



HAL
open science

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

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

► To cite this version:

Benjamin Galliot, Guillaume Wisniewski, Séverine Guillaume, Guillaume Jacques, Alexis Michaud. Facilitating NLP specialists' access to language archive materials: an update. Journées Jointes des Groupements de Recherche Linguistique Informatique, Formelle et de Terrain (LIFT) et Traitement Automatique des Langues (TAL), Nov 2022, Marseille, France. pp.109-118. hal-03846839

HAL Id: hal-03846839

<https://hal.science/hal-03846839>

Submitted on 14 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

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

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

ABSTRACT

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.

RÉSUMÉ

Faciliter l'accès des praticiens du traitement automatique des langues à des jeux de données de langues rares : un deuxième point d'étape

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.

KEYWORDS: benchmark datasets, computational language documentation, endangered languages.

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

1 Introduction

The deployment of automatic speech processing tools clearly has important implications for language documentation, at a time when the decline in linguistic diversity is accelerating (Kik et al., 2021), in parallel with the decline in biodiversity. Conversely, less-documented languages present a whole range of challenges to computational research, the interest of which is increasingly clearly perceived (Anastasopoulos et al., 2020). In this context, providing easily accessible, clearly versioned and user-friendly corpora of less-studied languages appears as a central necessity. This presentation, which follows a previous paper on the same theme (Galliot et al., 2021), presents our contribution to this endeavour. Following the example of the publication of the Mbochi (Bantu) corpus (Godard et al., 2018), we have deposited in Zenodo and in Hugging Face Datasets audio corpora (with transcriptions) of Japhug and Na, two minority languages of China.

These corpora have been used in automatic speech recognition work (Adams et al., 2018; Adams et al., 2021; Guillaume, Wisniewski, Macaire, et al., 2022; Macaire, 2021) and in interdisciplinary reflections associating NLP specialists and linguists (Guillaume, Wisniewski, Galliot, et al., 2022; Michaud et al., 2018; Michaud et al., 2020). The corpora are available online in an open archive, the Pangloss Collection¹ (Michaud et al., 2016), an open archive of (mostly) endangered languages, itself hosted by the Cocoon data repository², ensuring availability in open access.

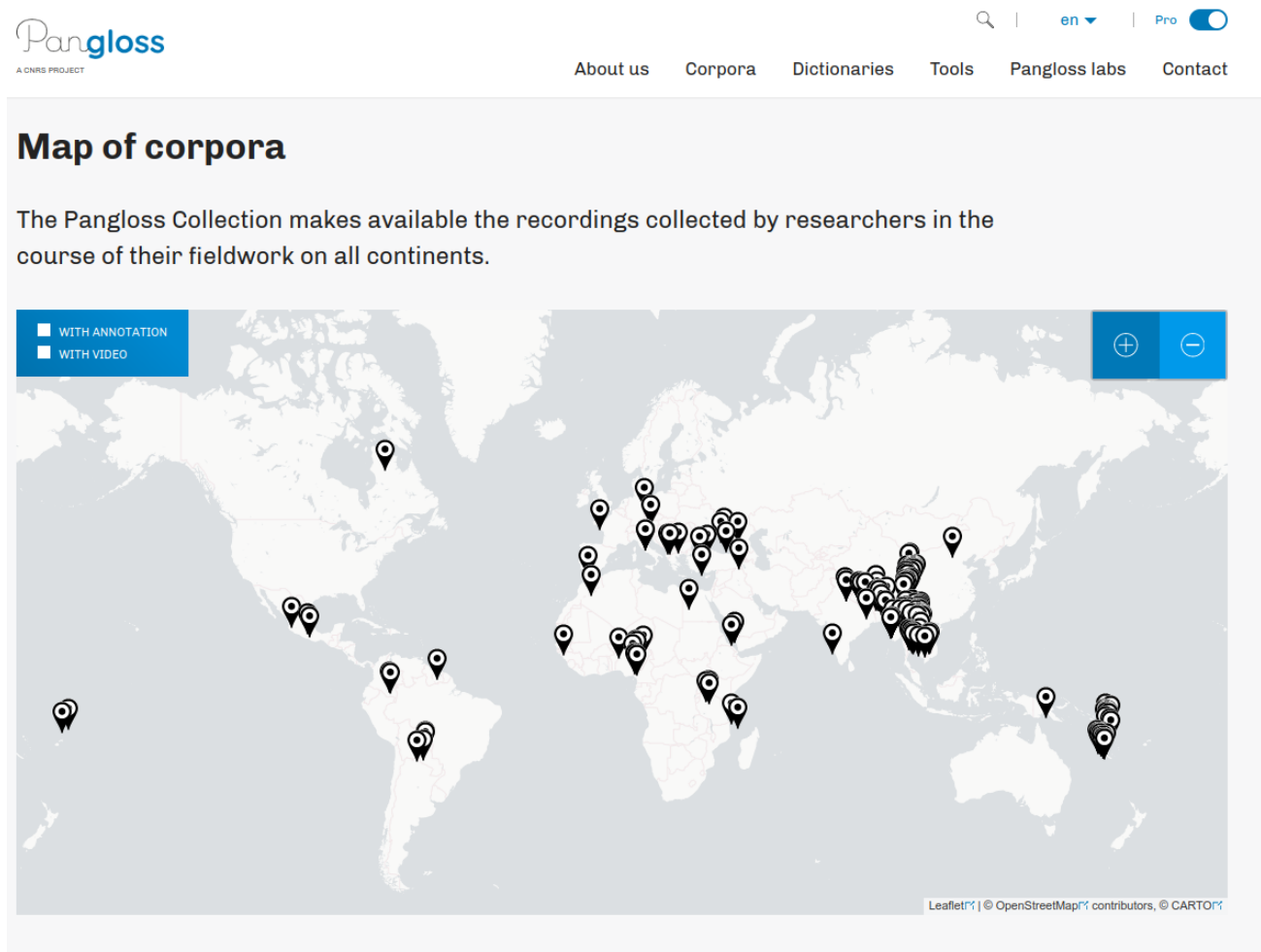


Figure 1: Pangloss – map of languages

¹<https://pangloss.cnrs.fr/>

²<https://cocoon.huma-num.fr/>

Home / Corpus / Yongning Na

Yongning Na (also known as *Narua* or *Mosuo* 摩梭语)

/nɑŋ-ɰwɿ/, i.e. the *Na language*, is spoken in an area straddling the boundary between the Chinese provinces of Yunnan and Sichuan, in the vicinity of lake Lugu (loɰɰyJ-hiJnɑɰmiɰ). The total number of speakers was estimated at about 40,000 on the basis of survey data from the late 1950s (He Jiren & Jiang Zhuyi 1985:107); the same figure is taken up by Yang Zhenhong (2009). The number of proficient speakers at present (2020) is much lower: language replacement is under way.

The Na are famous among anthropologists for their family structures, as a "society without marriage". Until the 21st century, much less attention was directed to their language, which has fascinating properties, such as its elaborate morphotology.

[See more](#)

Resources

| DOI | Type | Transcription(s) | Duration | Title | Researcher(s) | Speaker(s) |
|-----|------|------------------|----------|--|-----------------|---|
| | | | 00:08:37 | Sister: The sister's wedding (version 1) | Michaud, Alexis | Latami, Dashilame — lɑ-tʰɑ-imi- [æ-ɰwɿ- lɑJmɰ] — 拉它米 打史拉么 |
| | | | | | | Latami, Dashilame — |

[中文介绍](#)



Researchers

Michaud, Alexis



A sanctuary overlooking the plain of Yongning

Figure 2: Pangloss – page of a corpus: introduction and list of available multimedia resources

The transcriptions of these corpora are enriched and revised over the years, however, and new documents are added, so that referring to the open archive itself is not a sufficiently precise reference to achieve actual reproducibility of the experiments conducted on these data. Even though the primary data (the audio and video files) do not change, the transcriptions constitute “lukewarm data”, which get improved and modified slowly over time.

We have therefore made a deposit of a given state of these two corpora, making them accessible in a few clicks, in a stabilized form, in Zenodo (§3) and among the Hugging Face Datasets (§4). But the most important advance, in our view, is the creation of a script, *OutilsPangloss* (§2), to select among the corpora of the Pangloss Collection while maintaining a guarantee of full reproducibility (thanks to the versioning system of documents deposited in the archive).

2 A tool to build new datasets

The tool developed during the preparation of the corpora deposited in Zenodo and Hugging Face, entitled *OutilsPangloss*³, allows for putting together an open-ended range of new datasets. It consists of a toolbox in Julia language, which, among other functionalities, creates (sub)corpora of less-documented language data from the Pangloss Collection. The user fills in a YAML file (of which various examples are provided, see Figure 3a), specifying the names of the desired languages (as they

³<https://gitlab.com/lacito/outilspangloss>

appears in the Pangloss Collection). It can also provide a list of regular expressions for modifications in case processing of annotations is to be carried out (such as deletions or rearrangements of text blocks). It is possible to filter by speaker. Processing on the audio can also be set, to choose the sampling rate and bit-depth, and to separate the different tracks of multi-channel files into mono files (demultiplexing).

After harvesting (in Sparql), metadata checks (hashing, versions, etc.), data checks (mistake detection, etc.), and downloads, a general summary file (see Figure 3b) is generated in the target folder, alongside data and metadata folders (see Figure 4). The information contained in this summary file is sufficient to reproduce exactly, at any time, an experiment conducted with the dataset it describes.

```

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      spatialisation: "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"
44

```

(a) corpus configuration

```

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/Yongning320Na"
5  code_langue: "nru"
6  version: "1.0"
7  > modifications:--
17  nombre_patrimoines_culturels: 116
18  nombre_fichiers: 232
19  nombre_audios: 116
20  nombre_annotations: 116
21  ressources:
22  - identifiant: "CHO_cocoon-014b91c5-06a7-336c-8dc0-94b7fcec3ef"
23    type: "patrimoine culturel"
24    lien_métadonnées: "http://cocoon.huma-num.fr/pub/CHO_cocoon-014b91c5-06a7-336c-8dc0-94b7fcec3ef"
25    locuteur: "Latami, DashiLame"
26  enfants:
27  - identifiant: "WR_Record1_cocoon-014b91c5-06a7-336c-8dc0-94b7fcec3ef"
28    type: "audio original"
29    lien_métadonnées: "http://cocoon.huma-num.fr/pub/WR_Record1_cocoon-014b91c5-06a7-336c-8dc0-94b7fcec3ef"
30    lien_données: "http://purl.org/net/crdo/data/cocoon-014b91c5-06a7-336c-8dc0-94b7fcec3ef.version1"
31    clef_hachage: "4452db2de5627a61f89610179c28829"
32    version: 1
33    profondeur: 24
34    taux_échantillonnage: 44100
35    spatialisation: 1
36    enfants:
37    - type: "audio converti"
38      chemin_données: "Yongning Na - 16k-16/cocoon-014b91c5-06a7-336c-8dc0-94b7fcec3ef_C1.wav"
39      clef_hachage: "10ee9bd19926ec9de9ea86bce53a6688"
40      profondeur: 16
41      taux_échantillonnage: 16000
42      spatialisation: 1
43  - identifiant: "WR_Transcr1_cocoon-014b91c5-06a7-336c-8dc0-94b7fcec3ef"
44    type: "annotation originale"
45    lien_métadonnées: "http://cocoon.huma-num.fr/pub/WR_Transcr1_cocoon-014b91c5-06a7-336c-8dc0-94b7fcec3ef"
46    lien_données: "http://cocoon.huma-num.fr/data/michaud/masters/crdo-NRU-F4_TONEPEOPLES.xml"
47    version: 5
48    date_modification: "2022-02-16T16:41:58+01:00"
49    nature: "WORDLIST"
50    enfants:
51    - type: "annotation dupliquée"
52      chemin_données: "Yongning Na - 16k-16/cocoon-014b91c5-06a7-336c-8dc0-94b7fcec3ef_C1.xml"
53      clef_hachage: "b49da71311fd2643ab4da0ac9580cf6"
54  - identifiant: "CHO_cocoon-037c0b5f-33ba-34fb-83c6-39c26e4d09f8"
55    type: "patrimoine culturel"
56    lien_métadonnées: "http://cocoon.huma-num.fr/pub/CHO_cocoon-037c0b5f-33ba-34fb-83c6-39c26e4d09f8"
57    locuteur: "Latami, DashiLame"

```

(b) Versioned corpus (list of resources, metadata...)

Figure 3: OutilsPangloss : YML files

| 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 4: OutilsPangloss: local target folder of a corpus

| Nom | Taille |
|--|----------|
| cocoon-0a8404b4-49ee-3c80-964f-99320dad78f5_C1.wav | 2,8 Mo |
| cocoon-0a8404b4-49ee-3c80-964f-99320dad78f5_C1.xml | 14,2 ko |
| cocoon-0a8404b4-49ee-3c80-964f-99320dad78f5_C2.wav | 2,8 Mo |
| cocoon-0a8404b4-49ee-3c80-964f-99320dad78f5_C2.xml | 14,2 ko |
| cocoon-0ce38fef-4386-3446-b5b2-fc420d4cf148_C1.wav | 6,5 Mo |
| cocoon-0ce38fef-4386-3446-b5b2-fc420d4cf148_C1.xml | 43,4 ko |
| cocoon-0f6d7ca2-2fbd-3601-b3a5-b54232b43def_C1.wav | 34,5 Mo |
| cocoon-0f6d7ca2-2fbd-3601-b3a5-b54232b43def_C1.xml | 132,5 ko |
| cocoon-0f4576b5-e917-35a2-92f6-fed919660f5e_C1.wav | 4,9 Mo |
| cocoon-0f4576b5-e917-35a2-92f6-fed919660f5e_C1.xml | 44,0 ko |
| cocoon-1d91ab4c-fea9-3fd7-976f-e6753a255d1c_C1.wav | 24,0 Mo |
| cocoon-1d91ab4c-fea9-3fd7-976f-e6753a255d1c_C1.xml | 320,9 ko |
| cocoon-1fee21d5-125f-313d-aab7-121aa9c2cac5_C1.wav | 8,9 Mo |
| cocoon-1fee21d5-125f-313d-aab7-121aa9c2cac5_C1.xml | 172,5 ko |
| cocoon-1fee21d5-125f-313d-aab7-121aa9c2cac5_C2.wav | 8,9 Mo |
| cocoon-1fee21d5-125f-313d-aab7-121aa9c2cac5_C2.xml | 172,5 ko |
| cocoon-2bf8107a-6054-31b9-a355-38cb1f16706b_C1.wav | 17,3 Mo |
| cocoon-2bf8107a-6054-31b9-a355-38cb1f16706b_C1.xml | 84,3 ko |
| cocoon-2f5bb477-5267-38ea-bf77-f4e6dc258e9d_C1.wav | 27,7 Mo |
| cocoon-2f5bb477-5267-38ea-bf77-f4e6dc258e9d_C1.xml | 129,4 ko |
| cocoon-2f282ddc-cd08-32e0-b6b3-82f77cf0bf30_C1.wav | 14,5 Mo |
| cocoon-2f282ddc-cd08-32e0-b6b3-82f77cf0bf30_C1.xml | 53,8 ko |
| cocoon-2f282ddc-cd08-32e0-b6b3-82f77cf0bf30_C2.wav | 14,5 Mo |
| cocoon-2f282ddc-cd08-32e0-b6b3-82f77cf0bf30_C2.xml | 53,8 ko |
| cocoon-3b9c6235-e8a5-3927-8e35-d29664ed0a14_C1.wav | 37,9 Mo |

(a) processed resources

| Nom | Taille |
|--|--------------|
| cocoon-0a8404b4-49ee-3c80-964f-99320dad78f5_C1 | 33 éléments |
| cocoon-0ce38fef-4386-3446-b5b2-fc420d4cf148_C1 | 80 éléments |
| cocoon-0f6d7ca2-2fbd-3601-b3a5-b54232b43def_C1 | 200 éléments |
| cocoon-0f4576b5-e917-35a2-92f6-fed919660f5e_C1 | 81 éléments |
| cocoon-1d91ab4c-fea9-3fd7-976f-e6753a255d1c_C1 | 157 éléments |
| cocoon-1fee21d5-125f-313d-aab7-121aa9c2cac5_C1 | 73 éléments |
| cocoon-2bf8107a-6054-31b9-a355-38cb1f16706b_C1 | 141 éléments |
| cocoon-2f5bb477-5267-38ea-bf77-f4e6dc258e9d_C1 | 223 éléments |
| cocoon-2f282ddc-cd08-32e0-b6b3-82f77cf0bf30_C1 | 117 éléments |
| cocoon-3b9c6235-e8a5-3927-8e35-d29664ed0a14_C1 | 118 éléments |
| cocoon-3dd9341b-647a-33d3-b0fd-2b0ab0e7f638_C1 | 181 éléments |
| cocoon-3e084f49-332e-3744-88fc-cf5098405829_C1 | 20 éléments |
| cocoon-4a89f909-3102-3033-8afa-0b3bf27fe908_C1 | 115 éléments |
| cocoon-5a006333-133d-3952-94e0-cf774e0ca4b_C1 | 32 éléments |
| cocoon-6b82b30f-27b4-317b-82b9-42b2e79587b5_C1 | 110 éléments |
| cocoon-6bc8b80b-22f6-3986-92e3-4c6ed75ef0b5_C1 | 189 éléments |
| cocoon-6bd9a90b-95ac-374c-b12c-590fab22ad3c_C1 | 64 éléments |
| cocoon-6c87371d-f9ad-3496-83a4-5d5bf02b9388_C1 | 32 éléments |
| cocoon-7af967c5-d09c-332d-9bcd-9144c40e65b2_C1 | 79 éléments |
| cocoon-7d0ee91a-94ff-37f5-bccc-29c31fb7f4d_C1 | 120 éléments |
| cocoon-7d1f65c6-d83a-3e01-85b3-2e1d4cb901f9_C1 | 78 éléments |
| cocoon-7f9f8056-d111-3119-bafe-526f4160f0a4_C1 | 108 éléments |
| cocoon-8c729359-ac07-46bd-b293-59ac07b6db5_C1 | 33 éléments |
| cocoon-8e777be5-ab06-3fe9-b4f6-ac4123603ab7_C1 | 40 éléments |
| cocoon-8fd61350-7519-3b5d-a8b0-ba886d29b97_C1 | 176 éléments |

(b) processed resources for the Hugging Face format

Figure 5: OutilsPangloss : données

3 The Zenodo deposits: access links and explanations about technical choices

Links to the repository The Na and Japhug corpora were uploaded to Zenodo, where they constitute deposits 5336698 (Na) and 5521112 (Japhug), respectively. An entire corpus is identified by a DOI (*digital object identifier*): 10.5281/zenodo.5521112 for the [Japhug corpus](#), and 10.5281/zenodo.5336698 for the [Na corpus](#).

The same type of identifier has been deployed for the Pangloss Collection, the open archive where the corpora are deposited, but with a completely different granularity: a DOI for each document (Vasile et al., 2020), a choice which is suitable for linguists who wish to reference data at a highly specific level (a text and, within a text, a specific utterance) but does not give a handle on an entire corpus.

Audio files Audio files were downgraded to 16-bit, 16 kHz, mono (see Figure 5a). The logic behind this choice of parameters is that NLP experiments have requirements that differ from those of long-term archiving in the Pangloss Collection. The size of the two datasets deposited in Zenodo is compatible (as of today) with experiments carried out on a laptop computer: 1.8 GB for Na, 9.2 GB for Japhug.

Annotations The annotations are in the original format: XML structured according to a simple hierarchy (a text is composed of sentences, themselves composed of words, themselves composed of morphemes). An example is provided in Figure 6.

Some basic preprocessing has been carried out, so that users do not have to learn about the various conventions chosen by the depositors, which vary across corpora. In particular, when transcribing texts, it is not infrequent for language consultants and language workers (typically linguists) to make decisions that create an edit distance between the transcription and what is actually said in the recording. A convention used in the Na corpus is that added passages are placed within square brackets. Another convention is that passages that the language consultants wish to see removed from the “smoothed” transcription are placed between angle brackets. At preprocessing, the passages within square brackets were deleted, and the angle brackets were removed, thereby ensuring the closest match between audio and transcriptions.

| | |
|---|--|
| <pre> 1 <?xml version="1.0" encoding="utf-8"?> 2 <DOCTYPE WORDLIST SYSTEM "https://cocoon.huma-num.fr/schemas/Archive.did"> 3 <WORDLIST id="crdo-NRU_F4_TONE12" xml:lang="nru"> 4 <HEADER> 5 <TITLE xml:lang="en">The tones of compound nouns: body parts of animals, document 12. Speaker F4, year 2008; 6 with electroglottographic signal.</TITLE> 7 <SOUNDFILE href="Tone_BodyPartsOfAnimals.12_F4_2008_withEGG.wav"/> 8 </HEADER> 9 <NOTE xml:lang="en" message="All the compound nouns are framed in the sentence 'This is...: proximal demonstrative / 10 ts^w1 +target item +copula, /pij/. The demonstrative is realized [ts^w1] due to utterance-initial position. The 11 tone of the copula in context depends on what precedes."/><NOTE xml:lang="en" message="The elicitation was arranged 12 by head rather than by determiner: 'pig's skin' then 'tiger's skin', etc. This limits the repetitiveness of the 13 successive items, because the tone of the head has less influence on that of the compound than the tone of the 14 determiner: 'pig's skin', 'pig's intestines', 'pig fat' are more similar tonally than 'pig's skin', 'tiger's 15 skin', 'sheep's skin'... The present document, on the other hand, is arranged by determiner ('pig's skin', 'pig's 16 intestines', 'pig fat'...; then 'tiger's skin', 'tiger's intestines' and so on), as this seemed the more useful order 17 of presentation to study the tone system."/><NOTE xml:lang="fr" message="Toutes Les expressions sont précédées de / 18 ts^w1 /, réalisé [ts^w1], et suivies de la copule, /pij/, dont le ton en contexte dépend de ce qui précède."/ 19 ><NOTE xml:lang="fr" message="L'élitication est arrangée par tête (déterminé). Ce choix évite la monotonie tonale 20 d'un arrangement par déterminant, car le ton du déterminant a une plus grande influence sur celui du composé que le 21 ton de la tête. Le présent document, en revanche, est présenté par déterminant."/><w 22 id="Tone_BodyPartsOfAnimals.12_F4_2008_withEGG_001"> 23 <FORM ts^w1 boj-yw1 pij/><FORM> 24 <NOTE xml:lang="en" message="Determiner: boj^1"/> 25 <NOTE xml:lang="en" message="Head: wj^1"/> 26 <NOTE xml:lang="en" message="Input tones (separated by a space): LM LH"/> 27 <NOTE xml:lang="en" message="Output tone: LM"/> 28 <TRANSL xml:lang="fr">peau de porc (couenne de porc)</TRANSL> 29 <TRANSL xml:lang="zh">豬的皮</TRANSL> 30 <TRANSL xml:lang="en">pig's skin</TRANSL> 31 <AUDIO start="729.122" end="730.713"/> 32 </w> 33 <w id="Tone_BodyPartsOfAnimals.12_F4_2008_withEGG_002"> 34 <FORM ts^w1 boj-yw1 pij/><FORM> 35 <NOTE xml:lang="en" message="Determiner: boj^1"/> 36 <NOTE xml:lang="en" message="Head: by^1"/> 37 <NOTE xml:lang="en" message="Input tones (separated by a space): LM M"/> 38 <NOTE xml:lang="en" message="Output tone: LM"/> 39 <TRANSL xml:lang="fr">intestin de porc</TRANSL> 40 <TRANSL xml:lang="zh">豬的腸子</TRANSL> </pre> | <pre> 57 xmlns:eac-cpf="http://archivi.ibc.regione.emilia-romagna.it/ontology/eac-cpf/" 58 xmlns:foaf="http://xmlns.com/foaf/0.1/" 59 xmlns:lgdo="http://linkedgeodata.org/ontology/capital" > 60 <rdf:Description rdf:about="http://cocoon.huma-num.fr/pub/WR_Transcri1_cocoon-db3cf0e1-30bb-3225-b012-019252bb4fd"> 61 <dc:terms:accessRights>Freely accessible</dc:terms:accessRights> 62 <dc:terms:issued>rdf:datatype="http://pur1.org/dc/terms/W3CDTF">2013-09-21T21:18:16+02:00</dc:terms:issued> 63 <dc:identifier>oai:crdo.vjf.cnrs.fr:cocoon-5970d5ad-f863-33f5-adf9-71842837be9b</dc:identifier> 64 <dc:type>rdf:resource="http://www.language-archives.org/vocabulury/type/lexicon"/> 65 <dc:identifier>doi:10.34847/cocoon.5970d5ad-f863-33f5-adf9-71842837be9b</dc:identifier> 66 <dc:terms:replaces>rdf:resource="http://cocoon.huma-num.fr/data/michaud/masters/versions/crdo-NRU_F4_TONE12.v3.xml"/> 67 <dc:terms:license>rdf:resource="http://creativecommons.org/licenses/by-nc-sa/2.5"/> 68 <dc:identifier>hdl:10670/1.aralrux</dc:identifier> 69 <dc:type>rdf:resource="http://pur1.org/dc/dcmitype/Text"/> 70 <dc:terms:replaces>rdf:resource="http://cocoon.huma-num.fr/data/michaud/masters/versions/crdo-NRU_F4_TONE12.v2.xml"/> 71 <dc:subject>rdf:resource="http://lexvo.org/id/iso639-3/nru"/> 72 <foaf:primaryTopic>rdf:resource="http://cocoon.huma-num.fr/data/michaud/masters/crdo-NRU_F4_TONE12.xml"/> 73 <dc:ac:depositor>rdf:resource="http://viaf.org/viaf/120871374"/> 74 <dc:subject>Yongning Na</dc:subject> 75 <dc:language>rdf:resource="http://lexvo.org/id/iso639-3/nru"/> 76 <dc:identifier>doi:10.24397/PANGLOSS-000463</dc:identifier> 77 <dc:ac:speaker>rdf:resource="http://cocoon.huma-num.fr/pub/person_200#foaf:Person"/> 78 <dc:terms:created>rdf:datatype="http://pur1.org/dc/terms/W3CDTF">2008</dc:terms:created> 79 <dc:ac:researcher>rdf:resource="http://viaf.org/viaf/120871374"/> 80 <dc:terms:replaces>rdf:resource="http://cocoon.huma-num.fr/data/michaud/masters/versions/crdo-NRU_F4_TONE12.v1.xml"/> 81 <dc:identifier>doi:10.24397/PANGLOSS-000463</dc:identifier> 82 <dc:terms:modified>2022-02-16T16:39:24+01:00</dc:terms:modified> 83 <dc:terms:available>rdf:datatype="http://pur1.org/dc/terms/W3CDTF">2013-09-21</dc:terms:available> 84 <dc:contributor>rdf:resource="http://cocoon.huma-num.fr/pub/person_200#foaf:Person"/> 85 <dc:terms:conformsTo>rdf:resource="http://cho.cocoon-49aefa90-8c1f-3ba8-a099-0ebefc6a2aa"/> 86 <dc:publisher>rdf:resource="http://cocoon.huma-num.fr/pub/org_50#org:Organization"/> 87 <dc:contributor>rdf:resource="http://viaf.org/viaf/120871374"/> 88 <dc:identifier>ark:/87895/1.17-369643</dc:identifier> 89 <dc:ac:interviewer>rdf:resource="http://viaf.org/viaf/120871374"/> 90 <rdf:type>rdf:resource="http://www.europeana.eu/schemas/edm/WebResource"/> 91 <dc:format>rdf:datatype="http://pur1.org/dc/terms/IMT">text/xml</dc:format> 92 <dc:title>xml:lang="en">Tone: compound nouns. Body parts of animals 12, F4, 2008, with EGG</dc:title> 93 <dc:rights>Copyright (c) Michaud, Alexis</dc:rights> 94 <schema-org:version>4</schema-org:version> 95 </rdf:Description> 96 </rdf:RDF> </pre> |
|---|--|

(a) data

(b) metadata

Figure 6: XML files of annotation resources

4 The Pangloss corpus in Hugging Face Datasets

Hugging Face (HF) Datasets are currently a key hub for NLP researchers looking for easily usable corpora. The formatting of Pangloss corpora in HF format therefore seemed useful from the perspective of bringing corpora from less-studied languages to the attention of NLP researchers.

Preparing the data for Hugging Face Datasets format Each annotation file is divided into as many files as there are first-level units (sentences for texts, words for vocabulary lists). These files are named after these units, and are placed in folders named after the original resource (see Figure 5b). These data, structured in a tabular manner, contain some relevant metadata, such as the nature of the resource (word list or text) or the speaker. The data are then randomly partitioned. By default, the partitioning is carried out at the first level below the entire resource: hence sentences for texts, and words in the case of word lists. But it is possible to choose the level of the entire file instead. Three sets are thereby created: training, validation and test, thus three CSV files (see figure 7). Finally, all these data (audio files of the sentences and the three CSV files) are uploaded on a server.

Preparing the publicly available dataset A second part of the work consists in formally preparing the new datasets⁴. To achieve this, a precise description of each corpus is required. Particular attention is paid to the encoding of the language concerned⁵, a delicate point in the case of languages that are not standardized (or only slightly standardized), which constitute the core business of the Pangloss Collection. A Python script⁶ automates the creation of the dataset from the previously uploaded and accessible corpus archives. This script, of varying complexity depending on the dataset, is used to specify where the formatted archives are stored, the data types of each column, and their corresponding names. It is then possible to view the corpus data directly online and to listen to the audio segments (see Figure 8).

The choice made consists in creating one global dataset named Pangloss, of which the various

⁴Available here: <https://huggingface.co/datasets/Lacito/pangloss>

⁵<https://github.com/huggingface/datasets/issues/4881>

⁶<https://huggingface.co/datasets/Lacito/pangloss/blob/main/pangloss.py>

The screenshot shows the Hugging Face Datasets interface for the 'pangloss' dataset. At the top, there's a search bar and navigation links for Models, Datasets, Spaces, Docs, Solutions, and Pricing. Below that, the dataset name 'pangloss' is displayed with a 'like' button and a count of 2. There are several filters: Tasks (Automatic Speech Recognition, Fine-Grained Tasks: speech-recognition), Languages (jya, nru), Multilinguality (multilingual, translation), Size Categories (10K<n<100K), and Language Creators (expert-generated). Annotations Creators are listed as expert-generated, and Source Datasets are original. Licenses are cc-by-nc-sa-4.0. Below the filters, there are tabs for Dataset card, Files and versions, Community, and Settings.

Dataset Preview API

Subset: yong1288 Split: train

| path (string) | audio (audio) | sentence (string) | doctype (string) | translation:fr (string) | translation:en (string) | translation:zh (string) |
|--|---------------|--|------------------|--|---|-------------------------|
| "coco0n-d62e5852-a63f-3674-b09c-6b175a21cb81_C1/Tone_BodyParts0fAnimals_6_Ve..." | | "i, [shw Zmwzoi-bv pil" | "WORDLIST" | "intestin de poulain" | "colt's intestine" | "马驹子的肠子" |
| "coco0n-adf9c39f-e558-36ac-963a-e5b8f76e812_C1/F00D_SHORTAGE25061.wav" | | "t'hij, oiddi t'hv-i-v-i-rw-i, ..." | "TEXT" | "Alors, le père, il s'est tenu assis là, regardant..." | " " | " " |
| "coco0n-9245e7d0-3e19-3eac-8f3e-2dc3a04bbf43_C1/Sister3_S135.wav" | | "tojmij dā1-kvj-tswj ō -mvj, ..." | "TEXT" | "On frappe les poteaux; on donne un coup par ici, on..." | " " | "要打柱子: 这边打一下, 那边打一下, " |
| "coco0n-5152256a-dd16-345a-9325-73920210692c_C1/F4_TONETUTORIAL_PART3_025.wa..." | | "jo1-kil" | "WORDLIST" | " " | "to (a/the) sheep" | " " |
| "coco0n-3dd9341b-647a-33d3-b0fd-2b0ab0e7f638_C1/Sister_S090.wav" | | "njam-smjkwj hi1 lei-ɣwi-dzo1 ..." | "TEXT" | "chez nous autres, quand quelqu'un meurt," | "(when) (one of) our people dies, then" | "我们, 一个人去世的时候, 那么, " |
| "coco0n-6bd9a90b-95ac-374c-b12c-590fab22ad3c_C1/NumPlusCL_L3_Bund1e0fHay_1to..." | | "zvitshij soj-qdij" | "WORDLIST" | "43 grosses bottes (de paille...)" | "43 large bundles (of straw...)" | "43 抱 (麦秆...) " |

Figure 8: Page of the Pangloss corpus on Hugging Face Datasets

5 Perspectives: what is at stake in providing a full description of datasets without creating “hard copies”

The ongoing transition to Open Science carries a fundamental requirement for reproducibility of experiments, and the field of speech sciences and Automatic Language Processing is no exception (Garellek et al., 2020). Our hope is that practices will gradually evolve towards a description of datasets via metadata pointing to one master file hosted in an archive that guarantees both long-term conservation and 24/7 online availability. Describing the datasets in this way only takes a few kilobytes (Kb), whereas hosting each dataset as a “hard copy” would result in multiple and highly redundant deposits (in Zenodo or elsewhere) each of which amounts to gigabytes (GB).

Acknowledgments

Many thanks to the language consultants and friends for Japhug (in particular Tshendzin) and Na (in particular Mrs. Latami Dashilame and her son Latami Dashi). The present work is a contribution to the project “Computational language documentation by 2025” (ANR-19-CE38-0015-04) and to the Labex “Empirical foundations of linguistics” (ANR-10-LABX-0083). We thank the Institute for Linguistic Heritage and Diversity (ILARA) at *École pratique des hautes études*, the University of Queensland and the Australian Research Council Centre of Excellence for the Dynamics of Language for financial support to software development for language documentation.

References

- 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., KOUARATA, 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., NGUYÊN, 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., NGUYÊN, 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>
- KIK, A., ADAMEC, M., AIKHENVALD, A. Y., BAJZEKOVA, J., BARO, N., BOWERN, C., COLWELL, R. K., DROZD, P., DUDA, P., IBALIM, S., JORGE, L. R., MOGINA, J., RULI, B., SAM, K., SARVASY, H., SAULEI, S., WEIBLEN, G. D., ZRZAVY, J., & NOVOTNY, V. (2021). Language and ethnobiological skills decline precipitously in papua new guinea, the world's most linguistically diverse nation. *Proceedings of the National Academy of Sciences*, 118(22). <https://doi.org/10.1073/pnas.2100096118>

- 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., AOUNI, 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>