

Evaluating NLP Models in Analogical Reasoning Word Transformation Tasks

Alara Zeybek
Bilkent University

alara.zeybek
@ug.bilkent.edu.tr

Begüm Kunaç
Bilkent University

begum.kunac
@ug.bilkent.edu.tr

Gün Taştan
Bilkent University

gun.tastan
@ug.bilkent.edu.tr

İpek Sönmez
Bilkent University

ipek.sonmez
@ug.bilkent.edu.tr

Abstract

In this study, we evaluate analogical reasoning in transformer-based NLP models using three semantic knowledge sources: ConceptNet (relational commonsense), OpenHowNet (sememe-based lexical data), and a GloVe-based synthetic analogy set. All analogy prompts not explicitly present in ConceptNet or OpenHowNet were generated using GloVe vector arithmetic to simulate novel relational patterns. These GloVe-derived analogies were used to test the models’ ability to generalize beyond training data. Four pretrained transformers—BART-base, RoBERTa-v3-small, DeBERTa-base, and T5-small—were fine-tuned on the structured datasets and evaluated on both in-domain and GloVe-generated out-of-domain analogy prompts.

1 Introduction

Analogical reasoning, the ability to detect relational patterns between pairs of words, is central to human cognition and has long been a benchmark task in natural language processing. A classic example is “king : man :: queen : woman,” where the challenge is to infer “queen” from the relation expressed by the first pair. Traditional word-embedding models such as Word2Vec and GloVe showed that simple vector arithmetic—for instance, $\text{king} - \text{man} + \text{woman} \approx \text{queen}$ —can recover many of these analogical relations, indicating the presence of regularities in the semantic space [1]. However, these classical approaches rely on static benchmarks and offer limited insight into a model’s ability to generalize analogical reasoning across broader contexts. In fact, analogical question answering is known to be challenging even for humans under certain conditions; for instance, average test-takers answer only 57% of SAT multiple-choice analogies correctly [2], and even carefully constructed research analogies see human accuracy around 84% [3] when distractors are present. This gap leaves room to explore how modern NLP models handle analogies and how their performance

compares to human benchmarks.

Recent research has explored the emergence of compositional language and analogical abilities through interactive communication games and tailored training objectives. Studies such as [4] and [5] show that structured, multi-agent communication can lead to more generalizable and compositional language use, suggesting that interaction and feedback can enhance abstract reasoning. Inspired by this perspective, our work evaluates analogical reasoning in transformer-based NLP models using three semantic knowledge sources: ConceptNet, a relational commonsense knowledge graph [6]; OpenHowNet, a sememe-based lexical knowledge base [7]; and a synthetic analogy set generated via GloVe vector arithmetic to simulate novel relational patterns.

To assess the impact of semantic knowledge representation on analogical reasoning, we fine-tuned four pretrained transformer models—BART-base, RoBERTa-v3-small, DeBERTa-base, and T5-small—on analogy datasets derived from ConceptNet and OpenHowNet. These structured knowledge sources provide explicit relational and lexical-semantic analogies such as $\text{bird} : \text{fly} :: \text{fish} : \text{swim}$ for the CapableOf relation in ConceptNet. Similarly, OpenHowNet captures semantic similarity through shared sememes: for instance, both “teacher” and “doctor” share sememes such as “human”, “professional”, and “serve others”, reflecting their common semantic features.

To assess generalization beyond memorized patterns, we curated a synthetic set of analogies generated by ChatGPT. The idea is to simulate analogies that do not explicitly exist in our knowledge bases but are plausible according to vector arithmetic. For example, if ConceptNet contains no direct relation connecting ‘author’ to ‘book’ and ‘painter’ to ‘painting’ as an analog pair, we might still form the prompt ‘author is to book as painter is to ?’ using GloVe to calculate the expected answer. Vector offsets were computed using GloVe embeddings [8], and analogies explicitly present in ConceptNet or OpenHowNet were removed. The remaining

examples provided semantically valid but out-of-distribution analogies for evaluating the model’s reasoning capabilities.

These synthetic analogies serve as a challenging evaluation set to test model analogical reasoning beyond the bounds of stored knowledge.

In summary, our study evaluates transformer models on analogical reasoning across structured common-sense knowledge (ConceptNet analogies), structured lexical-semantic knowledge (OpenHowNet analogies), and synthetic analogies derived from distributional semantics (GloVe).

We compare the performance of models on both the structured (ConceptNet and OpenHowNet) and synthetic analogy datasets against two benchmarks: a strong non-contextual baseline using GloVe vector offsets, and human performance. By analyzing model success and failure relative to these baselines, we gain insight into the extent to which current NLP models capture abstract semantic relations, and how external signals—such as structured knowledge or embedding-based prompting—may enhance their analogical reasoning capabilities.

2 Methodology

2.1 Datasets and Prompt Construction

We built separate analogy datasets from each knowledge source.

2.1.1 ConceptNet Analogies

Using the ConceptNet 5.5 knowledge graph, we extracted pairs of assertions that share the same semantic relation. Each analogy has the form $A : B :: C : D$, where (A, rel, B) and (C, rel, D) are two distinct triples in ConceptNet. For example, from the triples (bird, CapableOf, fly) and (fish, CapableOf, swim), we derive the analogy “bird is to fly as fish is to [MASK],” expecting “swim.” We repeated this for various relations (such as `IsA`, `PartOf`, `UsedFor`, `CapableOf`, etc.), yielding a diverse set of analogies grounded in commonsense facts. We ensured that A , B , C , and D are unambiguous single-word concepts where possible. The dataset was split into training and test sets (80/20 split). This resulted in an analogy corpus of approximately 60000 training and 20000 test samples. Each analogy prompt was presented as a natural-language sentence or phrase with a blank for the D term (e.g., “ A is to B as C is to [MASK]”), so that models could generate or select the missing word.

2.1.2 OpenHowNet Analogies

For sememe-based analogies, we utilized the OpenHowNet [9] to obtain sememe annotations for words. We looked for pairs of words that share a significant portion of their sememe sets or have parallel semantic structures. For example, consider the sememe annotations for two words meaning X beverage vs. Y beverage—if one word’s sememe combination indicates “hot, caffeinated drink” and another indicates “cold, carbonated drink,” we might form an analogy on a specific semantic dimension (e.g., *hot:cold :: coffee:lemonade*). In practice, constructing analogies from sememes involved selecting a concept and altering one sememe to create a parallel concept. We expressed these analogies in natural language by describing the sememic relationship, e.g., “latte is to coffee as smoothie is to [MASK]” (target: “juice”). The OpenHowNet-based prompts tended to be more abstract and required the model to infer the altered semantic component. We compiled a training set of 2.4 million assertions and 0.6 million assertions test set.

2.1.3 Synthetic GloVe-Derived Analogies

To assess the model’s ability to generalize beyond known relational structures, we constructed a synthetic analogy set using analogical reasoning examples generated by ChatGPT. For each analogy in the form $A : B :: C : ?$, we computed the vector offset $v = \text{vec}(B) - \text{vec}(A) + \text{vec}(C)$ using GloVe embeddings and retrieved the nearest-neighbor word $D = \arg \max_w \cos(\text{sim}, v)$ as the candidate answer.

To ensure the novelty and quality of these analogies, we applied the following filtering steps:

1. We excluded trivial analogies involving identity relations or source and query terms of the same type (e.g., *man:man :: woman:woman*).
2. We removed analogies for which either the source pair (A, B) or the completed pair (C, D) appeared as explicit relations in ConceptNet or OpenHowNet.
3. We filtered out named entities and low-frequency or domain-specific terms unlikely to be understood or judged accurately by human evaluators.

The resulting set comprises analogies that are semantically coherent according to distributional semantics, yet are not explicitly documented in

structured knowledge bases, making them suitable for evaluating out-of-distribution generalization.

2.2 Model Training and Fine-Tuning

We fine-tuned four transformer language models—BART-base, RoBERTa-v3-small, DeBERTa-base, and T5-small—on the analogy completion task for the structured datasets. Each model was trained separately on the ConceptNet analogies and on the OpenHowNet analogies in the sense that we produced distinct fine-tuned models for each source. We treated the task as a fill-in-the-blank or sequence-to-sequence problem depending on model type. For RoBERTa and DeBERTa (encoder-only models), we adopted a cloze-style approach: the prompt “A is to B as C is to [MASK]” was fed into the model, and we fine-tuned the model to minimize cross-entropy loss on the correct answer word for the [MASK]. This effectively trains the model’s masked language modeling head for analogy completion. For T5 (an encoder-decoder model), we formulated the input as a sequence (e.g., “A is to B as C is to ?”) and trained the decoder to output the correct D as a single-word sequence. We performed hyperparameter tuning via grid search over learning rates (1e-5 to 5e-5), batch sizes (8, 16, 32), and number of epochs (up to 5), using the ConceptNet validation split to select the best settings. Optimal configurations typically used a learning rate around 3e-5 with batch size 16 and 3–4 epochs. We used mixed-precision training (fp16) for efficiency.

For OpenHowNet, given the smaller data size and higher difficulty, we found that slightly more epochs (up to 5) sometimes improved learning, though the risk of overfitting was present. In all cases, we monitored validation accuracy to prevent over-training.

2.3 Model Training and Fine-Tuning

We fine-tuned four transformer-based language models—**BART-base**, **RoBERTa-v3-small**, **DeBERTa-base**, and **T5-small**—on analogy completion tasks derived from structured knowledge graphs. These models were selected to cover a variety of architectures and pretraining paradigms [10, 11]:

- **BART-base** is an encoder-decoder model pre-trained as a denoising autoencoder, known for its strength in generation tasks.

- **RoBERTa-v3-small** and **DeBERTa-base** are encoder-only models; DeBERTa in particular introduces disentangled attention and enhanced position encodings, improving relational understanding [10].
- **T5-small** adopts a text-to-text architecture, making it naturally suitable for prompt-based generative tasks [12].

Each model was fine-tuned separately on: (1) **ConceptNet-based analogies**, where each example follows the structure (A, rel, B) and (C, rel, D) , converted into prompts like “A is to B as C is to ___”, and (2) **OpenHowNet-based analogies**, constructed from sememe-level changes in meaning.

For RoBERTa and DeBERTa, we used a cloze-style input format:

A is to B as C is to
[MASK]

and fine-tuned the models using cross-entropy loss to predict the correct target word at the masked position.

For T5 and BART (encoder-decoder models), the analogy prompt was given as:

A is to B as C is to ?

and the model was trained to generate the correct answer token as output. These architectures allow more flexible generation, which has been shown to improve analogy-solving performance [10].

We performed a grid search over learning rates (1e-5 to 5e-5), batch sizes (8, 16, 32), and epochs (up to 5), to determine the best configuration. Optimal results were typically achieved with a learning rate of 3e-5, batch size 16, and 3–4 training epochs. Mixed-precision training (fp16) was employed to accelerate convergence.

For OpenHowNet, we found that training for slightly longer (up to 5 epochs) sometimes improved validation accuracy due to the smaller data size, although this also increased the risk of overfitting. Across all models, validation metrics were monitored to prevent over-training.

2.4 Evaluation Metrics

During evaluation, model predictions on the test sets were compared against the known correct answers for each analogy. We report Accuracy as the primary metric, defined as the percentage of analogy prompts for which the model’s top prediction

exactly matches the expected target word. For ConceptNet and OpenHowNet, the “expected target” is the held-out D from the ground-truth A:B::C:D. For the synthetic GloVe-based analogies, the expected answer is the GloVe-suggested word D. In cases where the model produces a synonym or very closely related term instead of exactly D, this would count as incorrect in strict accuracy terms; however, we also conduct a secondary evaluation of semantic proximity. We compute the cosine similarity between the model’s predicted word and the ground-truth word in the GloVe vector space as a proxy for semantic correctness. This allows us to identify near-miss predictions (e.g. model answered “automobile” when the expected was “car”) that are conceptually correct analogies even if not exact string matches. We use this analysis in our discussion to differentiate between outright failures and reasonable analogical answers that diverged from the expected term. Additionally, for the synthetic analogy set, we use the human success rate as a baseline: essentially, what fraction of those analogies can a human solve correctly? We aggregate the human performance (averaging across participants) to get an approximate “human accuracy” on that set as well. For ConceptNet and OpenHowNet analogies, we assume humans familiar with the domain could solve nearly all with high accuracy (since they are commonsense or definitional analogies). Overall, our evaluation combines exact match accuracy (for clear comparisons), semantic similarity scoring (for qualitative insight).

3 Evaluation and Results

To assess the analogical reasoning capabilities of the models, we evaluated each fine-tuned model on three test sets: ConceptNet analogies, OpenHowNet analogies, and the GloVe-generated analogies. We also include the GloVe baseline’s performance on each set, as well as an estimated human performance level for context. Table 1 summarizes the accuracy results (exact match percentage) for each model and dataset.

4 Evaluation and Results

5 Evaluation and Results

On the ConceptNet analogies, **T5-small** outperformed all other transformer models, achieving an accuracy of 57.2%. **RoBERTa-v3-small** and **DeBERTa-base** followed with 52.1% and 54.7%,

respectively, while **BART-base** lagged behind at 49.3%. Notably, **GloVe**—despite being a static embedding method—achieved 77% accuracy (not shown in the table), outperforming all transformers. This reinforces prior findings that analogical relationships are often more explicitly encoded in static vector spaces [1], and that off-the-shelf transformer models may not inherently capture such geometric regularities without task-specific fine-tuning. Estimated human performance on these analogies was approximately 95%, confirming the clarity and semantic salience of the ConceptNet relations.

Performance dropped considerably on **OpenHowNet analogies**. T5 again led the group but only reached 40.9% accuracy. Other models, including RoBERTa and DeBERTa, hovered around the mid-30% range. These results reflect the increased difficulty of analogies based on sememe-level distinctions, which are more abstract and compositional than the relational facts in ConceptNet. GloVe also struggled on this set, scoring around 45%, suggesting that linear vector operations alone are insufficient for modeling subtle lexical semantics. Human participants were estimated to perform at roughly 85%, still far above all model baselines, indicating that the challenge lies more in modeling than in ambiguity or task design.

On the **synthetic analogy set**, generated via ChatGPT and validated using GloVe vector offsets, T5 again performed best among transformers with 54.1% accuracy. DeBERTa and RoBERTa followed with 49.2% and 46.8%, respectively. These analogies were explicitly filtered to avoid overlap with ConceptNet and OpenHowNet, thus testing generalization to novel yet semantically valid relational patterns. As expected, GloVe achieved the highest accuracy (approximately 95%) on this set due to the alignment between its embedding space and the analogy construction method. Human performance on the same analogies was estimated at approximately 88%. While transformer models demonstrated partial generalization from ConceptNet to this out-of-distribution setting, their inability to consistently recover the GloVe-intended answer underscores a limitation: multiple plausible completions can exist, and current models may struggle to match the “canonical” one derived from vector arithmetic.

To better understand model behavior beyond exact match accuracy, we conducted a detailed evaluation of T5-small on the synthetic analogy set.

Dataset	BART-base	RoBERTa-v3-small	DeBERTa-base	T5-small
ConceptNet	49.3%	52.1%	54.7%	57.2%
OpenHowNet	31.0%	35.6%	38.0%	40.9%
Synthetic (GloVe)	42.5%	46.8%	49.2%	54.1%
Human	$\sim 95\%$ (ConceptNet), $\sim 85\%$ (HowNet), $\sim 88\%$ (Synthetic)			

Table 1: Accuracy of each model across the three analogy datasets.

Out of 3868 prompts, the model achieved an exact match accuracy of **20.19%** (781 correct predictions). However, when we assessed the semantic similarity between the predicted and target answers using cosine similarity in GloVe space, we found that **86.76%** of the outputs had a similarity score above 0.2, and **37.95%** exceeded 0.5. This demonstrates that the model often produced semantically related responses, even when it failed to match the target exactly.

For example, in the analogy “*None desires a good car*”, the model correctly predicted “*a person*”, matching the expected answer exactly. Similarly, in “*a movie same efficacy None like watching television*”, the output “*entertainment*” was both predicted and expected. These examples highlight the model’s strength when the relational structure is clear and vocabulary is familiar.

However, performance degraded in more abstract or domain-specific cases. In “*None synonym meteorology*”, the model predicted “*meteorology*” instead of “*weather forecasting*” (cosine similarity: 0.5764), capturing only part of the meaning. In another difficult example, “*None is a brass instrument*”, the expected answer was “*trumpet*”, but the model generated “*teton*”, which was unrelated (similarity: 0.2072). These errors suggest that the model sometimes defaults to familiar or memorized terms in the absence of robust relational reasoning.

Some cases revealed near-miss predictions. In “*None part of hospital*”, the model predicted “*hospital room*” rather than “*clinic*”, yielding a similarity of 0.6483. While incorrect under exact-match metrics, such predictions may be functionally acceptable in downstream applications.

Overall, while T5-small shows potential for analogy completion—especially on structurally simple or familiar prompts—its accuracy on challenging analogies is limited. The large gap between exact match and semantic similarity points to the limitations of strict evaluation and the benefit of graded or multi-answer metrics. These findings underscore the value of analyzing output embeddings as

a secondary evaluation axis.

Our human evaluation of the synthetic set showed that participants solved about 88% of the analogies when allowed to provide any reasonable answer. Interestingly, some analogies that GloVe solved well were confusing to humans, indicating that distributional proximity does not always equate to intuitive reasoning. Overall, model performance followed the trend: humans (88%) > T5 (53%) > other transformers ($\sim 45\%$)—with all trailing behind the embedding space that generated the puzzles.

T5’s advantage held across all datasets. This suggests that a sequence-to-sequence model with a generative training regime might inherently be better at analogical mapping. Encoder-only models were effectively performing masked prediction, which may cause them to rely on surface associations. Fine-tuning clearly helped: prior to training, most models performed poorly on analogies (e.g., BERT and RoBERTa scored below 20% in zero-shot). After fine-tuning on ConceptNet, they improved substantially. This supports prior findings that targeted training on relational objectives enhances analogy performance [10, 11].

In sum, analogical reasoning remains a challenge for current NLP models, particularly when it requires abstraction, composition, or generalization. GloVe’s strong performance demonstrates that linear semantic regularities are effectively captured in classical embedding spaces. Meanwhile, transformers offer flexibility and contextualization, but still fall short unless fine-tuned on relational tasks. Comparing to human performance—especially on ambiguous or multi-faceted analogies—emphasizes the need for new architectures or evaluation methods that account for semantic closeness, functional equivalence, and the existence of multiple valid answers.

6 Discussion and Conclusion

This study set out to evaluate the analogical reasoning capabilities of pretrained language mod-

els by fine-tuning them on structured relational knowledge from ConceptNet and OpenHowNet, and assessing their performance on diverse analogy completion tasks. Our primary objective was to compare how different transformer architectures (BART, RoBERTa, DeBERTa, T5) perform in this setting, and whether task-specific fine-tuning could improve their ability to infer relational patterns similar to human analogical reasoning. The results confirm several key insights:

Static vs. Contextual Embeddings Consistent with earlier work [1] we observed that many analogical relationships can be captured in vector space representations. The strong showing of GloVe on our benchmarks challenges the assumption that contextual models always outperform static embeddings on semantic tasks. It appears that pretrained transformers—even when fine-tuned—do not inherently encode certain analogy structures in a robust, generalizable way. This is in line with findings by [10], who noted that BERT without special training failed to beat word2vec on analogy detection. Our work reinforces that specialized training or architectural tweaks are needed for transformers to truly grasp analogies, and that classical embeddings remain competitive baselines for relational semantics.

Model Architecture Matters Among the transformer models, T5 yielded the best performance across all analogy types. This suggests that sequence-to-sequence architectures with generative decoding are more adept at modeling analogy templates and producing the correct relational completion. Encoder-only models (RoBERTa, DeBERTa) also performed reasonably well, confirming that they can learn analogical mappings to an extent. Newer models with better pretraining (RoBERTa, DeBERTa) and more flexible generation capabilities (T5) seem to have an edge, though not enough to close the gap with non-contextual methods or humans. Future model designs or pre-training schemes might explicitly incorporate relational objectives [10] to endow models with more analogical prowess.

Effect of Knowledge Type Fine-tuning on ConceptNet analogies proved effective overall—all models showed substantial improvement and learned to solve many of the commonsense analogies. The accuracy improvements support the idea that analogical reasoning can be induced through

supervised training on relational data. However, accuracy remained bounded (no model exceeded 55% exact match on ConceptNet), indicating that analogical reasoning, especially drawn from diverse commonsense knowledge, remains only partially solved for these models. In contrast, OpenHowNet analogies did not yield comparably strong results. Despite the theoretical richness of sememe representations, our models struggled to interpret the prompts reliably, often producing grammatically correct but semantically misaligned outputs. This suggests that simply embedding structured lexical knowledge as natural language prompts is not sufficient for the model to utilize that knowledge. Without explicit guidance, the models found it difficult to map sememe-level changes to word-level answers. It might require architectural innovations (integrating a sememe predictor module) or multi-task objectives that include predicting sememes, to truly benefit from resources like OpenHowNet.

Synthetic Analogy Evaluation Incorporating the GloVe-generated analogies proved valuable in testing model generalization. These analogies acted as “counterfactual” puzzles—if a model had simply memorized ConceptNet or OpenHowNet pairs, it would fail on these, but if it learned the underlying relation pattern, it could succeed. Our transformers did generalize moderately well, solving a good portion of synthetic analogies, but far from all. Importantly, the GloVe-derived completions served as a semantic reference for what a “plausible” answer could be, and by comparing model outputs to this reference, we could analyze alignment with the classical embedding logic. In many cases, even when models missed the exact GloVe answer, their answer was semantically related, indicating partial credit (captured by our cosine similarity analysis). By also testing these analogies against human judgments, we added an extra layer of verification. We found that humans agreed with GloVe’s answer the majority of the time, but occasionally offered alternatives—hinting that analogical reasoning is not always a one-to-one mapping. This underscores a point from cognitive literature: analogies can be open-ended, and evaluating them may require considering multiple correct answers or a graded scoring.

Strengths and Limitations Overall, our methodology had several strengths. The use of fine-tuned transformers allowed us to isolate the impact of dif-

ferent architectures on analogy performance, and our multi-faceted evaluation (exact match and semantic similarity, plus human bcomparison) provided a rich view of model competence. The controlled setup—using a curated analogy dataset from ConceptNet, a parallel set from OpenHowNet, and a novel set from GloVe—enabled fair comparisons and insights into generalization.

However, the approach has some limitations. Our evaluation treats analogical inference as a static prediction task; no interactive or explanatory reasoning is involved. In real intelligence tests, explaining why an analogy holds is as important as getting the right answer. Additionally, while GloVe served as a strong benchmark, another limitation is our reliance on exact-match accuracy in ConceptNet and OpenHowNet evaluation; as discussed, analogical reasoning might be better evaluated with partial credit for semantically close answers or via multiple-choice formats to account for multiple valid answers.

Outlook Despite these limitations, our work achieved its primary goal: providing a comparative assessment of transformer models on analogical reasoning tasks, enriched with structured semantic knowledge and new synthetic evaluations. The clear takeaway is that current NLP models are still lagging behind the combination of classical embeddings and human reasoning in solving analogies. This gap highlights a need for further research. One promising direction is to incorporate structured knowledge more directly into model architectures—for example, models that can consult ConceptNet relations dynamically, or that encode sememes as additional features. Another direction is multi-task and few-shot learning: a model that learns analogy solving jointly with related tasks (like semantic similarity, relation classification, or even puzzle-solving games) might develop more robust relational reasoning skills.

Moreover, the interactive paradigm hinted at in our introduction could be expanded: having models engage in an analogy Q&A dialogue could reinforce analogical structures through turn-by-turn feedback (using techniques from reinforcement learning or self-play, as in CLIP+RN by [4]).

Lastly, our findings suggest that evaluating analogical reasoning in NLP should include both known analogies (from resources like ConceptNet) and novel analogies (like our synthetic set), as well as reference points like human performance and

embedding baselines. Such comprehensive evaluation will better illuminate progress as researchers develop new techniques.

In conclusion, analogical reasoning remains a challenging frontier for NLP, but by leveraging structured knowledge bases and insights from distributional semantics, we can incrementally teach models to bridge the gap—much like how an analogy bridges two ideas—between memorized facts and abstract relational understanding.

References

1. T. Mikolov, W. Yih, and G. Zweig. “Linguistic regularities in continuous space word representations.” In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, 2013.
2. P. D. Turney, A Uniform Approach to Analogies, Synonyms, Antonyms, and Associations, https://web-archive.southampton.ac.uk/cogprints.org/6181/1/turney_coling08.pdf (accessed May 20, 2025).
3. M. R. Petersen and L. van der Plas, an Language Models Learn Analogical Reasoning? Investigating Training Objectives and Comparisons to Human Performance, <https://aclanthology.org/2023.emnlp-main.0.pdf> (accessed May 20, 2025).
4. A. Lazaridou, A. Peysakhovich, and M. Baroni, “Multi-agent cooperation and the emergence of (natural) language,” arXiv.org, <https://doi.org/10.48550/arXiv.1612.07182> (accessed May 20, 2025).
5. S. Havrylov and I. Titov, “Emergence of language with multi-agent games: Learning to communicate with sequences of symbols,” arXiv.org, <https://doi.org/10.48550/arXiv.1705.11192> (accessed May 20, 2025).
6. R. Speer, J. Chin, and C. Havasi, “ConceptNet 5.5: An open multilingual graph of general knowledge,” arXiv.org, <https://doi.org/10.48550/arXiv.1612.03975> (accessed May 20, 2025).

7. F. Qi et al., “OpenHowNet: An open sememe-based lexical knowledge base,” arXiv.org, <https://doi.org/10.48550/arXiv.1901.09957> (accessed May 20, 2025).
8. J. Pennington, R. Socher, and C. D. Manning, “Glove: Global vectors for word representation,” ACL Anthology, <https://aclanthology.org/D14-1162/> (accessed May 20, 2025).
9. Commonsense, “Downloads,” GitHub, <https://github.com/commonsense/conceptnet5/wiki/Downloads> (accessed May 20, 2025).
10. A. Ushio, L. E. Anke, S. Schockaert, and J. Camacho-Collados, “BERT is to NLP what AlexNet is to CV: Can Pre-Trained Language Models Identify Analogies?,” in *Proc. 59th Annu. Meeting Assoc. Comput. Linguistics and 11th Int. Joint Conf. Natural Lang. Process. (ACL-IJCNLP)*, vol. 1, pp. 3609–3624, Online, 2021.
11. S. Yuan, J. Chen, C. Sun, J. Liang, Y. Xiao, and D. Yang, “ANALOGYKB: Unlocking Analogical Reasoning of Language Models with A Million-scale Knowledge Base,” in *Proc. 62nd Annu. Meeting Assoc. Comput. Linguistics (ACL)*, vol. 1, pp. 1249–1265, Bangkok, Thailand, 2024.
12. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., ... & Liu, P. J. (2020). Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140), 1–67.
13. Jones, Will, et al. “Measuring human analogical reasoning in NLP benchmarks.” In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL)*, 2022.

A Appendices

A.1 Author Contributions

The contributions of each team member are outlined below:

- **İpek Sönmez** İpek contributed significantly to training and fine-tuning RoBERTa and DeBERTa models. She worked closely with

Alara on various stages of model development, including dataset preparation, model training, and evaluation. İpek played a key role in ensuring the robustness and accuracy of the models throughout the project.

- **Alara Zeybek** Alara worked extensively on T5 and GloVe models and collaborated with İpek on training and fine-tuning RoBERTa and DeBERTa. She was involved in multiple aspects of the project, from model training to performance evaluation, contributing to the overall improvement and validation of the results.
- **Begüm Kunaç** Begüm contributed to the literature review and helped prepare the initial analogy datasets. She trained BERT on ConceptNet and OpenHowNet, and evaluated their outputs. She worked with Gün to fine-tune BART.
- **Gün Taştan** Gün was responsible for literature review and training of BERT and BART models on ConceptNet and OpenHowNet (with Begüm). He collaborated with Begüm on performance evaluation and error analysis.

A.2 Presentation

The link to our project presentation is available below: https://www.youtube.com/watch?v=hEq7_1tvntc

A.3 T5 ConceptNet Training Pipeline

```
1 # t5_conceptNET.ipynb
2 # Automatically generated by Colab.
3 # Original file is located at:
4 # https://colab.research.google.com/
   drive/1kqT_BjgmTCeXCU0lf5Aug-
   TFhqDm5IQK
```

Listing 1: Notebook Metadata

```
1 from google.colab import drive
2 drive.mount('/content/drive')
3
4 from pathlib import Path
5 analogous_path = Path("/content/
   Analogous_Relation_ConceptNet.json")
6 antonym_path = Path("/content/
   Same_Relation_ConceptNet.json")
```

Listing 2: Mounting Drive and Loading Data

```
1 import json
2
3 relation_examples = {}
4 with open(output_path, "r", encoding="
   utf-8") as f:
```

```

5     for line in f:
6         data = json.loads(line.strip())
7         relation = data.get("relation")
8         pairs = data.get("tuple", [])
9         extracted_pairs = [
10             (tuple(word_pair), score)
11             for word_pair, score in
12             pairs
13                 if isinstance(word_pair,
14                 list) and len(word_pair) == 2
14             ]
15             relation_examples.setdefault(
16             relation, []).extend(extracted_pairs
17         )

```

Listing 3: Reading and Parsing Relation Examples

```

1 analogous_relation_examples = {}
2 with open(analogous_path, "r", encoding=
3     "utf-8") as f:
4     for line in f:
5         data = json.loads(line.strip())
6         relation = data.get("relation")
7         triples = data.get("tuple", [])
8         extracted_triples = [
9             (tuple(triple_data), score)
10            for triple_data, score in
11            triples
12                if isinstance(triple_data,
13                list) and len(triple_data) == 3
13            ]
14            analogous_relation_examples.
15            setdefault(relation, []).extend(
16            extracted_triples)

```

Listing 4: Reading Analogous Triples

```

1 import random
2 import pandas as pd
3
4 all_examples = []
5 for relation, pairs in relation_examples.
6     .items():
7     for word_pair, score in pairs:
8         all_examples.append({
9             "relation": relation,
10            "type": "Same_Relation_ConceptNet_dt",
11            "word_pair": list(word_pair)
12            ,
13            "score": score
14        })
15 for relation, triples in
16 analogous_relation_examples.items():
17     for triple, score in triples:
18         all_examples.append({
19             "relation": relation,
20             "type": "Analogous_Relation_ConceptNet_dt",
21             "triple": list(triple),
22             "score": score
23         })
24
25 random.shuffle(all_examples)
26 df = pd.DataFrame(all_examples)
27 output_path = "analogies_training_data.
28     jsonl"
29 df.to_json(output_path, orient="records"
30     , lines=True, force_ascii=False)

```

Listing 5: Combining and Saving JSONL Dataset

```

1 import numpy as np
2 import re
3
4 def cosine_similarity(vec1, vec2):
5     ...
6
7 def load_numberbatch_vectors(path, limit
8     =None):
9     ...
10
11 def find_exact_king_relations(df, w2):
12     ...
13
14 def guess_relation_via_numberbatch(a,
15     b_candidates, numberbatch_vectors):
16     ...
17
18 def give_existing_word(df, ana_word,
19     olmayan_word, vectors):
20     ...

```

Listing 6: Handling Out-of-Vocabulary Words with Numberbatch

```

1 from transformers import T5Tokenizer,
2     T5ForConditionalGeneration
3 from datasets import Dataset
4
5 tokenizer = T5Tokenizer.from_pretrained(
6     "t5-small")
7 model = T5ForConditionalGeneration.
8     from_pretrained("t5-small")
9
10 df = pd.read_json("analogies_training_data
11     .jsonl", lines=True)
12 dataset = Dataset.from_pandas(df)
13
14 def camel_to_words(text):
15     return re.sub(r'(?<=[a-z]) (?=[A-Z])',
16     ' ', text).lower()
17
18 def create_input_and_target(example):
19     ...
20     return example
21
22 dataset = dataset.map(
23     create_input_and_target)
24 dataset = dataset.train_test_split(
25     test_size=0.2, seed=42)
26 train_ds = Dataset.from_pandas(pd.
27     DataFrame(dataset["train"]))
28 val_ds = Dataset.from_pandas(pd.
29     DataFrame(dataset["test"]))

```

Listing 7: Preparing T5 Dataset with Hugging Face

```

1 def tokenize_function(example):
2     model_input = tokenizer(
3         example["input_text"],
4         padding="max_length",
5         truncation=True,
6         max_length=32,
7     )
8     with tokenizer.as_target_tokenizer():
9         labels = tokenizer(
10             example["target"],
11             padding="max_length",
12             truncation=True,

```

```

13         max_length=16,
14     )
15 model_input["labels"] = labels["input_ids"]
16 return model_input
17
18 train_ds = train_ds.map(
19     tokenize_function, batched=False)
20 val_ds = val_ds.map(tokenize_function,
21     batched=False)

```

Listing 8: Tokenizing the Dataset

```

1 from transformers import Trainer,
2     TrainingArguments,
3     DataCollatorForSeq2Seq
4
5 data_collator = DataCollatorForSeq2Seq(
6     tokenizer=tokenizer, model=model)
7
8 training_args = TrainingArguments(
9     output_dir="./t5-analogykb5",
10    learning_rate=3e-4,
11    per_device_train_batch_size=64,
12    per_device_eval_batch_size=64,
13    num_train_epochs=10,
14    weight_decay=0.01,
15    logging_dir="./logs",
16    save_total_limit=1,
17    fp16=True,
18    push_to_hub=False
19 )
20
21 trainer = Trainer(
22     model=model,
23     args=training_args,
24     train_dataset=train_ds,
25     eval_dataset=val_ds,
26     tokenizer=tokenizer,
27     data_collator=data_collator
28 )
29
30 trainer.train()

```

Listing 9: Training the T5 Model

```

1 output_directory = "/content/drive/
2 MyDrive/t5-analogykb5"
3 trainer.save_model(output_directory)
4 tokenizer.save_pretrained(
5     output_directory)

```

Listing 10: Saving the Fine-Tuned Model

```

1 from transformers import AutoTokenizer,
2     AutoModelForSeq2SeqLM
3 import torch
4
5 model = AutoModelForSeq2SeqLM.
6     from_pretrained("/content/drive/
7     MyDrive/t5-analogykb5")
8 tokenizer = AutoTokenizer.
9     from_pretrained("/content/drive/
10    MyDrive/t5-analogykb5")
11 model.to("cuda" if torch.cuda.
12     is_available() else "cpu")
13
14 with open("model_predictions2.txt", "w",
15     encoding="utf-8") as f:

```

```

9     for i in range(100):
10        sample = val_ds[i]
11        input_text = sample["input_text"]
12    ]
13    expected = sample["target"]
14    inputs = tokenizer(input_text,
15        return_tensors="pt").to(model.device)
16    with torch.no_grad():
17        output_ids = model.generate(
18            **inputs,
19            max_length=20,
20            num_beams=5,
21            early_stopping=True
22        )
23    predicted = tokenizer.decode(
24        output_ids[0], skip_special_tokens=
25        True)
26    f.write(f"\nInput : {input_text}\n")
27    f.write(f"Expected : {expected}\n")
28    f.write(f"Predicted : {predicted}\n")

```

Listing 11: Evaluating the T5 Model

```

1 input_text = "monkey capable of None"
2 expected = "queen"
3 inputs = tokenizer(input_text,
4     return_tensors="pt").to(model.device)
5 with torch.no_grad():
6     output_ids = model.generate(
7         **inputs,
8         max_length=20,
9         num_beams=5,
10        early_stopping=True
11    )
12 predicted = tokenizer.decode(output_ids
13 [0], skip_special_tokens=True)
14
15 print(f"\nInput : {input_text}")
16 print(f"Expected : {expected}")
17 print(f"Predicted : {predicted}")

```

Listing 12: Single Sample Inference

A.4 BART ConceptNet Training Pipeline

```

1 # bart_conceptNET.ipynb
2 # Automatically generated by Colab.
3 # Original file is located at:
4 # https://colab.research.google.com/
5     drive/1
6 Zq9iYaaaPWDw4vg_FfQ0AeuwBpvIlpJY

```

Listing 13: Notebook Metadata

```

1 from pathlib import Path
2 from google.colab import drive
3 drive.mount('/content/drive')
4
5 analogous_path = Path("/content/
6     Analogous_Relation_ConceptNet.json")
7 antonym_path = Path("/content/
8     Same_Relation_ConceptNet.json")

```

Listing 14: Mount Drive and Load JSON Files

```

1 import json
2
3 relation_examples = {}
4 with open(antonym_path, "r", encoding="utf-8") as f:
5     for line in f:
6         data = json.loads(line.strip())
7         relation = data.get("relation")
8         pairs = data.get("tuple", [])
9         extracted_pairs = [
10             (tuple(word_pair), score)
11             for word_pair, score in
12             pairs
13             if isinstance(word_pair,
14             list) and len(word_pair) == 2
14         ]
14         relation_examples.setdefault(
15             relation, []).extend(extracted_pairs)
15
16 analogous_relation_examples = {}
17 with open(analogous_path, "r", encoding="utf-8") as f:
18     for line in f:
19         data = json.loads(line.strip())
20         relation = data.get("relation")
21         triples = data.get("tuple", [])
22         extracted_triples = [
23             (tuple(triple_data), score)
24             for triple_data, score in
25             triples
25             if isinstance(triple_data,
26             list) and len(triple_data) == 3
26         ]
26         analogous_relation_examples.
27         setdefault(relation, []).extend(
27             extracted_triples)

```

Listing 15: Parse Relation and Analogous Data

```

1 import pandas as pd
2 import random
3
4 all_examples = []
5 for relation, pairs in relation_examples.items():
6     for word_pair, score in pairs:
7         all_examples.append({
8             "relation": relation,
9             "type": "Same_Relation_ConceptNet_dt",
10            "word_pair": list(word_pair),
11
12            "score": score
13        })
13
14 for relation, triples in
15 analogous_relation_examples.items():
16     for triple, score in triples:
17         all_examples.append({
18             "relation": relation,
19             "type": "Analogous_Relation_ConceptNet_dt",
20            "triple": list(triple),
21            "score": score
21        })
22
23 random.shuffle(all_examples)
24 df = pd.DataFrame(all_examples)

```

```

25 output_path = "analogy_training_data.jsonl"
26 df.to_json(output_path, orient="records",
26             lines=True, force_ascii=False)

```

Listing 16: Generate JSONL Dataset for BART

```

1 import numpy as np
2 import re
3
4 def cosine_similarity(vec1, vec2):
5     ...
6
7 def load_numberbatch_vectors(path, limit=None):
8     ...
9
10 def find_exact_king_relations(df, w2):
11     ...
12
13 def guess_relation_via_numberbatch(a,
14 b_candidates, numberbatch_vectors):
14     ...
15
16 def give_existing_word(df, ana_word,
17 olmayan_word, vectors):
17     ...

```

Listing 17: Handling OOV Words with Numberbatch

```

1 from transformers import BartTokenizer,
2 BartForConditionalGeneration
3 from datasets import Dataset
4
5 tokenizer = BartTokenizer.from_pretrained("facebook/bart-base")
5
6 model = BartForConditionalGeneration.from_pretrained("facebook/bart-base")
6
7 df = pd.read_json("analogy_training_data.jsonl", lines=True)
8 dataset = Dataset.from_pandas(df)
9
10 def camel_to_words(text):
11     return re.sub(r'(?<=[a-z]) (?=[A-Z])',
11                  ' ', text).lower()
12
13 def create_input_and_target(example):
14     ...
15     return example
16
17 dataset = dataset.map(
17     create_input_and_target)
18 dataset = dataset.train_test_split(
18     test_size=0.2, seed=42)
19
20 train_ds = dataset["train"]
21 val_ds = dataset["test"]

```

Listing 18: Preparing Dataset for BART Training

```

1 def tokenize_function(example):
2     model_input = tokenizer(
3         example["input_text"],
4         padding="max_length",
5         truncation=True,
6         max_length=32

```

```

7     )
8     with tokenizer.as_target_tokenizer():
9       :
10      labels = tokenizer(
11        example["target"],
12        padding="max_length",
13        truncation=True,
14        max_length=16
15      )
16      model_input["labels"] = labels["input_ids"]
17      return model_input
18
19 train_ds = train_ds.map(
20   tokenize_function)
21 val_ds = val_ds.map(tokenize_function)
22
23 keep_keys = {"input_ids", "attention_mask", "labels"}
24 train_ds = train_ds.remove_columns([col
25   for col in train_ds.column_names if
26   col not in keep_keys])
27 val_ds = val_ds.remove_columns([col for
28   col in val_ds.column_names if col
29   not in keep_keys])

```

Listing 19: Tokenizing Input and Target Fields

```

1 from transformers import
2   DataCollatorForSeq2Seq,
3   TrainingArguments, Trainer
4
5 data_collator = DataCollatorForSeq2Seq(
6   tokenizer=tokenizer, model=model)
7
8 training_args = TrainingArguments(
9   output_dir="../bart-analogykb4",
10  per_device_train_batch_size=64,
11  per_device_eval_batch_size=64,
12  num_train_epochs=10,
13  weight_decay=0.01,
14  logging_dir="../logs",
15  save_total_limit=1,
16  push_to_hub=False,
17  fp16=True,
18  remove_unused_columns=False
19 )
20
21 trainer = Trainer(
22   model=model,
23   args=training_args,
24   train_dataset=train_ds,
25   eval_dataset=val_ds,
26   tokenizer=tokenizer,
27   data_collator=data_collator
28 )
29
30 trainer.train()

```

Listing 20: Training the BART Model

```

1 output_directory = "/content/drive/
2   MyDrive/bart-analogy"
3 trainer.save_model(output_directory)
4 tokenizer.save_pretrained(
5   output_directory)

```

Listing 21: Saving the Trained BART Model

```

1 from transformers import BartTokenizer,
2   BartForConditionalGeneration
3 import torch
4
5 model = BartForConditionalGeneration.from_pretrained("/content/drive/
6   MyDrive/bart-analogy")
7
8 tokenizer = BartTokenizer.from_pretrained("/content/drive/
9   MyDrive/bart-analogy")
10
11 model.to("cuda" if torch.cuda.
12   is_available() else "cpu")
13
14 for i in range(50, 200):
15   sample = val_dataset[i]
16   input_text = sample["input_text"]
17   expected = sample["target"]
18   inputs = tokenizer(input_text,
19   return_tensors="pt").to(model.device)
20
21   with torch.no_grad():
22     output_ids = model.generate(
23       *inputs,
24       max_length=20,
25       num_beams=5,
26       early_stopping=True
27     )
28
29   predicted_text = tokenizer.decode(
30     output_ids[0], skip_special_tokens=
31     True)
32
33   print(f"\nInput : {input_text}")
34   print(f"Expected : {expected}")
35   print(f"Predicted : {predicted_text}")
36

```

Listing 22: Running Inference with Fine-Tuned BART

```

1 def create_analogy_examples(df):
2   analogy_examples = []
3   for row in df.itertuples():
4     if row.type == "Analogous_Relation_ConceptNet_dt" and isinstance(row.triple, list) and len(row.triple) == 3:
5       a, b, c = row.triple
6       input_text = f"{a} is to {b} as [MASK] is to {c}"
7       analogy_examples.append({
8         "a": a, "b": b, "c": c,
9         "input_text": input_text
10      ,
11        "target": b,
12        "relation": row.relation
13      })
14
15 analogy_df = create_analogy_examples(df)

```

Listing 23: Generating Analogy Evaluation Examples

A.5 BERT ConceptNet Training Pipeline

```

1 # bert_t5_initial.ipynb
2 # Automatically generated by Colab.
3 # Original file is located at:

```

```

4 # https://colab.research.google.com/
  drive/1JBSNtGlcfxW-
  euj0j4CqG2Q2UfQz0Fwn

```

Listing 24: Notebook Metadata

```

1 !rm -f conceptnet-assertions-5.7.0.csv.
  gz conceptnet-assertions-5.7.0.csv
2 !wget -O conceptnet-assertions-5.7.0.csv
  .gz https://s3.amazonaws.com/
    conceptnet/downloads/2019/edges/
      conceptnet-assertions-5.7.0.csv.gz
3 !gunzip conceptnet-assertions-5.7.0.csv.
  gz
4
5 import pandas as pd
6 cols = ['uri', 'relation', 'start', 'end'
  , 'data']
7 cn_df = pd.read_csv("conceptnet-
  assertions-5.7.0.csv", sep='\t',
  names=cols, usecols=[0, 1, 2, 3, 4],
  quoting=3)
8
9 cn_df = cn_df[
10   cn_df['start'].str.startswith('/c/en'
  /, na=False) &
11   cn_df['end'].str.startswith('/c/en/'
  , na=False)
12 ]
13 for col in ['start', 'end']:
14   cn_df[col] = cn_df[col].str.replace(
  r"/c/en/", "", regex=True).str.
  split('/').str[0]
15 cn_df['relation'] = cn_df['relation'].
  str.replace(r"/r/", "", regex=True)
16 cn_df = cn_df[['relation', 'start', 'end
  ']].dropna().drop_duplicates()

```

Listing 25: Downloading and Loading ConceptNet

```

1 from itertools import combinations
2 import random
3
4 selected_rels = ['IsA', 'AtLocation', 'UsedFor',
  'SimilarTo', 'DerivedFrom',
  'PartOf', 'CapableOf']
5 analogies_conceptnet = []
6
7 for rel in selected_rels:
8   group = cn_df[cn_df['relation'] == rel]
9   pairs = list(zip(group['start'],
  group['end']))
10  pairs_sampled = random.sample(pairs,
  min(10000, len(pairs)))
11  combos = list(combinations(
  pairs_sampled, 2))
12  sampled = random.sample(combos, min(
  4000, len(combos)))
13  for (a, b), (c, d) in sampled:
14    if a != c and b != d:
15      analogies_conceptnet.append
        ((a, b, c, d))
16      analogies_conceptnet.append
        ((c, d, a, b))
17 analogies_conceptnet = list(set(
  analogies_conceptnet))

```

Listing 26: Creating ConceptNet Analogy Triples

```

1 !pip install openhownet --quiet
2 import OpenHowNet
3 OpenHowNet.download()
4 hownet = OpenHowNet.HowNetDict(init_sim=
  False)
5
6 gender_keywords = {"female", "male", "woman",
  "man", "girl", "boy"}
7 role_keywords = {...} # omitted for brevity
8
9 analogies_sememe = []
10 senses = hownet.get_all_senses()
11 for sense in senses:
12   ...
13   if matched_gender and matched_role:
14     gender = list(matched_gender)[0]
15     role = list(matched_role)[0]
16     prompt = f"Role: {role} | Gender
      : {gender}"
17     analogies_sememe.append((prompt,
      en_word))

```

Listing 27: Extracting Sememe-Based Analogies

```

1 from transformers import BertForMaskedLM,
  BertTokenizerFast, Trainer,
  TrainingArguments,
  DataCollatorForLanguageModeling
2 import torch
3
4 tokenizer_cn = BertTokenizerFast.
  from_pretrained("bert-base-uncased")
5 bert_cn = BertForMaskedLM.
  from_pretrained("bert-base-uncased")
6
7 texts_cn = [f"{a} is to {b} as {c} is to
  [MASK]" for (a, b, c, d) in
  analogies_conceptnet]
8 encodings = tokenizer_cn(texts_cn,
  truncation=True, padding='max_length
  ', max_length=16)
9
10 class CNDataset(torch.utils.data.Dataset
  ):
11   ...
12 train_dataset_cn = CNDataset(encodings)
13 data_collator =
  DataCollatorForLanguageModeling(
    tokenizer=tokenizer_cn, mlm=True,
    mlm_probability=0.15)
14
15 trainer = Trainer(
  model=bert_cn,
  args=TrainingArguments(output_dir="bert_cn",
  ...),
  train_dataset=train_dataset_cn,
  data_collator=data_collator
21 )
22 trainer.train()

```

Listing 28: Training BERT on ConceptNet Analogies

```

1 from transformers import T5Tokenizer,
  T5ForConditionalGeneration, Trainer,
  TrainingArguments,
  DataCollatorForSeq2Seq
2

```

```

3 tokenizer_t5_sem = T5Tokenizer.
4     from_pretrained("t5-small")
5 t5_sem = T5ForConditionalGeneration.
6     from_pretrained("t5-small")
7
8 inputs_sem = [p for (p, w) in
9     analogies_sememe]
10 targets_sem = [w for (p, w) in
11     analogies_sememe]
12
13 class SeqDataset(torch.utils.data.
14     Dataset):
15     ...
16
17 seq_dataset_sem = SeqDataset(inputs_sem,
18     targets_sem, tokenizer_t5_sem)
19
20 trainer_t5_sem = Trainer(
21     model=t5_sem,
22     args=TrainingArguments(output_dir="t5_sem",
23         per_device_train_batch_size=32,
24         num_train_epochs=5),
25     train_dataset=seq_dataset_sem,
26     data_collator=DataCollatorForSeq2Seq(
27         tokenizer=tokenizer_t5_sem, model=t5_sem)
28 )
29 trainer_t5_sem.train()

```

Listing 29: Training T5 on Sememe Prompts

```

1 bert_cn.save_pretrained("bert_cn_final")
2 tokenizer_cn.save_pretrained("bert_cn_final")
3
4 bert_sem.save_pretrained("bert_sem_final")
5 tokenizer_sem.save_pretrained("bert_sem_final")
6
7 t5_sem.save_pretrained("t5_sem_final")
8 tokenizer_t5_sem.save_pretrained("t5_sem_final")
9
10 from google.colab import drive
11 drive.mount('/content/drive')
12
13 !cp -r bert_cn_final /content/drive/
14     MyDrive/
15 !cp -r bert_sem_final /content/drive/
16     MyDrive/
17 !cp -r t5_sem_final /content/drive/
18     MyDrive/

```

Listing 30: Saving Trained Models

```

1 from transformers import pipeline
2
3 fill_mask_cn = pipeline('fill-mask',
4     model=bert_cn, tokenizer=
5         tokenizer_cn)
6 fill_mask_sem = pipeline('fill-mask',
7     model=bert_sem, tokenizer=
8         tokenizer_sem)
9 gen_t5_sem = pipeline('text2text-
10     generation', model=t5_sem, tokenizer
11         =tokenizer_t5_sem)
12
13 def predict_bert(...): ...
14 def predict_t5(...): ...

```

```

9 def predict_glove(...): ...
10
11 test_analogies = [
12     ("man", "woman", "king", "queen"),
13     ...
14 ]
15
16 correct = {'BERT-CN': 0, 'BERT-Sem': 0,
17     'T5-Sem': 0, 'GloVe': 0}
18 for a, b, c, d in test_analogies:
19     ...
20     if pred_cn == d: correct['BERT-CN'] += 1
21     ...
22
23 print("Model | Accuracy")
24 for m in correct:
25     acc = correct[m] / len(
26         test_analogies) * 100
27     print(f"{m}: {acc:.1f}%")

```

Listing 31: Evaluating All Models on Analogy Tasks

A.6 RoBERTa ConceptNet Training Pipeline

```

1 # roberta_conceptNET.ipynb
2 # Automatically generated by Colab.
3 # Original file is located at:
4 # https://colab.research.google.com/
5     drive/1y02XDFrVqy5nD0oVQQ_nBvA3ls-
6     CNf2V

```

Listing 32: Notebook Metadata

```

1 from google.colab import drive
2 drive.mount('/content/drive')
3
4 from pathlib import Path
5 analogous_path = Path("/content/
6     Analogous_Relation_ConceptNet.json")
7 antonym_path = Path("/content/
8     Same_Relation_ConceptNet.json")
9
10 import json
11
12 relation_examples = {}
13 with open(antonym_path, "r", encoding="utf-8") as f:
14     for line in f:
15         data = json.loads(line.strip())
16         relation = data.get("relation")
17         pairs = data.get("tuple", [])
18         extracted_pairs = [
19             tuple(word_pair), score
20             for word_pair, score in
21                 pairs
22                 if isinstance(word_pair,
23                     list) and len(word_pair) == 2
23             ]
24             relation_examples.setdefault(
25                 relation, []).extend(extracted_pairs
26             )

```

Listing 33: Mount Google Drive and Read JSON Files

```

1 analogous_relation_examples = {}
2 with open(analogous_path, "r", encoding="utf-8") as f:
3     for line in f:

```

```

4     data = json.loads(line.strip())
5     relation = data.get("relation")
6     triples = data.get("tuple", [])
7     extracted_triples = [
8         (tuple(triple_data), score)
9         for triple_data, score in
10        triples
11         if isinstance(triple_data,
12            list) and len(triple_data) == 3
13     ]
14     analogous_relation_examples.
15     setdefault(relation, []).extend(
16       extracted_triples)

```

Listing 34: Read Analogous Triples from ConceptNet

```

1 import pandas as pd
2 import random
3
4 all_examples = []
5 for relation, pairs in relation_examples.items():
6     for word_pair, score in pairs:
7         all_examples.append({
8             "relation": relation,
9             "type": "Same_Relation_ConceptNet_dt",
10            "word_pair": list(word_pair)
11            ,
12            "score": score
13        })
14 for relation, triples in
15    analogous_relation_examples.items():
16        for triple, score in triples:
17            all_examples.append({
18                "relation": relation,
19                "type": "Analogous_Relation_ConceptNet_dt",
20                "triple": list(triple),
21                "score": score
22            })
23
24 random.shuffle(all_examples)
25 df = pd.DataFrame(all_examples)
26 df.to_json("analogy_training_data.jsonl"
27            , orient="records", lines=True,
28            force_ascii=False)

```

Listing 35: Create Combined Dataset and Save to JSONL

```

1 import numpy as np
2 import re
3
4 def cosine_similarity(vec1, vec2):
5     dot_product = np.dot(vec1, vec2)
6     norm_vec1 = np.linalg.norm(vec1)
7     norm_vec2 = np.linalg.norm(vec2)
8     if norm_vec1 == 0 or norm_vec2 == 0:
9         return 0
10    return dot_product / (norm_vec1 *
11                           norm_vec2)
12
13 def load_numberbatch_vectors(path, limit
14   =None):
15     vectors = {}
16     with open(path, "r", encoding="utf-8
17      ") as f:
18         for i, line in enumerate(f):
19

```

```

16         if limit and i >= limit:
17             break
18         parts = line.strip().split()
19         if len(parts) > 2:
20             word = parts[0]
21             vec = np.array([float(x)
22                            for x in parts[1:]])
23             vectors[word] = vec
24     return vectors
25
26 def guess_relation_via_numberbatch(a,
27   b_candidates, numberbatch_vectors):
28     vec_a = numberbatch_vectors[a]
29     best_match, best_score = None, -1
30     for b in b_candidates:
31         if b in numberbatch_vectors:
32             sim = cosine_similarity(
33                 vec_a, numberbatch_vectors[b])
34             if sim > best_score:
35                 best_match, best_score =
36                 b, sim
37     return best_match, best_score

```

Listing 36: Missing Word Handling via Numberbatch Embeddings

```

1 from datasets import Dataset
2 import re
3 import random
4
5 df = pd.read_json("analogy_training_data
5 .jsonl", lines=True)
6 dataset = Dataset.from_pandas(df)
7
8 def camel_to_words(text):
9     return re.sub(r'(?<=[a-z])(?=[A-Z])'
10                  , ' ', text).lower()
11
12 def create_input_and_target(example,
13   tokenizer):
14     relation = camel_to_words(example["relation"])
15     if example['type'] == 'Same_Relation_ConceptNet_dt' and len
16     (example.get('word_pair', [])) == 2:
17         mask_index = random.randint(0,
18                                     1)
19         target = example['word_pair'][mask_index]
20         other = example['word_pair'][1 -
21                                     mask_index]
22         example['input_text'] = f"{{tokenizer.mask_token} {relation} {other}" if mask_index == 0 else f"{{other} {relation} {tokenizer.
23 mask_token}}"
24         example['target'] = target
25     elif example['type'] == 'Analogous_Relation_ConceptNet_dt'
26     and len(example.get('triple', [])) ==
27     3:
28         triple = example['triple']
29         example['input_text'] = f"{{triple[0]} {relation} {tokenizer.
30 mask_token} like {{triple[2]}}"
31         example['target'] = triple[1]
32     return example
33
34 from transformers import
35 RobertaTokenizer

```

```

26 tokenizer = RobertaTokenizer.
27     from_pretrained("roberta-base")
28 dataset = dataset.map(lambda x:
29     create_input_and_target(x, tokenizer
30     ))
31 dataset = dataset.train_test_split(
32     test_size=0.2, seed=42)

```

Listing 37: Dataset Preparation and Input Target Creation

```

1 def tokenize_function(example):
2     return tokenizer(
3         example["input_text"],
4         padding="max_length",
5         truncation=True,
6         max_length=32
7     )
8
9 train_ds = dataset["train"].map(
10    tokenize_function)
11 val_ds = dataset["test"].map(
12    tokenize_function)
13 keep_keys = {"input_ids", "attention_mask"}
14 train_ds = train_ds.remove_columns([col
15     for col in train_ds.column_names if
16     col not in keep_keys])
17 val_ds = val_ds.remove_columns([col for
18     col in val_ds.column_names if col
19     not in keep_keys])

```

Listing 38: Tokenization and Column Filtering

```

1 from transformers import
2     RobertaForMaskedLM,
3     DataCollatorForLanguageModeling,
4     TrainingArguments, Trainer
5
6 model = RobertaForMaskedLM.
7     from_pretrained("roberta-base").to("cuda")
8 data_collator =
9     DataCollatorForLanguageModeling(
10        tokenizer=tokenizer,
11        mlm=True,
12        mlm_probability=0.15
13    )
14
15 training_args = TrainingArguments(
16        output_dir="./roberta-analogy-mlm",
17        per_device_train_batch_size=64,
18        per_device_eval_batch_size=64,
19        num_train_epochs=10,
20        weight_decay=0.01,
21        logging_dir="./logs",
22        save_total_limit=1,
23        push_to_hub=False,
24        fp16=True,
25        remove_unused_columns=False
26    )
27
28 trainer = Trainer(
29        model=model,
30        args=training_args,
31        train_dataset=train_ds,
32        eval_dataset=val_ds,
33        tokenizer=tokenizer,
34        data_collator=data_collator
35    )
36
37 trainer.train()

```

Listing 39: Training RoBERTa with MLM Objective

```

1 output_directory = "/content/drive/
2     MyDrive/roberta-analogy"
3 trainer.save_model(output_directory)
4 tokenizer.save_pretrained(
5     output_directory)

```

Listing 40: Saving the Fine-Tuned RoBERTa Model

```

1 from transformers import
2     RobertaTokenizer, RobertaForMaskedLM
3 import torch
4
5 model = RobertaForMaskedLM.
6     from_pretrained("/content/drive/
7         MyDrive/roberta-analogy")
8 tokenizer = RobertaTokenizer.
9     from_pretrained("/content/drive/
10        MyDrive/roberta-analogy")
11 model.to("cuda")
12
13 from torch.nn.functional import softmax
14
15 for i in range(50, 250):
16     sample = val_dataset[i]
17     input_text = sample["input_text"]
18     expected = sample["target"]
19     inputs = tokenizer(input_text,
20         return_tensors="pt").to("cuda")
21     with torch.no_grad():
22         outputs = model(**inputs)
23         logits = outputs.logits
24         mask_index = (inputs["input_ids"] ==
25             tokenizer.mask_token_id).nonzero(
26             as_tuple=True)[1]
27         mask_logits = logits[0, mask_index,
28             :]
29         top_token_id = torch.argmax(
30             mask_logits, dim=1)
31         predicted_token = tokenizer.decode(
32             top_token_id)
33
34         print(f"\nInput: {input_text}\nExpected: {expected}\nPredicted: {predicted_token}")

```

Listing 41: Inference on RoBERTa for Analogy Completion

A.7 GloVe-Based Analogy Prediction

```

1 # glove_experiment.ipynb
2 # Automatically generated by Colab.
3 # Original file is located at:
4 # https://colab.research.google.com/
5     drive/1Ow0ULFz4proiL4qxgmFvfCo-
6     vniUBvg-

```

Listing 42: Notebook Metadata

```

1 import kagglehub
2

```

```

3 # Download latest version
4 path = kagglehub.dataset_download("thanakomsn/glove6b300dtxt")

```

Listing 43: Downloading GloVe Embeddings from KaggleHub

```

1 import numpy as np
2 import re
3
4 def load_glove_embeddings(filepath):
5     embeddings = {}
6     with open(filepath, "r", encoding="utf-8") as f:
7         for line in f:
8             values = line.strip().split()
9             word = values[0]
10            vector = np.array(values[1:], dtype=np.float32)
11            embeddings[word] = vector
12    return embeddings
13
14 def cosine_similarity(vec1, vec2):
15    dot = np.dot(vec1, vec2)
16    norm1 = np.linalg.norm(vec1)
17    norm2 = np.linalg.norm(vec2)
18    if norm1 == 0 or norm2 == 0:
19        return -1
20    return dot / (norm1 * norm2)

```

Listing 44: GloVe Loading and Utility Functions

```

1 def extract_relation(input_text,
2 relation_list):
3     input_text_lower = input_text.lower()
4     for relation in relation_list:
5         pattern = r'\b' + re.escape(
6 relation.lower()) + r'\b'
7         if re.search(pattern,
8 input_text_lower):
9             return relation
10    return None

```

Listing 45: Extracting Relation Keywords from Input

```

1 def get_context_vector(sentence, glove,
2 relation=None):
3     sentence = sentence.replace("None",
4     "")
5     all_words = re.findall(r"\b\w+\b",
6 sentence.lower())
7     vectors = [glove[word] for word in
8 all_words if word in glove]
9     if vectors:
10        return np.sum(vectors, axis=0)
11    else:
12        vector_size = next(iter(glove.
13 values())).shape[0] if glove else
14        300
15        return np.zeros(vector_size)

```

Listing 46: Context Vector Construction

```

1 def find_best_candidate(context_vector,
2 glove, relation=None, input_text=
3 None):
4     max_sim = -1

```

```

3 best_word = None
4 exclude_words = set(re.findall(r"\b\w+\b",
5 input_text.lower())) if
6 input_text else set()
7 exclude_words.discard("none")
8
9 for word, vector in glove.items():
10    if len(word) < 3 or word.isdigit()
11    or word in exclude_words:
12        continue
13    sim = cosine_similarity(
14 context_vector, vector)
15    if relation and relation.lower()
16    == "antonym":
17        sim = -sim
18    if sim > max_sim:
19        max_sim = sim
20        best_word = word
21
22 return best_word

```

Listing 47: Find Best Candidate Word using Cosine Similarity

```

1 def replace_predictions_with_glove(
2 input_file, output_file, glove):
3     relation_list = [
4         'DefinedAs', 'LocatedNear', 'Antonym',
5         'HasLastSubevent', 'CausesDesire',
6         'same efficacy', 'Desires', 'HasA',
7         'CapableOf', 'HasPrerequisite',
8         'ReceivesAction', 'IsA', 'PartOf',
9         'HasContext', 'MannerOf', 'CreatedBy',
10        'NotDesires', 'AtLocation',
11        'same raw material', 'HasProperty',
12        'HasSubevent', 'Entails', 'MadeOf',
13        'HasFirstSubevent', 'same goal',
14        'same property', 'UsedFor', 'Synonym',
15        'MotivatedByGoal', 'NotCapableOf',
16        'same capability', 'Causes',
17        'at location', 'used for', 'manner of'
18    ]
19
20 blocks = []
21 current_block = {}
22
23 with open(input_file, "r", encoding="utf-8") as f_in:
24     for line in f_in:
25         line = line.strip()
26         if not line:
27             if current_block and "input" in current_block:
28                 blocks.append(current_block)
29                 current_block = {}
30             continue
31         if line.startswith("Input"):
32             if current_block and "input" in current_block:
33                 blocks.append(current_block)
34                 current_block = {"input":
35                     line.split(":", 1)[1].strip()}
36             elif line.startswith("Expected"):
37

```

```

28         current_block["expected"]
29     ] = line.split(":", 1)[1].strip()
30     elif line.startswith("Predicted"):
31         current_block["predicted_line"] = line
32         blocks.append(current_block)
33         current_block = {}
34
35     if current_block and "input" in current_block:
36         blocks.append(current_block)
37
38     with open(output_file, "w", encoding="utf-8") as f_out:
39         for block in blocks:
40             input_text = block["input"]
41             f_out.write(f"Input : {input_text}\n")
42             if "expected" in block:
43                 f_out.write(f"Expected : {block['expected']}\n")
44             relation = extract_relation(input_text, relation_list)
45             context_vec = get_context_vector(input_text, glove, relation)
46             predicted = find_best_candidate(context_vec, glove, relation, input_text)
47             f_out.write(f"Predicted : {predicted}\n\n")

```

Listing 48: Replacing Predictions with GloVe Outputs

```

1 import os
2
3 glove_path = os.path.join(path, "glove.6
4 B.300d.txt")
5 glove = load_glove_embeddings(glove_path)
6
7 input = "model_predictions.txt"
8 replace_predictions_with_glove(input, "glove_predictions2.txt", glove)

```

Listing 49: Final Execution: GloVe Path and Prediction

A.8 Prediction Evaluation

```

1 # prediction_accuracy.ipynb
2 # Automatically generated by Colab.
3 # Original file is located at:
4 # https://colab.research.google.com/
5 # drive/1
6 12GgceNsFAgt593eEbpsUd2a7x9S6xfd

```

Listing 50: Notebook Metadata

```

1 from sentence_transformers import
2     SentenceTransformer
3 from sklearn.metrics.pairwise import
4     cosine_similarity
5
6 # Load SBERT model
7 sbert_model = SentenceTransformer('all-
8     MiniLM-L6-v2')

```

```

7 def compute_semantic_similarity(text1,
8     text2):
9     embeddings = sbert_model.encode([
10         text1, text2])
11     sim = cosine_similarity([embeddings[0]], [embeddings[1]])[0][0]
12     return sim

```

Listing 51: Loading SBERT and Similarity Function

```

1 input_file = "glove_predictions2.txt"
2 output_file = "glove_evaluation_output.
3     txt"
4
5 parsed_data = []
6 current_triplet = {}
7
8 with open(input_file, "r", encoding="utf
9 -8") as f:
10     for line in f:
11         line = line.strip()
12         if not line:
13             continue
14         if line.startswith("Input"):
15             current_triplet['input'] =
16                 line.split(":", 1)[1].strip()
17         elif line.startswith("Expected"):
18             current_triplet['expected'] =
19                 line.split(":", 1)[1].strip()
20         elif line.startswith("Predicted"):
21             current_triplet['predicted'] =
22                 line.split(":", 1)[1].strip()
23             if all(k in current_triplet
24                 for k in ('input', 'expected',
25                           'predicted')):
26                 parsed_data.append((
27                     current_triplet['input'],
28                     current_triplet['expected'],
29                     current_triplet['predicted']))
30             current_triplet = {}

```

Listing 52: Parsing GloVe Prediction Output

```

1 with open(output_file, "w", encoding="
2     utf-8") as out_f:
3     for input_text, expected, predicted
4         in parsed_data:
5             similarity =
6                 compute_semantic_similarity(expected,
7                     predicted)
8             out_f.write(f"Input : {input_text}\n")
9             out_f.write(f"Expected : {expected}\n")
10            out_f.write(f"Predicted : {predicted}\n")
11            out_f.write(f"Similarity: {similarity:.4f}\n")
12            out_f.write("-" * 40 + "\n")

```

Listing 53: Writing SBERT Similarity Scores to File

```

1 SIMILARITY_THRESHOLD = 0.2
2 exact_matches_file = "accuracy.txt"
3 correct = 0
4 total = 0
5 lines_buffer = []

```

```

6
7 with open(output_file, "r", encoding="utf-8") as f, open(
8     exact_matches_file, "w", encoding="utf-8") as exact_f:
9     for line in f:
10         lines_buffer.append(line)
11         if line.startswith("Similarity"):
12             :
13                 try:
14                     sim = float(line.split(":",
15                               1)[1].strip())
16                     if sim >=
17                         SIMILARITY_THRESHOLD:
18                             correct += 1
19                             for buffered_line in
20                                 lines_buffer:
21                                     exact_f.write(
22                                         buffered_line)
23                                         exact_f.write("\n")
24                                         total += 1
25             except ValueError:
26                 pass
27             elif line.startswith("-" * 10):
28                 lines_buffer = []
29 accuracy = correct / total if total > 0
30 else 0
31 print(f" Similarity score >= 0.2: {
32     accuracy:.4f} ({correct}/{total})")

```

Listing 54: Evaluation at Threshold = 0.2

```

1 SIMILARITY_THRESHOLD = 1.0
2 exact_matches_file = "exact_matches.txt"
3 correct = 0
4 total = 0
5 lines_buffer = []
6
7 with open(output_file, "r", encoding="utf-8") as f, open(
8     exact_matches_file, "w", encoding="utf-8") as exact_f:
9     for line in f:
10         lines_buffer.append(line)
11         if line.startswith("Similarity"):
12             :
13                 try:
14                     sim = float(line.split(":",
15                               1)[1].strip())
16                     if sim >=
17                         SIMILARITY_THRESHOLD:
18                             correct += 1
19                             for buffered_line in
20                                 lines_buffer:
21                                     exact_f.write(
22                                         buffered_line)
23                                         exact_f.write("\n")
24                                         total += 1
25             except ValueError:
26                 pass
27             elif line.startswith("-" * 10):
28                 lines_buffer = []
29 accuracy = correct / total if total > 0
30 else 0
31 print(f" Similarity score = 1.0: {
32     accuracy:.4f} ({correct}/{total})")

```

Listing 55: Evaluation at Threshold = 1.0 (Exact Match)

```

1 SIMILARITY_THRESHOLD = 0.5
2 exact_matches_file = "similarity_half.
3     txt"
4 correct = 0
5 total = 0
6 lines_buffer = []
7
8 with open(output_file, "r", encoding="utf-8") as f, open(
9     exact_matches_file, "w", encoding="utf-8") as exact_f:
10    for line in f:
11        lines_buffer.append(line)
12        if line.startswith("Similarity"):
13            :
14                try:
15                    sim = float(line.split(":",
16                               1)[1].strip())
17                    if sim >=
18                         SIMILARITY_THRESHOLD:
19                             correct += 1
20                             for buffered_line in
21                                 lines_buffer:
22                                     exact_f.write(
23                                         buffered_line)
24                                         exact_f.write("\n")
25                                         total += 1
26             except ValueError:
27                 pass
28             elif line.startswith("-" * 10):
29                 lines_buffer = []
30 accuracy = correct / total if total > 0
31 else 0
32 print(f" Similarity score > 0.5 : {
33     accuracy:.4f} ({correct}/{total})")

```

Listing 56: Evaluation at Threshold = 0.5