



Universidad Politécnica de Madrid

Escuela Técnica Superior de Ingenieros Informáticos

Máster Universitario en Digital Innovation

Master Thesis

Easy-to-Read (E2R) Adaptation of Figurative Language

Author: **Jelle van Lieshout**

Supervisor: **Dr. Mari Carmen Suárez-Figueroa**

Madrid, MM YYYY

Abstract

Figurative language, such as idioms and metaphors, presents a significant barrier to text accessibility and downstream Natural Language Processing tasks. While Large Language Models (LLMs) have shown promise in general language understanding, their ability to reliably handle figurative expressions in zero-shot scenarios remains under-explored compared to task-specific supervised models. This thesis investigates the effectiveness of agentic LLM pipelines for the detection, interpretation, and Easy-to-Read (E2R) adaptation of figurative language.

We propose a multi-stage agentic workflow that explicitly decomposes the task into detection, explanation, literal replacement, and self-verification. Using the SemEval-2022 Task 2 dataset for idioms and the VU Amsterdam Metaphor Corpus for metaphors, we evaluate the system's performance against both monolithic single-prompt LLM baselines and established supervised benchmarks.

We hypothesize that an agentic decomposition significantly improves the quality of literal replacements compared to standard prompting, particularly for complex metaphors that rely on abstract source-target mappings. Furthermore, we explore the trade-offs between semantic fidelity and readability in the generated replacements. This work aims to demonstrate that structured, observable agentic reasoning enables reliable figurative language interpretation and transformation without the need for extensive task-specific training, thereby contributing to more accessible and inclusive digital content.

Resumen

...

Acknowledgement

...

Contents

Abstract	A
Resumen	B
Acknowledgement	C
1 Introduction	1
2 Background & Related Work	2
2.1 Figurative language	2
2.2 Datasets	2
2.3 Large Language Models & Agentic Systems	2
2.4 Literature Review	3
3 System Design & Methodology	4
3.1 Problem formulation	4
3.2 Agentic pipeline	4
3.3 Baselines	4
4 RQ1: Detection & Interpretation	5
4.1 Experiments	5
4.2 Metrics	5
5 RQ2: Agentic Decomposition	6
5.1 Experiments	6
5.2 Metrics	6
6 RQ3: Replacement Quality	7
6.1 Idioms	7
6.2 Metaphors	7
6.3 Analysis	7
7 RQ4: Observability & Error Analysis	8
8 Discussion & Limitations	9
9 Conclusion & Future Work	10
10 Annex	12

List of Tables

List of Figures

List of Equations

1 Introduction

- Motivation: figurative language as a barrier to accessibility, simplification, and downstream NLP
- Limitations of supervised, task-specific approaches
- Rise of LLMs and agentic reasoning
- Contributions
 - Zero/one-shot/few-shot agentic system for figurative language handling
 - Unified treatment of idioms and metaphors (is this the case?)
 - Empirical evaluation of replacement quality [1]

2 Background & Related Work

2.1 Figurative language

- Idioms vs metaphors (linguistic + cognitive distinction)
- Prior work on detection, classification, and simplification

2.2 Datasets

- SemEval-2022 Task 2 (idioms) [11]
 - This shared task focuses on multilingual idiomaticity detection and sentence embedding.
 - It involves a binary classification subtask to determine whether a sentence contains an idiomatic expression.
 - A Python client has been developed to interface with this dataset for training and evaluation purposes.
- VU Amsterdam Metaphor Corpus (VUAMC) [8]
 - The VUAMC is a large, annotated corpus of English texts, manually tagged for metaphorical language use.
 - It is generally available via the Oxford Text Archive (currently down, legacy link accessible).
 - The corpus is available as XML, providing rich metadata for each token.
 - A Python client has been developed to interface with the dataset, enabling easy access and filtering of metaphorical instances.
 - Goal: Compare annotation methodology with other datasets and explore possibilities for merging.
- Supporting datasets (MAGPIE, TroFi, MOH)
- Conceptual resources (MetaNet)

2.3 Large Language Models & Agentic Systems

- Zero-shot, one-shot and few-shot prompting
- Agentic decomposition (planning, reflection, self-verification)
- Observability (LangSmith-style traces)

2.4 Literature Review

Relevant papers to include:

- [9] Mari Carmen Suárez-Figueroa, Isam Diab, Edna Ruckhaus, and Isabel Cano. First steps in the development of a support application for easy-to-read adaptation. *Universal Access in the Information Society*, 23:365–377, 3 2024
- [10] Mari Carmen Suárez-Figueroa, Isam Diab, Álvaro González, and Jesica Rivero-Espinosa. Towards an automatic easy-to-read adaptation of morphological features in spanish texts. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 14142 LNCS:176–198, 2023
- [2] Alexander Vladislavovitch Dmitrijev, Elena Sergeevna Krupnova, and Anastasia Aleksandrovna Protopopova. Metaphors and analogies in the context of large language models. In *Scenarios, Fictions, and Imagined Possibilities in Science, Engineering, and Education*, volume 1203 LNNS, pages 326–341. Springer Science and Business Media Deutschland GmbH, 2024
- [7] Arthur Neidlein, Philipp Wiesenbach, and Katja Markert. An analysis of language models for metaphor recognition. *COLING 2020 - 28th International Conference on Computational Linguistics, Proceedings of the Conference*, pages 3722–3736, 2020
- [6] Emmy Liu, Chenxuan Cui, Kenneth Zheng, and Graham Neubig. Testing the ability of language models to interpret figurative language. *NAACL 2022 - 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference*, pages 4437–4452, 4 2022
- [5] Huiyuan Lai and Malvina Nissim. A survey on automatic generation of figurative language: From rule-based systems to large language models. *ACM Computing Surveys*, 56, 5 2024
- [3] Nicholas Ichien, Dušan Stamenković, and Keith J. Holyoak. Large language model displays emergent ability to interpret novel literary metaphors. *Metaphor and Symbol*, 39:296–309, 10 2024
- [4] Kaidi Jia, Yanxia Wu, Ming Liu, and Rongsheng Li. Curriculum-style data augmentation for llm-based metaphor detection. arXiv:2412.02956 [cs.CL], 12 2024

3 System Design & Methodology

3.1 Problem formulation

- Input: raw text
- Output: literalized, meaning-preserving text + intermediate explanations
- Challenge: Metaphors often encode meaning that cannot be fully literalized without loss; replacing them requires significant nuance and linguistic understanding.

3.2 Agentic pipeline

- Figurative span detection:
 - Input: Raw text (e.g., “Yesterday it was raining cats and dogs.”)
 - Output: JSON structure containing `metaphor_spans` (text, character start/end positions, confidence score).
- Interpretation / explanation
- Literal replacement generation
- Self-verification & revision
- Self-learning through addition to RAG?

3.3 Baselines

- Monolithic single-prompt LLM
- Detection-only prompting
- Naïve paraphrasing

4 RQ1: Detection & Interpretation

To what extent can a zero-/one-shot agentic LLM pipeline accurately detect and interpret idiomatic and metaphorical expressions in context?

4.1 Experiments

- Idioms: SemEval-2022 (detection / idiomaticity) [11]
- Metaphors: VUA (token-level metaphor detection) [8]
- Note: Initial prototyping focuses on the detection module as a foundation for subsequent replacement tasks.
- Compare:
 - Agentic pipeline
 - Single-prompt LLM
 - Reported supervised benchmarks (from literature)

4.2 Metrics

- Accuracy / F1 (detection)
- Qualitative interpretation correctness

5 RQ2: Agentic Decomposition

Does agentic decomposition improve figurative language handling compared to monolithic prompting?

5.1 Experiments

- Same inputs, different architectures
- Ablation:
 - no explanation step
 - no self-verification
 - full agentic pipeline

5.2 Metrics

- Replacement quality (see Chapter 6)
- Error types
- Failure traceability

6 RQ3: Replacement Quality

How effectively can figurative expressions be replaced with literal, meaning-preserving alternatives?

6.1 Idioms

- SemEval-2022 Subtask B-style evaluation
- Semantic similarity to gold paraphrases

6.2 Metaphors

- Custom evaluation protocol

6.3 Analysis

- Idioms vs metaphors
- Trade-off: readability vs semantic fidelity

7 RQ4: Observability & Error Analysis

How observable and debuggable are agentic systems compared to end-to-end prompting?

- Case studies using LangSmith traces
- Error localization
- Correlation between explanation quality and outcome quality

8 Discussion & Limitations

- What agentic systems can/cannot do
- Where metaphors fundamentally resist literalization
- Limitations:
 - Conceptual: metaphors encoding meaning that cannot be fully literalized; replacement oversimplifying nuance
 - Evaluation limits: Subjectivity of human evaluation; no gold standard for metaphor replacement
 - Model dependence: Results vary across LLM providers
 - Scope limits: Sentence/MWE level focus; no discourse-wide tracking; primarily English

9 Conclusion & Future Work

- Summary of findings
- Implications for NLP accessibility
- Directions: multilinguality, discourse-level metaphors, human-in-the-loop

References

- [1] A. Author. Sample article title. *Journal Name*, 1(1):1–10, 2024.
- [2] Alexander Vladislavovitch Dmitrijev, Elena Sergeevna Krupnova, and Anastasia Aleksandrovna Protopopova. Metaphors and analogies in the context of large language models. In *Scenarios, Fictions, and Imagined Possibilities in Science, Engineering, and Education*, volume 1203 LNNS, pages 326–341. Springer Science and Business Media Deutschland GmbH, 2024.
- [3] Nicholas Ichien, Dušan Stamenković, and Keith J. Holyoak. Large language model displays emergent ability to interpret novel literary metaphors. *Metaphor and Symbol*, 39:296–309, 10 2024.
- [4] Kaidi Jia, Yanxia Wu, Ming Liu, and Rongsheng Li. Curriculum-style data augmentation for llm-based metaphor detection. arXiv:2412.02956 [cs.CL], 12 2024.
- [5] Huiyuan Lai and Malvina Nissim. A survey on automatic generation of figurative language: From rule-based systems to large language models. *ACM Computing Surveys*, 56, 5 2024.
- [6] Emmy Liu, Chenxuan Cui, Kenneth Zheng, and Graham Neubig. Testing the ability of language models to interpret figurative language. *NAACL 2022 - 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference*, pages 4437–4452, 4 2022.
- [7] Arthur Neidlein, Philipp Wiesenbach, and Katja Markert. An analysis of language models for metaphor recognition. *COLING 2020 - 28th International Conference on Computational Linguistics, Proceedings of the Conference*, pages 3722–3736, 2020.
- [8] Gerard J Steen, Aletta G Dorst, J Berenike Herrmann, Anna A Kaal, and Tina Krennmayr. VU amsterdan metaphor corpus, 2010. Oxford Text Archive.
- [9] Mari Carmen Suárez-Figueroa, Isam Diab, Edna Ruckhaus, and Isabel Cano. First steps in the development of a support application for easy-to-read adaptation. *Universal Access in the Information Society*, 23:365–377, 3 2024.
- [10] Mari Carmen Suárez-Figueroa, Isam Diab, Álvaro González, and Jesica Rivero-Espinosa. Towards an automatic easy-to-read adaptation of morphological features in spanish texts. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 14142 LNCS:176–198, 2023.
- [11] Harish Tayyar Madabushi, Edward Gow-Smith, Marcos Garcia, Carolina Scarton, Marco Idiart, and Aline Villavicencio. SemEval-2022 task 2: Multilingual idiomaticity detection and sentence embedding. In Guy Emerson, Natalie Schluter, Gabriel Stanovsky, Ritesh Kumar, Alexis Palmer, Nathan Schneider, Siddharth Singh, and Shyam Ratan, editors, *Proceedings of the 16th International Workshop on Semantic Evaluation (SemEval-2022)*, pages 107–121, Seattle, United States, July 2022. Association for Computational Linguistics.

10 Annex

...