

Background and Motivation

Analogies in Human Cognition

- Analogy-making in human cognition and AI.
- Analogies play an important role across many areas.







- **3. Automatic Filtering and Labeling**
- Estimation: analogies are < 30% of the candidates data.
- We use part of our annotated data as few-shot examples for our **filtering** model.
 - Inputs: two paragraphs, their subjects, similar relations. Label: how many workers labeled it as an analogy (0-3).
- **Goals**:

 \bigcirc

Target

 $t_1, t_2, ..., t_n$

- To identify the most probable analogous candidates to show our annotators.
- Potentially replace the human-in-the-loop and achieve a **fully automated pipeline**.



A cork is *stuck* inside an empty wine bottle.

A Baby is stuck inside the birth canal

Existing Analogy Resources

- **Surprisingly, few** analogy resources exist today.
- We believe this **lack of data** hinders progress in computational analogy.
- Most resources focus on **word-analogies** (man:king is like woman:queen).
- Sentence-level analogies. Jiayang et al. (2023)- dataset of 24K story pairs.
- Full paragraph-level analogies. Stories from cognitive-psychology.

The Structure Mapping Theory (SMT), (Gentner, 1983)

- Analogy is a mapping from entities in **base B** to entities in $b_1, b_2, ..., b_n$ target T, relying on relational similarity, not object attributes.
- **Example**: analogy between electrical circuit & water pump. Mappings for example: electrons \rightarrow water, wire \rightarrow pipe (electrons move through wires like water flows in pipes). ground potential

Approach **ParallelPARC (Parallel Paragraph Creator) Pipeline** Random split Human annotation of a random sample -----



4. Human Validation

- We show annotators both the **most likely analogous candidates** (as predicted) by the model), but also the **least likely candidates**.
- 3 annotators per sample. Strict setting: positive if all 3 agree it is an analogy.
- We randomly gave annotators small batches to label until **310** positives.
- Annotators' agreement is 78.6%, where random chance is 25%.

<u>Filtering models' predictions vs. workers' majority vote</u>

- Accuracy of **85.1%**, f1-score of **83.4%**.
- 79.5% precision, predicting high likelihood of an analogy (> 30%)

5. Distractors Generation (Challenging Negatives)



- **Motivation:** In addition to the the analogies, our aim is to create negatives. • **Formulation**. Let B and T be two analogous paragraphs. We create distractor T' that keeps first-order relations of T, but changes the higher-order relations – i.e., relations between first-order relations (e.g., cause and effect, or temporal dependencies). **How?** To create T', we find two dependent events in T such that one must precede the other, and switch their order.
- **Generation.** GPT-4 with two separate prompts:
- Finding & Replacing two dependent events (one-shot).
- \circ Writing a coherent T' (few-shot).

Evaluation.		Base:	Target (Analogy):	Target (Distractor):				
\sim	GPT4 - 89% accuray.	How do bats use echolocation?	How do submarines use sonar?	How do submarines use sonar?				
0		(Natural Sciences)	(Engineering)	(Engineering)				
0	We create distractors for	Bats use echolocation to navigate and	Submarines use sonar technology to	Submarines interpret the echo to				
		find food. They emit high frequency	detect objects in the water.	determine the distance and size of				
	both gold and silver sets.	sound waves that bounce off of objects	They emit sound waves, which	the object. After interpreting the				
		in their environment.	travel through the water and bounce	echo, they emit sound waves, which				
		The bats then receive the echoes and in-	off the objects.	travel through the water and bounce				
		terpret the information to locate their	The sound waves are then received	off the objects. These sound waves				
		prey and navigate their surroundings.	back as an echo. Submarines in-	are then received back as an echo.				
		Submarines interpret the echo to deter-	terpret the echo to determine the	Finally, submarines use sonar tech-				
		mine the distance and size of the object.	distance and size of the object.	nology to detect objects in the water.				

Our ProPara-Logy Generated Dataset

Base	Target	Similar Relations
Title: How does a solar panel work?	Title: How does photosynthesis occur?	(solar energy, <i>powers</i> , electric current)
Domain: Engineering	Domain: Natural Science	(sunlight, <i>powers</i> , chemical reactions)
Paragraph: solar energy powers an	Paragraph: Photosynthesis occurs	
electric current within a solar panel.	when sunlight powers chemical	(photovoltaic cells, convert, energy)
The photovoltaic cells within the	reactions within the chloroplasts of a	(chloroplasts, tranform, energy)
panel convert the energy from the	plant. The chloroplasts are able to	
sun into electricity. The electrical	transform the energy from the sunlight	(electrical wires, spread, power)
wires then spread this power	into usable energy for the plant. This	(plants, distribute, nutrients)
throughout the panel. The electric	energy is then used to produce	
current is then used to power	nutrients for the plant, which are then	
whatever the panel is connected to.	distributed throughout the plant.	

<u>1. Analogy Candidates Generation</u>



- **Goal:** to generate analogy candidates from **diverse** scientific domains.
- How? We employed GPT-3- high-quality results at a very reasonable cost.
- (1): GPT tends to repeat itself. (2): GPT creates analogies of similar topics.
- (1): Seed GPT with B instead of asking it to generate both B and T.
- (2): Broad target domains: Eng., Natural, Social, and Biomedical Science.
- Using a single prompt for the task X
- Using two separate prompts V
- Finding an analogous subject, and similar relations.

Evaluating Humans and LLMs on ProPara-Logy Benchmark

<u>Binary Classification Task.</u> To decide whether the processes are analogous.

The target paragraph could be:

• Analogy (positive) / Random (simple negative) / Distractor (challenging negative)

Multiple Choice Task. Given a base paragraph B, along with 4 candidate

paragraphs, the task is to identify the paragraph that is most analogous to B.

- **Basic:** includes one analogous paragraph and **3 random paragraphs**.
- Advanced: includes challenging distractors. Ο

Research Questions:

- **RQ1:** What is the performance of humans and models?
 - Humans achieve better performance than models (~13% gap on both tasks)!
 - **GPT4** achieves the best accuracy out of the models!
- **RQ2:** Is the automatically-generated "silver set" (without human validation) useful for training models?
 - The training of FlanT5-small on the silver-set significantly improved its Performance! Ο
- Generating a paragraph in natural language (given subject, and relations). We include Similar relations, in addition to paragraphs, subjects & domains. • RQ3: Can the distractors fool humans and models? In total: 4,288 candidates.

2. Human Annotation Task



- We now annotate a small portion of the **candidates data**.
- **Goal**: to estimate % of analogies & use the annotated data to train models.
- Given two paragraphs (B, T), corresponding subjects, domains & similar relations. The task: to decide whether the paragraphs are analogous and the similar relations are correct.
 - **YES** (close / far) **analogy**. Ο
 - **NO** "for further inspection" (dissimilar relations, misinformation, cyclic vs. non-cyclic process, other)

- - The challenging distractors confuse LLMs, but not humans!

Binary Classification Task

Multiple Choice Task

Row	Settings	Method	Overall	Per Target Type		
				Positives (50%)	Negatives (50%)	
				Analogy	Random	Distractor
1		Random Guess	50	50	50	50
2		GPT4	79.5	95.2	92.9	34.8
3		ChatGPT	68.2	53.5	96.8	69.0
4	Zero-shot	Gemini Pro	73.9	79.7	100	36.1
5		FlanT5-XXL	61.1	28.1	100	88.4
6		FlanT5-XL	59.7	25.1	100	88.4
7		FlanT5-small	49.3	0	97.4	100
8		Humans	79	58	100	100
9		GPT4 (in-context)	78	86.5	98.1	40.7
10	Guided	FlanT5-small (fine-tune)	74.4	87.1	96.1	27.1
11		Humans	92.5	95	100	80

Basic Advanced Settings Method 25 Random Guess 25 83.2 95.5 GPT4 74.2 ChatGPT Gemini Pro FlanT5-XXL 75.2 FlanT5-XL 55.5 FlanT5-small 32.9 32.9 Guided Humans 100

We hope researchers will use the pipeline in domains where analogies have shown promise, and that this work will inspire more NLP work on analogies, leading to new tasks and benchmarks!