

About this book

What?

This textbook is a graduate-level introduction to proof-based (rather than modeling-based) computational linguistics for linguists. The focus is on how a computationally informed perspective can illuminate aspects of language that theoretical linguists and cognitive scientists care about. At concerted effort has been made to keep the material as approachable as possible without sacrificing exactness. To get the most out of the book, readers have to be willing to engage with mathematical notation, but each unit is designed in a modular fashion so that the less mathematically inclined can skip the parts they find too tedious.

Why?

Like most textbooks, this one was written because of a picky instructor who wasn't quite happy with the existing options. It grew out of my lecture notes for *Computational Linguistics 2* at Stony Brook University's Department of Linguistics. Within linguistics, computational linguistics tends to be geared towards model-building and simulations: MaxEnt learners, corpus-based techniques, and so on. These topics are covered in a different course at Stony Brook, though, with Computational Linguistics 2 exploring the proof-based side of the field: formal language theory, subregular complexity, logic, string and tree automata/transducers, learnability, parsing theory, algebra, plus some algorithms and data structures. There is a number of excellent textbooks that cover a subset of these topics, but they all fell short in some respect:

- Marcus Kracht's *The Mathematics of Language* is a treasure trove and still one of the most rewarding textbooks ever written, but it is too hard and theoretical for the average linguistics student.
- *Speech and Language Processing* by Jurafsky and Martin is an excellent reference book, but it is clearly aimed at computer science students. The focus is on engineering techniques rather than investigating linguistic questions from a computational perspective.
- András Kornai's *Mathematical Linguistics* presents a more appropriate compromise between linguistic inquiry and applications, but is still too far removed from linguistics and cognition for my taste.
- *Mathematical Methods in Linguistics* by Partee, ter Meulen & Wall is a classic, but as the title implies it mostly focuses on mathematics, in particular those areas that are needed for semantics (set theory, algebra, and predicate logic). Its final

chapter on formal language theory is a commendable inclusion, but does not make a strong case for why linguists should care about this perspective.

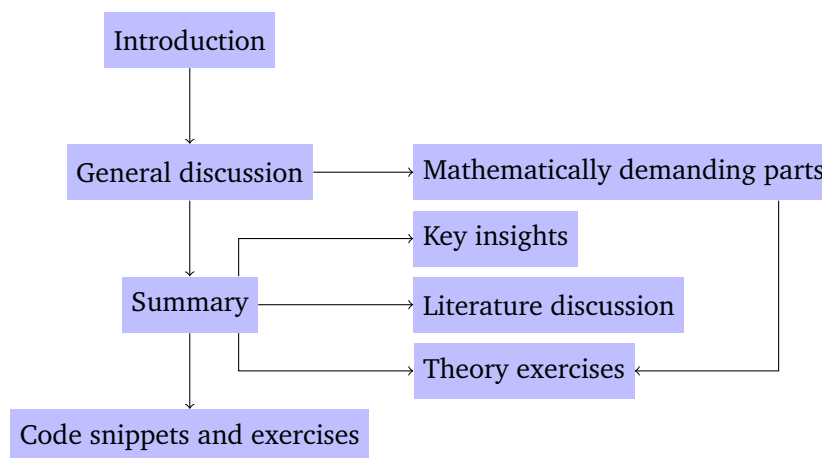
- Ed Keenan and Larry Moss have produced an impressive and very comprehensive introduction to mathematical linguistics with *Mathematical Structures in Language* (and I'm not just saying that because of the many fond memories I have of TAing for the course that their book sprang from). But overall it is a mathematical methods book, not an introduction to proof-based computational linguistics as a means to study language.
- Some subtopics have dedicated textbooks, e.g. Laura Kallmeyer's *Parsing Beyond Context-Free Grammars*, *Categorial Grammar* by Glyn Morrill, *The Logic of Categorial Grammars* by Richard Moot and Christian Retoré, or my colleague Jeff Heinz's book *Grammatical Inference for Computational Linguistics*, co-authored with Colin de La Higuera and Menno van Zaanen. I have happily used them in special topics courses, but they proved difficult to sample from for a broad introduction course like Computational Linguistics 2.

I am sure I have forgotten a few others, all of them deserving of a shout-out here. But the bottom-line is that none fit my vision of proof-based computational linguistics as a way of studying language. So here we are.

How?

Computational linguistics is a very interdisciplinary field, and this is also reflected in the students that have taken the course over the years: from advanced Ph.D. students in linguistics to philosophy and psychology grad students as well as computer scientists with no background in linguistics at all. It is impossible to serve all of them equally well, so I decided to keep the focus on what would benefit and interest linguists the most. But I did make minor additions to include other groups — that is why the linguistically trained reader will sometimes encounter brief sections on very basic material they obviously know about already. I hope they won't feel patronized.

Even the group of linguistics students is far from homogeneous, though, in particular at Stony Brook. Depending on which classes they have taken before, their mathematics and programming background may range from non-existent to far above average. For this reason the textbook has a partly non-linear design:



Every reader should be able to complete the main path from the introduction to the summary. While some math is unavoidable even in those general sections, the more demanding parts are put in a separate section. All the math that is required to follow the main discussion is explained in separate background boxes:

Background

These boxes occur throughout the text and introduce basic concepts of mathematics and computer science.

Exercises are included without solution so that they can be reused by instructors in their course. Each exercise is assigned one of four difficulty levels, indicated by asterisks:

[**Exercise 0.1** The lowest difficulty. Every student should be able to solve this by simply applying a key technique of the unit in a mechanical fashion.]

[**Exercise 0.2*** Requires some effort and creativity. The student cannot simply follow a recipe from the unit but has to build on their understanding of the material to come up with a solution of their own.]

[**Exercise 0.3**** Very difficult for students, and instructors might have to sit down for a minute or two to figure out the general strategy. Writing up the full solution might take quite a while, and/or the solution involves a lot of lateral, out-of-the-box thinking.]

[**Exercise 0.4***** Open research problem. If somebody figures out the answer, they should write a paper about it.]

The book also includes optional discussions of programming issues. Make not mistake, though, this is not a programming textbook. Nonetheless it is often instructive to analyze how exactly formal ideas can be translated into concrete code. I have opted not to put these code snippets directly in the text as they may distract those readers who want to focus exclusively on the interplay of linguistics and abstract computation. Instead, the relevant code snippets are discussed at the very end of the respective unit, after the summary.

A few more remarks on the coding sections are in order. I decided to go with Python instead of pseudo code. Pseudo code is difficult to follow for the uninitiated. Python, on the other hand, enjoys rapidly growing popularity in linguistics and the sciences. Even a reader with a general programming background but no prior experience with Python in particular should be able to understand the code snippets.

To further increase the approachability of the code sections, I have adopted a beginner-friendly coding style. I err very much on the verbose side with line comments and docstrings (which document the purpose and overall design of each function). I also make a deliberate effort to avoid advanced techniques such as decorators and generators. Error handling, which is very important for production-level software, is completely omitted. Similarly, efficiency is considered only to the extent that it illustrates a key pedagogical point. Quite simply, the code sections are of little interest on their own and are designed only to support the concepts covered in each unit.

Depending on which parts of a unit one skips, the book can range from a light-weight introduction that conveys the key intuitions to a very formal *tour de force*.