

### Les données orales du CRDO

Michel Jacobson

Direction des Archives de France



#### Plan

- 1. Les types d'objet manipulés par le CRDO
- 2. Les fonctions de pré-ingestion
- 3. Les fonctions d'accès



### 1. Les types d'objets manipulés

- Pour que les ressources soient éligibles il faut qu'il s'agisse:
  - D'enregistrements de la parole
  - Accompagnés éventuellement d'annotations
  - Dont le but est d'étudier la langue ou le langage
  - Constitués dans le cadre d'une démarche scientifique
  - Par ou pour la communauté SHS



### Le périmètre

- Un enregistrement de parole:
  - Types de production orale: Récits, conversation, élicitations, textes lus, logatones...
  - Type de locuteurs: Vous et moi, des enfants avant 2 ans, des personnes atteintes de trouble de la parole (pas actuellement dans l'entrepôt)
  - Différentes langues. Il existe environ 6000 langues (ref. cataloge ethnologue) dont une infime partie possède une tradition d'écriture.
  - Type de situation: interactions au travail, interviews, expériences en chambre sourde...



### Le périmètre

- Un enregistrement de parole:
  - Ca ne veut pas dire enregistrement audio
    - Langues des signes
    - Langage enfantin
    - Interactions dans un groupe
  - En général les enregistrements sont audio ou vidéo mais on rencontre aussi des enregistrements de mesures physiologiques (electroglotogramme, endoscopie, IRMf...)
  - Ces enregistrements peuvent être accompagnés d'annotations (transcriptions, traductions, indications scénographiques, analyses syntaxiques, prosodiques, etc.)



### Description des ressources

- Les enregistrements
  - Audio
    - format = way, codage = pcm, fréquence: >= 44.1 Khz, taille de l'échantillon >= 16 bits
  - Vidéo
    - En attente des prescriptions du CINES
- Les annotations
  - format = xml, encodage = utf-8, différents schémas et dtd
  - format text seul, encodage ascii ou utf-8
  - format jpeg (pour les scans de manuscripts)
- Des métadonnées
  - Format xml, encodage utf-8, schémas d'OLAC



### Exemple d'une ressource type

Ethnologue > Web version > Country index > Asia > China > Jiarong

#### **Jiarong**

#### A language of China

Récit de 7 minutes, enregistré en 2005, transcrits en 2007 par Guillaume Jacques (CNRS/CRLAO). Langue Japhug. Chine (Sichuan, Maerkang, Bar-khams)

ISO 639-3: <u>jya</u>

Population 83,000 (1999 Sun Hongkai). 25,000 monolinguals. Ethnic population: 151,197 including 139,000 in Situ

Jiarong, 12,197 in Chabao and Sidaba (1993 Lin).

Region North central Sichuan. Situ is in the traditional territory of four chieftaincies: Zhuokeji, Suomo, Songgang,

Dangba, Chabao is in the northeastern corner of Maerkang county, at Longerjia, Dazang, and Shaerzong townships in Chabao District. Sidaba is in Caodeng, Kangshan, and Ribu townships in Sidaba District of

Maerkang County. Some outlying Sidaba communities are to the Rongan townships, at the southwesten corner of the Aba Counter River between Wuyi and Shill townships in Rangtang County, sonfluence of the Seda and Duke rivers in Seda County.

Alternate names

Jyarung, Gyarong, Gyarung, Rgyarong, Chiarong, Jarong

Dialects Chabao (Dazang, Northeastern Jiarong), Sidaba (Caodeng, No

Subdialects of Situ are: Maerkang, Lixian, Jinchuan, Xiaojin; of Western and Northern are fairly similar and differ greatly from E<sub>lesa</sub>

Eastern and Northern Jiarong, 60% between Western and Nort

Classification Sino-Tibetan, Tibeto-Burman, Tangut-Qiang, rGyarong

Language use Vigorous. All domains. All ages. Positive language attitude, but use of the language. 56,000 use Chinese, 950 Tibetan, 50 Qiariauhati

Radio programs, Dictionary,

Language Radio programs. Dio development

Comments Part of the Tibetan nationality. SOV; phonologically and lexically similar to Pumi and Qiang; complex consonant clusters; limited

riverine. Agriculturalists: apples, pears; lumbermen. Traditional

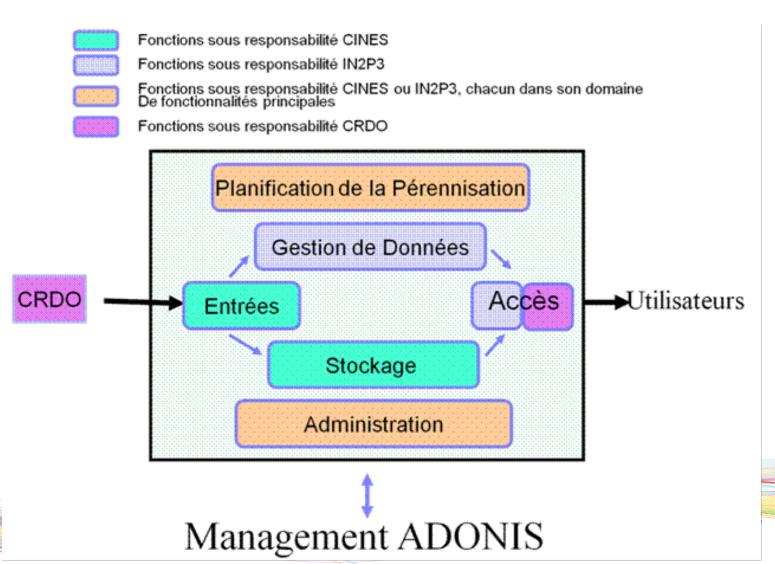
Démo





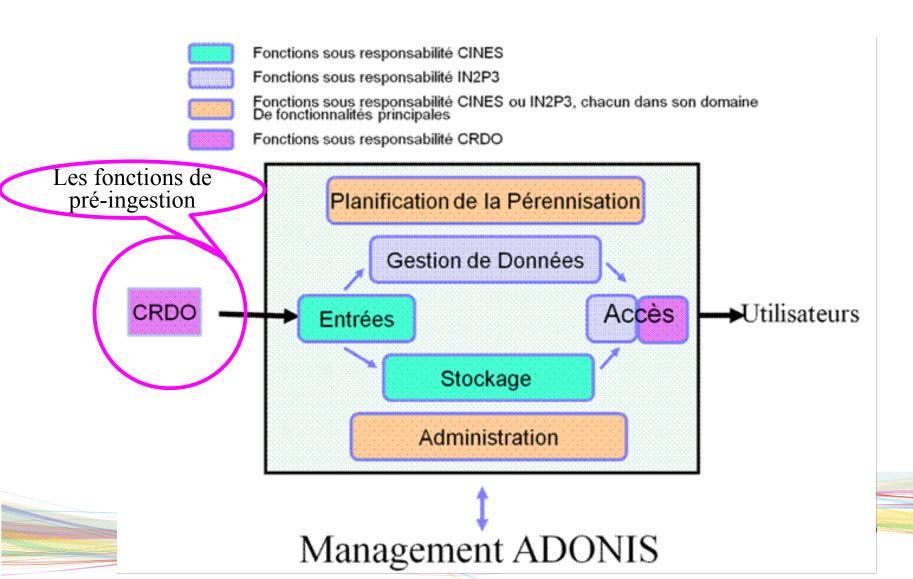


### La nouvelle organisation





#### 2. Les fonctions d'entrée





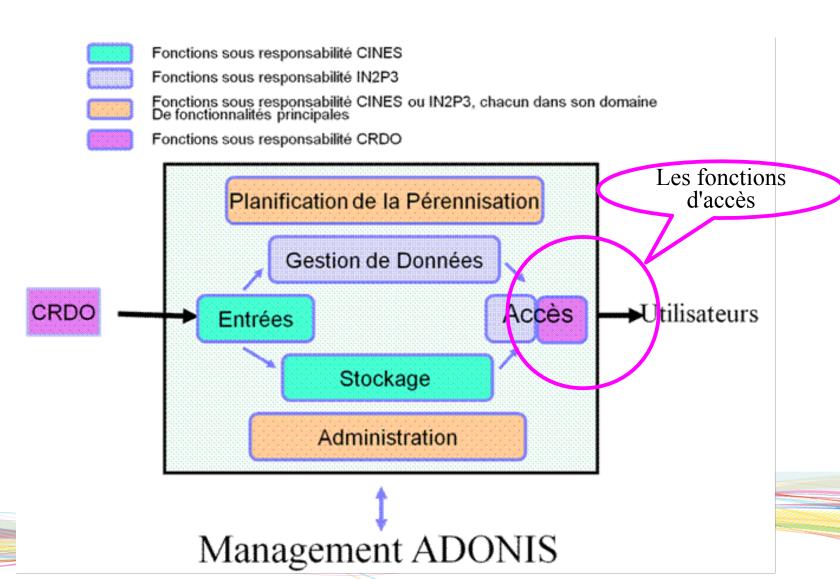
### Pré-ingestion

- Interface avec les producteurs pour
  - Aide à l'édition des données (outils d'édition, formations)
  - Aide à la normalisation des données
  - Aide à la construction des métadonnées
- Prévalidation des données
  - un outil a été crée pour les fichiers audio: vérification du format, du codage, de la qualité utilisée...
  - Xerces pour la validation des fichier xml
  - Validation de l'exactitude et de la complétude des métadonnées par des documentalistes
- Constitution et versement du SIP au CINES lorsque tout est ok





#### 3. Les fonctions d'accès





### 3. Les fonctions d'accès

- Des moteurs de recherche
- Des outils de consultation
- ✓ La Diffusion

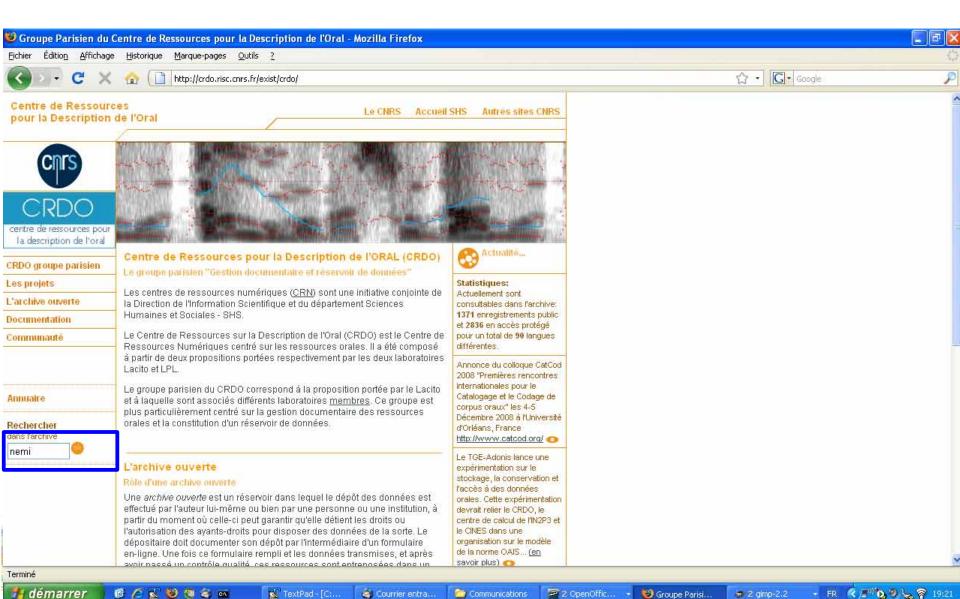


#### 3.1. Des moteurs de recherche

- Moteurs de recherche accéssibles sur le site web du CRDO
  - Par simple mot-clef
  - Par catégories Dublin-core
  - Par une carte géographique (googlemaps)
  - Par un axe temporel (simile timeline)
- La recherche peut être effectuée par les moteurs de recherche des fournisseurs de services qui ont moissonné l'entrepôt du CRDO
  - En général exploite les catégories Dublin-Code (OAIster) ou OLAC (LinguistList)

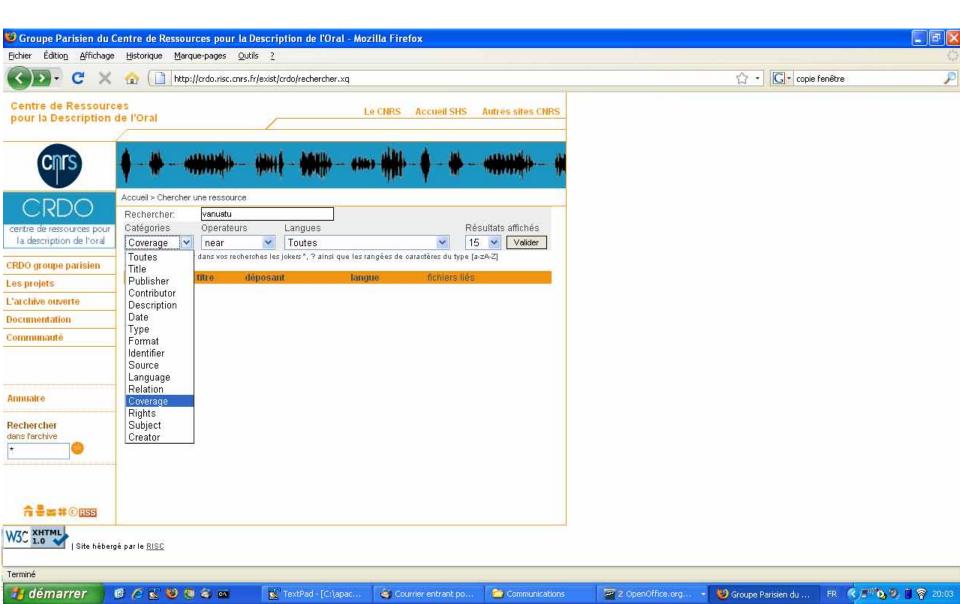


### Recherche par mot-clef



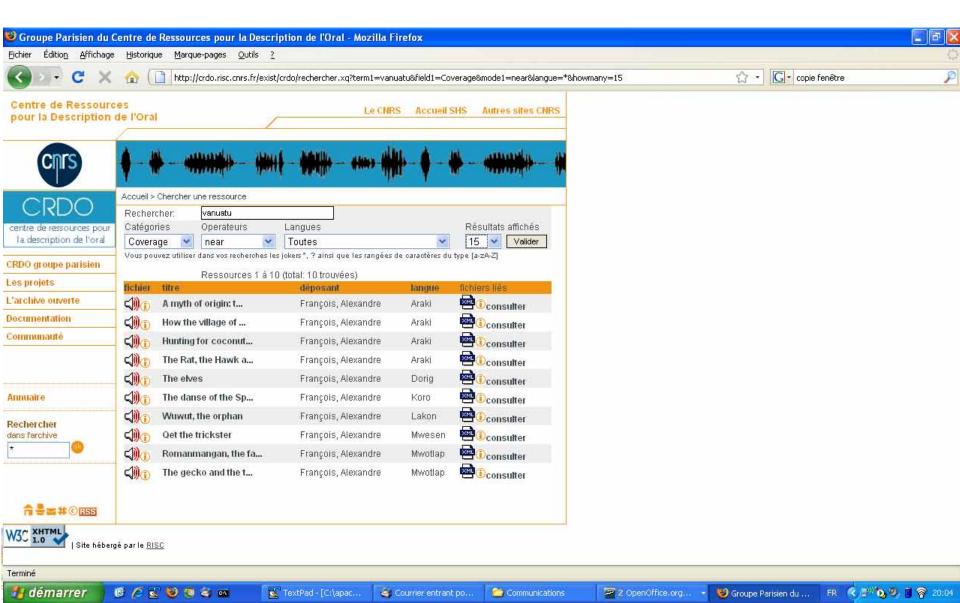


### Recherche par categorie D-C





### Résultats d'une recherche



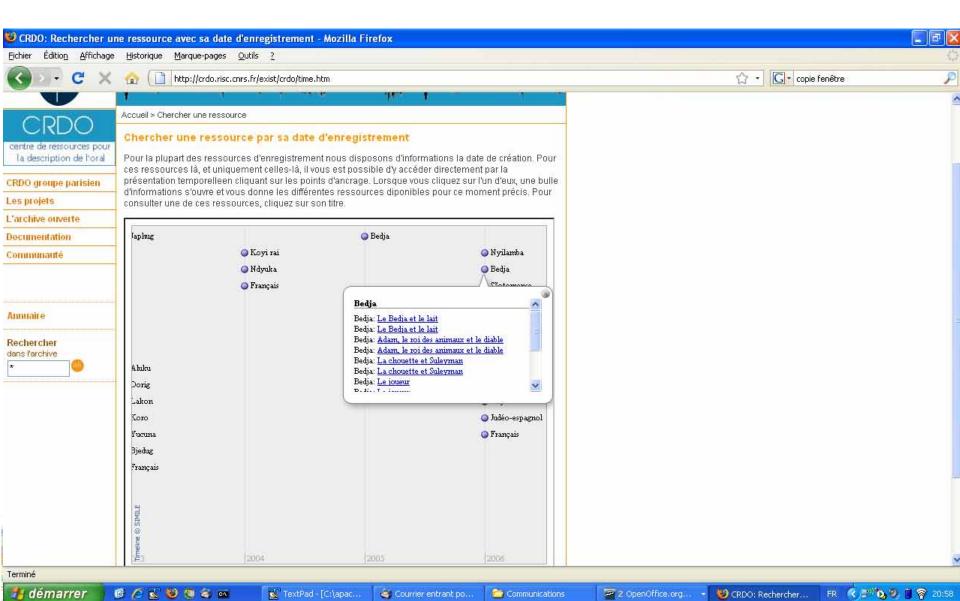


### Recherche par point d'enquête





### Recherche par date de création



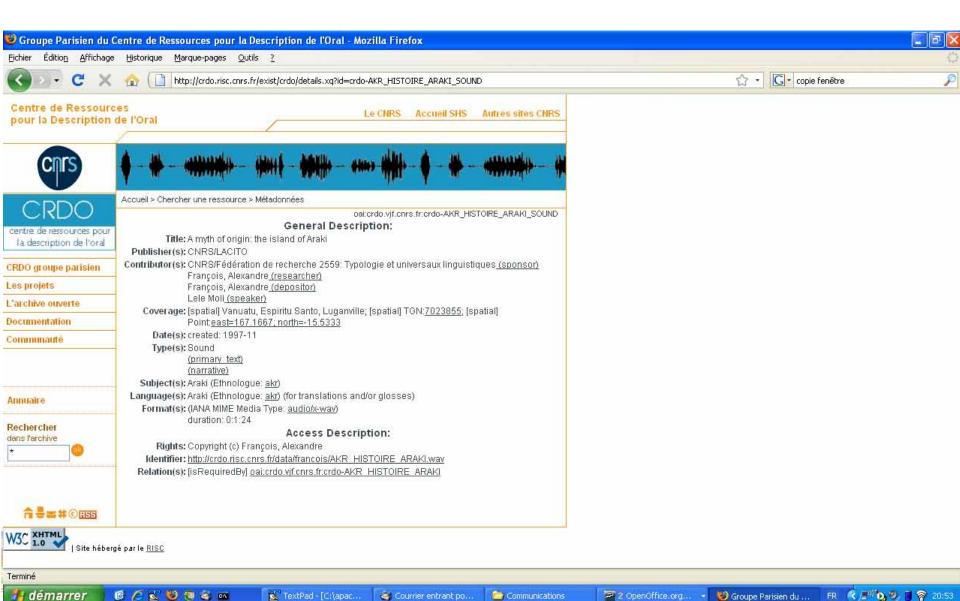


#### 3.2. Des outils de consultation

- Feuilles de styles xslt pour présenter les transcriptions et les métadonnées
- Outil de synchronisation (plugin + javascript)
- Outil de découpage d'extrait sonore (servlet java)



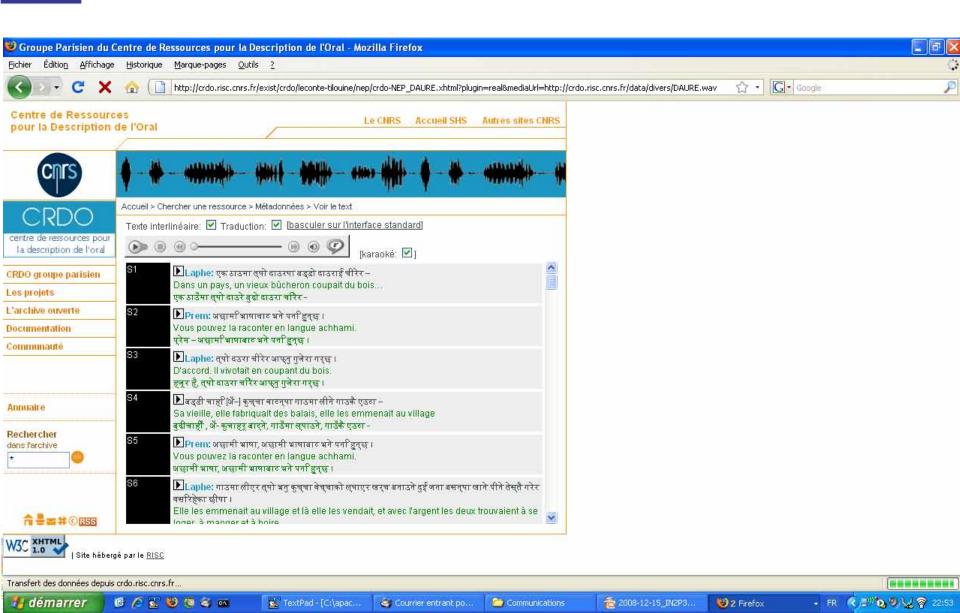
### Consultation des métadonnées





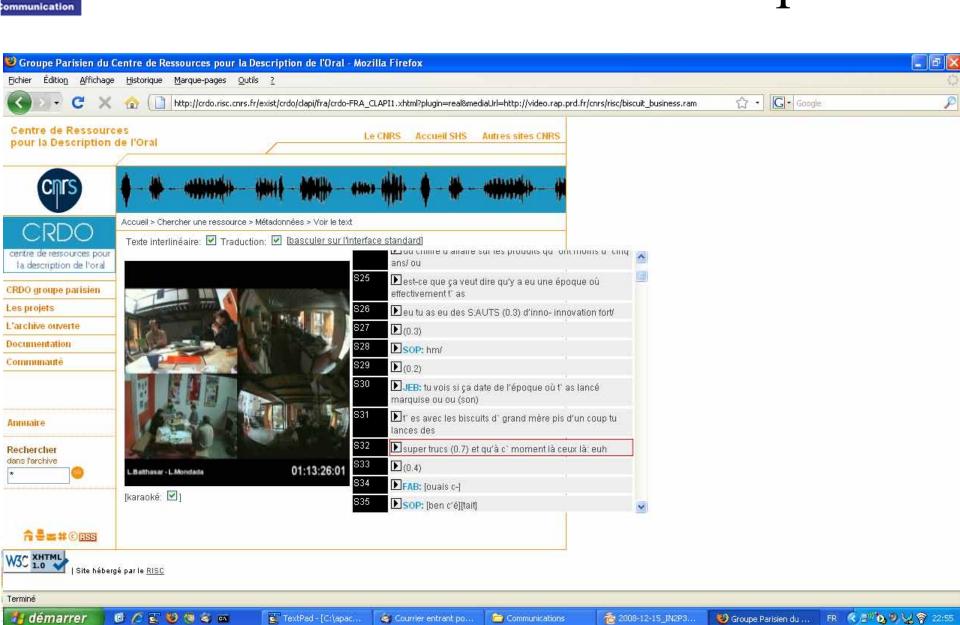
ommunication

### Consultation audio + transcription





### Consultation vidéo + transcription





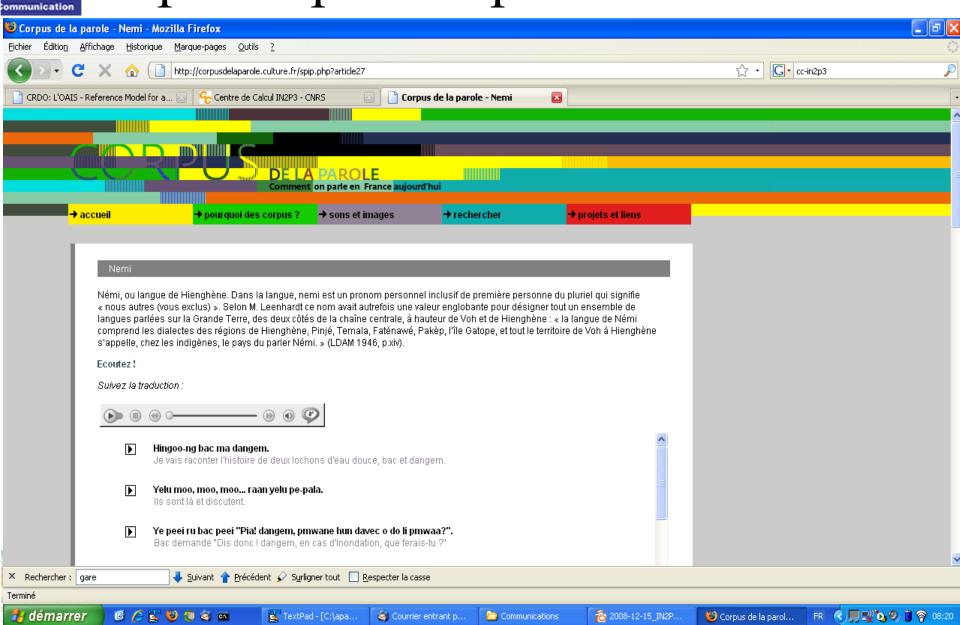
### 3.3. Publication référencement

- Le CRDO s'est organisé sous la forme d'un entrepôt OAI. Le protocole OAI-PMH est implémenté sous forme d'une sevlet java.
- Des fournisseurs de ressources moissonnent régulièrement l'entrepôt et offrent des moteurs de recherche dans un ensemble plus large d'entrepôts
- Des portails web moissonnent régulièrement l'entrepôt pour offrir des consultations particulieres des ressources ou des vues partielles de l'entrepôt
  - Corpus de la parole: portail sur les langues de France
  - Dallith: portail sur les langues de la zone himalayenne (archives du CRDO + Univ of Virginia)
  - Lacito: portail sur les ressources des chercheurs de ce labo



Cultur

## http://corpusdelaparole.culture.fr/





# http://dallith.vjf.cnrs.fr/

