

Why natural language is the right vehicle for complex reasoning



Greg Durrett

NYU

April 14, 2022



Reasoning about text: entailment

premise

hypothesis

A dog is chasing a cat | *Two animals are running* — Transformer → **P entails H**

A dog is chasing a cat | *Two animals are sitting* — Transformer → **P contradicts H**

There's **some** reasoning happening, in a purely latent way:

X chasing \Rightarrow X is running

X chasing \Rightarrow X is not sitting

A dog is chasing a cat | *Three animals are running* — Transformer → **P entails H**

...but when the model's latent reasoning is flawed, it's hard to diagnose

Example: multi-hop QA where systems only do single-hop reasoning



Where latent reasoning breaks down

Can we just improve latent reasoning models? Better data, debiasing, contrastive learning, ...



Applies to **Canada** only (new dataset with this context: SituatedQA; Michael J.Q. Zhang and Eunsol Choi, 2021)

Are airsoft and BB guns the same? It's complicated!

End-to-end models don't model these nuances well.

We need justified reasoning in addition to answers.



Contrast: Theorem Provers

A dog is chasing a cat

$\exists d. \exists c. \exists e. \text{dog}(d) \wedge \text{cat}(c) \wedge \text{chase}(e) \wedge$
 $\text{agent}(e, d) \wedge \text{patient}(e, c)$

Two animals are running

$\exists a. \exists b. \text{animal}(a) \wedge \text{animal}(b) \wedge$
 $\text{running}(a) \wedge \text{running}(b)$



Theorem Prover

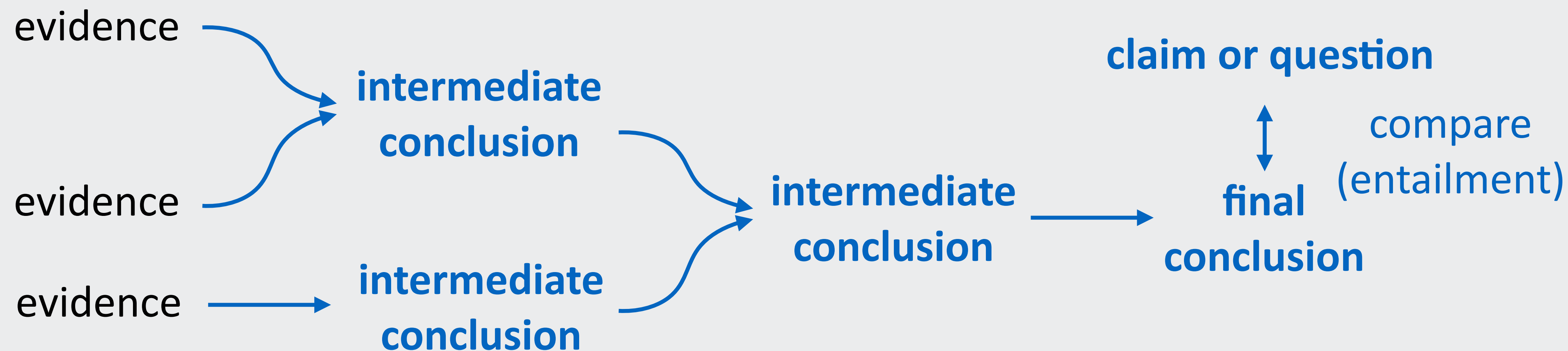
$\forall x. \forall e. \text{chase}(e) \wedge \text{agent}(e, x) \Rightarrow \text{running}(x)$
 $\forall x. \forall e. \text{chase}(e) \wedge \text{patient}(e, x) \Rightarrow \text{running}(x)$

$\forall x. \text{dog}(x) \Rightarrow \text{animal}(x)$
 $\forall x. \text{cat}(x) \Rightarrow \text{animal}(x)$

- ▶ Advantage: Articulates explicit intermediate reasoning states
- ▶ Disadvantage: requires high-coverage semantic formalism, parser into that formalism, and background knowledge hard to learn from data



Our vision



Use pre-trained models to do reasoning **directly in natural language**

- Combine **logical inference** (modus ponens, ...) and **lexical inference** (paraphrasing, ...)

Natural logic,
theorem provers

This approach

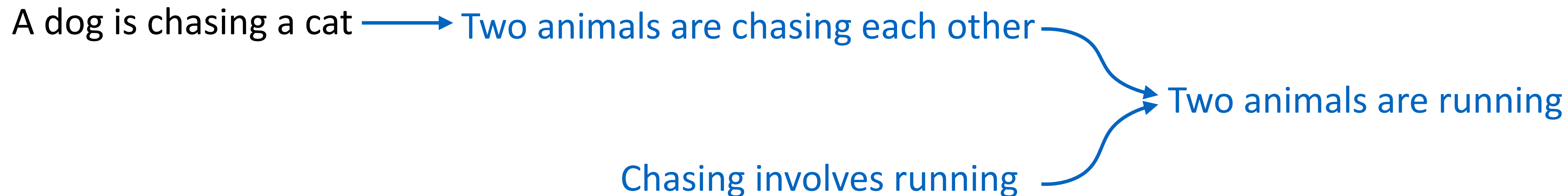
End-to-end models,
chain-of-thought

(Bill MacCartney and Manning, 2007;
Hai Hu et al. “MonaLog”, 2020, ...)

(Maxwell Nye et al., 2021;
Jason Wei et al., 2022)



Why natural language?



- ▶ **Expressive.** Text *is already* a broad-coverage semantic representation
- ▶ **Flexible.** Approaches operating over text can synthesize pre-trained models, Wikipedia, commonsense knowledge bases, ...
- ▶ **Interpretable** reasoning chains
- ▶ But: we need the **right data** and need to ensure our models are doing sound reasoning



Outline

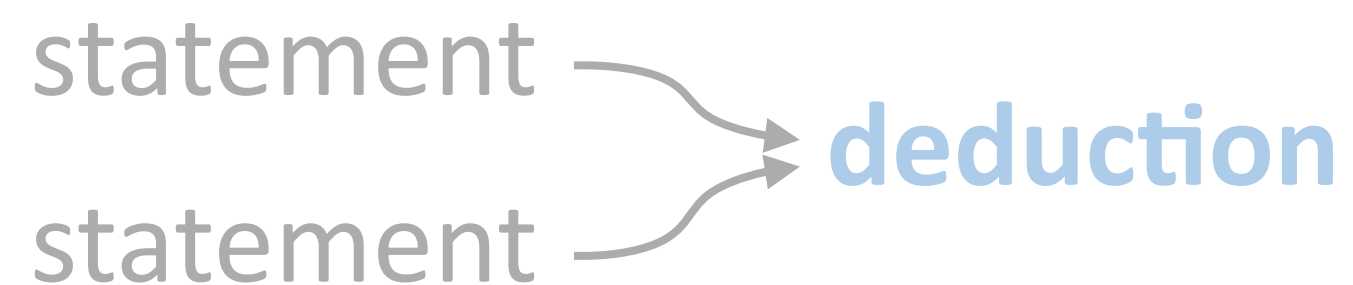
Entailment to verify QA

Jifan Chen, Eunsol Choi, GD. EMNLP-Findings21
Can NLI Models Verify QA Systems' Predictions?



Logically manipulating statements

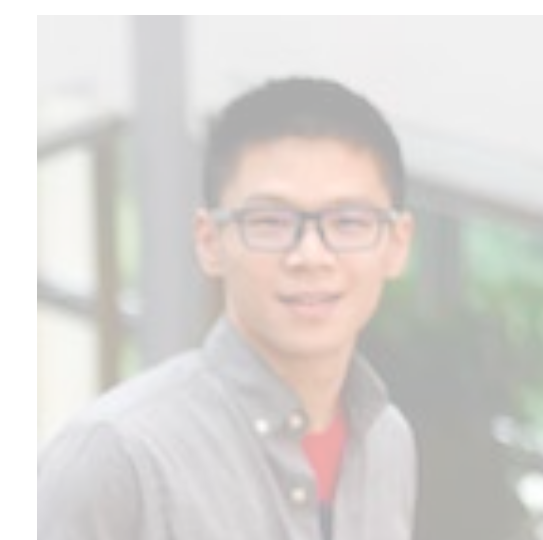
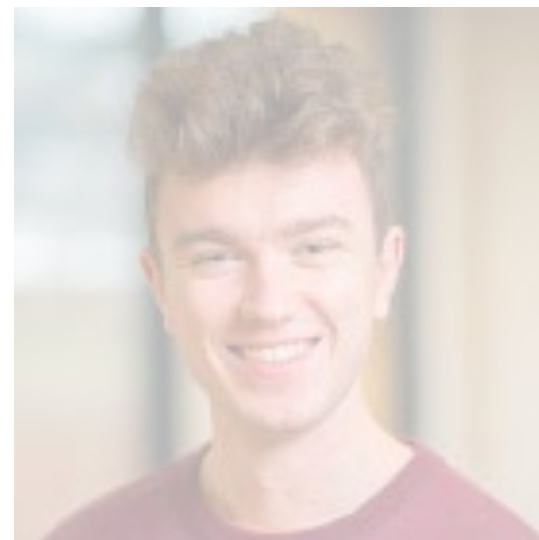
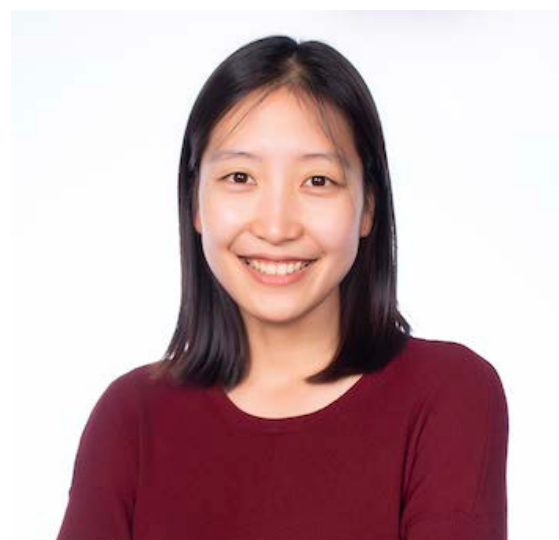
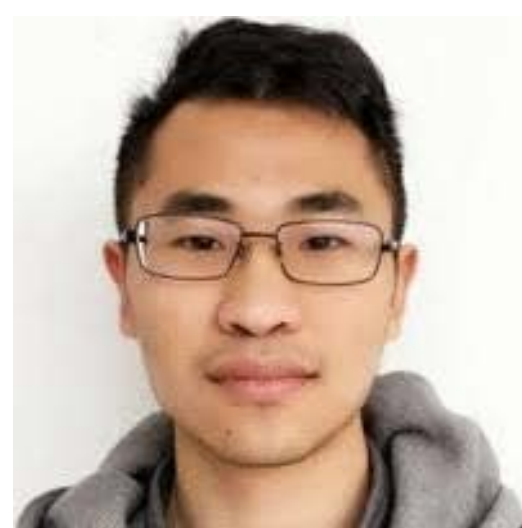
Kaj Bostrom, Xinyu Zhao, Swarat Chaudhuri, GD. EMNLP21
Flexible Generation of Natural Language Deductions



Kaj Bostrom, Zayne Sprague, Swarat Chaudhuri, GD. In submission.
Natural Language Deduction through Search over Statement Compositions

Improving diverse generation (if time)

Jiacheng Xu, GD. NAACL22.
Massive-scale Decoding for Text Generation using Lattices





Verifying Reading Comprehension

Ted Danson



RoBERTA QA model

Who plays the bad guy in the Good Place?

The first season of the fantasy comedy television series The Good Place [...] The series focuses on Eleanor Shellstrop (Kristen Bell), a woman who wakes up in the afterlife and is introduced by Michael (Ted Danson) to a Heaven-like utopia [...]

Assume a base QA system with a latent reasoning process. Can we check the answer?

- ▶ Can better determine if question is unanswerable
- ▶ Can improve confidence, “selective QA setting” (Amita Kamath et al. 2020)
- ▶ Can validate presuppositions in the question (Najoung Kim et al., 2021)



Our Method

The series focuses on Eleanor Shellstrop (Kristen Bell) , a woman who wakes up in the afterlife and is introduced by Michael (Ted Danson) [...]

Decontextualization

Eunsol Choi et al. (2021)

**standalone
statement (premise)**

*The series **The Good Place** focuses on Eleanor Shellstrop (Kristen Bell) , a woman who wakes up in the afterlife and is introduced by Michael (Ted Danson) [...]*

Who plays the bad guy in the Good Place?

**Question-to-statement
conversion**

Dorottya Demszky et al. (2018)
(upgraded to use T5-3B)

Ted Danson

NLI model

false

***Ted Danson** plays the bad guy in the Good Place.*

hypothesis

(in this case: right for the wrong reasons)



System

Verifier

Answer sentence

Decontextualization

Q+A

Question-to-statement
conversion

NLI model

Confidence

QA System

Answer

RoBERTA QA model

Question

Context

→ Verifier → Confidence

- ▶ Can use these confidence values to reject low-confidence answers



Results: Unanswerable Questions

- ▶ Train QA system on (En) SQuAD 1.1 (answers every question), run on SQuAD 2.0 (contains unanswerable questions), use the verifier to reject bad answers

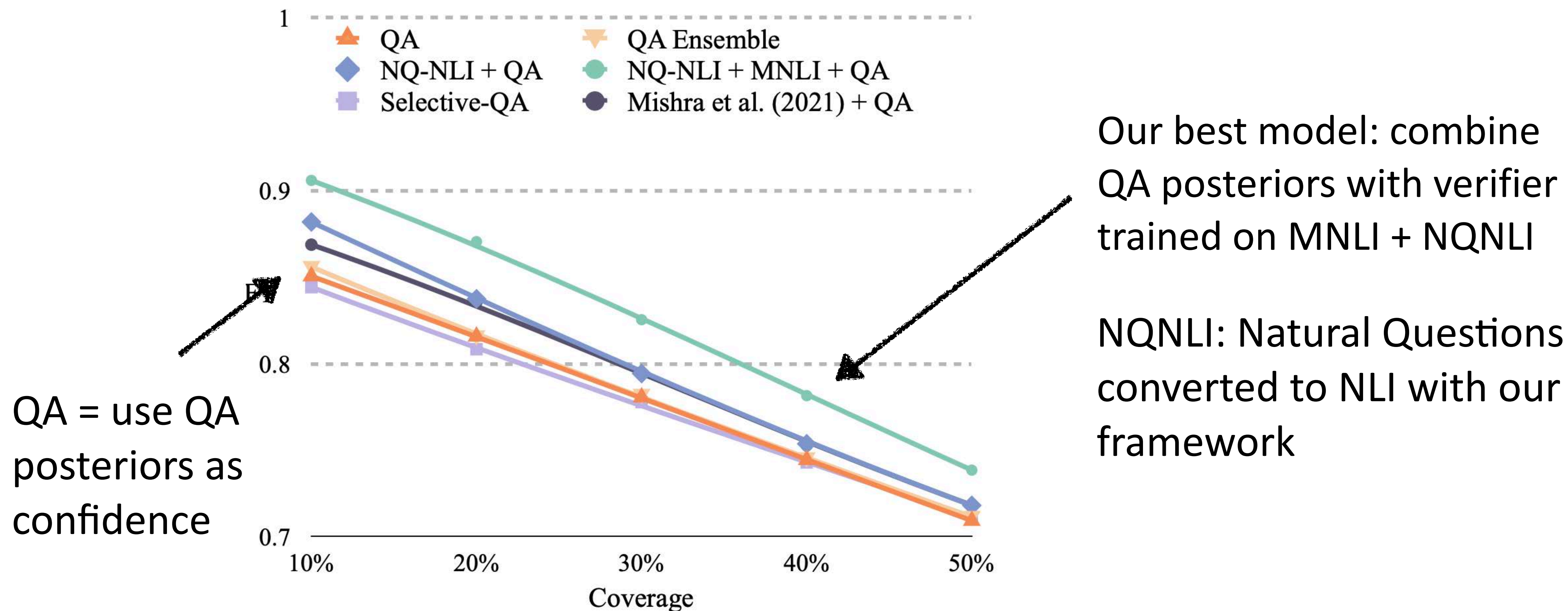


- ▶ A RoBERTa MNLI model can reject **78.5%** of unanswerables, accept **82.5%** of answerables. Good zero-shot performance (NLI model is not trained on SQuAD + pipeline is not optimized end-to-end)
- ▶ Can use MNLI off-the-shelf here because we have single-sentence premises
(unlike Anshuman Mishra et al., 2021; Wenpeng Yin et al., DocNLI, 2021)



Results: Selective QA (5 datasets)

- Base QA: BERT-Large on (En) NaturalQuestions. Target: NQ + 4 out-of-domain (En) sets.



NLI works well as an answer verifier



Errors

- ▶ **QA conversion / decontextualization errors are rare.** NL manipulation is great with big models like T5-3B!
- ▶ **Entailment errors** are more common. But sometimes, the entailment model disagrees with the QA dataset **and is right**

Reformulated Q+A: ***John von Neumann** developed the central processing unit (CPU).*

Context: *On June 30, 1945, before ENIAC was made, mathematician **John von Neumann** distributed the paper entitled First Draft of a Report on the EDVAC. It was the outline of a stored-program computer that would eventually be completed in August 1949.*

- ▶ John von Neumann is marked as the gold answer (debatable), NLI model disagrees
- ▶ Manipulation of these examples makes it easy to evaluate reasoning. **Is this evidence really sufficient to validate the answer?**



Takeaways

- ▶ We can manipulate question-answer pairs and evidence sentences **in natural language** and use NLI to check QA answers
- ▶ Manipulation was highly reliable, and the two operations we had were sufficient to allow us to employ a pre-existing model (NLI)
- ▶ NLI can improve calibration for QA and lets us audit both our models and datasets



Outline

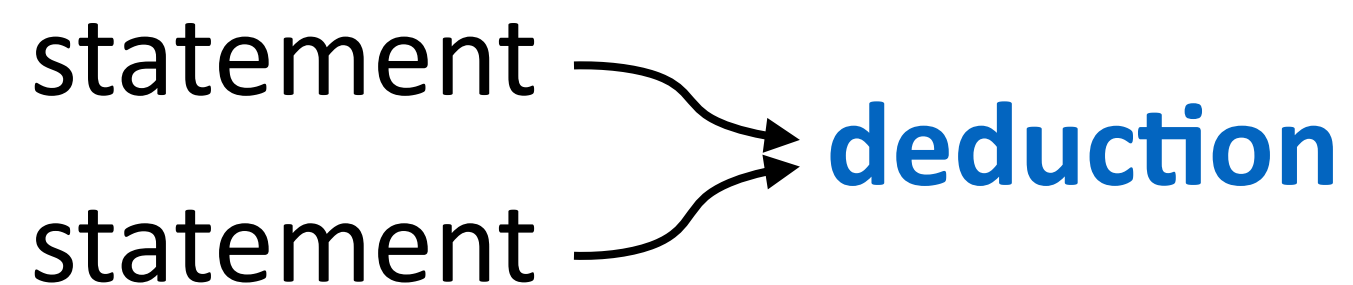
Entailment to verify QA

Jifan Chen, Eunsol Choi, GD. EMNLP-Findings21
Can NLI Models Verify QA Systems' Predictions?



Logically manipulating statements

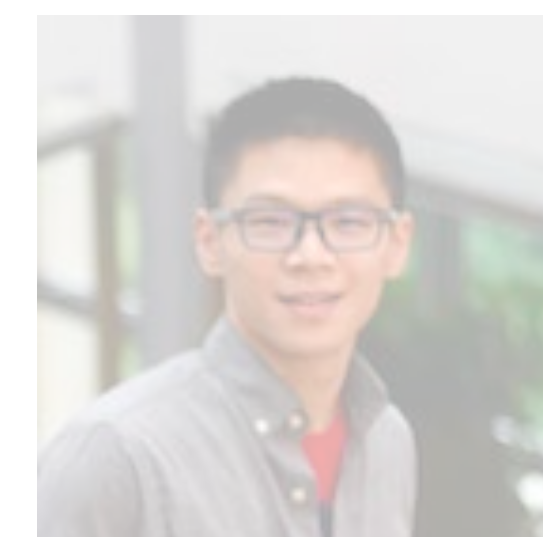
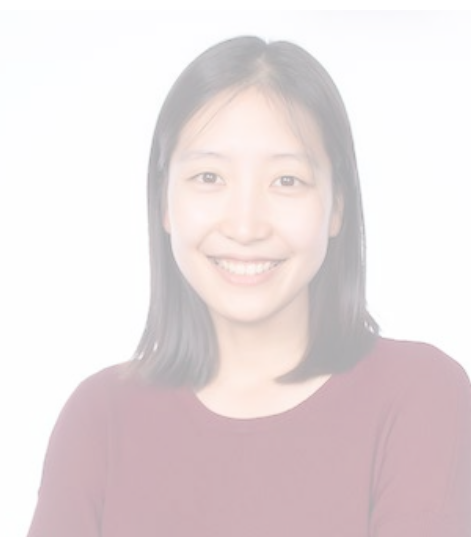
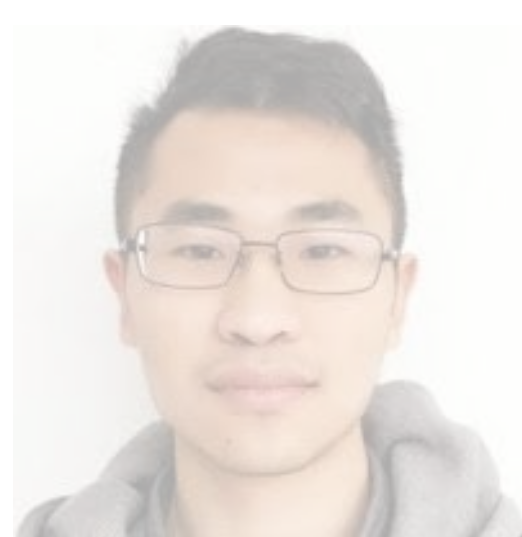
Kaj Bostrom, Xinyu Zhao, Swarat Chaudhuri, GD. EMNLP21
Flexible Generation of Natural Language Deductions



Kaj Bostrom, Zayne Sprague, Swarat Chaudhuri, GD. In submission.
Natural Language Deduction through Search over Statement Compositions

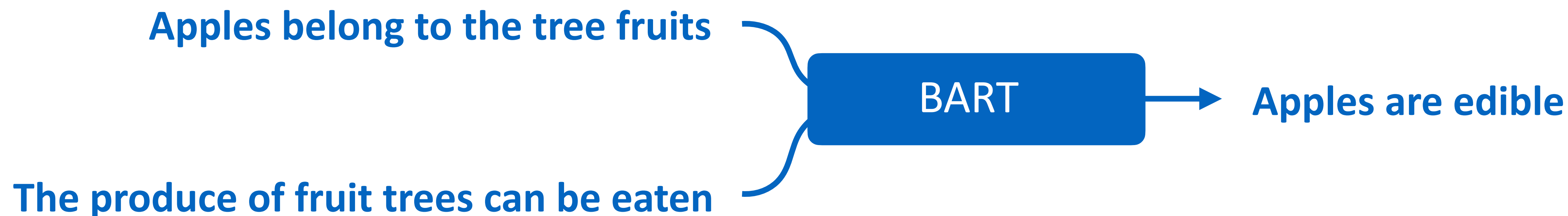
Improving diverse generation

Jiacheng Xu, GD. NAACL22.
Massive-scale Decoding for Text Generation using Lattices





Natural Language Deduction



Natural language deduction: place a distribution over the set of valid (and useful) conclusions

1. Can we automate collecting this kind of data at scale?
2. Can we chain these inferences together into multiple steps?



Our Approach

We characterize inference as a blend of two processes:

Logical inference

Fruits are edible. Apples are a fruit.
→ Apples are edible.

Invariant w.r.t. lexical content
Easy to describe with a concise set of rules
Hard to learn distributionally



Automatic template-based data generation

Lexical inference

edible \leftrightarrow can be eaten

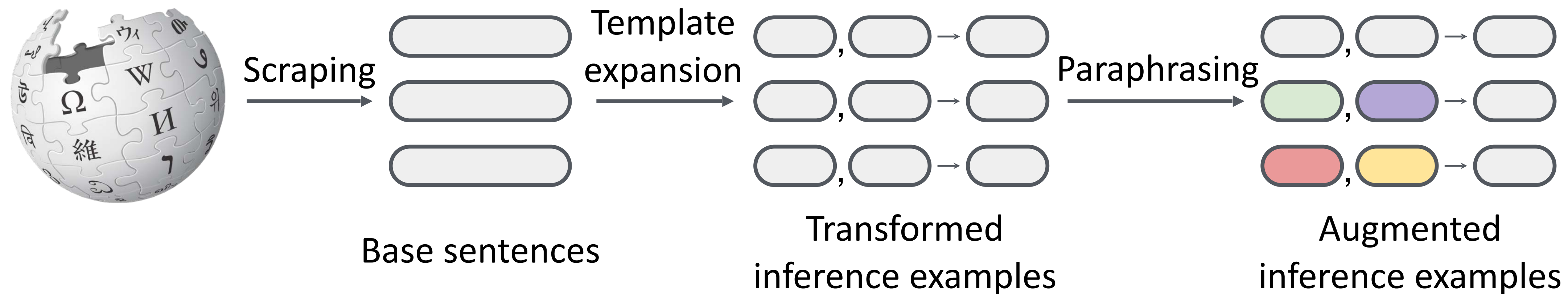
Idiosyncratic, depends on words
Hard to describe with a concise set of rules
Can be learned distributionally



Transfer learning from a pre-trained LM
Data augmentation with paraphrasing



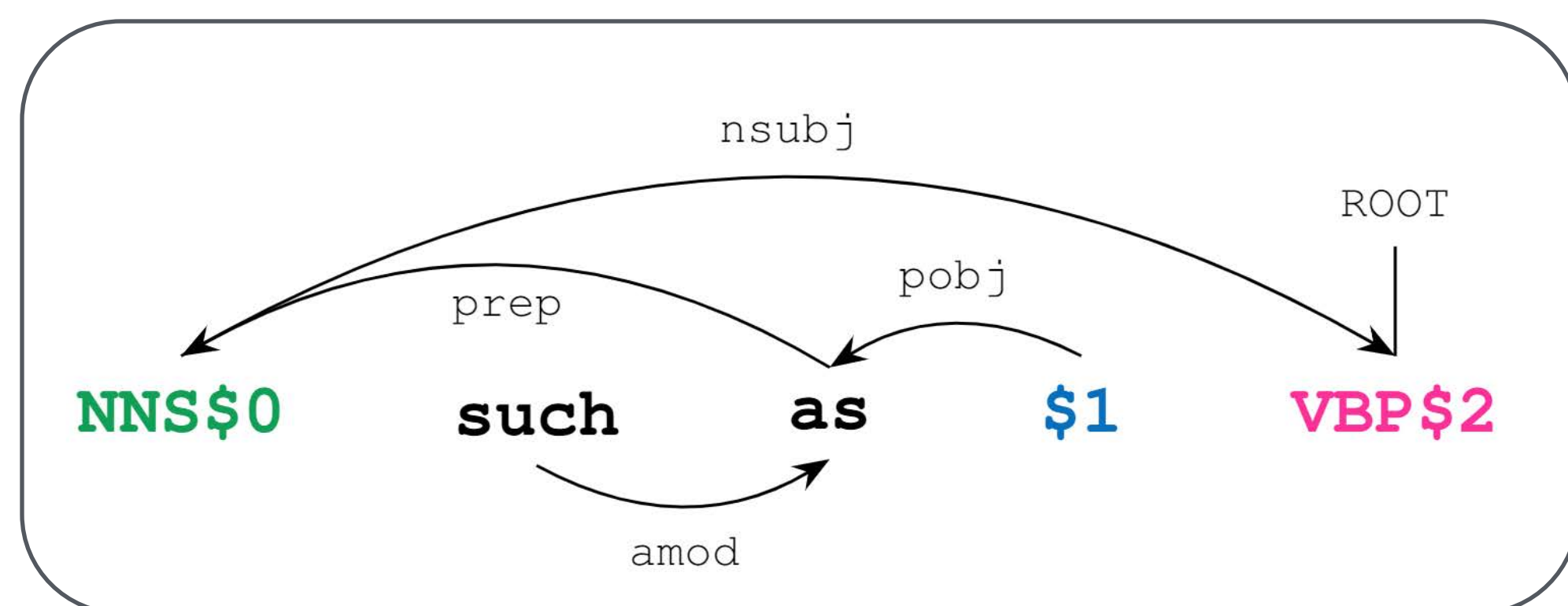
Data Generation





Data Generation

1. Source sentence scraping using Hearst patterns (6 patterns for this type of *substitution* reasoning)



dependency parsing (spaCy)

Dependency
index

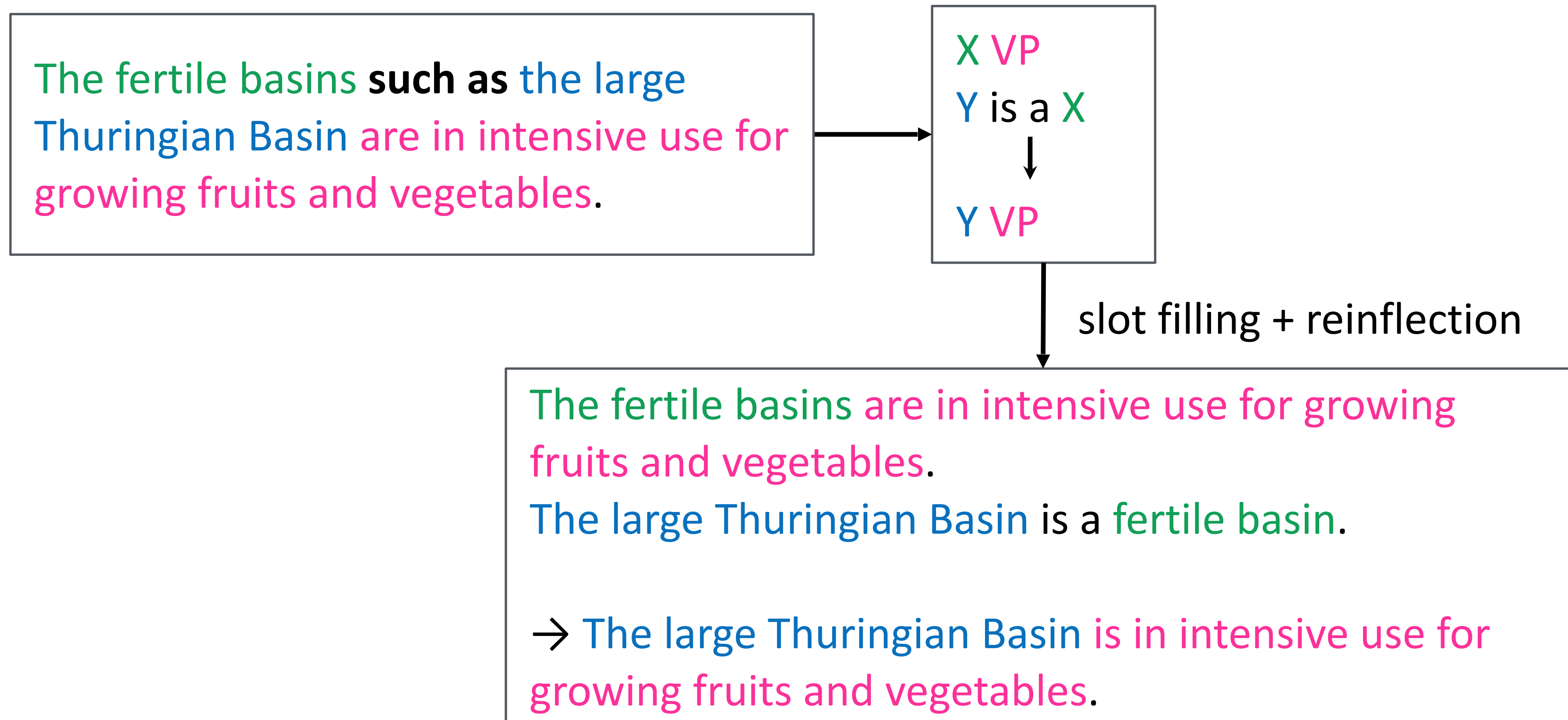
retrieval

The fertile basins such as the large Thuringian Basin are in intensive use for growing fruits and vegetables.



Data Generation

2. Template expansion (1 template per Hearst pattern)





Data Generation

3. Paraphrasing

The fertile basins are in intensive use for growing fruits and vegetables.

The large Thuringian Basin is a fertile basin.

→ The large Thuringian Basin is in intensive use for growing fruits and vegetables.

Paraphrasing adds noise, but only to the input. We find the model **still does sound** reasoning and can handle more lexical variation

Automatic paraphrasing model (PEGASUS)

Growing fruits and vegetables requires a lot of fertile basins.

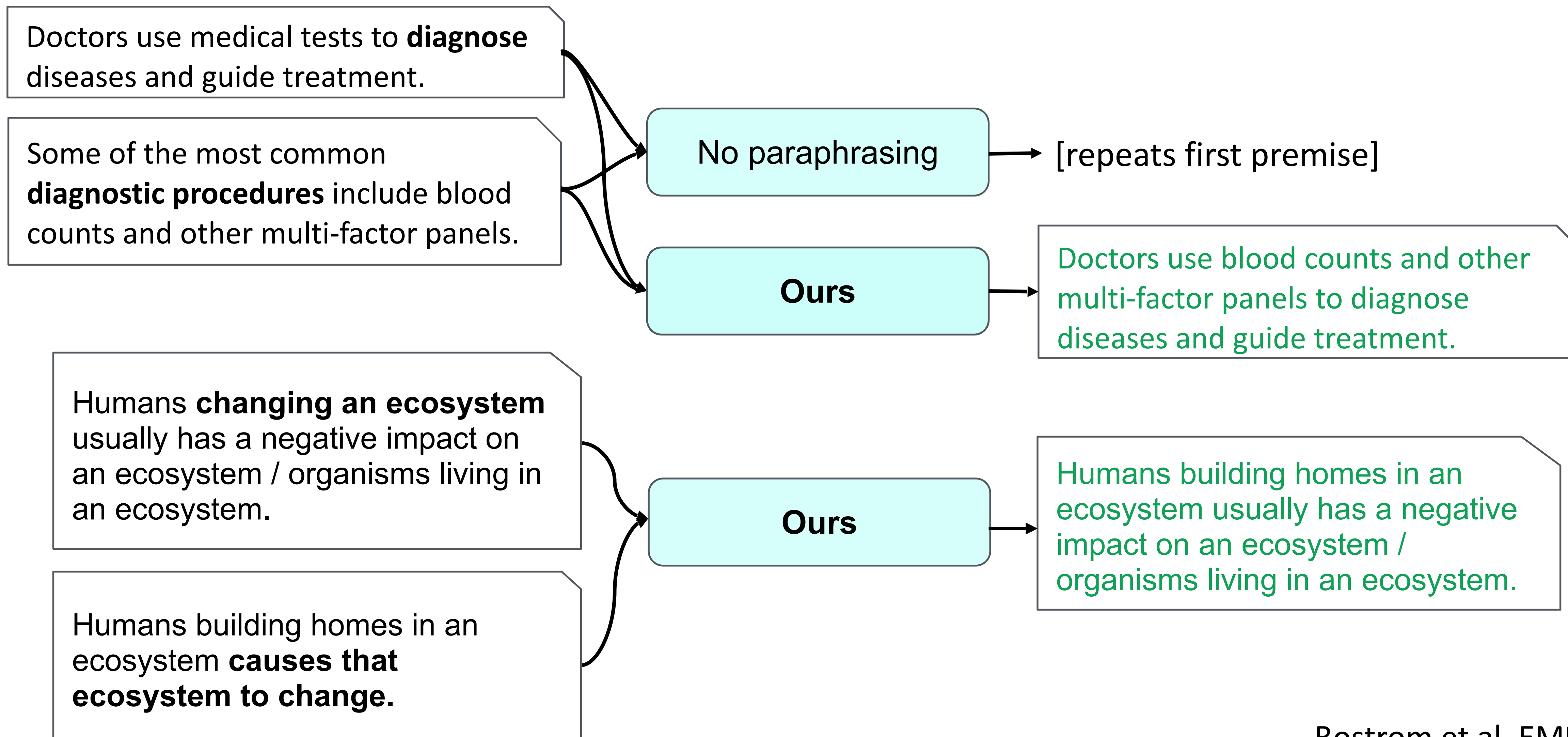
The Thuringian Basin is fertile.

→ The large Thuringian Basin is in intensive use for growing fruits and vegetables.



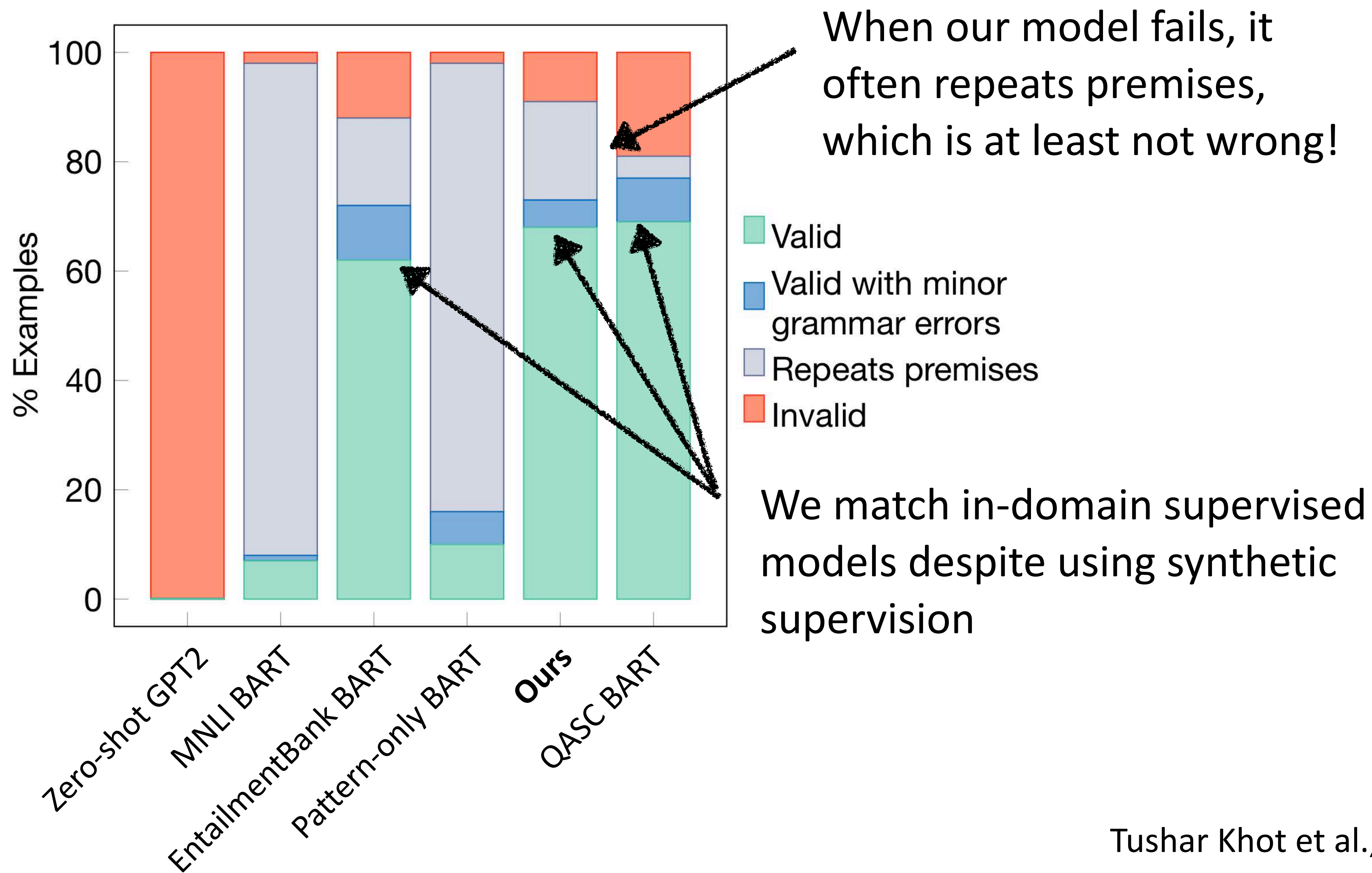
Examples

- ▶ BART trained on 126k examples we automatically collect from Wikipedia





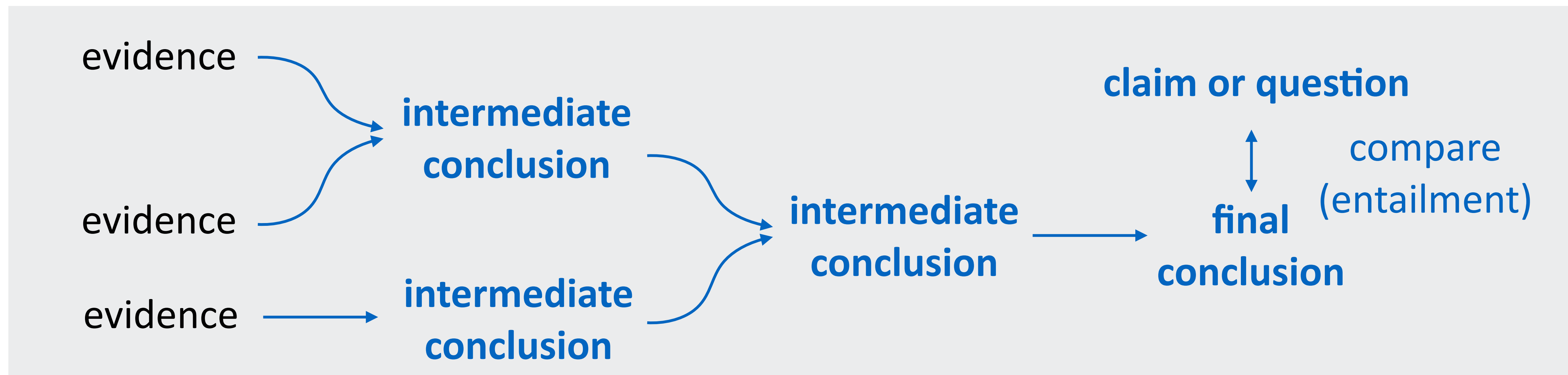
Results: QASC Human Eval





Multi-step Deduction

- ▶ We showed something that works well for single steps. Let's return to our goal...



- ▶ Goal: dynamically apply operations (including deduction and decontextualization) to give the conclusions we need
- ▶ This is a hard search problem: intermediate states are potentially all natural language sentences that can be reasonably generated from the evidence



Multi-step Deduction: Setup

evidence

- S_1 Paper is recyclable.
- S_2 Recyclable means old material can be converted into new material
- S_3 Cardstock is a type of paper

hypothesis

Old cardstock can be turned into new cardstock.

- ▶ Collection of evidence (science domain, taken from EntailmentBank) and hypothesis
- ▶ Can we prove the hypothesis deductively using our generative step model? (not just throwing everything into a discriminative model)



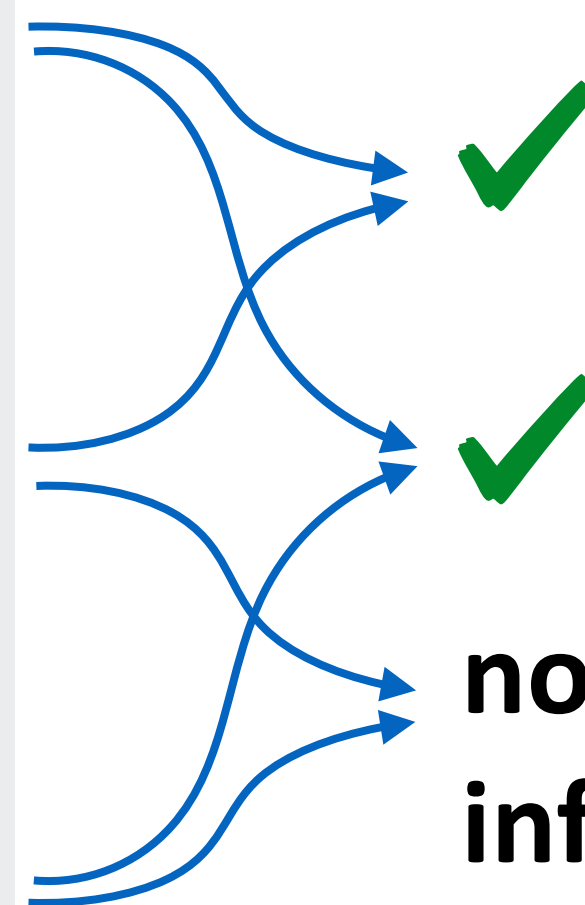
Search Heuristics

evidence

S_1 Paper is recyclable.

S_2 Recyclable means old material can be converted into new material

S_3 Cardstock is a type of paper



hypothesis

Old cardstock can be turned into new cardstock.

- ▶ Search frontier of pairs (s_i, s_j) of sentences we can combine — how to prioritize?
- ▶ Goal-conditioned heuristic: learn a model $g(s_i, s_j, h)$ — how likely will combining s_i and s_j eventually lead to h ? Requires training on EntailmentBank
- ▶ **Search and deduction are decoupled.** Search conditions on the hypothesis, but the deduction itself uses only the premises

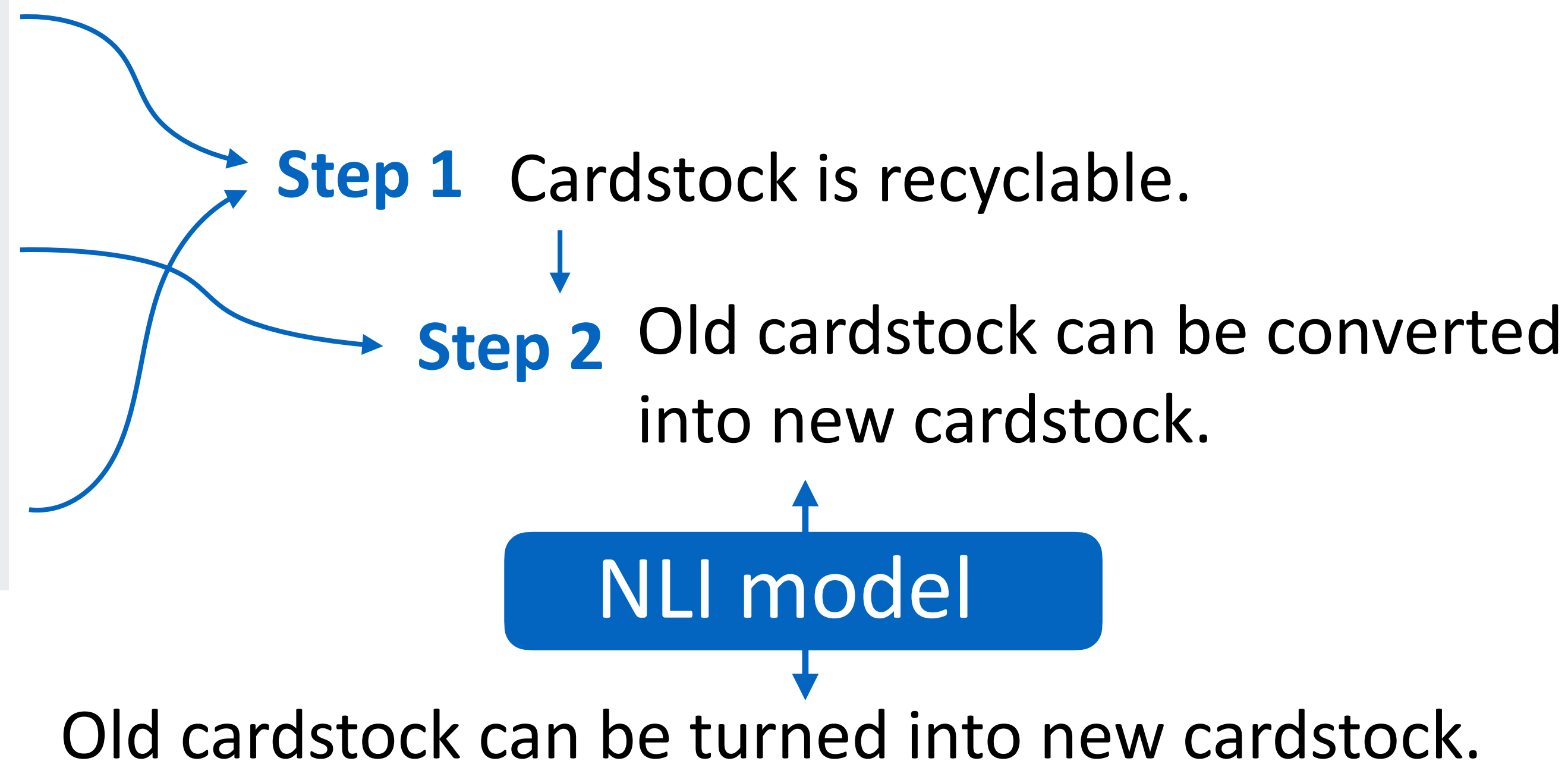


Multi-Step Deduction

evidence

S_1 Paper is recyclable.
 S_2 Recyclable means old material can be converted into new material
 S_3 Cardstock is a type of paper

hypothesis



- Repeatedly prove statements and expand the search space, then check if each entails the claim with an NLI model (fine-tuned for this domain)



Evaluation: EntailmentBank

- Input: 25 English premises and a potentially true hypothesis.
Goal: classify hypothesis as true/false based on premises

premises (+ 22 distractors)

A planet rotating causes cycles of day and night on that planet.

Earth rotating on its tilted axis occurs once per day.

Earth is a kind of planet.

true goal statement

Goal: The earth rotating on its tilted axis causes the cycles of day and night on earth.



random distractor goal

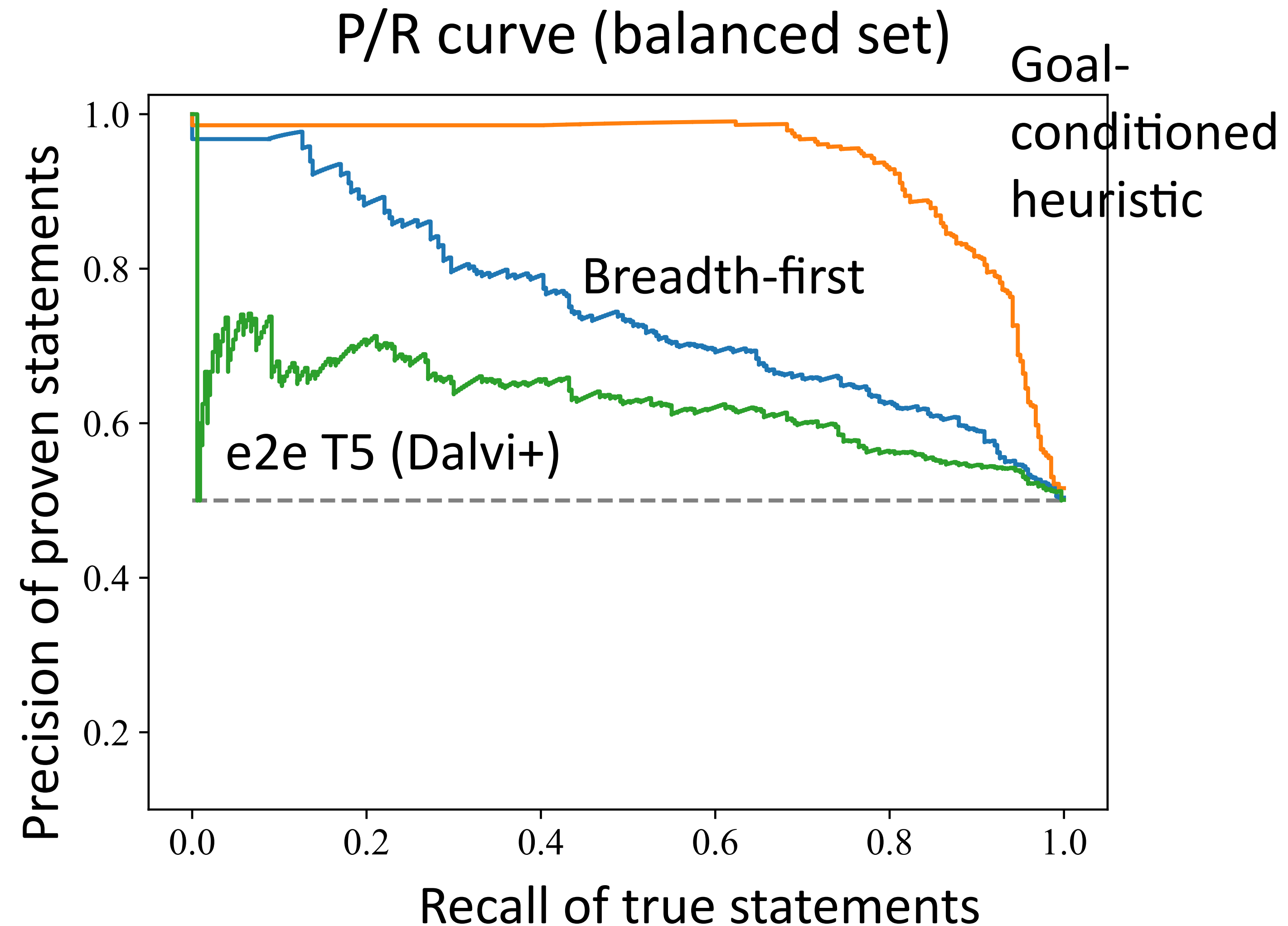
Goal: Looking at the moon has less of a negative impact on the eyes.





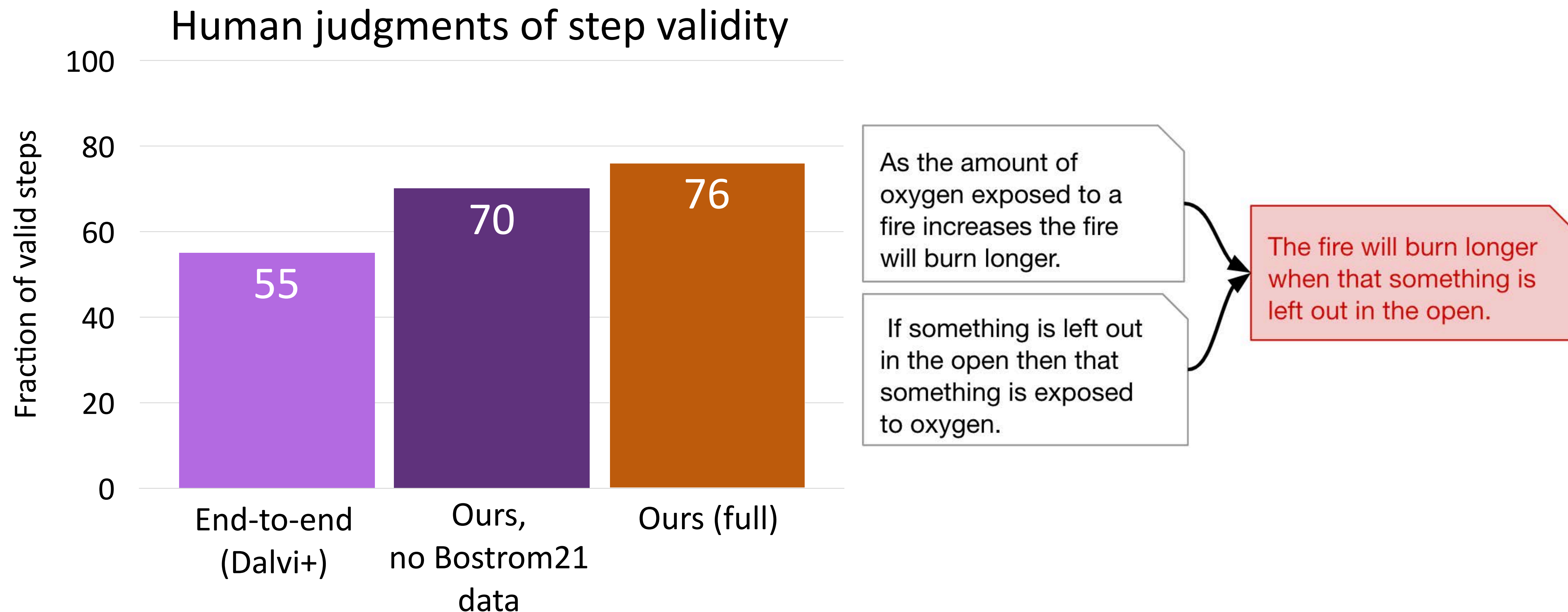
Results

- ▶ Input: 25 premises and a potentially true hypothesis. Goal: classify hypothesis as T/F
- ▶ Baseline: a pure end-to-end T5 model (Dalvi et al., 2021). Rank outputs by generation probability and apply a threshold to classify
- ▶ Separating concerns of search and deduction is important, and a good heuristic is important





Results: Individual Steps



- ▶ Our model is substantially better than end-to-end T5
- ▶ Some gray area about what's an error or not



Takeaways

- ▶ Our deduction models can capture broad-domain reasoning patterns **with little human training signal**, no logical forms
- ▶ Our models are **expressive** (can represent statements across several datasets) and are **flexible**
- ▶ A multi-step reasoning system founded on our deduction principles outperforms a pure end-to-end approach. **Structuring the reasoning this way helps!**
- ▶ Ongoing work: learning a backward model to do abductive inference, be able to hypothesize missing premises
- ▶ Ongoing work: take a step towards symbolic components in the model



Outline

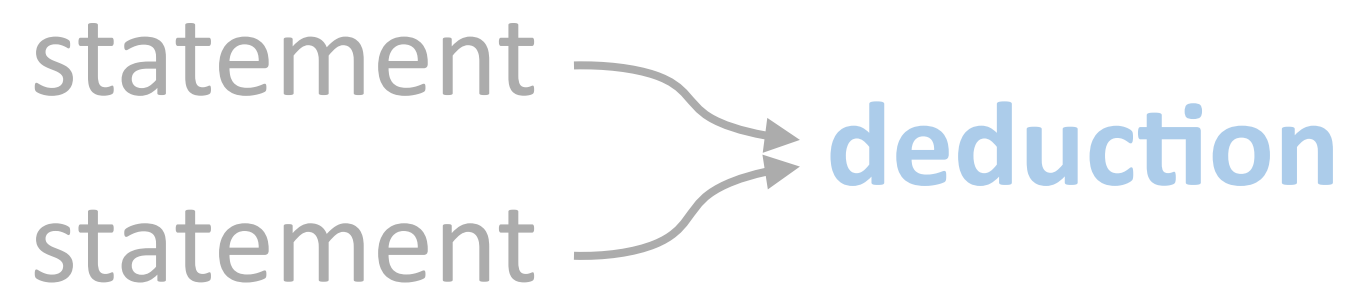
Entailment to verify QA

Jifan Chen, Eunsol Choi, GD. EMNLP-Findings21
Can NLI Models Verify QA Systems' Predictions?



Logically manipulating statements

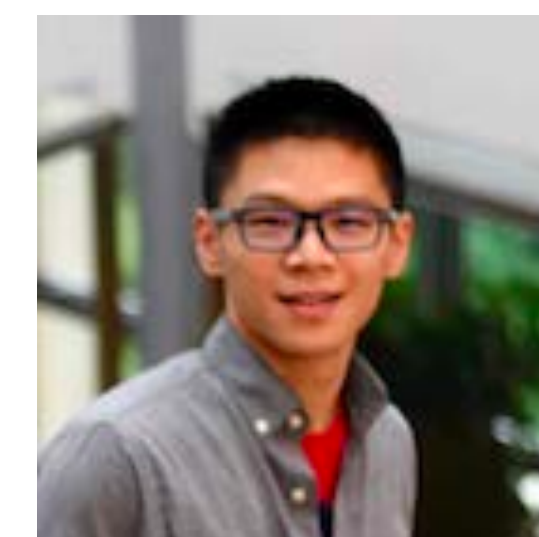
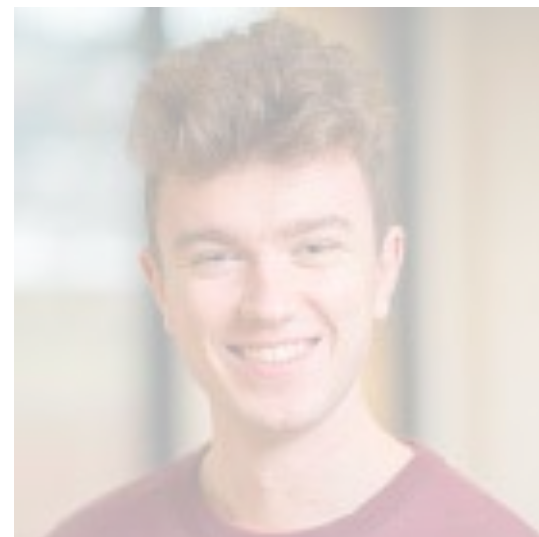
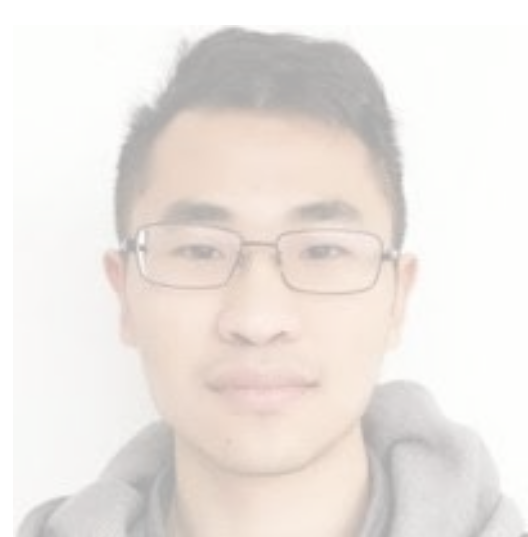
Kaj Bostrom, Xinyu Zhao, Swarat Chaudhuri, GD. EMNLP21
Flexible Generation of Natural Language Deductions



Kaj Bostrom, Zayne Sprague, Swarat Chaudhuri, GD. In submission.
Natural Language Deduction through Search over Statement Compositions

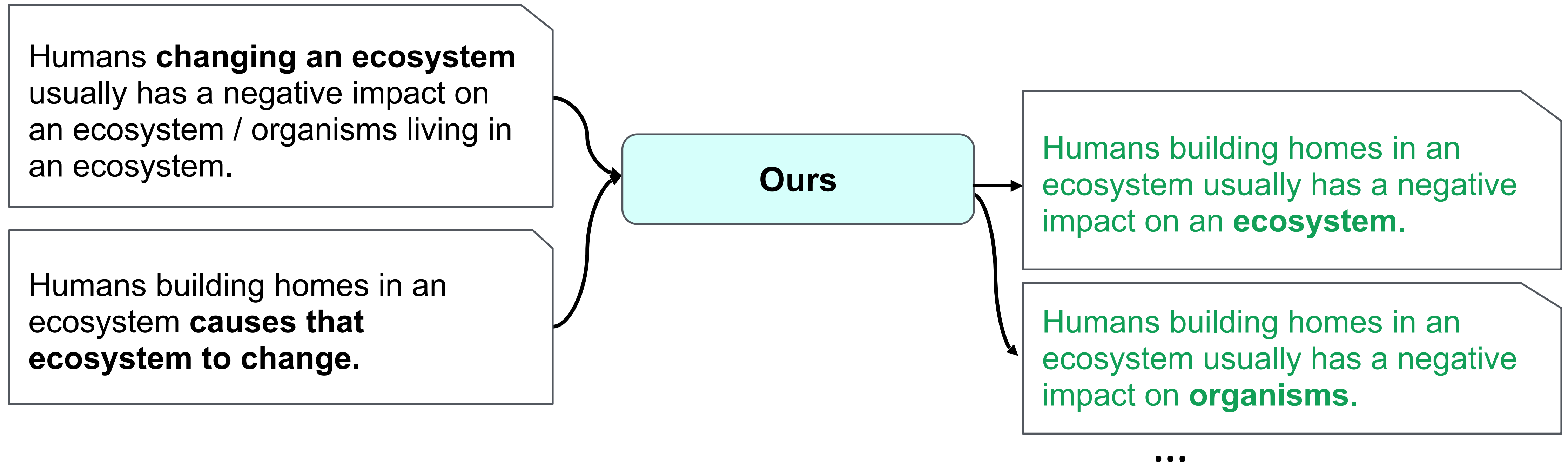
Improving diverse generation

Jiacheng Xu, GD. NAACL22.
Massive-scale Decoding for Text Generation using Lattices





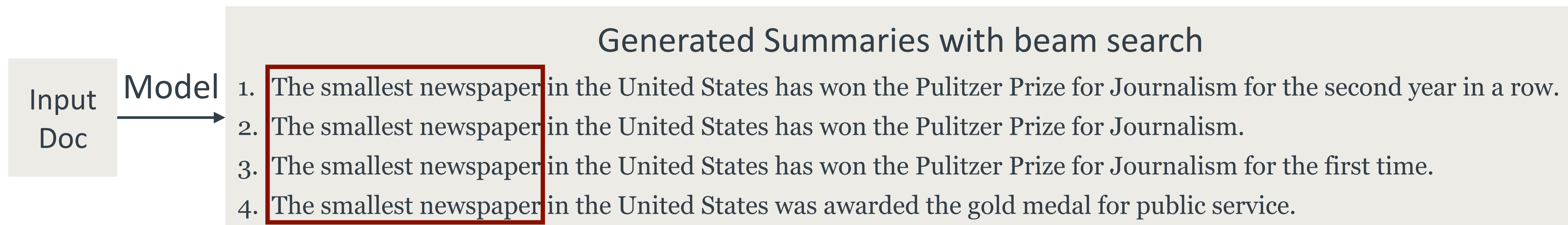
Advancements in Generation



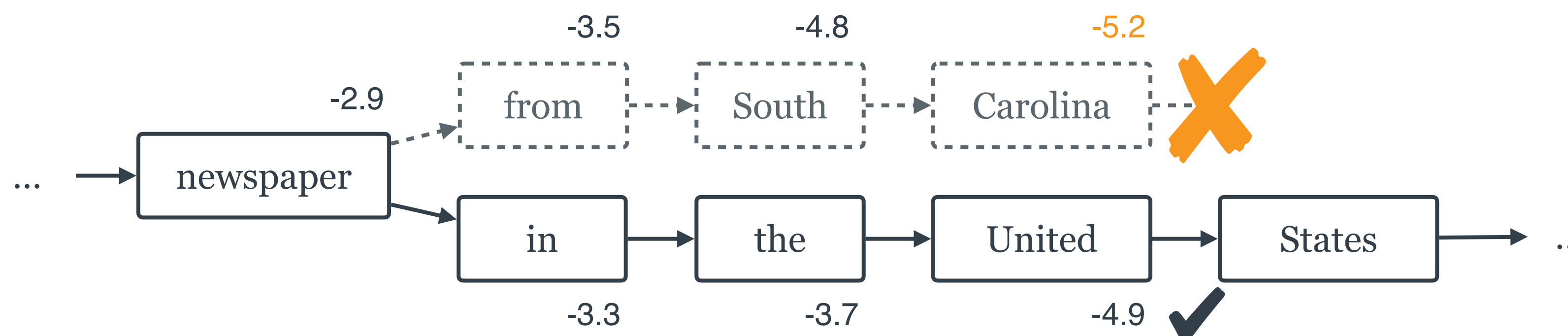
- ▶ Multiple correct generations — we're not sure which ones might be useful
- ▶ How can we access **as many generation candidates as possible?**



Getting Diverse Summaries



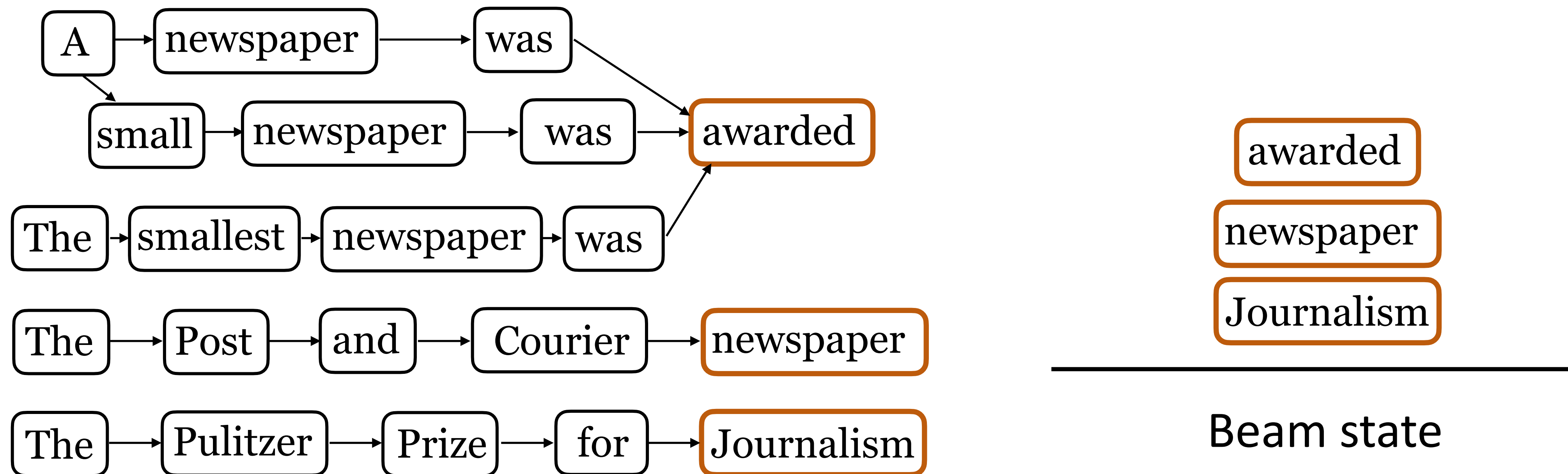
- ▶ Top summaries are similar and wrong! (not the smallest). **Too much redundancy**
- ▶ Other useful states were explored (info about location) but **pruning eliminated them**



- ▶ We're going to fix **two problems** with beam search to improve diversity



Reducing Redundancy with Recombination



- ▶ Hypotheses are stored in a lattice. Beam search now operates over nodes in this lattice
- ▶ Expanding a node continues all of the hypotheses ending in that node



Hypothesis Recombination

Recombine partial generated hypotheses A and B if:

- ▶ The last n tokens of A and B are the same ($n = 3$ or 4)
- ▶ A and B are roughly the same length

Prefix A: A small newspaper was awarded

Prefix B: A newspaper was awarded



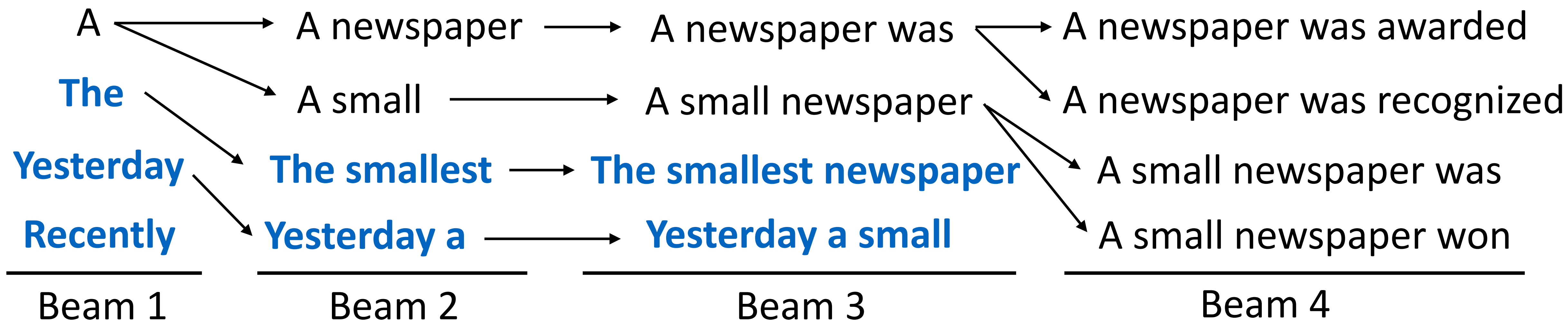
Assumption: if these criteria are met,
the rest of the summary will be similar:
 $P(\mathbf{y} \mid \text{document}, \mathbf{A}) \approx P(\mathbf{y} \mid \text{document}, \mathbf{B})$

- ▶ For summarization: we find that when this heuristic applies, **~70%** of time the greedy completion of the summary is exactly the same.
- ▶ When these distributions match, merging states in the lattice is completely okay!



Reducing Pruned States

- ▶ **Beam search wastes time** — most expanded hypotheses are eventually discarded. (Makes sense if you want the one-best, but not to get diverse options)

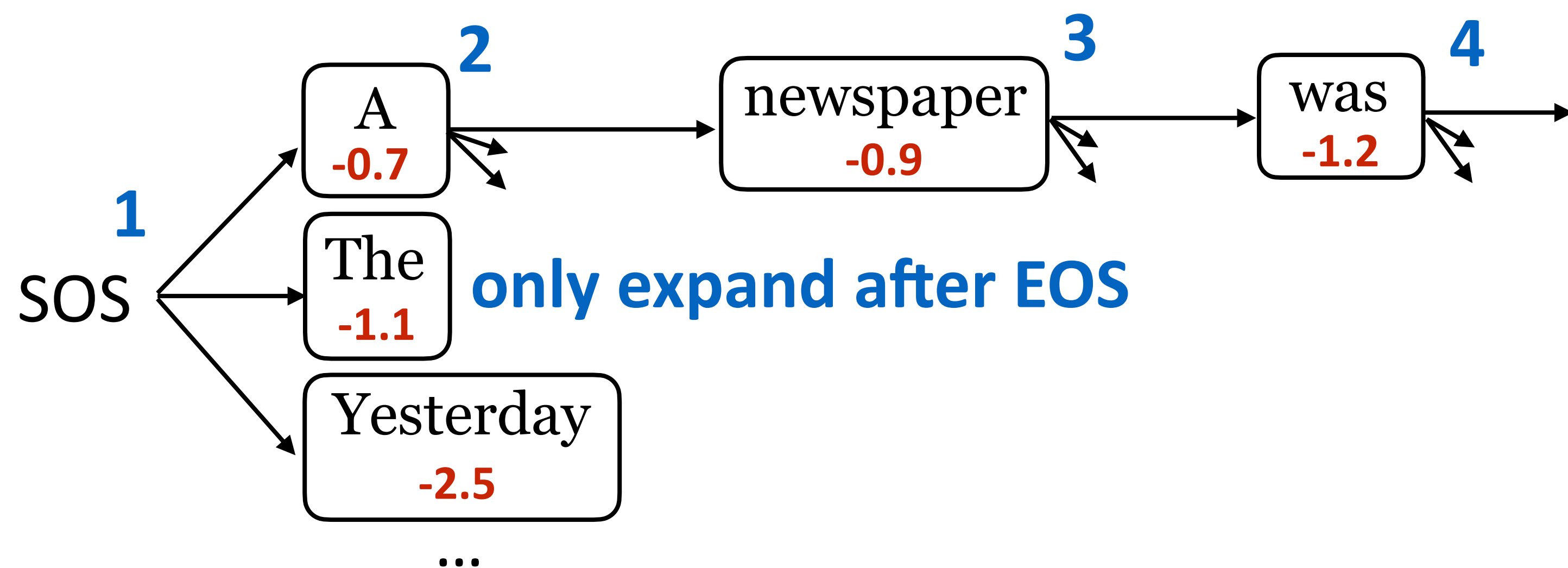


- ▶ Every **blue step** ultimately got pruned, even though these could be good summaries



Reduce Pruning with BFS/DFS

- ▶ We want a search algorithm where **every explored state** is on some finished path
- ▶ Use a modified best-first search with a depth-first stage: greedily expand each node until an EOS token is reached
- ▶ Expand by model score, shown in **red** below



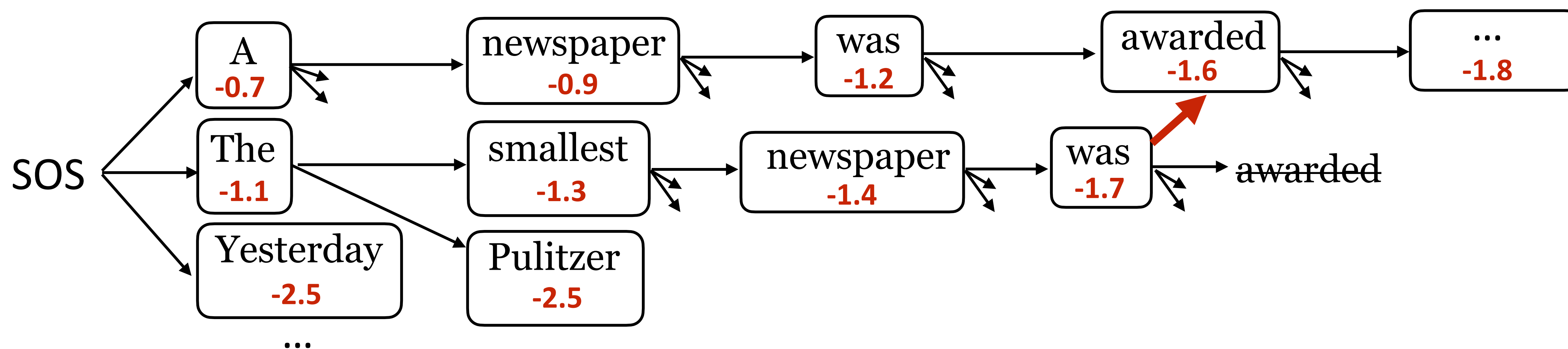
(expansion order)

...but for our depth-first stage, we continue expanding until we reach EOS. This greedy path is typically high-scoring.



Reduce Pruning with BFS/DFS

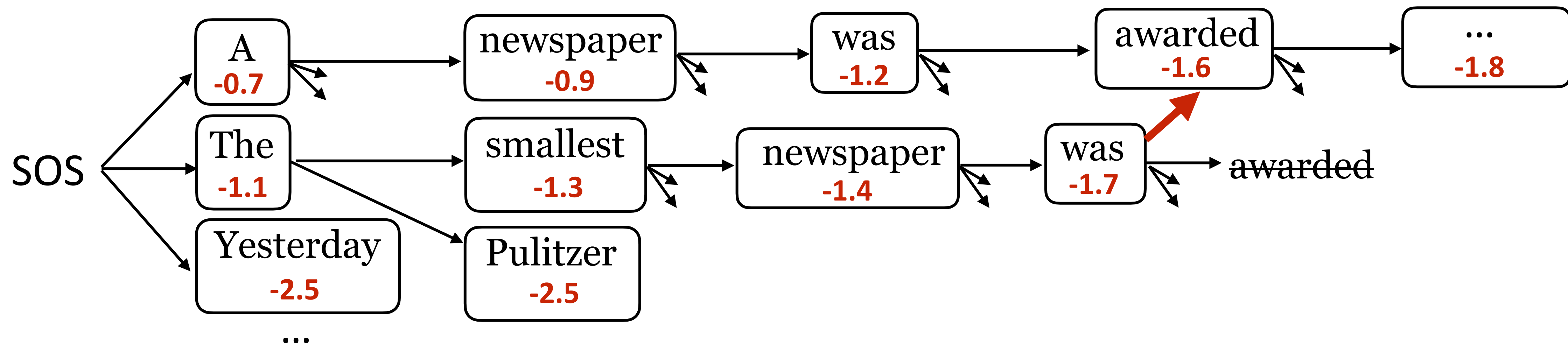
- Effective when combined with recombination: subsequent paths that are explored may overlap with earlier ones

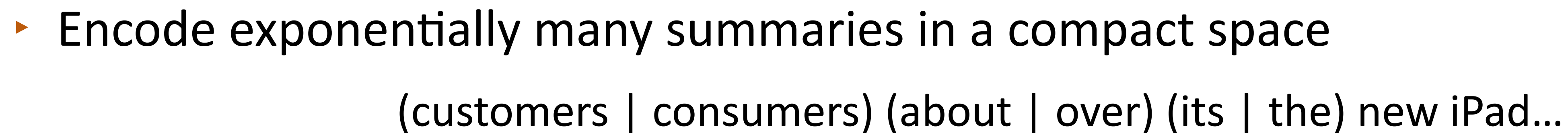
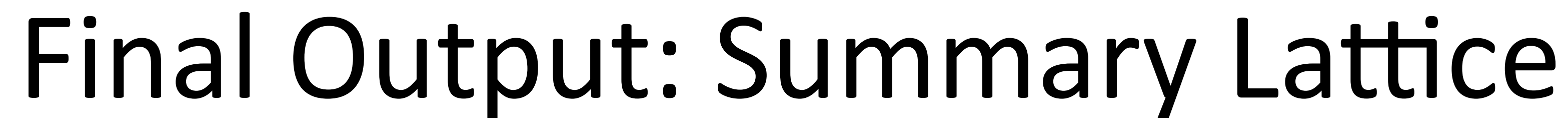




Overall Algorithm

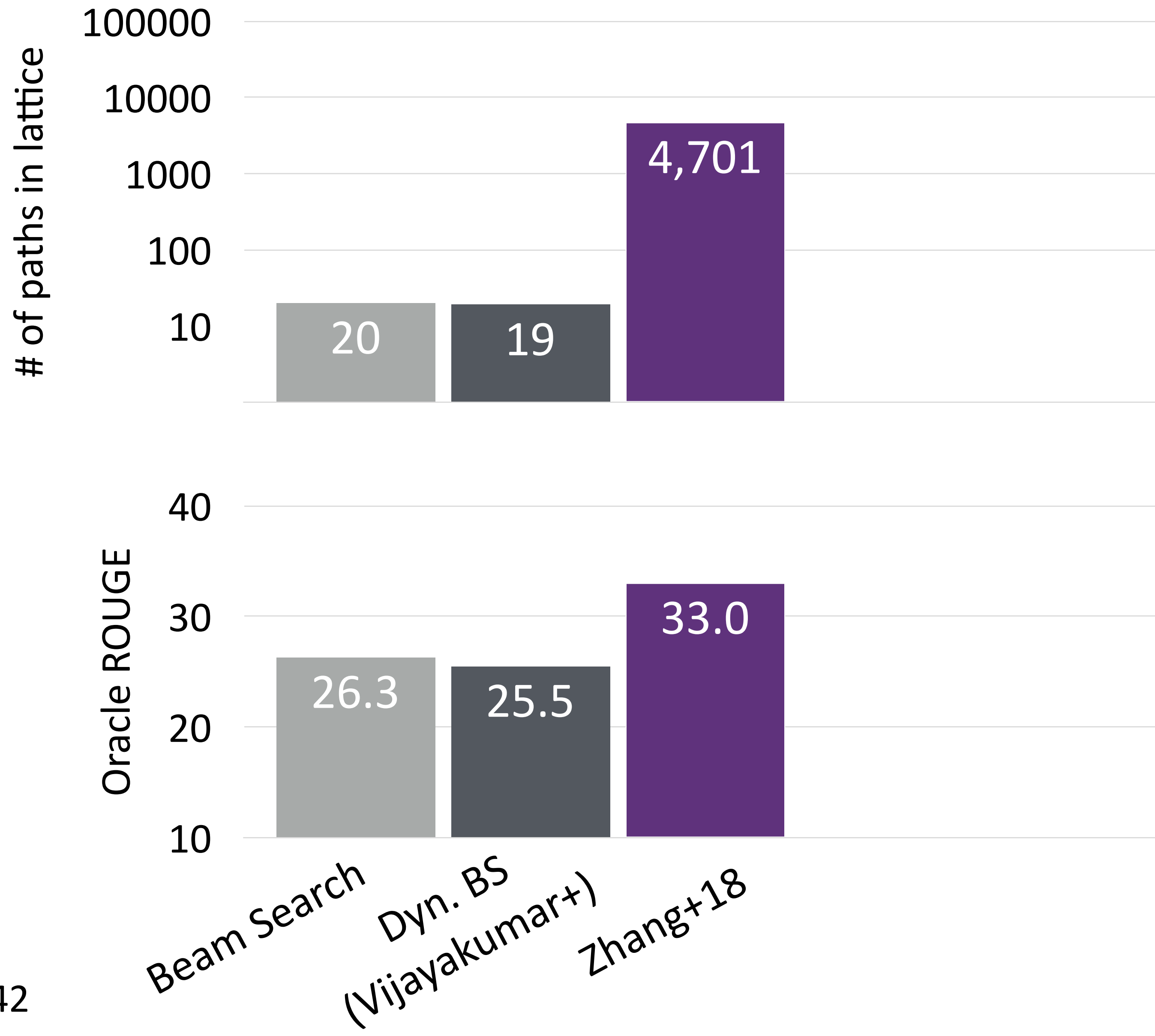
- ▶ Explore paths with best-first/depth-first search
- ▶ Merge states when redundancy is identified
- ▶ Construct a **lattice** of many possible options





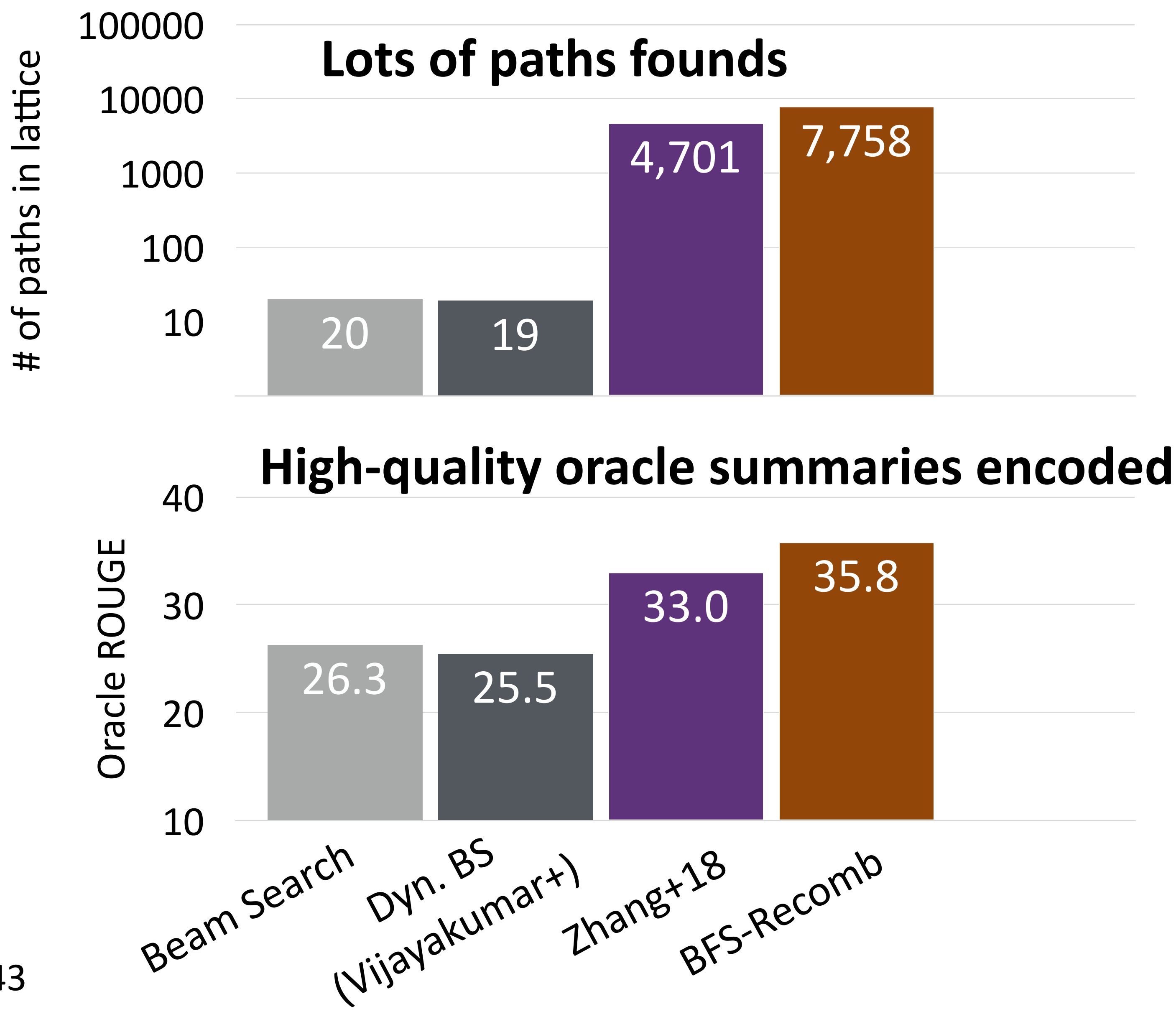


Results



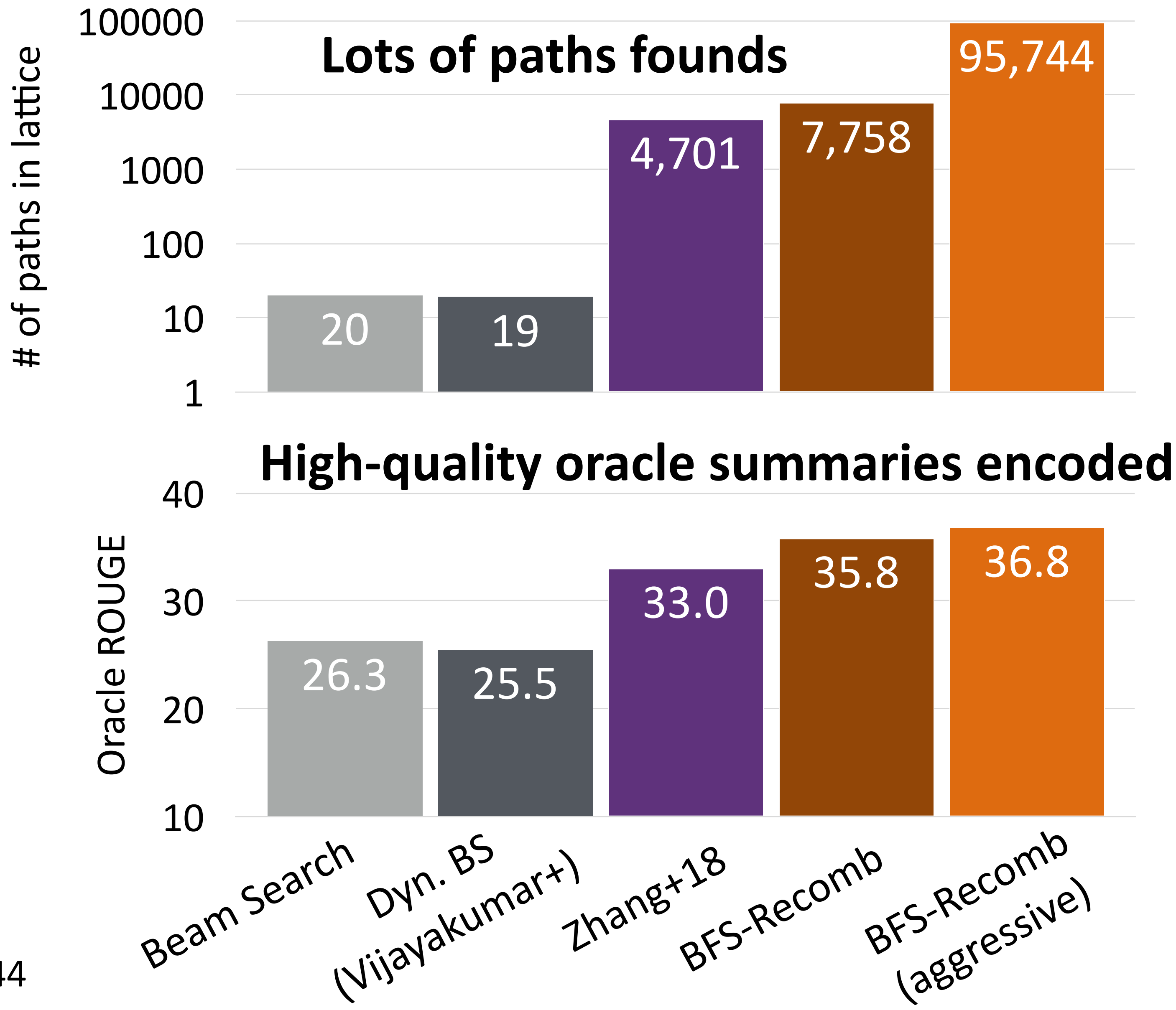


Results



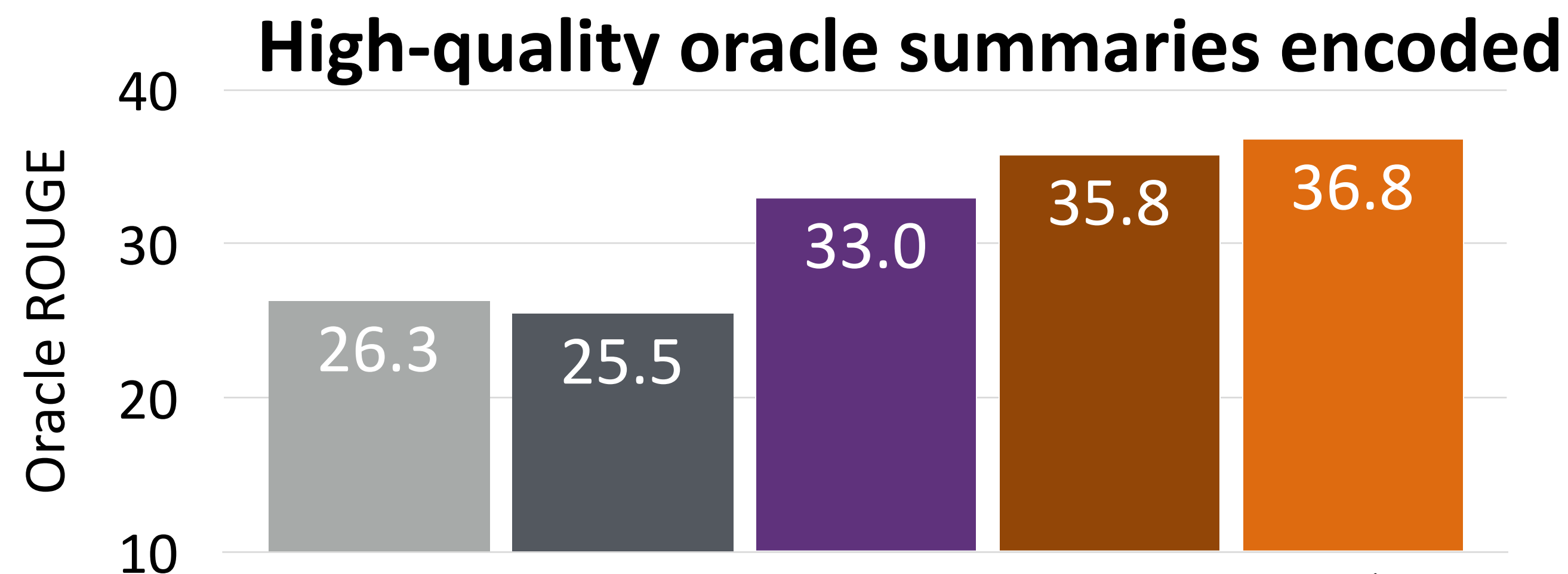
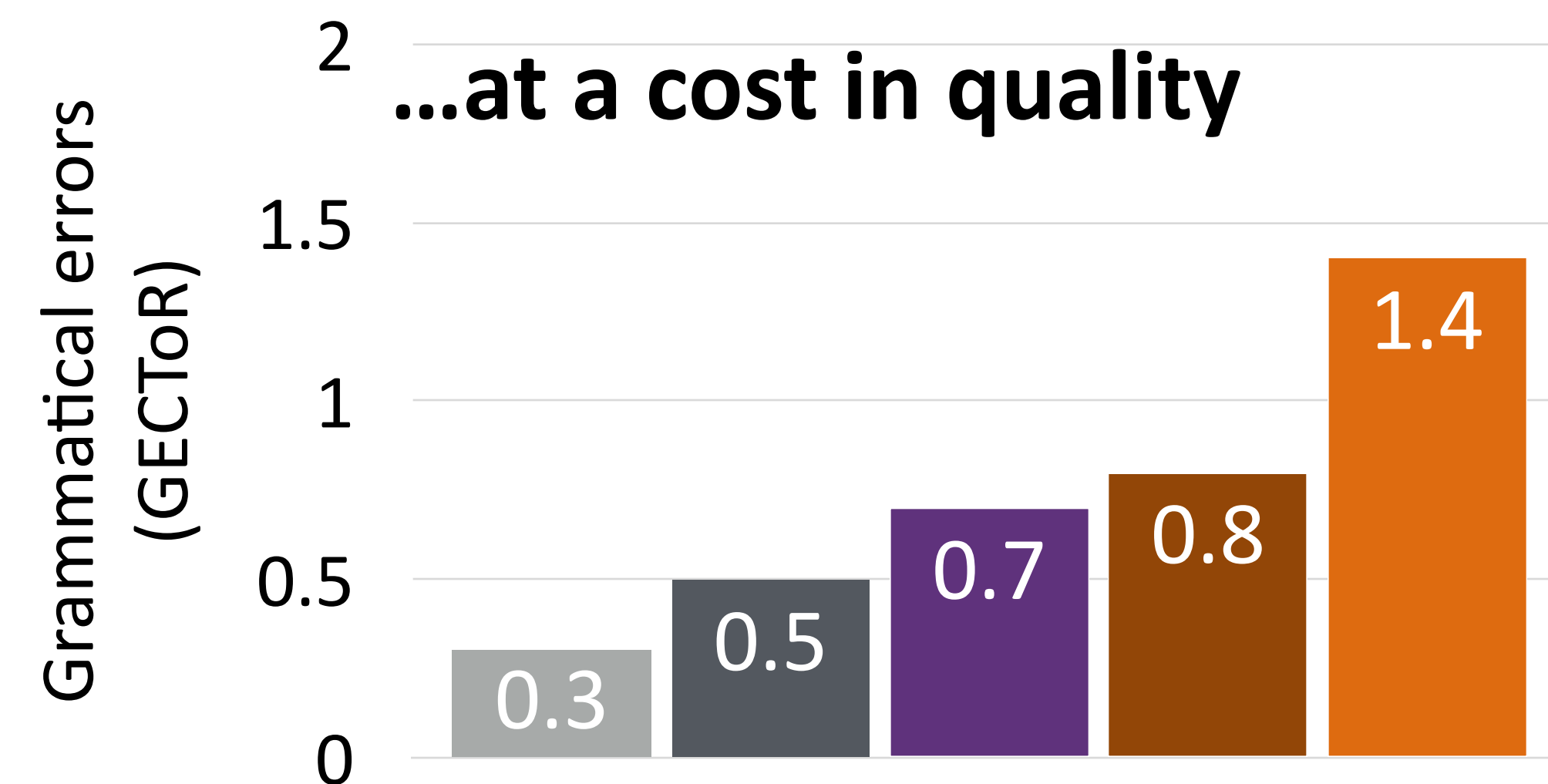
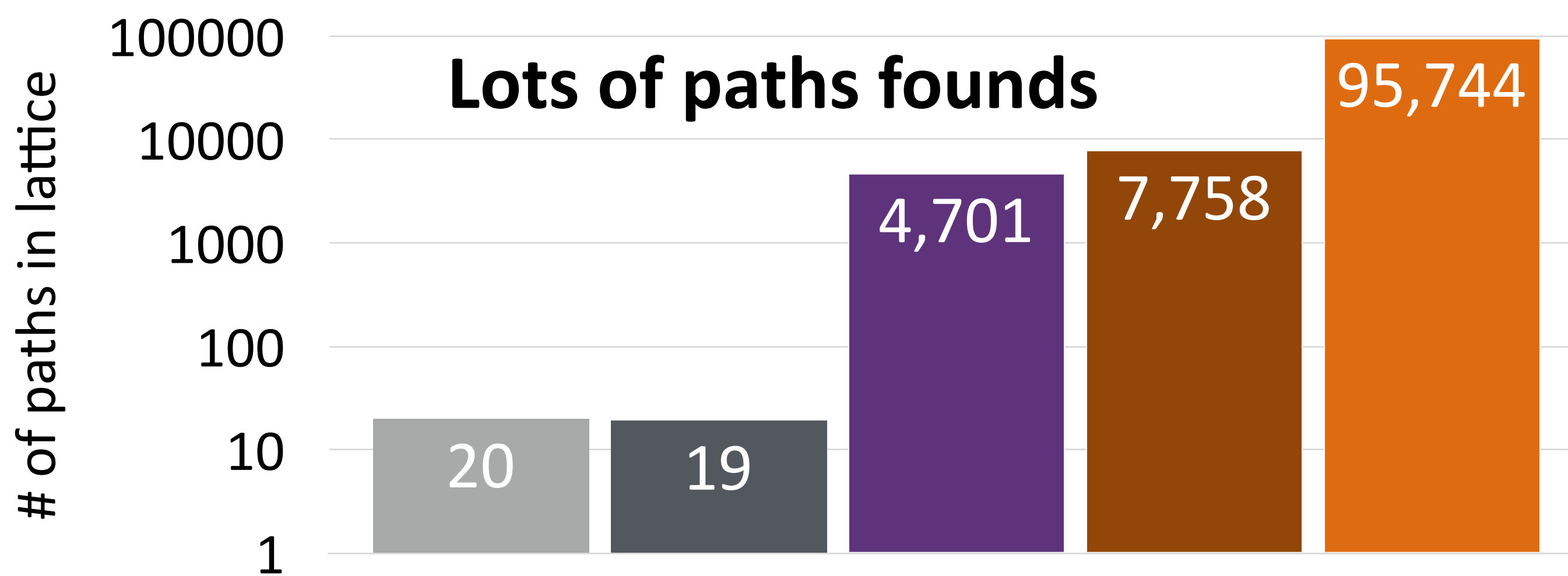


Results





Results



- ▶ Our most aggressive merging does introduce some grammatical errors. Better merging heuristics could help with this.



Goals for Lattices

Our generation systems can already encode lots of good options. We just need to be able to efficiently find them and encode them!

- ▶ Can we rerank our generated deductions and pick out the good ones?
- ▶ Can users control + correct the system on-the-fly, with the system learning those corrections?

Applications: factuality, controllable dialogue, diverse paraphrasing, and more!

Enumerating all valid deductions + a strong proof engine = effective search over natural language proofs using present-day models?



Outline

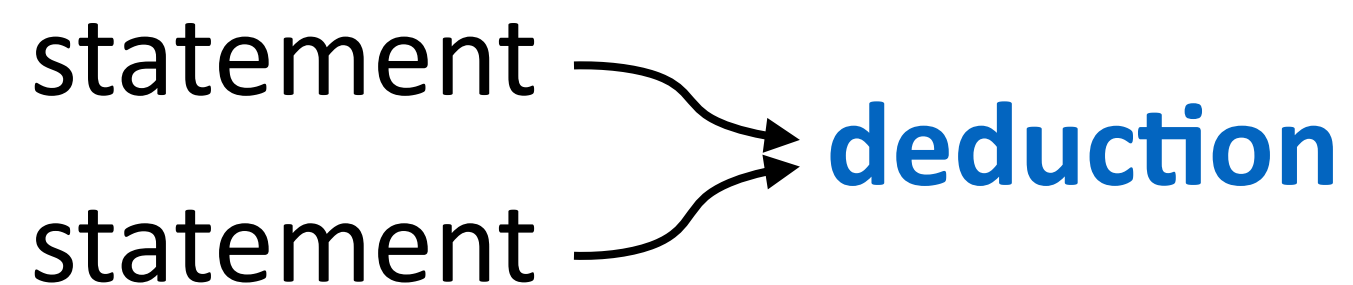
Entailment to verify QA

Jifan Chen, Eunsol Choi, GD. EMNLP-Findings21
Can NLI Models Verify QA Systems' Predictions?



Logically manipulating statements

Kaj Bostrom, Xinyu Zhao, Swarat Chaudhuri, GD. EMNLP21
Flexible Generation of Natural Language Deductions



Kaj Bostrom, Zayne Sprague, Swarat Chaudhuri, GD. In submission.
Natural Language Deduction through Search over Statement Compositions

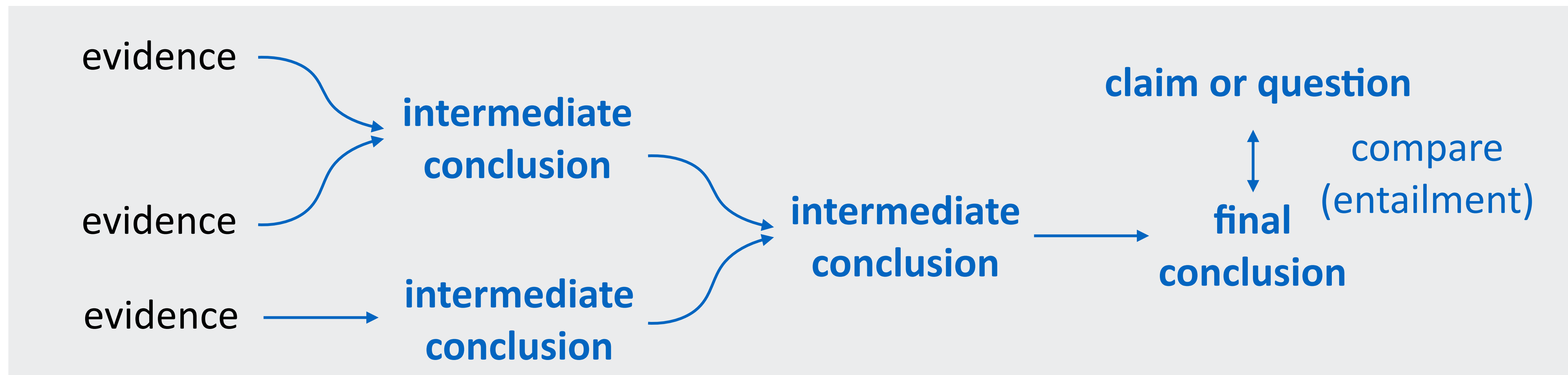
Improving diverse generation

Jiacheng Xu, GD. NAACL22.
Massive-scale Decoding for Text Generation using Lattices





Path to Multi-step Reasoning



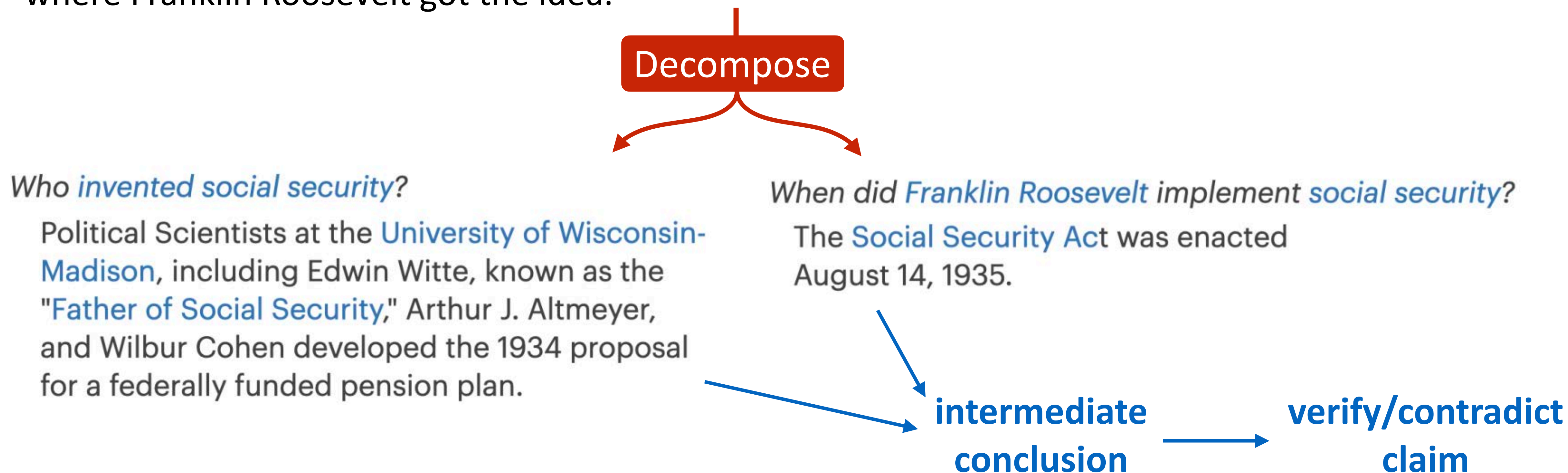
With better generation techniques and advances in pre-trained models, these steps will become easier and easier!



Goals

- Can we use it to do better, more explainable fact-checking?

Claim: Social Security was basically invented at the University of Wisconsin-Madison; that's where Franklin Roosevelt got the idea.



Angela Fan et al. *Generating Fact-Checking Briefs*, EMNLP 2020

Jifan Chen, Aniruddh Sriram, Eunsol Choi, GD. In preparation.



Goals

- Can we materialize reasoning about entities?

Claim: Harry Potter can teach classes on how to fly on a broomstick.



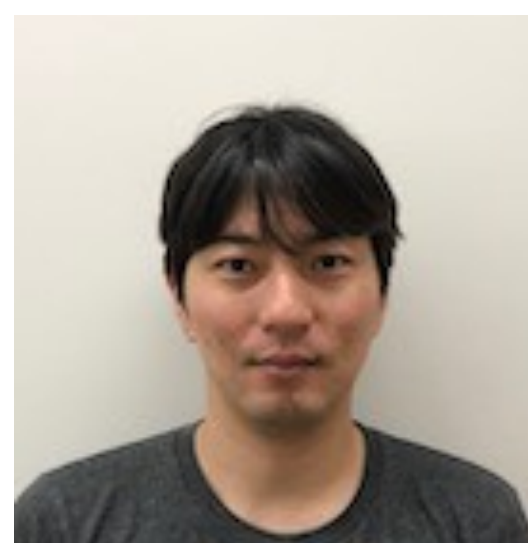
Harry Potter is a wizard ...
He plays Quidditch while riding
on a broomstick.

+



Someone who's good at
something can teach it.

Large annotated dataset (13k total claims about entities). Can textual reasoning help materialize an explanation?



Yasumasa Onoe, Michael J.Q. Zhang, Eunsol Choi, GD. *NeurIPS Datasets 2021*
CREAK: A Dataset for Commonsense Reasoning over Entity Knowledge



Conclusion

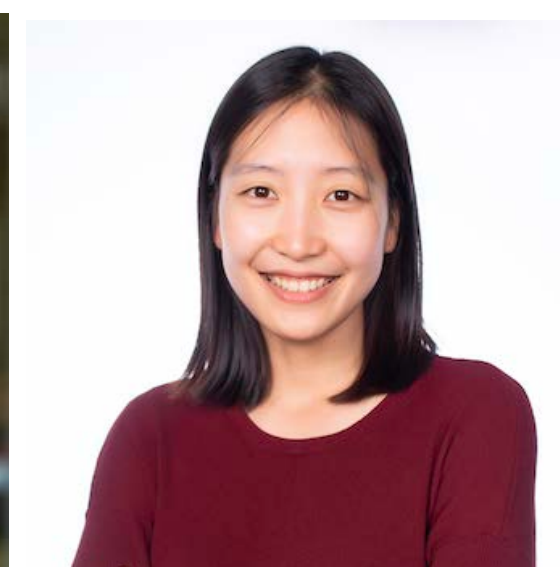
- ▶ Complex reasoning problems can be tackled using natural language to represent intermediate conclusions
- ▶ Natural language is an **expressive, flexible, and interpretable** vehicle for reasoning
- ▶ New datasets and better models are dramatically improving our ability to manipulate text (PaLM). **Making logical inferences in text is increasingly becoming viable.**



Acknowledgments



Bloomberg
arm



amazon

Walmart 

Thanks!



Results: EntailmentBank Human Eval

