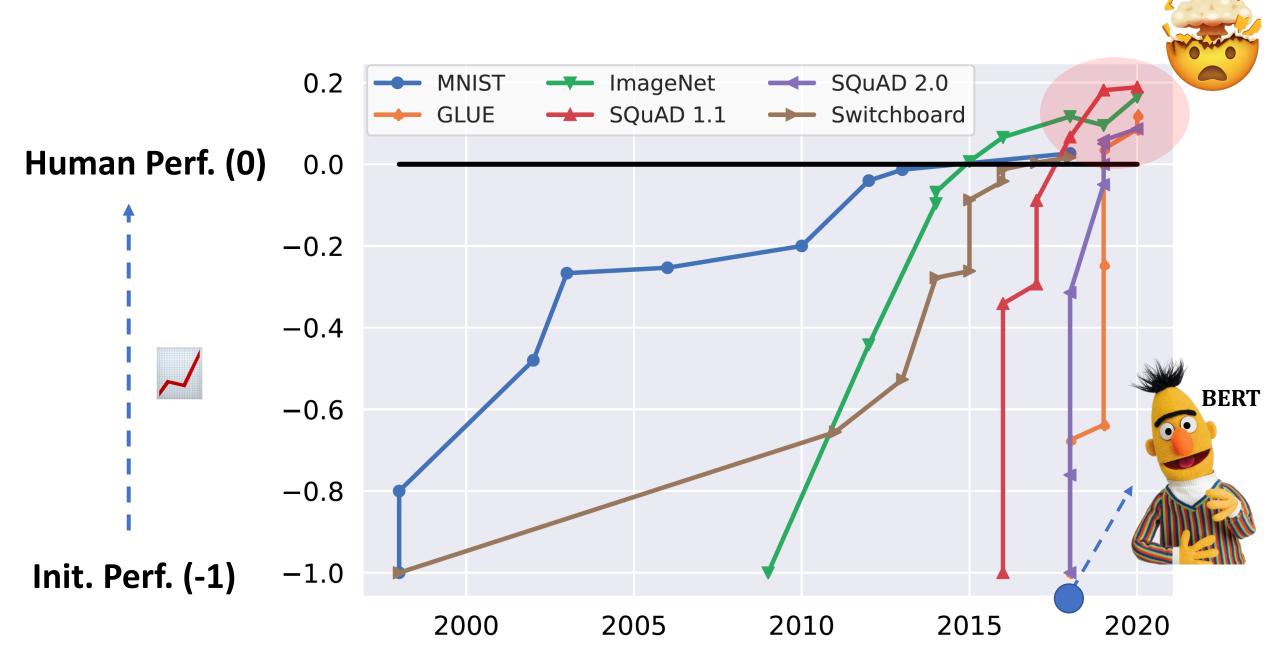# Reflex or Reflect

## When Do Language Tasks Need Slow Reasoning?

Xiang Ren

Associate Professor, CS & ISI

Viterbi Early Career Chair

University of Southern California

http://inklab.usc.edu

REFLEXIVE

REFLECTIVE

Human Perf. (0)

Init. Perf. (-1)

Legend: MNIST, GLUE, ImageNet, SQuAD 1.1, SQuAD 2.0, Switchboard

BERT

Initial baseline performance and human performance are normalized to **-1** and **0** respectively (Credit: Kiela et al., 2021). **2**
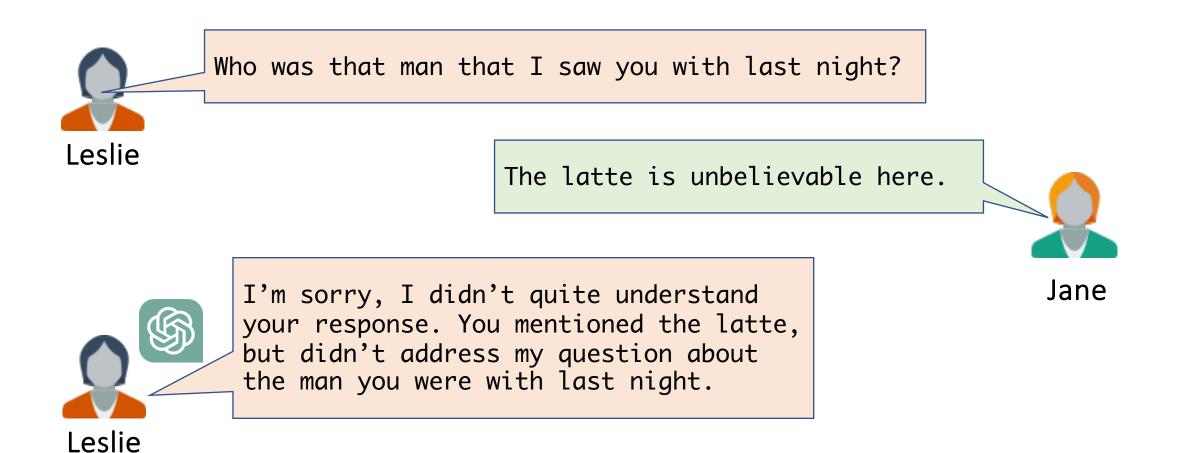
| Benchmark | GPT-4 Evaluated few-shot | GPT-3.5 Evaluated few-shot | LM SOTA Best external LM evaluated few-shot | SOTA Best external model (includes benchmark-specific training) |
|---|---|---|---|---|
| MMLU Multiple-choice questions in 57 subjects (professional & academic) | 86.4% 5-shot | 70.0% 5-shot | 70.7% 5-shot U-PaLM | 75.2% 5-shot Flan-PaLM |
| HellaSwag Commonsense reasoning around everyday events | 95.3% 10-shot | 85.5% 10-shot | 84.2% LLAMA (validation set) | 85.6% ALUM |
| AI2 Reasoning Challenge (ARC) Grade-school multiple choice science questions. Challenge-set. | 96.3% 25-shot | 85.2% 25-shot | 84.2% 8-shot PaLM | 85.6% ST-MOE |
| WinoGrande Commonsense reasoning around pronoun resolution | 87.5% 5-shot | 81.6% 5-shot | 84.2% 5-shot PALM | 85.6% 5-shot PALM |
| HumanEval Python coding tasks | 67.0% 0-shot | 48.1% 0-shot | 26.2% 0-shot PaLM | 65.8% CodeT + GPT-3.5 |
| DROP (f1 score) Reading comprehension & arithmetic. | 80.9 3-shot | 64.1 3-shot | 70.8 1-shot PaLM | 88.4 QDGAT |

https://openai.com/research/gpt-4

3

Your Magical AI-generated World

# On My Wishlist: Reading the Air

Leslie and Jane are chatting at a coffee shop.

Leslie: Who was that man that I saw you with last night?

Jane: The latte is unbelievable here.

Leslie: I'm sorry, I didn't quite understand your response. You mentioned the latte, but didn't address my question about the man you were with last night.

# On My Wishlist: Indirect Speech

Adam and Bill are working on a project in Bill's room. Bill opens the window to get some fresh air. A cold breeze blows in.

Adam: Is the window open?

Bill: Yes, I just opened it.

# On My Wishlist: Indirect Speech

Adam and Bill are working on a project in Bill's room. Bill opens the window to get some fresh air. A cold breeze blows in.

Is the window open?

Adam

Bill

- Adam feels the breeze and would like to be warmer
- Adam probably wants to close the window
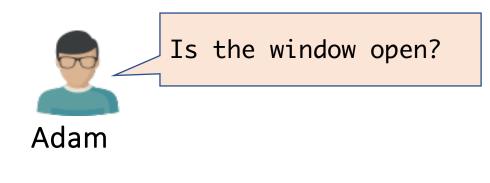- (but Adam didn't want to seem rude)

# On My Wishlist: Indirect Speech

Adam and Bill are working on a project in Bill's room. Bill opens the window to get some fresh air. A cold breeze blows in.

Adam: Is the window open?

**Adam**

Bill: Is it too cold? Do you want me to close it?

**Bill**
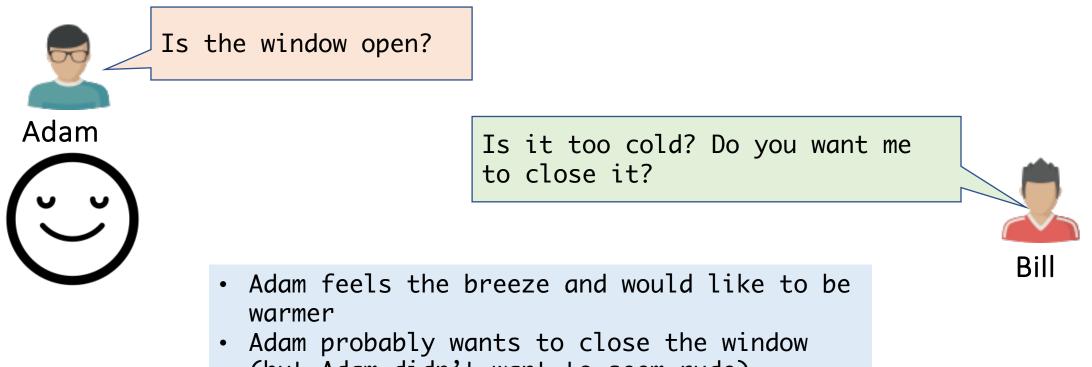
- Adam feels the breeze and would like to be warmer
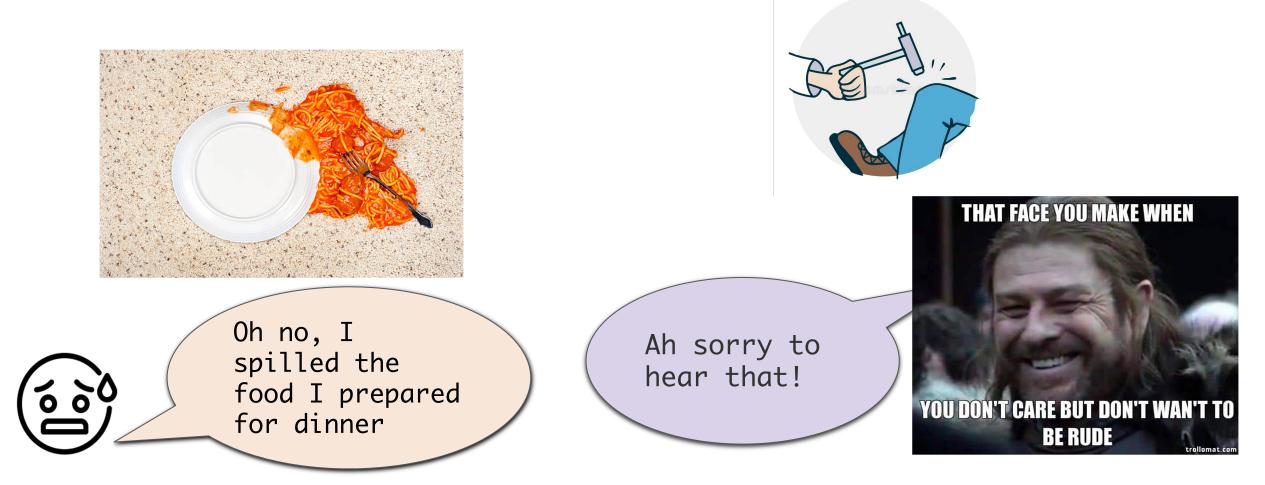- Adam probably wants to close the window
- (but Adam didn't want to seem rude)

# Muscle-*Reflex* Style Language

"*Reflect*" Style Language

# STUDIES IN THE WAY OF WORDS

## PAUL GRICE

Paul Grice's Maxims on *cooperative principles*

Communication is a **collaborative** effort with **intents** and people tend to "*minimize the total effort spent*". [Least collaborative effort]

Due to least collaborative effort, we need to **make inferences** to draw conclusions about the speaker's **intentions, emotion states, and experiences.** [Build Common Ground]

PAUL GRICE
Philosopher and Linguist
*Siobhan Chapman*

# "*Reflect*" Style Language

# "*Reflect*" Style Language

## Why Challenging?

- Often implicit in training corpora → more prone to generate *shallow* replies

- Appropriate answers require *slow reasoning* about others' true intents and common sense

pasta place?

# How do we reply in conversations?

*I'm going to sing in front of hundreds tomorrow...*

# How do we reply in conversations?

*I'm going to perform in a piano recital tomorrow...*

*Performing in front of audience can cause anxiety*

Deep breaths, you'll do great!

*Recalling & Combining common sense* with information expressed in NL to *make inferences*

Producing *consistent* inferences amidst *logically-equivalent yet linguistically-varied* paraphrases

Clark, H. H., & Brennan, S. E. (1991). *Grounding in communication.*

# RICA: Robust Inference on Commonsense Axioms

➢ Test model's robustness against linguistic variations
➢ Focus on implicit commonsense inferences
➢ Scalable probe set construction process

## RICA: Evaluating Robust Inference Capabilities Based on Commonsense Axioms

**Pei Zhou**    **Rahul Khanna**    **Seyeon Lee**    **Bill Yuchen Lin**    **Daniel Ho**

**Jay Pujara**    **Xiang Ren**

Department of Computer Science and Information Sciences Institute
University of Southern California

{peiz,rahulkha,seyeonle,yuchen.lin,hsiaotuh,jpujara,xiangren}@usc.edu

Commonsense Logic to Probe:

*A.Size < B.Size → P(A in Container) > P(B in Container)*

pridag    fluberg

pridag    fluberg

**A pridag is smaller than a fluberg,** so it is **[MASK]** to put a pridag into a box than a fluberg.

**A fluberg is smaller than a pridag,** so it is **[MASK]** to put a pridag into a box than a fluberg.

easier (86.6%)
harder (1.1%) ✅

easier (87.2%)
harder (1.3%) ❌

17

# RICA: Robust Inference on Commonsense Axioms

- Examples:
  - **Original**: "A is heavier than B, so A is <better> at sinking than B."
  - **Negation**: "A is heavier than B, so A is **not** <worse> at sinking than B."
  - **Entity Swap**: "**B** is heavier than **A**, so A is <worse> at sinking than B."
  - **Antonym**: "A is heavier than B, so A is <worse> at **floating** than B."
  - …

# RICA: Robust Inference on Commonsense Axioms

- Masked word prediction task: **Choose <better> or <worse>**:

  - **Original**: "A is heavier than B, so A is <MASK> at sinking than B."

  - **Perturb1**: "A is heavier than B, so A is **not** <MASK> at sinking than B."

  - **Perturb2**: "**B** is heavier than **A**, so A is <MASK> at sinking than B."

  - **Perturb3**: "A is heavier than B, so A is <MASK> at **floating** than B."

  - …

(Zhou et al., EMNLP'21)

# **Results:** Human-Curated Set

- **Random-guessing** like performance on *all settings* for all models.

- Training on similar data does **not** help achieve real robustness

Average Accuracy

| | |
|---|---|
| **Human** | **91.7%** |

| | |
|---|---|
| Zero-Shot | BERT etc. |
| Low-Resc. | BERT etc. |
| High-Resc. | BERT etc. |
| Noisy 100k | BERT etc. |

**~50%**

# **Results:** How About Fancy New LLMs?

RICA still remains challenging to LLMs

- Larger models tend to perform better for T5-family models

- GPT-family models seem less magical
  - Bidirectional attention better captures logic with perturbations?

### Average Accuracy on **Zero-Shot Prompting**

| Model | Accuracy |
|-------|----------|
| Human | 91.7% |
| BERT et al. | 50% |
| Flan-T5-Base | 50% |
| Flan-T5-11B | 64% |
| Flan-UL2-20B | 68% |
| GPT3.5-Turbo (ChatGPT) | 54% |
| GPT4 | 64% |

# Analysis: Positivity Bias

- Heavy bias towards positive-valence words such as "*more*", "*better*".

- Fine-tuning on RICA mitigates the imbalance issue (but still fails)

Average Accuracy without Fine-Tuning

| | |
|---|---|
| Human (both positive and negative) | 91.7% |
| Pos. Words — BERT etc. | 87.2% |
| Neg. Words — BERT etc. | 12.5% |

Average Accuracy after Fine-tuning

| | |
|---|---|
| Pos. Words — BERT etc. | |
| Neg. Words — BERT etc. | ~50% |

# Scaling is the Way Going Forward!



GPT3

Scaling Laws for Neural Language Models

$L = (D/5.4 \cdot 1$

$L = (C_{min}/2.3 \cdot 10^8)^{-0.050}$

Zero-shot   One-shot   Few-shot

Natural Language Prompt

No Prompt

175B Params

13B Params

1.3B Params

Examples in Context (K)

# Bing, Bard, and ChatGPT: AI chatbots are rewriting the internet

How we use the internet is changing fast, thanks to the advancement of AI-powered chatbots that can find information and redeliver it as a simple conversation.

# Does Scaling Always Work?

| Many tasks like this | Any |
|---|---|



**Zhengping Zhou and Yuhui Zhang**, for *NeQA: Can Large Language Models Understand Negation in Multi-choice Questions?*

This task takes an existing multiple-choice dataset and negates a part of each question to see if language models are sensitive to negation. The authors find that smaller language models display approximately random performance whereas the performance of larger models become significantly worse than random.

## Modus Tollens, by Sicong Huang and Daniel Wurgaft (Third Prize)

**TL;DR** This task shows strong inverse scaling on almost all models and represents a simple logical reasoning task (*modus tollens*) that might be expected to show regular scaling. Inverse scaling trends hold across both pretrained LMs and LMs finetuned with human feedback via RL from Human Feedback (RLHF) and Feedback Made Easy (FeedME).

**Robustness on logical reasoning?**

# Language-based Deductive Reasoning

**fact1:** Charlie is blue.
**fact2:** Charlie is round.
**fact3:** Erin is kind.
**fact4:** Dave is round.

**rule1:** If someone is blue then they are kind.
**rule2:** Round, kind people are white.

**statement:** Charlie is white.

**Theory**

**Input**: Facts + Rules (theory), Statement

**Output**: Entailment label

- ➢ **True**: Theory → statement is True

- ➢ **False**: Theory → negation of the statement is True

- ➢ **Unknown**: No conclusion

# Can ChatGPT do Deductive Reasoning?



For a given theory and statement, tell if the statement is provable using the theory. If it is provable, then output "True". If the negation of the statement is provable, then output "False". If nothing can be said about the statement, then output "Unknown".

Theory: Sam is tall. Tom is smart. Tall people are good. Tall people are not blue.
Statement: Sam is blue.

Output: False. ✔

Sure, it can get it right sometimes, but ...

# Can ChatGPT do Deductive Reasoning?



For a given theory and statement, tell if the statement is provable using the theory. If it is provable, then output "True". If the negation of the statement is provable, then output "False". If nothing can be said about the statement, then output "Unknown".

Theory: Sam is tall. Tom is smart. Tall people are good. Tall people are blue.
Statement: Sam is blue.

False. ✗

**… *not* robust to negation within the theory..**

# Can ChatGPT do Deductive Reasoning?



For a given theory and statement, tell if the statement is provable using the theory. If it is provable, then output "True". If the negation of the statement is provable, then output "False". If nothing can be said about the statement, then output "Unknown".

Theory: Sam is tall. Tom is smart. Tall people are good. Tall people are blue and not kind.
Statement: Sam is kind.

Unknown. ❌

The theory provides information about Sam being tall, but it doesn't provide any direct information about whether Sam is kind or not. Therefore, we cannot determine the truth value of the statement "Sam is kind" based on the given theory alone.

**Erroneous reasoning given the theory...**

# Robust Reasoning: Lexical Perturbation

# Robust Reasoning: Lexical Perturbation

**True**, because Sam is tall and tall people are good.

**Unknown**!

**False**, because Sam is kind and all kind people are good.

## FaiRR: Faithful and Robust Deductive Reasoning over Natural Language, ACL 2022

**Input**: Sam is tall. Tom is smart. All tall people are good.

**Conclusion**: Sam is good.

**Input**: Alex is tall. Tom is smart. All tall people are good.

**Conclusion**: Alex is good.

**Input**: Sam is kind. Tom is smart. All kind people are good.

**Conclusion**: Sam is good.

# RobustLR: Logical Perturbation

Sam is tall. Tom is smart. Tall people are good. Tall people are blue.

Sam is good. **True**

# RobustLR: Logical Perturbation

Sam is tall. Tom is smart. <mark>Tall people are good.</mark> Tall people are blue.

Sam is good. **True**

➢ Logical Equivalence **Contraposition**
(A → B ≣ ~B → ~A)

Sam is tall. Tom is smart. <mark>A person who's not good is also not tall.</mark> Tall people are blue.

Sam is good. **True**

# RobustLR: Logical Perturbation

➤ Logical Equivalence **Contraposition**

   (A → B ≡ ~B → ~A)

➤ Logical Equivalence **Distributive**

   (A → B; A → C ≡ A → B AND C)

# RobustLR: Logical Perturbation

Sam is tall. Tom is smart. Tall people are good. Tall people are blue.

Sam is good. **True**

➢ Logical Equivalence **Contraposition**

(A → B ≡ ~B → ~A)

➢ Logical Equivalence **Distributive**

(A → B; A → C ≡ A → B AND C)

➢ Logical **Contrast**

(A → B  vs  A → B & C, etc.)

Sam is tall. Tom is smart. Tall people are good. Tall people are blue.

Sam is good. **True**
Sam is kind. **Unknown**

Sam is tall. Tom is smart. Tall people are good and not kind. Tall people are blue.

Sam is good. **True**
Sam is kind. **False**

# RobustLR: Dataset generation process



Facts
- Tall(Sam)
- Kind(Ana)
- Red(Bob)

Rules

Conclusions

1. Sample some predicates

2. Label the predicates as **valid** and **invalid**

3. Break down into multiple levels

4. Starting from level 1, select predicates from lower level, such that a valid rule is formed

# RobustLR: Dataset generation process



**Facts**

Tall(Sam)   Kind(Ana)   Red(Bob)

0   1   2

Blue(Sam)

3   4

~Big(Bob)

**Rules**

5   6

Kind(Bob)

~Red(Ana)

Good(Sam)

~Tall(Bob)

7   8

**Conclusions**

1. Sample some predicates

2. Label the predicates as **valid** and **invalid**

3. Break down into multiple levels

4. Starting from level 1, select predicates from lower level, such that a valid rule is formed

Can control the degree of the rule, #negations, multiple proof graphs, etc., in a flexible manner

**10k+ test Instances**

**50k+ training instances**

f1: Charlie is tall.
r1: Erin is kind, if Charlie is tall.
statement: Erin is kind.
Label: *True*

**Original Theory**

f1: Charlie is tall.
r1: Erin is kind, if Charlie is tall *or round*.
statement: Erin is kind.
Label: *True*

**Disjunction Contrast**

f1: Charlie is tall.
r1: Erin is kind, if Charlie is tall *and round*.
statement: Erin is kind.
Label: *Unknown*

**Conjunction Contrast**

f1: Charlie is tall.
r1: *If Erin is not kind, then Charlie is not tall.*
statement: Erin is kind.
Label: *True*

**Contrapositive Equivalence**

# Results - Machine vs Human

Macro F1



*Training a RoBERTa architecture from scratch

# Results - Machine vs Human

Macro F1



*Finetune a pretrained checkpoint

# Results - Machine vs Human

Macro F1



| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1.00 — | | | | | | | |
| 0.75 — | | | | | | | |
| 0.50 — | | | | | | | |
| 0.25 — | | | | | | | |
| 0.00 — | | | | | | | |

Scratch · RoBERTa · T5-Large · T5-3B · T5-11B · GPT-3 · GPT-3.5 · Human

*6-shot in-context learning

# Results - Machine vs Human

Macro F1



1.00

0.75

0.50

0.25

0.00

Scratch   RoBERTa   T5-Large   T5-3B   T5-11B   GPT-3   GPT-3.5   Human

*7 CS graduates annotating a subset of the data

# Results - Machine vs Human

Macro F1



**Logical Contrast** ■  **Logical Equivalence** ■

Scratch, RoBERTa, T5-Large, T5-3B, T5-11B, GPT-3, GPT-3.5, Human

Training from scratch fails!

Pretrained knowledge is crucial

# Results - Machine vs Human

Macro F1



Model size is not a very significant factor, but T5 > RoBERTa!

# Results - Machine vs Human

Macro F1



**■ Logical Contrast   ■ Logical Equivalence**

GPT3/3.5 performance is worse than finetuned models!

# Results - Machine vs Human

Macro F1



**Logical Contrast**    **Logical Equivalence**

The performance gap is **low** for humans

→ more robust reasoning!

# Results - Variation with Logical Operators



Macro F1

RoBERTa-Large | T5-Large | T5-3B | T5-11B

Difficulty level

Negation > Conjunction > Disjunction

# Related Works

**P1:** David, Jack and Mark are colleagues in a company. David supervises Jack, and Jack supervises Mark. David gets more salary than Jack.

**Q:** *What can be inferred from the above statements?*
   A. Jack gets more salary than Mark.
   B. David gets the same salary as Mark.
   C. One employee supervises another who gets more salary than himself.
✓ **D. One employee supervises another who gets less salary than himself.**

**P2:** Our factory has multiple dormitory areas and workshops. None of the employees who live in dormitory area A are textile workers. We conclude that some employees working in workshop B do not live in dormitory area A.

**Q:** *What may be the missing premise of the above argument?*
   A. Some textile workers do not work in workshop B.
   B. Some employees working in workshop B are not textile workers.
✓ **C. Some textile workers work in workshop B.**
   D. Some employees living in dormitory area A work in the workshop B.

**LogiQA**

---

*(Input Facts:)* Alan is blue. Alan is rough. Alan is young.
Bob is big. Bob is round.
Charlie is big. Charlie is blue. Charlie is green.
Dave is green. Dave is rough.

*(Input Rules:)* Big people are rough.
If someone is young and round then they are kind.
If someone is round and big then they are blue.
All rough people are green.

Q1: Bob is green. True/false? **[Answer: T]**
Q2: Bob is kind. True/false? **[F]**
Q3: Dave is blue. True/false? **[F]**

**RuleTaker**

---

**Question:** How might eruptions affect plants?
**Answer:** They can cause plants to die

**Hypothesis**
H (hypot): Eruptions can cause plants to die

**Text**
sent1: eruptions emit lava.
sent2: eruptions produce ash clouds.
sent3: plants have green leaves.
sent4: producers will die without sunlight
sent5: ash blocks sunlight.

*or* Corpus

**Entailment Tree**
H (hypot): Eruptions can cause plants to die

int1: Eruptions block sunlight.
sent4: producers will die without sunlight.

sent2: eruptions produce ash clouds.
sent5: ash blocks sunlight.

**Entailment Bank**

---

**RICA**

1. **Base Predicates**
- Property(A,p)
- Relation(A,B,r)
- Comparator(x,y)

2. **Logical Template**
Rel(A,B,r) →
Comp(Prop(A,p), Prop(B,p))

3. **Knowledge Table**

| Relation | Property |
|----------|----------|
| Lawyer | Knowledge of Law |
| Doctor | Takes care of people |
| ... | ... |

4. **Created Axiom**
Rel(A,B, *lawyer*) →
Comp(Prop(A, *knowledge of law*), Prop(B, *knowledge of law*))

*Perturbation Functions*

5. **Commonsense Statement Set**
A is B's *lawyer*, so A is *more knowledgeable about law* than B
B is A's *lawyer*, so A is *not more knowledgeable about law* than B
A is B's *lawyer*, so A is *less clueless about law* than B
A is B's *lawyer*, so B is *less informed on the law* than A

*Text Conversion Module*

...

*Replace A and B with Novel Entities:* A → *prindag* B → *fluberg*

---

**CLUTRR**

**Kristin** and her son **Justin** went to visit her mother **Carol** on a nice Sunday afternoon. They went out for a movie together and had a good time.

Q: How is **Carol** related to **Justin** ?

A: Carol is the **grandmother** of Justin

# "*Reflect*" Style Language Reasoning

# We Need Slower and Deeper Language Reasoning

- Paul Grice's Maxims on *cooperative principles*

- Herbert H Clark: *Common ground*

- Jens Allwood: *Linguistic Communication as Action and Cooperation*
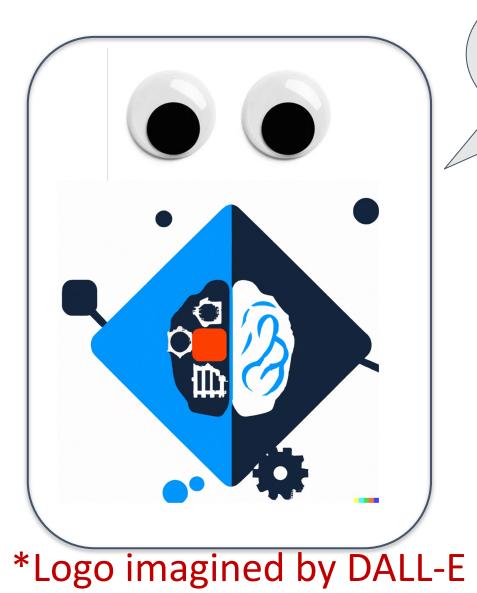
# We Need Slower and Deeper Language Reasoning

★ Communication is a **collaborative** effort with **intents** and people tend to "*minimize the total effort spent*". [Least collaborative effort]

★ Effective communications require "*reaching **mutual beliefs and knowledge** among participants called grounding*". Common sense serves a critical role in building such knowledge [Common Ground]

★ Due to least collaborative effort, we need to **make inferences to draw conclusions about the speaker's intentions, emotion states, and experiences.** [Build Common Ground]