

Diskursmarker

in

schriftlichem & akustischem

Diskurs

BACHELORVERTEIDIGUNG

Johanna Sacher, 4.2.2021



Diese Arbeit liefert Evidenz für Unterschiede zwischen

- oral-akustischem und literat-schriftlichem Diskurs
- geskripteten und improvisierten oral-akustischen Texten
- interaktiven und passiven oral-akustischen Texten

Im Folgenden wird die unterschiedliche Verwendung von **Diskursmarkern** in den genannten Textsorten nachgewiesen

- oral-akustische Texte nutzen mehr Diskursmarker als literat-schriftliche
- improvisierte Texte nutzen mehr Diskursmarker als geskriptete
- interaktive Texte nutzen mehr Diskursmarker als passive

BEGRIFF

Diskurs

Einheit von Sprache, länger als ein einzelner Satz

Quelle

BEGRIFF

Diskursmarker

Wörter wie *and*, *but* und *so*

- keine inhaltliche Bedeutung
- signalisieren Beziehungen zwischen Diskurssegmenten
- Wegweiser im Text

Literat vs. Oral

Literat – Konzept für das **schriftliche** Medium

- » literat-schriftliche Texte (LS)
- » Readability

Oral – Konzept für das **akustische** Medium

- » oral-akustische Texte (OA)
- » Listenability

Medien	Konzepte	
	literat	oral
schriftlich	Stummes lesen eines Zeitungsartikels	Stummes lesen eines YouTube Kommentars
akustisch	Anhören eines vorgelesenen Zeitungsartikels	Persönliches Gespräch

MOTIVATION

Wieso dieses Thema?

Verwendung von Sprachassistenten zum Vorlesen von z.B. Zeitungsartikeln



By **Dennis Overbye**

Published Jan. 19, 2021 Updated Jan. 20, 2021

Astronomers are searching the cosmic lost-and-found for one of the biggest, baddest black holes thought to exist. So far they haven't found it.

In the past few decades, it has become part of astronomical lore, if not quite a law, that at the center of every luminous city of light, called a galaxy, lurks something like a hungry Beelzebub, a giant black hole into which the equivalent of millions or even billions of suns have disappeared. The bigger the galaxy, the more massive the black hole at its center.

Zeitungsartikel wurde geschrieben, um gelesen zu werden

⇒ vorgelesen ggf. nicht mehr so gut verständlich

Wie können Texte so formuliert werden,
dass sie über beide Medien funktionieren?



Welche Faktoren erhöhen die **Listenability** eines Textes?

Listenability \neq Readability

- kaum Forschung zu Listenability
 - » kaum Methoden zur Messung
- Readability Measures nur für linguistische Features

- Kurze Sätze
- Einfache Wörter
- Zahlen runden
- Koordination / Bindewörter / Diskursmarker

DISKURSMARKER

Begriff

- Begriff ist nicht eindeutig definiert
- Verschiedene Begriffe in Benutzung

Funktionale Definiton nach Bruce Fraser

- DM ist ein lexikaler Ausdruck
- In $\langle S1 \ S2 \rangle$ muss ein DM Teil von S2 sein
- DM trägt nicht zur semantischen Bedeutung der Sequenz bei, sondern signalisiert eine Relation zwischen S1 und S2

I love the Shire. But I begin to wish, somehow, that I had gone, too.

I love the Shire. I begin to wish, somehow, that I had gone, too.

You are the master of Bag End now. And also, I fancy, you'll find a golden ring.

You are the master of Bag End now. I fancy, you'll find a golden ring.

You are the master of Bag End now. You'll find a golden ring.

Kriterien des EnDimLex

- DM ist ein lexikaler Ausdruck und kann nicht flektiert werden
- DM signalisiert eine zweiseitige Relation, deren Argumente abstrakte Objekte sind
- Argumente können in Klausel-, Satz- oder Phrasenstruktur ausgedrückt werden

Weitere Bedingungen

- feststehender, nicht modifizierbarer Ausdruck
 - » nicht: *for this reason (for this **exact** reason)*
- darf nicht semantisch kombinierbar sein
 - » nicht: *particularly if*
 - » feststehende Phrasen sind ok: *even if*

DISKURSMARKER

Vergleich

Funktionale Definition (Fraser)		EnDimLex-Kriterien
lexikaler Ausdruck	✓	lexikaler Ausdruck
	→	kann nicht flektiert werden
signalisiert Relation zwischen Diskurssegmenten	✓	signalisiert zweiseitige Relation zwischen Klauseln, Sätzen oder Phrasen
trägt nicht zur Bedeutung des Satzes bei	←	
meistens Adverbien, Konjunktionen, Präpositionalphrasen	✓	meistens Adverbien, Konjunktionen, Präpositionalphrasen

DM setzen sich aus verschiedenen anderen Wortgruppen zusammen

» erschwert automatische Erkennung

Bilbo won the ring. *As a result*, Gollum was very angry. (*Diskursmarker*)

Gollum was very angry *as a result* of Bilbo winning the ring. (*Adverb*)

DISKURSMARKER

Zusammenfassung

- keine inhaltliche Bedeutung
- signalisieren Beziehungen zwischen Diskurssegmenten
- setzen sich aus verschiedenen Wortgruppen zusammen
- automatische Erkennung schwierig

DISKURSMARKER

Bedeutungsgruppen

Gibt verschiedene Ansätze, DM anhand ihrer Bedeutung in Klassen aufzuteilen

Fraser

CONTRASTIVE MARKERS Kontrast zwischen S1 und S2

but, *alternatively, although, even so, still, yet, ...*

ELABORATIVE MARKERS Genauere Ausführung von S1 in S2

and, *also, besides, for instance, moreover, similarly, ...*

INFERENTIAL MARKERS S2 kann aus S1 gefolgert werden

so, *consequently, therefore, thus, ...*

EnDimLex

COMPARISON Vergleich

but, although, in contrast, still, while, yet, ...

CONTINGENCY Folgern, Möglichkeiten aufzeigen

so, for, because, given, in case, whatever, ...

EXPANSION Hinzufügen eines Aspektes

and, also, besides, finally, instead, rather, ...

TEMPORAL Zeitlicher Bezug

afterwards, as, before, next, thereafter, ...

EnDimLex	Fraser	Funktion	Beispiele
Comparison	Contrastive	Vergleich, Kontrast	<i>but, although, still, yet, ...</i>
Contingency	Inferential	Folgern	<i>so, for, because, thus, ...</i>
Expansion	Elaborative	Ausführen, Illustrieren	<i>and, also, besides, ...</i>
Temporal	Discourse Structure Markers	Zeitlich in Bezug setzen	<i>as, before, next, ...</i>

DM können in mehrer Klassen gleichzeitig fallen:

Sam and Pippin crouched behind a large tree-bole, **while** Frodo crept back a few yards towards the lane.

(Temporal & Comparison)

Since they were all hobbits, and were trying to be silent, they made no noise that even hobbits would hear.

(Contingency)

I came also upon two others, but they turned away southward. **Since** then I have searched for your trail.

(Temporal)

TEXTSORTEN

Diskursarten

- oral-akustisch
- literat-schriftlich

TEXTSORTEN

Genres

- News
- Discussion
- Science/Education
- Documentary
- Presentation

TEXTSORTEN

Konversationsarten

- Dialog
- Monolog
- Kooperativer Monolog
- Rede

TEXTDATEN

Corpora

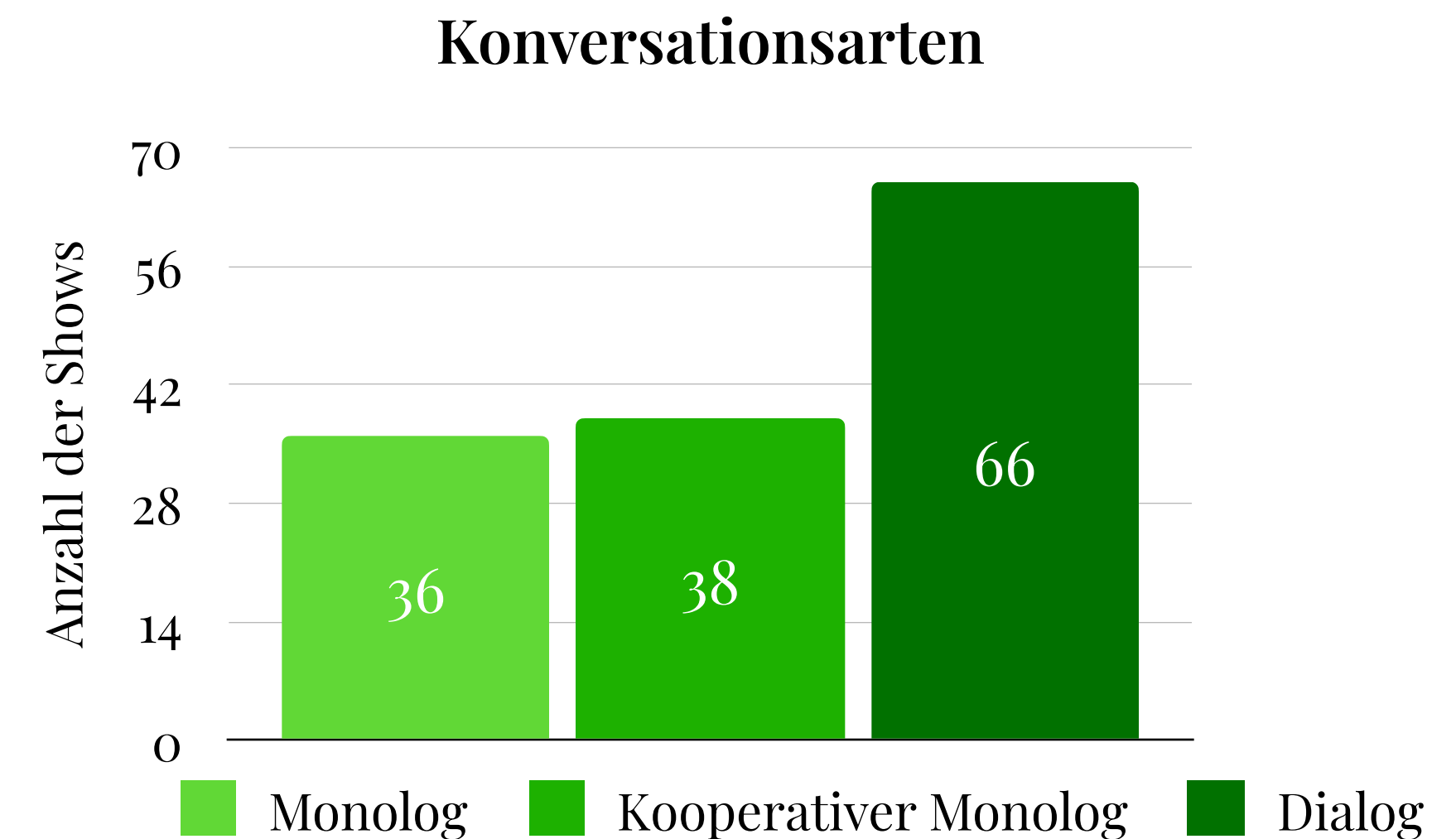
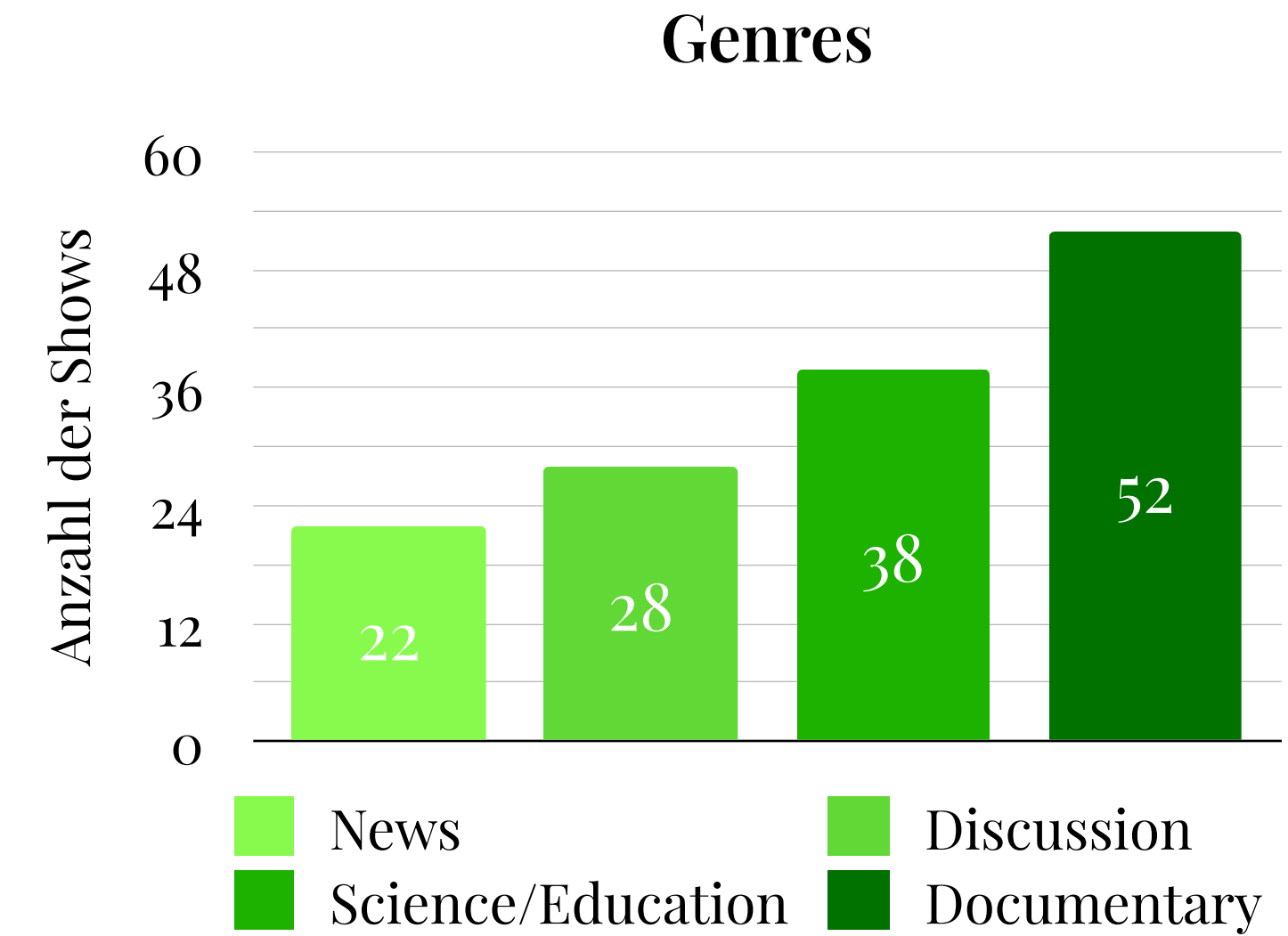
akustische Corpora mit Transkripten von Audiomaterial &

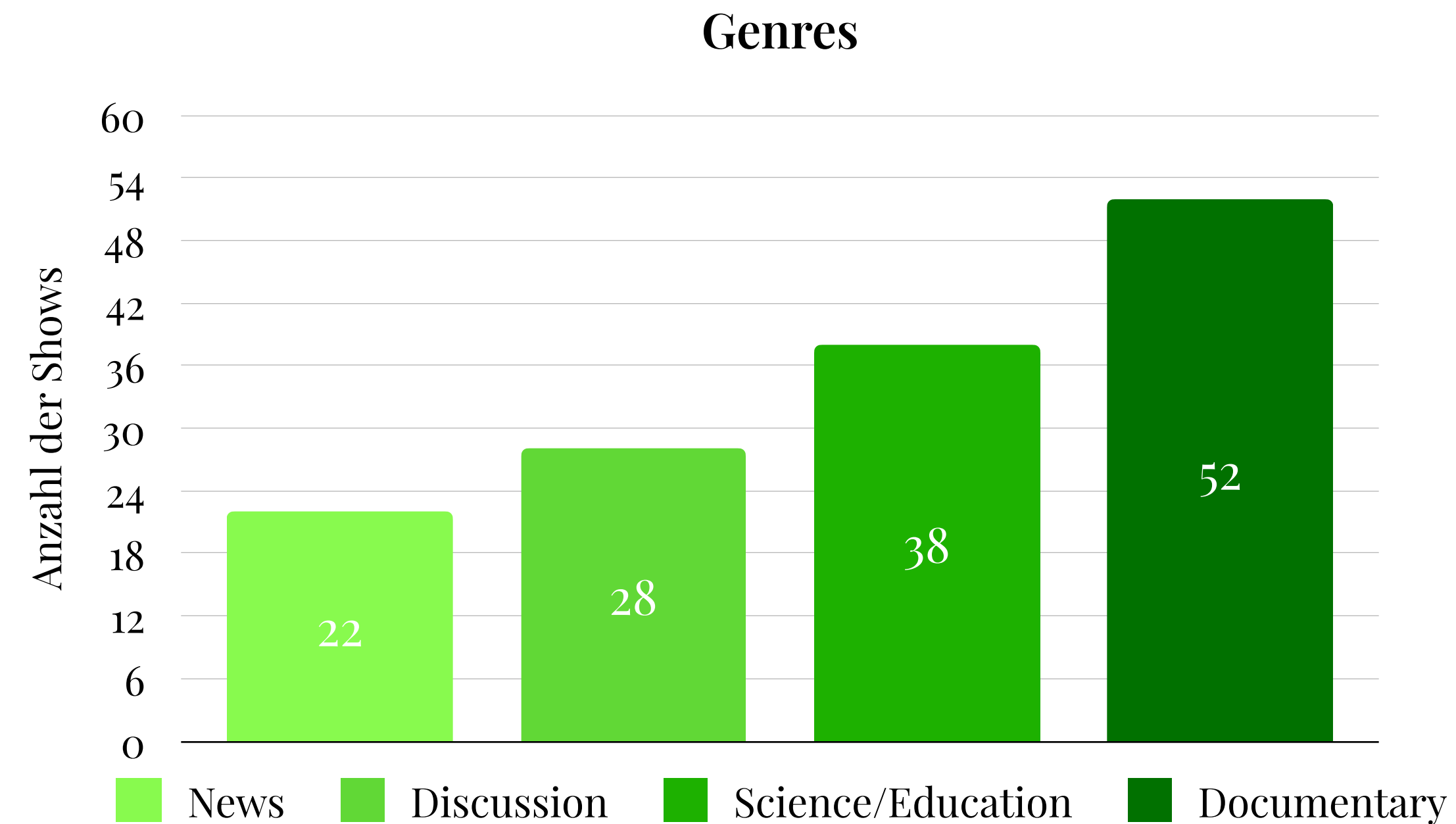
schriftliche Corpora mit ursprünglich schriftlichem Material

- kostenlos
- groß
- qualitativ hochwertig
- nachrichtenähnlich

CORPUS	Spotify Podcast Corpus	TED-LIUM 3 Corpus
DATEN	<ul style="list-style-type: none">• fast 60.000 Stunden transkribiertes Audiomaterial• verschiedenste Produzenten• WER: 18,1 %	<ul style="list-style-type: none">• 1.983 TED-Talks• ca. 4 Mio. Wörter• WER: 6,7 %

CORPUS	Spotify Podcast Corpus
DATEN	<ul style="list-style-type: none"> • fast 60.000 Stunden transkribiertes Audiomaterial • verschiedenste Produzenten • WER: 18,1 %
NUTZBAR	<ul style="list-style-type: none"> • 140 Shows, 2.782 Episoden • ca. 17 Mio. Wörter



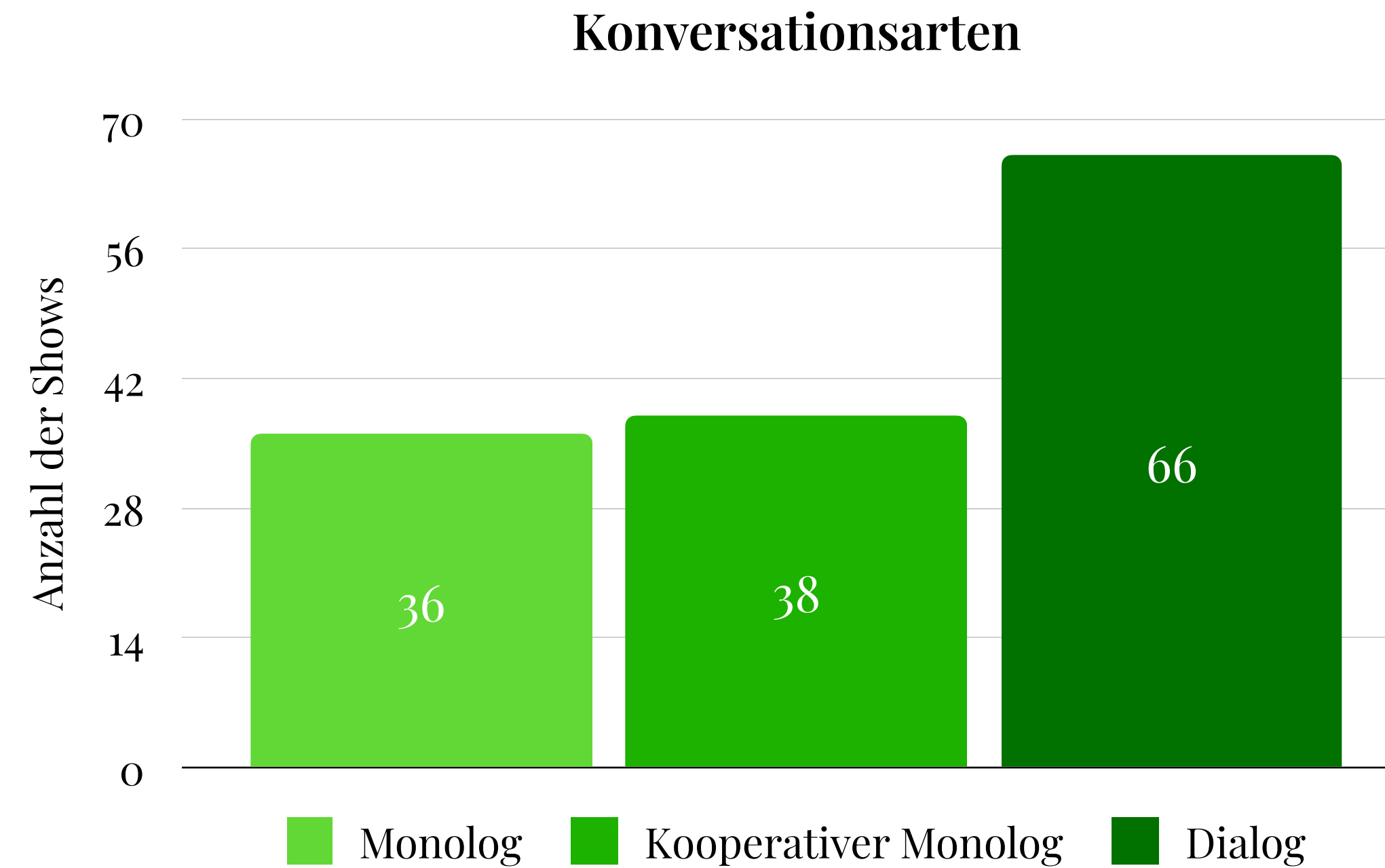


NEWS Fokus auf Nachrichten

DISCUSSION Fokus auf Diskussionen und Meinungsaustausch

SCIENCE/EDUCATION übermitteln Wissen

DOCUMENTARY geskriptet, gut recherchiert, zu einem bestimmten Thema



MONOLOG Hauptsächlich eine Person spricht

KOOPERATIVER MONOLOG mehrere Personen sprechen zum gleichen Thema, aber nicht miteinander

DIALOG mindestens zwei Personen reden miteinander

CORPUS	TED-LIUM 3 Corpus
DATEN	<ul style="list-style-type: none">• 1.983 TED-Talks• ca. 4 Mio. Wörter• WER: 6,7 %
NUTZBAR	<ul style="list-style-type: none">• Alle Talks

Genre

Presentation
100 %

Konversationsart

Rede
100 %

Genre



PRESENTATION vor einer Menge Zuhörer nach einem vorbereiteten Skript präsentiert

Konversationsart



REDE vor einer Menge Zuhörer nach einem
vorbereiteten Skript gehalten

CORPUS	Spotify Podcast Corpus	TED-LIUM 3 Corpus
DATEN	<ul style="list-style-type: none">• fast 60.000 Stunden transkribiertes Audiomaterial• verschiedenste Produzenten• WER: 18,1 %	<ul style="list-style-type: none">• 1.983 TED-Talks• ca. 4 Mio. Wörter• WER: 6,7 %
NUTZBAR	<ul style="list-style-type: none">• 140 Shows, 2.782 Episoden• ca. 17 Mio. Wörter	<ul style="list-style-type: none">• Alle Talks

CORPUS	New York Times Corpus
DATEN	<ul style="list-style-type: none">• 1,8 Mio. Nachrichtenartikel der New York Times• ca. 1,1 Mrd. Wörter
NUTZBAR	<ul style="list-style-type: none">• Alles

CORPUS	New York Times Corpus	Gigaword Corpus
DATEN	<ul style="list-style-type: none">• 1,8 Mio. Nachrichtenartikel der New York Times• ca. 1,1 Mrd. Wörter	<ul style="list-style-type: none">• Newswire Textdaten• aus 7 Quellen• ca. 4 Mrd. Wörter
NUTZBAR	<ul style="list-style-type: none">• Alles	<ul style="list-style-type: none">• Alles

TYP	Akustische Corpora		Schriftliche Corpora	
CORPUS	Spotify	TED-LIUM 3 Corpus	New York Times	Gigaword
NUTZBAR	<ul style="list-style-type: none"> • 140 Shows, 2.782 Episoden • ca. 17 Mio. Wörter 	<ul style="list-style-type: none"> • 1.983 TED-Talks • ca. 4 Mio. Wörter 	<ul style="list-style-type: none"> • 1,8 Mio. Artikel • 1,1 Mrd. Wörter 	<ul style="list-style-type: none"> • ca. 4 Mrd. Wörter

Diskursmarker im Text erkennen

Diskursmarker wurden mit Hilfe eines einfachen
String-Matching Verfahren mit den Texten gematched

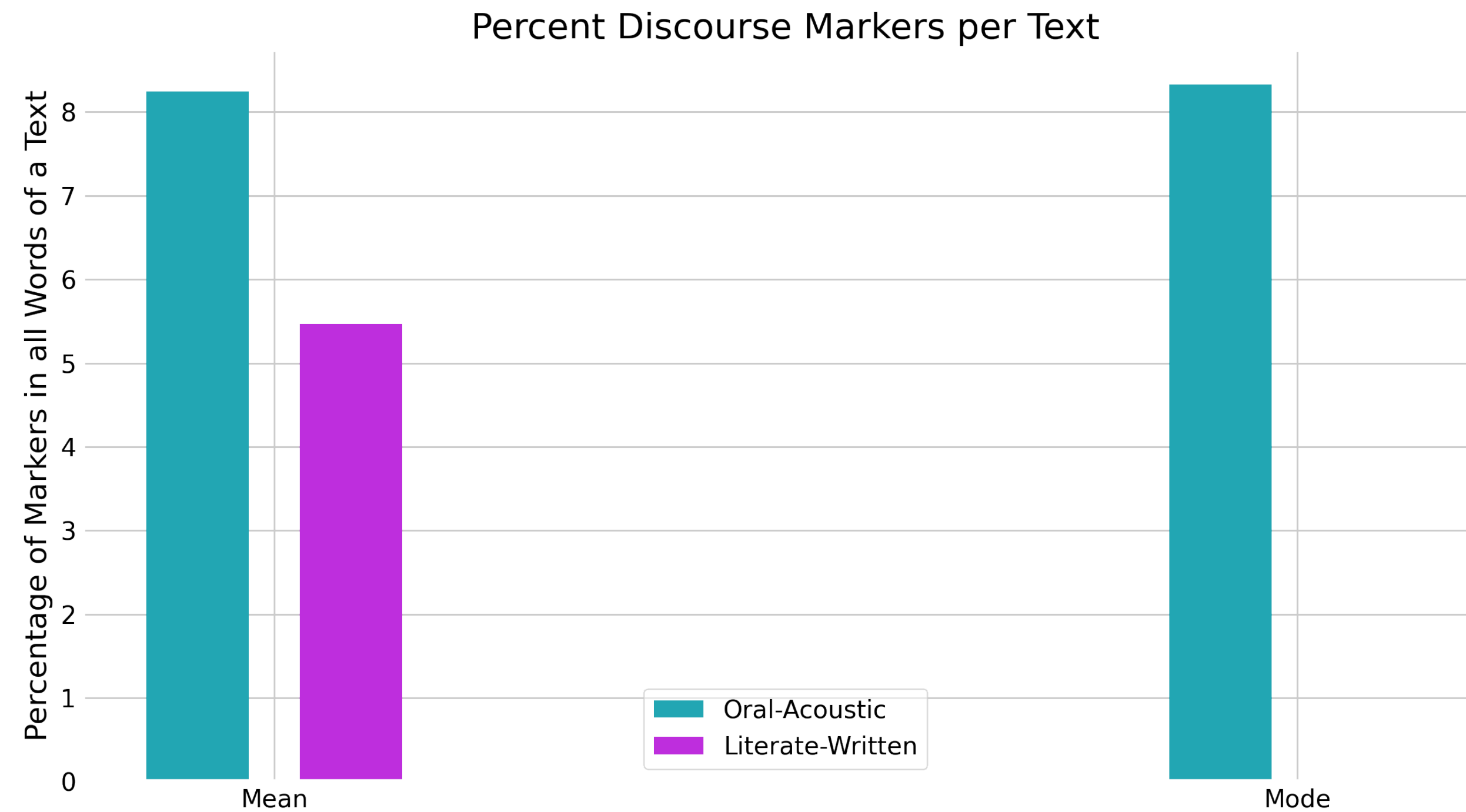
FRAGEN

1. Welche Textsorten stützen sich besonders auf Diskursmarker?
2. An welchen Positionen im Text stützen sich die jeweiligen Textsorten besonders auf Diskursmarker?
3. An welchen Positionen im Satz stützen sich die jeweiligen Textsorten besonders auf Diskursmarker?
4. Auf welche Klassen von Diskursmarkern stützen sich die jeweiligen Textsorten besonders?
5. Welche Diskursmarker werden innerhalb der jeweiligen Klassen besonders genutzt?

AUSWERTUNG

1. Generelle Verteilung

- oral-akustische Texte nutzen mehr DM als literat-schriftliche
- improvisierte OA Texte nutzen mehr DM als geskriptete



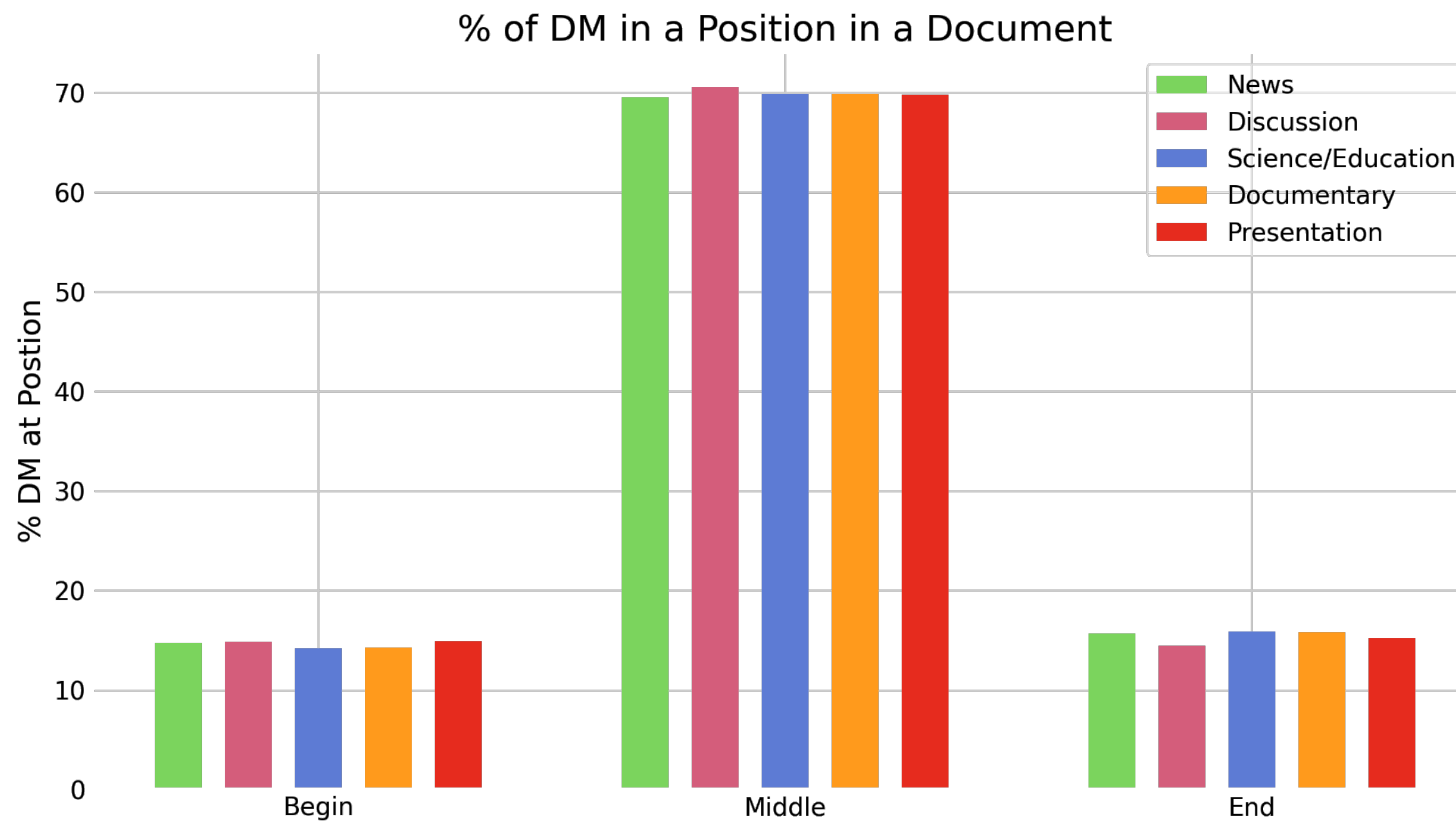
P-Wert < 0.001

» OA nutzt mit einer Effektgröße (EG)

von 1,63 mehr DM als LS

2. Textpositionen

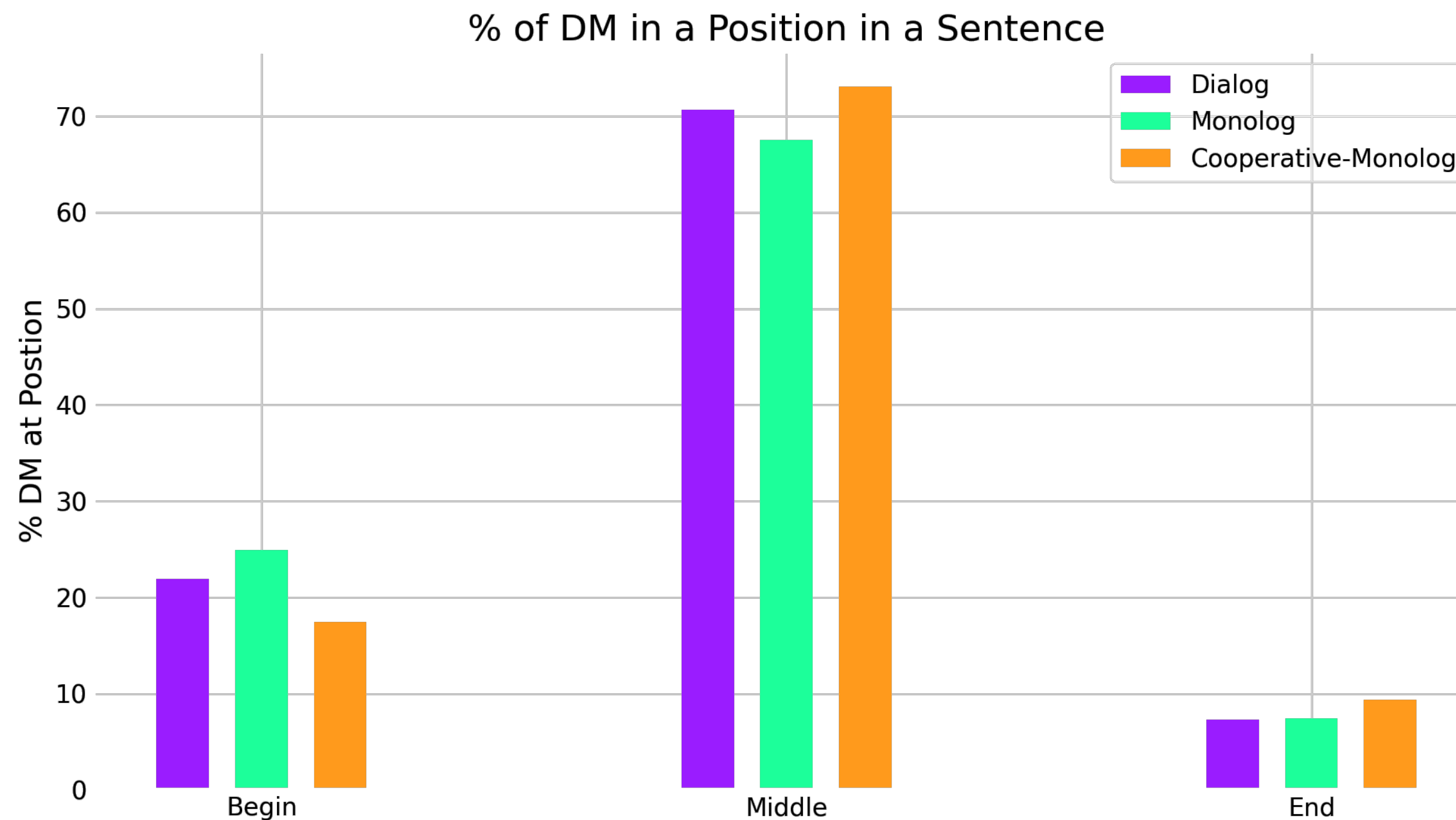
- oral-akustische Texte nutzen am Anfang und am Ende mehr DM als literat-schriftliche, LS mehr in der Mitte
- interaktionslastigere Genres (Discussion, Presentation) nutzen am Anfang und in der Mitte mehr DM, sachlichere (Science, Documentary) am Ende
- improvisierte (Dialoge) und interaktive (Reden) Texte nutzen am Anfang mehr DM, geskriptete (Koop. Monolog) am Ende



Discussion und Presentation
nutzen am Anfang und in der Mitte
mehr DM,
Documentary und Science am Ende

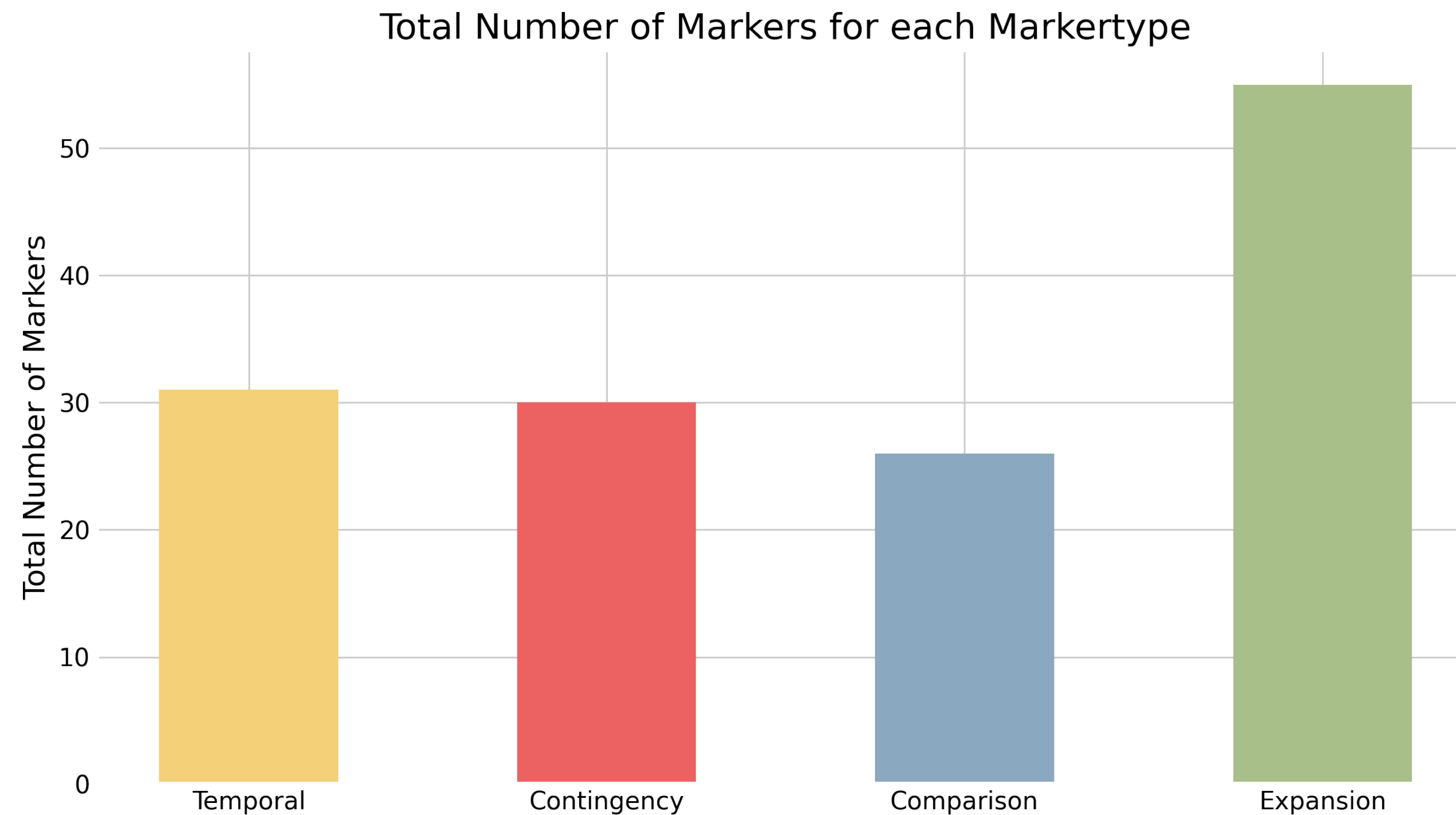
3. Satzpositionen

- oral-akustische Texte nutzen am Anfang mehr DM als literarische, LS mehr in der Mitte
- sachliche Genres (News, Science) nutzen am Anfang mehr DM und weniger in der Mitte
- interaktionslastigere Texte (Dialoge) nutzen am Anfang mehr DM und am Satzende weniger, passive (Koop. Monolog) mehr am Satzende



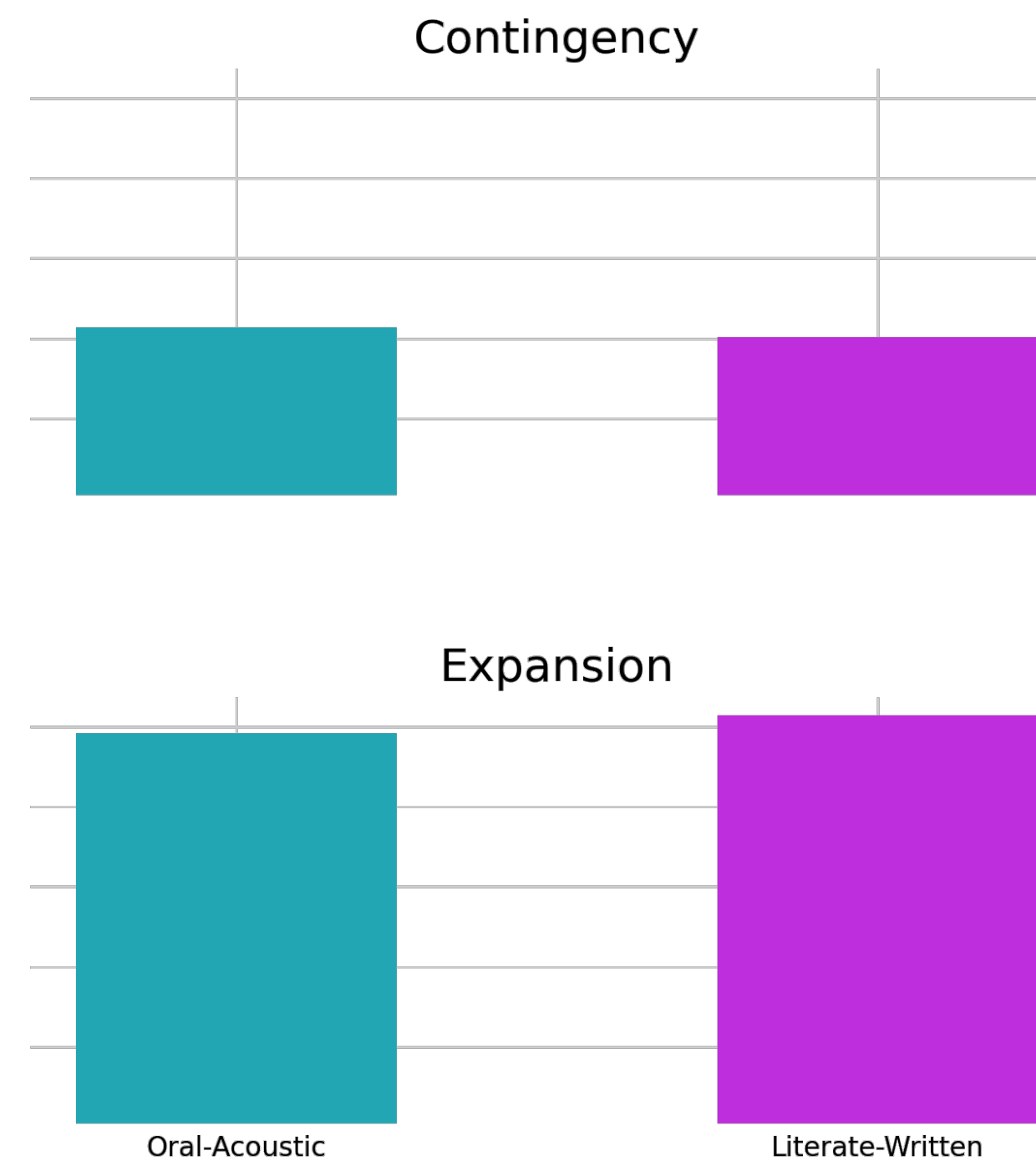
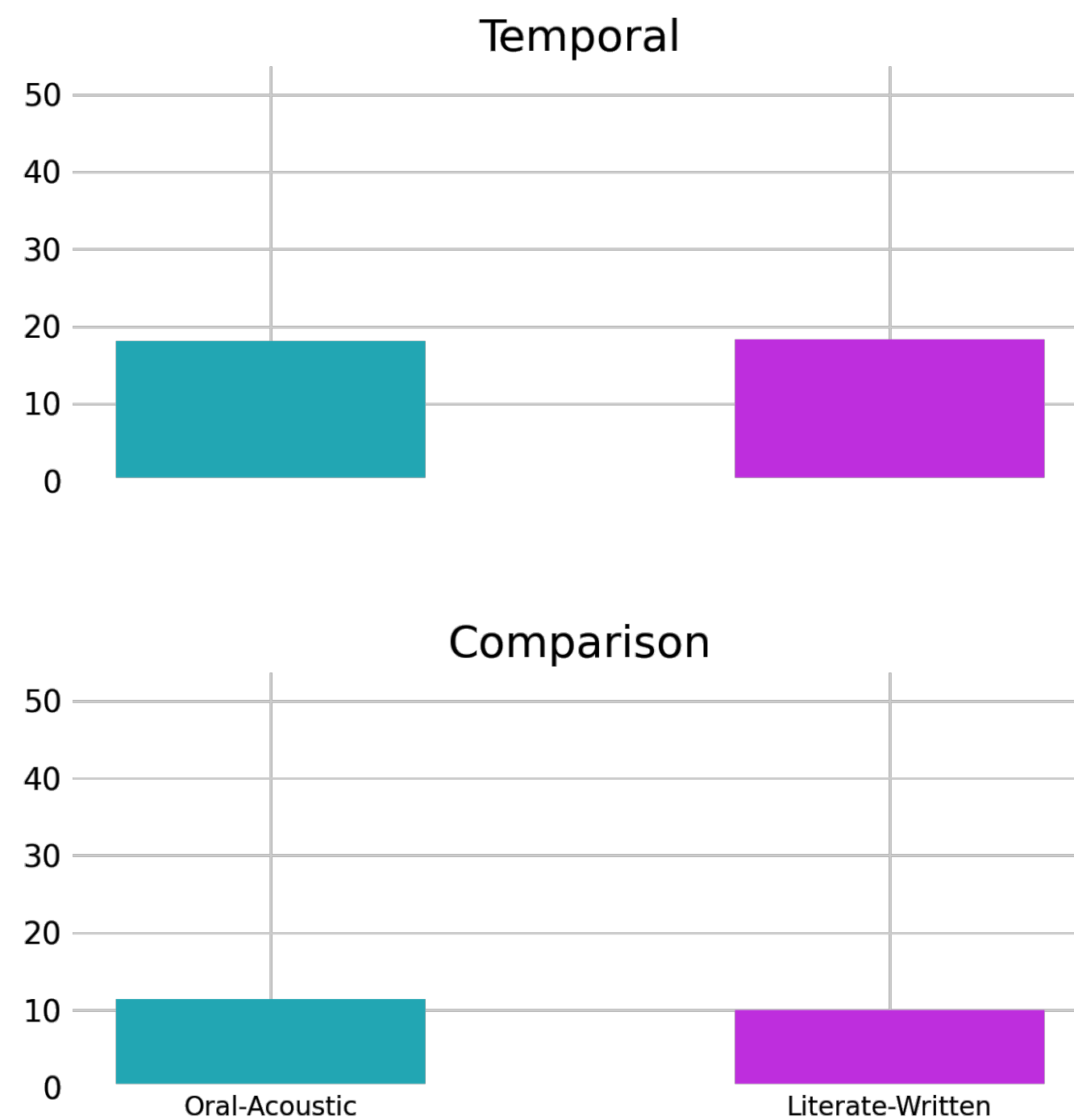
Dialoge und Monologe nutzen am
Satzanfang mehr DM,
Kooperative Monologe dafür am
Satzende

4. Diskursmarker-Klassen



größter Anteil an allen DM in
allen Texten: **EXPANSION** DM
(u.a. *and*)

Share of Markertypes in all Markers of a Texttype (%)



OA Texte nutzen mehr

COMPARISON und CONTINGENCY

DM als LS Texte

5. Häufigste Diskursmarker der Klassen

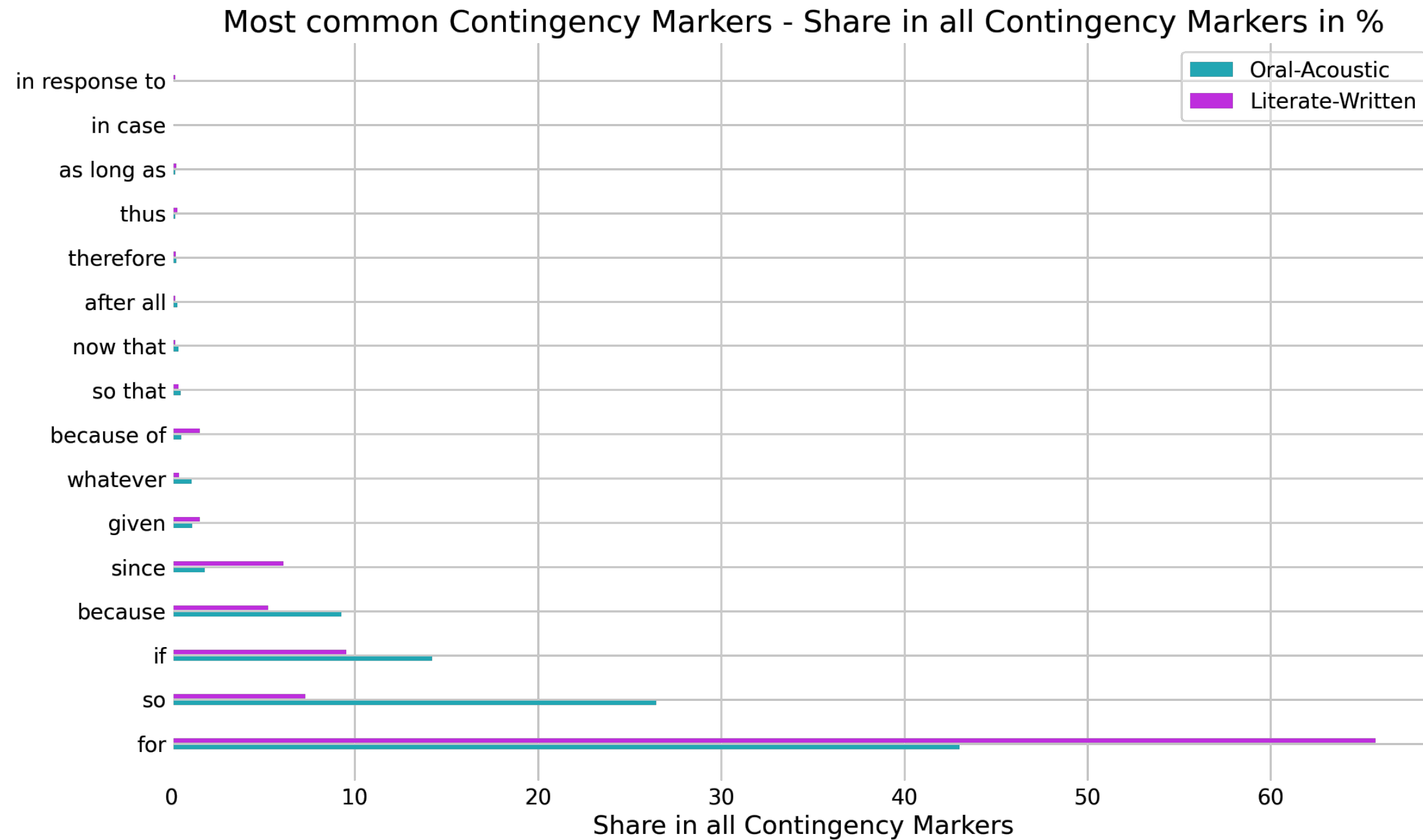
as TEMPORAL

and EXPANSION

but COMPARISON

so, for CONTINGENCY

- umgangssprachlichere DM eher in improvisierten OA Texten
- gehobener DM eher in LS Texten und geskripteten OA Texten



3 häufigste DM:

> 80% aller Contingency DM

häufiger in OA:

because, now that, whatever

häufiger in LS:

in response to, since

Offene Fragen

- Positionen der DM-Klassen im Satz
- Aussortieren von Homographen
- alternative DM-Liste nutzen
- Listenability aller betrachteten Texte messen
und mit DM Nutzung in Zusammenhang setzen

ZUSAMMENFASSUNG

- Diskursmarker sind nicht eindeutig definiert
 - wegweisende Wörter im Text, ohne inhaltlichen Aspekt
 - » Problem: Mehrdeutigkeit (Homographie)
 - Einteilung in Bedeutungsgruppen möglich
 - » Problem: Mehrdeutigkeit (mehrere Gruppen möglich)
- Diskursmarker werden in verschiedenen Textsorten unterschiedlich eingesetzt
 - oral-akustische Texte nutzen mehr Diskursmarker als literat-schriftliche
 - improvisierte Texte nutzen mehr Diskursmarker als geskriptete
 - interaktive Texte nutzen mehr Diskursmarker als passive
 - » Auswirkung auf Listenability