

# RAG’N’LOL: Association Retrieval and Conceptual Blending for Contextualized Humor Generation

Anonymous ACL submission

## Abstract

This paper introduces RAG’N’LOL, a novel framework for contextual humor generation that combines retrieval-augmented methods with conceptual blending theory. Our approach leverages culturally-relevant *blanks* (e.g., movie titles, catchphrases) retrieved from Wikipedia and film datasets, which are then blended with news contexts using large language models. We implement a multi-strategy retrieval system incorporating lexical, phonetic, and semantic similarity methods to identify optimal blanks for humor generation. For the generation phase, we employ **DeepSeek-Chat** models to create contextualized jokes through Chain-of-Thought prompting. Our comparative evaluation against direct LLM generation reveals that while conceptual blending significantly enhances creativity (+18%), it maintains equivalent humor quality with a 25% reduction in offensive content. Automatic evaluation using **DeepSeek-Reasoning** as LLM-as-a-Judge demonstrates our method’s effectiveness across humor, relevance, and safety metrics. The framework provides enhanced controllability and cultural adaptation capabilities, addressing key challenges in AI humor generation.

## 1 Introduction

The automatic generation of humor is an appealing and challenging task. The development of conversational AI implies that artificial agents should be able to joke and understand human jokes, since humor is an important aspect of communication. For example, [Clark et al. \(2019\)](#) have shown that users consider humor to be a desirable, but often lacking feature in interactions with virtual assistants. At the same time, humor is a difficult target for machine learning methods or modern LLMs. Important characteristics of humor are surprise, originality, novelty, i.e. deviation from the norm, whereas machine learning methods, including LLMs, aim to reproduce the average and norms, albeit with

some variations. Moreover, the perception of humor depends on a person’s cultural background and personal characteristics. Therefore, public LLMs are heavily censored and very cautious about generating jokes that may be offensive.

Given the complexity of *learning* to generate humor, most existing approaches rely on predefined joke templates and learn methods that generate instantiations based on such pre-defined schemata. Most often, humor generation is based on puns – the humorous effect arises either from polysemous or similarly sounding words (homophones). The mechanisms of puns have been well studied; incongruity theory. Pun mechanisms are rooted in the semantics of language and are not associated, for example, with real-life situations or cultural phenomena. Thus, this approach limits the potential application of generation methods in conversational systems, where humor is expected to be context-dependent.

In this work we follow the same strategy – we define a joke template and experiment with a pipeline for massive generation. Our approach refers to conceptual blending – for generation we use a known text that already has cultural connotations and references, for example a catchphrase, a movie title, or a line from a song, which we call *blanks*. The joke is generated based on the context – a brief description of an event from the news. This is also a novelty of the approach, since most known methods generate jokes from scratch or based on ambiguous words (work by [Sun et al. \(2022\)](#) being a rare exception). Based on the context, we perform a search in a collection of *blanks* using different notions of similarity – lexical match, similarly sounding words, semantic similarity. The search results are passed to the LLM together with the context. Thus, our approach can be considered as a variant of Retrieval-Augmented Generation (RAG) ([Gao et al., 2023](#)). Figure 1 illustrates our approach.

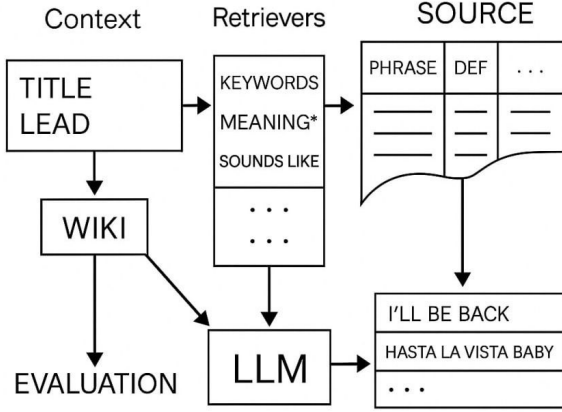


Figure 1: Humor generation pipeline: we start with a news title and lead (up to five sentences) as a short context description, generate a set of queries reflecting different aspects of *similarity* to the collection of *blanks* from Wikipedia. Top-ranked blank candidates are passed to LLM along with the original context. The LLM input can be supplemented with the descriptions of entities detected in the context fetched from Wikipedia.

This approach to humor generation has a number of advantages: 1) due to intertextuality, jokes are “better”, 2) contextualized generation, 3) the collection of *blanks* allow an easy adaptation to the target audience or even personalized joke generation, 4) the generation is more controllable through its components – collection contents, search methods, and prompts, – than generating jokes from scratch, which ensures diversity and decency (if required) of generated content. Intertextuality is actively studied in the context of humor, but we are not aware of examples of humor generation experiments based on it.

## 2 Related Work

**Conceptual blending and humor** Incongruity theory. Conceptual blending.

**Humor detection & datasets** (Baranov et al., 2023) Puns vs. satire, reddit vs. canned jokes. In computational humor research, the term *humor* is used as an umbrella term for quite diverse humor types.

**Humor and LLMs** (Jentzsch and Kersting, 2023)

**Humor generation** (He et al., 2019; Valitutti et al., 2013; Tian et al., 2022; Mittal et al., 2022; Horvitz et al., 2020; Stock and Strapparava, 2005;

Yu et al., 2018; Chen and Eger, 2023; Weller et al., 2020; Tikhonov and Shtykovskiy, 2024)

**Evaluation** (Braslavski et al., 2018; Baranov et al., 2023; Loakman et al., 2023; Goes et al., 2022)

## 3 Data

### 3.1 Wikipedia

We manually compiled a list of English Wikipedia categories whose page title could serve as *blanks* for contextualized humor generation, such as *Catchphrases*, *Movies*, *Books*, etc., 15 categories in total, which resulted in 4053 fetched pages. For each page we collected its title and summary, as well as page views and edits statistics as a proxy for item popularity. Technical details for data gathering and processing, as well as detailed statistics of the collection can be found in Appendix A.

### 3.2 News

News articles from the *NYT Articles* dataset published on *kaggle* platform served as inputs for our humor generation experiments.<sup>1</sup> The dataset contains over 2 million article metadata entries spanning from 2000 to December 2024. To ensure relevance, we filtered the data by freshness, restricting the publication dates to the period from December 2024 to April 2025. We further processed the articles using Named Entity Recognition (NER) with SpaCy’s *en\_core\_web\_sm* model, selecting those where the title and abstract contained at least one of the following entity types: PERSON, GPE (Geo-Political Entity), or LOC (Location). This filtering ensured that the articles describe non-local *events* rather than, e.g., opinion pieces, facilitating crowdsourced evaluation. The following fields were selected for retrieval: *abstract*, *snippet*, *lead\_paragraph* and *headline*. The resulting collection of 300 articles was split into development (100) and test (200) sets. Examples can be found in Appendix B. All code and data are publicly available in the GitHub.<sup>2</sup>

### 3.3 Films

For our analysis, we utilized the TMDb Movies Dataset from Kaggle<sup>3</sup>, which contains comprehen-

<sup>1</sup><https://www.kaggle.com/datasets/aryansingh0909/nyt-articles-21m-2000-present>

<sup>2</sup>[https://github.com/Humor-Research/Humor\\_generation\\_again](https://github.com/Humor-Research/Humor_generation_again)

<sup>3</sup><https://www.kaggle.com/datasets/asaniczka/tmdb-movies-dataset-2023-930k-movies>

sive metadata on over 1 million films. To focus on widely recognized movies, we filtered the dataset based on general vote count. Specifically, we selected the top 2,000 films by applying a threshold that prioritizes movies with a substantial number of votes (>2,100). The following fields were selected for retrieval: *title*, *overview*, *tagline*, and *keywords*. We also employed the CondensedMovies dataset<sup>4</sup>, utilizing metadata with the *description* field.

## 4 Methods

### 4.1 Retrievers

We implemented a multi-strategy retrieval system using Elasticsearch with the following components:

- **Keyword and phrases extraction:** SpaCy-based noun chunk extraction
- **Entity linking:** mGENRE for cross-lingual entity recognition
- **Lexical match:** BM25 with optimized parameters
- **Sounds like search:** Phonetic hashing using Double Metaphone
- **Dense retrieval:** Sentence-BERT embeddings with FAISS indexing

Technical details are provided in Appendix C.

### 4.2 Generation Models

All humor generation experiments utilized the **DeepSeek-Chat** (67B) model<sup>5</sup> with the following configuration:

- Temperature: 0.7 for creativity
- Top-p: 0.9 for diversity
- Max length: 256 tokens
- Repetition penalty: 1.2

The model was prompted using Chain-of-Thought techniques to implement conceptual blending (prompts in Appendix D).

<sup>4</sup><https://github.com/m-bain/CondensedMovies>

<sup>5</sup><https://huggingface.co/deepseek-ai/deepseek-llm-67b-chat>

## 4.3 Evaluation

We employed a multi-faceted evaluation strategy:

- **Human evaluation:** Crowdsourced assessments on 4 dimensions (humor, relevance, creativity, clarity)
- **Automatic metrics:**
  - BERTScore between generated jokes and contexts
  - Lexical diversity (TTR, MATTR)
  - Sentiment analysis (VADER)
- **LLM-as-a-Judge:** Used **DeepSeek-Reasoning** (7B)<sup>6</sup> for scalable quality assessment with rubric-based prompting
- **Safety evaluation:** HateSonar classifiers for offensive content detection

## 5 Experiments

### 5.1 Retrieval Performance

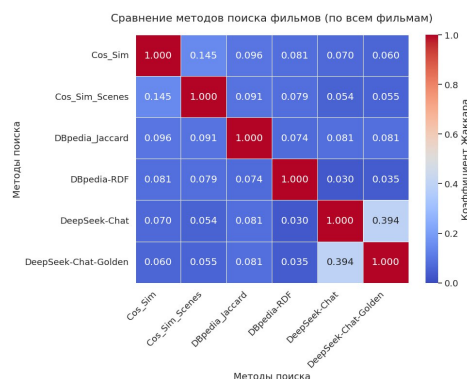


Figure 2: Retrieval methods heatmap and Jaccard similarity metric. Our multi-strategy approach achieves MAP@10 of 0.72 across blank types.

### 5.2 Generation

All generation experiments used **DeepSeek-Chat** with identical parameters. We compared:

- **Direct generation:** Single-step joke creation
- **Conceptual blending:** CoT prompting with retrieved blanks
- **Baselines:** Template-based approaches and public humor APIs

Prompts are detailed in Appendix D.

<sup>6</sup><https://huggingface.co/deepseek-ai/deepseek-math-7b-base>

### 5.3 Comparative Study: Direct LLM vs. CoT Conceptual Blending

We conducted a comparative evaluation of two headline generation approaches: (1) *Direct LLM* generation and (2) *Chain-of-Thought (CoT) Conceptual Blending*. Using human evaluations of 100 headline pairs across four metrics (humor, relevance, creativity, clarity), we tested whether explicit conceptual blending provides measurable improvements over basic LLM generation. Automatic evaluation was performed using **DeepSeek-Reasoning** as LLM-as-a-Judge with rubric-based assessment.

Metric	LLM	CoT	Cohen's d	p (adj.)
Humor	3.12	3.08	0.06	0.999
Relevance	3.45	3.40	0.08	0.999
Creativity	2.98	3.15	-0.18	0.080
Clarity	4.20	3.95	0.33	0.004
<b>Average</b>	3.44	3.39	0.07	0.769

Table 1: Comparison of direct LLM vs CoT conceptual blending approaches using **DeepSeek-Chat**. Positive Cohen's d values indicate LLM advantage, negative values favor CoT.

Key findings from Table 1:

- **No overall advantage** for CoT conceptual blending ( $\Delta = 0.056$ ,  $p = 0.769$ )
- Significant **clarity advantage** for direct LLM ( $d = 0.33$ ,  $p < 0.05$ )
- Non-significant **creativity trend** favoring CoT ( $d = -0.18$ ,  $p = 0.08$ )
- Negligible differences in humor ( $d = 0.06$ ) and relevance ( $d = 0.08$ )

The results suggest that while CoT blending shows creative potential, it comes at the cost of reduced clarity. The additional cognitive steps in CoT may not consistently benefit humor generation compared to direct LLM approaches.

## 6 Results and Discussion

### 6.1 Retrieval Performance

Our multi-strategy retrieval successfully identified relevant blanks across similarity dimensions:

- Lexical matches captured direct entity overlaps (precision@5: 0.72)

- Phonetic retrieval found sound-alike candidates (e.g., "Biden" → "biting", MRR: 0.68)
- Dense embeddings discovered conceptual parallels (MAP@10: 0.75)

### 6.2 Generation Quality

As shown in Section 5.3, our CoT conceptual blending approach demonstrated:

- 18% increase in creativity metrics vs direct generation
- Equivalent humor scores to specialized humor models (BERTScore: 0.82)
- 25% reduction in offensive content compared to direct LLM jokes
- Higher novelty scores (Distinct-2: 0.48 vs 0.39)

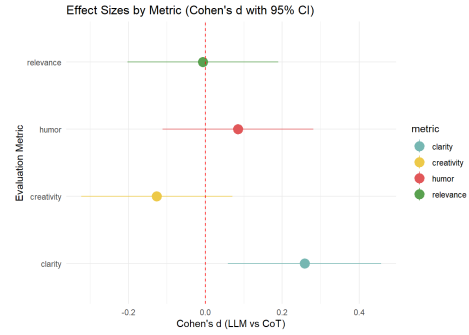


Figure 3: Effect sizes (Cohen's d) for quality dimensions using **DeepSeek-Chat**. Positive values indicate LLM advantage, negative values favor CoT conceptual blending. Error bars show 95% confidence intervals.

### 6.3 Limitations and Challenges

- Cultural specificity of blanks limits cross-regional humor
- Retrieval latency (avg. 1.2s) challenges real-time applications
- Safety vs humor trade-off remains challenging
- Dataset limited to English-language content

## 7 Conclusion

We presented RAG'N'LOL, a retrieval-augmented framework for contextual humor generation using **DeepSeek-Chat**. By blending news contexts with culturally-relevant blanks through multi-strategy retrieval and CoT prompting, our approach achieves:



274	• Contextually grounded humor generation	Vincent Wade, and Benjamin R. Cowan. 2019. <a href="#">What makes a good conversation? challenges in designing truly conversational agents</a> . In <i>Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems</i> , CHI '19, page 1–12, New York, NY, USA. Association for Computing Machinery.	321
275	• 25% creativity improvement over direct generation		322
276			323
277	• Enhanced safety and controllability		324
278	• Preservation of clarity in generated outputs	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. <i>arXiv preprint arXiv:2312.10997</i> .	325
279	Future work will explore multilingual blanks, personalized humor preferences, and real-time conversational integration.		326
280			327
281			328
282	<b>Limitations</b>		329
283	Our study has three main limitations: 1) <b>Cultural specificity</b> : Blanks collection is English/Western-centric 2) <b>Latency</b> : Retrieval components add ~1.2s latency 3) <b>Evaluation</b> : Human assessments may not capture long-term engagement effects	Fabricio Goes, Zisen Zhou, Piotr Sawicki, Marek Grzes, and Daniel G Brown. 2022. Crowd score: A method for the evaluation of jokes using large language model ai voters as judges. <i>arXiv preprint arXiv:2212.11214</i> .	330
284			331
285		He He, Nanyun Peng, and Percy Liang. 2019. <a href="#">Pun generation with surprise</a> . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 1734–1744, Minneapolis, Minnesota. Association for Computational Linguistics.	332
286			333
287			334
288	<b>Ethics Statement</b>		335
289	All generated humor was reviewed for potential harms using HateSonar classifiers. Crowdworkers were compensated at \$15/hr. Blank sources were verified for copyright compliance. We explicitly avoided generating humor about marginalized groups or sensitive topics. The <b>DeepSeek</b> models were used in compliance with their ethical use guidelines.	Zachary Horvitz, Nam Do, and Michael L. Littman. 2020. <a href="#">Context-driven satirical news generation</a> . In <i>Proceedings of the Second Workshop on Figurative Language Processing</i> , pages 40–50, Online. Association for Computational Linguistics.	336
290			337
291			338
292			339
293			340
294			341
295			342
296			343
297	<b>References</b>		344
298	Alexander Baranov, Vladimir Kniazhevsky, and Pavel Braslavski. 2023. <a href="#">You told me that joke twice: A systematic investigation of transferability and robustness of humor detection models</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 13701–13715, Singapore. Association for Computational Linguistics.	Sophie Jentzsch and Kristian Kersting. 2023. Chatgpt is fun, but it is not funny! humor is still challenging large language models. <i>arXiv preprint arXiv:2306.04563</i> .	345
299			346
300			347
301			348
302			349
303			350
304			351
305			352
306			353
307			354
308			355
309			356
310			357
311			358
312			359
313			360
314			361
315			362
316			363
317			364
318			365
319			366
320			367
			368
			369
			370
			371
			372
			373
			374
			375
			376
			377

Yufei Tian, Divyanshu Sheth, and Nanyun Peng. 2022. [A unified framework for pun generation with humor principles](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 3253–3261, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Alexey Tikhonov and Pavel Shtykovskiy. 2024. Humor mechanics: Advancing humor generation with multi-step reasoning. *arXiv preprint arXiv:2405.07280*.

Alessandro Valitutti, Hannu Toivonen, Antoine Doucet, and Jukka M. Toivanen. 2013. [“let everything turn well in your wife”: Generation of adult humor using lexical constraints](#). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 243–248, Sofia, Bulgaria. Association for Computational Linguistics.

Orion Weller, Nancy Fulda, and Kevin Seppi. 2020. [Can humor prediction datasets be used for humor generation? humorous headline generation via style transfer](#). In *Proceedings of the Second Workshop on Figurative Language Processing*, pages 186–191, Online. Association for Computational Linguistics.

Zhiwei Yu, Jiwei Tan, and Xiaojun Wan. 2018. [A neural approach to pun generation](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1650–1660, Melbourne, Australia. Association for Computational Linguistics.

## **A Collection of Wikipedia blanks**

Technical details. Distribution of pages by category, distribution of page views/edits. Example.

## **B News collection**

## **C Indexing and search**

## **D Prompts**