# Using the Tromsø Recommendations to cite data in language work

**Helene N. Andreassen (@n_andreassen), Andrea L. Berez-Kroeker, Lauren Gawne** (@superlinguo), Philipp Conzett (@PhilippConzett), Koenraad De Smedt, Christopher Cox, Lauren B. Collister (@parnopaeus)

**bit.ly/TRecsICLDC21 #lingdata**

Andrea, Helene & Lauren

Montréal, September 2017

Photo: private

# Overview: The Tromsø Recommendations For Citation of Research Data in Linguistics

Language and linguistics datasets are often not cited well.

The people involved in creating language data are not receiving proper recognition.

The Tromsø Recommendations detail how to cite data - and people! - in language work.

# What do we mean by "language work"?

Leonard 2017:

"an umbrella expression to include language **documentation**, **description**, **teaching**, **advocacy**, and **resource development**" (2017:16)

...everything that ICLDC participants do!

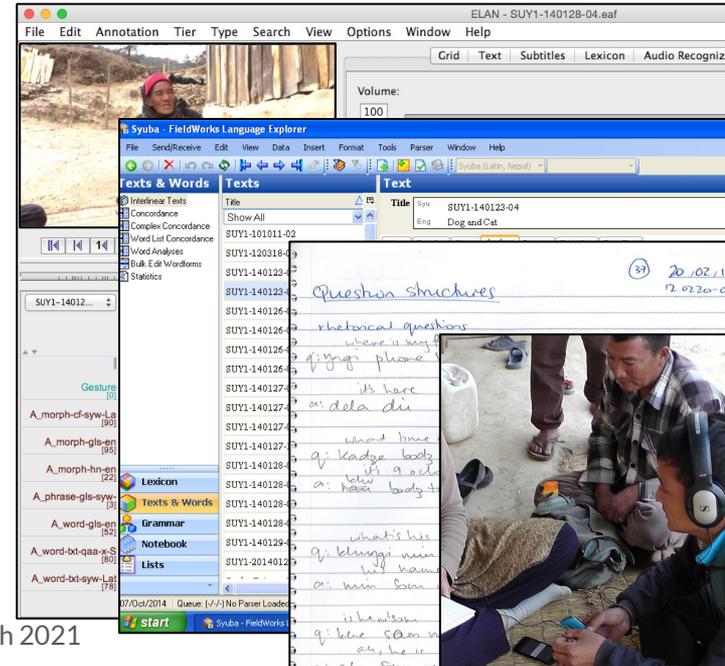# What do we mean by "data" in language work?

**Data** in language work means all samples of language:

Recordings and written language

Words, sentences, verbal art, storytelling, oration, song, etc.

**People** create and contribute to language data.

Data in language work is **precious** because language is about **people**.

# What do we mean by "citing" data?

**People** are speakers, learners, elders, youth, teachers, helpers, parents, grandparents, researchers, culture-bearers, translators, poets, authors and more.

**People** deserve to be thanked - credited - acknowledged for their contributions in language work.

**Citation** is one way to do that.

We do this for publications **All. The. Time.**

# Typical citations (of publications, not data)

Good, Jeff. 2011. Data and language documentation. In Peter K. Austin & Julia Sallabank (eds.), *The Cambridge handbook of endangered languages*, 212–234. Cambridge: Cambridge University Press.

Haspelmath, Martin & Susanne Maria Michaelis. 2014. Annotated corpora of small languages as refereed publications: A vision. *Diversity linguistics comment*. http://dlc.hypotheses.org/691 (accessed 10 January 2017).

(Acknowledgment given to **author** and **editors**)

# Hmmm.... we've been here before



**Putting practice into words: Fieldwork methodology in grammatical descriptions**

ICLDC 4 February 26-March 1, 2015 Honolulu

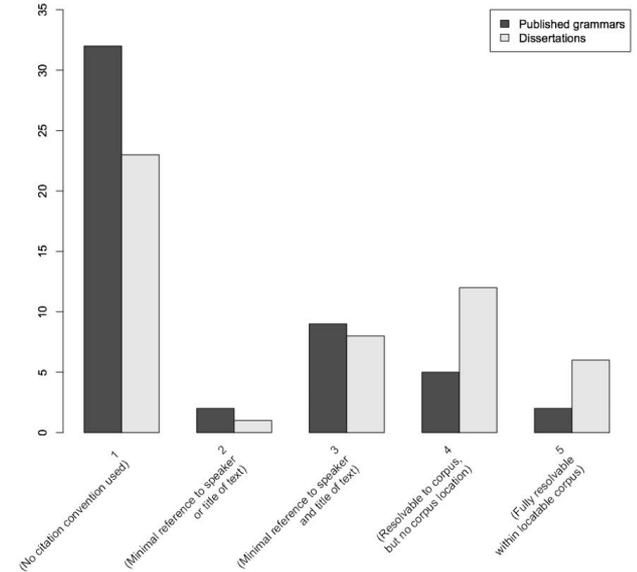Lauren Gawne(1), Barbara F. Kelly(2), Andrea Berez(3), & Tyler Heston(3)

(1) NTU Singapore, (2) The University of Melbourne, (3) The University of Hawaii

ICLDC 2015: language documenters aren't that great at citing **data**.

(Gawne et al. 2015, Gawne et al. 2017)

Which means we aren't giving proper credit to the **people** involved in language work.

Why not?

Linguistics doesn't have a history of requiring the citation of data.

(even though we *do* have a history of requiring the citation of *publications*)

# Documenters don't really do this for data

In 2015/2017 we were concerned about the effects of this on our **science**.

At this ICLDC today we are concerned about the effects of this on **people** who are being left out.

People serve a lot of **roles** in language work.



Gawne et al. 2017

# Many roles for people in language work

| | | | |
|---|---|---|---|
| Author | Translator | Illustrator | Participant |
| Editor | Recorder | Participant | Depositor |
| Speaker | Data inputter | Interviewer | Developer |
| Signer | Consultant | Compiler | Sponsor |

So **why** aren't we giving people credit for these roles through proper citation?

# **Because we don't know how!**

Part of the problem is that we don't know **how** or **why** to cite data (Berez-Kroeker et al. 2018).

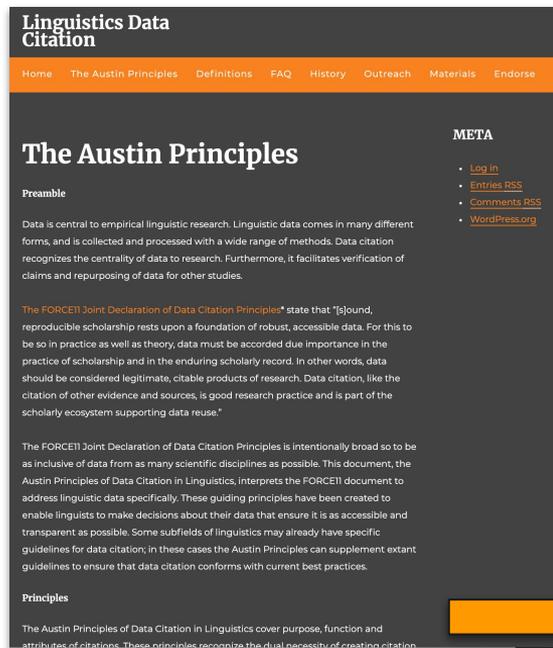2017: Research Data Alliance Linguistic Data Interest Group (link below!)

Two key publications:

For the WHY: The Austin Principles of Data Citation (link below!)

For the HOW: Tromsø Recommendations for Citation of Research Data in Linguistics (link below!)

# Austin Principles of Data Citation in Linguistics
## www.linguisticsdatacitation.org

**Linguistics Data Citation**

Home   The Austin Principles   Definitions   FAQ   History   Outreach   Materials   Endorse

## The Austin Principles

**Preamble**

Data is central to empirical linguistic research. Linguistic data comes in many different forms, and is collected and processed with a wide range of methods. Data citation recognizes the centrality of data to research. Furthermore, it facilitates verification of claims and repurposing of data for other studies.

The FORCE11 Joint Declaration of Data Citation Principles* state that "[s]ound, reproducible scholarship rests upon a foundation of robust, accessible data. For this to be so in practice as well as theory, data must be accorded due importance in the practice of scholarship and in the enduring scholarly record. In other words, data should be considered legitimate, citable products of research. Data citation, like the citation of other evidence and sources, is good research practice and is part of the scholarly ecosystem supporting data reuse."

The FORCE11 Joint Declaration of Data Citation Principles is intentionally broad so to be as inclusive of data from as many scientific disciplines as possible. This document, the Austin Principles of Data Citation in Linguistics, interprets the FORCE11 document to address linguistic data specifically. These guiding principles have been created to enable linguists to make decisions about their data that ensure it is as accessible and transparent as possible. Some subfields of linguistics may already have specific guidelines for data citation; in these cases the Austin Principles can supplement extant guidelines to ensure that data citation conforms with current best practices.

**Principles**

The Austin Principles of Data Citation in Linguistics cover purpose, function and attributes of citations. These principles recognize the dual necessity of creating citation

**META**

- Log In
- Entries RSS
- Comments RSS
- WordPress.org

## 2. Credit and Attribution

*Data citations should facilitate giving scholarly credit and normative and legal attribution to all contributors to the data, recognizing that a single style or mechanism of attribution may not be applicable to all data.*

In linguistics, citations should facilitate readers retrieving information about who contributed to the data, and how they contributed, when it is appropriate to do so. One way to do this is through citations that list individual contributors and their roles. Another way is by using citations that link to metadata about contributors and their roles.

Model: FORCE11 Joint Declaration of Data Citation Principles

# The Tromsø recs

- Minimal and expanded templates for in-text citations & bibliographic references
- Explanation of elements in clear terms
- Examples from real linguistic data
- Highlights issues that are important to linguistic data

# The Tromsø recommendations for citation of research data in linguistics

Developed through asynchronous meetings of the LDIG, plus invited input from VIPs

**Aim:** Practical and concise advice for data citation, with consideration of the variety of linguistic data

**Intended audience:** Editors of linguistic publications, researchers, and repositories.

# The Tromsø recommendations - Outline

Includes:

- Recommendations for in-text citation and bibliographic reference
- Full data set and specific example citation
- Examples of citation using real data
- Flexibility to fit with journal style guidelines

# Of note here: Roles

The T-Recs allow you to give credit to many people and explain the role they played.

For example: the data **collector:**

Adelaar, Alexander (Collector). 2005. *Ma'anyan narratives* (AA4). PARADISEC. https://doi.org/10.4225/72/56E979455A05E.

Also, **researchers, depositors, speakers, consultants, interviewers**…

Hauk, Bryn (Researcher, Depositor), Omar P'ap'ashvili (Speaker) & Rezo Orbetishvili (Consultant). 2018. BH2-074. In Batsbi (Tsova-Tush). Kaipuleohone University of Hawaii Digital Language Archive. http://hdl.handle.net/10125/58935.

Krauss, Michael E. (Interviewer), Jeff Leer (Interviewer) & Anna Nelson Harry (Speaker). 1975. Interview with Anna Nelson Harry. In Krauss Eyak Recordings, item ANLC0082. Alaska Native Language Archive. https://www.uaf.edu/anla/.

Even for **in-line citations**:

(Hauk 2018: BH2-081, 00:00:01–00:00:03, Rezo Orbetishvili (Speaker))

# There are many lists of Roles you can use

Some standard lists of contributor roles include

CASRAI

DataCite

OLAC role vocabulary

# Next: Data citation in your work

Working towards normalising the practice of citing linguistic data

Relevant to everyone who works with linguistic data

# Language workers: Cite your data

Build data citation into projects from the beginning

Data citation is distinct from, but closely related to, making underlying data available. Citation co-exists with ethical approaches to archiving and access

Cite other people's linguistic data if you use it in your work

# Supervisors & project managers: Encourage best practice

Introduce students to best practice in the field (cf. Pawley 2014)

**Example:** Data citation and archiving have been expectations at University of Hawaii since Fall 2013. Included in PhD student handbook.

# Publishers: Make citation an expectation

The Trømso Recommendations can be adopted by any journal or publication

**Example:** The [Australian Journal of Linguistics](#) guidelines include the Generic Style Rules for Linguistics, the Leipzig Glossing Rules, the Austin Principles & the Tromsø Recommendations:

"For research based on original fieldwork or archival documentary materials, authors must provide the sources and provenance of data, as well as the methods used to collect it, including the time period and locations in which fieldwork was conducted".

# Data managers: Encourage citation

Data managers can provide training and support to encourage citation.

**Example:** PARADISEC provides a "cite as" field on all pages of the archive, giving a formatted citation to the relevant level of granularity.

**Cite as** Lauren Gawne (collector), 2009. *Kagate (Nepal)*. Collection SUY1 at catalog.paradisec.org.au [Open Access]. https://dx.doi.org/10.4225/72/56E976A071650

# Normalising data citation in language work

Language documentation and reclamation is about **people**. **People** make language records.

Proper citation of linguistic records (data) gives credit to everyone involved.

The Tromsø Recommendations provide practical examples for how to cite linguistic data.

# References

- Andreassen, H. N. (2020). The Tromsø Recommendations for Citation of Research Data in Linguistics: Collaboration illustrated. *Ravnetrykk* 39: 114-122.
- Andreassen, H. N., Berez-Kroeker, A. L., Collister, L., Conzett, P., Cox, C., Smedt, K.D., McDonnell, B. and  and the Research Data Alliance Linguistic Data Interest Group. (2019). Tromsø recommendations for citation of research data in linguistics (Version 1). R*esearch Data Alliance*. DOI: 10.15497/RDA00040
- Berez-Kroeker, A.L., L. Gawne, S. Kung, B.F. Kelly, T. Heston, G. Holton, P. Pulsifer, D. Beaver, S. Chelliah, S. Dubinsky, R.P. Meier, N. Thieberger, K. Rice & A. Woodbury. 2018. Reproducible research in linguistics: A position statement on data citation and attribution in our field. *Linguistics* 56(1): 1–18.
- Berez-Kroeker, A.L., H.N. Andreassen, L. Gawne, G. Holton, S. Smythe Kung, P. Pulsifer, L.B. Collister, The Data Citation and Attribution in Linguistics Group, & the Linguistics Data Interest Group. 2018. *The Austin Principles of Data Citation in Linguistics.* Version 1.0. http://site.uit.no/linguisticsdatacitation/austinprinciples/
- Data Citation Synthesis Group: *Joint Declaration of Data Citation Principles*. Martone M. (ed.) San Diego CA: FORCE11; 2014 https://doi.org/10.25490/a97f-egyk
- Gawne, L., B.F. Kelly, A.L. Berez- Kroeker & T. Heston. 2015. Putting practice into words: Fieldwork methodology in grammatical descriptions. 4th International Conference on Language Documentation and Conservation. Hawaii: February 26-March 1.
- Gawne, L., B.F. Kelly, A.L. Berez- Kroeker & T. Heston. 2017. Putting practice into words: The state of data and methods transparency in grammatical descriptions. Language Documentation & Conservation 11: 157-189.
- Leonard, W.Y. (2017. Producing language reclamation by decolonising 'language'. In W. Y. Leonard & H. De Korne (Eds), *Language Documentation and Description*, vol 14. London: EL Publishing. pp. 15-36.
- Pawley, A. 2014. Grammar writing from a dissertation advisor's perspective. In T. Nakayama & K. Rice *The Art and Practice of Grammar Writing*, 7-23. University of Hawaii Press.

# Acknowledgments

Many thanks to the participants in the Data Citation and Attribution project, the Data Science for All of Linguistics project, members of the RDA LDIG, and attendees at our previous workshops, courses and presentations for fruitful discussion.
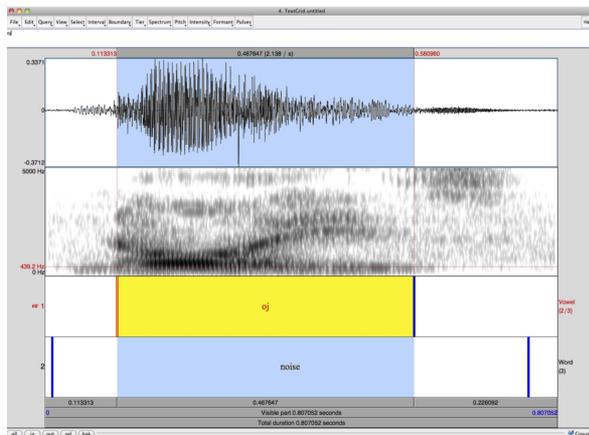
**Slides:** **bit.ly/TRecsICLDC21**

# Background:
# What is linguistic data?



21st century

Jansegers & Gries
2017:10



http://ase.tufts.edu/psychology/psy
cholinglab/asl-lex/visualization.html



Styler
2017:54

Ak'a-*ggem*        ayag-llru-uq

already-INFER leave-PAST-3

"It seems he already left."

Payne 1997:253

# Background: A long-noticed problem

1994: Editor of *Language*, top journal in the field found many cases where use of data was problematic

"...so frequently, in fact, that the assumption that the **data in accepted papers is reliable** began to look questionable"

(Thomason 1994:409)

Exhortation to use data carefully,
Describe and cite sources well,
Say how data was collected.

# Overview of presentation

Background: Citing linguistic data
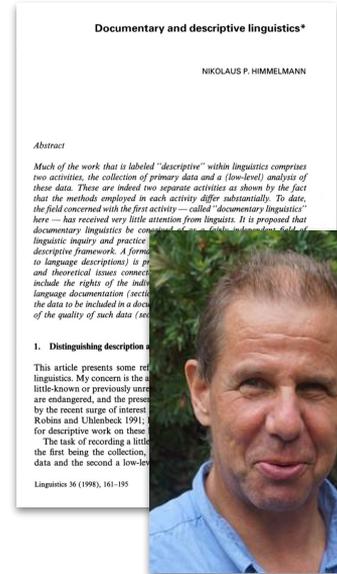
Outline: The main features of the Tromsø Recommendations

Next: Building support for data citation in your work and community

# Background: A long-noticed problem

"It is simply a feature of a scientific enterprise to make one's primary data accessible to further scrutiny"

(Himmelmann 1998: 165)

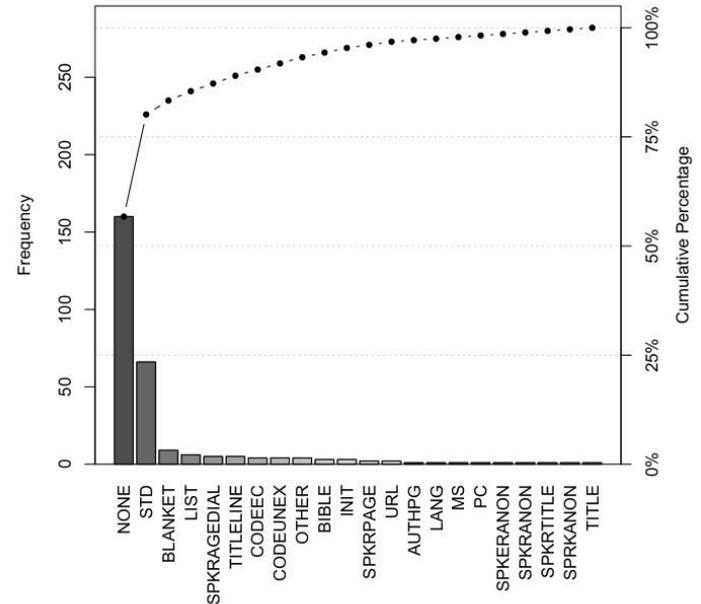See also Gawne & Berez-Kroeker (2018)

# Background: Linguists don't cite data (much)

Data in publications **don't generally have citations**

>       (cf Berez-Kroeker et al. 2017)

If they do, citation only vaguely linked to the actual data set, making reproducible research very hard.



Citation convention frequencies: All journals

# Background: LDIG

**Related LDIG publications**

2018: Open access position paper on reproducibility in linguistics.
Most downloaded article of the journal.

To appear: The Open Handbook of Linguistic Data Management, MIT Press Open (Berez-Kroeker, McDonnell, Koller & Collister, eds.).
13 chapters on conceptual foundations of data management for linguistics and best practices. 50 short data management use cases. Appr. 90 authors from four continents.

**Abstract:** This paper is a position statement on reproducible research in linguistics, including data citation and attribution, that represents the collective views of some 41 colleagues. Reproducibility can play a key role in increasing

*Corresponding author: **Andrea L. Berez-Kroeker,** Department of Linguistics, University of Hawai'i at Mānoa, 1890 East West Road, Moore 569, Honolulu, HI 96822, USA, E-mail: andrea.berez@hawaii.edu
**Lauren Gawne,** Department of Languages and Linguistics, SOAS University of London, London WC1H 0XG, UK; La Trobe University, Melbourne, VIC 3086, Australia, E-mail: l.gawne@latrobe.edu.au
**Susan Smythe Kung,** Archive of the Indigenous Languages of Latin America, University of Texas at Austin, Austin, TX 78712, USA, E-mail: skung@austin.utexas.edu
**Barbara F. Kelly,** Department of Languages and Linguistics, The University of Melbourne, Parkville, VIC 3010, Australia, E-mail: b.kelly@unimelb.edu.au
**Tyler Heston,** Payap University, Chiang Mai 50000, Thailand, E-mail: tylerheston@earthlink.net
**Gary Holton,** Department of Linguistics, University of Hawai'i at Mānoa, 1890 East West Road, Moore 569, Honolulu, HI 96822, USA, E-mail: holton@hawaii.edu
**Peter Pulsifer,** National Snow and Ice Data Center, Boulder, CO 80303, USA, E-mail: pulsifer@nsidc.org
**David I. Beaver,** Department of Linguistics, University of Texas at Austin, Austin, TX 78712, USA, E-mail: dib@utexas.edu
**Shobhana Chelliah,** Department of Linguistics, University of North Texas, Denton, TX 76203,

# Background: LDIG

2017: Research Data Alliance [Linguistic Data Interest Group](#) founded

First publication: [The Austin Principles of Data Citation](#)
Explains the importance of *why* to cite data, but not *how*.

Data is central to empirical linguistic research. Linguistic data comes in many different forms, and is collected and processed with a wide range of methods. Data citation recognizes the centrality of data to research. Furthermore, it facilitates verification of claims and repurposing of data for other studies.

# Background References

- Andreassen, H. N., Berez-Kroeker, A. L., Collister, L., Conzett, P., Cox, C., Smedt, K.D., McDonnell, B. and and the Research Data Alliance Linguistic Data Interest Group. (2019). Tromsø recommendations for citation of research data in linguistics (Version 1). R*esearch Data Alliance.* DOI: 10.15497/RDA00040
- Berez-Kroeker, A.L., H.N. Andreassen, L. Gawne, G. Holton, S. Smythe Kung, P. Pulsifer, L.B. Collister, The Data Citation and Attribution in Linguistics Group, & the Linguistics Data Interest Group. 2018. *The Austin Principles of Data Citation in Linguistics.* Version 1.0. http://site.uit.no/linguisticsdatacitation/austinprinciples/
- Berez-Kroeker, A.L., L. Gawne, B.F. Kelly & T. Heston. 2017. *A survey of current reproducibility practices in linguistics journals, 2003-2012.* link.
- Berez-Kroeker, A.L., L. Gawne, S. Kung, B.F. Kelly, T. Heston, G. Holton, P. Pulsifer, D. Beaver, S. Chelliah, S. Dubinsky, R.P. Meier, N. Thieberger, K. Rice & A. Woodbury. 2018. Reproducible research in linguistics: A position statement on data citation and attribution in our field. *Linguistics* 56(1): 1–18.
- Berez-Kroeker, A.L., B. McDonnell, E. Koller & L.B. Collister. In prep. *The Open Handbook of Linguistic Data Management.* MIT Press Open.
- Gawne, L. & A.L. Berez-Kroeker. 2018. Reflections on reproducible research. In B. McDonnell, A.L. Berez-Kroeker & G. Holton. (Eds.) *Reflections on Language Documentation 20 Years after Himmelmann 1998*, pp. 22–32. University of Hawaiʻi Press.
- Himmelmann, N. P. (1998). Documentary and descriptive linguistics. Linguistics, 36, 161-196.
- Jansegers, M. & S.Th. Gries. 2017. Towards a dynamic behavioral profile: A diachronic study of polysemous *sentir* in Spanish. *Corpus Linguistics and Linguistic Theory* 1-43.
- Payne, T. E. 1997. *Describing Morphosyntax: A Guide for Field Linguists.* Cambridge University Press.
- Styler, W. 2017. *Using Praat for Linguistic Research*, version 1.8.1. Online: http://savethevowels.org/praat. Accessed 20 November 2019.
- Thomason, S. 1994. The editor's department. *Language* 70: 409-413.