# Entity-Based Knowledge Conflicts in Question Answering
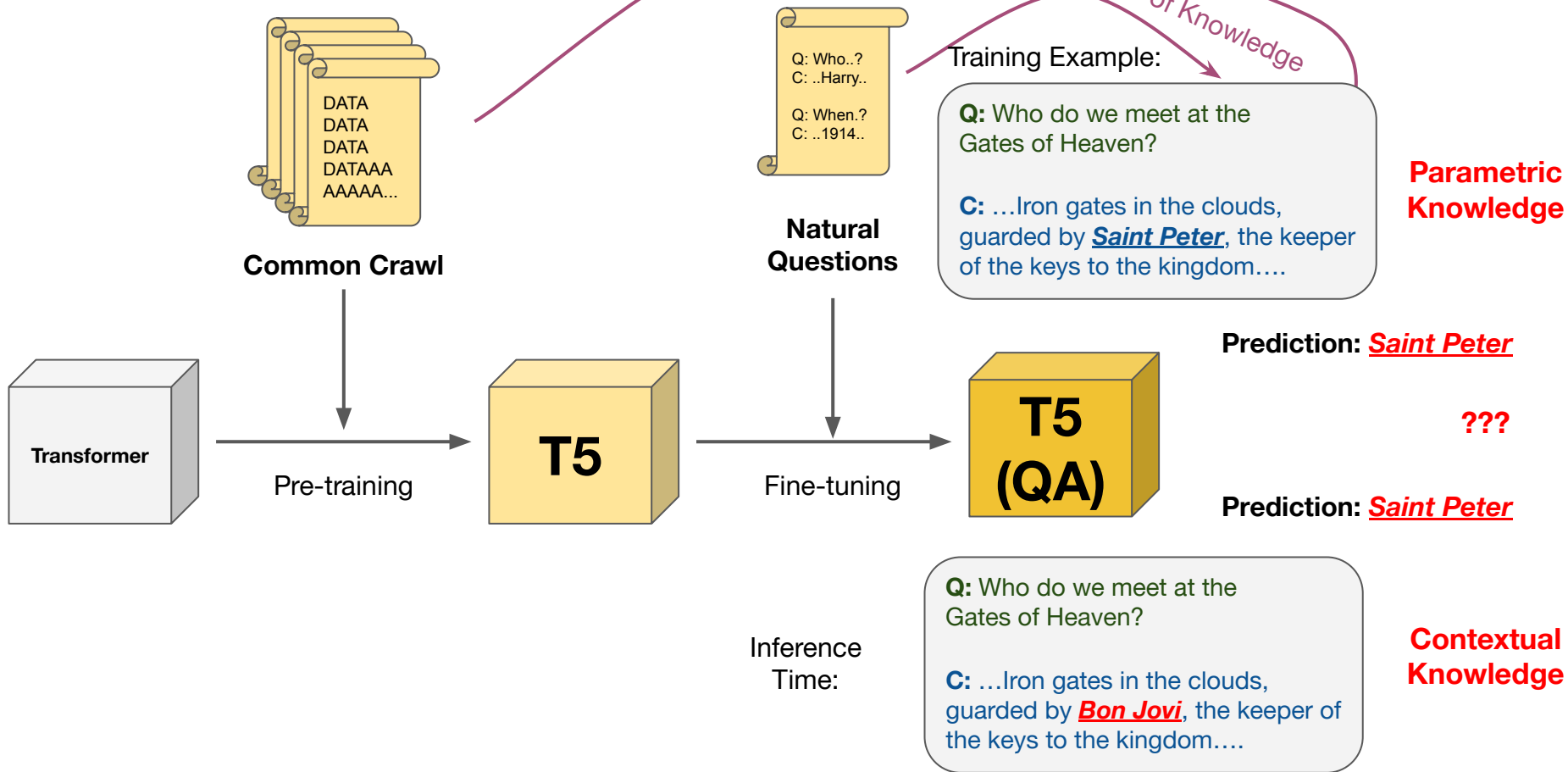
Shayne Longpre*, Kartik Perisetla*, Anthony Chen*, Nikhil Ramesh, Chris DuBois, Sameer Singh

(* = Equal contribution)

What is a <span style="color:red">contextual-parametric *knowledge conflict*</span>?

**Knowledge Conflicts**

Common Crawl

DATA
DATA
DATA
DATAAA
AAAAA...

Natural
Questions

Q: Who..?
C: ..Harry..

Q: When.?
C: ..1914..

Transformer

Pre-training

T5

Fine-tuning

T5
(QA)

Source of Knowledge

Training Example:

**Q:** Who do we meet at the Gates of Heaven?

**C:** …Iron gates in the clouds, guarded by *Saint Peter*, the keeper of the keys to the kingdom….

**Parametric Knowledge**

**Prediction:** *Saint Peter*

**???**

**Prediction:** *Saint Peter*

Inference Time:

**Q:** Who do we meet at the Gates of Heaven?

**C:** …Iron gates in the clouds, guarded by *Bon Jovi*, the keeper of the keys to the kingdom….

**Contextual Knowledge**

# Why do we care if models *ignore* the context?

**Why do we care _which_ _knowledge_ models use?**

1. *Static* knowledge v. *Temporal* knowledge → Generalization

2. *Interpretability* of a prediction

3. Context _grounding_ mitigates **hallucination**, **bias**, **stochastic parroting**

**Summary of Findings**

1. QA Dataset → **Substitution Framework** → Knowledge Conflicts

2. **Benchmark behaviour** (parametric vs contextual) → lots of **hallucination**!!!

3. **Factors**: (1) model size, (2) quality of retriever at training, (3) popularity of entities

4. Mitigate this behaviour → improves **generalization**.

# Substitution Framework

**Original Example**

**Q:** Who do we meet at the Gates of Heaven?

**C:** ...Iron gates in the clouds, guarded by **_Saint Peter_**, the keeper of the keys to the kingdom….

**Alias Substitution**

**C:** ...Iron gates in the clouds, guarded by **_Peter the Apostle_**, the keeper of the keys to the kingdom….

<> Policy: *Wikidata alias of original answer* <>

**Corpus Substitution**

**C:** ...Iron gates in the clouds, guarded by **_Bon Jovi_**, the keeper of the keys to the kingdom….

<> Policy: *Sample* PERSON *from training set* <>

Types: [*PER, LOC, ORG, DAT, NUM*]

**Human Assessment**

# 98%

**Fluency** of original Natural Questions examples
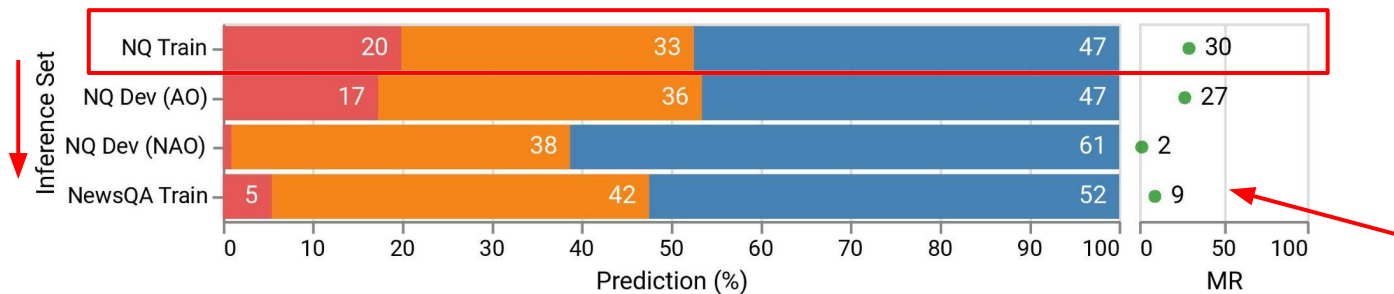
# 84%

**Fluency** of *Corpus Substituted* Natural Questions examples

**Model Behaviour**

$$p_o \qquad\qquad p_s \qquad\qquad p_t$$

Predict **<u>Original</u>** Answer     Predict **<u>Substitute</u>** Answer     Predict **<u>Other</u>** Answer

$$M_R = \frac{p_o}{p_o + p_s}$$

**"Memorization Ratio"**

# Model Behaviour

**Train: Natural Questions**
**Test: Corpus Substitution**



**Train: NewsQA**
**Test: Corpus Substitution**

# Model Behaviour

**Train: NQ**
**Test: Alias Substitution**



**_Extractive_ QA Model**
**Train: NQ**
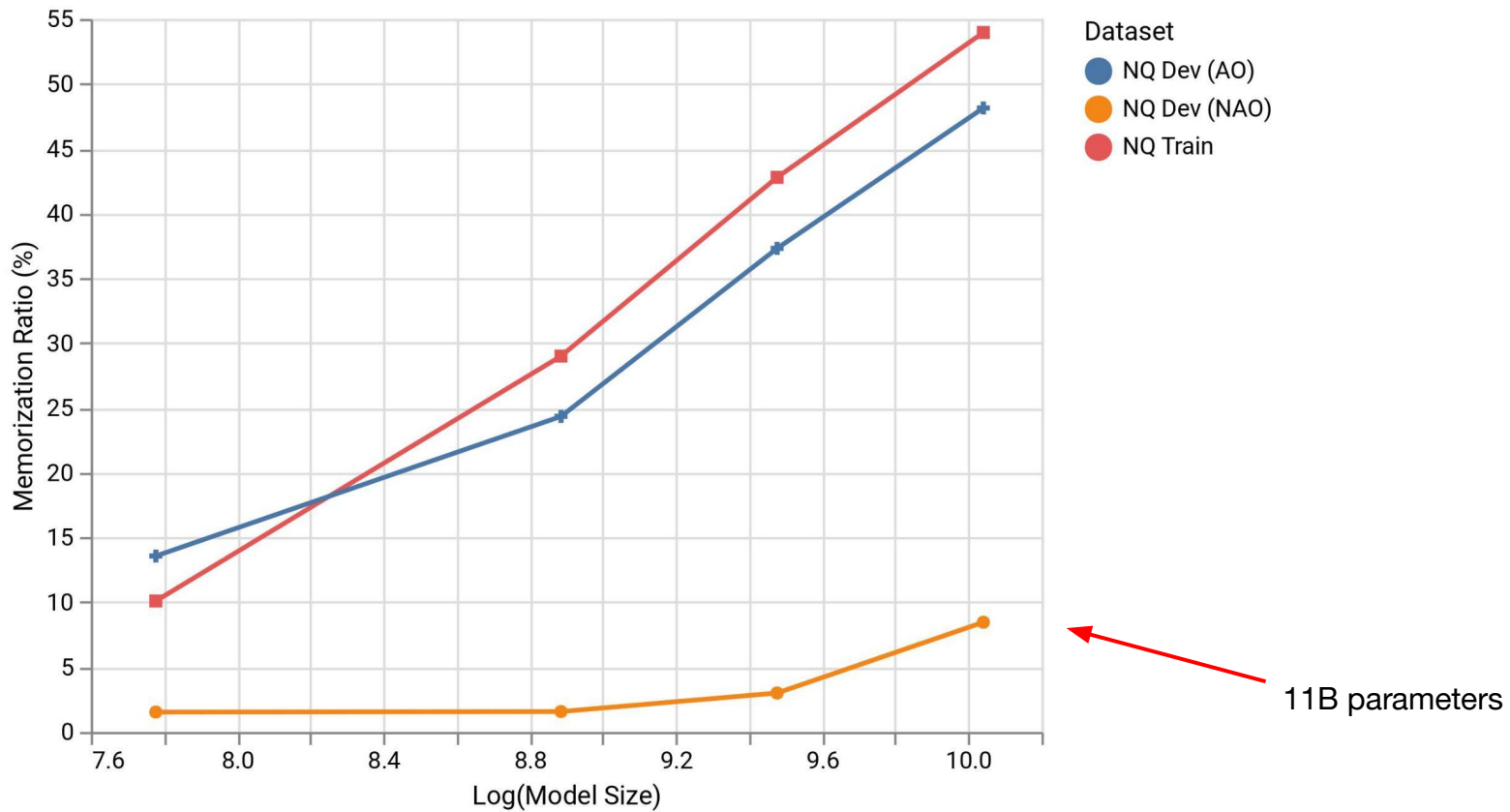**Test: Corpus Substitution**

Takeaway:

Parametric preference over context is prevalent, and contradictions cause confusion/instability in predictions.

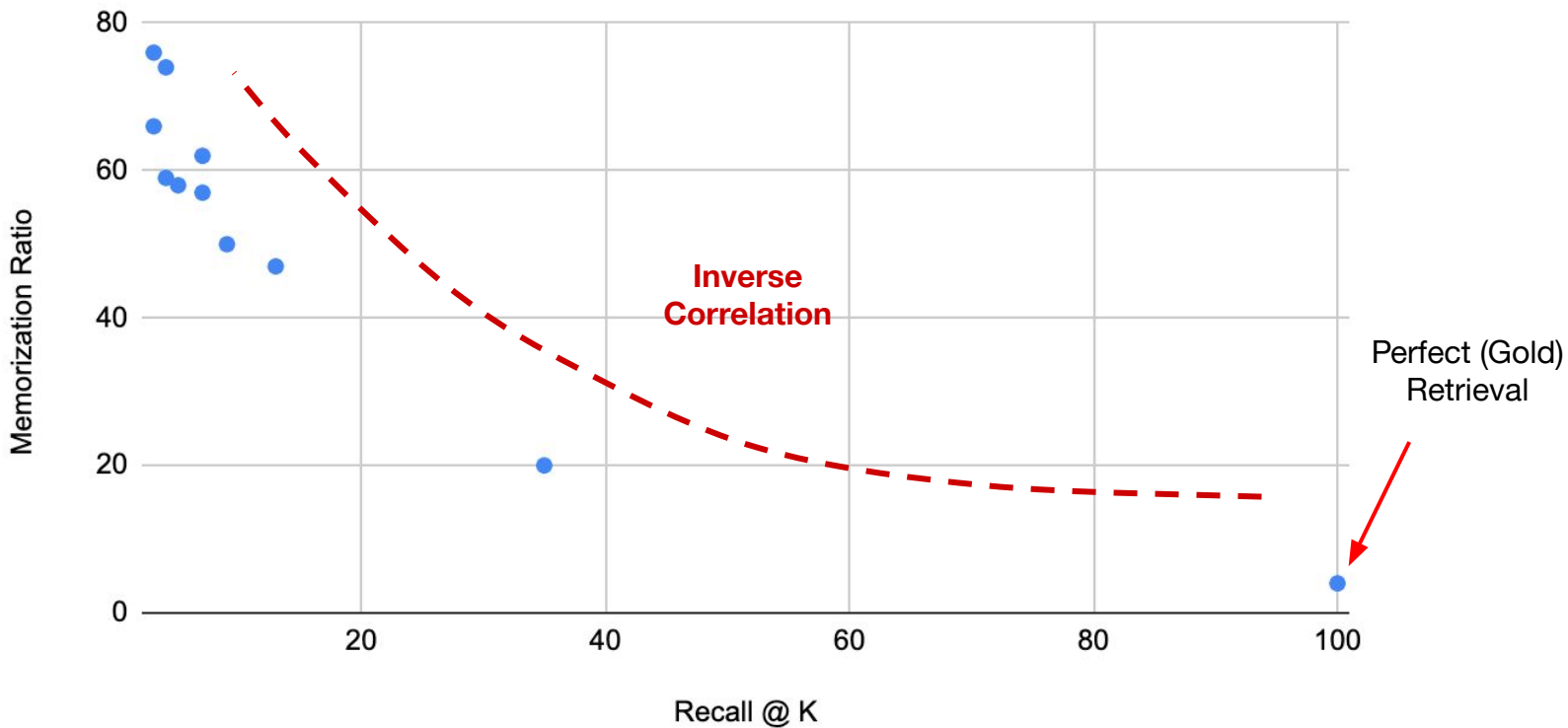# What *Factors* affect the Memorization Ratio?

# Factor 1: Model Size

Takeaway:

As model capacity grows, it relies more heavily on memorized information (even from pre-training).

# Factor 2: Retriever Quality during Training



Recall@K vs. Memorization Ratio

Takeaways:

Reader models ignore context when retrievers are poor.

Only trust context when retrievers are near-perfect.
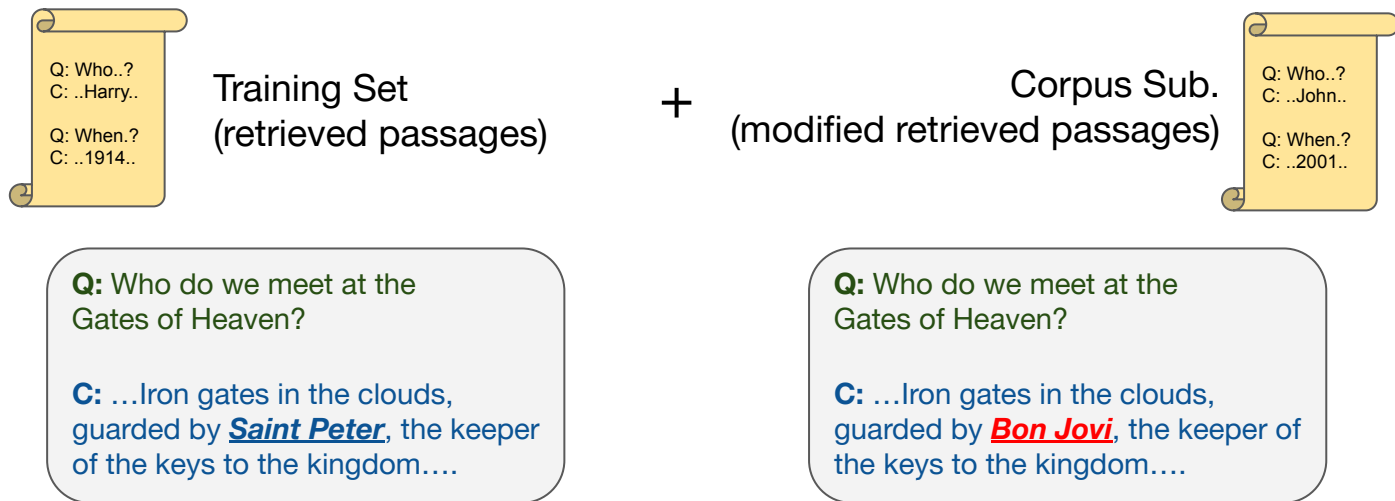
# How can we *mitigate* memorization/hallucination?

# Mitigating Memorization

**Recall Key Insight:** training with perfect retrieval → low Memorization Ratio

But…
- We don't have unlimited gold passage labels to train on
- And SOTA QA models need to train on the same retriever they will use at test time…

**Solution:** Train on: [1] (fallible) model-retrieved passages + [2] corpus-substitution version.

Q: Who..?
C: ..Harry..

Q: When.?
C: ..1914..

Training Set
(retrieved passages)

+

Corpus Sub.
(modified retrieved passages)

Q: Who..?
C: ..John..

Q: When.?
C: ..2001..

**Q:** Who do we meet at the Gates of Heaven?

**C:** …Iron gates in the clouds, guarded by ***Saint Peter***, the keeper of the keys to the kingdom….

**Q:** Who do we meet at the Gates of Heaven?

**C:** …Iron gates in the clouds, guarded by ***Bon Jovi***, the keeper of the keys to the kingdom….

**Mitigating Memorization**

| Inference Set | $M_R$ | $EM$ ($\Delta$) |
|---|---|---|
| NQ TRAIN | $29.5 \rightarrow 2.6$ | $70.9 \rightarrow 64.9$ (-5.0) |
| NQ DEV (AO) | $27.1 \rightarrow 1.9$ | $62.7 \rightarrow 64.2$ (+1.5) |
| NQ DEV (NAO) | $1.5 \rightarrow 0.0$ | $32.9 \rightarrow 40.0$ (+7.1) |
| NEWSQA | $9.3 \rightarrow 0.6$ | $21.4 \rightarrow 25.8$ (+4.4) |

# Thank you!

Please don't hesitate to reach out!
- Email: slongpre@mit.edu
- Repository: https://github.com/apple/ml-knowledge-conflicts
- Paper: https://arxiv.org/abs/2109.05052