# Cross-Dialect Information Retrieval: Information Access in Low-Resource and High-Variance Languages

Robert Litschko,   Oliver Kraus,   Verena Blaschke,   Barbara Plank

COLING 2025 • Abu Dhabi

The 31st International Conference on Computational Linguistics

CIS

nlp

mcml
Munich Center for Machine Learning

🔍 München ("Munich")

# 🔍 München ("Munich")



Wikipedia
https://de.wikipedia.org › wiki › München ⋮

## München

Sie ist mit gut 1,5 Millionen Einwohnern die bevölke
Gemeinde Deutschlands und mit 4.861 Einwohnern
Geschichte Münchens · Altstadt (München) · Landk

# 🔍 München ("Munich")

What about culture-specific knowledge that can often be found in dialect Wikis?



Wikipedia
https://de.wikipedia.org › wiki › München ⋮

**München**

Sie ist mit gut 1,5 Millionen Einwohnern die bevölke...
Gemeinde Deutschlands und mit 4.861 Einwohnern
Geschichte Münchens · Altstadt (München) · Landk...

# 🔍 München ("Munich")



**Boarische Wikipedia**
https://bar.wikipedia.org › wiki · Translate this page ⋮

**Minga**

In **Minga** sogt ma München. **Minga** is mid mehra wia 1
Stod vo Bayern und hinta Berlin und Hamburg d'drittgre

# München ("Munich")



Alemannische Wikipedia
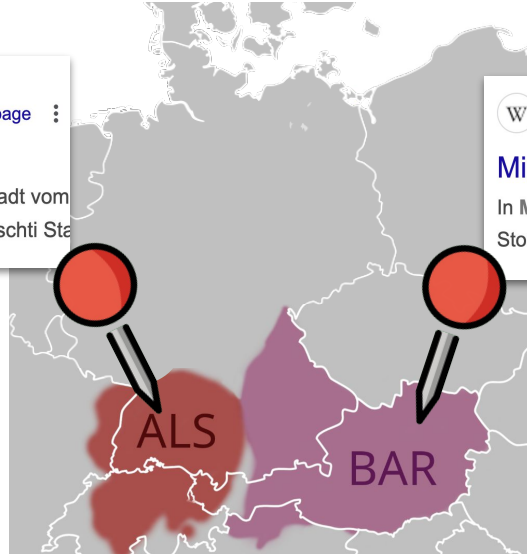https://als.wikipedia.org › München  · Translate this page  ⋮

**Münche**

**Münche** (hd. **München**, bar. Minga) isch d Hauptstadt vom
Bayern un mit über 1,4 Millione Iiwohner au die gröschti Sta

Boarische Wikipedia
https://bar.wikipedia.org › wiki  · Translate this page  ⋮

**Minga**

In **Minga** sogt ma München. **Minga** is mid mehra wia 1
Stod vo Bayern und hinta Berlin und Hamburg d'drittgre

ALS

BAR

München ("Munich")

Alemannische Wikipedia
https://als.wikipedia.org › München · Translate this page

Münche

Münche (hd. München, bar. Minga) isch d Hauptstadt vom
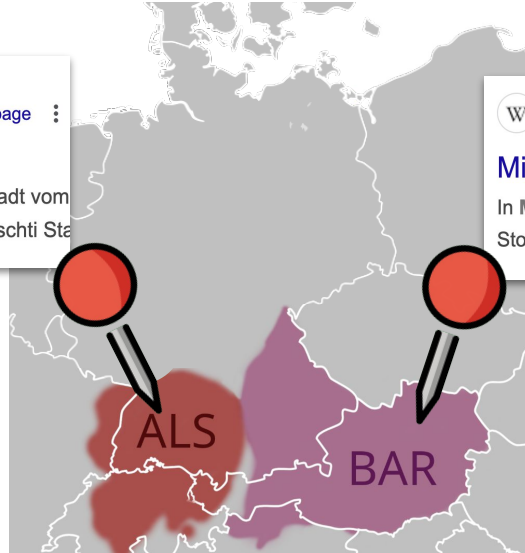Bayern un mit über 1,4 Millione Iiwohner au die gröschti Sta

Boarische Wikipedia
https://bar.wikipedia.org › wiki · Translate this page

Minga

In Minga sogt ma München. Minga is mid mehra wia 1
Stod vo Bayern und hinta Berlin und Hamburg d'drittgre

ALS

BAR

Mincke    Mincha
Minche
Müncha                Münchu
Minke    Münchè
                Münschen
Minchen    Münchä
Minga    Mìncha
        Minchä

Minchn        Minkhn
        Münch'n
Minkcha        Münchn
    Minkn
        Mingna

**High Lexical variation** due to regional word choices and different pronunciations.

Alemannische Wikipedia
https://als.wikipedia.org › München · Translate this page

Münche

**Münche** (hd. **München**, bar. Minga) isch d Hauptstadt vom
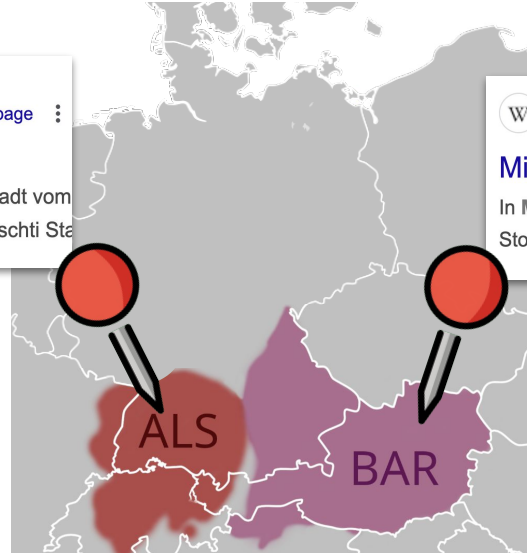Bayern un mit über 1,4 Millione Iiwohner au die gröschti Sta

Boarische Wikipedia
https://bar.wikipedia.org › wiki · Translate this page

Minga

In **Minga** sogt ma München. **Minga** is mid mehra wia 1
Stod vo Bayern und hinta Berlin und Hamburg d'drittgre

Mincke
Mincha
Minche
Müncha
Münchu
Münchè
Minke
Münschen
Minchen
Münchä
Minga
Mìncha
Minchä

ALS
BAR

Minchn
Minkhn
Münch'n
Minkcha
Münchn
Minkn
Mingna

High Lexical variation due to regional word choices and different pronunciations.

Alemannische Wikipedia
https://als.wikipedia.org › München · Translate this page

**Münche**

**Münche** (hd. **München**, bar. Minga) isch d Hauptstadt vom
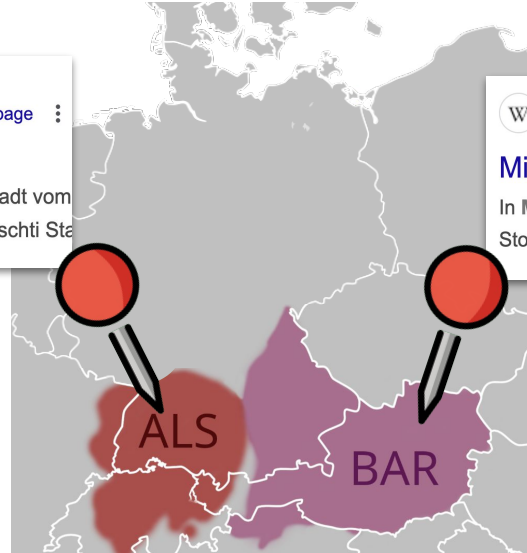Bayern un mit über 1,4 Millione Iiwohner au die gröschti Sta

Boarische Wikipedia
https://bar.wikipedia.org › wiki · Translate this page

**Minga**

In **Minga** sogt ma München. **Minga** is mid mehra wia 1
Stod vo Bayern und hinta Berlin und Hamburg d'drittgre

Mincke   Mincha
Minche
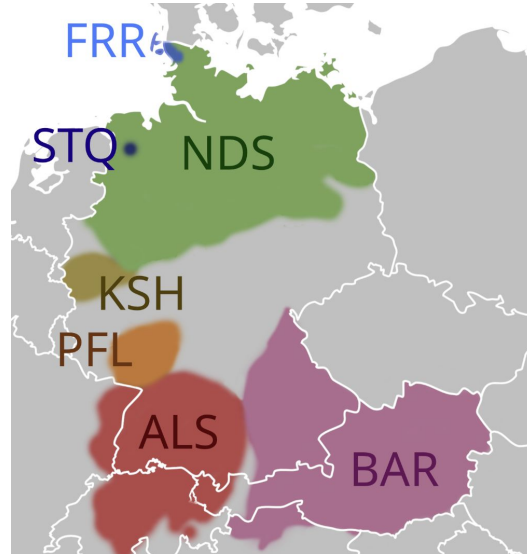Müncha   Münchu
Münchè
Minke   Münschen
Minchen   Münchä
Minga   Mìncha
Minchä

Minchn   Minkhn
Münch'n
Minkcha
Münchn
Minkn
Mingna

Lexical retrieval falls short: Normalizers do not exists for most dialects.

icons by OpenMoji - CC BY-SA 4.0

9

💡 High Lexical variation due to regional word choices and different pronunciations.



Low German (**nds**)
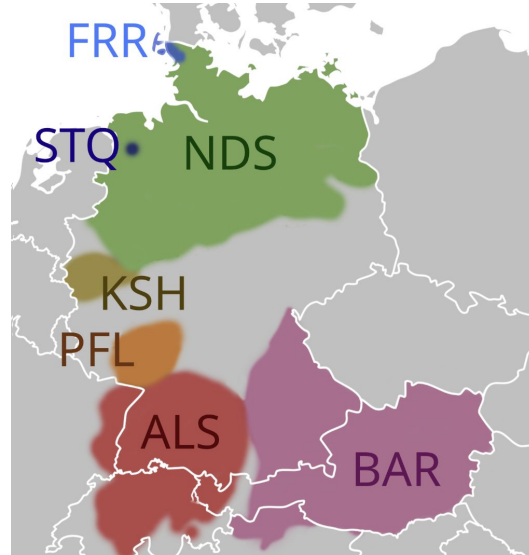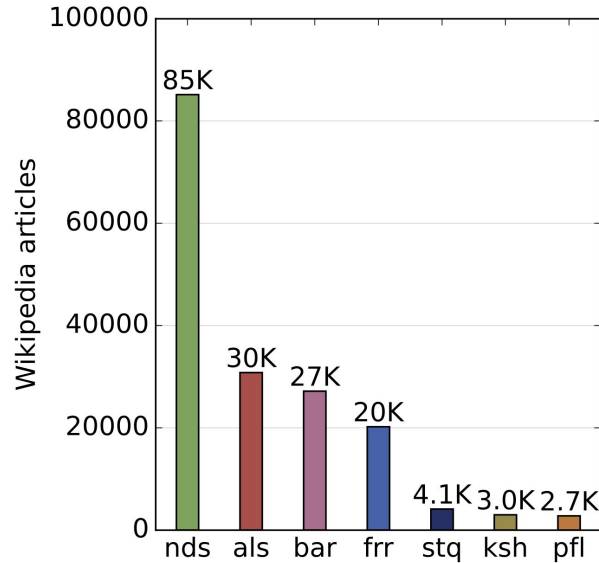Alemannic (**als**)
Bavarian (**bar**)
North Frisian (**frr**)
Saterfrisian (**stq**)
Ripuarian (**ksh**)
Rhine Franconian (**pfl**)

💡 **High Lexical variation** due to regional word choices and different pronunciations.
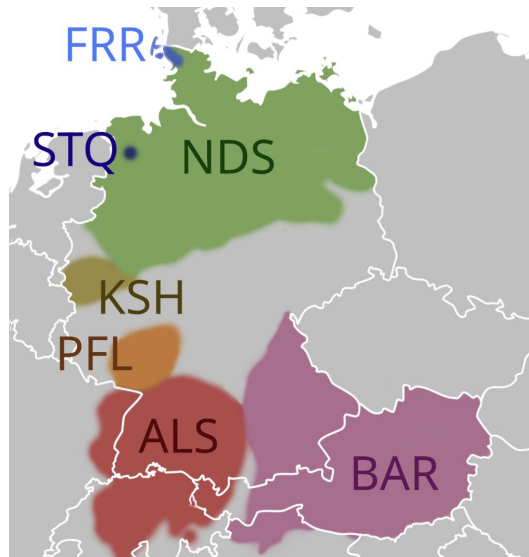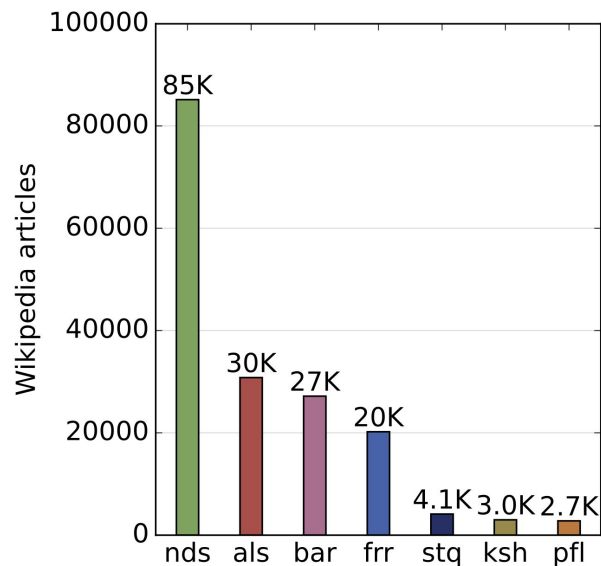


Low German (**nds**)
Alemannic (**als**)
Bavarian (**bar**)
North Frisian (**frr**)
Saterfrisian (**stq**)
Ripuarian (**ksh**)
Rhine Franconian (**pfl**)

Standard German: 2.9M Wiki articles

**High Lexical variation** due to regional word choices and different pronunciations.

Low German (**nds**)
Alemannic (**als**)
Bavarian (**bar**)
North Frisian (**frr**)
Saterfrisian (**stq**)
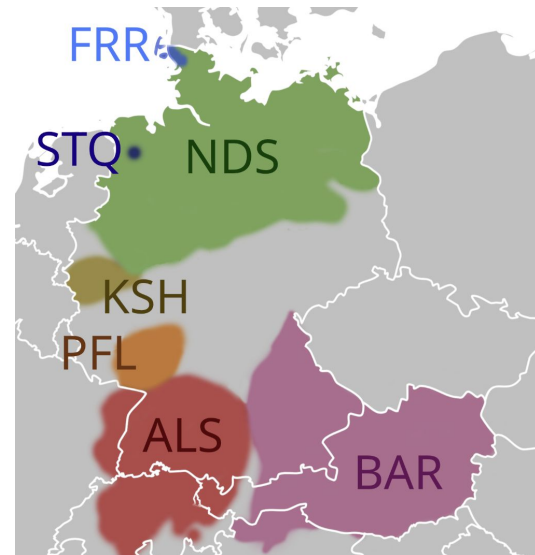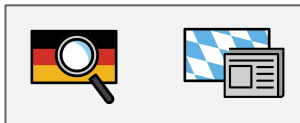Ripuarian (**ksh**)
Rhine Franconian (**pfl**)

**Low-resource**: Very Limited resources data to train neural retrieval models.

# Contribution

- New task: Cross-dialect information retrieval

- New dataset: WikiDIR

- Dialect variation dictionaries

- Evaluation of IR models on WikiDIR

**Example**

# Agenda

1. **Motivation**

2. WikiDIR dataset

3. Dialect dictionaries

4. Models

5. Results

# Agenda

1. Motivation

2. **WikiDIR dataset**

3. Dialect dictionaries

4. Models

5. Results

# Dataset Pipeline



Wikipedia
Dump of
dialect Y

## Minga

# Dataset Pipeline

Wikipedia Dump of dialect Y

titles

Query $q_i$

## Minga

文A 190 Sproochen ⌄

Leesen | Werkeln | Am Gwëntext werkeln | Gschicht åhschaun | Sunstigs ⌄

Der Artikl is im Dialekt **Mingarisch** gschriem worn.
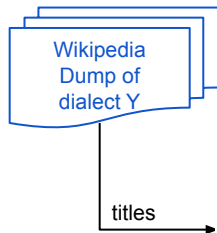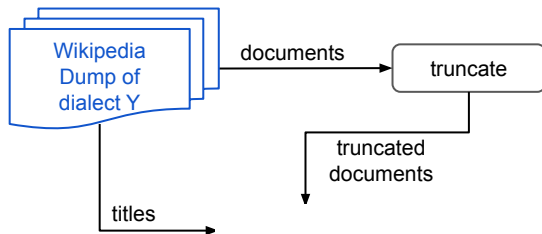
Fia andane Bedeitunga schau: Minga (Begriffsklearung).

**Minga** (amtli: **München**) Aussproch: [ˈmɪŋ(ː)ɐ] is d'Haptstod vo Bayern. In da Umgebung (20–30 km) hoaßt ma s'Minga oda oft aa oafach *d'Stod*. In Minga sogt ma München. Minga is mid mehra wia 1,5 Milliona Eihwohna d'gresste Stod vo Bayern und hinta Berlin und Hamburg d'drittgresste Stod vo Deitschland.[2] D'Stod g'head zua d'wichtigstn Wirtschofts-, Vakeas- und Kuituazentren vo Eiropa. Minga is aa da Vawoitungssitz vom Regiarungsbeziak Owabayern.

Minga is in da ganzen woid aa zwengs da Wiesn und am Hofbraihaus bekannt. Dazua hods no vui andane Sengswiadigkeitn wias Glocknspui am Rathaus am Marienplotz, d'Residenz und s'Schloss Nymphenburg. Z'Minga gibt 's aa an Hauffa Museen, wias Deitsche Museum, oda d'oide, d'neie und d'Pinakothek vo da Modeane.

| Woppn | Deitschlandkoatn |
|---|---|
| | |

| Basisdotn | |
|---|---|
| Bundesland: | Bayern |
| Regiarungsbeziak: | Owabayern |

17

based on CLIRMatrix (Sun and Duh, 2018)

# Dataset Pipeline

Wikipedia Dump of dialect Y

documents → truncate

truncated documents

titles

Query $q_i$

Corpus $\mathcal{D}$

## Minga

文A 190 Sproochen ∨

Artikl    dischkrian                        Leesen    Werkeln    Am Gwëntext werkeln    Gschicht åhschaun    Sunstigs ∨

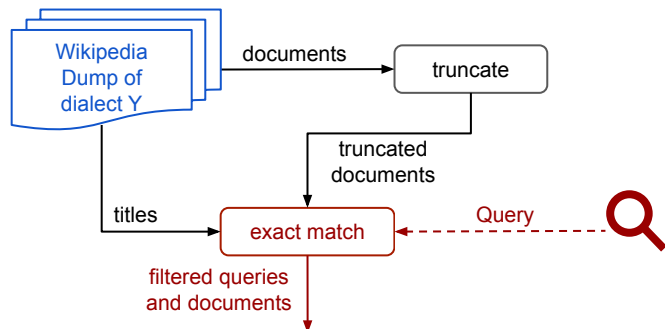Der Artikl is im Dialekt **Mingarisch** gschriem worn.

Fis andane Bedeitunga schau: Minga (Begriffsklearung)

**Minga** (amtli: **München**) Aussproch: [ˈmɪŋ(ː)ɐ] is d'Haptstod vo Bayern. In da Umgebung (20–30 km) hoaßt ma s'Minga oda oft aa oafach *d'Stod*. In Minga sogt ma München. Minga is mid mehra wia 1,5 Milliona Eihwohna d'gresste Stod vo Bayern und hinta Berlin und Hamburg d'drittgresste Stod vo Deitschland.[2] D'Stod g'head zua d'wichtigstn Wirtschofts-, Vakeas- und Kuituazentren vo Eiropa. Minga is aa da Vawoitungssitz vom Regiarungsbeziak Owabayern.

Minga is in da ganzen woid aa zwengs da Wiesn und am Hofbraihaus bekannt. Dazua hods no vui andane Sengswiadigkeitn wias Glocknspui am Rathaus am Marienplotz, d'Residenz und s'Schloss Nymphenburg. Z'Minga gibt 's aa an Hauffa Museen, wias Deitsche Museum, oda d'oide, d'neie und d'Pinakothek vo da Modeane.

| Woppn | Deitschlandkoatn |
|---|---|
| | |
| **Basisdotn** | |
| Bundesland: | Bayern |
| Regiarungsbeziak: | Owabayern |

# Dataset Pipeline



$$\mathcal{D}_{\mathrm{rel}}^{q_i} = \{d_j \in \mathcal{D} \mid d_j \text{ contains } q_i\}$$

# Dataset Pipeline

Wikipedia Dump of dialect Y

documents → truncate

truncated documents

titles

exact match ← - - - Query 🔍

filtered queries and documents

Pyserini

BM25 scores

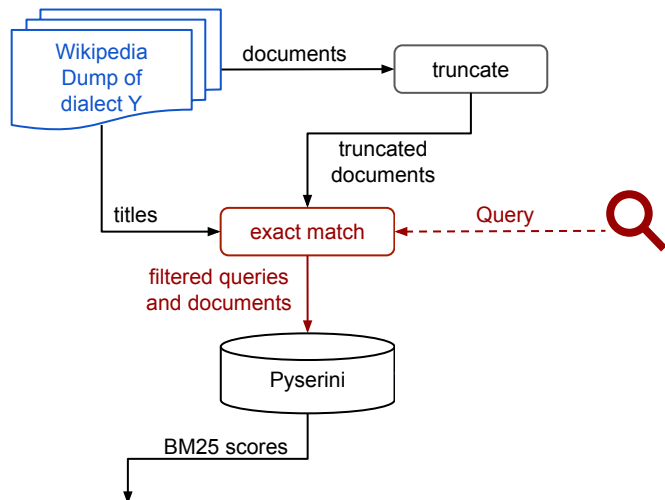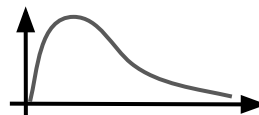Query $q_i$          Corpus $\mathcal{D}$

$$\mathcal{D}_{\text{rel}}^{q_i} = \{d_j \in \mathcal{D} \mid d_j \text{ contains } q_i\}$$

**Lexical Similarity Scores**

all (q,d)-pairs

based on CLIRMatrix (Sun and Duh, 2018)

# Dataset Pipeline



Wikipedia Dump of dialect Y

documents → truncate

truncated documents

titles

exact match ← Query

filtered queries and documents

Pyserini

BM25 scores

Jenks

Discrete Relevance Labels

Query $q_i$

Corpus $\mathcal{D}$

$$\mathcal{D}_{\text{rel}}^{q_i} = \{d_j \in \mathcal{D} \mid d_j \text{ contains } q_i\}$$

Monolingual Relevance Labels
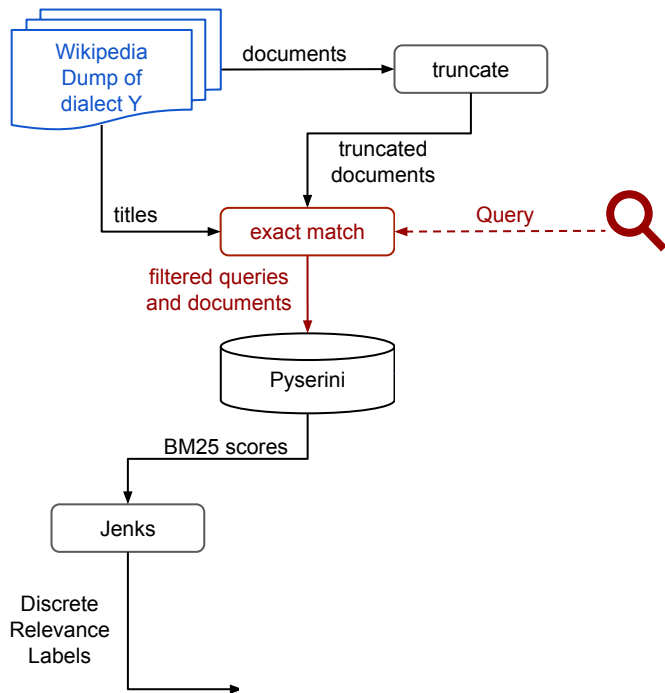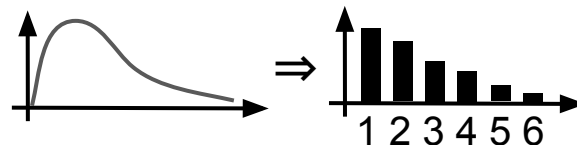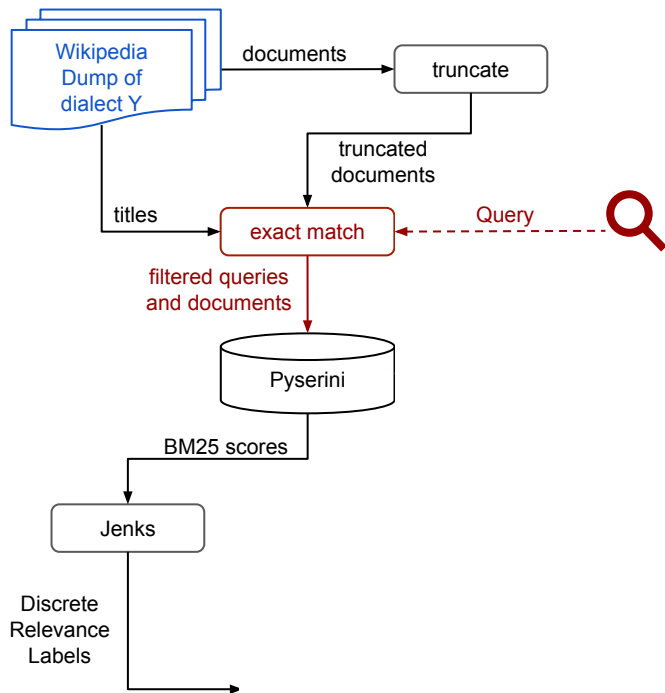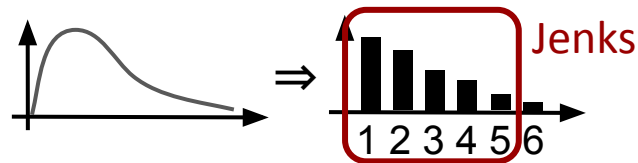
$\Rightarrow$

1 2 3 4 5 6

based on CLIRMatrix (Sun and Duh, 2018)

# Dataset Pipeline



Wikipedia Dump of dialect Y

documents → truncate

truncated documents

titles → exact match ⬸ Query 🔍

filtered queries and documents

Pyserini

BM25 scores

Jenks

Discrete Relevance Labels

Query $q_i$   Corpus $\mathcal{D}$

$$\mathcal{D}^{q_i}_{\mathrm{rel}} = \{d_j \in \mathcal{D} \mid d_j \text{ contains } q_i\}$$

Monolingual Relevance Labels

⇒   Jenks

1 2 3 4 5 6

# Dataset Pipeline



Wikipedia Dump of dialect Y

documents → truncate → truncated documents

titles → exact match ← Query

filtered queries and documents → Pyserini

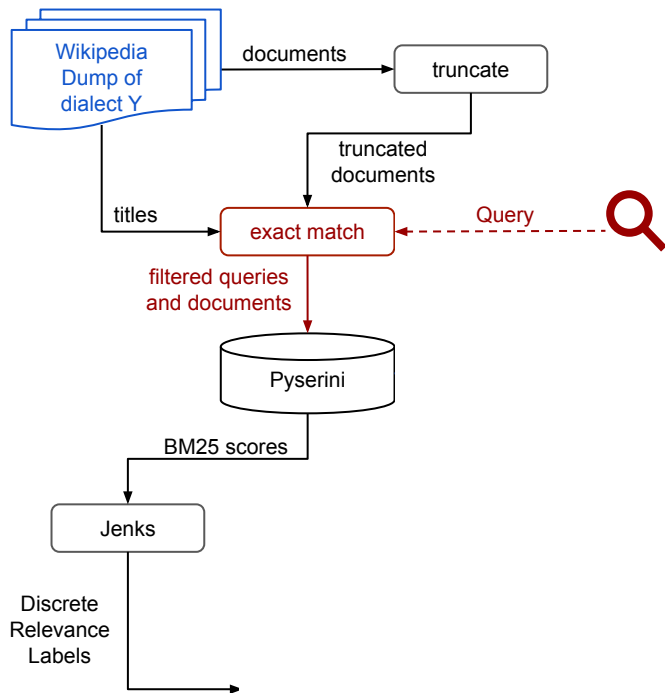BM25 scores → Jenks → Discrete Relevance Labels

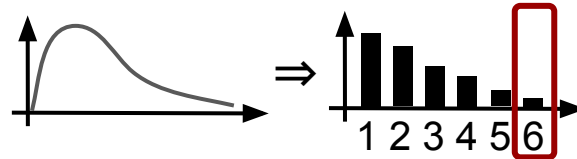Query $q_i$     Corpus $\mathcal{D}$

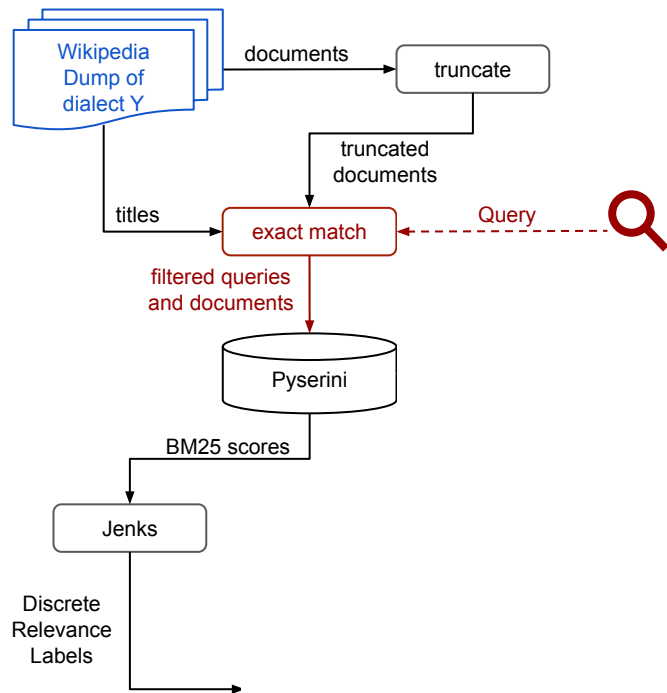$$\mathcal{D}^{q_i}_{\text{rel}} = \{d_j \in \mathcal{D} \mid d_j \text{ contains } q_i\}$$

Monolingual Relevance Labels

$\Rightarrow$    1 2 3 4 5 6    same article

based on CLIRMatrix (Sun and Duh, 2018)

23

# Dataset Pipeline



Wikipedia
Dump of
dialect Y

documents → truncate

truncated
documents

titles

exact match ← Query

filtered queries
and documents

Pyserini

BM25 scores

Jenks

Discrete
Relevance
Labels

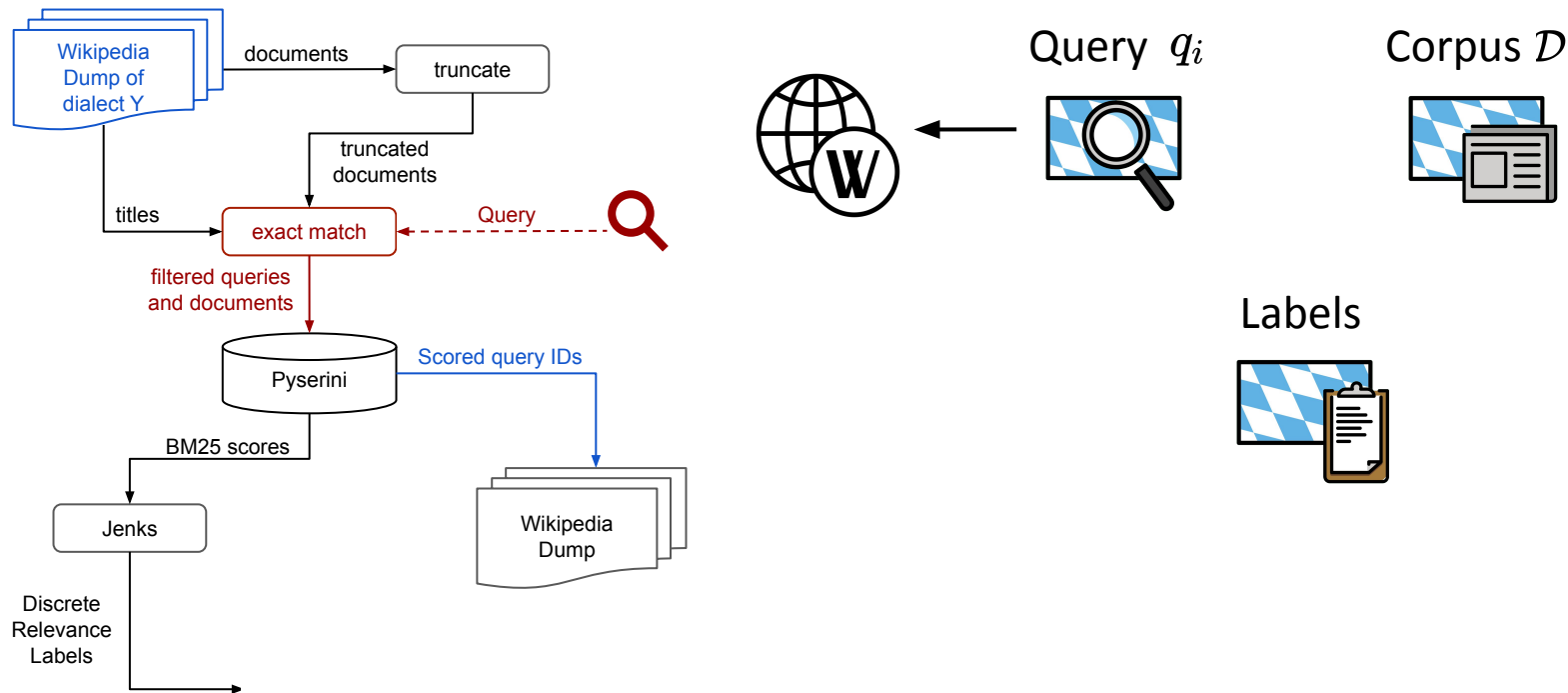Query $q_i$   Corpus $\mathcal{D}$

Labels

based on CLIRMatrix (Sun and Duh, 2018)

# Dataset Pipeline

# Dataset Pipeline

# Dataset Pipeline



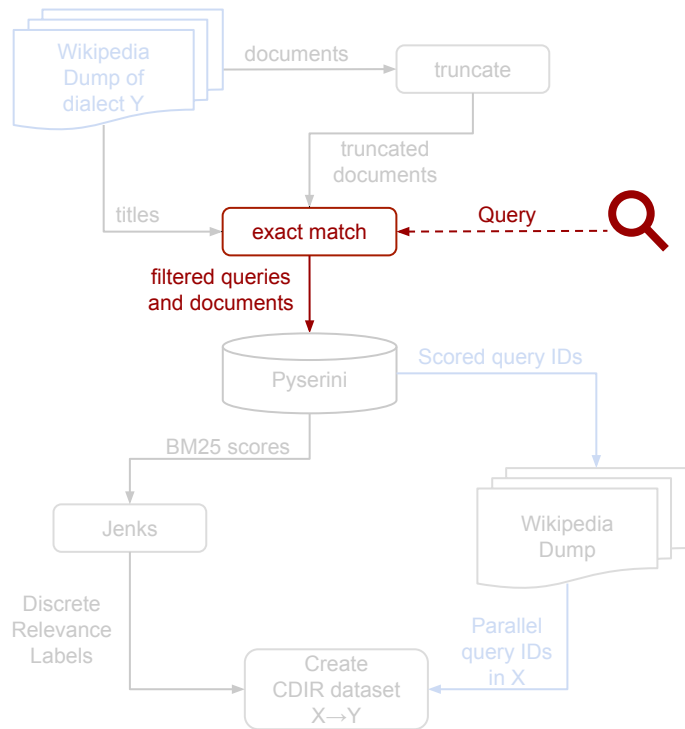based on CLIRMatrix (Sun and Duh, 2018)

# Dataset Pipeline

Wikipedia Dump of dialect Y → documents → truncate

truncated documents

titles → exact match ← Query 🔍

filtered queries and documents

Pyserini → Scored query IDs

BM25 scores

Jenks

Wikipedia Dump

Discrete Relevance Labels

Create CDIR dataset X→Y ← Parallel query IDs in X
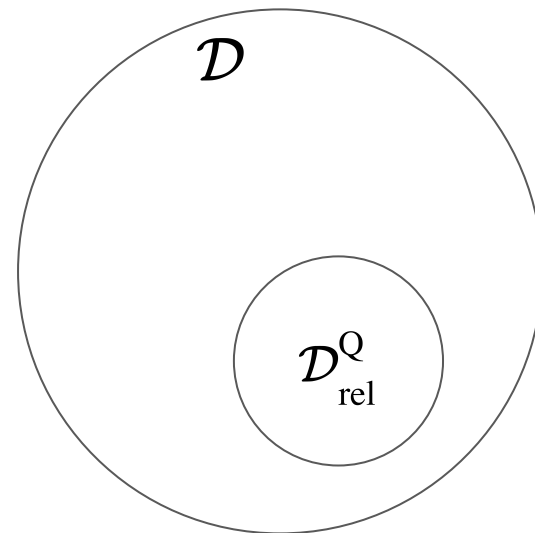
based on CLIRMatrix (Sun and Duh, 2018)

Train    Dev    Test

**Set of rel. docs.**

All documents that contain a query.

$\mathcal{D}$

$\mathcal{D}^{Q}_{rel}$

# Dataset Pipeline



Wikipedia Dump of dialect Y

documents → truncate

truncated documents

titles → exact match ← dialect variations

filtered queries and documents

Pyserini → Scored query IDs

BM25 scores

Jenks

Discrete Relevance Labels

Wikipedia Dump

Parallel query IDs in X

Create CDIR dataset X→Y

Train | Dev | Test | An.

## Analysis Split

All documents that contain a query **or any of its dialect variations.**

$\mathcal{D}$

$\mathcal{D}_{\text{rel}}^{\text{Variants}}$

$\mathcal{D}_{\text{rel}}^{\text{Q}}$

based on CLIRMatrix (Sun and Duh, 2018)

# Dataset Pipeline



Where do **dialect variations** come from?

based on CLIRMatrix (Sun and Duh, 2018)
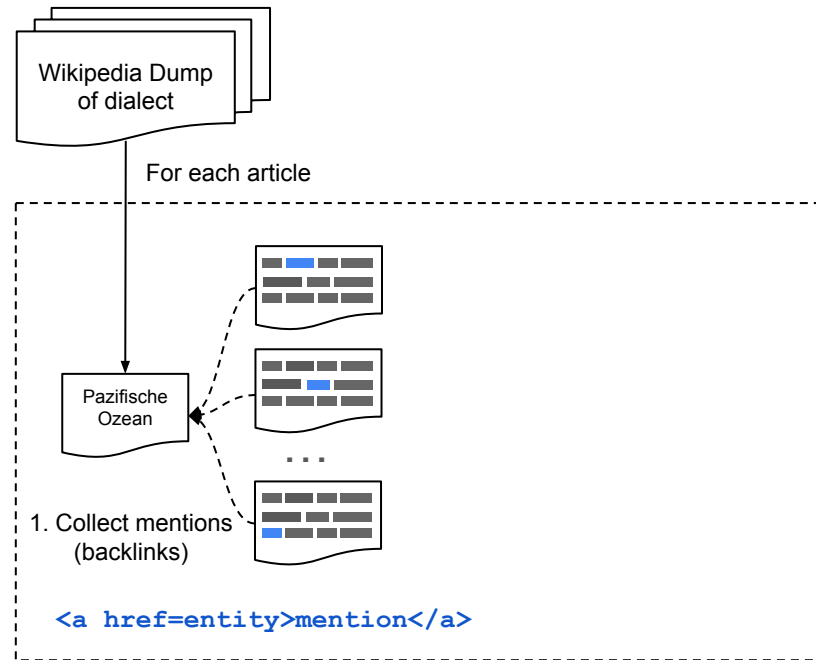
# Agenda

1. Motivation

2. WikiDIR Dataset

3. **Dialect dictionaries**
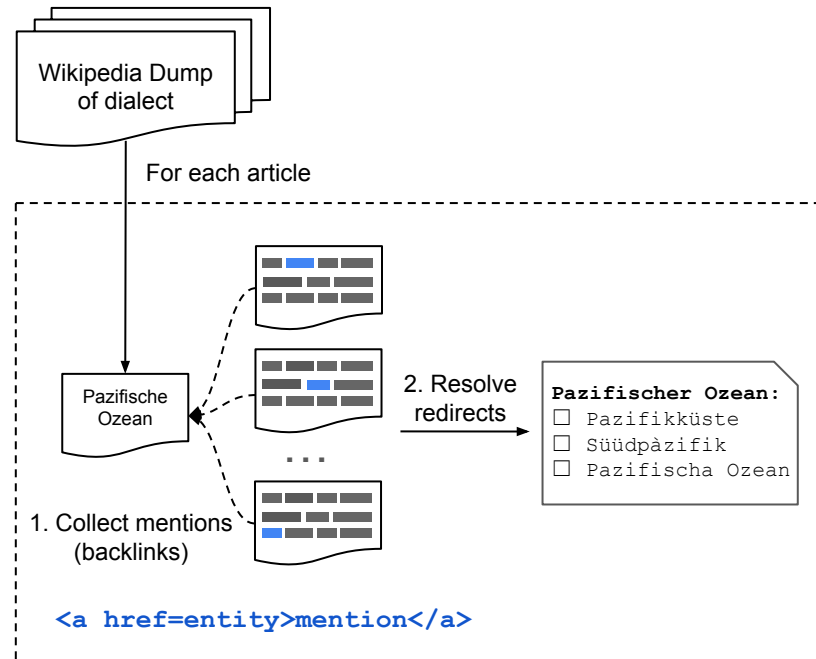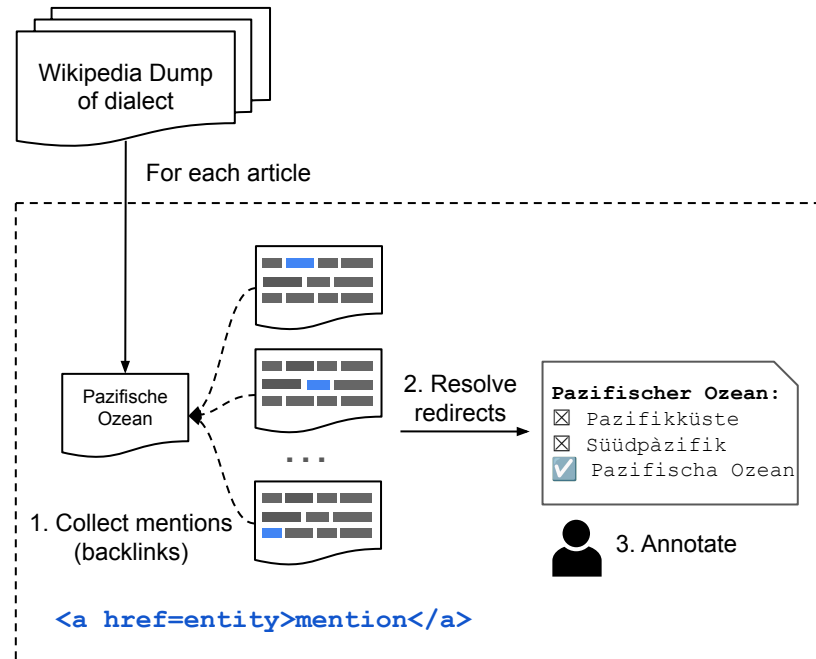
4. Models

5. Results

# Dialect variation dictionaries

Wikipedia Dump
of dialect

For each article

Pazifische
Ozean

# Dialect variation dictionaries



Wikipedia Dump
of dialect

For each article

Pazifische
Ozean

. . .

1. Collect mentions
(backlinks)

`<a href=entity>mention</a>`

# Dialect variation dictionaries

# Dialect variation dictionaries

# Dialect variation dictionaries



Wikipedia Dump of dialect

For each article

4. Add dialect variations

Pazifische Ozean

2. Resolve redirects

**Pazifischer Ozean:**
☒ `Pazifikküste`
☒ `Süüdpàzifik`
☑ `Pazifischa Ozean`

1. Collect mentions (backlinks)

. . .

3. Annotate

`<a href=entity>mention</a>`

# Example Record (Bavarian dictionary)

```
{
 "de_id": "3215",

 "de_title": "München",

 "dial_id": "12259",

 "dial_title": "Minga",

 "variants": ["Münch'n", "Minkcha", "Minkn", "Minchn", "Mingna", "Minkhn", "Münchn"]

}
```

# Agenda

1. Motivation

2. WikiDIR dataset

3. Dialect dictionaries

4. **Models**

5. Results

# Models

**Baseline:** BM25 (Robertson, 1995)

# Models



RankGPT (Llama 3.1)

============ LLM-RERANKING ============

**system:** You are RankGPT, an intelligent assistant that can rank passages based on their relevancy to the query.

**user**: I will provide you with num passages, each indicated by number identifier []. Rank them based on their relevance to query: {{query}}.
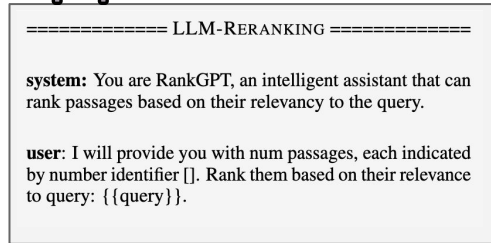
(Sun et al., 2023)
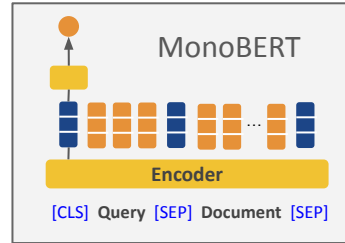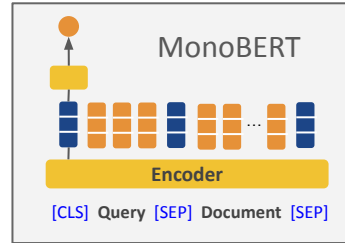
**Baseline:** BM25 (Robertson, 1995)

# Models



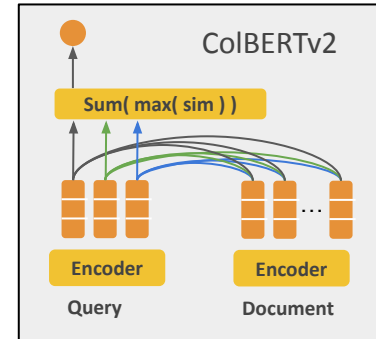RankGPT (Llama 3.1)

============ LLM-RERANKING ============

**system:** You are RankGPT, an intelligent assistant that can rank passages based on their relevancy to the query.

**user**: I will provide you with num passages, each indicated by number identifier []. Rank them based on their relevance to query: {{query}}.

(Sun et al., 2023)



MonoBERT

Encoder

[CLS] Query [SEP] Document [SEP]

(Nogueira et al., 2019)
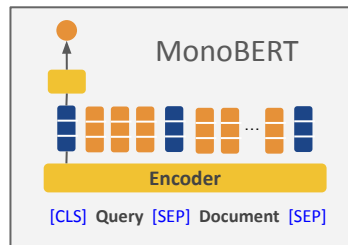
**Baseline:** BM25 (Robertson, 1995)

# Models



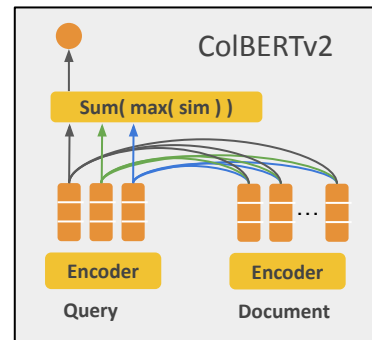RankGPT (Llama 3.1)

============ LLM-RERANKING ============

**system:** You are RankGPT, an intelligent assistant that can rank passages based on their relevancy to the query.

**user**: I will provide you with num passages, each indicated by number identifier []. Rank them based on their relevance to query: {{query}}.

(Sun et al., 2023)

MonoBERT

**Encoder**

[CLS] **Query** [SEP] **Document** [SEP]

(Nogueira et al., 2019)

ColBERTv2

**Sum( max( sim ) )**

**Encoder** **Encoder**

**Query** **Document**

(Santhanam et al., 2022)

**Baseline:** BM25 (Robertson, 1995)
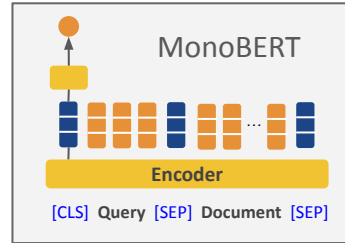
# Models

## Rerank top 100



RankGPT (Llama 3.1)
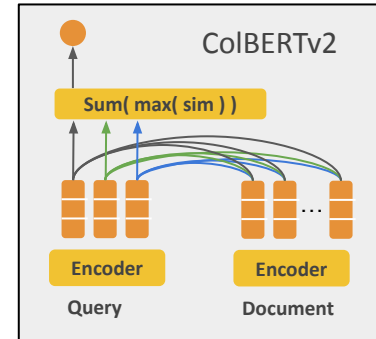
============ LLM-Reranking ============

**system:** You are RankGPT, an intelligent assistant that can rank passages based on their relevancy to the query.

**user:** I will provide you with num passages, each indicated by number identifier []. Rank them based on their relevance to query: {{query}}.

(Sun et al., 2023)

MonoBERT

Encoder

[CLS] Query [SEP] Document [SEP]

(Nogueira et al., 2019)

ColBERTv2

Sum( max( sim ) )

Encoder          Encoder

Query          Document

(Santhanam et al., 2022)

**Baseline:** BM25 (Robertson, 1995)
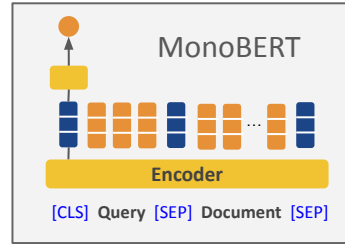
# Models

## Retrieval

RankGPT (Llama 3.1)

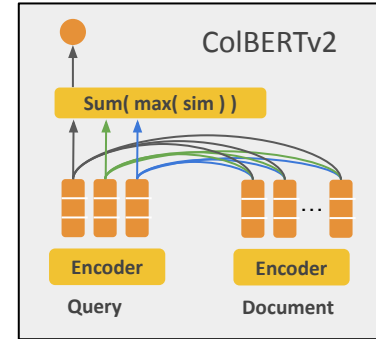============ LLM-RERANKING ============

**system:** You are RankGPT, an intelligent assistant that can rank passages based on their relevancy to the query.

**user:** I will provide you with num passages, each indicated by number identifier []. Rank them based on their relevance to query: {{query}}.
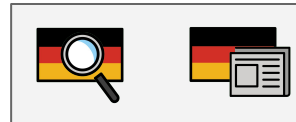
(Sun et al., 2023)

MonoBERT

Encoder

[CLS] Query [SEP] Document [SEP]

(Nogueira et al., 2019)

ColBERTv2

Sum( max( sim ) )

Encoder          Encoder

Query            Document

(Santhanam et al., 2022)

**Baseline:** BM25 (Robertson, 1995)

# Models



RankGPT (Llama 3.1)

============= LLM-Reranking =============

**system:** You are RankGPT, an intelligent assistant that can rank passages based on their relevancy to the query.

**user**: I will provide you with num passages, each indicated by number identifier []. Rank them based on their relevance to query: {{query}}.
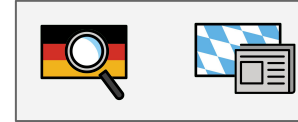
(Sun et al., 2023)

MonoBERT

Encoder

[CLS] Query [SEP] Document [SEP]

(Nogueira et al., 2019)

ColBERTv2

Sum( max( sim ) )

Encoder          Encoder

Query            Document

(Santhanam et al., 2022)
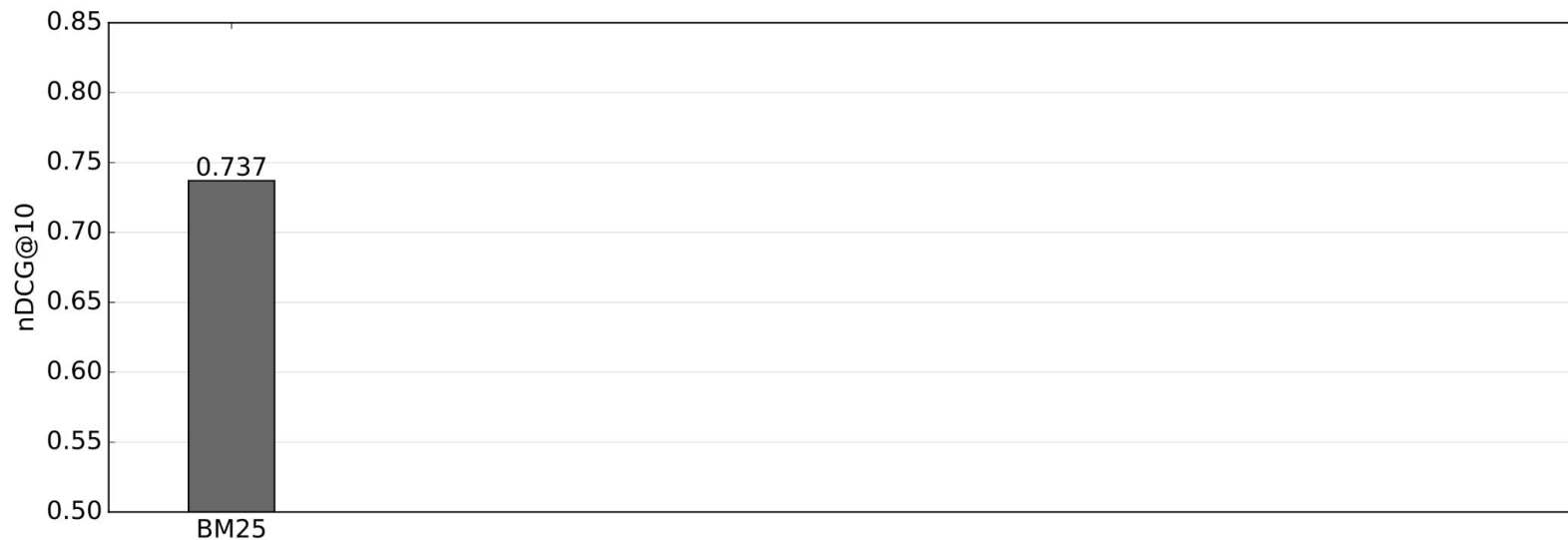
Zero-shot Transfer          Fine-tuning
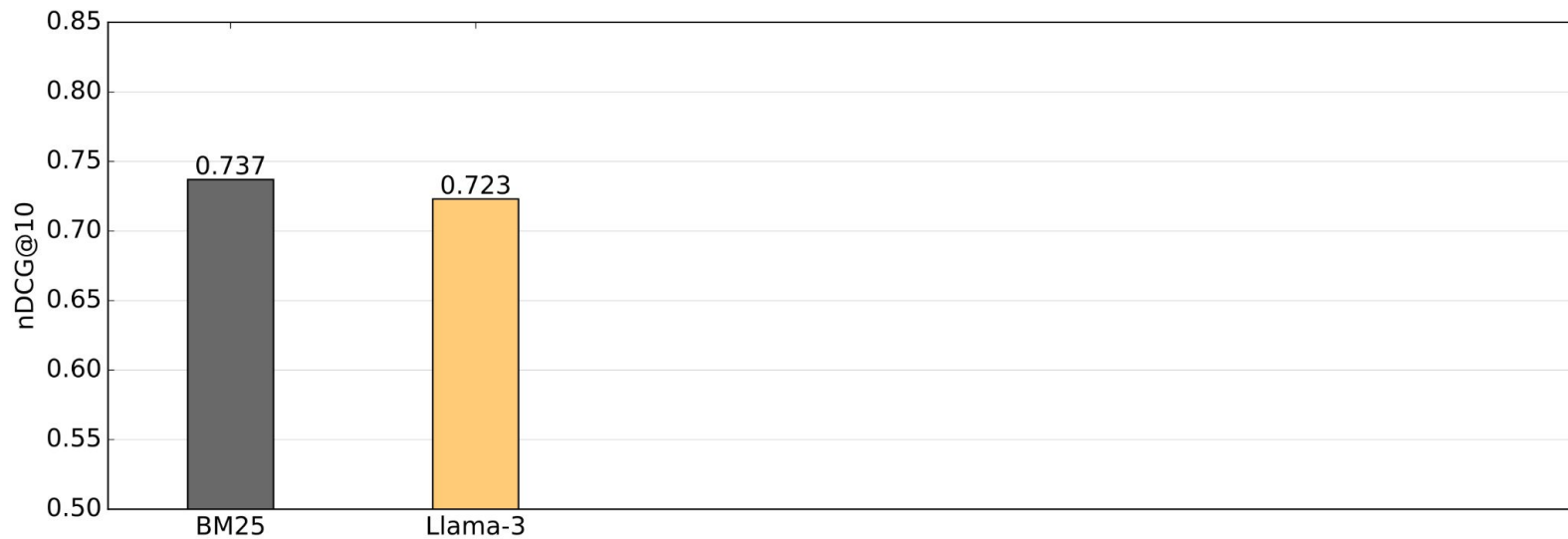
**Baseline:** BM25 (Robertson, 1995)

# Agenda

1. Motivation

2. WikiDIR dataset

3. Dialect dictionaries

4. Models

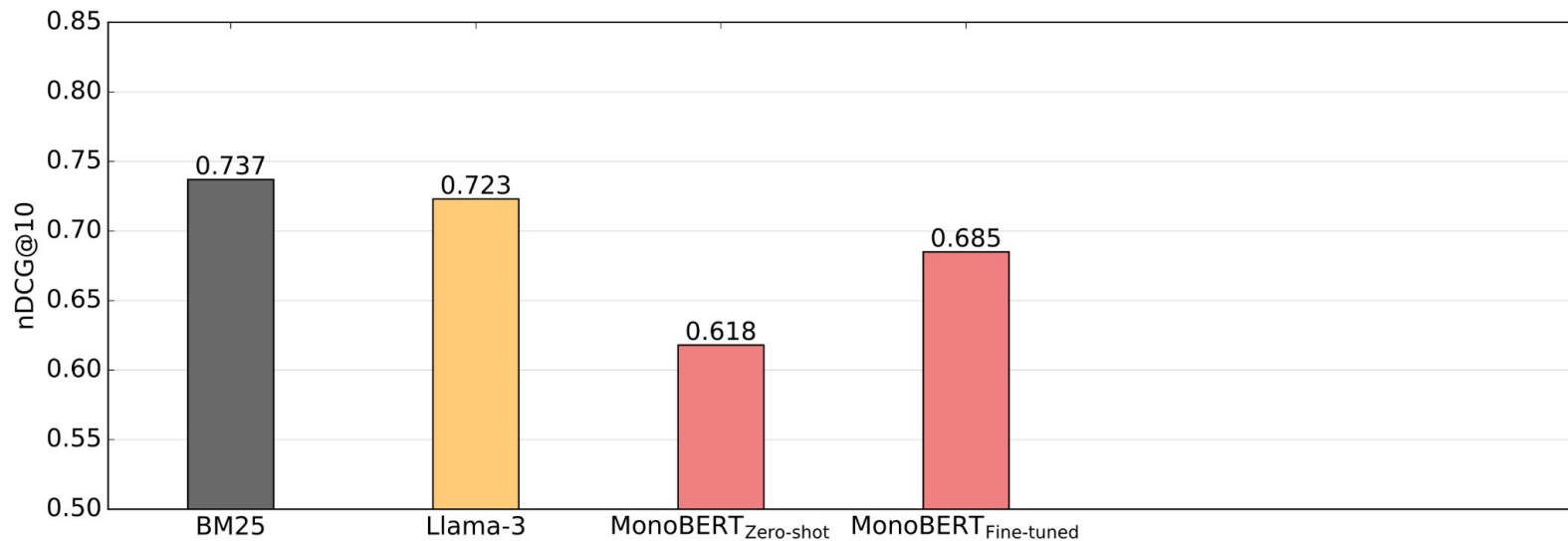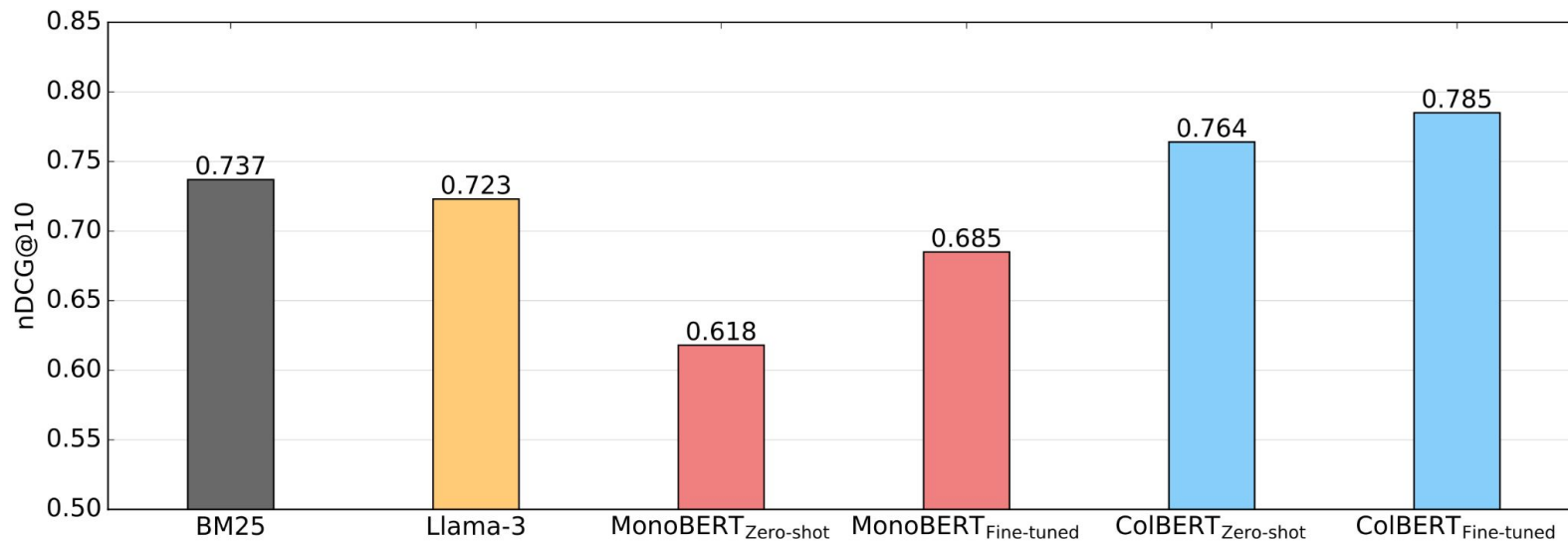5. **Results**

# Main results



*average over 7 dialects

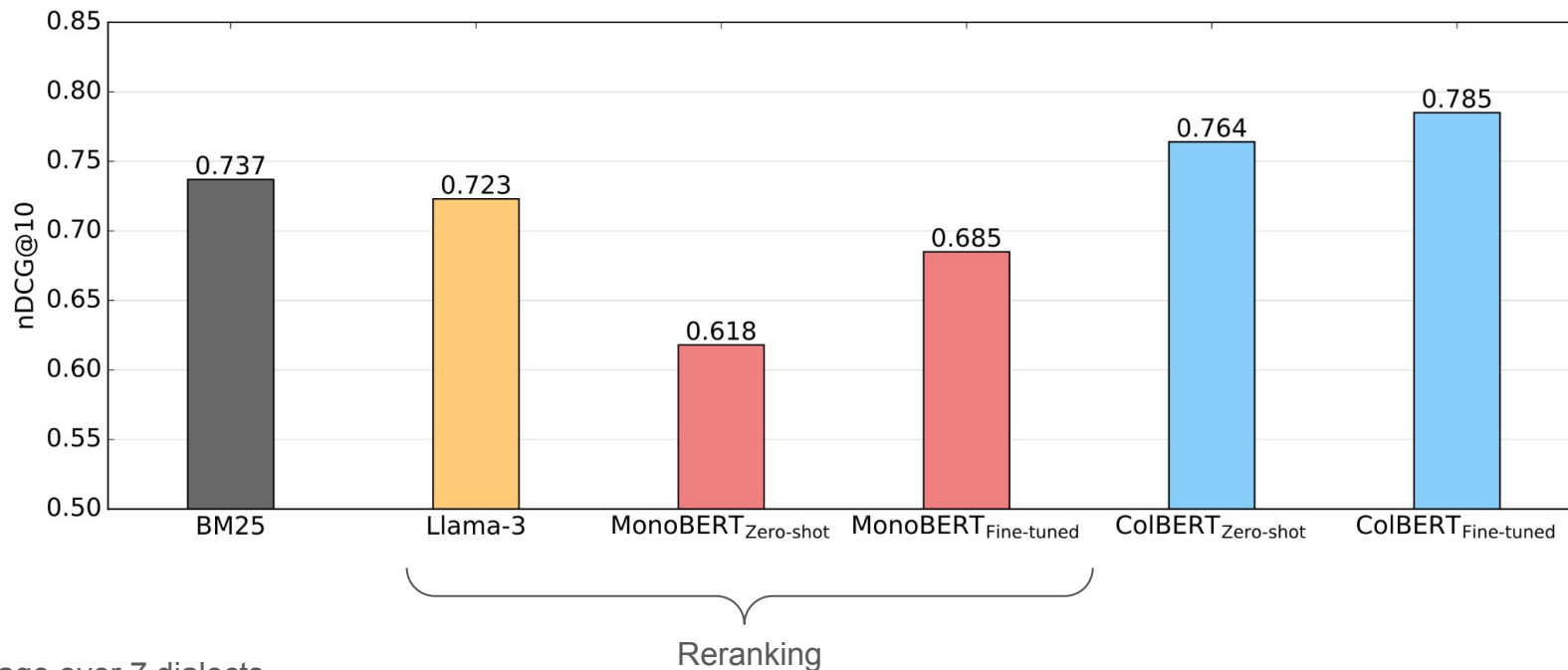# Main results



*average over 7 dialects

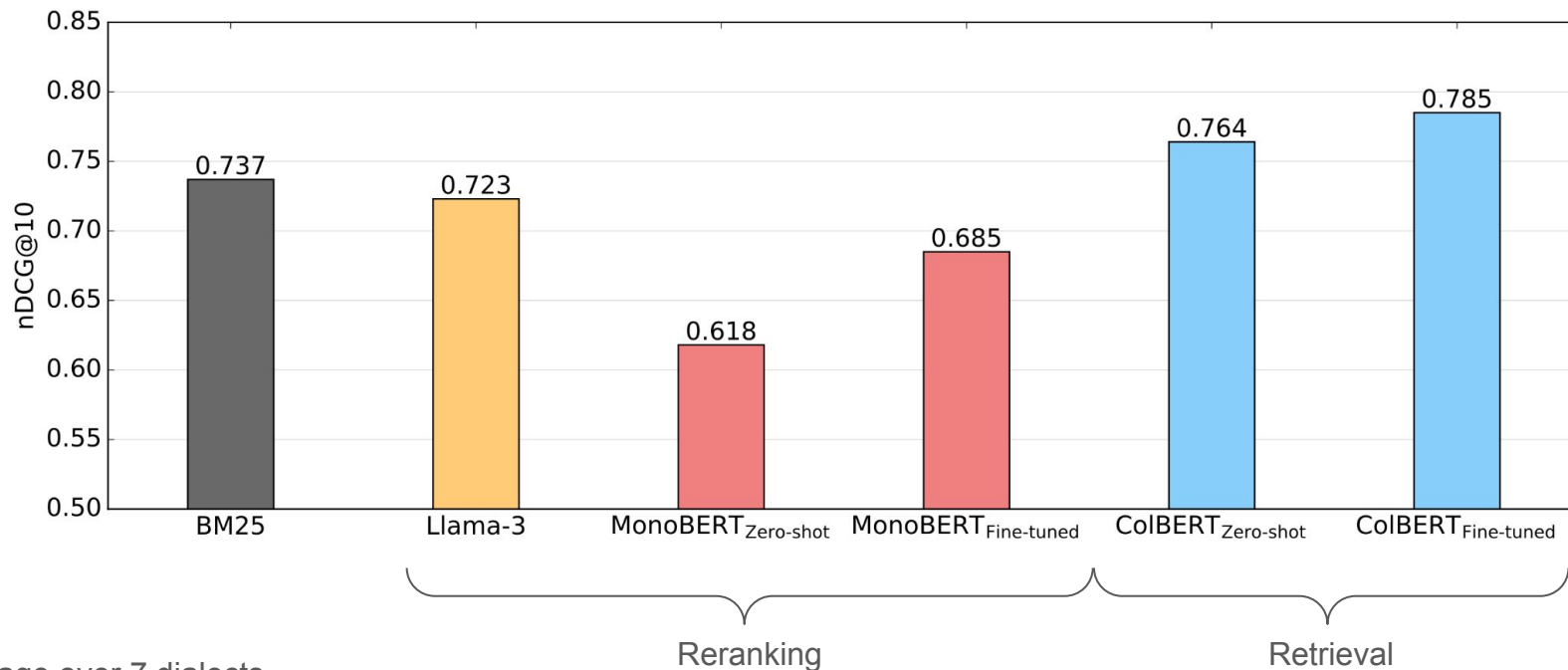# Main results



*average over 7 dialects

# Main results

*average over 7 dialects

# Main results
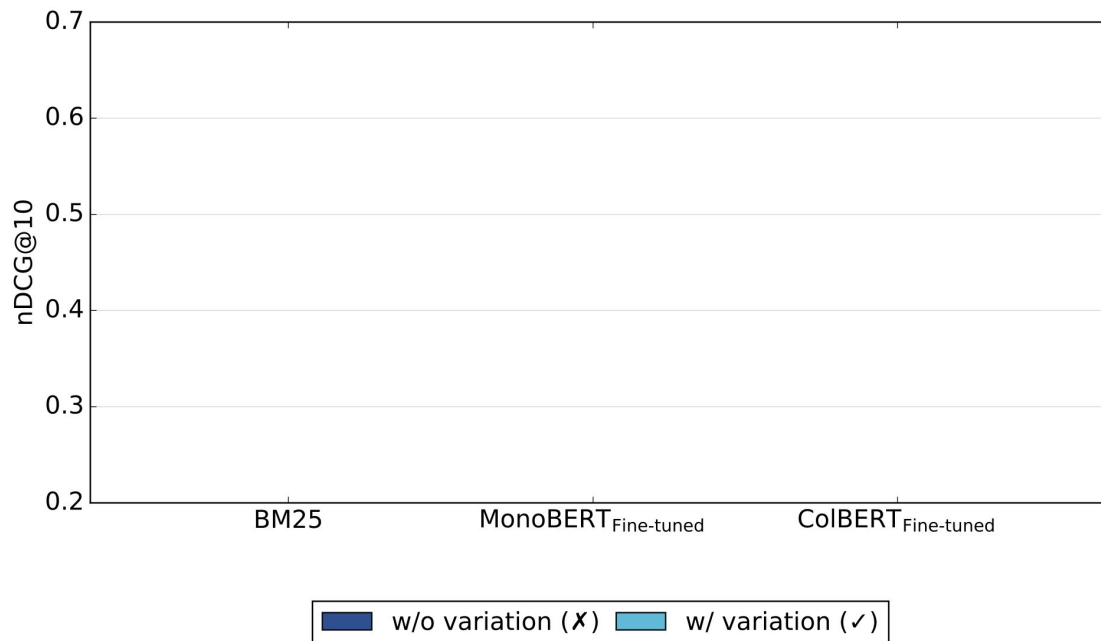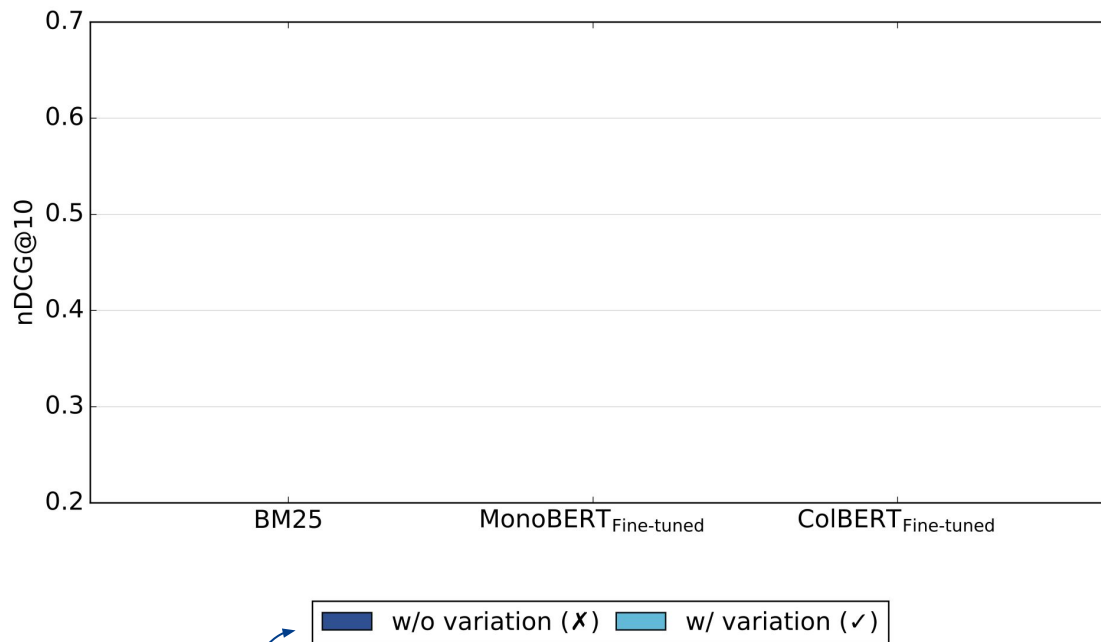
*average over 7 dialects

# Main results

*average over 7 dialects
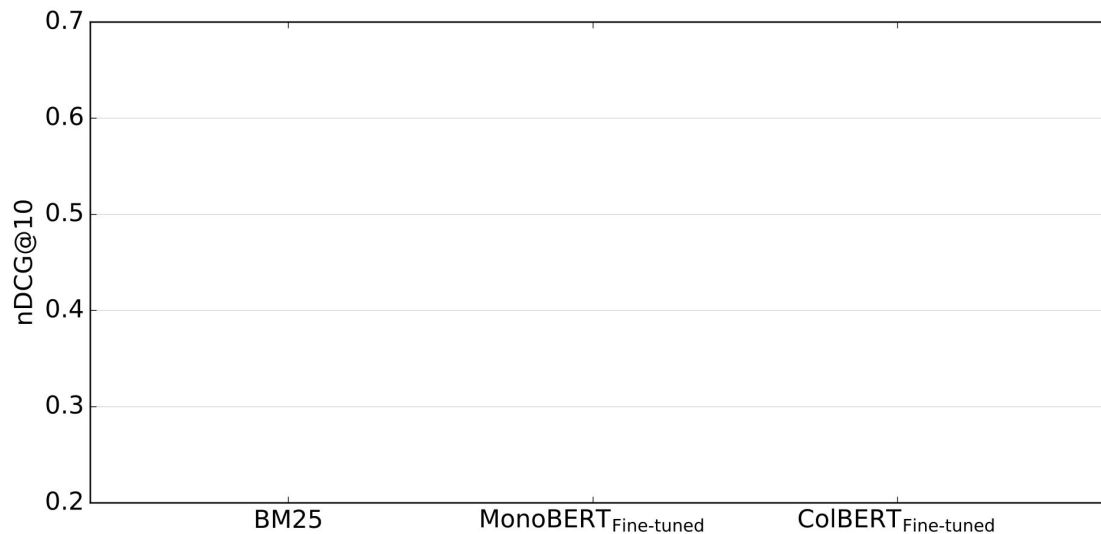
# Dialect variation results



*average over 5 dialects

# Dialect variation results



*average over 5 dialects

Exclude documents
containing variations

# Dialect variation results



nDCG@10 chart with x-axis categories: BM25, MonoBERT_{Fine-tuned}, ColBERT_{Fine-tuned}. Legend: ■ w/o variation (✗)  ■ w/ variation (✓)
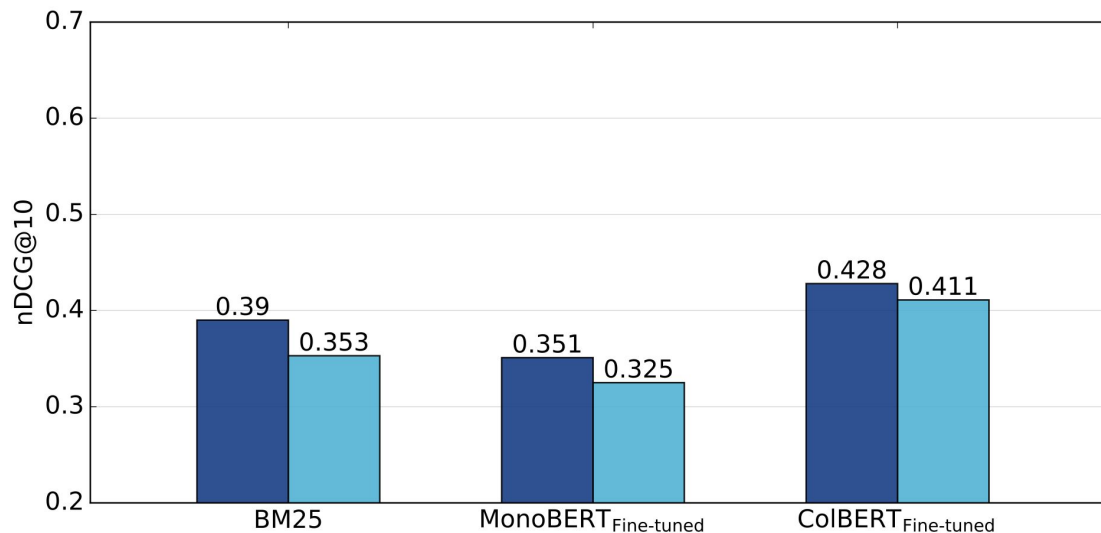
*average over 5 dialects

Exclude documents
containing variations

full analysis split

# Dialect variation results

*average over 5 dialects

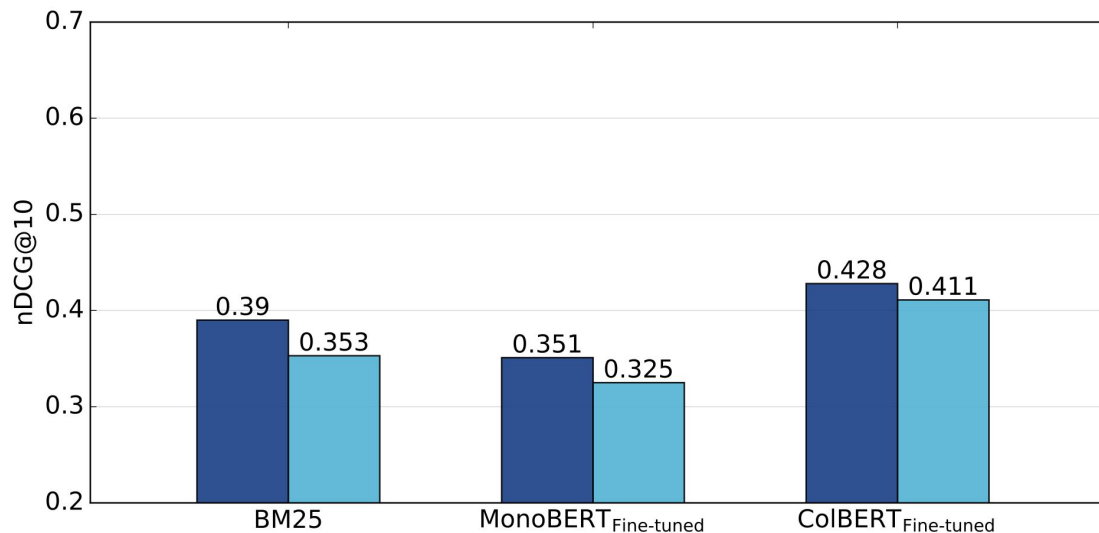Exclude documents containing variations

full analysis split

# Document translation results

Can we use LLMs to close the dialect gap?

Document transl.
Dialect → DE



*average over 5 dialects

Exclude documents containing variations
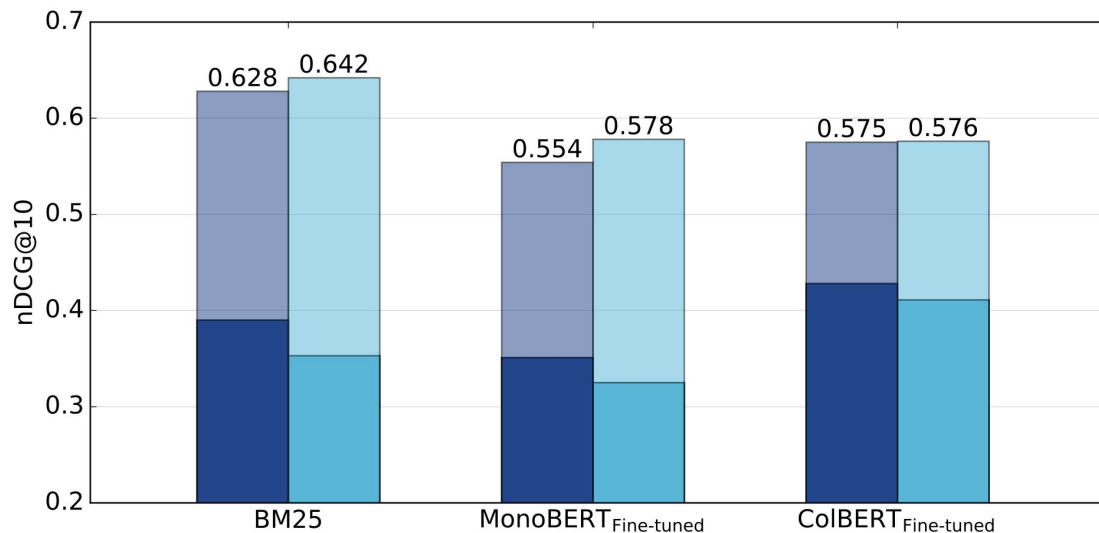
full analysis split

# Document translation results

Can we use LLMs to close the dialect gap?

Document transl.
Dialect → DE



nDCG@10

- BM25: 0.628, 0.642
- MonoBERT_{Fine-tuned}: 0.554, 0.578
- ColBERT_{Fine-tuned}: 0.575, 0.576

■ w/o variation (✗)   ■ w/ variation (✓)

Exclude documents containing variations

full analysis split

*average over 5 dialects
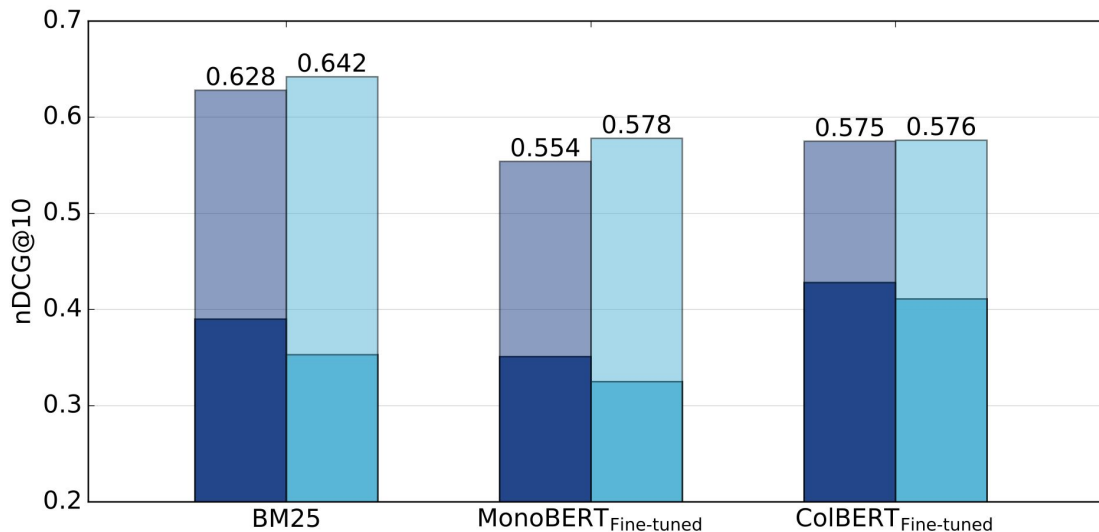
# Document translation results

Can we use LLMs to close the dialect gap?

Document transl.
Dialect → DE

There are still large gaps!



nDCG@10

| | 0.628 | 0.642 |
| BM25 | | |
| | 0.554 | 0.578 |
| MonoBERT_{Fine-tuned} | | |
| | 0.575 | 0.576 |
| ColBERT_{Fine-tuned} | | |

w/o variation (✗)    w/ variation (✓)

Exclude documents containing variations

full analysis split

*average over 5 dialects

# Conclusion

- We introduce WikiDIR, a cross-dialect information retrieval dataset.

- We release dialect variation dictionaries for German dialects.

- More results and analyses in the paper.

**GitHub**

# Conclusion

- We introduce WikiDIR, a cross-dialect information retrieval dataset.

- We release dialect variation dictionaries for German dialects.

- More results and analyses in the paper.

CDIR is novel and challenging task!

→ Low-resource

→ High-Variance

The gaps are still large.

**GitHub**