

# Analysis and Generation of Verbal Humor in Portuguese

## Paper type: Doctoral Consortium

**Marcio Lima Inácio**

Centre for Informatics and Systems of the University of Coimbra (CISUC)  
Univeristy of Coimbra  
Polo II, Pinhal de Marrocos, 3030-290 Coimbra, Portugal  
{mlinacio}@dei.uc.pt

### Abstract

Computational Creativity is the field of research focused on formalizing, discussing, and developing computational systems capable of producing artifacts which would be considered as creative by humans. For Natural Language Processing, working with creativity is specially challenging and intriguing, since it raises questions about figurative language, pragmatics, ambiguity, and other topics which are still open for fruitful research. From the many tasks involving language within the Computational Creativity area, this project focuses on the automatic recognition and generation of humorous texts (jokes, riddles, puns, and the like). The core of this project will be to integrate knowledge from different levels of linguistic investigation – such as Phonology, Semantics, Pragmatics, and others – in order to produce, analyze and evaluate texts which can be considered more funny by human readers. It will also be important to measure how important each of the information types are for computational tasks involving humor. As a side effect, it is expected to develop resources, tools, and methods which could be valuable for research not only in Computational Creativity but also in Natural Language Processing as a whole, specially for the Portuguese language, the main target of this project, which is considerably less researched than other languages such as English or Mandarin Chinese.

### Introduction

There is a natural path on Natural Language Processing (NLP), as argued by Cambria and White (2014), which follows the commonly defined layers of linguistic knowledge – from morphology through syntax towards semantics and pragmatics – since interpreting figurative language, context and intentions are a part in human-to-human interaction. In this way, computational systems capable of producing and dealing with such information are essential to achieve the ultimate goal of full natural language-based human-computer interaction (Hempelmann 2008); furthermore, those systems may provide important insights for psycholinguists to better understand how humans themselves deal with such phenomena (Reyes, Rosso, and Buscaldi 2012; Reyes Pérez 2013; Cambria and White 2014).

Humorous texts pose a challenge for NLP, as their interpretation is influenced by different aspects studied across

multiple disciplines, such as linguistics, psychology, philosophy and sociology (Reyes, Rosso, and Buscaldi 2012). The specific case of humor is also important in some applications, as in identifying humor, irony or sarcasm in interactions with a user (Valitutti et al. 2016), differentiating satirical news from “true” ones (Yang, Mukherjee, and Dragut 2017), or even to make the communication with people more pleasant and human-like (Amin and Burghardt 2020). In this context, NLP borders on the Computational Creativity (CC) field of research. CC researchers deal with the processing – generation and interpretation – of any kind of activity or artifact which would be interpreted as being creative by unbiased observers (Colton and Wiggins 2012), e.g. poetry or humor as instances of linguistic artifacts, whose automatic creation has a strong relation with the Natural Language Generation (NLG) field of NLP (Gatt and Krahmer 2018).

As discussed by different authors in the literature, there is still much to advance within the field of CC, more specifically on humor processing within NLP. Amin and Burghardt (2020) mention that only a minority of research on Computational Humor Generation is grounded explicitly in some sort of humor theory, be it linguistic, psychological, or sociological. The authors also observed that there is no standardized methodology to evaluate the quality of humor-related computational models. The absence of such evaluation hinders the development of systems, as it could help to design techniques capable of distinguishing creativity from nonsense (He, Peng, and Liang 2019).

Some authors also advocate to the need of grounding systems on some sort of existing theory, because standard methods for generation, e.g. Neural Networks, are created to follow specific patterns, which is the direct opposite concept of creativity (He, Peng, and Liang 2019). Current neural attempts are still reported to not produce humor in general, despite creating more complex and sophisticated outputs (Amin and Burghardt 2020). Furthermore, there is some arguing about how much those types of models can learn about more subtle aspects of language (Bender and Koller 2020).

Another important aspect to be taken into account is that NLP research is generally concentrated on a small set of languages, such as English and Mandarin Chinese (Bender 2019). Computational Humor research in Portuguese is limited to only a few works (Gonçalo Oliveira, Clemêncio, and Alves 2020; Gonçalo Oliveira and Rodrigues 2018;

Gonçalo Oliveira, Costa, and Pinto 2016; Mendes and Gonçalo Oliveira 2020), even when expanding the interpretation to some related phenomena like irony (Wick-Pedro and Vale 2020) or satires (Wick-Pedro et al. 2020). Therefore, there is much to be advanced with this language, whose investigation on humor can also produce important resources and methods to be adopted on other tasks, such as standard NLG, Natural Language Understanding (NLU), lexical knowledge bases, sentiment analysis and others.

## Related Work

The works most related to this project are those focused on computational humor research, specially for the Portuguese language. Costa, Gonçalo Oliveira, and Pinto (2015) presented in 2015 the first work on automatic humor generation for Portuguese, which was later carried on by Gonçalo Oliveira, Costa, and Pinto (2016), creating memes based on some templates and linguistic features. Later, in 2018, an automatic riddles generation method was developed by Gonçalo Oliveira and Rodrigues (2018). The authors also used templates to create the riddles based on lexical semantic relations.

More recently, there has also been some research on recognizing humorous texts in Portuguese (Clemêncio 2019; Gonçalo Oliveira, Clemêncio, and Alves 2020) using n-gram features, Named Entity Recognition, sentiment analysis, and other linguistic characteristics to train Machine Learning models. Another generation approach is the TECO system, created by Mendes and Gonçalo Oliveira (2020), which uses word embeddings for creating humorous headlines for news articles.

Some corpora are also available for dealing with humor and figurative language in Portuguese. Gonçalo Oliveira, Clemêncio, and Alves (2020) presented a corpus of one-liners – short jokes – alongside non-humorous content with similar structure. There is also a corpus for irony detection in tweets created by Wick-Pedro and Vale (2020). Later, the same authors published a corpus including satirical news and their corresponding non-satirical ones (Wick-Pedro et al. 2020).

On a more historical note, computational humor generation can be traced back to the late decade of 1990, with systems such as JAPE (Binsted and Ritchie 1994; 1997) and LIBJOG (Raskin and Attardo 1994). Another notable methods from the literature are HAHACronym (Stock and Strapparava 2003) and STANDUP (Binsted et al. 2006; Ritchie et al. 2006).

## Objectives

The objective of this PhD Thesis Project will be to take advantage of existing linguistic and psycholinguistic theories of humor, as they are believed to overcome some of the limitations presented by current systems – such as low humorosity, a trade-off between funniness and linguistic complexity (Amin and Burghardt 2020), and low grammaticality (He, Peng, and Liang 2019) – with a special focus on the Portuguese language.

In order to apply those theories, it will be necessary to develop methods for “mapping” the abstract linguistic concepts used in the literature into some computer-tractable artifacts, such as rules, metrics, lexicons, and others. Those techniques will be employed not only on the recognition (or analysis), but also on the creation of humor.

Regarding the functioning of such systems, we can focus on using non-humorous content as input and develop methods to manipulate it in order to eventually create humor; this will be relevant not only because of the methods per se, but also because it can be a suitable way of assessing the validity of some underlying linguistic theory and creating new insights about humor that can be useful to different areas. The nature of the non-humorous content is yet to be specified, but it could be news, interactions of the user with some chatbot, social media posts, commonly used expressions, and others.

There is also a specific interest on evaluating and creating methods for benefiting from the advantages of Neural Language Models, such as their capability of generating good quality and linguistically complex text, despite their known weaknesses on dealing with creativity (He, Peng, and Liang 2019) and more abstract linguistic knowledge (Bender and Koller 2020). Tying to conciliate both approaches seems an interesting and important path for research.

With regards to humor recognition, we could use or develop some corpora to classify or even determine a level of funniness for jokes. Another option is to dive into identifying satirical news, which can also be fruitful to the field of fake news detection (Wick-Pedro et al. 2020).

Besides generating and recognizing humor, it is of interest to develop techniques for these computational models to be able of explaining why some texts may be considered funny by some readers, how they reached this conclusion, and also what are the decisions taken towards generating some humorous text, which is an interesting path to follow given the current trend of transparency in Artificial Intelligence (Ehsan et al. 2021).

## Expected Results

As a result, it is expected that this PhD project produces some methods for creating, analyzing, recognizing, evaluating, and explaining verbal humor automatically, based on existing and developing theories of humor and figurative language. We expect in the end that the findings and systems developed not only help the final user in creating their texts, but also may contribute to the scientific community (specially the Portuguese-speaking one) with interesting artifacts, resources and tools to be analyzed and with insights for further research on humor and how it is created and perceived by people and machines.

## Acknowledgements

This work is funded by national funds through the FCT - Foundation for Science and Technology, I.P., within the scope of the project CISUC - UID/CEC/00326/2020 and by European Social Fund, through the Regional Operational Program Centro 2020. The author is funded by Foundation

for Science and Technology (FCT), Portugal, under the grant UI/BD/153496/2022.

## References

- Amin, M., and Burghardt, M. 2020. A survey on approaches to computational humor generation. In *Proceedings of the 4th Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature*, 29–41. Online: International Committee on Computational Linguistics.
- Bender, E. M., and Koller, A. 2020. Climbing towards NLU: On Meaning, Form, and Understanding in the Age of Data. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 5185–5198. Online: Association for Computational Linguistics.
- Bender, E. M. 2019. The #BenderRule: On Naming the Languages We Study and Why It Matters.
- Binsted, K., and Ritchie, G. 1994. An implemented model of punning riddles. In Hayes-Roth, B., and Korf, R. E., eds., *Proceedings of the 12th National Conference on Artificial Intelligence*, volume 1, 633–638. Seattle: AAAI Press / The MIT Press.
- Binsted, K., and Ritchie, G. 1997. Computational rules for generating punning riddles. *Humor - International Journal of Humor Research* 10(1).
- Binsted, K.; Bergen, B.; Coulson, S.; Nijholt, A.; Stock, O.; Strapparava, C.; Ritchie, G.; Manurung, R.; Pain, H.; Waller, A.; and O’Mara, D. 2006. Computational Humor. *IEEE Intelligent Systems* 21(2):59–69.
- Cambria, E., and White, B. 2014. Jumping NLP Curves: A Review of Natural Language Processing Research. *IEEE Computational Intelligence Magazine* 9(2):48–57.
- Clemêncio, A. 2019. *Reconhecimento Automático de Humor Verbal*. MSc, Universidade de Coimbra, Coimbra.
- Colton, S., and Wiggins, G. A. 2012. Computational creativity: The final frontier? In *Proceedings of the 20th European Conference on Artificial Intelligence, ECAI’12*, 21–26. NLD: IOS Press.
- Costa, D.; Gonçalves Oliveira, H.; and Pinto, A. M. 2015. “In reality there are as many religions as there are papers” – First Steps Towards the Generation of Internet Memes. In Toivonen, H.; Colton, S.; Cook, M.; and Ventura, D., eds., *Proceedings of the Sixth International Conference on Computational Creativity (ICCC 2015)*, 300–307. Park City, Utah: Brigham Young University.
- Ehsan, U.; Liao, Q. V.; Muller, M.; Riedl, M. O.; and Weisz, J. D. 2021. Expanding Explainability: Towards Social Transparency in AI systems. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, 1–19. Yokohama Japan: ACM.
- Gatt, A., and Kraemer, E. 2018. Survey of the State of the Art in Natural Language Generation: Core tasks, applications and evaluation. *Journal of Artificial Intelligence Research* 61:65–170.
- Gonçalo Oliveira, H., and Rodrigues, R. 2018. Explorando a Geração Automática de Adivinhas em Português. *Linguamática* 10(1):3–18.
- Gonçalo Oliveira, H.; Clemêncio, A.; and Alves, A. 2020. Corpora and baselines for humour recognition in Portuguese. In *Proceedings of the 12th Language Resources and Evaluation Conference*, 1278–1285. Marseille, France: European Language Resources Association.
- Gonçalo Oliveira, H.; Costa, D.; and Pinto, A. M. 2016. One does not simply produce funny memes! - Explorations on the Automatic Generation of Internet humor. In Pachet, F.; Cardoso, A.; Corruble, V.; and Ghedini, F., eds., *Proceedings of the Seventh International Conference on Computational Creativity*, 238–245. Paris: Sony CSL Paris, France.
- He, H.; Peng, N.; and Liang, P. 2019. Pun Generation with Surprise. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, volume 1, 1734–1744. Minneapolis: Association for Computational Linguistics.
- Hempelmann, C. F. 2008. Computational humor: Beyond the pun? In *The Primer of Humor Research*, number 8 in Humor Research. Berlin, New York: Victor Raskin. 333–360.
- Mendes, R., and Gonçalves Oliveira, H. 2020. Amplifying the range of news stories with creativity: Methods and their evaluation, in Portuguese. In *Proceedings of the 13th International Conference on Natural Language Generation*, 252–262. Dublin, Ireland: Association for Computational Linguistics.
- Raskin, J. D., and Attardo, S. 1994. Non-literality and non-bona-fide in language: An approach to formal and computational treatments of humor. *Pragmatics & Cognition* 2(1):31–69.
- Reyes Pérez, A. 2013. Linguistic-based Patterns for Figurative Language Processing: The Case of Humor Recognition and Irony Detection. *Procesamiento del Lenguaje Natural*.
- Reyes, A.; Rosso, P.; and Buscaldi, D. 2012. From humor recognition to irony detection: The figurative language of social media. *Data & Knowledge Engineering* 74:1–12.
- Ritchie, G. D.; Manurung, R.; Pain, H.; Waller, A.; and O’Mara, D. 2006. The STANDUP interactive riddle builder. *IEEE Intelligent Systems* 21(2):67–69.
- Stock, O., and Strapparava, C. 2003. HAHAcronym: Humorous Agents for Humorous Acronyms. *Humor - International Journal of Humor Research* 16(3).
- Valitutti, A.; Doucet, A.; Toivanen, J. M.; and Toivonen, H. 2016. Computational generation and dissection of lexical replacement humor. *Natural Language Engineering* 22(5):727–749.
- Wick-Pedro, G., and Vale, O. A. 2020. Comentcorpus: descrição e análise de ironia em um corpus de opinião para o português do Brasil. *Cadernos de Linguística* 1(2):01–15.
- Wick-Pedro, G.; Santos, R. L. S.; Vale, O. A.; Pardo, T. A. S.; Bontcheva, K.; and Scarton, C. 2020. Linguistic Analysis Model for Monitoring User Reaction on Satirical News for Brazilian Portuguese. In Quaresma, P.; Vieira, R.;

Aluísio, S.; Moniz, H.; Batista, F.; and Gonçalves, T., eds., *Computational Processing of the Portuguese Language*, volume 12037. Cham: Springer International Publishing. 313–320.

Yang, F.; Mukherjee, A.; and Dragut, E. 2017. Satirical News Detection and Analysis using Attention Mechanism and Linguistic Features. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 1979–1989. Copenhagen, Denmark: Association for Computational Linguistics.