

Can NLI Models Verify QA Systems' Predictions?





Jifan Chen, Eunsol Choi and Greg Durrett

The University of Texas at Austin



Unreliable Predictions of QA systems

- ▶ Current QA systems try to return a most-plausible answer to users:


who invented the first central processing unit (cpu) ×  

[All](#) [Images](#) [News](#) [Shopping](#) [Videos](#) [More](#) [Tools](#)

About 880,000 results (1.08 seconds)

physicist Federico Faggin

Italian physicist Federico Faggin invented the first commercial CPU. It was the Intel 4004 released by Intel in 1971.



[https://stu](https://study.com/who-invented-the-first-cpu.html)
[Who in](https://study.com/who-invented-the-first-cpu.html) [ly.com](https://study.com)

First commercial CPU != First CPU

[?](#) About featured snippets • [Feedback](#)



Unreliable Predictions of QA systems

- ▶ Current QA systems try to return a most-plausible answer to users:

who invented the first centr

Google

When was Vancouver's last earthquake?

All News Images Maps Shopping More Tools

About 23,200,000 results (0.68 seconds)

1946

The **1946** Vancouver Island earthquake struck Vancouver Island on the coast of British Columbia, Canada, on June 23 at 10:15 a.m. with a magnitude estimated at 7.0 M_s and 7.5 M_w .

...

1946 Vancouver Island earthquake.

UTC time	1946-06-23 17:13:24
Local time	10:15 a.m.
Magnitude	7.0 M_s 7.5 M_w
Depth	15 km (9.3 mi)
Epicenter	49.62°N 125.26°W

8 more rows

No information about "last" is given



Unreliable Predictions of QA systems

- ▶ Current QA systems try to return a most-plausible answer to users:

who invented the first central processing unit (cpu) | **Google** | When was Vancouver's last earthquake?

All News Images Maps Shopping More

About 23,200,000 results (0.68 seconds)

About 880,000 results (1.08 seconds)

1946

The **1946** Vancouver Island earthquake struck Vancouver Island on the coast of British Columbia, Canada, on June 23 at 10:15 a.m. with a magnitude estimated at 7.0 M_s and M_w .

Italian physicist **Federico Faggin** invented the first commercial CPU. It was the Intel 4004 released by Intel in 1971.

UTC time **1946-06-23 17:13:24**

Local time **10:15 a.m.**

Magnitude **7.0 M_s 7.5 M_w**

Depth **15 km (9.3 mi)**

Epicenter **49.62°N 125.26°W**

8 more rows

2011-09-09	Vancouver Island
2010-06-23	Central Canada
2009-11-17	Queen Charlotte Islands, BC
2009-07-07	Baffin Bay
2008-01-05	Queen Charlotte Islands, BC
2007-10-09	The Nazko region
2004-11-02	Vancouver Island, BC
2004-07-19	Vancouver Island



NLI as a verifier

- ▶ Idea: Use Natural Language Inference (NLI) to verify whether an answer can be **truly** entailed from its corresponding context (Harabagiu and Hickl, 2006; Peñas et al.2008; Yin et al. 2021; Mishra et al. 2021).

Question: What is the revolution period of Venus in earth days?
Answer: 243 days

Question conversion ▼

Hypothesis: The **revolution** period of Venus in earth days is 243 days.

Context: **Venus** is the second planet from the Sun... **It** has the longest rotation period (243 days)...

▼ **Decontextualization**

Premise: Venus has the longest **rotation** period (243 days)...

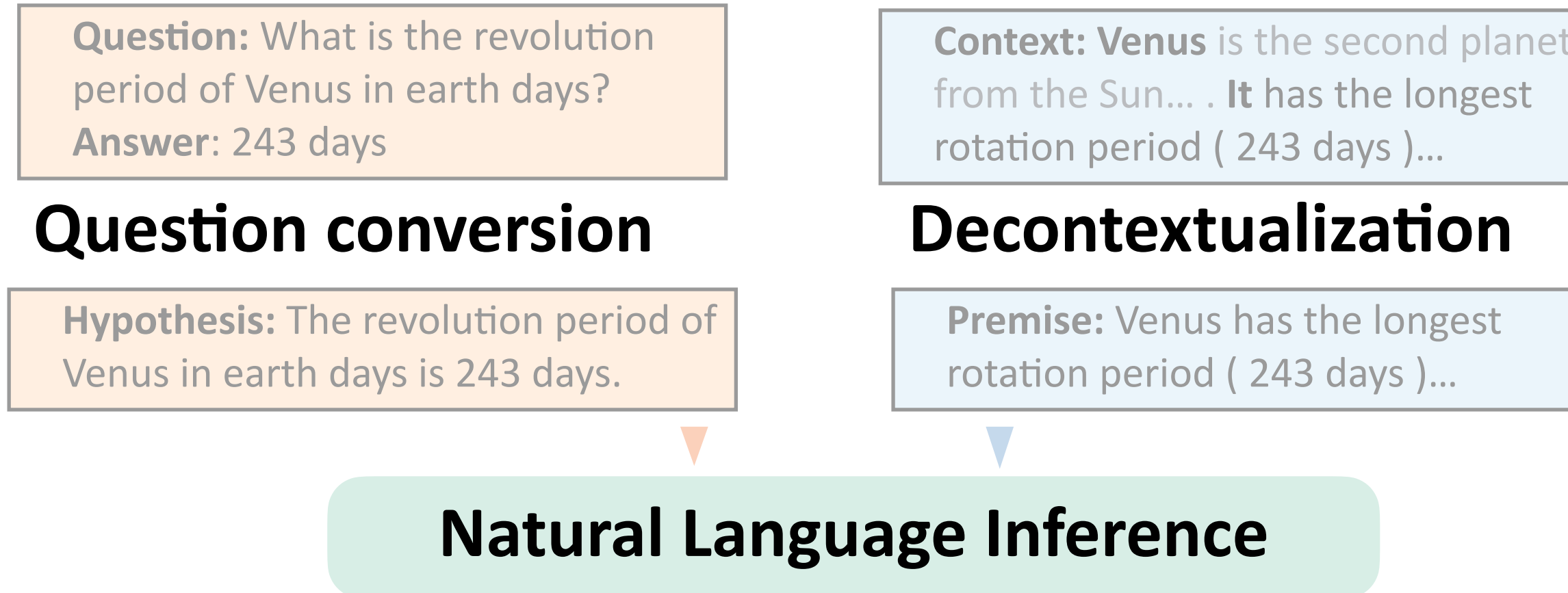
▼ **Natural Language Inference**

▶ **Not entailed**



Outline

1) NLI as a QA Verifier



2) Experiments

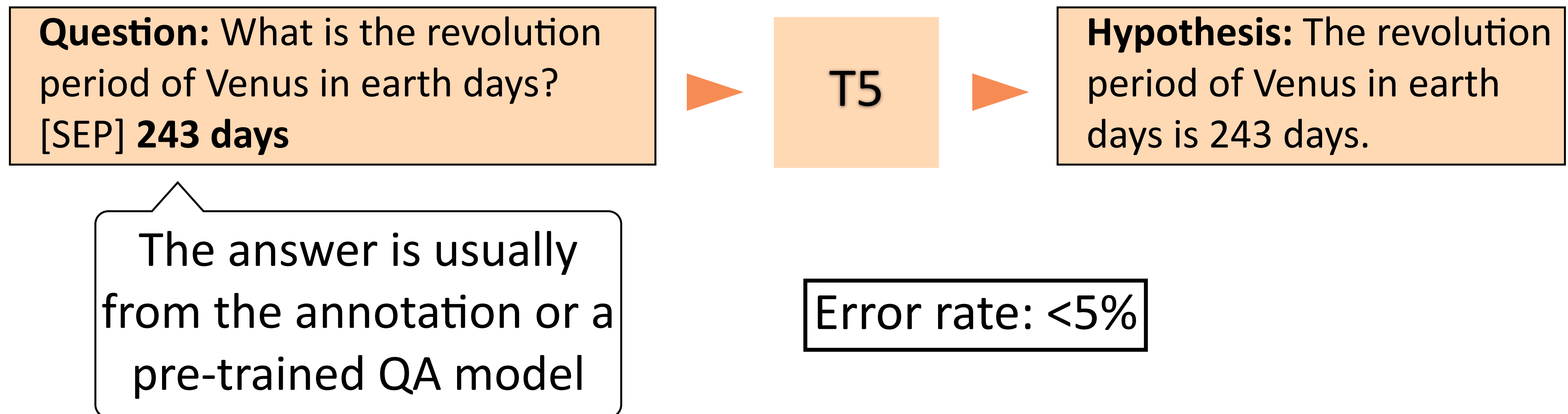
- ▶ Rejecting unanswerable questions
- ▶ Improving the prediction confidence
- ▶ Disagreement between NLI and QA

3) Takeaways



Question Conversion

- ▶ Demszky et al. (2018) explored a rule-based conversion system.
- ▶ This work: A T5-based model converting a (question, answer) pair to a statement trained on the annotations from Demszky et al. (2018).





Decontextualization

- ▶ The whole paragraph contains too much information
- ▶ A T5-based model to rewrite (name completion, NP/pronoun swap, bridging) a sentence to be interpretable out of context if feasible (Choi et al. 2021).

Error rate: <10%

Venus is the second planet from the Sun , orbiting it every 224.7 Earth days . **It** has the longest rotation period (243 days) of any planet in the Solar System and rotates in the opposite direction to most other planets ...

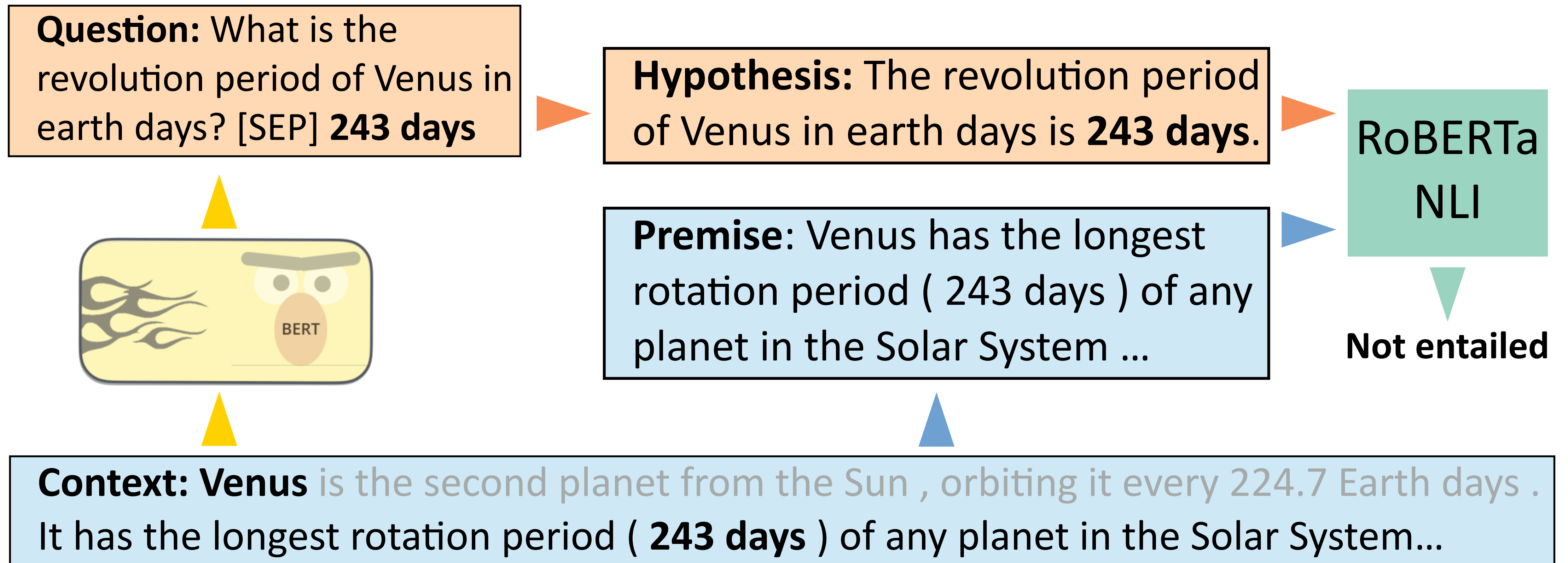
T5

Venus has the longest rotation period (243 days) of any planet in the Solar System and rotates in the opposite direction to most other planets ...



NLI model

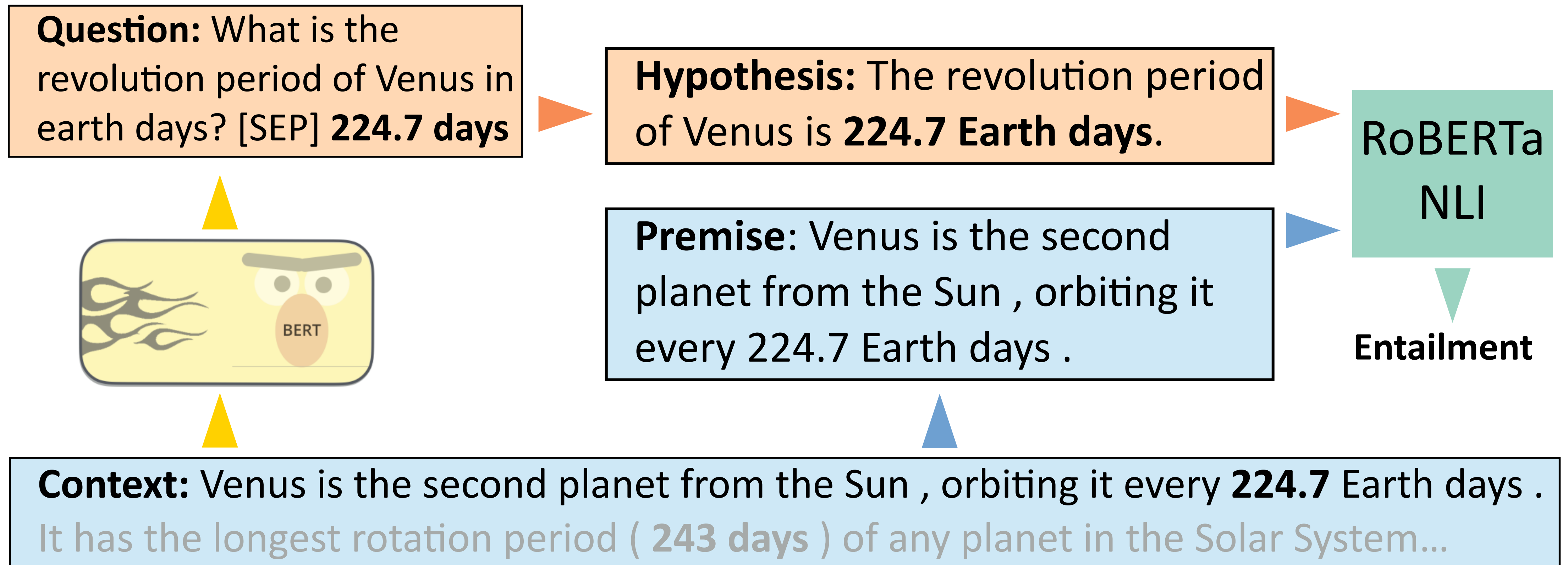
- ▶ A Roberta-based model trained with the (premise, hypothesis) pairs from QA datasets: gold answer + context as positive, non-gold answers in the top-k predictions as negatives





NLI model

- ▶ A Roberta-based model trained with the (premise, hypothesis) pairs from QA datasets: gold answer + context as positive, non-gold answers in the top-k predictions as negatives





Outline

1) NLI as a QA Verifier

Question: What is the revolution period of Venus in earth days?
Answer: 243 days

Question conversion

Hypothesis: The revolution period of Venus in earth days is 243 days.

Context: Venus is the second planet from the Sun... . It has the longest rotation period (243 days)...

Decontextualization

Premise: Venus has the longest rotation period (243 days)...

Natural Language Inference

2) Experiments

- ▶ Rejecting unanswerable questions
- ▶ Improving the prediction confidence
- ▶ Disagreement between NLI and QA

3) Takeaways



Rejecting unanswerable questions

- ▶ Experimental setup:
 - Train a QA model on SQuAD1.1 (every question is answerable) and test it on SQuAD2.0 (contains unanswerable questions).
 - Using an NLI model pre-trained on MNLI to verify the predictions from the SQuAD1.1 model (always gives an answer).
- ▶ Results: MNLI model successfully rejects **78.5%** of the unanswerable examples and accepts **82.5%** of the answerable examples.
 - Notice that MNLI is totally out-of-domain regarding the task of QA



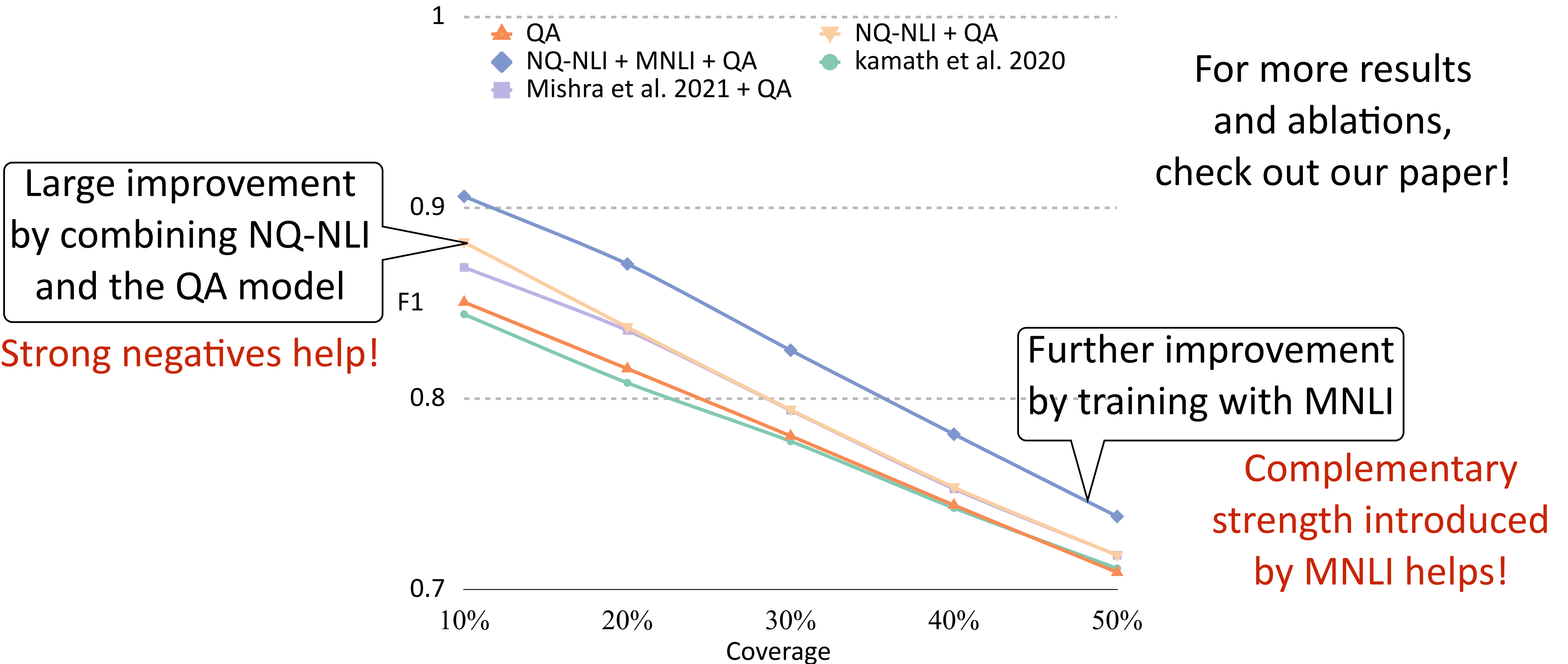
Improving the prediction confidence

- ▶ Selective QA: If our model can choose to answer only the k percentage of examples it is most confident about (the coverage), what F1 can it achieve?
 - Examples are ranked by the confidence score of a model
 - QA -> posterior probability of the answer span
 - NLI -> posterior probability associated with “Entailment”
- ▶ Experimental setup:
 - Train a QA model on Natural Questions (Kwiatkowski et al. 2019) and test it on NQ and 4 out-of-domain datasets.
 - Train a NLI model using the generated NLI pairs from Natural Questions and use it to verify the predictions from the previous step



Improving the prediction confidence

► Results:





Disagreement between NLI and QA

- ▶ Right answer for the wrong reason (25% in NQ):

Question: When was Clash Royale released in the US?

Answer: The game Clash Royale was released globally on **March 2 , 2016**

Globally —> US? Need further evidence

Question: Who plays the bad guy in the Good Place?

Answer: 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**) to The Good Place " ...

Is Michael the bad guy? Need to check



Takeaways

- ▶ Existing QA datasets encourage models to return answers when the context does not actually contain sufficient information
 - Fully verifying the answer is a challenging direction
- ▶ The proposed approach is helpful to locate the information mismatches between the question and the supporting context

Thank you!

Q&A