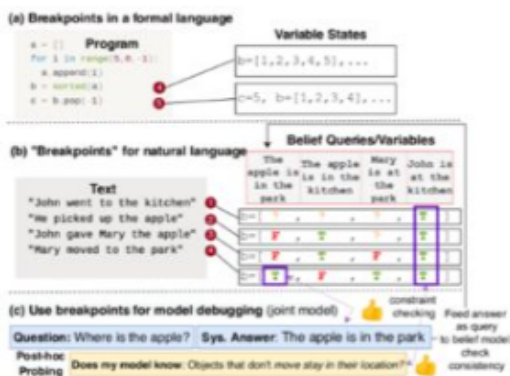# Breakpoint Transformers for Modeling and Tracking Intermediate Beliefs

Kyle Richardson[2]❄, Ronen Tamari[1]*❄, Oren Sultan[1],
Reut Tsarfaty[3], Dafna Shahaf[1] and Ashish Sabharwal[2]

[1]Hebrew University of Jerusalem, Israel, [2]Allen Institute of Artificial Intelligence, Seattle, USA, [3]Bar–Ilan University, Israel
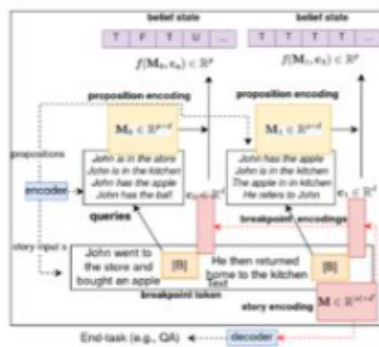* Work begun during internship at AI2  ❄ Equal contribution

## Tracking intermediate Beliefs: Motivation

- Wish your language model (LM) had **breakpoints** you could inspect to probe its intermediate semantic representations?
- Breakpoints in programming are vital for code interpretability: **allow inspection** of program state at **intermediate points throughout execution**
- We develop a new idea of "natural language breakpoints" that can be used to probe LM encodings of input texts



(a) Breakpoints in a formal language — Variable States

(b) "Breakpoints" for natural language — Belief Queries/Variables

(c) Use breakpoints for model debugging (joint model)

## Breakpoint Transformers (BPTs): Modeling Approach

- Breakpoints are simply a special token **[B]** inserted after each sentence
- Breakpoint encoding can then be queried against natural language proposition $p$ to obtain $\{T,F,?\}$ prediction
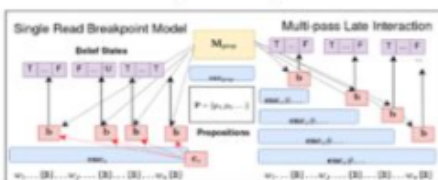- Represent summary of model "beliefs" at that point in text



### Optimization

Semantic logic loss translates to cross-entropy loss:

$$\bigwedge_{b_j \in B^{(i)}} \bigwedge_{p_k \in P_j^{(i)}} Y_{j,k}^{(i)}(b_j, p_k) = \sum_{j=1}^{m} \sum_{k=1}^{t} -\log \Pr[y_{j,k}]$$

- BPT compared against standard late-interaction baseline (Sentence Transformers*) - BPT single read enables efficient scaling



*Reimers & Gureyvich (2019)

## Datasets: Annotating Intermediate State



| Task | Example Stories | Breakpoint Propositions |
|---|---|---|

bAbI, CLUTRR: synthetic | TRIP (Storks et. al, 2021): human-authored

## Experimental Results

### CLUTRR
- BPTs show improved prediction acc., consistency and training efficiency



### bAbI
- BPTs can accurately predict hundreds of relations across long stories jointly (compared to **multi-pass** baseline)



### TRIP



Tiered 3-task eval — Which story is more plausible? A

- BPT show up to 20-30% improvement against RoBERTA-based (RoB) approach of Storks et. al (2021)
- BPT needed no additional arch. adaptation, RoB tailored arch specifically for TRIP

| Split | Model | Task 1 (Plaus.) | Task 2 (Consist.) | Task 3 (Verif.) |
|---|---|---|---|---|
| Dev | RoB | 73.6 | 22.4 | 10.6 |
| | BPT-base | 81.99 ±0.91 | 58.07 ±0.76 | 36.44 ±0.53 |
| Test | RoB | 72.9 | 19.1 | 9.1 |
| | BPT-base | 80.55 ±1.20 | 53.83 ±1.65 | 32.37 ±0.27 |

*example figure from Storks et. al (2021)

## Discussion

- BPTs are modular extension of Transformers: added to existing models without harming performance
- BPTs improve model interpretability, easily applicable to narrative/procedural text understanding tasks
- Limitations & future work:
  ○ Systematic generalization: BPTs inherit limitations of pre-trained LMs
  ○ Causal relation between breakpoints and generated outputs is unclear, can possibly be enhanced by new joint consistency losses

Code and experiments available on GitHub! https://github.com/allenai/situation_modeling