

Annotation  
mit  
Sprachmodellen  
zur Sprachmodellierung  
mit  
Annotation



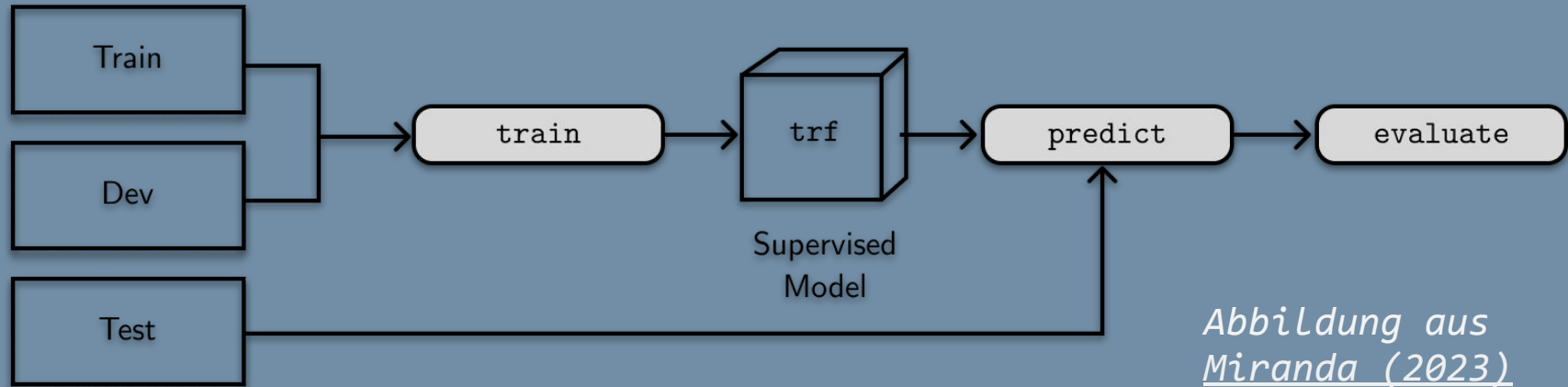
# Annotation von Textdaten

- **Annotation:** Auszeichnung von (Text)Einheiten mit Labels / Verknüpfung von Sprachdaten mit deskriptiven oder analytischen Notationen.  
Beispiel.
- **Nutzen** in verschiedenen Anwendungsbereichen: Informationsextraktion, Überprüfung linguistischer Theorien, Training von sprachverarbeitenden Anwendungen
- **Annotationsschema:**
  - Definiert gültige Auszeichnungen / Annotationen (Labels), die sprechend gewählt werden sollen, um auffindbar zu sein
  - Definiert mit welchen Einheiten (Annotation Units) die Labels assoziiert werden sollen
  - Beispiel PoS-Tags.

# Annotation von Textdaten

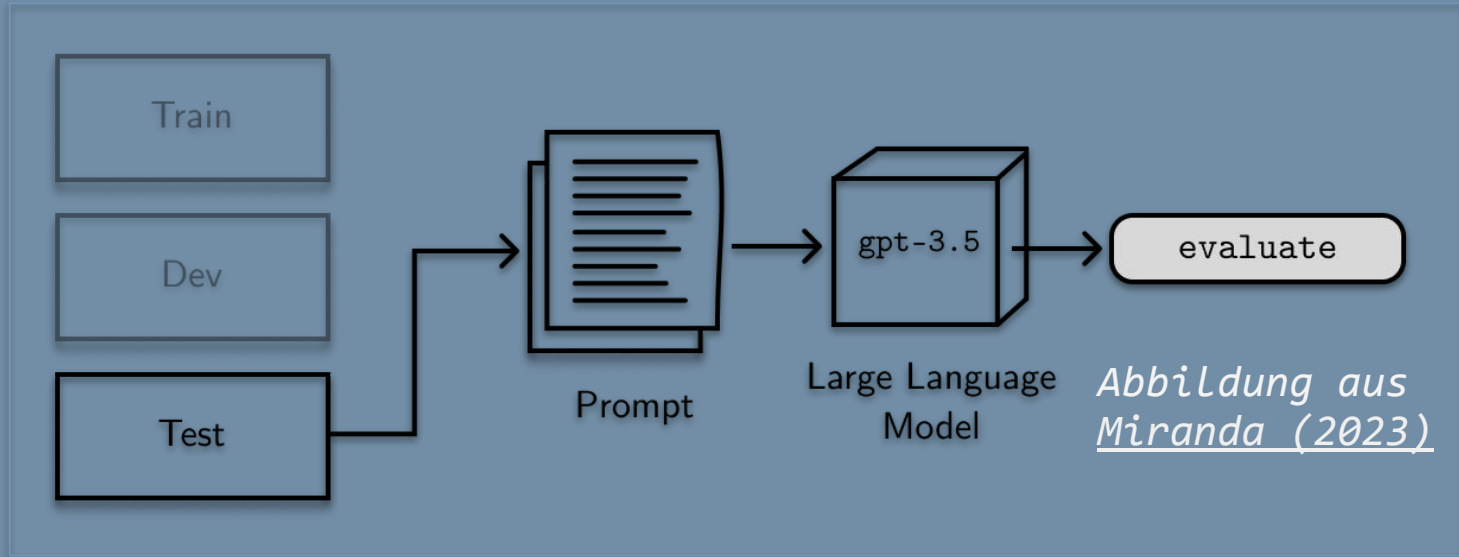
- **Richtlinien zur Tokenisierung:** Legen fest, wie Einheiten (Annotation Units, wie z.B. Wörter, Phrasen) aus den Rohdaten (Text, Audioaufzeichnung etc.) gewonnen werden.
- **Richtlinien zur Annotation:** Legen fest, wie bei der Zuweisung von Labels zu Einheiten vorgegangen werden soll (ggfs. zyklisch entwickelt).
- **Inter-Annotator-Agreement:** Übereinstimmung von Annotationen zwischen verschiedenen (meist menschlichen) Annotatoren. Wie über unterschiedliche Metriken ermittelt (z.B. Cohen's Kappa, Krippendorff's Alpha)
- **Annotation Tools:** Software zu manueller (Beispiel: INCEPTION) oder automatischer (Beispiel: Stanza) Annotation oder Mischformen daraus.

# Maschinelles Lernen auf annotierten Daten



- Standard-Workflow Maschinelles Lernverfahren: Erst trainieren, dann finetunen, dann testen, dann evaluieren.
- Um Daten zu annotieren, benötigt man annotierte Daten!
- Bsp: Klassifikation von Abschnitten

# Zero-Shot-Learning mit LLMs



- Workflow zur Annotation mit LLMs
- Um Daten zu annotieren, benötigt man lediglich Testdaten.

# Few-Shot-Learning mit LLMs

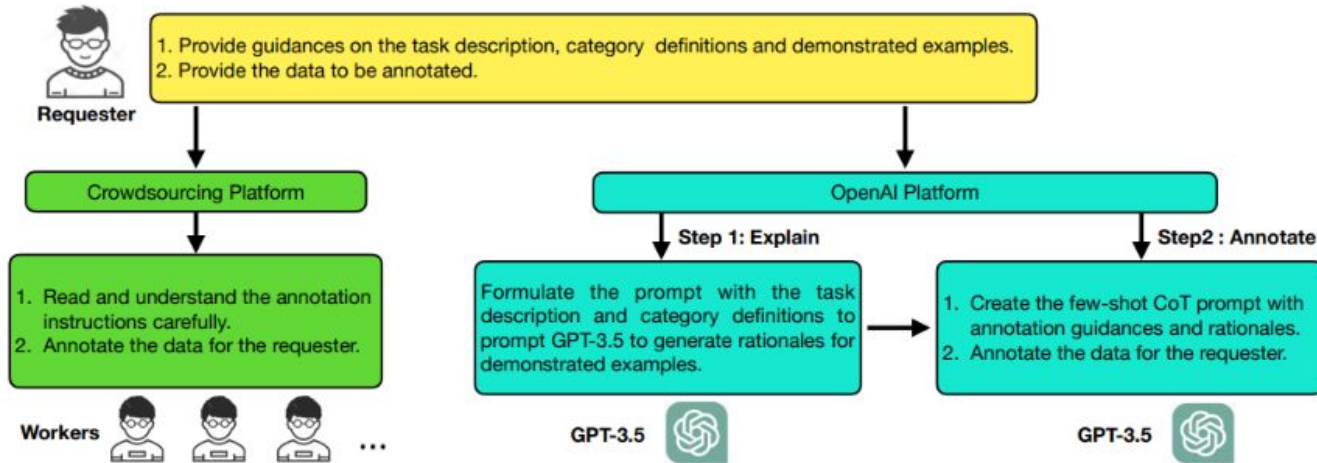


Figure 1: On the left is the annotation process used by crowdsourced workers, while on the right is AnnoGPT's process. AnnoGPT mimics the manual annotation process, with the exception that it generates explanations for each example before annotation. This ensures that each demonstrated example is accompanied by helpful explanations, making the annotation guidelines more informative and useful."

*Abbildung aus  
He et al. (2023)*

# Mögliche Tasks (Vorschläge, kann ergänzt werden)

- Klassifikation von Kurznachrichten (andere Textsorten möglich)
  - Erkennung von Spam, HateSpeech, Themenrelevanz
  - Kategorisierung in festgelegte Topics
  - Erkennung der Polarität einer Einstellung (für, neutral, gegen)
- Resolution von Anaphern
  - Identifikation sämtlicher Bezüge zu Person/Ereignis X in einem Text
  - Präferierte Lesarten für ambige Pronomen/Anaphern (Bsp. Der Wein steht auf dem Tisch. Er/der/dieser/jener ist alt)

# Aufgabe bis zum 27.11. (siehe ILIAS)

Suchen Sie sich ein Thema, zu dem Sie ein Sprachmodell Daten annotieren lassen können.

Spezifizieren Sie dabei folgendes:

1. Annotationsaufgabe (gerne mit Beispiel)
2. Datenbasis (Textsorte, Sprache, Quelle)
3. Sprachmodell / Interface
4. (Nur grob, wird noch besprochen) Experimentdesign - wie wollen Sie vorgehen und wie die Ergebnisse beurteilen?