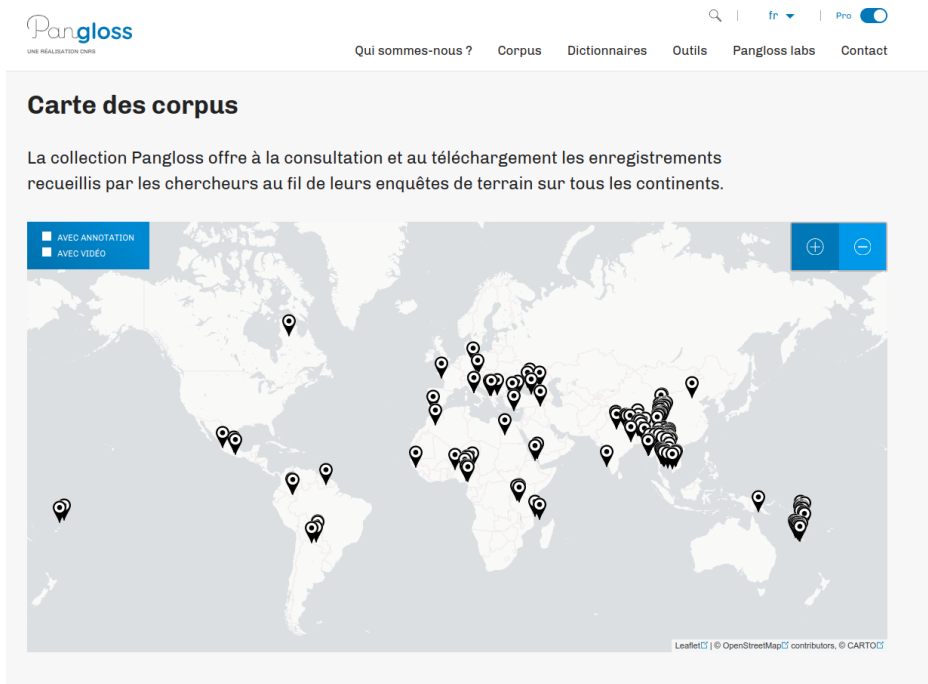


FIGURE 1 – Pangloss – carte des langues

Néanmoins, les transcriptions de ces corpus sont enrichies et revues au fil des années, et de nouveaux documents s'y ajoutent, de sorte que renvoyer à l'archive ouverte elle-même ne constitue pas une

1. <https://pangloss.cnrs.fr/>

2. <https://cocoon.huma-num.fr/>

Accueil / Corpus / na de Yongning

Na de Yongning (aussi appelé *narua* et *mosuo* 摩梭语)

La langue na (endonyme : /nɑ̃J-zwɛ˥/, 'na'+langue) est parlée à la frontière des provinces chinoises du Yunnan et du Sichuan, aux abords du lac Lugu (lo˥˥vJ-hiJnɑ̃˥mi˥).


Elle appartient au groupe naish de la famille sino-tibétaine, qui comprend également le naxi et le lazé.

En 2010, il lui a été accordé une entrée à part, sous le nom translittéré de Narua, dans l'inventaire des langues du Summer Institute of Linguistics (code: nru). Jusque-là, le na était considéré comme un dialecte de la langue naxi.

La grande majorité des documents offerts ici à la consultation et au téléchargement ont été recueillis dans la plaine de Yongning (nom en na : /ɦiɦiJ-djiJmi˥/), et proviennent du parler décrit dans une monographie parue en 2017 ([disponible en ligne](#)).

[Voir la suite](#)

Ressources

DOI	Type	Transcription(s)	Durée	Titre	Chercheur(s)	Locuteur(s)
			00:08:37	Le mariage de la sœur (version 1)	Michaud, Alexis	Lafami, Dashilame — lo˥˥q̄˥˥ tmi˥- [æ ɦ̄w̄ɪ- loJmɪ] — 拉它米 打史拉么

🔄 中文介绍



Chercheurs

Michaud, Alexis



Sanctuaire en contrehaut de la plaine de

FIGURE 2 – Pangloss – page d'un corpus : introduction, et liste des ressources

référence suffisamment précise pour parvenir à la reproductibilité des expériences menées sur ces données. Il s'agit, sinon de « données chaudes », du moins de « données tièdes », qui évoluent lentement au fil du temps.

Nous avons donc effectué un dépôt d'un état donné de ces deux corpus, ainsi rendu accessibles en quelques clics, sous une forme stabilisée, dans Zenodo (§2) puis parmi les jeux de données Hugging Face (§3). Mais l'avancée la plus importante à nos yeux est la création d'un script, OutilPangloss (§4), pour panacher à son aise parmi les corpus de la collection Pangloss tout en conservant une garantie de reproductibilité (grâce au système de versionnage dont bénéficient les documents déposés en archive).

2 Les dépôts Zenodo : liens d'accès et choix techniques

Lieu de dépôt Les deux corpus na et japhug ont été téléversés sur Zenodo, dont ils constituent respectivement les dépôts 5336698 (na) et 5521112 (japhug). Un corpus entier est identifié par un DOI (*digital object identifier*) : 10.5281/zenodo.5521112 pour le **japhug**, et 10.5281/zenodo.5336698 pour le **na**.

Le même type d'identifiant a été déployé pour la collection Pangloss, l'archive ouverte où sont déposés les corpus, mais avec une granularité tout à fait différente : un DOI pour chaque document (VASILE et al., 2020), ce qui est bien adapté pour les linguistes qui souhaitent faire référence aux données avec une granularité fine (un texte et, à l'intérieur d'un texte, un énoncé précis) mais ne donne pas prise sur

un corpus entier.

Fichiers audio Les fichiers audio ont été dégradés en 16 bit, 16 kHz, mono. La logique qui préside à la constitution de ces corpus versionnés pour expériences de TAL s'éloigne de celle de l'archivage pérenne dans la collection Pangloss. La taille des deux jeux de données déposés dans Zenodo est compatible (au jour d'aujourd'hui) avec des expériences menées sur un ordinateur portable : 1,8 Go pour le na, 9,2 Go pour le japhug.

Annotations Les annotations sont dans le format d'origine : du XML organisé selon une hiérarchie simple (un texte est composé de phrases, composées de mots, composés de morphèmes). Un exemple est fourni en Figure 3.

```
1 <?xml version="1.0" encoding="utf-8"?>
2 <IDCTYPE KWOLIST SYSTEM "https://ccocoon.huma-num.fr/schemas/Archive.dtd">
3 <KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
4 <HEAD>
5 <TITLE xml:lang="en">The tones of compound nouns: body parts of animals, document 12. Speaker F4, year 2008:
6 with electrolinguistic signals.</TITLE>
7 <SOUNDFILE href="Tone_BodyPartofAnimals_12_F4_2008_wiTHEGG.wav"/>
8 </HEAD>
9 <NOTE xml:lang="en" message="All the compound nouns are framed in the sentence 'this is...': proximal demonstrative /
10 is'na| / target item <sup>hoj</sup>we| pi|j| / . The demonstrative is realized [is'na] due to utterance-initial position. The
11 tone of the copula in context depends on what precedes."/>
12 <NOTE xml:lang="en" message="Determiner: boi?"/>
13 <NOTE xml:lang="en" message="Head: na?"/>
14 <NOTE xml:lang="en" message="Input tones (separated by a space): LM LM"/>
15 <NOTE xml:lang="en" message="Output tone (LH)"/>
16 <FRAMBL xml:lang="en">peau de porc (couenne de porc)/FRAMBL
17 <FRAMBL xml:lang="en">NRJ/FRAMBL
18 <FRAMBL xml:lang="en">pig's skin/FRAMBL
19 </KWOLIST id="c20-c122" xml:lang="fr">
20 </>
21 <?xml version="1.0" encoding="utf-8"?>
22 <IDCTYPE KWOLIST SYSTEM "https://ccocoon.huma-num.fr/schemas/Archive.dtd">
23 <KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
24 <HEAD>
25 <TITLE xml:lang="en">The tones of compound nouns: body parts of animals 12. Speaker F4, year 2008:
26 with electrolinguistic signals.</TITLE>
27 <SOUNDFILE href="Tone_BodyPartofAnimals_12_F4_2008_wiTHEGG.wav"/>
28 </HEAD>
29 <NOTE xml:lang="en" message="All the compound nouns are framed in the sentence 'this is...': proximal demonstrative /
30 is'na| / target item <sup>hoj</sup>we| pi|j| / . The demonstrative is realized [is'na] due to utterance-initial position. The
31 tone of the copula in context depends on what precedes."/>
32 <NOTE xml:lang="fr" message="Élicitation est arrangée par tête (détérminé). Ce choix existe la nomenclature tonale
33 d'un arrangement par déterminant, car le ton du déterminant a une plus grande influence sur celui du composé que le
34 ton de la tête. Le présent document, en revanche, est présenté par déterminant."/>
35 <NOTE xml:lang="en" message="Toutes les expressions sont précédées de
36 is'na| /, 'na|is [is'na], et suivies de la copule, /pi|j|, dont le ton en contexte dépend de ce qui précède."/>
37 <NOTE xml:lang="fr" message="L'élicitation est arrangée par tête (détérminé). Ce choix existe la nomenclature tonale
38 d'un arrangement par déterminant, car le ton du déterminant a une plus grande influence sur celui du composé que le
39 ton de la tête. Le présent document, en revanche, est présenté par déterminant."/>
40 <NOTE xml:lang="en" message="Determiner: boi?"/>
41 <NOTE xml:lang="en" message="Head: na?"/>
42 <NOTE xml:lang="en" message="Input tones (separated by a space): LM LM"/>
43 <NOTE xml:lang="en" message="Output tone (LH)"/>
44 <FRAMBL xml:lang="en">peau de porc (couenne de porc)/FRAMBL
45 <FRAMBL xml:lang="en">NRJ/FRAMBL
46 <FRAMBL xml:lang="en">pig's skin/FRAMBL
47 </KWOLIST id="c20-c122" xml:lang="fr">
48 </>
49 </IDCTYPE KWOLIST SYSTEM "https://ccocoon.huma-num.fr/schemas/Archive.dtd">
50 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
51 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
52 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
53 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
54 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
55 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
56 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
57 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
58 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
59 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
60 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
61 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
62 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
63 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
64 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
65 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
66 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
67 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
68 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
69 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
70 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
71 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
72 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
73 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
74 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
75 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
76 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
77 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
78 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
79 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
80 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
81 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
82 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
83 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
84 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
85 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
86 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
87 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
88 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
89 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
90 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
91 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
92 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
93 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
94 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
95 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
96 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
97 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
98 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
99 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
100 </KWOLIST id="c10-crd0-NRUL_F4_TOME12" xml:lang="fr">
```

(a) données

(b) métadonnées

FIGURE 3 – Fichiers XML des ressources d'annotations

Un prétraitement élémentaire a été effectué, afin de ne pas imposer aux utilisateurs de devoir prendre connaissance d'un certain nombre de conventions choisies par les déposants. En particulier, lors de la transcription de textes, il arrive que des retouches soient apportées, qui éloignent la transcription de ce qui a été dit sur l'enregistrement ; les passages ajoutés sont signalés par des crochets [] , et les passages que les consultants linguistiques souhaitent voir retranchés de la transcription « lissée » sont placés entre chevrons <> . Lors du prétraitement, les premiers ont été effacés, et les seconds allégés de leurs chevrons, afin qu'audio et transcription coïncident.

3 Les corpus au format Hugging Face

Les jeux de données Hugging Face (HF) constituent actuellement un lieu central pour les TAListes à la recherche de corpus aisément utilisables. Le formatage des corpus au format HF est donc apparu opportuniste dans la perspective qui consiste à porter des données de langues rares à l'attention de chercheurs et chercheurs en TAL.

Préparation de l'archive formatée Chaque fichier d'annotations est divisé en autant de fichiers qu'il comporte d'unités de premier niveau (phrases pour les textes, mots pour les listes de vocabulaire). Ces fichiers portent le nom de ces unités, et sont placés dans des dossiers du nom de la ressource

d'origine. Ces données, structurées de manière tabulaire, sont partitionnées en trois jeux (entraînement, validation et test), donc trois fichiers CSV.

	chemin_audio	nature	locuteur	forme	traduction.fr	traduction.zh	traduction.en	chemin_audio_segment
5966	cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1/Tone_BodyPartsOfA	WORDLIST	Latami, Dashiame	tsʰui qʰimi-ŋilpi pil	oreilles de renard	狐狸的耳朵	fox's ears	./langues/corpus/Yongning Na/corpus HFT Yongning Na ~ 16k-16/Yongning Na/cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1/Tone_BodyPartsOfA
5967	cocoon-71bb9ae8-6794-3509-aba1-648f583a1e5e_C1/NumPlusCL_13_Bun	WORDLIST	Latami, Dashiame	ŋwri-qal	5 grosses bottes (de paille...)	5包 (麦秆...)	5 large bundles (of straw...)	./langues/corpus/Yongning Na/corpus HFT Yongning Na ~ 16k-16/Yongning Na/cocoon-71bb9ae8-6794-3509-aba1-648f583a1e5e_C1/NumPlusCL
5968	cocoon-fef3d28-7b82-371e-be90-c02850f95e6_C1/HOUSEBUILDING_2	TEXT	Latami, Dashiame	tsʰui tsʰui soɪ-bæŋ-ŋil tsʰui-dao tsʰui-ŋil-dao tsʰui-qʰwri-doi-tsoŋ-pil pil-zoi ŋyʰimi-qoŋ-ŋui lae... 4ilb-ŋil-zoi lei-qʰwɛŋ ŋwri dzoŋ ŋil-qoŋ tsʰui-ŋil-ŋil-zoi ŋyʰ-qʰwɛŋ	Alors, ces trois sortes (la locutrice inclut les feuilles de nœve dans la liste, acceptant élargement la suggestion qui lui a été faite)... Comme on se dit: "Cette année, on va construire une maison", au neuvième mois, euh... on va déterrer les radis; ce qu'on a planté sur les terres de la famille, on le déterre!			./langues/corpus/Yongning Na/corpus HFT Yongning Na ~ 16k-16/Yongning Na/cocoon-fef3d28-7b82-371e-be90-c02850f95e6_C1/HOUSEBUILDING
5969	cocoon-4419730d-5bad-387b-9a15-47a5bb5418ae_C1/Tone_BodyPartsOfA	WORDLIST	Latami, Dashiame	tsʰui polloɪ-byɪ qʰitʰsæɪ pil	gorge de bélier	公绵羊的喉咙	ram's throat	./langues/corpus/Yongning Na/corpus HFT Yongning Na ~ 16k-16/Yongning Na/cocoon-4419730d-5bad-387b-9a15-47a5bb5418ae_C1/Tone_BodyPartsOfA
5970	cocoon-285cfefb-5ba1-3c90-b6c7-ce7e9a2197_C1/NumPlusCL_MH2_P4	WORDLIST	Latami, Dashiame	qʰitʰŋyɪ-kyj	99+classificateur des personnes/hommes	99个人	99+classifier for people/persons	./langues/corpus/Yongning Na/corpus HFT Yongning Na ~ 16k-16/Yongning Na/cocoon-285cfefb-5ba1-3c90-b6c7-ce7e9a2197_C1/NumPlusCL

(a) données pour le format Hugging Face

identifiant_annotation	chemin_annotation	identifiant_audio	chemin_audio	identifiant_forme	nature	locuteur	début	fin	forme	traduction.fr	traduction.zh	traduction.en	chemin_audio_segment
9399	cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.xm	Yongning Na ~ 16k-16/cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.xm	cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.b012-019252bb4f4d	Tone_BodyPartsOfAnim	WORDLIST	Latami, Dashiame	988.0	1003.0	tsʰui qʰimi-ŋilpi pil	oreilles de renard	狐狸的耳朵	fox's ears	./langues/corpus/Yongning Na/corpus HFT Yongning Na ~ 16k-16/Yongning Na/cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1/Tone_E
9400	cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.xm	Yongning Na ~ 16k-16/cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.xm	cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.b012-019252bb4f4d	Tone_BodyPartsOfAnim	WORDLIST	Latami, Dashiame	1078.0	1080.0	tsʰui qʰimi-ŋilpi pil	taille de renard	狐狸的腰	fox's waist	./langues/corpus/Yongning Na/corpus HFT Yongning Na ~ 16k-16/Yongning Na/cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1/Tone_E
9401	cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.xm	Yongning Na ~ 16k-16/cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.xm	cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.b012-019252bb4f4d	Tone_BodyPartsOfAnim	WORDLIST	Latami, Dashiame	1170.0	1172.0	tsʰui ŋyʰim-ŋyɪ pil	yeux de renard	狐狸的眼睛	fox's eyes	./langues/corpus/Yongning Na/corpus HFT Yongning Na ~ 16k-16/Yongning Na/cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1/Tone_E
9402	cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.xm	Yongning Na ~ 16k-16/cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.xm	cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.b012-019252bb4f4d	Tone_BodyPartsOfAnim	WORDLIST	Latami, Dashiame	1227.0	1229.0	tsʰui qʰimi-sæŋŋil pil	cou de renard	狐狸的脖子	fox's neck	./langues/corpus/Yongning Na/corpus HFT Yongning Na/cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1/Tone_E
9403	cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.xm	Yongning Na ~ 16k-16/cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.xm	cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.b012-019252bb4f4d	Tone_BodyPartsOfAnim	WORDLIST	Latami, Dashiame	208.37	210.5	tsʰui zwætzoɪ-yuɪ pil	peau de poulain	马驹子的皮	colt's skin	./langues/corpus/Yongning Na/corpus HFT Yongning Na ~ 16k-16/Yongning Na/cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1/Tone_E
9404	cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.xm	Yongning Na ~ 16k-16/cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.xm	cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1.b012-019252bb4f4d	Tone_BodyPartsOfAnim	WORDLIST	Latami, Dashiame	280.25	282.25	tsʰui zwætzoɪ-byɪ pil	intestin de poulain	马驹子的肠子	colt's intestine	./langues/corpus/Yongning Na/corpus HFT Yongning Na ~ 16k-16/Yongning Na/cocoon-db3cf0e1-30bb-3225-b012-019252bb4f4d_C1/Tone_E

(b) données au format Hugging Face (après partitionnement aléatoire)

FIGURE 4 – OutilsPangloss : fichiers CSV

Certaines métadonnées pertinentes, comme la nature de la ressource (liste de mots ou texte), apparaissent également dans le format tabulaire. Les données sont ensuite aléatoirement partitionnées (par défaut au niveau des phrases, mais il est possible de choisir le niveau des fichiers) en trois jeux (entraînement, validation et test), donc trois fichiers CSV. Enfin, toutes ces données (fichiers audios des phrases et 3 fichiers CSV) sont téléversées sur un serveur.

Préparation du jeu de données public La seconde partie consiste à préparer formellement les nouveaux jeux de données³. Pour ce faire, il faut préparer une description précise de chaque corpus. Une attention particulière est portée à l'encodage de la langue concernée⁴, point délicat dans le cas des langues pas ou peu standardisées qui constituent le cœur de métier de la collection Pangloss. Un

3. Disponibles ici : <https://huggingface.co/datasets/Lacito/pangloss>

4. <https://github.com/huggingface/datasets/issues/4881>

script⁵ Python automatise la création du jeu de données depuis les archives des corpus préalablement téléversées et accessibles. Ce script, de complexité variable selon les jeux de données, sert notamment à préciser où sont stockées les archives formatées, les types de données de chaque colonne et leurs noms correspondants. Il est ensuite possible de visualiser en ligne directement les données des corpus et d'écouter les précieux segments audio : voir Figure 5.

Il a été choisi de créer un jeu de données global du nom de Pangloss, dont les diverses langues (japhug⁶, na de Yongning⁷) constituent un sous-jeu.

The screenshot shows the Hugging Face interface for the 'Lacito pangloss' dataset. The 'Dataset Preview' section is active, showing a subset named 'yong1288' with a 'train' split. Below this, a table displays several rows of data. Each row includes a file path, an audio player, a sentence in the source language, a doctype (e.g., 'WORDLIST' or 'TEXT'), and translations in French, English, and Chinese.

path (string)	audio (audio)	sentence (string)	doctype (string)	translation:fr (string)	translation:en (string)	translation:zh (string)
"cocoon-062e5982-a63f-3674-b09c-0a375a21c0d1_C1/Topic_BodyParts0EAnIma1s_6_Ve..."		"īi, t̄ʰm̄ ʒweizoi-bɔi p̄i."	"WORDLIST"	"intestin de poulain"	"colt's intestine"	"马驹子的肠子"
"cocoon-adf9c39f-e558-36ac-963a-e58f876e812_c1/F00D_SHORTAGE25061.wav"		"t̄ʰīj̄, əd̄ōi t̄ʰvi-ji-ŋ̄i, ..."	"TEXT"	"Alors, le père, il s'est tenu assis là, regardant..."	".."	".."
"cocoon-9245e7f00-3e19-3eac-8E3a-2dc3a040a4d3_C1/gisiter3_5135.wav"		"tōm̄j̄ ɔ̄l-kv̄j- t̄ʰvi-ji-ŋ̄i, ..."	"TEXT"	"On frappe les poteaux: on donne un coup par ici, un..."	".."	"要打柱子: 这边打一下, 那边打一下."
"cocoon-5152256a-dd16-345a-9325-73920219692c_c1/F4_TONETUTORIAL_PART3_025.wa..."		"jōj-kīl"	"WORDLIST"	".."	"to (a/the) sheep"	".."
"cocoon-3a093d18-647a-33a3-802a-2ab0a0e7f638_c1/gisiter_5090.wav"		"n̄jēl-suk̄j̄ h̄ī le-t̄gal-d̄āl ..."	"TEXT"	"chez nous autres, quand quelqu'un meurt,"	"(when) (one of) our people dies, then"	"我们, 一个人去世了的时候, 那么..."
"cocoon-6bd9a90b-95ac-374c-d12c-598fab22ad3c_c1/MuPPlusCL_L3_Bund1e0H̄ay_1to..."		"ʒ̄ēt̄s̄t̄s̄ī sōj- ɔ̄j̄ "	"WORDLIST"	"43 grosses bottes (de paille...)"	"43 large bundles (of straw...)"	"43 抱 (麦秆...)"

FIGURE 5 – Hugging Face : page des corpus Pangloss

Il est dès lors possible en deux lignes⁸ de charger le jeu de données voulu et de commencer à les utiliser.

4 Un outil pour constituer de nouveaux jeux de données

L'outil élaboré à l'occasion de la préparation des corpus déposés dans Zenodo et Hugging Face, intitulé OutilsPangloss⁹, consiste en une boîte à outils divers (en langage Julia) servant notamment à créer des (sous-)corpus de langues rares de Pangloss. L'utilisateur-riche remplit un fichier YAML (dont différents exemples sont fournis), en indiquant notamment le nom de la langue (telle qu'elle figure sur Pangloss). Elle peut également fournir une liste d'expressions rationnelles de modifications si des traitements

5. <https://huggingface.co/datasets/Lacito/pangloss/blob/main/pangloss.py>

6. <https://huggingface.co/datasets/Lacito/pangloss/viewer/japh1234/train>

7. <https://huggingface.co/datasets/Lacito/pangloss/viewer/yong1288/train>

8. `import datasets puis`

`datasets.load_dataset("pangloss", "japh1234")` pour le japhug et

`datasets.load_dataset("pangloss", "yong1288")` pour le na de Yongning.

9. <https://gitlab.com/lacito/outilspangloss>

sur les annotations sont à faire (telles que suppressions ou réarrangements de blocs de textes). Il est possible de filtrer par locuteur pour les sous-corpus. Des traitements sur l'audio peuvent également être paramétrés, pour choisir le taux d'échantillonnage et la profondeur, et séparer les différentes pistes des fichiers multicanaux en fichiers mono (démultiplexage).

Après le moissonnage (en Sparql), les vérifications des données par les métadonnées (hachage, versions, etc.) et les téléchargements, un fichier récapitulatif général (données.yml) se trouvera dans le dossier cible, aux côtés de dossiers données et métadonnées : voir Figure 6. Les informations contenues dans ce fichier récapitulatif suffisent à reproduire exactement, à tout moment, une expérience menée avec le jeu de données qu'il décrit.

Nom	Taille
corpus HFT Yongning Na – 16k-16	3 éléments
données	433 éléments
données.yml	416,7 ko
métadonnées	750 éléments
ressources obsolètes	2 éléments
Yongning Na – 16k-16	348 éléments
Yongning Na (annotated).yml	133,3 ko
Yongning Na (annotated, converted).yml	220,6 ko
Yongning Na (annoté).yml	138,5 ko
Yongning Na (annoté, converti).yml	230,7 ko

FIGURE 6 – OutilsPangloss : dossier cible local d'un corpus

```

1  Langue: Yongning Na
2  > modifications:--
9  corpus:
10  - chemin: "../langues/corpus/Yongning Na"
11  nom complet: "Corpus de Na de Yongning complet"
12  commentaire: "Données complètes."
13  Langue: fr
14  sous-corpus:
15  - récapitulatifs:
16  - nom de fichier: "Na de Yongning (annoté)"
17  nom complet: "Corpus de Na de Yongning (complet)"
18  commentaire: "Données annotées brutes."
19  Langue: fr
20  - nom de fichier: "Yongning Na (annotated)"
21  nom complet: "Yongning Na corpus (complete)"
22  commentaire: "Raw annotated data."
23  Langue: en
24  sous-chemin: "Yongning Na - 16k-16"
25  traitements:
26  - audio:
27  - taux d'échantillonnage: 16000
28  - profondeur: 16
29  - spatialisations: "séparer"
30  récapitulatifs:
31  - nom de fichier: "Na de Yongning (annoté, converti)"
32  nom complet: "Corpus de Na de Yongning (complet)"
33  commentaire: "Données annotées brutes, conversion audio réalisée par Sox 14.4.2."
34  Langue: fr
35  - nom de fichier: "Yongning Na (annotated, converted)"
36  nom complet: "Yongning Na corpus (complete)"
37  commentaire: "Raw annotated data, audio conversion performed by Sox 14.4.2."
38  Langue: en
39  conversions:
40  - format: "HuggingFace-Transformers"
41  sous-chemin: "corpus HFT Yongning Na - 16k-16"
42  nom de fichier: "Na de Yongning (annoté, converti)"
43  nom du corpus: "yong1288"

```

(a) configuration d'un corpus

```

1  Langue: "Yongning Na"
2  nom: "Corpus de Yongning Na (complet)"
3  commentaire: "Données annotées brutes, conversion audio réalisée par Sox 14.4.2."
4  lien_documentation: "https://pangloss.cnrs.fr/corpus/Yongning420na"
5  code_langue: "fr"
6  version: "1.0"
7  > modifications:--
17  nombre_patrimoine_culturels: 116
18  nombre_fichiers: 232
19  nombre_audios: 116
20  nombre_annotations: 116
21  ressources:
22  - Identifiant: "CMO_cocoon-014b91c5-06a7-336c-8dc0-94b7fcca3ef"
23  type: "patrimoine_culturel"
24  lien_métadonnées: "http://cocoon.huma-num.fr/pub/CMO_cocoon-014b91c5-06a7-336c-8dc0-94b7fcca3ef"
25  locuteur: "Latani, DashiLame"
26  enfants:
27  - Identifiant: "Mw_Transcri_cocoon-014b91c5-06a7-336c-8dc0-94b7fcca3ef"
28  type: "audio original"
29  lien_métadonnées: "http://cocoon.huma-num.fr/pub/Mw_Transcri_cocoon-014b91c5-06a7-336c-8dc0-94b7fcca3ef"
30  lien_données: "http://purl.org/net/crodo/data/cocoon-014b91c5-06a7-336c-8dc0-94b7fcca3ef_version"
31  clef_hachage: "4452b242e967461f8910179c28a9"
32  version: 1
33  profondeur: 24
34  taux_echantillonnage: 64100
35  spatialisations: 1
36  enfants:
37  - type: "audio converti"
38  chemin_données: "Yongning Na - 16k-16/cocoon-014b91c5-06a7-336c-8dc0-94b7fcca3ef_c1.wav"
39  clef_hachage: "1dee6bd19926c9de9a8b8ce3a6688"
40  profondeur: 18
41  taux_echantillonnage: 16000
42  spatialisations: 1
43  Identifiant: "Mw_Transcri_cocoon-014b91c5-06a7-336c-8dc0-94b7fcca3ef"
44  type: "annotation originale"
45  lien_métadonnées: "http://cocoon.huma-num.fr/pub/Mw_Transcri_cocoon-014b91c5-06a7-336c-8dc0-94b7fcca3ef"
46  lien_données: "http://purl.org/net/crodo/data/cocoon-014b91c5-06a7-336c-8dc0-94b7fcca3ef"
47  version: 9
48  note_modification: "2022-02-16T16:41:56+01:00"
49  nature: "WORLDSIT"
50  enfants:
51  - type: "annotation dupliquée"
52  chemin_données: "Yongning Na - 16k-16/cocoon-014b91c5-06a7-336c-8dc0-94b7fcca3ef_c1.xml"
53  clef_hachage: "4e8d01211d2434848d8c9590815"
54  Identifiant: "CMO_cocoon-037c0b5f-33ba-34fd-836c-39c2664809f8"
55  type: "patrimoine_culturel"
56  lien_métadonnées: "http://cocoon.huma-num.fr/pub/CMO_cocoon-037c0b5f-33ba-34fd-836c-39c2664809f8"
57  locuteur: "Latani, DashiLame"

```

(b) corpus versionné (liste des ressources, métadonnées...)

FIGURE 7 – OutilsPangloss : fichiers YML

5 Perspectives : enjeux d’une description exhaustive des jeux de données sans copie *en dur*

La transition en cours vers la Science ouverte porte une exigence fondamentale de reproductibilité des expériences, et le domaine des sciences de la parole et du Traitement Automatique des Langues ne fait pas exception (GARELLEK et al., 2020). Notre espoir serait que les pratiques s’orientent à l’avenir vers une description des jeux de données via des métadonnées renvoyant vers un unique hébergement des données dans une archive pérenne. En effet, décrire de cette façon le jeu de données qu’on a utilisé ne prend que quelques kilooctets (Ko), tandis qu’un dépôt *en dur* de chaque bouquet de données multiplierait des dépôts (dans Zenodo ou ailleurs) dont chacun se compte en gigaoctets (Go).

Remerciements

Un grand merci aux collègues et amis consultants de langue japhug (en particulier Tshendzin) et na (en particulier M^{me} Latami Dashilame et son fils Latami Dashi). Le présent travail est une contribution au projet « La documentation computationnelle des langues à l’horizon 2025 » (ANR-19-CE38-0015-04) ainsi qu’au Labex « Fondements empiriques de la linguistique » (ANR-10-LABX-0083). Nous remercions l’Institut des langues rares (ILARA) de l’École pratique des hautes études, l’Université du Queensland et l’*Australian Research Council Centre of Excellence for the Dynamics of Language* pour le soutien financier apporté au développement d’outils logiciels pour la documentation linguistique.

Références

ADAMS, O., COHN, T., NEUBIG, G., CRUZ, H., BIRD, S., & MICHAUD, A. (2018). Evaluating phonemic transcription of low-resource tonal languages for language documentation. *Proceedings of LREC 2018 (Language Resources and Evaluation Conference)*, 3356–3365. <https://halshs.archives-ouvertes.fr/halshs-01709648>

ADAMS, O., GALLIOT, B., WISNIEWSKI, G., LAMBOURNE, N., FOLEY, B., SANDERS-DWYER, R., WILES, J., MICHAUD, A., GUILLAUME, S., BESACIER, L., COX, C., APLONOVA, K., JACQUES, G., & HILL, N. (2021). User-friendly automatic transcription of low-resource languages: Plugging ESPnet into Elpis. *Proceedings of ComputEL-4: Fourth Workshop on the Use of Computational Methods in the Study of Endangered Languages*. <https://halshs.archives-ouvertes.fr/halshs-03030529>

ANASTASOPOULOS, A., COX, C., NEUBIG, G., & CRUZ, H. (2020). Endangered languages meet Modern NLP. *Proceedings of the 28th International Conference on Computational Linguistics: Tutorial Abstracts*, 39–45. <https://doi.org/10.18653/v1/2020.coling-tutorials.7>

GALLIOT, B., WISNIEWSKI, G., GUILLAUME, S., MICHAUD, A., ROSSATO, S., NGUYÊN, M.-C., & FILY, M. (2021). Deux corpus audio transcrits de langues rares (japhug et na) normalisés en vue d’expériences en traitement du signal. *Journées scientifiques du Groupement de recherche “Linguistique informatique, formelle et de terrain” (GDR LIFT)*. <https://halshs.archives-ouvertes.fr/halshs-03475436>

- GARELLEK, M., GORDON, M., KIRBY, J., LEE, W.-S., MICHAUD, A., MOOSHAMMER, C., NIEBUHR, O., RECASENS, D., ROETTGER, T. B., SIMPSON, A., et al. (2020). Toward open data policies in phonetics: What we can gain and how we can avoid pitfalls. *Journal of Speech Science*, 9(1), 3–16.
- GODARD, P., ADDA, G., ADDA-DECKER, M., BENJUMEA, J., BESACIER, L., COOPER-LEAVITT, J., KOARATA, G. N., LAMEL, L., MAYNARD, H., & MUELLER, M. (2018). A very low resource language speech corpus for computational language documentation experiments. *Proceedings of LREC 2018 (Language Resources and Evaluation Conference)*, 3366–3370.
- GUILLAUME, S., WISNIEWSKI, G., GALLIOT, B., NGUYỄN, M.-C., FILY, M., JACQUES, G., & MICHAUD, A. (2022). Plugging a neural phoneme recognizer into a simple language model: a workflow for low-resource settings. *Interspeech 2022 - 23rd Annual Conference of the International Speech Communication Association*. <https://doi.org/10.5281/zenodo.5521111>
- GUILLAUME, S., WISNIEWSKI, G., MACAIRE, C., JACQUES, G., MICHAUD, A., GALLIOT, B., COAVOUX, M., ROSSATO, S., NGUYỄN, M.-C., & FILY, M. (2022). Fine-tuning pre-trained models for Automatic Speech Recognition: experiments on a fieldwork corpus of Japhug (Trans-Himalayan family). *ComputEL-5 5th Workshop on Computational Methods for Endangered Languages (ComputEL-5)*. <https://halshs.archives-ouvertes.fr/halshs-03647315>
- MACAIRE, C. (2021). *Recognizing lexical units in low-resource language contexts with supervised and unsupervised neural networks* (Research Report). LACITO (UMR 7107). <https://hal.archives-ouvertes.fr/hal-03429051>
- MICHAUD, A., ADAMS, O., COHN, T., NEUBIG, G., & GUILLAUME, S. (2018). Integrating automatic transcription into the language documentation workflow: Experiments with Na data and the Persephone toolkit. *Language Documentation and Conservation*, 12, 393–429. <http://hdl.handle.net/10125/24793>
- MICHAUD, A., ADAMS, O., COX, C., GUILLAUME, S., WISNIEWSKI, G., & GALLIOT, B. (2020). La transcription du linguiste au miroir de l'intelligence artificielle : réflexions à partir de la transcription phonémique automatique. *Bulletin de la Société de Linguistique de Paris*, 116(1). <https://halshs.archives-ouvertes.fr/halshs-02881731/>
- MICHAUD, A., GUILLAUME, S., JACQUES, G., MAC, D.-K., JACOBSON, M., PHAM, T.-H., & DEO, M. (2016). Contribuer au progrès solidaire des recherches et de la documentation : la Collection Pangloss et la Collection AuCo. *Journées d'Etude de la Parole 2016*, 1, 155-163. <https://halshs.archives-ouvertes.fr/halshs-01341631>
- VASILE, A., GUILLAUME, S., AOUINI, M., & MICHAUD, A. (2020). Le Digital Object Identifier, une impérieuse nécessité? L'exemple de l'attribution de DOI à la Collection Pangloss, archive ouverte de langues en danger. *I2D - Information, données & documents*, 2, 156-175. <https://halshs.archives-ouvertes.fr/halshs-02870206>