# From Models to Microtheories: Distilling a Model's Topical Knowledge for Grounded Question Answering

Nathaniel Weir, Bhavana Dalvi Mishra, Orion Weller, Oyvind Tafjord,
Sam Hornstein, Alexander Sabol, Peter Jansen, Benjamin Van Durme, Peter Clark

Ai2 · ICLR International Conference On Learning Representations · JOHNS HOPKINS UNIVERSITY

## Motivation

Chain-of-thought materializes knowledge for **one** question. What about a whole **topic** of questions?

Can we create a method for extracting and reasoning over a discrete representation of an LLM's beliefs?

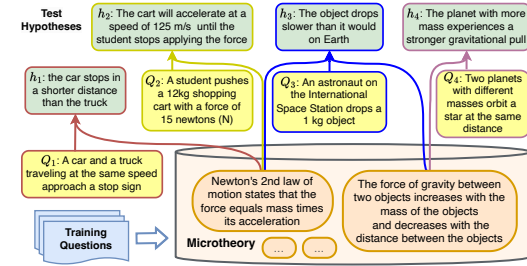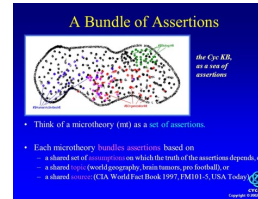E.g. what does GPT-4 know about mechanical physics?

## Inspiration

Humans develop and teach each other general *microtheories* about a topic space e.g. you learned the "magnetism" microtheory from your school teacher

A microtheory serves as **representational buffer** that allows us to operationalize our understanding of an area. Can we apply the same idea to neural reasoning algorithms?
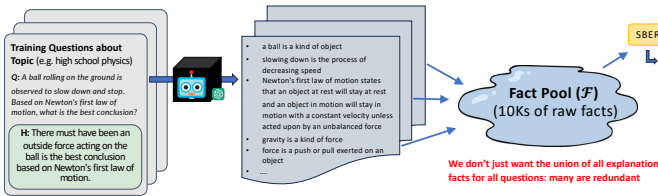
We define a **microtheory** as a concise, generalizable set of core facts about a domain

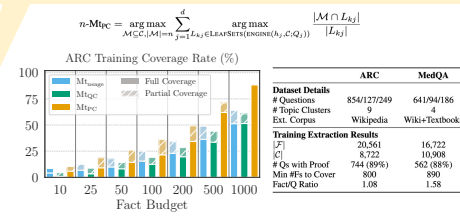The **Cyc project** used a Similar idea in the 90s/00s


A Bundle of Assertions

the Cyc KR, as a set of assertions

- Think of a microtheory (mt) as a set of assertions.
- Each microtheory bundles assertions based on —
  - a shared set of assumptions on which the truth of the assertions depends, or
  - a shared topic (world geography, brain tumors, pro football), or
  - a shared source (CIA World Fact Book 1997, FM101-5, USA Today).

Check out the paper!

**Test Hypotheses**

$h_2$: The cart will accelerate at a speed of 125 m/s until the student stops applying the force

$h_3$: The object drops slower than it would on Earth

$h_4$: The planet with more mass experiences a stronger gravitational pull

$h_1$: the car stops in a shorter distance than the truck

$Q_2$: A student pushes a 12kg shopping cart with a force of 15 newtons (N)

$Q_3$: An astronaut on the International Space Station drops a 1 kg object

$Q_4$: Two planets with different masses orbit a star at the same distance

$Q_1$: A car and a truck traveling at the same speed approach a stop sign

Newton's 2nd law of motion states that the force equals mass times its acceleration

The force of gravity between two objects increases with the mass of the objects and decreases with the distance between the objects

Training Questions ⟹ Microtheory

## Construction Pipeline

Step 1: Extract Large, Redundant Fact Pool



**Training Questions about Topic** (e.g. high school physics)
Q: A ball rolling on the ground is observed to slow down and stop. Based on Newton's first law of motion, what is the best conclusion?
H: There must have been an outside force acting on the ball is the best conclusion based on Newton's first law of motion.

- a ball is a kind of object
- slowing down is the process of decreasing speed
- Newton's first law of motion states that an object at rest will stay at rest and an object in motion will stay in motion with a constant velocity unless acted upon by an unbalanced force
- gravity is a kind of force
- force is a push or pull exerted on an object
- ....

**Fact Pool (F)** (10Ks of raw facts)

We don't just want the union of all explanation facts for all questions: many are redundant

## Usage Optimization

Step 2: Distill out semantic redundancies and find $n$ best coverage facts

SBERT Deduplication → NLI Condensation → Usage Optimization (Linear Program) → Microtheory

1. Using GPT-4 entailment engine, find all combinations of facts from condensed fact pool $\mathcal{C}$ that are an "argumentative basis" for each train hypothesis
$$L_i = [\text{fact}_1, \text{fact}_2, \text{fact}_3, \dots]$$
2. Linear program finds optimal $n$ facts for max partial coverage (PC) of a vector for each hypothesis
(We also tried full question coverage (QC) and sorting by question usage frequency, those work less well)

$$n\text{-Mt}_{PC} = \underset{\mathcal{M} \subseteq \mathcal{C}, |\mathcal{M}| = n}{\arg\max} \sum_{j=1}^{d} \underset{L_{kj} \in \text{LEAFSETS}(\text{ENGINE}(h_j, \mathcal{C}; Q_j))}{\arg\max} \frac{|\mathcal{M} \cap L_{kj}|}{|L_{kj}|}$$

### ARC Training Coverage Rate (%)



| Dataset Details | ARC | MedQA |
|---|---|---|
| # Questions | 854/127/248 | 641/94/186 |
| # Topic Clusters | 9 | 4 |
| Ext. Corpus | Wikipedia | Wiki+Textbooks |
| **Training Extraction Results** | | |
| $|\mathcal{F}|$ | 20,561 | 16,722 |
| $|\mathcal{C}|$ | 8,722 | 10,908 |
| # Qs with Proof | 744 (89%) | 562 (88%) |
| Min #Fs to Cover | 800 | 890 |
| Fact/Q Ratio | 1.08 | 1.58 |

## Subgraph from an entailment network for ~750 hypotheses

Unselected Fact Pool Facts · Entailment Tree Nodes · Selected Microtheory Facts · Training Hypotheses · Question Contexts



the force of gravity between two objects increases with the mass of the objects and decreases as the distance between the objects increases

the second law of motion states that the force acting on an object is equal to the mass of that object times its acceleration

Greater force is required to stop an object with greater mass given that the speed is the same

**The algorithm finds the core, generalizable facts**

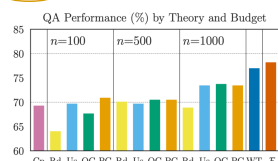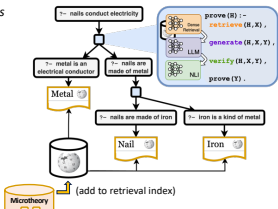## Application to Entailment Engine

*Does this hypothesis follow from this knowledge corpus? (and with what confidence?)*

- Hooks up to any text retrieval corpus
- Uses LLMs/classifiers for generating + verifying decompositions
- Accepts a hypothesis if it can fully ground it to the corpus
  - Identifies different "argumentative bases" for supporting hypothesis
- As you give it more (useful) facts, it should reason better



nails conduct electricity
prove(H) :-
retrieve(H, X),
metal is an electrical conductor
nails are made of metal
generate(H, X, Y),
verify(H, X, Y),
prove(Y).
Metal
Nail · Iron
nails are made of iron
iron is a kind of metal
Microtheory (add to retrieval index)

### Engine QA Evaluation

Adding 1000-fact Mts improves QA w/Engine by 4%

Adding full fact pool ($\mathcal{F}$) improves QA by 9%

QA Performance (%) by Theory and Budget
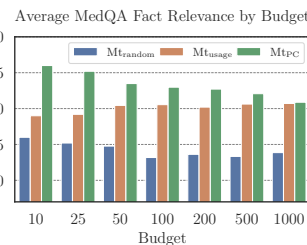

$n=100$ · $n=500$ · $n=1000$

Figure 7: Question Answering performance on ARC topical test questions using Wikipedia plus various Mts (**R**andom, **Us**age, **QC**, **PC**) as knowledge sources.

## Relevance Evaluation

**Human:**
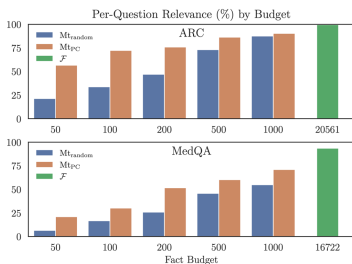We recruited senior medical students to annotate relevance of MedQA microtheory facts

Microtheory optimized for partial coverage (PC) has higher scores than random selection or sorting by usage frequency
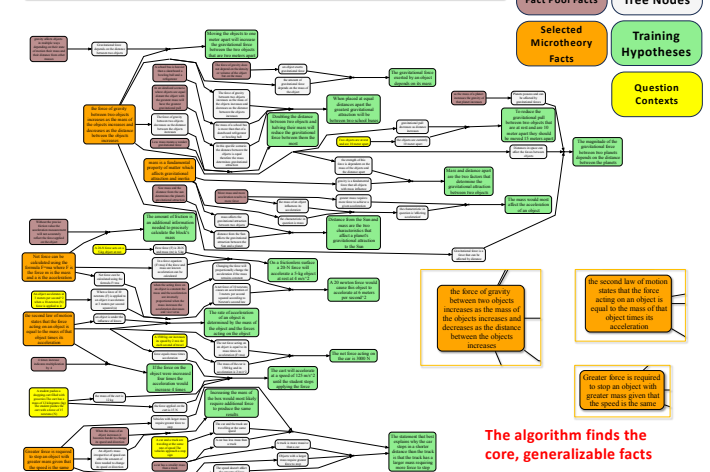

Average MedQA Fact Relevance by Budget
Mt_random · Mt_usage · Mt_PC

**Automatic:**
We asked GPT-4 whether *at least one* fact in a microtheory is useful for answering a test question correctly

PC-optimized microtheory has substantially higher relevance rate than random


Per-Question Relevance (%) by Budget
ARC · MedQA
Mt_random · Mt_PC · $\mathcal{F}$

## Takeaways

We present a method for materializing a model's latent, topical knowledge into a discrete **microtheory** articulating the model's core, reusable knowledge about the topic

We explore ways to identify the most core $n$ facts to explain training hypotheses from a large knowledge pool

Microtheories can improve a entailment engine's ability to ground hypotheses and correctly answer questions