

LDSS seminar, 12<sup>th</sup> July 2022

# Language documentation, language description, and language revitalization: how and why?

Peter K. Austin

Department of Linguistics  
SOAS, University of London

© 2022 Peter K. Austin

Creative commons licence

Attribution-NonCommercial-NoDerivs

CC BY-NC-ND

[www.peterkaustin.com](http://www.peterkaustin.com)

# General issues for discussion

1. Collecting and analysing language research materials (a corpus)?
2. What do we mean by 'language documentation', 'language description', 'language revitalisation'?
3. How do you analyse a corpus?
4. How do you archive a corpus?
5. How can a corpus be used? – for description, revitalisation, other purposes?

# What is a corpus (plural corpora)?

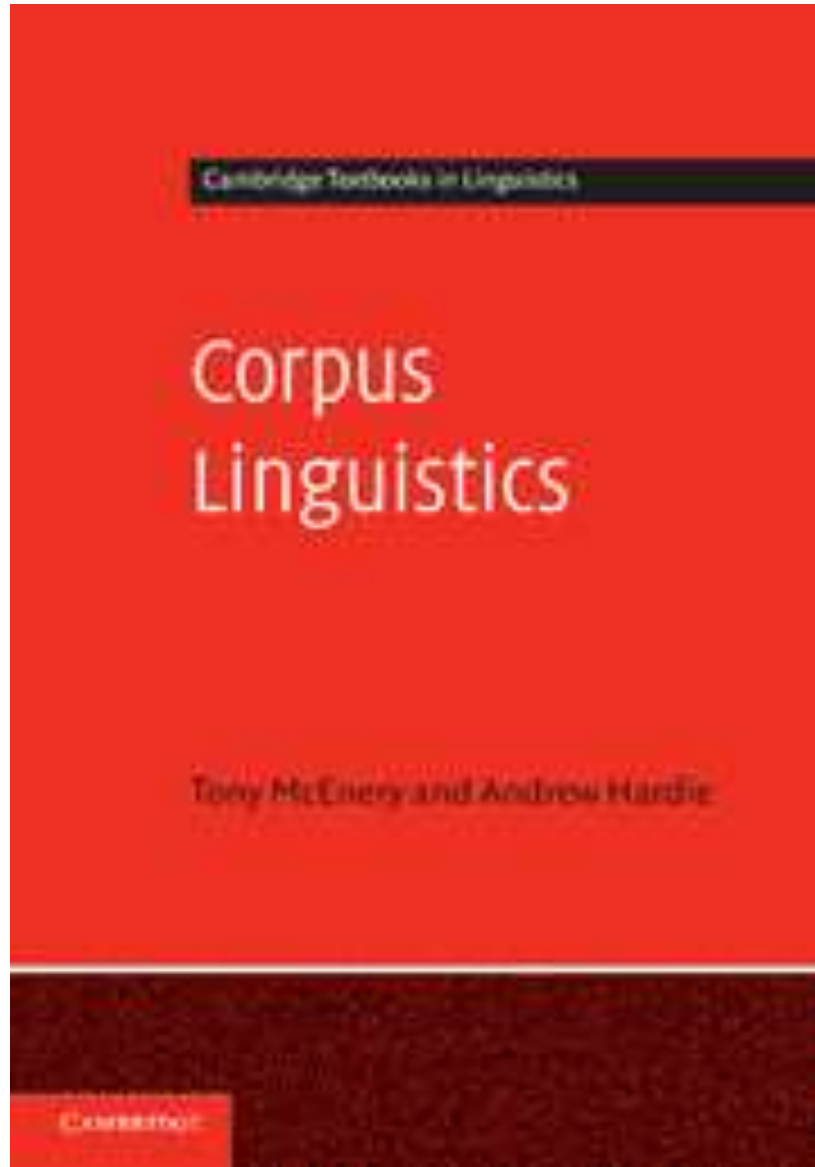
**Traditional definition:** A collection of linguistic data, either compiled as written texts or as a transcription of recorded speech. The main purpose of a corpus is to verify a hypothesis about language - for example, to determine how the usage of a particular sound, word, or syntactic construction varies. Corpus linguistics deals with the principles and practice of using corpora in language study. A computer corpus is a large body of machine-readable texts.

(Crystal, David. 1992. *An Encyclopedic Dictionary of Language and Languages*. Oxford: Blackwell.)

## **Note:**

1. emphasis on written text;
2. primarily quantitative analysis;
3. software tools;
4. particular collection is justified by research hypothesis (goals)

# A good introduction



Tony McEnery & Andrew Hardie. 2011. *Corpus Linguistics: Method, Theory, and Practice*. Cambridge: Cambridge University Press.

# How do you create a corpus?

**Traditional approach:** collect a set of texts (already born digital or by scanning print publications):

1. a sample or collection which is **representative** with regards to the research hypothesis;
2. there may be an attempt to **balance** the corpus to represent various non-linguistic variables (e.g. x% novels, y% poetry, z% conversation)
3. with a defined size (e.g. 1 million words), and content (e.g. English literature).
4. e.g. British National Corpus (1980-90, 1 million words, multiple written genres) <https://www.english-corpora.org/bnc/>
5. e.g. frTenTen: Corpus of French Web (current, 10 billion words, European Canadian and African French) <https://www.sketchengine.eu/frtnten-french-corpus/>

# Language description

- Study of language as a **system** separated from its actual use by speakers and the social-political-cultural-economic conditions of use
- Requires abstraction and search for general principles (phonology, morphology, syntax, semantics, pragmatics)
- Requires idealisation and “cleaning up” recordings of actual use
- Data collection often involves **elicitation** through surveys or interviews or **experiments**
- Studying a language the researcher does not speak is often done via **translation** or asking for **speaker judgements**
- The records of interview or survey are **not** of interest in themselves, but just a way to accumulate “**the data**” for analysis

## For discussion

1. Why do we want to describe languages?
2. Give some examples of language description
3. Who is the audience for language description?



# Language description

## 1. **Goals** of description:

- to present language structures for others to understand;
- to identify common features and differences across languages (typology);
- to understand how the human mind works (psychology, neurophysiology);
- to understand how humans interact and express personal, social and cultural relationships

## 2. Analysis is often highly structured and formal and written in an abstract **metalanguage**

## 3. the **audience** for description is typically other researchers, and distributed in books or articles (grammars, dictionaries, maps, graphs, narratives, text collections)

# A new approach: language documentation

- “concerned with the **methods, tools, and theoretical underpinnings** for compiling a **representative** and **lasting multipurpose** record of a natural language or one of its varieties” (Himmelman 1998)
- Features:
  - *Focus on primary data*
  - *Accountability*
  - *Long-term storage and preservation of primary data (archiving)*
  - *Interdisciplinary teams*
  - *Cooperation with and direct involvement of the speech community*
- Outcome is **annotated and translated corpus** of archived representative materials on a language, cf. TLA/Dobes, ELAR
- Woodbury (2003, 2011) ‘transparent records of a language’

# What's new in documentation?

- **Data focus** – this is Himmelman's “primary data” and includes **audio, video, still images**, and **text**, but also **structured data** derived from processed materials (transcribed, translated, annotated digital files). A collection of such material is called a **corpus**. See Himmelman 2012. We discuss problems with this later.
- **Accountability** – we expect the materials (“primary” and analysed) to be made available to others. Some have argued for **reproducibility**, i.e. the possibility of recreating the researcher's analytical steps to see if the outcome is the same (or different). See Berez-Kroeker et al 2017. We discuss problems with this later.
- **Preservation** -- long-term storage in safe archival facilities where the data and analysis (corpora) can be safeguarded for the long term (including refreshing data formats to take into account changing software)
- **New software tools** – data and analysis is stored in digital files and access is mediated via computer software (e.g. Praat, ELAN, FLEx)

## For discussion

- Why do we want to document languages/dialects by collecting instances of language performance (use in social-cultural-political-economic context)?
- If languages/dialects are disappearing, what is the point of studying them?

# Metadata and meta-documentation

In order to organise, manage, understand, and analyse the corpus we need:

- **Metadata** – data about the data – several types at file-level (see Nathan & Austin 2004):
  - *cataloguing* — title, speakers, collectors, time and place of recording, language name etc.
  - *descriptive* — information about content, relationship to other resources etc.
  - *structural* — what structural devices and patterns exist in the document etc.
  - *technical* — performance and preservation information, description of formats etc.
  - *administrative* — work log, responsibilities, access protocol statements etc.
- **Meta-documentation** – metadata at the project level: goals, corpus theory, data collection and analytical methods, stakeholders, ethics (informed consent), access and use

This is also important for archiving (see below)

# Data collection methods

## 1. **Elicitation** (interviewing):

- Translation (L1 → L2, L2 → L1)
- Speaker judgements: “can you say xxxxxx?”, “does yyy sound/mean the same as xxx?”

## 2. **Narratives** (telling stories)

## 3. **Conversation** (2 or more participants)

## 4. **Experimentation** – puzzles, games, other tasks (video description, image description = use of “stimuli”)

## 5. **Participant observation** – spending time with speakers observing language use and attempting to use the language oneself

Each method has strengths and weaknesses

# Metadata collection and management

1. Can be done manually (pen and paper) but best done electronically for ease of storage, searching, sharing, and restructuring
2. Plain text files, spreadsheet (Excel), or database (Access, MySQL)
3. Dedicated metadata software – SayMore (stand alone), CIMDI maker (online) etc.

Archives may have preferences on metadata tools and formats – check when you are first designing your project and planning your corpus

# Secondary corpora ('legacy materials')

- It is rarely the case that first-hand research is carried out on languages or communities that have never been documented before, so typically there already exists material in some form, in missionary or traveller reports, government records, or from previous linguistic or anthropological researchers. With careful use, these **legacy materials** can provide valuable information to contemporary researchers and communities, and assist language recovery or revitalisation
- In some cases there are no contemporary fluent speakers and legacy materials are the richest or only sources for description and revitalisation
- Sometimes, field research in communities is not possible due to danger from violence, e.g. civil war or gangs, or from disease, including pandemics like Ebola and Covid-19



## Secondary corpora example

- Peter Austin *Toolbox* databases (lexicon, glossed and annotated texts) from S.A. Wurm's 1955-57 handwritten notes of Australian languages
- Retype original, add metadata on sources (speaker, recorder, fieldnotes location), abbreviation definitions, date of last edit
- Add sentence analysis: phonemicization, morpheme glossing, part of speech, free translation in English, notes, link to lexicon (lexnum), link to abbreviations
- Add lexicon: headword, gloss, definition, scientific name, scientific name source, picture, semantic relations (synonym, antonym, cf), notes, cognates, example sentence link (text, free translation)

Epi Talar (Mac Talar) (R) (M.L.) daddja layambuni  
wanju inipi punju - ε

ε - ε, (0 00 / 000000 - 000)

a'ei, wala to daddja layambuni inipi punju

minju inipi punju ε

banimama lea lam daddja lambuna  
wanima gagulunu

buppani punju ε, ε ε ε ε ε ε ε ε ε ε

punju ε

inipida jalida (punju ε ε ε ε ε ε ε ε ε ε)

malyangapa_lexicon	
\x	<b>dhadja</b>
\a	
\u	
\xnum	029
\g	MI
\ps	vtr
\ge	bite
\def	bite
\sd	Actions
\eth	
\sci	
\sci_source	
\nt	
\syn	
\ant	
\cf	<b>dhaba</b>
\se	
\gcf	
\ety	
\rec	SW
\sp	HQ
\x_ref	011
\v	<b>dhadjangarndambunyi gunyungu</b>
\xe	The dog bit me
\wnum	425· 426· 427· 428· 429· 430· 431

malyangapa_notes	
\snum	017
\text	<b>wanyu yinigi gunyuyi, dhadjalangambuni</b>
\morpheme	<i>wanyu yinigi gunyu -yi dhadja -langa -mbu -ni</i>
\gloss	bad that dog -emph bite -might -3sg.A -2sg.P
\ps	n dem n -suff vtr -suff -suff -suff
\xnum	028 035 002 -034 029 -049 -031 -050
\free_translation	This dog is bad, it will bite you
\reference	SW1/2As01
\recorder	SW
\speaker	HQ
\note	wannju i:nigi gunju-?: daddja la?ambuni
\date	11/Sep/2020

Malyangapa_abbreviations	
\abb	<b>vtr</b>
\mng	transitive verb
\type	<b>sub-category</b>
\nt	transitive verbs are a sub-category of verbs; they take a transitive subject argument in ergative case and a transitive object argument in accusative case.
\cf	<b>vi, vdi</b>
\gr	
\date	11/Sep/2020

# Archiving

- An archive is a trusted repository with a collection policy and a commitment to:
  - appraise the value of certain materials
  - preserve selected items
  - make known their existence
  - enable access to them (or their ‘content’)
- Archives have a catalogue that presents metadata (data about the data in the archive), often in a standardized format, some have **finding aids**
- Archives have access management protocols
- Many funders now require that projects archive their materials

## Henke & Berez-Kroeker (2016: 411)

“It is difficult to imagine a contemporary practice of language documentation that does not consider among its top priorities the **digital preservation** of endangered language materials. Nearly all handbooks on documentation contain chapters on it; conferences hold panels on it; funding agencies provide money for it; and even this special issue evinces the **central role of archiving** in endangered language work. In fact, archiving language data now stands as a regular and normal part of the field linguistics workflow (e.g., Thieberger & Berez 2011).” [emphasis added]

**Note:** there is a free online course about archiving at <https://archivingforthefuture.teachable.com/>, but it does not cover how to use other people’s collections or legacy materials

# Important

## A website is **NOT** an archive!!!



1. Websites are volatile – they come and go, and rarely have the institutional support like an archive does
2. The files on websites can become obsolete and no longer accessible; archives plan for ‘forward migration’ of file formats
3. Access to websites cannot typically be controlled to the fine degree that archives allow
4. Anyone can put anything on the web – archives involve collection policies, selection, and curation (quality control)



# Archive types

## 1. Classified according to the types of material:

- **Physical** (analogue) archives – contain paper records, tape recordings, physical objects, e.g. Smithsonian Institution, British Library, Bibliothèque nationale de France
- **Digital** archives – contain digital files only: audio-visual, text, still images, maps, e.g. ELAR, TLA, AILLA (see DELAMAN for a list)
- **Mixed** archives – contain analogue and digital materials, e.g. AIATSIS, CLA, ANLA

## 2. Classified according to scope:

- **International** – world-wide or multi-country coverage, e.g. ELAR, TLA, BL, BNdeF, AILLA, Pangloss
- **National** – cover one country, e.g. AIATSIS
- **Regional** – cover an area in a country, e.g. CLA, ANLA
- **Local** – cover a town or community, e.g. local museums
- **Personal** – records of an individual or family

# Large international digital – ELAR

 Endangered Languages Archive

Show deposits:  Curated  In-process  Forthcoming

Find a deposit:  
List  
Map

Help  
Home

ARCADIA  
 SOAS

POWERED BY Google



Map data ©2013 MapLibre



# Large international digital – TLA DOBES



# Archive access management



- Universal – resource available to all, e.g. online
- Register – resource available to registered users
- Closed – resource not generally available (embargoed, “black box”)
- Strict – resource available to users who have been given *individual* access rights for that resource

# Language revitalisation

- efforts to increase **language vitality** by taking action to:
  - increase the domains of use of a language and/or
  - increase the number of speakers (often in the context of reversing language shift) both adults and children
- older than language documentation (serious work began in 1970s and 1980s among Maori, Native American groups and others)
- Speech/language community members are often more interested in revitalisation than documentation
- Often assumed revitalisation = formal language learning (school lessons, immersion)
- Many communities are now using corpora to support language learning

# Approaches to language revitalisation

- Often assumed revitalisation = formal language learning (school lessons, immersion, bilingual education)
- Other methods may be easier to establish and to identify resources for:
  - Language landscape
  - Language nests
  - Master-apprentice programmes
  - Cultural immersion and informal learning: camps, seeing and doing
- Many communities are now trying to use descriptive and corpora materials to support language learning – we discuss issues with this further in tomorrow's seminar

## For discussion

Does it make any sense to talk about “revitalisation” in Italy?

Should we distinguish between “real languages”, e.g. Ladin, Arberesh, and “dialects”, e.g. Abruzzese?

How?

# How can an online archived corpus be used – for description, revitalisation, other purposes?

- Online corpora might seem at first sight to be great potential sources of data to replicate the research of others, to explore topics not covered in the analysis so far, or for language learning and revitalisation
- However, there are often numerous problems with using archived corpora:
  - Problems of principle concerning epistemology of the content
  - Problems of principle concerning scope and goals of the data collectors
  - Problems of practice concerning accessing the corpus and using it

# Archive content and interfaces

Wasson et al (2016: 669): ‘In their presentations, the archivists provided a rich list of problems that might be encountered by users of language archives. The most frequently mentioned items were:

- A lack of contextual information at the deposit level, or metadata
- Incomplete materials—missing annotation, missing translations
- Inadequate search/browse functions
- Problems with the interface/information display
- Users may be frustrated when they don’t have access to data; it may be hard for the archivist to get hold of a collection owner to request access for a user
- Technology issues—outdated, broken scripts, Flash/Java problems, etc.
- Interface language(s) may not [be] ... spoken by would-be users’

# Corpora epistemology issues

- Corpus materials do not exist in a vacuum but have their own socio-cultural context within which they were created, their “**social lives**” (Dobrin & Schwartz 2021), reflecting the nature and history of people (researchers, consultants, community), relationships, interactions, assumptions, goals (Christensen 2018)
- These are often **implicit** and need to be understood in order to make sense of the materials, however they are rarely documented or made explicit by researchers as meta-documentation (documentation of the documentation, cf. Austin 2013), so research on the corpus and its creation needs to be undertaken (see Austin 2017)



# Corpora epistemology issues

- **biography** of creator(s): prior language knowledge and/or study and/or exposure, their teachers/mentors/correspondents, how/when they learnt the language, how long they worked on the language and at what point in their careers, how the work was funded and with what goals, whether there were previous studies of the language or the community that they could have had access to, who was the consultant, what prior experience, who did the value-adding (transcription, translation)
- aspects of **historical period**: kind and impact of contact between communities, including colonialists, and influential descriptive categories and formats known to author(s), e.g. traditional grammar based on Latin or Greek models, structural linguistics, IPA

# Form, content, and interpretation issues

- Audio may be poor quality, noisy, difficult to hear (e.g. multiparty conversations)
- Video may be poorly recorded or unwatchable, poor audio, partially record people out of frame
- Hand-written or typed text can be difficult to read or interpret, with crossing out, abbreviations, obscurities (requires philological analysis)
- In digital text files, characters may be mismapped or missing due to font problems, tabbed text may not align, structured text may be uninterpretable if the structure definition is missing
- Retranscriptions should link back to documents or files on which they are based (so we can retrace the steps), XML representations can create various different outputs ('diplomatic edition', edited (clean) edition), e.g. Dawes MS



Book B, Page 8

◀ Prev

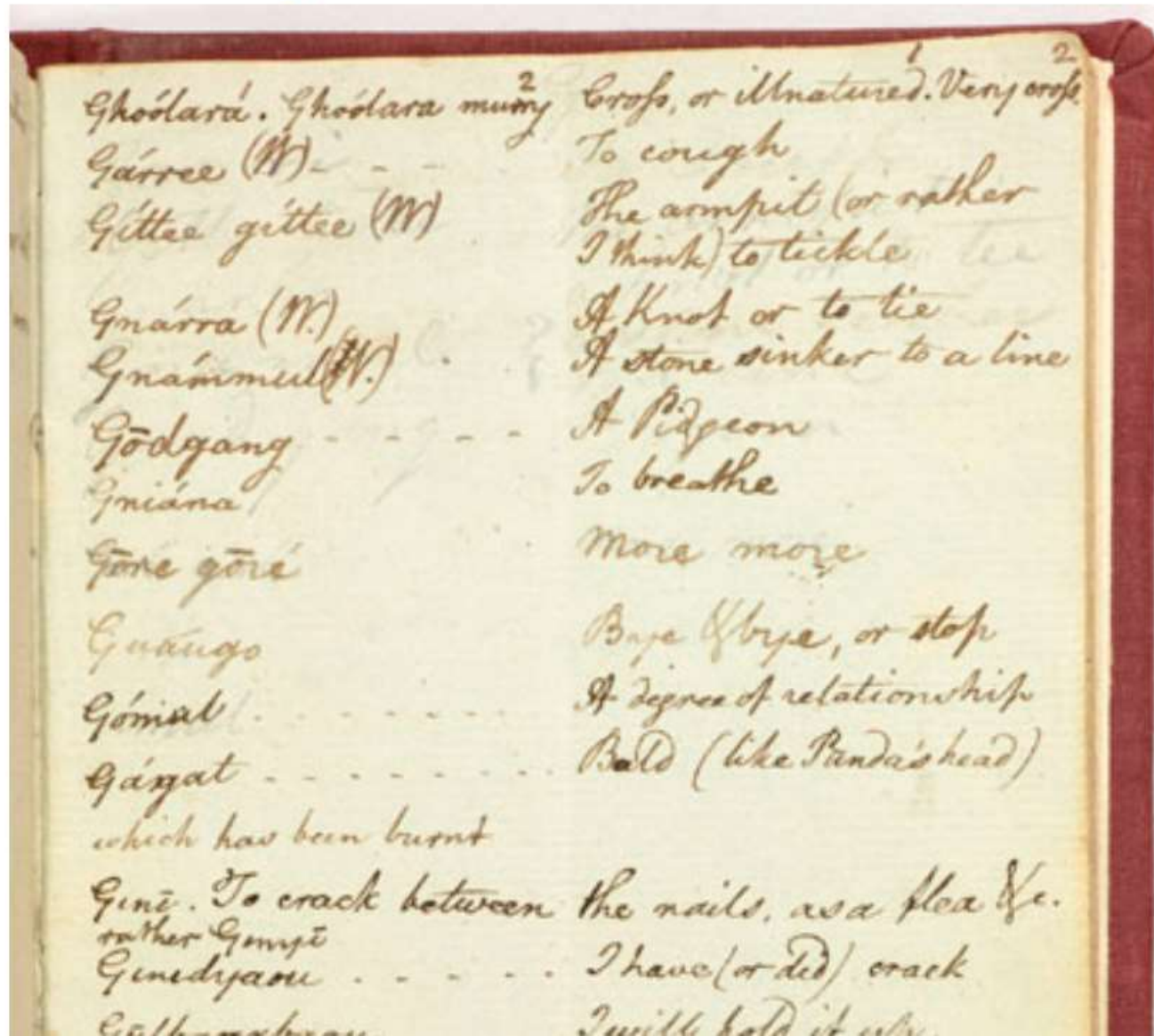
Contents ▲

Next ▶

Manuscripts Book A ▼

Manuscripts Book B ▼

Manuscripts Book C ▼



Ghoólara <sup>1</sup> . Ghoólara murry <sup>2</sup>	Cross, or illnatured <sup>1</sup> . Very cross <sup>2</sup> .
Gárree (W)	To cough
Gíttee gíttee (W)	The armpit (or rather I think) to tickle
Gnárra (W.)	A knot or to tie
Gnámmul (W.)	A stone sinker to a line
Gödgang	A Pidgeon
Gniána	To breathe
Göre göré	More more
Guáugo	Bye & bye, or stop
Gómül	A degree of relationship
Gáñat	Bald (like Pūnda's head)
which has been burnt	
Ginī. To crack between the nails, as a flea &c.	
rather Ginyi	I have (or did) crack
Gindyaou	I will hold it up
Gonanjúye desiring to wear one of Patyegaran's pettycoats: I told her it was too long for her; on which she said Gūlbarabaou which Patye explained as above. –	
Gwára buráwá	The wind is fallen.



## Book B, Page 8

◀ Prev

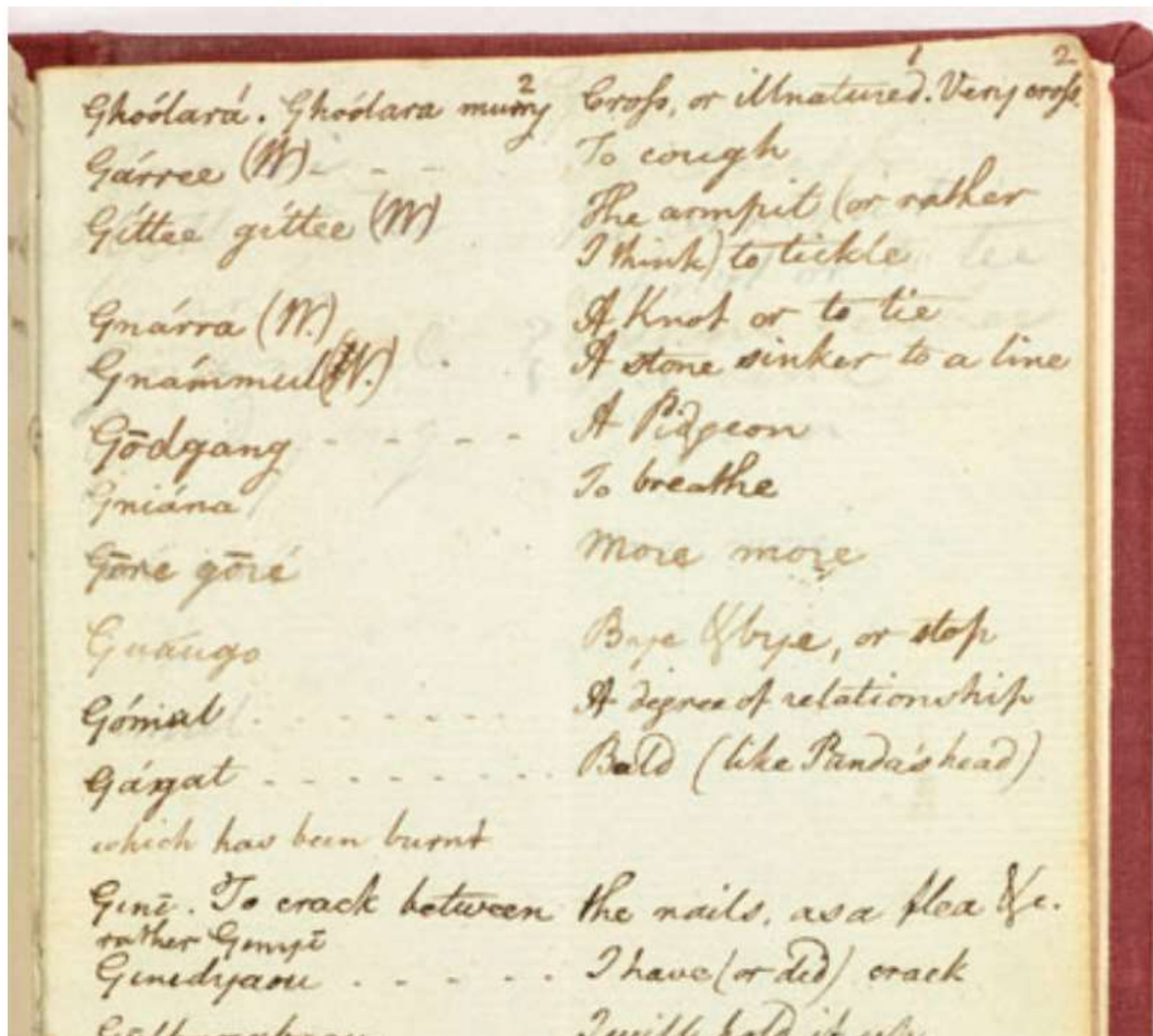
Contents ▲

Next ▶

Manuscripts Book A ▼

Manuscripts Book B ▼

Manuscripts Book C ▼



*Ghoólara*<sup>1</sup>. *Ghoólara murry*<sup>2</sup>

*Gárrée* (W)

*Gíttee gíttee* (W)

*Gnárra* (W.)

*Gnámmul* (W.)

*Gōdgang*

*Gniána*

*Gōre gōré*

*Guángo*

*Gómúl*

*Gáŋat*

*which has been burnt*

*Ginī*. *rather Ginyī To crack between the nails, as a flea etc.*

*Ginedyaou*

*Gūlbaŋabaou*

*Gonaŋulye* *desiring to wear one of Patyegarāŋ's pettycoats: I told her it was too long for her; on which she said Gūlbaŋabaou which Patye explained as above. –*

*Gwára buráwā*

*Cross, or illnatured*<sup>1</sup>. *Very cross*<sup>2</sup>.

*To cough*

*The armpit (or rather I think) to tickle*

*A knot or to tie*

*A stone sinker to a line*

*A Pidgeon*

*To breathe*

*More more*

*Bye & bye, or stop*

*A degree of relationship*

*Bald (like Pūnda's head)*

*I have (or did) crack between the nails*

*I will hold it up*

*The wind is fallen.*

# Form, content, and interpretation issues

- implicitly structured materials, e.g. using typography or page layout to distinguish analytical categories or kinds of information, can be made more useful by encoding the structure separately from the form, e.g. by XML markup, or database model
- structure is not always computable from typography and may need to be manually added (cf. quotes from Nathan in Austin 2017): over-use of quotation marks, for multiple purposes and often redundantly, unclear scoping, spelling errors
- cryptic or incorrect glossing, because author(s) could not understand language consultant's accent or pronunciation, or because the semantics of the source language terms were misunderstood (the "gavagai problem")
- changing interpretations over time, and the author(s) hearing what they think they heard rather than what is in the recording or was dictated by speakers due to analytical decisions

Table 1 Semantic problems in text materials, from Crowley and Austin (2005).

	Wordlist meaning	Correct meaning
1.	Pronunciation problems	
	heart	hot
	wet	sweat
	moths	boss
	dung, shit	tongue
2.	Meaning problems	
a.	generic versus specific	
	grass	vegetation
	boy	uninitiated youth
	beard	hair
	day	now
	thumb	your hand
	girl	female
b.	related word	
	thighs	buttocks
	cloud	sky
	woman	wife
	hair	head
	frown	blind
	spider	to bite
	dig	drink

# Form, content, and interpretation issues

- Inappropriate content (taboo, sacred) for various audiences
- Dated content using expressions that are no longer acceptable or now inappropriate, e.g. personal remarks about the ancestors of living persons
- Over-distinguishing or under-distinguishing crucial contrasts, in phonology (voicing, aspiration, vowel quality or quantity), morphology (ergativity), or syntax (vocative case, applicative, switch-reference)
- Lack of sociolinguistic context: who says what to whom when and where?
- Relationship between corpus forms and contemporary usage – issues about what is “right”, especially for shifting languages undergoing change

# Stakeholder issues

- projects typically have many stakeholders who may have different kinds of interests in the materials collected and the analyses created
- Issues of control, consultation, and decision-making are important when deciding what kind of documentary material to include in any corpus and how it can be used
- For legacy materials possible mismatches between past situations and the present
  - current membership of a contemporary ‘community’ may not coincide with past membership
  - people who provided legacy materials may not even now be viewed as rightful members of a given group and therefore their information may be deprecated



# Stakeholder issues

- Unclear agreements, if any, between original collector and the community or particular individuals at the time (and whether these agreements were documented) as well as the relationship between any such agreements and arrangements that are currently being negotiated between contemporary researchers and other stakeholders, e.g. Austin told not to distribute copies of Wurm's materials without permission of current Aboriginal group who self-identify as descendants
- Best to clarify if possible before creating and using the corpus, especially for North America and Australia

# Conclusions

- Creating and analysing corpora can be very rewarding and enable various exciting kinds of linguistic and cultural research to be done
- However, working with corpora involves dealing with often **complex issues** about the form, content, context, and use of materials and analyses arising from them
- Good corpus management principles and practices (e.g. file naming, folder structure, backup, software tools) will make life easier. It is essential to build in archiving plans from the beginning of a project
- maximising opportunities to use a corpus requires thinking about data entities, data types and relationships, and being **explicit** about them in the project design and application (e.g. in database design or XML tagging)
- very important role for **metadata** and **meta-documentation**
- by creating good meta-documentation now we can reduce legacy data problems for future researchers

# Conclusions

- there are many **opportunities** for researchers to add substantial value to corpus materials, and create **secondary** corpora, especially if they are able to work with other historical sources and/or contemporary knowledge holders to elucidate them and the context surrounding their creation, analysis and current status
- careful work with corpora can also be very **rewarding** for researchers and communities, especially for unique documents on languages/varieties or areas of knowledge that are no longer available, and that can serve as **important sources** for language support and revitalisation
- Thank you for your attention

# Abbreviations

AIATSIS	Australian Institute of Aboriginal and Torres Strait Islander Studies
AILLA	Archive of the Indigenous Languages of Latin America (UTexas Austin)
ANLA	Alaskan Native Languages Archive
APS	American Philosophical Society
BL	British Library
BNdeF	Bibliothèque nationale de France
CLA	California Languages Archive (UC Berkeley)
DELAMAN	Digital Endangered Languages and Musics Archives Network
ELAR	Endangered Languages Archive (SOAS University of London)
SI	Smithsonian Institution
TLA	The Language Archive (MPI Nijmegen)

# References

Austin, Peter K. 2013. Language documentation and meta-documentation. In Mari Jones & Sarah Ogilvie (eds.) *Keeping Languages Alive: Documentation, Pedagogy and Revitalization*, 3-15. Cambridge: Cambridge University Press.

Austin, Peter K. 2016. Language documentation 20 years on. In Martin Pütz & Luna Filipović (eds.) *Endangered Languages and Languages in Danger: Issues of ecology, policy and human rights*, 147-170. Amsterdam: John Benjamins.

Austin, Peter K. 2017. Language documentation and legacy text materials. *Asian and African Languages and Linguistics* 11, 23-44.

Berez-Kroeker, Andrea L., Lauren Gawne, Susan Smythe Kung, Barbara F. Kelly, Tyler Heston, Gary Holton, Peter Pulsifer, David I. Beaver, Shobhana Chelliah, Stanley Dubinsky, Richard P. Meier, Nick Thieberger, Keren Rice & Anthony C. Woodbury. 2017. Reproducible research in linguistics: A position statement on data citation and attribution in our field. *Linguistics* 56(1), 1-18.

# References

- Christen, Kim. 2018. Relationships, not records: Digital Heritage and the Ethics of Sharing Indigenous Knowledge Online. In Jenetry Sayers. (ed.) *Routledge Companion to Media Studies and Digital Humanities*, 403-412. London: Routledge.
- Dobrin, Lise & Saul Schwartz. 2021. The social lives of linguistic field materials. *Language Documentation and Description* 21.
- Gawne, Lauren, Barbara F. Kelly, Andrea L. Berez-Kroeker & Tyler Heston. 2017. Putting practice into words: The state of data and methods transparency in grammatical descriptions. *Language Documentation & Conservation* 11, 157–189.
- Henke, Ryan & Andrea L. Berez-Kroeker. 2016. A Brief History of Archiving in Language Documentation, with an Annotated Bibliography. *Language Documentation & Conservation* 10, 411-457.
- Himmelmann, Nikolaus P. 1998. Documentary and descriptive linguistics. *Linguistics* 36, 161–195.
- Himmelman, Nikolaus P. 2012. Linguistic Data Types and the Interface between Language Documentation and Description. *Language Documentation and Conservation* 6, 187-207.
- McEnnery, Tony & Andrew Hardie. 2011. *Corpus Linguistics: Method, Theory, and Practice*. Cambridge: Cambridge University Press.

# References

Nathan, David & Peter K. Austin. 2004. Reconceiving metadata: language documentation through thick and thin. *Language Documentation and Description* 2, 179-187.

Nathan, David, Susannah Rayner & Stuart Brown (eds.) 2009. *William Dawes : notebooks on the Aboriginal language of Sydney : a facsimile version of the notebooks from 1790-1791 on the Sydney language written by William Dawes and others*. London: SOAS. (see also [www.williamdawes.org](http://www.williamdawes.org))

Rice, Keren. 2006. et the language tell the story? The role of linguistic theory in writing grammars. In Felix K. Ameka, Alan Charles Dench & Nicholas Evans (eds.) *Catching Language: The Standing Challenge of Grammar Writing*, 235-268. Berlin: Mouton de Gruyter.

Thieberger, Nicholas & Andrea L. Berez. 2011. Linguistic data management. In Nicholas Thieberger (ed.) *The Oxford handbook of linguistic fieldwork*, 90-118. Oxford: Oxford University Press.

Warner, N., Q. Luna, and L. Butler. 2007. Ethics and revitalization of dormant languages: The Mutsun language. *Language Documentation & Conservation* 1(1), 58-76.

Warner, N., Q. Luna, L. Butler & H. van Volkinburg. 2009. Revitalization in a scattered language community: Problems and methods from the perspective of Mutsun language revitalization. *International Journal of the Sociology of Language* 198, 135-148.

# References

Woodbury, Anthony C. 2003. Defining documentary linguistics. *Language Documentation & Description* 1, 35-51.

Woodbury, Anthony C. 2011. Language documentation. In Peter K. Austin & Julia Sallabank (eds.), *Cambridge handbook of endangered languages*, 159–186. Cambridge: Cambridge University Press