

Deux corpus audio transcrits de langues rares (japhug et na) normalisés en vue d'expériences en traitement du signal

Benjamin Galliot² Guillaume Wisniewski³ Séverine Guillaume²
Laurent Besacier¹ Guillaume Jacques⁵ Alexis Michaud² Solange Rossato¹
Minh-Châu Nguyễn^{1,2} Maxime Fily²

(1) Laboratoire d'Informatique de Grenoble (LIG), Unité Mixte de Recherche 5217 CNRS - Université Grenoble Alpes - Grenoble INP - Institut national de recherche en informatique et en automatique (INRIA)
(2) Langues et Civilisations à Tradition Orale (LACTO), Unité Mixte de Recherche 7107 CNRS - Sorbonne Nouvelle - Institut National des Langues et Civilisations Orientales (INALCO)
(3) Laboratoire de Linguistique Formelle (LLF), Unité Mixte de Recherche 7110 CNRS - Université de Paris
(4) Naver Labs Europe, Grenoble
(5) Centre de Recherches Linguistiques sur l'Asie Orientale (CRLAO), Unité Mixte de Recherche 8563 CNRS - École des Hautes Études en Sciences Sociales - Institut National des Langues et Civilisations Orientales
b.gillyon@gmail.com, Guillaume.Wisniewski@univ-paris-diderot.fr, severine.guillaume@cnrs.fr, laurent.besacier@univ-grenoble-alpes.fr, rgyalrongskad@gmail.com, alexis.michaud@cnrs.fr, Solange.Rossato@univ-grenoble-alpes.fr, minhchau.ntm@gmail.com, maxime.fily@gmail.com

RÉSUMÉ

Deux corpus audio transcrits de langues « rares » (langues minoritaires de Chine : japhug et na) sont proposés comme corpus de référence pour des expériences en traitement automatique des langues. Les données, collectées et transcrites au fil d'enquêtes de terrain en immersion, s'élevaient à un total de 1907 minutes d'audio transcrit en japhug et de 209 minutes en na. Nous décrivons les traitements effectués pour les mettre à disposition sous une forme aisément accessible et utilisable, et présentons un outil qui permet d'assembler divers jeux de données de la collection Pangloss (archive ouverte de langues rares) en assurant la reproductibilité des expériences menées sur ces données.

ABSTRACT

Two Very-Low-Resource Language Speech Corpora for Experiments in NLP : Japhug and Na

Two audio corpora of minority languages of China (Japhug and Na), with transcriptions, are proposed as reference data sets for experiments in Natural Language Processing. The data, collected and transcribed in the course of immersion fieldwork, amount to a total of 1,907 minutes in Japhug and 209 minutes in Na. By making them available in an easily accessible and usable form, we hope to facilitate the development and deployment of state-of-the-art NLP tools for the full range of human languages. We present a tool for assembling datasets from the Pangloss Collection (an open archive) in a way that ensures full reproducibility of experiments conducted on these data.

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

KEYWORDS: Benchmark datasets, Computational Language Documentation, low-resource languages, endangered languages.

Difficulté identifiée : accès des TAListes aux corpus de langues rares. Archive ouverte Cocoon (Collection de Corpus Oraux Numériques) : portails internet pour consultation plutôt que téléchargement en masse.

The screenshot shows the Pangloss website interface. At the top, there's a navigation bar with 'Qui sommes-nous?', 'Corpus', 'Dictionnaires', and 'Contact'. Below that, a search bar and a 'Recherche' button. The main content area displays the title 'Souvenir d'enfance : jeûner pendant la moisson' and a description: 'Le locuteur raconte un souvenir de moisson, pendant le ramadan.' There are options for 'Options d'affichage' and 'Lecture en continu'. A map shows the location of the recording. On the right, there's a sidebar with 'Téléchargements' (File media original (wav), File media compressed (mp3), Annotation 1 (xml)) and 'Métadonnées' (Resource (xml), Annotation 1 (xml)).

Besoin de reproductibilité des expériences.

Exemple d'utilisation : explorations en Reconnaissance Automatique de la Parole.

Voir exposé oral de Macaire et al.

The screenshot shows a GitHub repository page for 'macairececile / internship_lacito_2021'. It displays the repository name, a 'Public' badge, and a list of contributors. The main content is a code snippet from a Python file, showing imports for 'argparse', 'difflib', 'itertools', 'os', 'path', 'librosa', 'torchaudio', 'datasets', and 'pandas'.



Zenodo : téléchargement en quelques clics.

Formats : selon pratiques courantes en traitement du signal (audio dégradé en 16 kHz, 16-bit) et du texte (transcriptions et annotations en XML ; simplifications mineures pour confort d'utilisation

The screenshot shows a Zenodo record page for 'Japhug for Natural Language Processing: a single-speaker audio corpus with transcriptions'. It includes the title, authors (Jacques, Guillaume; Galliot, Benjamin; Guillaume, Séverine), and a description of the dataset. The page shows 30 views and 5 downloads. It is indexed in OpenAIRE. The publication date is September 22, 2021. The DOI is 10.5281/zenodo.5521112. The license is Creative Commons Attribution-NonCommercial-ShareAlike 2.0 France. The page also shows a list of files for download, including audio files and XML annotations.

Outil générique pour création de jeux de données intégralement documentés : Outils-Pangloss

Ce programme sobrement intitulé **OutilsPangloss** consiste en une boîte à outils divers servant notamment à créer des (sous-)corpus de langues rares de **Pangloss**.

Installation

Ce programme est codé en Julia et téléchargera si besoin est un moteur XSLT 2/3 (Saxon), ainsi Java devrait aussi être installé sur l'ordinateur hôte. Du fait de sa nature, il est nécessaire d'avoir une connexion internet.

Utilisation

Création de corpus

Tout ce dont a besoin un utilisateur est de remplir un fichier YAML (dont différents exemples se trouvent dans le dossier `exemples`) dont les éléments importants sont les suivants :

- le nom de la langue du corpus qui se trouve sur Pangloss (« Japhug », « Yongning Na », par exemple) ;
- la liste des graphèmes (notamment complexes) ;
- la liste des expressions rationnelles de modifications si des traitements sur les annotations sont à faire (suppressions ou réarrangements de blocs de textes...) ;
- les informations sur les corpus et sous-corpus, notamment :
 - le filtre de locuteur pour le sous-corpus ;
 - la langue du fichier récapitulatif associé (français : fr ou anglais : en) ;
 - les traitements à réaliser sur l'audio, notamment :
 - le taux d'échantillonnage (typiquement 16000 Hz) ;