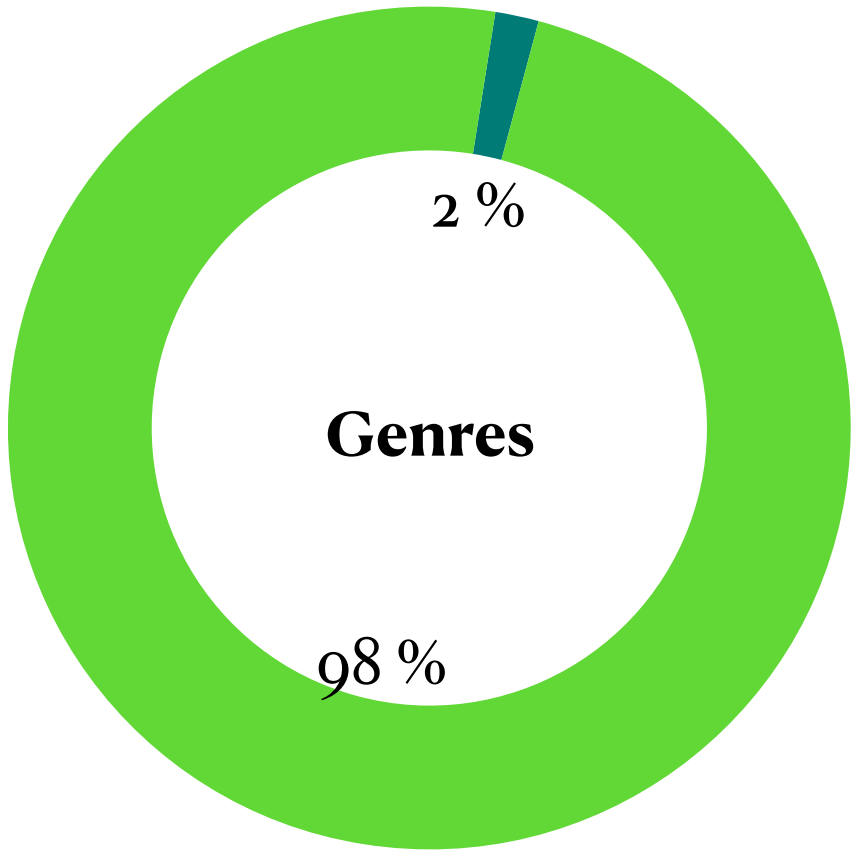# Spotify Podcast Dataset

- 100.000 podcast episodes with aligned ASP transcripts

- more than 47.000 hours of transcribed audio

- automatic transcripts

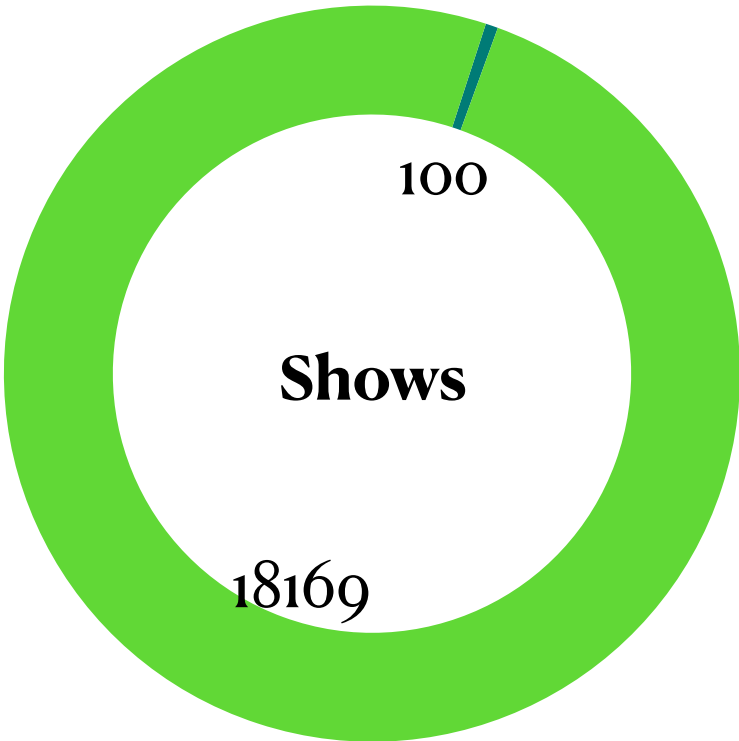- Word Error Rate: 18.1%

- has punctuation

## General

### Genres

Total Genres · Relevant Genres
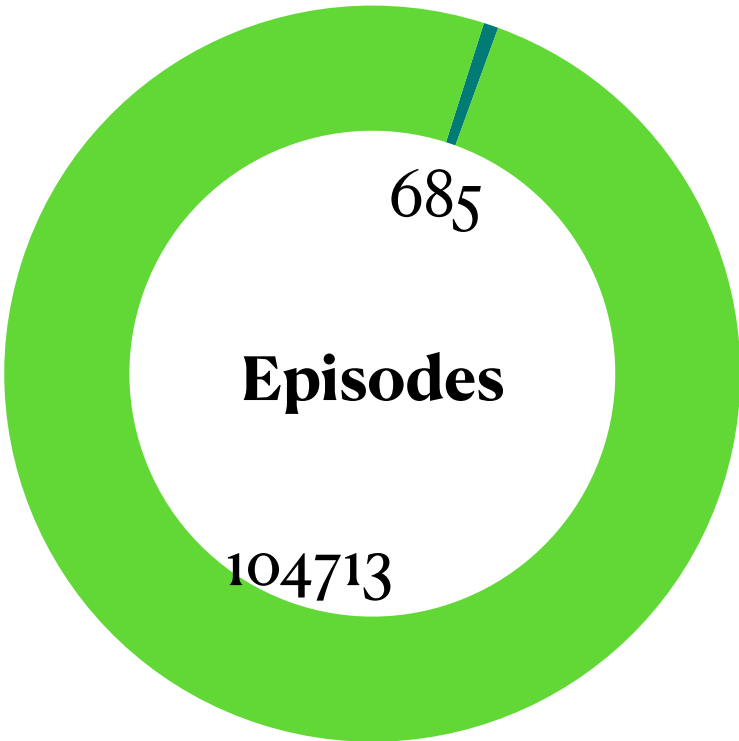
2 %

**Genres**

98 %

## Relevant Genres

- Business News
- Daily News
- News
- news
- Sports News
- Tech News

### Shows

Total Shows · Relevant Shows

100

**Shows**

18169

### Episodes

Total Episodes · Relevant Episodes

685

**Episodes**

104713

## Episode data

|  | min | average | max |
|---|---|---|---|
| minutes | 1 | 31,6 | 305 |
| words | 11 | 5728 | 43504 |

Quelle: The Spotify Dataset Paper

# Relevant Shows



Daily News
2 %

Business News
3 %

Tech News
3 %

News
73 %

Sports News
19 %

news
1 %

● Business News    ● Daily News
● News             ● news
● Sports News      ● Tech News

# Relevant Data

## Relevant Episodes



Business News
2 %

Tech News
4 %

News
73 %

Sports News
16 %

news
4 %

● Business News    ● Daily News
● News             ● news
● Sports News      ● Tech News

\* The List of relevant shows contains shows that may have
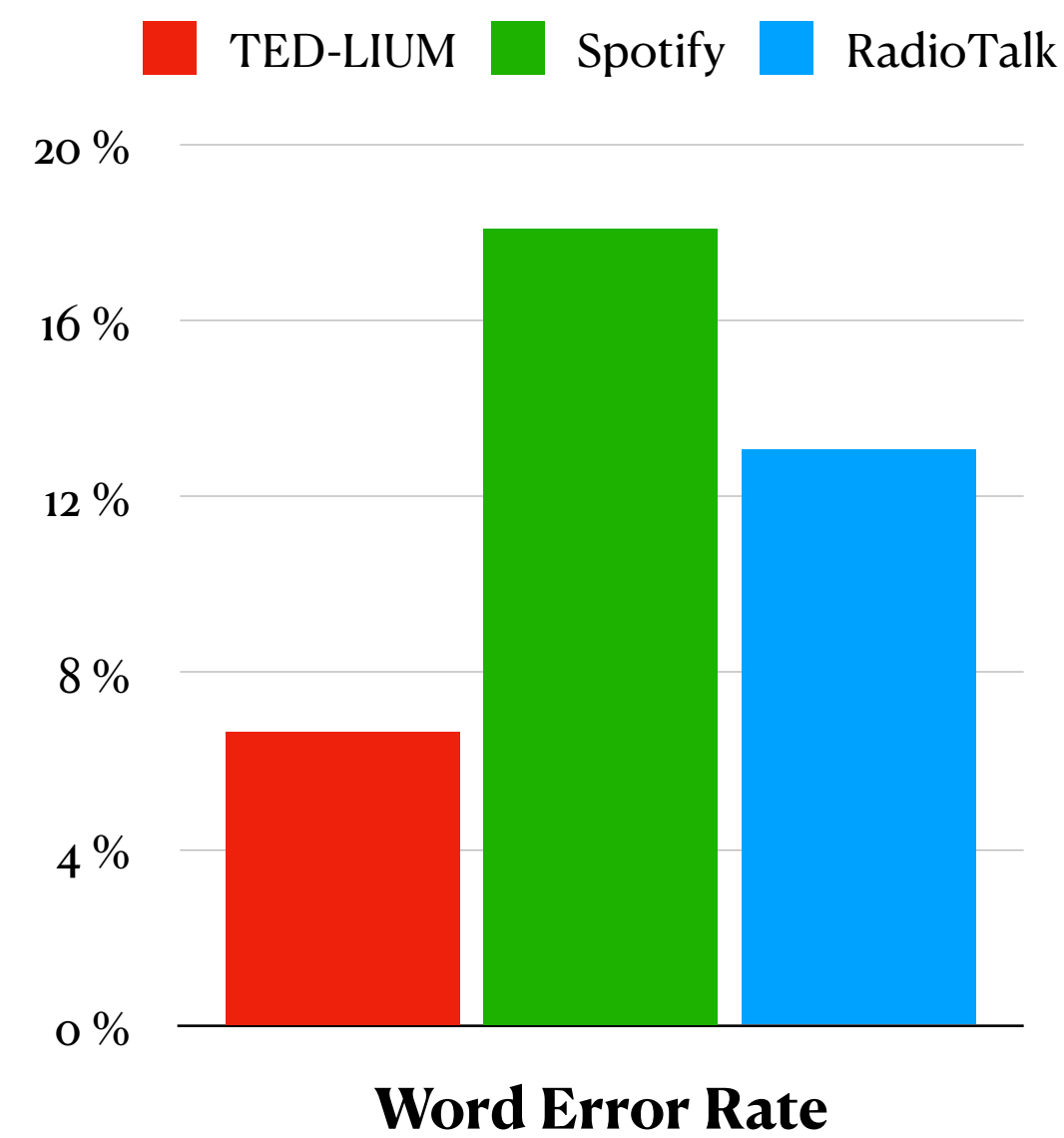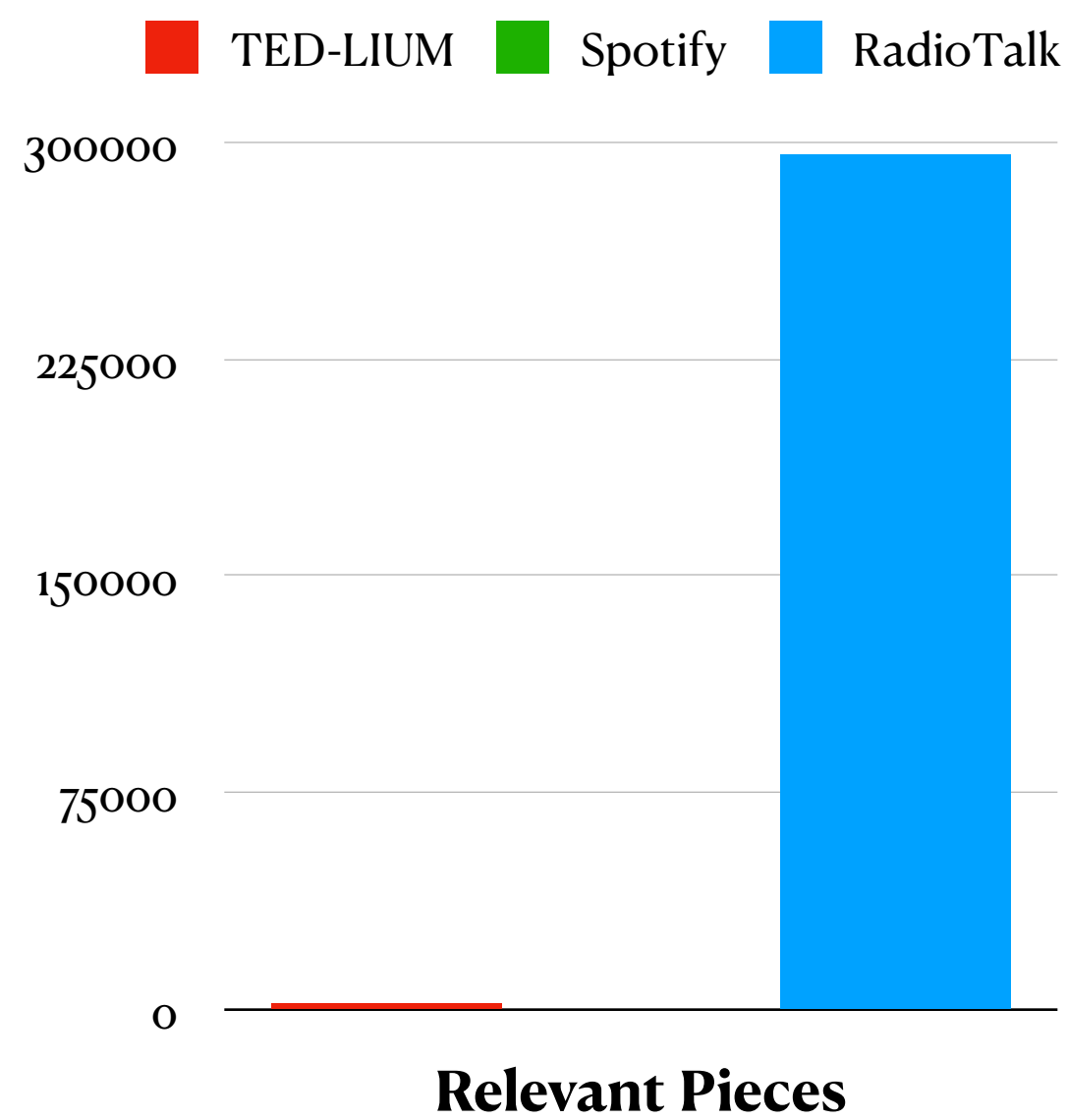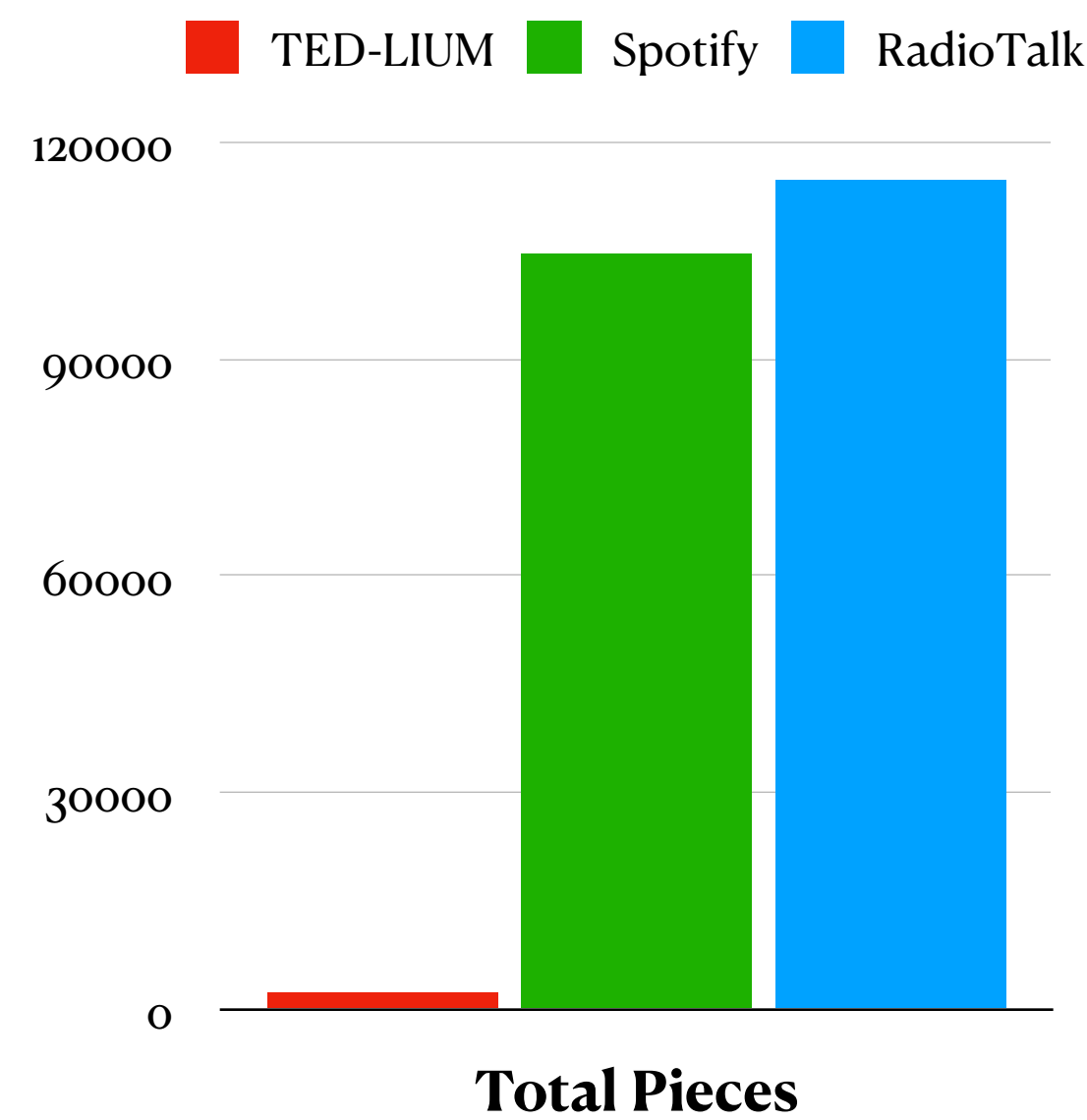multiple relevant genres and are therefore listet multiple times.

- We can use approx. 685 Episodes from 100 news-related shows

- The already collected Data contains a list of all relevant shows incl. their respective show_uri

- the transcription files are named after the show_uri

- 452 hours of transcribed speech

- 2351 speeches

- 4.9 M words

- automatic transcripts

- Word Error Rate: 6.7%

- no punctuation

- no genre information

# Comparison

**Total Pieces** — **Relevant Pieces** — **Word Error Rate**

Checklist

| Category | TED-LIUM | Spotify | RadioTalk |
|---|---|---|---|
| Punctuation | ☐ | ✔ | ☐ |
| Genre-Info | ☐ | ✔ | ☐ |
| Contains News | ☐ | ✔ | ✔ |
| Contains Opinion Pieces | ✔ | ✔ | ✔ |
| „Good enough" Transcripts | ✔ | ✔ | ☐ |
| Enough relevant Data | ✔ ★ | ✔ | ✔ |
| Human-made Transcripts | ☐ | ☐ | ☐ |
| RSS-Feed | ☐ | ✔ | ☐ |
| Free | ✔ | ✔ | ✔ |
| Variety of Producers | ✔ | ✔ | ✔ |

★ according to Ortmann and Dipper (2019) speeches have the same *index of orality* as the texts we want to produce:

- many participants = -1
- monologue = -1
- synchronous production = -1
- asynchronous reception = 1

—> **Index of orality** = **-2**