

# Can Large Language Models compete with specialized models in Lexical Semantic Change Detection?

Frank D. Zamora-Reina  
Felipe Bravo-Márquez  
**University of Chile, CENIA and IMFD**

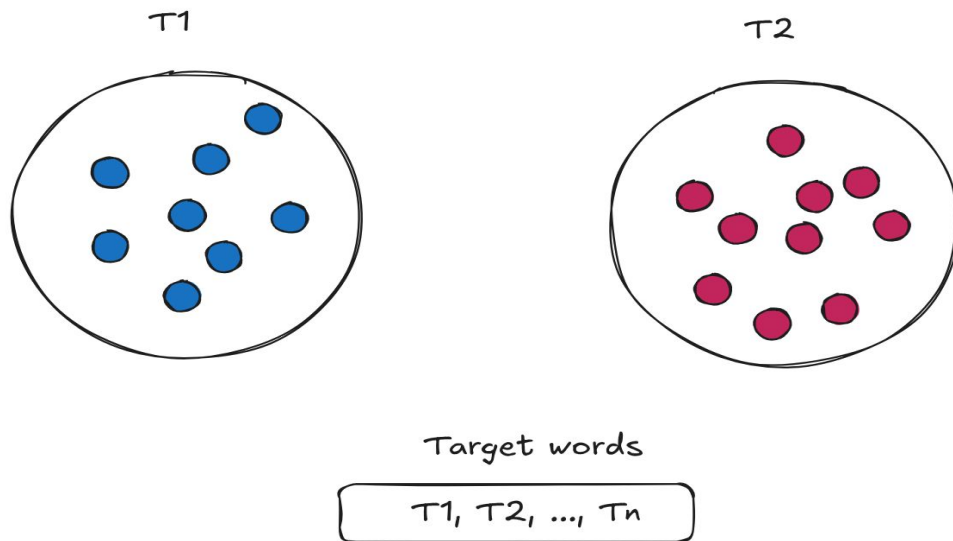
Dominik Schlechtweg  
**University of Stuttgart**

Nikolay Arefyev  
**University of Oslo**

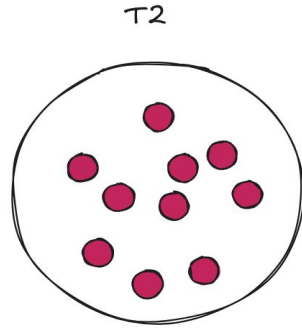
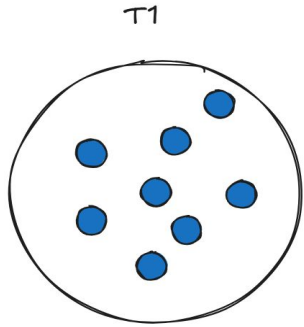
# Research questions

- [RQ1] Can automatically optimized prompts yield better results for the LSCD task than manually crafted prompts designed through prompt engineering?
- [RQ2] Can LLMs solve the Graded Change LSCD task well? Can these results surpass the WiC models reported as state-of-the-art?
- [RQ3] Can LLMs outperform state-of-the-art LSCD models at the annotation level?

# Lexical Semantic Change Detection



# Lexical Semantic Change Detection



Target words

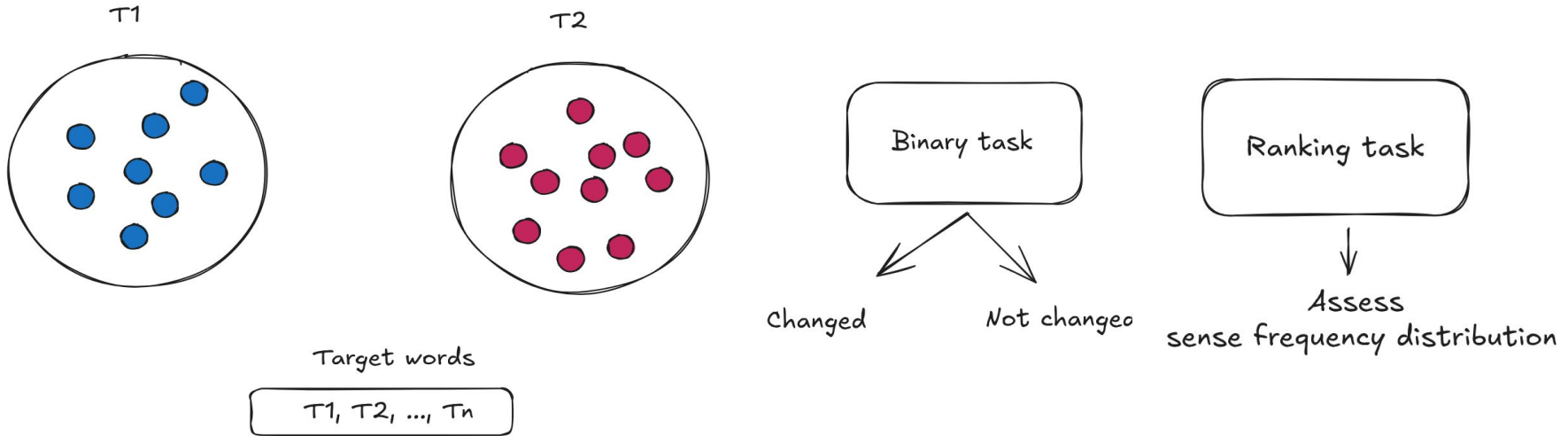
T1, T2, ..., Tn

Binary task

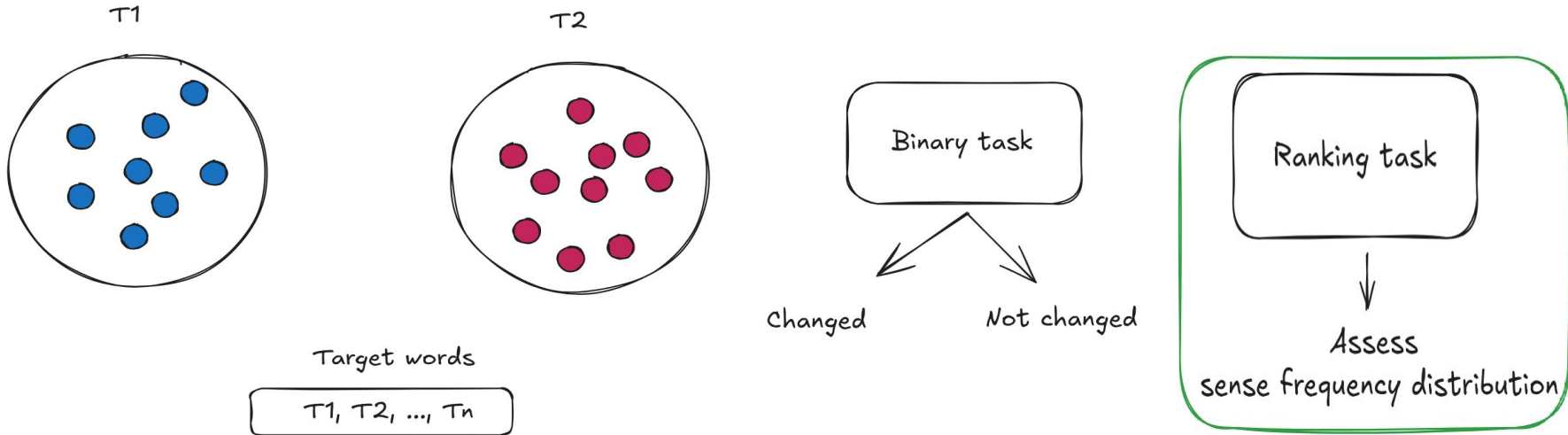
Changed

Not changed

# Lexical Semantic Change Detection



# Lexical Semantic Change Detection



# Metrics

- SPR\_WiC - Spearman correlation between the annotations provided by human annotators and models.
- SPR\_LSCD - Spearman correlation between gold graded change scores and model predictions

## Data

- DWUG EN (Schlechtweg et al., 2020)
  - 37 target words (~23k usage pairs)
  - (1810 - 1860) - (1960 - 2010)
- DWUG ES (Zamora-Reina et al., 2022)
  - 60 target words (~27k usage pairs)
  - (1810 - 1906) - (1994 - 2020)

(Schlechtweg et al., 2021)



# Data

**first usage:** His parents had left a lot of money in the **bank** and now it was Measle's, but a judge had said that Measle was too young to get it.

**second usage:** Sherrell, it is said, was sitting on the **bank** of the river close by, and as soon as the men had disappeared from sight he jumped on board the schooner.

**target word:** bank

annotation



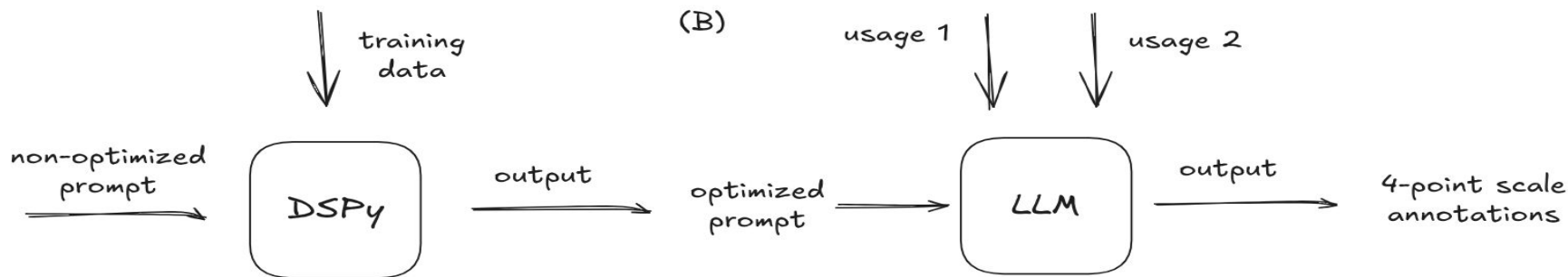
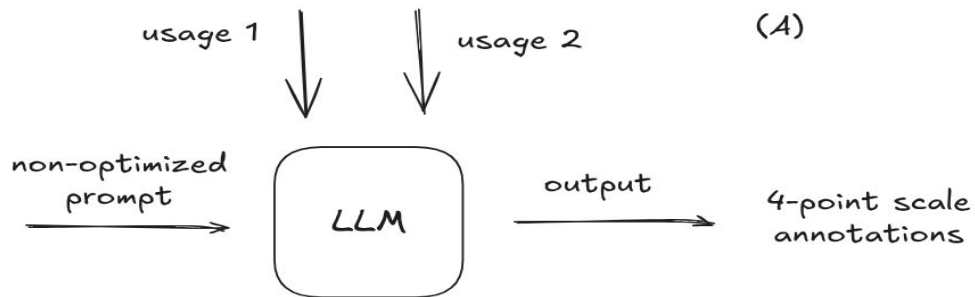
- 4: Identical
- 3: Closely Related
- 2: Distantly Related
- 1: Unrelated

(Schlechtweg et al., 2018)

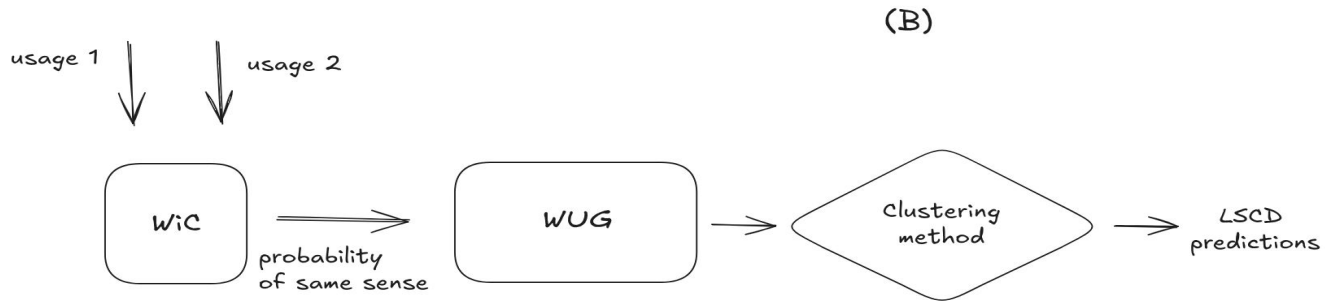
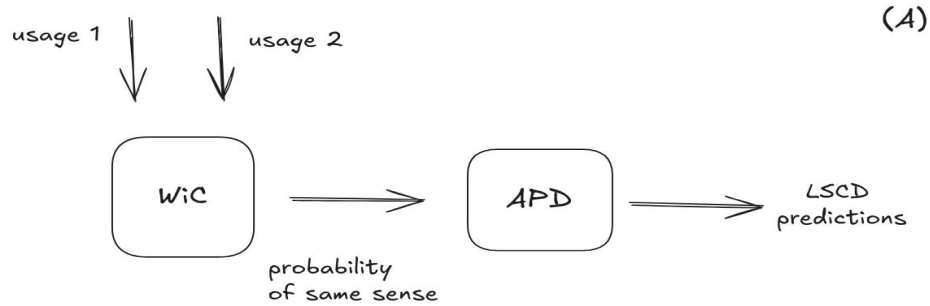
# Models

- LLMs
  - Llama 3.1 (8B)
  - Mixtral 8x7B (Jiang et al., 2024)
  - Llama 3.3 (70B)
- WiC models
  - DeepMistake (Arefyev et al., 2021)
    - SOTA for Russian and Spanish datasets
    - MCL, enMCL, MCL-> es (fine-tuned on various WiC datasets across multiple languages)
  - XL-LEXEME (Casotti et al., 2023)
    - SOTA for English, Swedish, German datasets

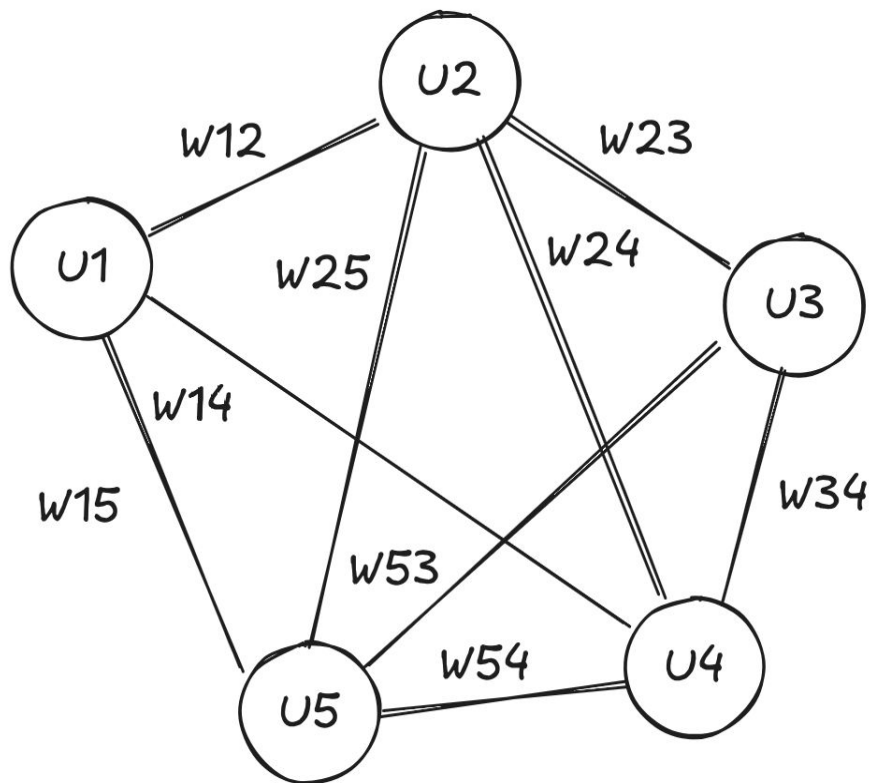
# Experimental Setup - LLMs



# Experimental setup - WiC models



## Word Usage Graph



$U_n$  — usage

$w_{ij}$  — weight

$w_{ij} \rightarrow [1-4], P(i)$

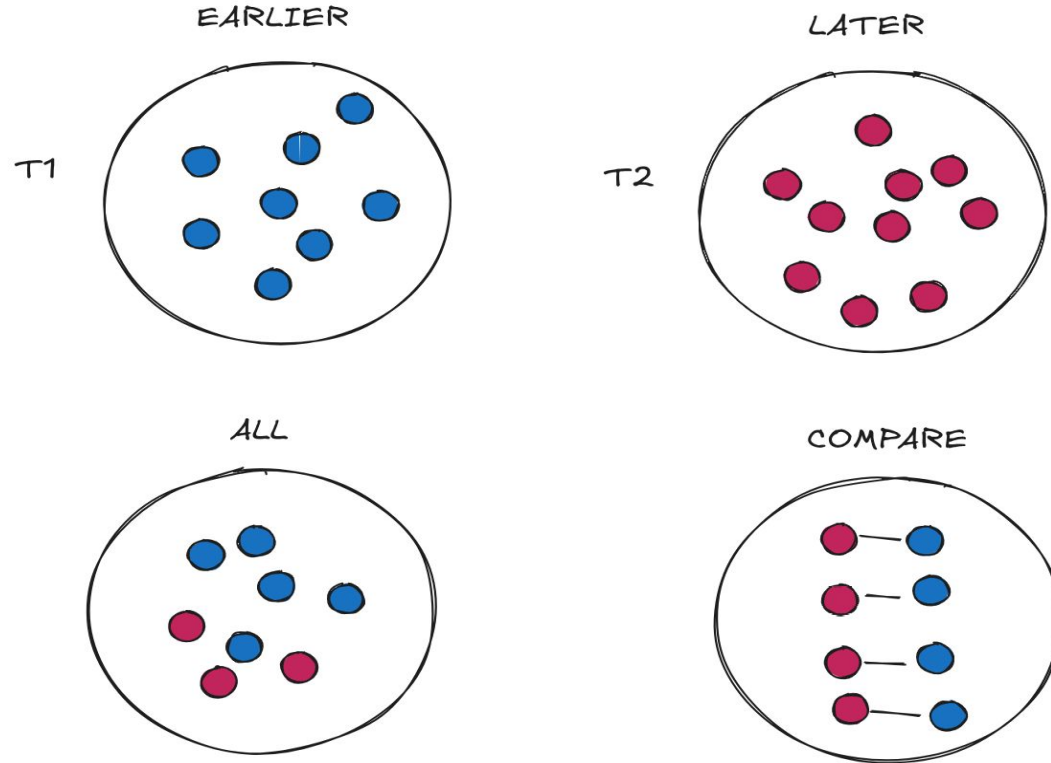
# Clustering algorithms

Spectral clustering (von Luxburg, 2007)

Agglomerative clustering (Jain and Dubes, 1988)

Weighted Stochastic Block Model (Peixoto, 2019)

# Lexical Semantic Change Detection



# Prompts

- optimized prompt from Yadav et al. (2024)
  - translate the prompt into Spanish
  - extend the prompt using samples from the DUREl framework (Schlechtweg et al., 2018)
- optimized prompts further using DSPy framework (Khattab et al., 2023)



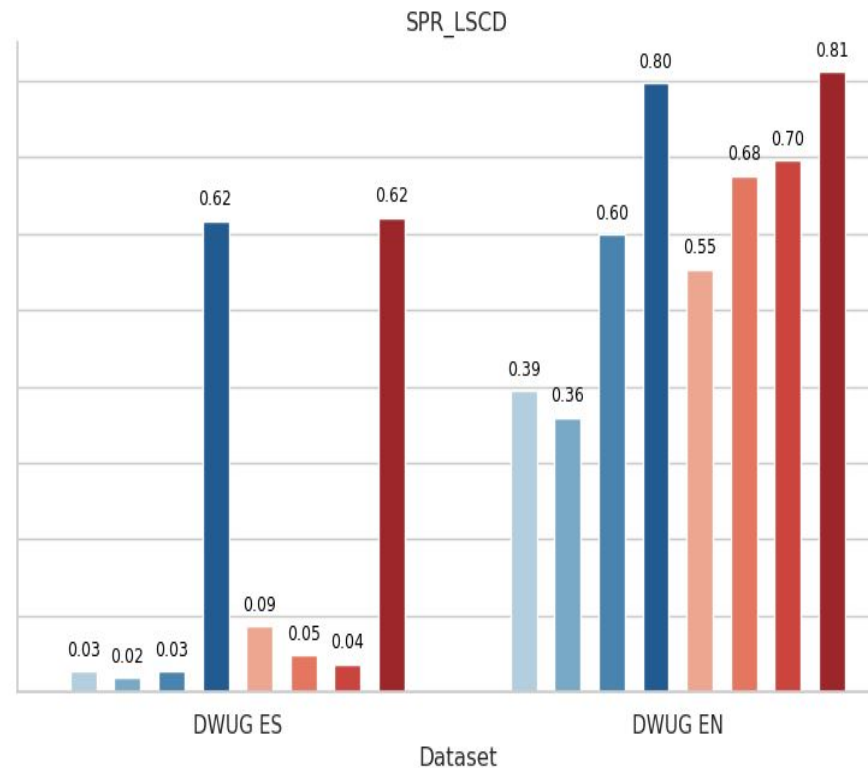
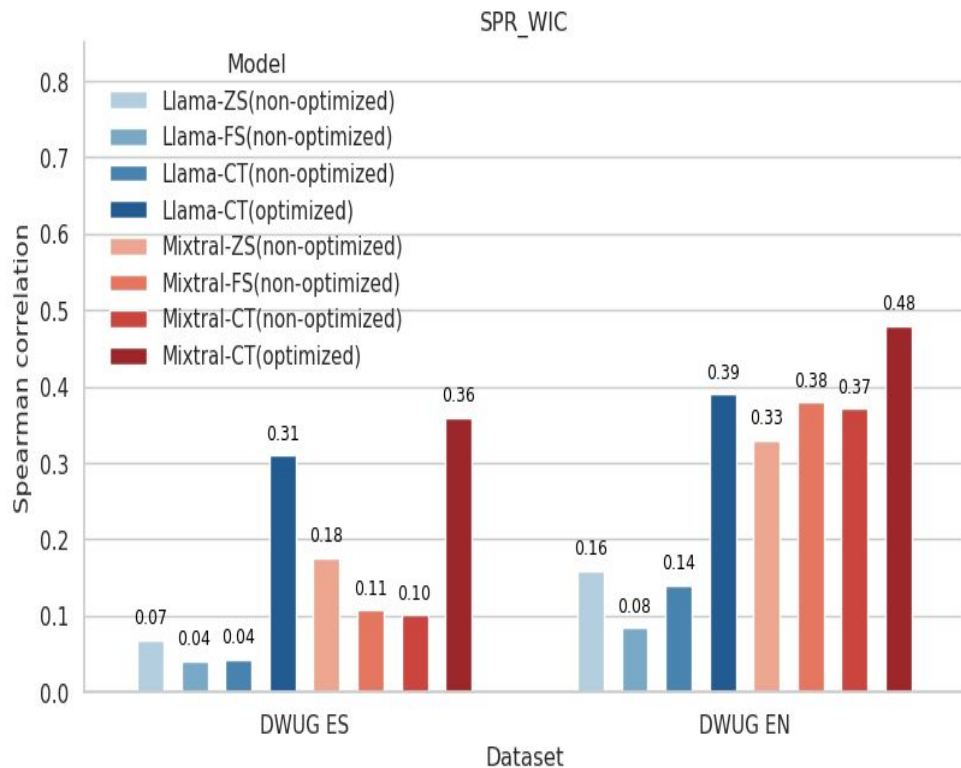
# Prompt optimization

Dataset-Prompts	Llama 3.1		Mixtral		Llama 3.3	
	NOP	OP	NOP	OP	NOP	OP
DWUG ES - PrS	26.8	29.5	32.6	31.22	32.33	39.77
DWUG ES - PrE	26.5	<b>35.5</b>	33.7	<b>34.80</b>	40.88	<b>46.33</b>
DWUG EN - PrS	25.75	31.0	32.8	<b>40.75</b>	37.25	45.5
DWUG EN - PrE	28.75	<b>37.25</b>	33.25	38.75	35.75	<b>49.25</b>

## RQ1

- Can automatically optimized prompts yield better results for the LSCD task than manually crafted prompts designed through prompt engineering?

## Non-optimized Prompts vs. Optimized Prompts



## RQ2

- Can LLMs solve the Graded Change LSCD task well? Can these results surpass the WiC models reported as state-of-the-art?

## Specialized WiC models vs LLMs in SPR\_LSCD



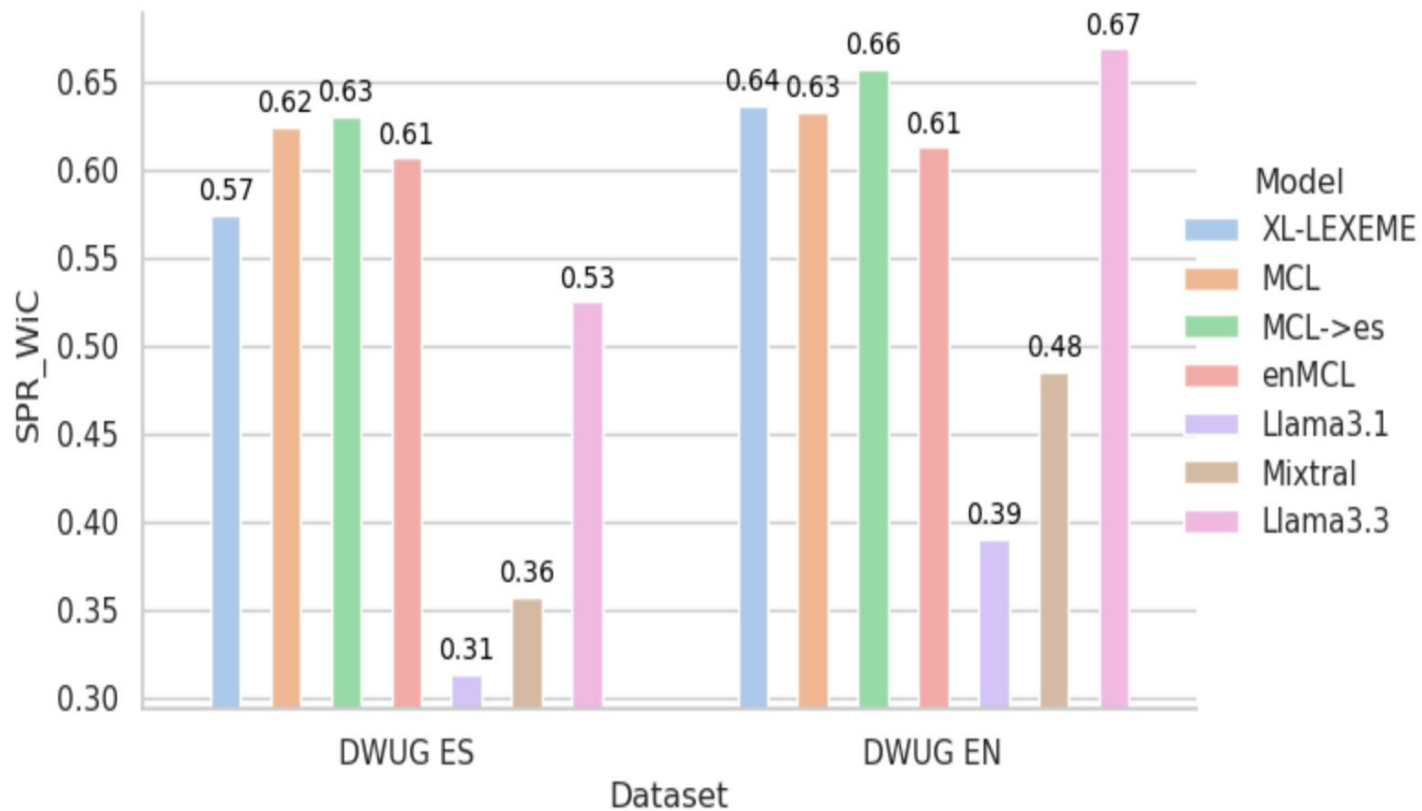
## WiC + WUG + Clustering methods - cross validation

Methods	Models	Spr_LSCD (ES)	ARI (ES)	Spr_LSCD (EN)	ARI (EN)
WSBM	Llama 3.1	.369 ± .204	.356 ± .116	<b>.845 ± .097</b>	.152 ± .067
SC		.271 ± .461	.102 ± .097	.014 ± .411	-.03 ± .01
AC		.478 ± .286	.073 ± .058	.205 ± .540	-.01 ± .03
APD		.636 ± .236	-	.645 ± .368	-
WSBM	Mixtral	.454 ± .180	.380 ± .104	.776 ± .219	.161 ± .07
SC		.565 ± .141	.092 ± .049	-.171 ± .492	-.03 ± .01
AC		.414 ± .075	.068 ± .027	-.04 ± .525	-.003 ± .02
APD		.567 ± .332	-	.612 ± .280	-
WSBM	Llama 3.3	.659 ± .181	<b>.502 ± .09</b>	.729 ± .241	.183 ± .08
SC		.507 ± .231	.294 ± .05	.302 ± .436	.124 ± .113
AC		.423 ± .184	.228 ± .05	.195 ± .273	-.01 ± .01
APD		.676 ± .195	-	.752 ± .227	-
WSBM	DeepMistake	<b>.727 ± .206</b>	.397 ± .074	.730 ± .212	.231 ± .212
SC		.561 ± .140	.355 ± .036	.520 ± .436	<b>.273 ± .115</b>
AC		.457 ± .320	.341 ± .054	.433 ± .245	.215 ± .128
APD		.653 ± .250	-	.638 ± .292	-
WSBM	XL-LEXEME	.630 ± .377	.452 ± .095	.686 ± .200	.152 ± .059
SC		.484 ± .215	.318 ± .043	.491 ± .176	.137 ± .063
AC		.426 ± .255	.292 ± .087	.143 ± .367	.02 ± .024
APD		.566 ± .354	-	.814 ± .199	-
WSBM	Random Baseline	-.199 ± .310	-.02 ± .184	-.111 ± .254	-.05 ± .147

## RQ3

- Can LLMs outperform state-of-the-art LSCD models at the annotation level?

Specialized WiC models vs LLMs in SPR\_WiC





## Conclusions

- recent prompt optimization techniques are crucial for achieving better results on the Graded Change LSCD task, as demonstrated by the performance of Llama3.3:70B
- medium-sized LLMs such as Mixtral:8x7B and Llama3.1:8B still underperform compared to smaller and faster specialized LSCD models in both the DWUG EN and DWUG ES datasets
  - in addition to optimization techniques, the size of the model also significantly influences the results

Thanks