



# 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)

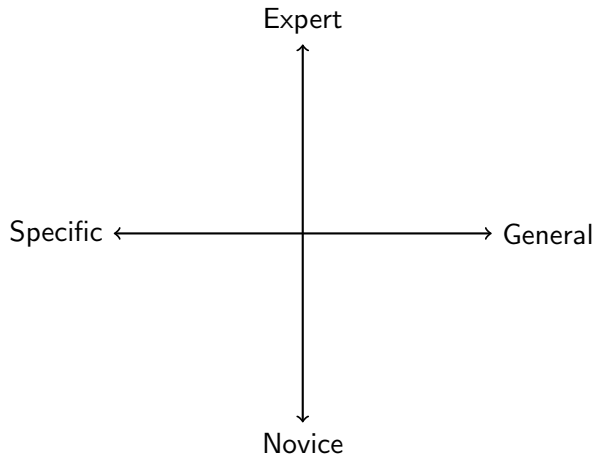


# Dissemination of Research Results

- ▶ 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

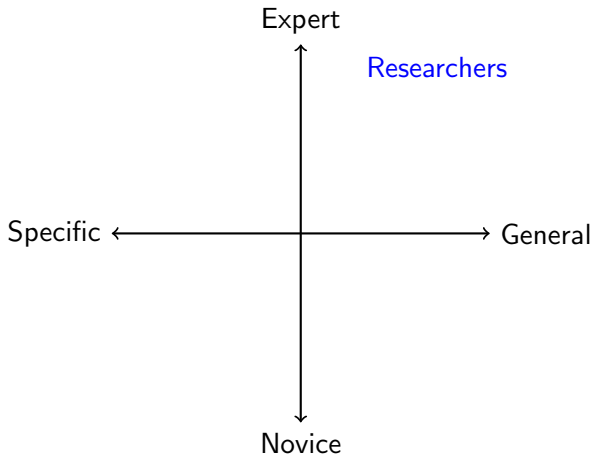


# The Receiver



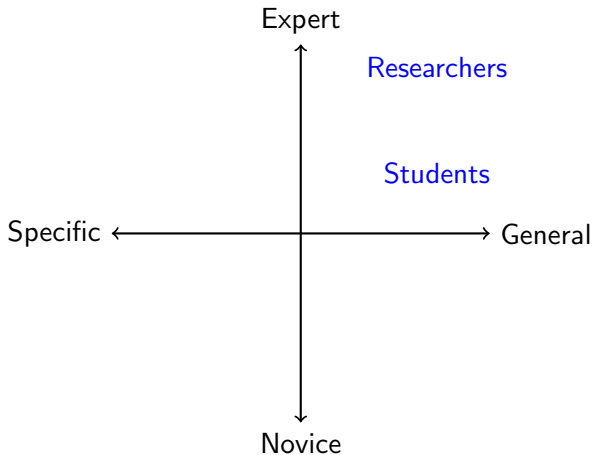


# The Receiver



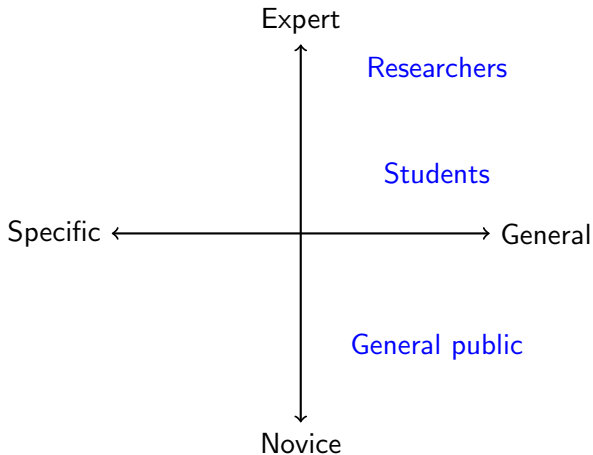


# The Receiver



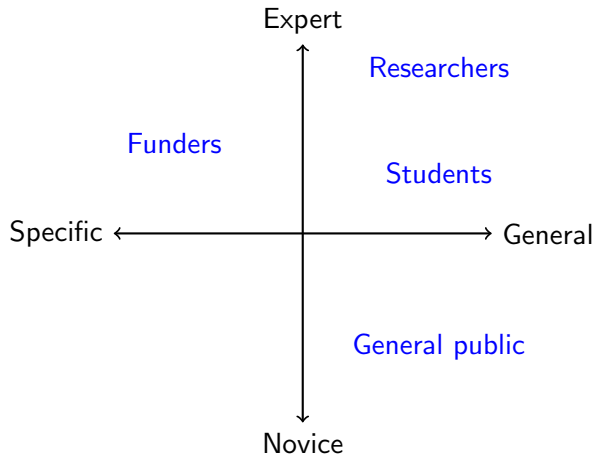


# The Receiver



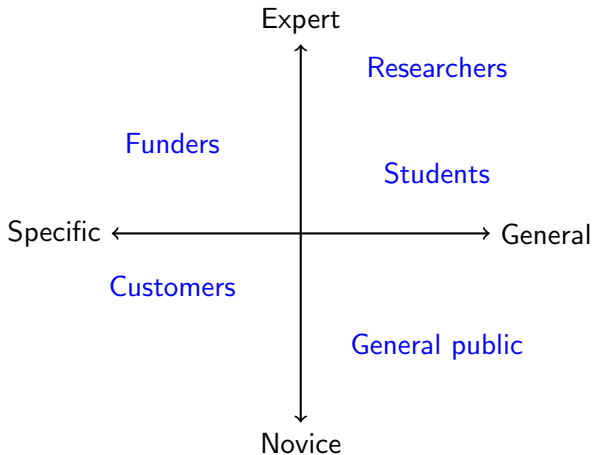


# The Receiver





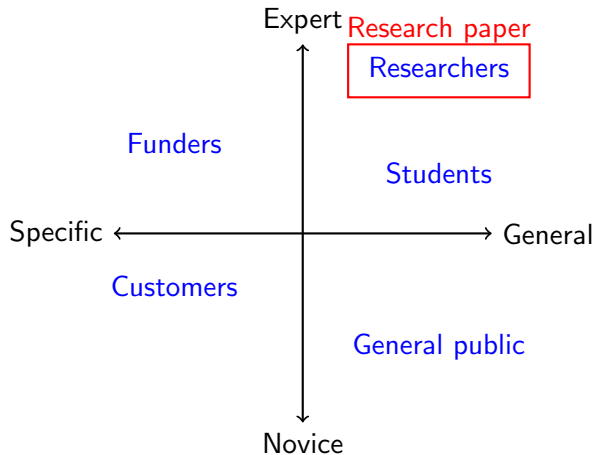
# The Receiver





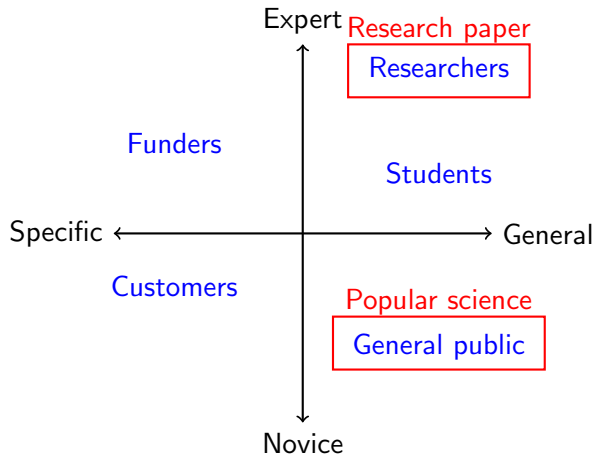


# The Receiver



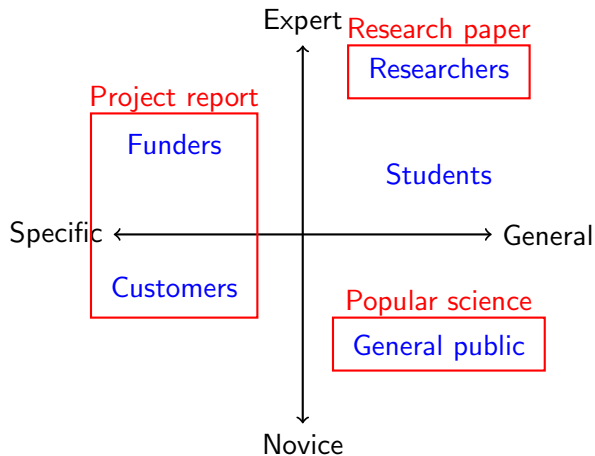


# The Receiver



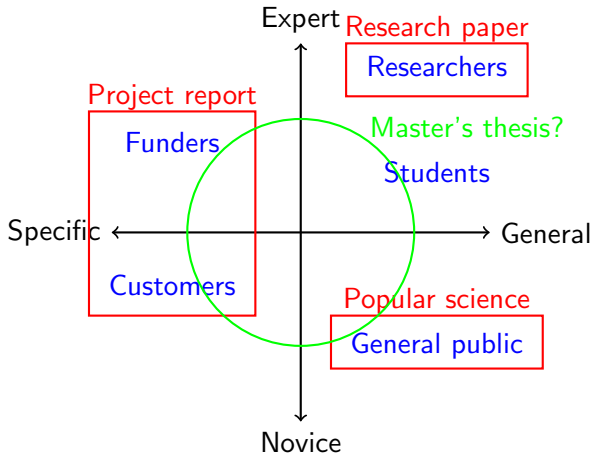


# The Receiver



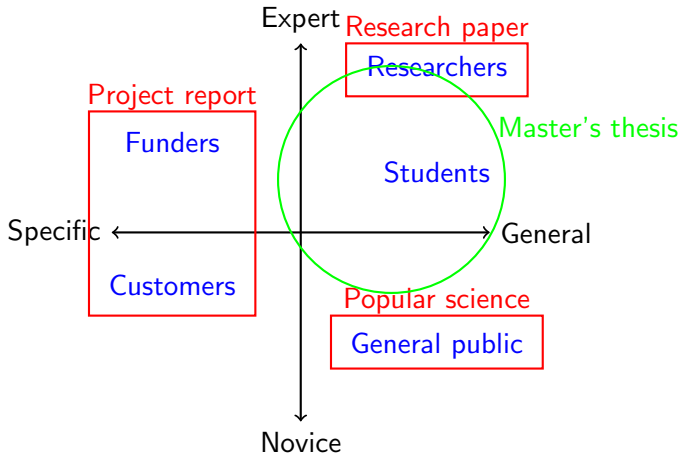


# The Receiver





# 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

## 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 Garschev 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

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

Taylor Berg-Karantatzis, Alexandros Bouras-COx, John Doehren, and Dan Klein. 2010. Patches unperceivable: learning with features. In *Proceedings of NAACL-HLT*. Suttons Brothers and Erwin Mann.

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 Chakraborty, Sharan Goelwater, and Mark Steinhilber. 2010. Two decades of unperceivable 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-IJL*.

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-IJL*.

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 hidden marker models for part-of-speech tagging with incomplete tag dictionaries. In *Proceedings of EMNLP-CoNLL*.

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

Philippe Kories. 2005. *Samuel: 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-based 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 Srinivasan. 1995. Building a large annotated corpus of English: the Penn treebank. *Computational Linguistics*, 19(2).

Tabita Nauwer, Benjamin Seyfar, Jacob Elmanita, 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 Razi and Kevin Knight. 2008. Minimalist models for supervised part-of-speech tagging. In *Proceedings of ACL-IJL*.

Main text in numbered sections

Acknowledgments (optional)

References (alphabetical by last name)



# The Anatomy of a TACL 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



## 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{konratyukstraka19udify,  
  title="75 Languages, 1 Model: Parsing Universal Dependencies  
    Universally",  
  author="Dan Kondratyuk and Milan Straka",  
  year=2019,  
  note = "{\it arXiv preprint arXiv:1904.02099v3}",  
}
```

Note: do NOT cite arXiv article if there is a published version of it!



## 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}
```



# 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 17: final version of paper, taking reviews into account
- ▶ If you miss a deadline (AVOID!)
  - ▶ First version: January 17
  - ▶ Second version: February 21
  - ▶ You may present during the workshop (Jan 15)
- ▶ Reviewing and first version should be handed in via EasyChair  
- more information will come!