



# Language Technology: Research and Development

Dissemination of Research Results

Sara Stymne

Uppsala University  
Department of Linguistics and Philology  
[sara.stymne@lingfil.uu.se](mailto:sara.stymne@lingfil.uu.se)

Partially based on slides from Joakim Nivre

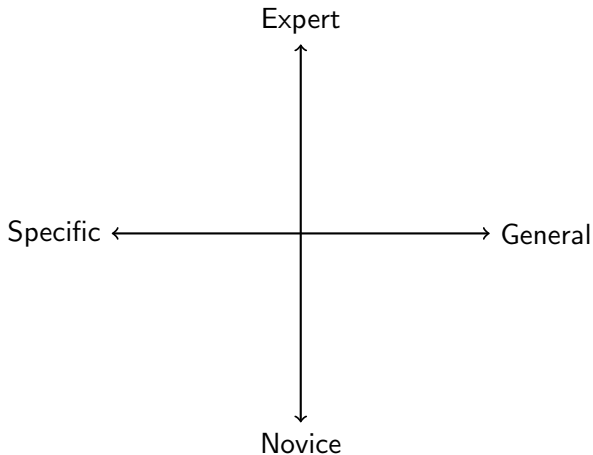


# Dissemination of Research Results

- ▶ Why?
  - ▶ Submit results for critical review
  - ▶ Inform other researchers, users, society
  - ▶ Satisfy requirements from funders or customers
  - ▶ Promote research career – publish or perish
- ▶ To whom?
  - ▶ Other researchers
  - ▶ Potential users
  - ▶ Students
  - ▶ The general public
  - ▶ Funding bodies
  - ▶ Customers

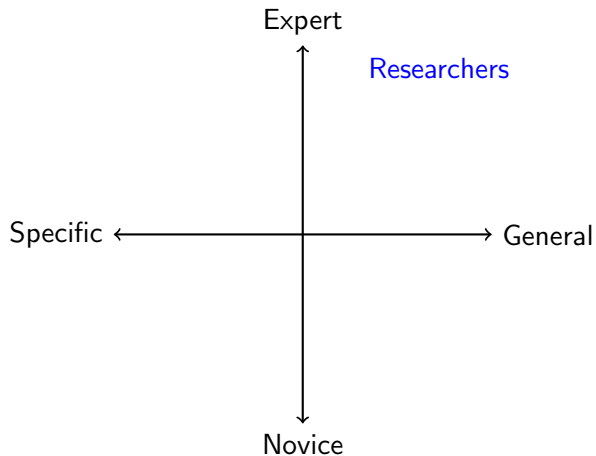


# The Receiver



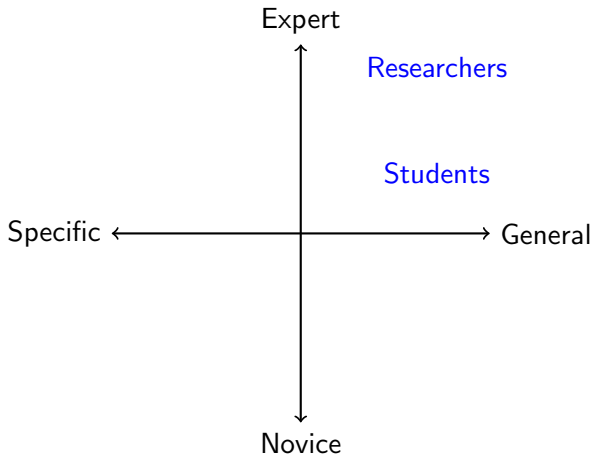


# The Receiver



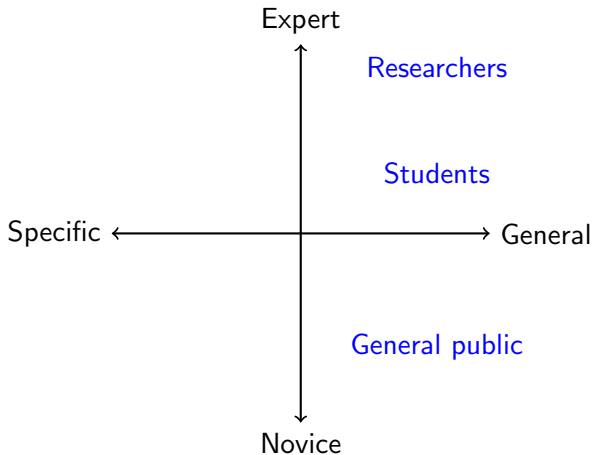


# The Receiver



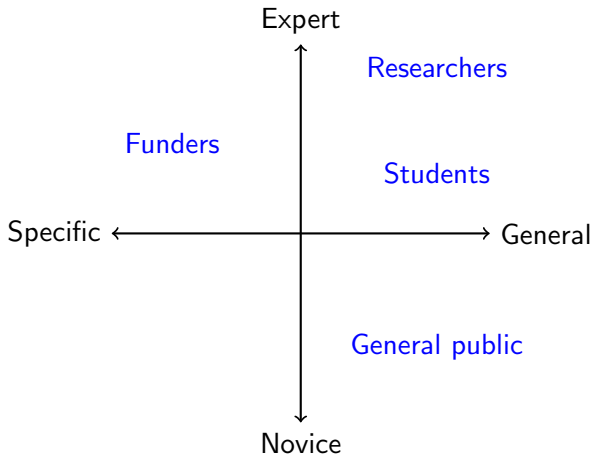


# The Receiver



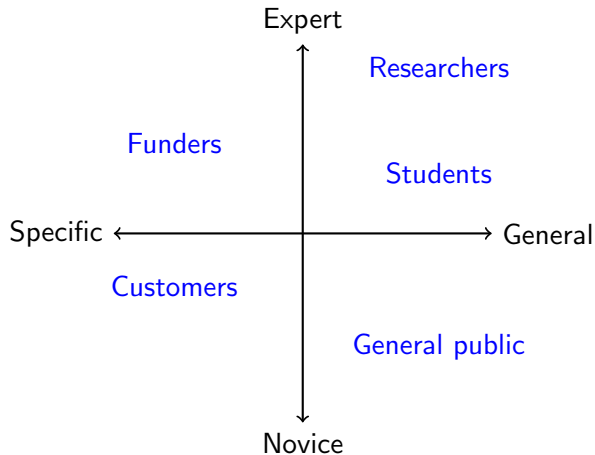


# The Receiver





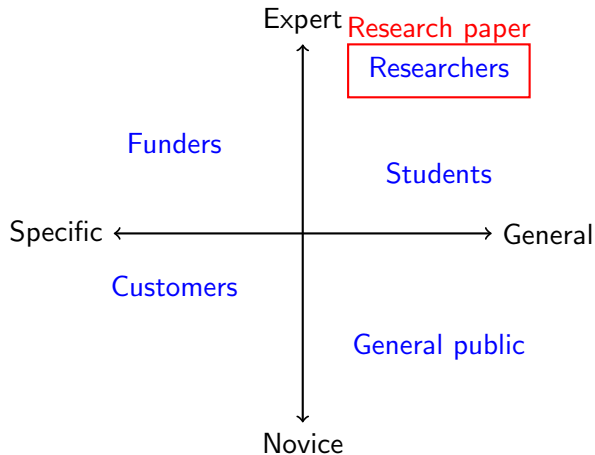
# The Receiver





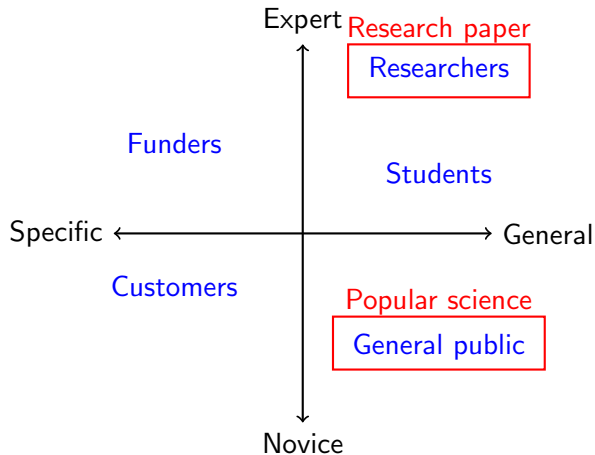


# The Receiver



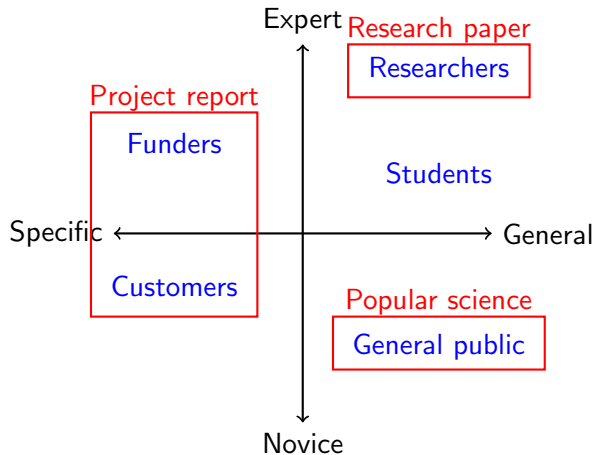


# The Receiver



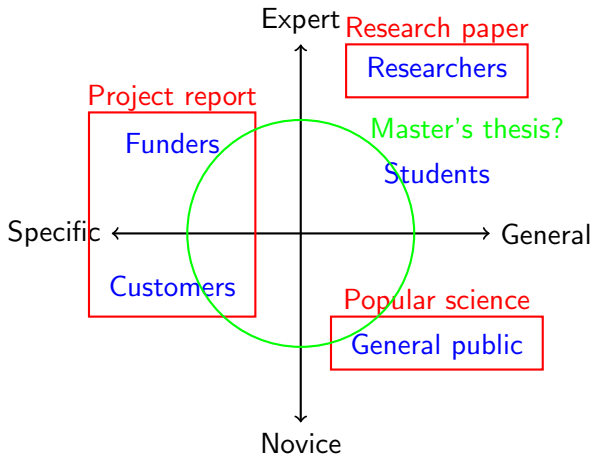


# The Receiver



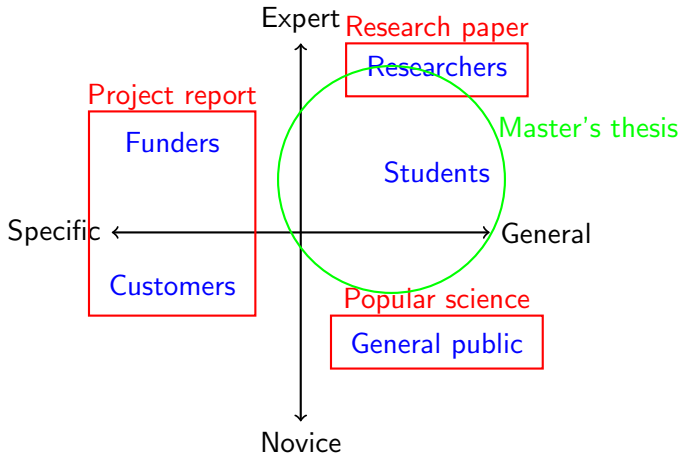


# The Receiver





# The Receiver





# How?

## Written:

1. Publications (indexed and archived)
2. Internal reports (public or confidential)
3. Digital archives, web pages, etc.

## Oral:

1. Lectures (especially at conferences)
2. Demonstrations, posters, discussions, etc.
3. Internal meetings (seminars, workshops)



## Written Genres – Single Topic

### Papers (short)

1. Journal article – refereed and approved by editorial board
2. Conference paper – often but not always refereed
3. Technical report – usually not refereed

### Monographs (long)

1. Book – standards of refereeing depends on publisher
2. Thesis – refereed in examination, may or may not be published



## Written Genres – Other

### Collections

1. Conference proceedings – collection of conference papers
2. Edited volume – book with different chapter authors

### Meta-genres

1. Survey or handbook article
2. Review in scientific journal
3. Bibliography
4. Abstract





# Oral Genres

## Lecture

- ▶ Presentation by 1 person followed by discussion (large group)
  1. Conference talk (15–30 min)
  2. Invited talk (45–90 min)

## Seminar

- ▶ Presentation or introduction by 1 or more persons with more or less continuous discussion (small group)

## Panel

- ▶ Short presentations on a set topic from a selected group of persons with questions and opinions from the audience



# Mixed Genres

## Poster

- ▶ Written presentation displayed on poster board
- ▶ Oral interaction with interested audience
- ▶ Sometimes combined with short talk (1–5 min)

## Demonstration

- ▶ System demonstration (or similar)
- ▶ Oral interaction with interested audience
- ▶ Sometimes combined with poster



# Requirements on Scientific Reports

- ▶ Ethics:
  - ▶ Sensitive information requires permission and anonymization
- ▶ Accessibility:
  - ▶ Reports should be understandable by target audience
- ▶ Novelty and relevance:
  - ▶ Results should be novel, original, unpublished
  - ▶ Relevance to research area should be made clear
- ▶ Quality:
  - ▶ Claims clearly stated and possible to challenge (falsifiability)
  - ▶ Claims supported by arguments and/or evidence (justification)
  - ▶ Claims not misleading (e.g., by withholding information)



# Scientific Writing

Writing takes time (to learn)

- ▶ Practice makes perfect – write a lot!
- ▶ Writing requires rewriting – start early!

Scientific writing is a standardized genre

- ▶ Collect good examples – and study them!
- ▶ Copy structure and formulations – but not content!



# The Structure of Scientific Publications



# The Structure of Scientific Publications

**Pre-matter:** Title page (abstract, preface, contents)

**Post-matter:** References (appendices, indexes)



# The Structure of Scientific Publications

**Pre-matter:** Title page (abstract, preface, contents)

**Introduction:** What is the problem/question?  
Why is it relevant/interesting?

**Conclusion:** What is the solution/answer?  
Where do we go from here?

**Post-matter:** References (appendices, indexes)



# The Structure of Scientific Publications

**Pre-matter:** Title page (abstract, preface, contents)

**Introduction:** What is the problem/question?  
Why is it relevant/interesting?

**Body:** What has been done before?  
How is the problem tackled?  
What are the results?

**Conclusion:** What is the solution/answer?  
Where do we go from here?

**Post-matter:** References (appendices, indexes)





# The Main Theme

The research question

- ▶ is stated in the introduction
- ▶ is related to previous research
- ▶ motivates the approach taken
- ▶ determines the selection of results
- ▶ is revisited in the conclusion



# The Anatomy of a TACL Style Article

**Token and Type Constraints for Cross-Lingual Part-of-Speech Tagging**  
 Oscar Täckström<sup>1,2\*</sup> Dipanjan Das<sup>1</sup> Slav Petrov<sup>1</sup> Ryan McDonald<sup>1</sup> Jaakko Niivre<sup>2</sup>

<sup>1</sup>Swedish Institute of Computer Science  
<sup>2</sup>Department of Linguistics and Philology, Uppsala University  
 Google Research, New York

oscar@aic.a.se  
 {dipan.jand|slav|ryan.mcd}@google.com  
 jaakko.nivre@lingfil.uu.se

**Abstract**

We consider the construction of part-of-speech taggers for resource-poor languages. Recently, manually constructed tag dictionaries from Wiktionary and dictionaries proposed via biword have been used as type constraints to increase the accuracy of automatic data in this setting. In this paper, we show that additional noise constraints can be projected from a resource-rich source language to a resource-poor target language via word-aligned biword. We present several models in this regard, in particular a partially observed conditional random field (model), where capital tokens and type tags are used as a special signal for training. We compare across eight previously studied Indo-European languages, our model achieves a 27% relative error reduction over the prior state of the art. We further present successful results on seven additional languages from different families, empirically demonstrating the applicability of capital tokens and type constraints across a diverse set of languages.

**1 Introduction**

Supervised part-of-speech (POS) taggers are available for more than twenty languages and achieve accuracies of around 95% on in-domain data (Petrov et al., 2012). Thanks to their efficiency and robustness, supervised taggers are routinely employed in many natural language processing applications, such as syntactic and semantic parsing, named-entity recognition and machine translation. Unfortunately, the resources required to train supervised taggers are expensive to create and unlikely to exist for the majority of written languages. The necessity of building NLP tools for these resource-poor languages has been part of the motivation for research on unsupervised learning of POS taggers (Chenoukova et al., 2010).

In this paper, we instead take a weakly supervised approach towards this problem. Recently, learning POS taggers with type-level tag dictionaries constraints has gained popularity. Tag dictionaries, mostly projected via word-aligned biword, have bridged the gap between poorly unsupervised and fully supervised taggers, resulting in an average accuracy of over 85% on a benchmark of eight Indo-European languages (Das and Petrov, 2013; Li et al., 2012). Further improvement upon this result by employing Wiktionary<sup>1</sup> as a tag dictionary source, resulting in the hitherto best published result of almost 95% on the same setup.

Although the aforementioned weakly supervised approaches have resulted in significant improvements over fully unsupervised approaches, they have not exploited the benefits of state-of-the-art cross-lingual projection methods, which are possible with word-aligned biword between a target language of interest and a resource-rich source language, such as English. This is the setting we consider in this paper (32). While prior work has traditionally considered both token- and type-level projections across word-aligned biword for estimating the model parameters of generative tagging models (Chenoukova and Nivre, 2010; Xu and Hwa, 2005, inter alia), a key observation underlying the present work is that token- and type-level information offer different and complementary signals. On the one hand, high confidence token-level projections offer precise constraints on a tag in a particular context. On the other hand, manually con-

\*Work primarily carried out while at Google Research.

<sup>1</sup><https://www.wiktionary.org/>.

3

Title page: title, authors, affiliations

Abstract: self-contained summary

Main text in numbered sections



# The Anatomy of a TACL Style Article

## 6 Conclusions

We considered the problem of constructing multilingual POS taggers for resource-poor languages. To this end, we explored a number of different models that combine token constraints with type constraints from different sources. The best results were obtained with a partially observed CRF model that effectively integrates these complementary constraints. In an extensive empirical study, we showed that this approach substantially improves on the state of the art in this context. Our best model significantly outperformed the second best model on 10 out of 15 evaluated languages, when tested on identical data sets, with an insignificant difference on 3 languages. Compared to the prior state of the art (Li et al., 2012), we observed a relative reduction in error by 25%, averaged over the eight languages common to our studies.

## Acknowledgments

We thank Alexander Bush for help with the hypergraph framework that was used to implement our models and Klaus Macherey for help with the bi-text extraction. This work benefited from many discussions with Yves Goldberg, Keith Hall, Rasmus Geisler and Hao Zhong. We also thank the editor and the three anonymous reviewers for their valuable feedback. The first author is grateful for the financial support from the Swedish National Graduate School of Language Technology (GSLT).

## References

Azou Aheille, Lionel Clement, and Françoise Tissotani. 2003. Building a Truthbank for French. In A. Aheille, editor, *Textbank: Building and Using Parallel Corpora*, chapter 10. Kluwer.

Taylor Berg-Kirkpatrick, Alexander Buchsart-Ciik, John Dieffen, and Dan Klein. 2010. Patches unpermitted: learning with features. In *Proceedings of NAACL-HLT*. Salt Lake, Utah: Morgan Kaufmann.

Siddhartha Dandekar and Ervin Shadmehr. 2008. CoNLL-X shared task on multilingual dependency parsing. In *Proceedings of CoNLL*.

Stanley F Chen. 2003. Conditional and joint models for grapheme-to-phoneme conversion. In *Proceedings of Eurospeech*.

Chiranjit Choudhury, Sharan Goelwater, and Mark Steinhilber. 2010. Two decades of unpermitted POS induction: How far have we come? In *Proceedings of EMNLP*.

Dipayan Das and Slav Petrov. 2011. Unsupervised part-of-speech tagging with bilingual graph-based projections. In *Proceedings of ACL-IJLW*.

Arthur P Dempster, Nan M Laird, and Donald B Rubin. 1977. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society Series B*, 39.

John DeNero and Klaus Macherey. 2011. Model-based aligner combination using dual decomposition. In *Proceedings of ACL-IJLW*.

Bruce Edlin and Robert J. Tibshirani. 1993. *An Introduction to the Bootstrap*. Chapman & Hall, New York, NY, USA.

Victoria Papan and Steven Abney. 2005. Asymmetrically inducing a part-of-speech tagger by projecting from multiple source languages across aligned corpora. In *Proceedings of ACL*.

Das Garenna and Jason Baldridge. 2012. Type-supervised token marker models for part-of-speech tagging with incomplete tag dictionaries. In *Proceedings of EMNLP-CoNLL*.

Yves Goldberg, Shari Adin, and Michael Elhadad. 2008. EM can find pretty good HMM POS-tagger values given a good start. In *Proceedings of ACL-IJLW*.

Philippe Korien. 2005. *Empirical: A parallel corpus for statistical machine translation*. In *MT Summit*.

John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of ICML*.

Shen Li, John Ong, and Ben Taskar. 2012. Wiki-style supervised part-of-speech tagging. In *Proceedings of EMNLP-CoNLL*.

Dong C. Liu and Jorge Nocedal. 1989. On the limited memory BFGS method for large scale optimization. *Mathematical Programming*, 45.

Michael J. Marcus, Mary Ann MacChieveze, and Bruce Santner. 1995. Building a large annotated corpus of English: the Penn treebank. *Computational Linguistics*, 19(2).

Tabita Nauwer, Benjamin Seyfar, Jacob Eisenstein, and Regina Barzilay. 2009. Multilingual part-of-speech tagging: Two unsupervised approaches. *AACL*, 36.

Jackie Niwe, Johan Hall, Sander Kiffin, Ryan McDonald, Henk Nilsson, Sebastian Riedel, and Denis Yuret. 2007. The CoNLL 2007 shared task on dependency parsing. In *Proceedings of EMNLP-CoNLL*.

Slav Petrov, Dipayan Das, and Ryan McDonald. 2012. A universal part-of-speech tagset. In *Proceedings of LREC*.

Sajid Ravi and Kevin Knight. 2008. Minimalist models for supervised part-of-speech tagging. In *Proceedings of ACL-IJLW*.

Main text in numbered sections

Acknowledgments (optional)

References (alphabetical by last name)



# The Anatomy of a TAACL Style Article

## Introduction

- ▶ State the research problem and relate it to previous research
- ▶ Give a synopsis of the rest of the article

## Related work

- ▶ Model 1: After introduction, before contributions
- ▶ Model 2: After contributions, before conclusion

## Contributions

- ▶ Theory → Method → Results → Discussion

## Conclusion

- ▶ Evaluate contributions, point to new research directions



## Supplementary materials

It has become increasingly common to provide additional materials along with a research paper

- ▶ Appendix – additional materials that does not fit in the main paper, e.g.:
  - ▶ Surveys
  - ▶ Guidelines
  - ▶ Sample of data sets
  - ▶ Proofs
  - ▶ Long tables (not containing key material)
  - ▶ Technical details such as hyper-parameters
- ▶ Code
- ▶ Data sets



## Supplementary materials

It has become increasingly common to provide additional materials along with a research paper

- ▶ Appendix – additional materials that does not fit in the main paper, e.g.:
  - ▶ Surveys
  - ▶ Guidelines
  - ▶ Sample of data sets
  - ▶ Proofs
  - ▶ Long tables (not containing key material)
  - ▶ Technical details such as hyper-parameters
- ▶ Code
- ▶ Data sets

We do not expect you to have appendices in your term paper, but we welcome links to your code and data (in case you performed some annotation)



## References

- ▶ Language technology mostly uses the Harvard system
  - ▶ Author-year citations in text
  - ▶ Alphabetical list of references at the end (no footnotes)
- ▶ Citations in the text:
  - ▶ Parenthetical: Translation is hard (Smith, 2012).
  - ▶ Syntactic: Smith (2012) claims that translation is hard.
  - ▶ More than two authors:
    - ▶ In text, use et al.  
Parsing is hard (Anderson et al., 2010).  
Anderson et al. (2010) claims that parsing is hard.
    - ▶ All authors in reference list  
Anderson, P., Svensson, G, Lind, W. and Sund, T. 2017.  
Parsing is hard. . . .



## Reference List

- ▶ Reference list including all (and only) works cited in the text:
  - ▶ **Journal article:** author, year, title, *journal*, volume, number, pages
  - ▶ **Conference paper:** author, year, title, *proceedings*, pages, location
  - ▶ **Book chapter:** author, year, title, *book*, editors, publisher, pages
  - ▶ **Book:** author, year, *title*, publisher
  - ▶ **Technical report:** author, year, title, organization
  - ▶ **Thesis:** author, year, title, type of thesis, school
- ▶ Important: BE CONSISTENT!





## Bibtex example – journal article

```
@article{songetal2019semantic,  
  title = "Semantic Neural Machine Translation Using {AMR}",  
  author = "Song, Linfeng and Gildea, Daniel and Zhang, Yue  
    and Wang, Zhiguo and Su, Jinsong",  
  journal = "Transactions of the Association for Computational  
    Linguistics",  
  volume = "7",  
  year = "2019",  
  pages = "19–31",  
}
```



## Bibtex example – conference article

```
@inproceedings{rahimietal2019massively,  
  title = "Massively Multilingual Transfer for {NER}",  
  author = "Rahimi, Afshin and Li, Yuan and Cohn, Trevor",  
  booktitle = "Proceedings of the 57th Annual Meeting of  
    the Association for Computational Linguistics",  
  year = "2019",  
  address = "Florence, Italy",  
  pages = "151–164",  
}
```



## Bibtex example – arXiv article

```
@misc{deWynterPerryOptimal,  
  title="Optimal Subarchitecture Extraction for {BERT}",  
  author="Adrian de Wynter and Daniel J. Perry",  
  year=2020,  
  howPublished = "\it arXiv preprint arXiv:2010.10499v1",  
}
```



## Bibtex example – arXiv article

```
@misc{deWynterPerryOptimal,  
  title="Optimal Subarchitecture Extraction for {BERT}",  
  author="Adrian de Wynter and Daniel J. Perry",  
  year=2020,  
  howPublished = "\it arXiv preprint arXiv:2010.10499v1",  
}
```

Note: do NOT cite an arXiv article if there is a published version of it! Double check each time you have an arXiv article in your reference list, since many of them eventually get published!



## Bibtex example – book

```
@Book{MS99statmet,  
  author = {Christopher D. Manning and Hinrich Sch\"utz},  
  title = {Foundations of Statistical Natural Language  
    Processing},  
  publisher = {MIT Press},  
  year = 1999,  
  address = {Cambridge, Massachusetts, USA}  
}
```



## Bibtex example – book chapter

```
@InCollection{Lude11corpus,  
  author = {Anke L\"udeling},  
  title = {Corpora in Linguistics: Sampling and Annotations},  
  booktitle = {Going Digital, Evolutionary and Revolutionary  
    Aspects of Digitization},  
  pages = {220–243},  
  publisher = {The Nobel Foundation},  
  year = 2011,  
  editor = {Karl Grandin},  
}
```



## Using bibtex bibliography

```
%style file  
\bibliographystyle{tacl2018}
```

```
%Name of your bibtex file:  
\bibliography{myRefs.bib}
```



## Presenting non-English examples

- ▶ Often useful to show language examples
- ▶ You cannot expect all readers to know all languages
- ▶ Help the reader understand your key point(s) with the example!
- ▶ Standard presentation
  - ▶ Foreign example
  - ▶ Word-by-word gloss
  - ▶ 'Translation'





## Swedish example

- (1) I förrgår gick pennan av  
In past yesterday went the pen off  
'The day before yesterday, the pen broke'

Or more detailed:

- (2) I förrgår gick pennan av  
In past yesterday go.PAST pen.DEF.SG off  
'The day before yesterday, the pen broke'



## More on non-English examples

- ▶ Think about why you include the example
  - ▶ What does the reader need to know in order to understand the point?
- ▶ Depending on this answer, the amount of detail can vary
  - ▶ Translation only
  - ▶ Translation and simple gloss
  - ▶ Translation and detailed gloss
- ▶ For detailed glosses: Leipzig Glossing Rules
- ▶ Latex packages: Covington, gb4e, linguex, . . .



## Ethics: plagiarism

- ▶ Always attribute references to ideas of others!
- ▶ Quotations
  - ▶ Use quotations marks, give reference and page number
  - ▶ Do not overuse. Mainly useful for definitions, et.c.
- ▶ Paraphrasing
  - ▶ Describing the work of others in your own words. Give reference
  - ▶ Changing a few words, tense, et.c. is not enough.
  - ▶ Tip: do not look at the paper you are paraphrasing, write from memory (and then double check)
- ▶ Images:
  - ▶ Often under copyright. If so: DO NOT COPY!
  - ▶ If permissive license: you can copy, and give reference
  - ▶ Otherwise: draw your own variant, and give reference



# Ethics: data manipulation

- ▶ Fabrication
  - ▶ Making up false results
  - ▶ Manipulating experiments
- ▶ Cherry picking / suppressing evidence
  - ▶ Only presenting results that supports your hypothesis
  - ▶ Setting up the study so that the experiments are not representative
  - ▶ Not (attempting to) controlling for confounding variables

DO NOT DO THIS!



# Giving Oral Presentations

Preparation is the key

- ▶ Think through what you want to say
- ▶ Formulate key passages in concrete sentences
- ▶ Prepare audiovisual aids (if relevant)

Practice makes perfect

- ▶ Rehearse the presentation (many times)
- ▶ Time the presentation and note any disfluencies
- ▶ Modify and rehearse until fluent



## The Structure of Oral Presentations

Oral presentations are basically structured as written reports but

- ▶ typically contain less material due to time constraints (especially the background part)
- ▶ are often less formal and detailed due to real-time processing (the big picture instead of the formal details)
- ▶ can be more repetitive due to memory limitations (get the take-home message across)

The discussion part:

- ▶ Listen to the question
- ▶ Answer the question – if you can



## Audiovisual Aids

Slides provide support for the presentation

- ▶ Key points and important concepts
- ▶ Graphical illustrations (and sound if relevant)
- ▶ Material that is hard to present orally (equations, examples)

But remember

- ▶ Not too much information (or too small fontsize) on one slide
- ▶ Not running text (to be read aloud)
- ▶ Slides should support presentation, not vice versa



## Geoff Pullum's Golden Rules



- ▶ Don't ever begin with an apology
- ▶ Don't ever underestimate the audience's intelligence
- ▶ Respect the time limits
- ▶ Don't survey the whole damn field
- ▶ Remember that you're an advocate, not the defendant
- ▶ Expect questions that will floor you





## Requirements for your course papers

- ▶ Follow the TACL guidelines
- ▶ Use the TACL Latex templates
- ▶ 4–7 pages of content + references
- ▶ Content is text + tables + figures
  - ▶ Only references and acknowledgement allowed on additional pages



## Deadlines and submissions

- ▶ December 13: First version of **full paper**
- ▶ January 13: final version of paper, taking reviews into account
- ▶ If you miss a deadline (AVOID!)
  - ▶ First version: January 13
  - ▶ Second version: February 17
  - ▶ You may present during the workshop, if your project is ready to present (Jan 12)
- ▶ Reviewing and first version will be handled through a conference management system – more information will come!



## Final seminar

- ▶ January 12: all day
- ▶ Will be held on Campus
- ▶ Online presentations will only be allowed if there are exceptional circumstances (for instance travel bans). If this is the case for you, you need to contact Sara in advance, and get approval.
- ▶ There will probably be two types of sessions
  - ▶ Plenary
  - ▶ Half class



## Next group seminar

- ▶ Special theme: ethics
- ▶ Reading:
  - ▶ Main reading: Hovy, D. and Spruit, S. L. *The Social Impact of Natural Language Processing*. ACL 2016.
  - ▶ Additional reading: Bender, E. M., Gebru, T., McMillan-Major, A. and Shmitchell, S. *On the dangers of stochastic parrots. Can language models be too big?* FAccT'21
- ▶ Reflect on ethical issues related to your projects
- ▶ Also feel free to think about positive consequences of your project!



## Coming up

- ▶ Latex tutorial, 16/11
  - ▶ Mainly help with any Latex issues you may have
  - ▶ Some exercises for practice
  - ▶ Only on Campus
- ▶ For those of you who have not yet passed the proposal, remember to hand in the main proposal, and the popular science abstract by November 4!
  - ▶ Make sure that you follow the instructions (template, length, content)
- ▶ Second take-home exam for those who did not pass the first exam, and signed up for it:
  - ▶ Handed out: November 4, 08:00 in Studium
  - ▶ Deadline: November 11, 23:59