

# Research Methods

## Transcriptions & Annotations (ELAN Tool)

---

Amandine Decker (*[amandine.decker@gu.se](mailto:amandine.decker@gu.se)*)



UNIVERSITY OF  
GOTHENBURG

# Table of contents

## 1. Transcriptions and Annotations

Examples of Transcription / Annotation Uses

Definition(s)

What is an Annotation Campaign?

## 2. One Transcription / Annotation Tool: ELAN

ELAN?

Examples

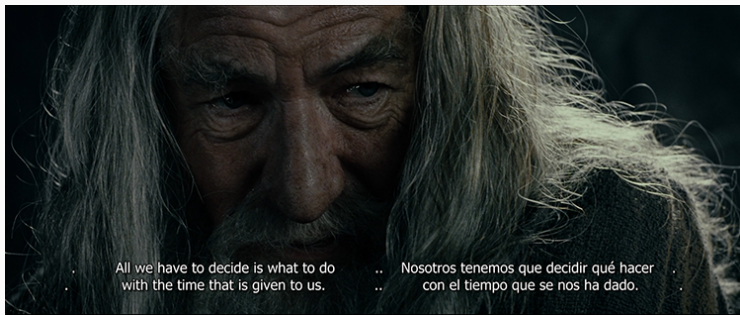
Core Concept

Application and Assignment

# Transcriptions and Annotations

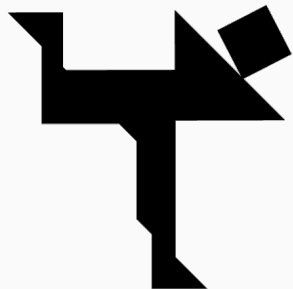
---

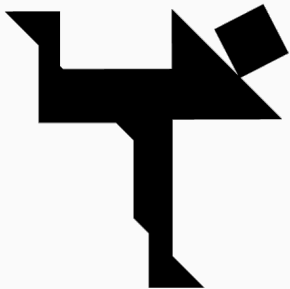
# Application for daily life



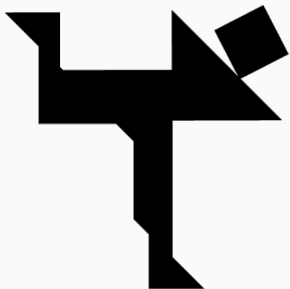
---

Image from <https://bonigarcia.dev/dualsub/>



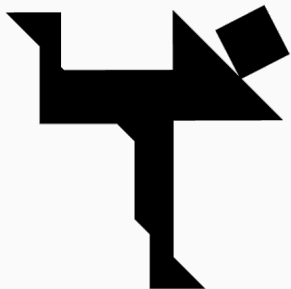


“The guy with one foot on the ground and one leg up.”



“The guy with one foot on the ground and one leg up.”

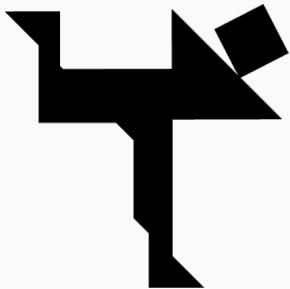
“The ice skater guy.”



“The guy with one foot on the ground and one leg up.”

“The ice skater guy.”

“Someone trying to throw a ball.”

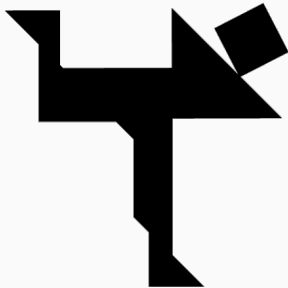


“The guy with one foot on the ground and one leg up.”

“The ice skater guy.”

“Someone trying to throw a ball.”

“Someone dancing or in balance on one foot.”



“The guy with one foot on the ground and one leg up.”

“The ice skater guy.”

“Someone trying to throw a ball.”

“Someone dancing or in balance on one foot.”

“A yoga posture.”

# Application for Language Technologies

The image shows a software interface for Named Entity Recognition. At the top, there is a black header bar with the text "Named Entity Recognition" in white, preceded by a small yellow circle icon. Below the header, there are four colored boxes: "Name" (blue), "Date" (red), "Designation" (green), and "Subject" (yellow). The main area contains a text snippet: "John McCarthy who was born on September 4, 1927 was an American computer scientist and cognitive scientist. He was one of the founders of the discipline of artificial intelligence. He co-authored the document that coined the term 'Artificial intelligence' (AI), developed the programming language family Lisp, significantly influenced the design of the language ALGOL .....". The words "John McCarthy", "September 4, 1927", "American", "cognitive scientist", and "'Artificial intelligence' (AI)" are highlighted with colored boxes corresponding to the categories in the header.

**Name**   **Date**   **Designation**   **Subject**   **Named Entity Recognition**

John McCarthy who was born on September 4, 1927 was an American computer scientist and cognitive scientist. He was one of the founders of the discipline of artificial intelligence. He co-authored the document that coined the term "Artificial intelligence" (AI), developed the programming language family Lisp, significantly influenced the design of the language ALGOL .....

---

Image from *https://www.amygb.ai/blog/what-is-named-entity-recognition-in-nlp*

# Annotation?

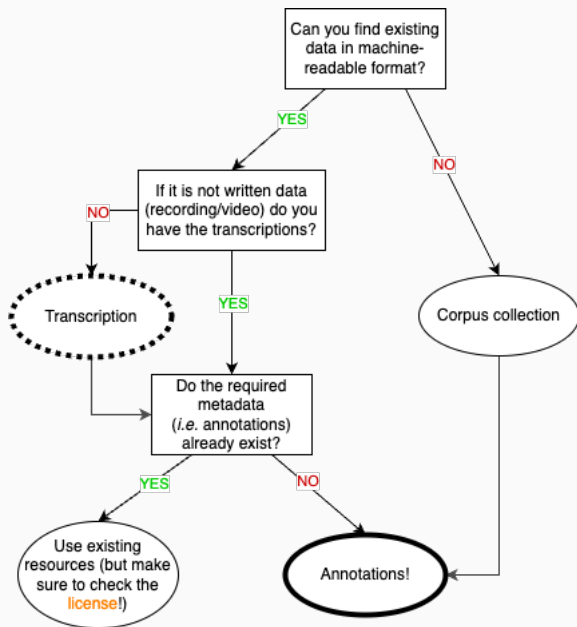
[Leech, 1997]

“[corpus annotation] can be defined as the practice of adding interpretative, linguistic information to an electronic corpus of spoken and/or written language data. ‘Annotation’ can also refer to the end-product of this process”

[Bird and Liberman, 2001]

“‘Linguistic annotation’ covers any descriptive or analytic notations applied to raw language data. The basic data may be in the form of time functions - audio, video and/or physiological recordings - or it may be textual.”

# Define the Goal



# Choose the Annotation Unit(s)

Whole Document

Paragraph

Sentence

Span

Lexical Unit



"All men have the stars," he answered, "but they are not the same things for different people. For some, who are travellers, the stars are guides. For others they are no more than little lights in the sky. For others, who are scholars, they are problems. For my businessman they were wealth. But all these stars are silent. You--you alone--will have the stars as no one else has them--"

"All men have the stars," he answered, "but they are not the same things for different people."

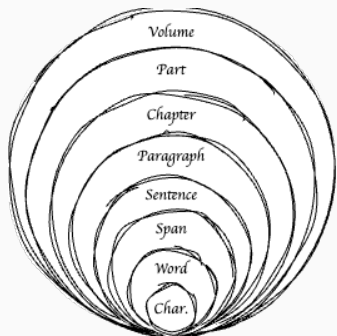
"All men have the stars," he answered, "but they are not the same things for different people."

"All men have the stars," he answered, "but they are not the same things for different people."

Inspired from the slides of Valerio Basile and Enrica Troiano on Data Perspectivism at ESSLLI 2023.

# Choose the Annotation Unit(s)

E.g. Recognise if parts of narrative texts describe a fact or a character's beliefs:



## Alice's Adventures in Wonderland, Chapter II

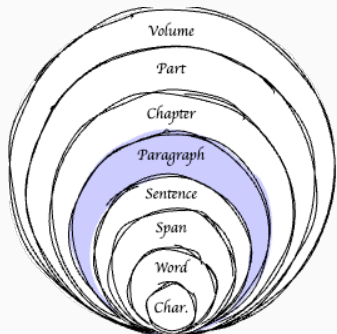
'I'm sure I'm not Ada,' she said, 'for her hair goes in such long ringlets, and mine doesn't go in ringlets at all; and I'm sure I can't be Mabel, for I know all sorts of things, and she, oh! she knows such a very little! Besides, she's she, and I'm I, and—oh dear, how puzzling it all is!

---

Inspired from the slides of Valerio Basile and Enrica Troiano on Data Perspectivism at ESSLI 2023.

# Choose the Annotation Unit(s)

E.g. Recognise if parts of narrative texts describe a fact or a character's beliefs:



## Alice's Adventures in Wonderland, Chapter II

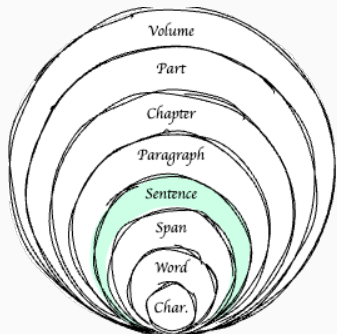
'I'm sure I'm not Ada,' she said, 'for her hair goes in such long ringlets, and mine doesn't go in ringlets at all; and I'm sure I can't be Mabel, for I know all sorts of things, and she, oh! she knows such a very little! Besides, she's she, and I'm I, and—oh dear, how puzzling it all is!

---

Inspired from the slides of Valerio Basile and Enrica Troiano on Data Perspectivism at ESSLLI 2023.

# Choose the Annotation Unit(s)

E.g. Recognise if parts of narrative texts describe a fact or a character's beliefs:



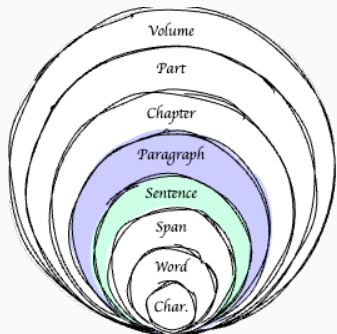
## Alice's Adventures in Wonderland, Chapter II

'I'm sure I'm not Ada,' she said, 'for her hair goes in such long ringlets, and mine doesn't go in ringlets at all; and I'm sure I can't be Mabel, for I know all sorts of things, and she, oh! she knows such a very little! Besides, she's she, and I'm I, and—oh dear, how puzzling it all is!

Inspired from the slides of Valerio Basile and Enrica Troiano on Data Perspectivism at ESSLLI 2023.

# Choose the Annotation Unit(s)

E.g. Recognise if parts of narrative texts describe a fact or a character's beliefs:



## Alice's Adventures in Wonderland, Chapter II

'I'm sure I'm not Ada,' she said, 'for her hair goes in such long ringlets, and mine doesn't go in ringlets at all; and I'm sure I can't be Mabel, for I know all sorts of things, and she, oh! she knows such a very little! Besides, she's she, and I'm I, and—oh dear, how puzzling it all is!

Inspired from the slides of Valerio Basile and Enrica Troiano on Data Perspectivism at ESSLLI 2023.

## How to ask for the target information?

→ The guidelines are instructions to apply your coding scheme to the data. They include:

- ▶ Goal;
- ▶ How to use the tagset:
  - Where to apply the tags (*i.e.*, the annotation units);
  - Examples of good/bad uses;
  - Grey areas and how to deal with ambiguities.

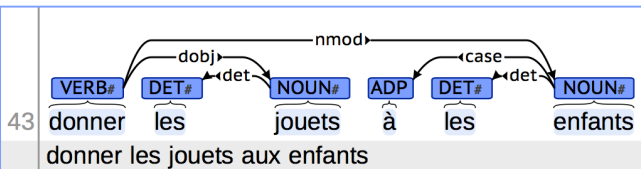
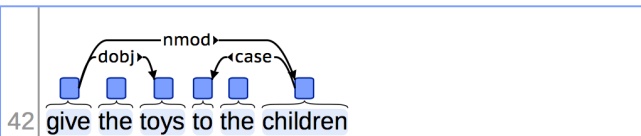
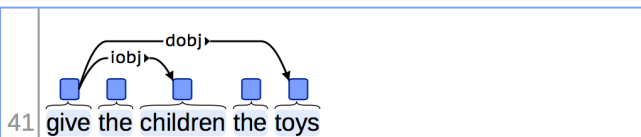
Be **informative** and keep things **clear** and **simple**!

---

Inspired from the slides of Valerio Basile and Enrica Troiano on Data Perspectivism at ESSLLI 2023.



# Choose the Right Tool

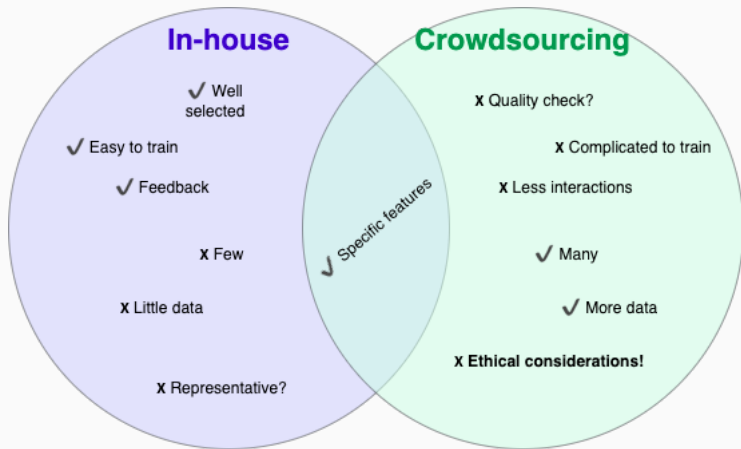


# Choose the Right Tool

The screenshot displays the ELAN 5.1 software interface for the file 'amy\_andrea.eaf'. The top menu bar includes 'File', 'Edit', 'Annotation', 'Tier', 'Type', 'Search', 'View', 'Options', 'Window', and 'Help'. Below the menu is a toolbar with tabs for 'Grid', 'Text', 'Subtitles', 'Lexicon', 'Comments', 'Recognizers', 'Metadata', and 'Controls'. A dropdown menu shows '< select a tier >'. The main workspace is a grid with columns for 'Nr', 'Annotation', 'Begin Time', 'End Time', and 'Duration'. A playback control bar at the bottom shows the current time as 00:40:59.310 and a selected segment from 00:40:59.057 to 00:40:59.883. The audio waveform 'amy.wav' is visible at the top of the grid. The grid contains several tiers of annotations:

- %com@ADD (147)**: A red bar spanning from approximately 00:40:56.000 to 00:41:03.000.
- \*ADD (14891)**: A blue bar containing the text 'no?' from 00:40:59.057 to 00:41:00.000, and another blue bar containing 'oh yeah' from 00:41:02.000 to 00:41:03.000.
- %act@ADD (151)**: A pink bar spanning from approximately 00:40:56.000 to 00:41:03.000.
- %com@BRI (151)**: A yellow bar spanning from approximately 00:40:56.000 to 00:41:03.000.
- \*BRI (14891)**: A grey bar containing the text 'so (.) nothing really new.' from 00:40:56.000 to 00:41:00.000, 'she quit her job and' from 00:41:00.000 to 00:41:02.000, and 'yeah. she was I/ you know with' from 00:41:02.000 to 00:41:03.000.
- Question\_Type (152)**: A white bar containing 'YN' from 00:40:59.057 to 00:41:00.000.
- Feature (144)**: An empty tier.
- Answer\_Type (112)**: An empty tier.
- Quoted (126)**: An empty tier.
- Complexity (16)**: An empty tier.

# Recruit Participants (Annotators)



---

Inspired from the slides of Valerio Basile and Enrica Troiano on Data Perspectivism at ESSLLI 2023.

## How to assess the quality of each annotation?

→ Assumption = Some annotators produce better annotations overall. We find them by:

- ▶ Having a common **subset** of data annotated by everyone;
- ▶ Comparing all the annotators' answers **pairwise**;  
→ Annotators with very low agreement are sometimes ignored.
- ▶ **larger the common subset ↔ better assessment of the agreement**  
**BUT less annotated data in total**

## What you should remember:

- ▶ Define your **goal** clearly;
- ▶ **Guidelines** are crucial for the annotators;
- ▶ Transcribing and annotating is **time-consuming** and **expensive!**  
→ Do not create a resource if it already exists.

# One Transcription / Annotation Tool: ELAN

---

## 1. Transcriptions and Annotations

Examples of Transcription / Annotation Uses

Definition(s)

What is an Annotation Campaign?

## 2. One Transcription / Annotation Tool: ELAN

ELAN?

Examples

Core Concept

Application and Assignment

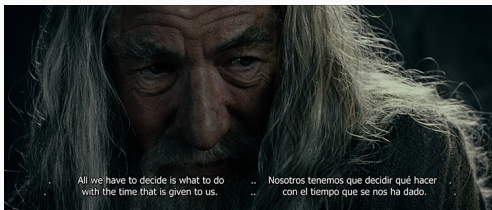
ELAN makes **time-aligned** transcripts:

---

Slides inspired from the course '2019 Linguistic Institute Course 353: Digital Methods in Language Documentation' licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

# What Does it Do?

ELAN makes **time-aligned** transcripts:



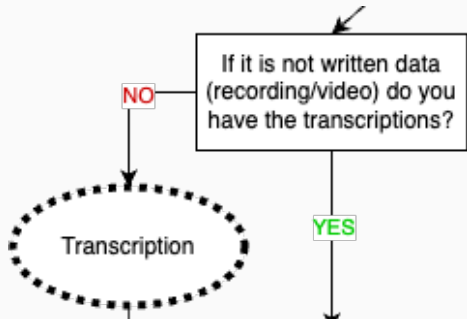
Match **text** (transcription  
and/or annotations)  
with sections of an audio  
or video **recording**.

---

Slides inspired from the course '2019 Linguistic Institute Course 353: Digital Methods in Language Documentation' licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

# What Does it Do?

ELAN makes **time-aligned** transcripts:



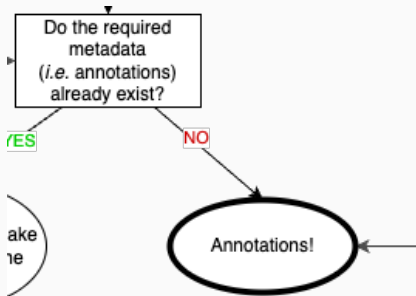
Match **text** (transcription and/or annotations) with sections of an audio or video **recording**.

---

Slides inspired from the course '2019 Linguistic Institute Course 353: Digital Methods in Language Documentation' licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

# What Does it Do?

ELAN makes **time-aligned** transcripts:



Match **text** (transcription and/or annotations) with sections of an audio or video **recording**.

---

Slides inspired from the course '2019 Linguistic Institute Course 353: Digital Methods in Language Documentation' licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

# What Does it NOT Do?

It is **not** an audio/video editor:

- ▶ Cannot edit media;
- ▶ If you change the media (with another software) you **must** change the ELAN file as well;

It is **not** a text editor:

- ▶ Output is *plain text*;
- ▶ Cannot add **embellishments**.

---

Slides inspired from the course '2019 Linguistic Institute Course 353: Digital Methods in Language Documentation' licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

## Example 1: One language, one speaker

(Morality, provided by RadioLab <http://www.wnyc.org/shows/radiolab>)

Josh: “How do people make this judgment? Forget whether or not these judgments are right or wrong, just, what’s going on in the brain that makes people distinguish so naturally and intuitively between these two cases? Which, from an actuarial point of view, are very very very similar if not identical...”

---

Slides inspired from the course ‘2019 Linguistic Institute Course 353: Digital Methods in Language Documentation’ licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

## Example 2: One language, multiple speakers

(Excerpt: Television)

ALVIN: yeah, I haven't –

like admittedly,

I haven't,

(SNIFF) (0.3)

It's funny, I haven't watched those in years, you know,

PETER: Yeah.

ALVIN: so I've thought, (0.7) I've thought,  
it might be fun to see them again.

---

Slides inspired from the course '2019 Linguistic Institute Course 353: Digital Methods in Language Documentation' licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

# Examples of Applications

## Example 3: Two languages, richer linguistic info

(Pear Story, Kannada. Speaker: Keshava Subramanya)

**KAN Sentence:** li kathe obba, haṇṇu maaliya bagge, mattu, obba huḍugana bagge ide.

---

<b>Intonation Unit</b>	li	kathe	obba,	haṇṇu	maaliya
<b>Morphemes</b>	ii	kathe	obba	haṇṇu	maali -ya
<b>Gloss</b>	this	story	one	fruit	gardener- SRC

---

bagge,	mattu,	obba	huḍugana	bagge	ide.
bagge	mattu	obba	huḍuga -na	bagge	ide
about	and	someone.M	boy -GEN	about	be.3SM

---

**Free translation:** ‘This story is about a fruit farmer and a boy.’

Slides inspired from the course ‘2019 Linguistic Institute Course 353: Digital Methods in Language Documentation’ licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

ELAN enables you to add as many **layers** of annotations as you want

→ Everything that is the same “kind” of information is part of the same **tier**:

- ▶ All the transcription of a monologue (Example 1);
- ▶ All the utterances of one speaker in a dialogue (Example 2);
- ▶ All the glosses of the Kannada words in Example 3.

---

Slides inspired from the course ‘2019 Linguistic Institute Course 353: Digital Methods in Language Documentation’ licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

# Concept of Tier

You are filming a 2-character puppet show designed to teach English speakers sentences in Irish:

- ▶ How many tiers do you need?
- ▶ What are they?



# Concept of Tier

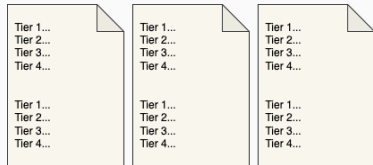
1. Fox's utterances in Irish
2. Translations of Fox's sentences into English.
3. Archer's utterances in Irish
4. Translations of Archer's sentences into English.



# Concept of Tier

Printed transcripts have a **fixed** page width;

→ Forced to **split content** of single tier across many lines.



# Concept of Tier

But ELAN is **time-based**!

Tiers are **continuous** and **contiguous** to the media timeline:



# Concept of Tier

But ELAN is **time-based!**

Multiple continuous tiers can capture **speaker overlap:**



## A Practice Example

Install ELAN:

*<https://archive.mpi.nl/tla/elan/download>*;

Download the **recording** we will work on (Canvas);

Open the recording in ELAN and:

- ▶ **Transcribe** it;
- ▶ Annotate the questions with **question types** (Wh-question, Yes-No-Question, Other);
- ▶ Consider that the **goal** of this experiment is to gather empirical data to **determine the frequency at which people use the different question types**;

You can find help on how to use ELAN on **these slides**.





# Assignment

Write a paragraph on the transcription/annotation task you just did:

- ▶ How many tiers did you use and what were they?
- ▶ Why did you choose these tiers?
- ▶ Was the research question specific enough to make decisions regarding the tiers and the annotations? If not what additional information could have been helpful?
- ▶ What difficulties did you encounter?
- ▶ How much time did it take? Do you feel that it is a lot based on the length of the recording and the difficulty of the task?

*It does not need to be very long, just show that you understand the concept of tier and how to take advantage of it for transcription and annotations. Also show that you are aware of the difficulties when it comes to creating transcriptions and defining an annotation task.*

Send your paragraph by email to [amandine.decker@gu.se](mailto:amandine.decker@gu.se)

-  Berez-Kroeker, Andrea and Christopher Cox (2019). *Linguistic Institute Course 353: Digital Methods in Language Documentation – Days 3&4: ELAN Lesson 1*. Lecture Slides.
-  Bird, Steven and Mark Liberman (2001). **“A formal framework for linguistic annotation”**. In: *Speech Communication* 33.1. Speech Annotation and Corpus Tools, pp. 23–60.
-  Dostal, Mateja (2021). *Business English learner speech corpus SAPS*. Slovenian language resource repository CLARIN.SI.
-  Leech, Geoffrey (1997). **“Corpus annotation : Linguistic information from computer text corpora”**. In: Longman, Londres, Angleterre. Chap. Introducing corpus annotation, pp. 1–18.