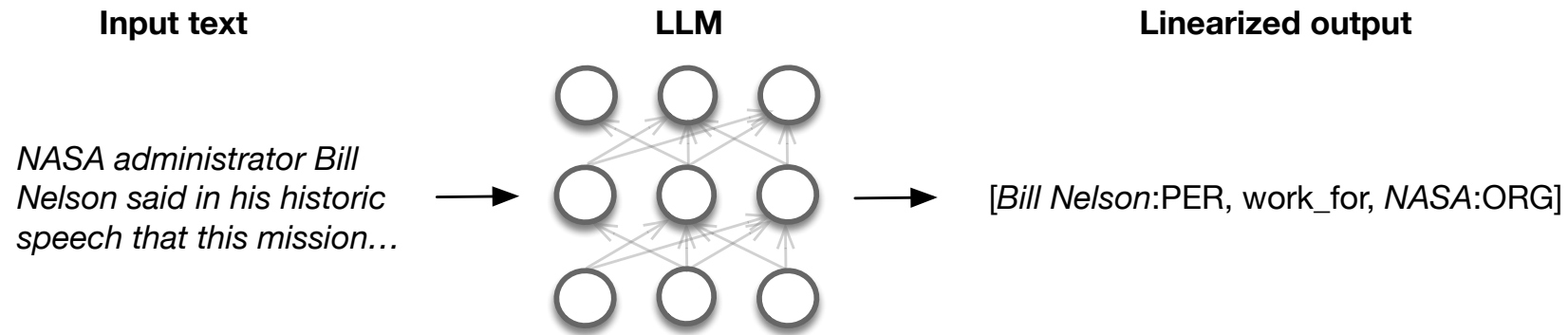


Evaluating AI systems

Evaluating AI systems (is hard)

Relation extraction with LLMs



Revisiting Relation Extraction in the era of Large Language Models

Somin Wadhwa Silvio Amir Byron C. Wallace
Northeastern University,
{wadhwa.s,s.amir,b.wallace}@northeastern.edu



PICO extraction + evidence inference w/LLMs

Input

*Patients receiving **aspirin** experienced **headaches** with comparable **duration** but significantly lower **reported pain** compared to those receiving **placebo**.*

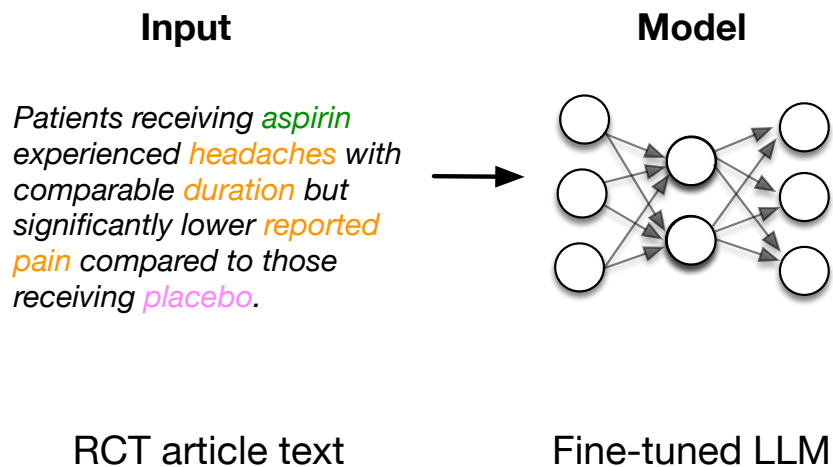
RCT article text

Jointly Extracting Interventions, Outcomes, and Findings from RCT Reports with LLMs

Somin Wadhwa, Jay DeYoung, Benjamin Nye, Silvio Amir, Byron C. Wallace



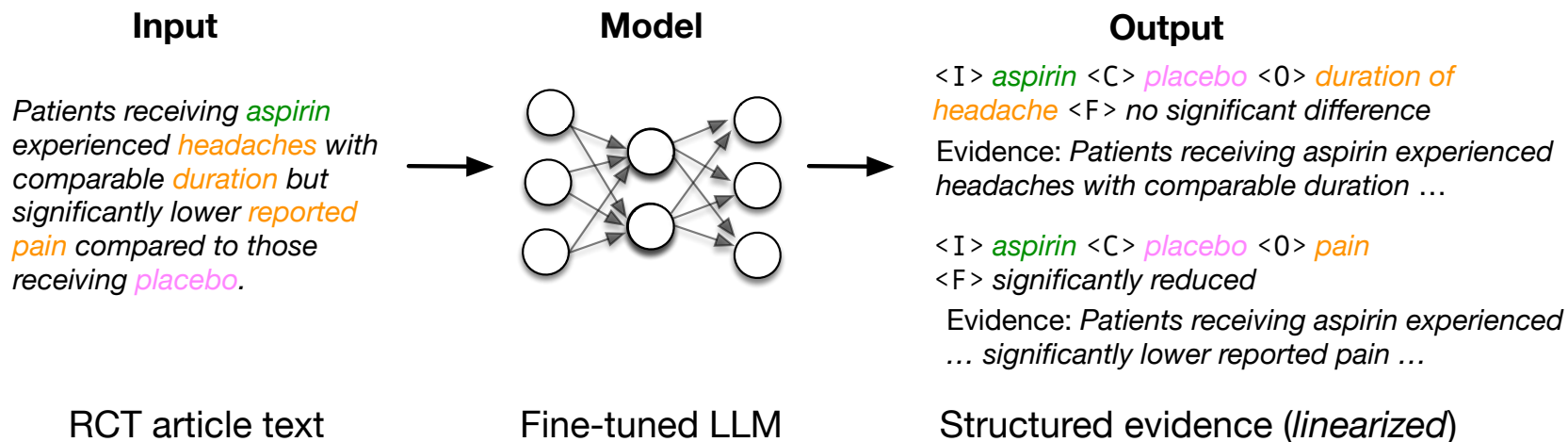
PICO extraction + evidence inference w/LLMs



Jointly Extracting Interventions, Outcomes, and Findings from RCT Reports with LLMs

Somin Wadhwa, Jay DeYoung, Benjamin Nye, Silvio Amir, Byron C. Wallace

PICO extraction + evidence inference w/LLMs



Jointly Extracting Interventions, Outcomes, and Findings from RCT Reports with LLMs

Somin Wadhwa, Jay DeYoung, Benjamin Nye, Silvio Amir, Byron C. Wallace

Evaluation is hard

Abstract snippet: Canagliflozin increased urinary glucose excretion in a dose-dependent manner and produced statistically significant reductions in body weight compared with placebo (least squares mean percent changes from baseline of -2.2%, -2.9%, -2.7%, and -1.3% with canagliflozin 50, 100, and 300 mg and placebo; $P < 0.05$ for all comparisons). Overall adverse event (AE) rates were similar across groups. Canagliflozin was associated with higher rates of genital mycotic infections in women, which were generally mild and led to few study discontinuations. Osmotic diuresis-related AE rates were low and similar across groups.

Reference: [canagliflozin, body weight, placebo, Canagliflozin increased urinary glucose excretion in a dose-dependent manner and produced statistically significant reductions in body weight compared with placebo., canagliflozin [LABEL] significantly decreased [OUT] body weight [COMP] placebo]

Generated: [canagliflozin, body weight reduction, placebo, Canagliflozin increased urinary glucose excretion in a dose-dependent manner and produced statistically significant reductions in body weight compared with placebo., canagliflozin [LABEL] significantly increased [OUT] body weight reduction [COMP] placebo]

Evaluating summaries



COVID-19 [population] × Chloroquine [interventions] × Start typing a Population, Intervention, Comparator, or Outcome (PICO)

Showing 13 results

Automatically generated summary (β!): There is currently insufficient evidence to support the routine use of HCQ for the prophylaxis of SARS-CoV-2 infection in healthcare personnel. Further randomised controlled trials are needed to determine whether HCQ is of benefit to healthcare personnel and their carers, and to compare HCQ with other antiviral therapies.

All (168) Published articles (13) Preprints (6) Registered trials (149)

Get large/high quality trials first

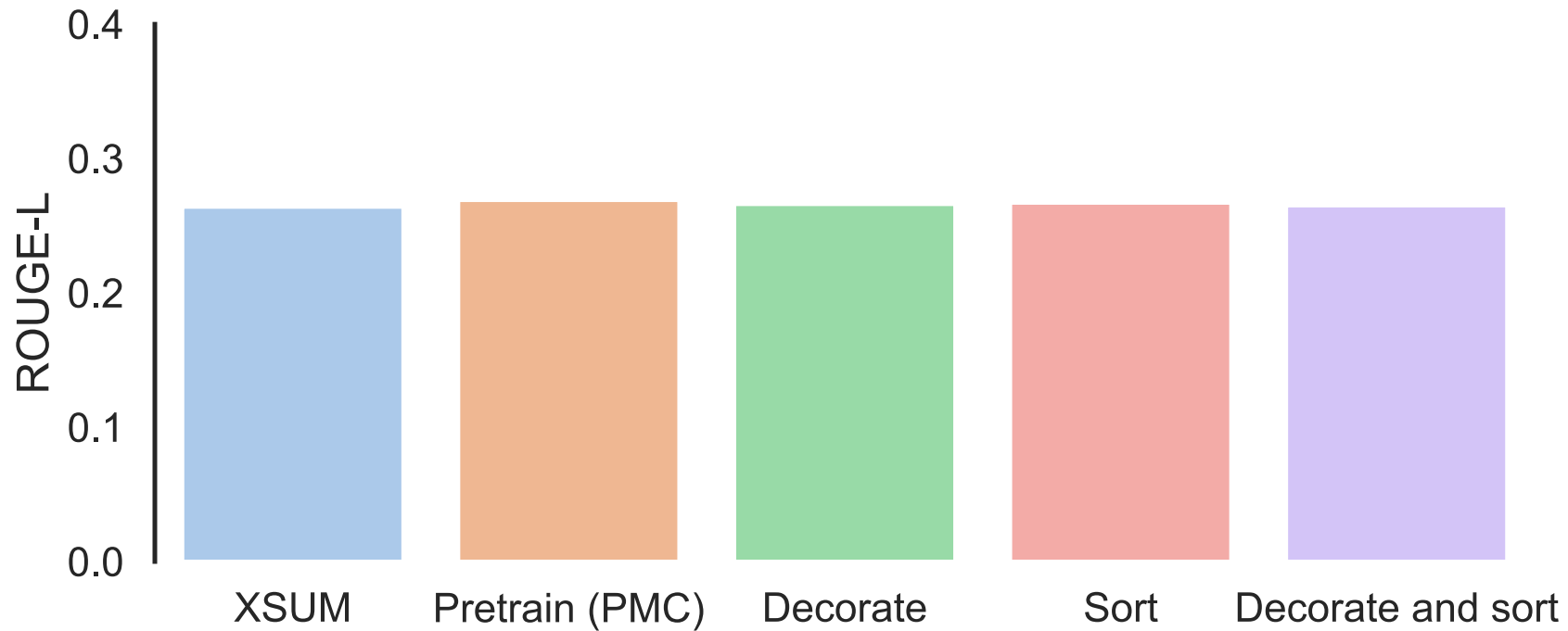
Newest first



In the olden days (~5 years ago)

Vaccines against SARS-CoV-2 cannot be recommended for routine clinical practice at this time. There is a need for well-designed RCTs with long-term follow-up to evaluate the efficacy and safety of vaccines against this disease in healthy adults ...

ROUGE (basically equivalent)



But this says nothing about how accurate these summaries are.

Annotate x +

127.0.0.1:5000

System Summary

Magnesium sulphate may reduce the incidence of eclampsia in women with mild to moderate preeclampsia. However, there is insufficient evidence to assess the effect of magnesium sulphate on other important outcomes, such as perinatal mortality and neurodevelopmental outcome. Further randomised controlled trials are needed to determine the optimal dose and route of administration, the optimal duration of prophylaxis, and the cost-effectiveness of this intervention.

Is the System Summary relevant to the topic *Magnesium sulphate and other anticonvulsants for women with pre-eclampsia*?

○ ○ ○

Mostly off topic (does not seem to address the key question). Moderately on topic, but contains seemingly irrelevant information as well. Strongly focusses on this topic.

Now we would now like you to assess the "semantic plausibility" of the text (without regard for the reference review or source abstracts). This text is ...

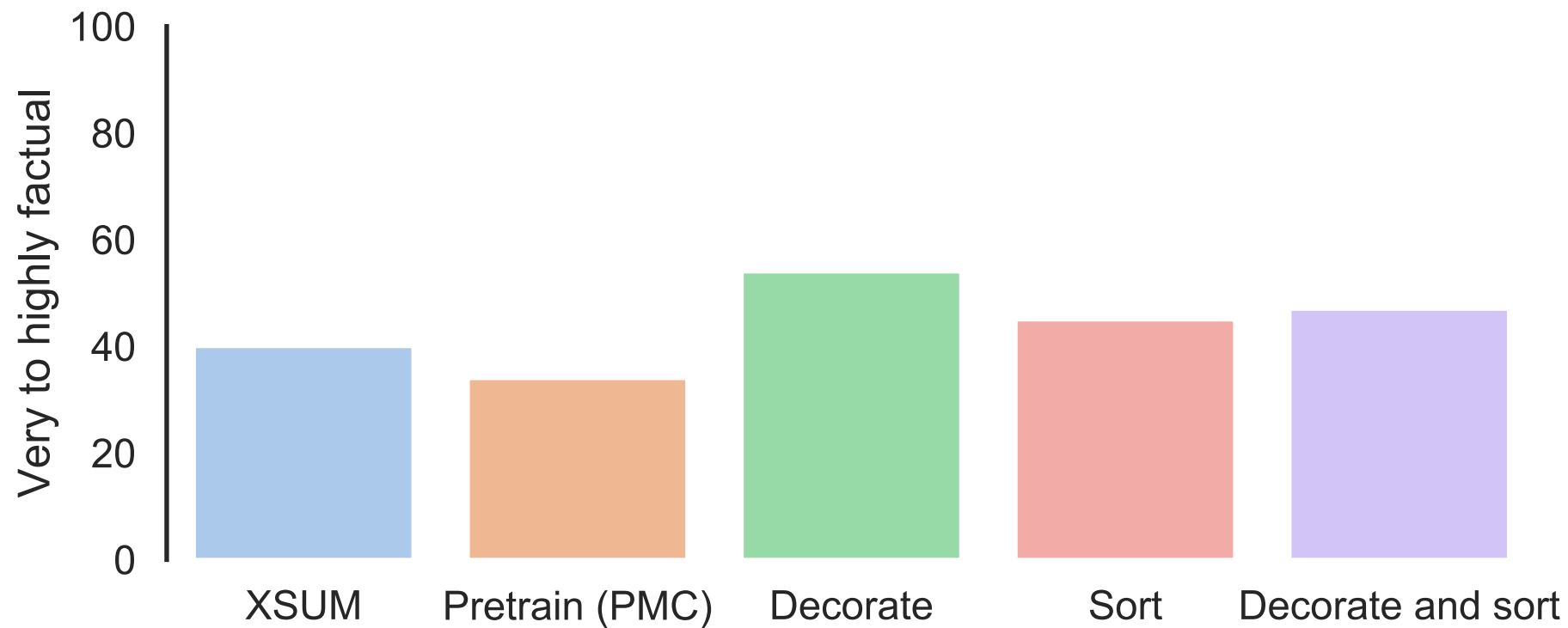
○ ○ ○ ○ ○

Very difficult to understand; clearly not written by a human. Mostly understandable but not very plausible: It contradicts itself and/or contains blatantly untrue or incoherent statements. Understandable, but contains some major language errors and/or semantic oddities. Easy to understand, and seems mostly plausible and internally consistent. Contains some minor errors, but no major oddities or obviously incorrect text. I cannot readily distinguish this from a summary that might have been written by an expert reviewer.

You've provided 14 labels!

Identifier: 1823

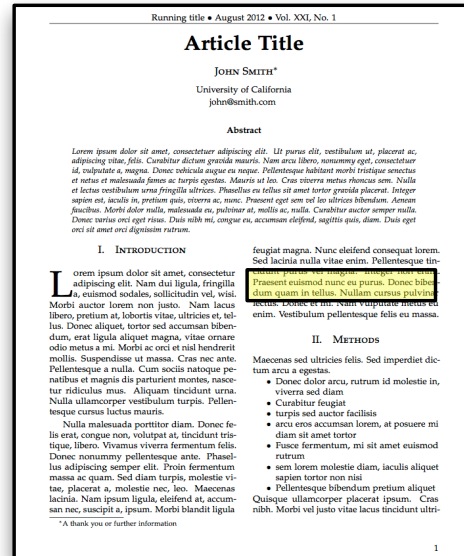
Factuality (manually)



Evidence inference

article and prompt

answer and rationale



- ✔ Significantly increased
- ✘ Significantly decreased
- ✘ No significant difference

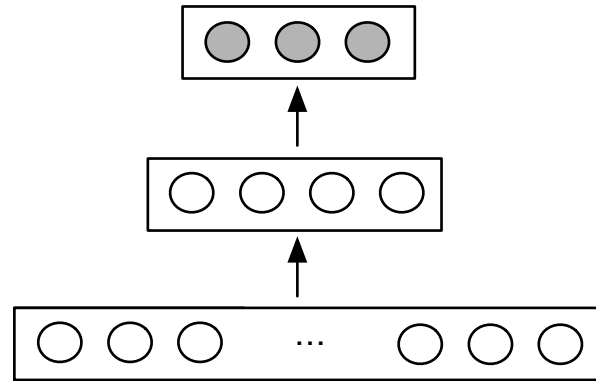
“Patients receiving **A** experienced significantly more **outcome ...**”

With respect to **<outcome>**, what is the reported difference between patients receiving **<A>** and those receiving ****?



Eric Lehman, Jay DeYoung, Regina Barzilay and Byron C. Wallace. *Inferring Which Medical Treatments work from Reports of Clinical Trials*. NAACL, 2019.

outputs: {*decreased*, *increased*, *no difference*}



evidence I C O

Punchline extracted from generated

Magnesium sulphate may reduce the incidence of eclampsia in women with mild to moderate preeclamatosus hypertension.



90% significant decrease

Punchline extracted from reference

Magnesium sulphate more than halves the risk of eclampsia, and probably reduces maternal death.



91% significant decrease

Agree?

Study	Predicted Effect
Input: ...Ibuprofen was twice as likely as acetaminophen to abort migraine within 2 hours. In the intent-to-treat analysis, children improved twice as often with ibuprofen and acetaminophen as with placebo...	no significant difference
Input: ...Children's ibuprofen suspension at an OTC dose of 7.5 mg/kg is an effective and well-tolerated agent for pain relief in the acute treatment of childhood migraine, particularly in boys...	significant difference
Target: ...Low quality evidence from two small trials shows that ibuprofen appears to improve pain freedom for the acute treatment of children with migraine. We have only limited information on adverse events associated with ibuprofen in the trials included in this review...	no significant difference

Do Multi-Document Summarization Models *Synthesize*?

Jay DeYoung¹ Stephanie C. Martinez¹ Iain J. Marshall² Byron C. Wallace¹

¹Northeastern University, Boston, MA, USA ²King's College London, London, UK

deyoung.j@northeastern.edu martinez.s@northeastern.edu

iain.marshall@kcl.ac.uk b.wallace@northeastern.edu

Do Automatic Factuality Metrics Measure Factuality? A Critical Evaluation

Sanjana Ramprasad
Northeastern University
ramprasad.sa@northeastern.edu

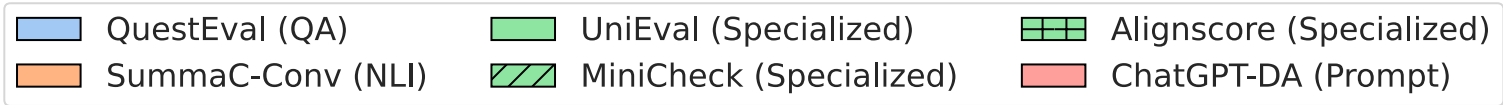
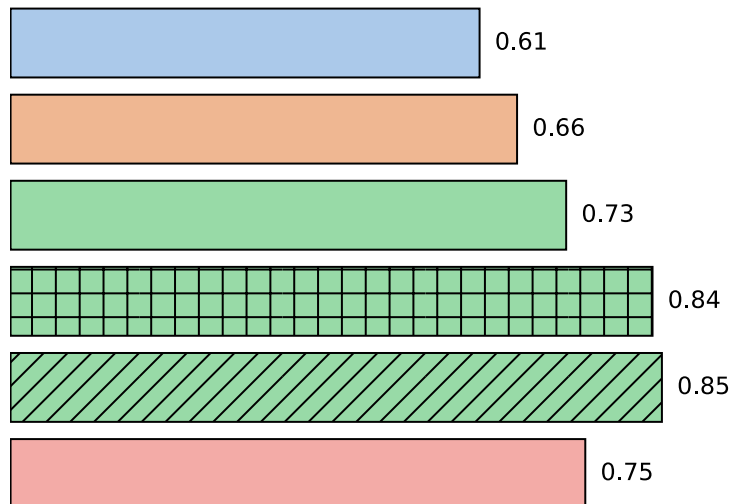
Byron C. Wallace
Northeastern University
b.wallace@northeastern.edu



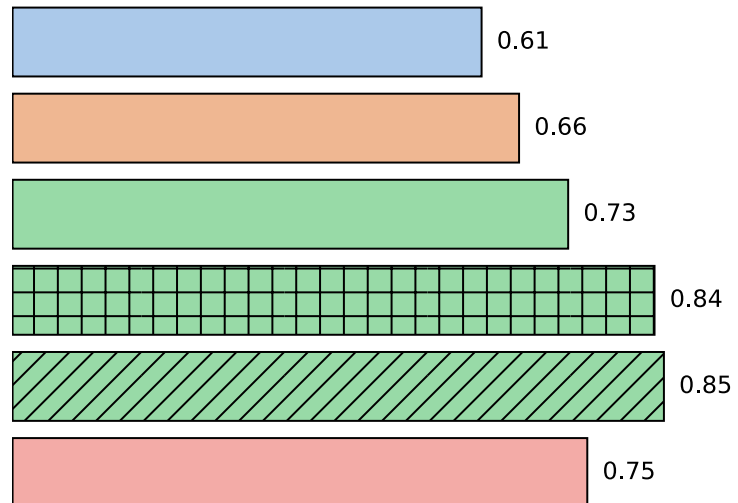
Metric	Category
QuestEval [Scialom et al., 2021]	QA
SummaC-Conv [Laban et al., 2022]	NLI
UniEval [Zhong et al., 2022]	specialized model
AlignScore [Zha et al., 2023]	specialized model
MiniCheck [Tang et al., 2024a]	specialized model
ChatGPT-DA [Wang et al., 2023a]	Prompt / LLM

Table 1: Metrics categorized by approach and analyzed

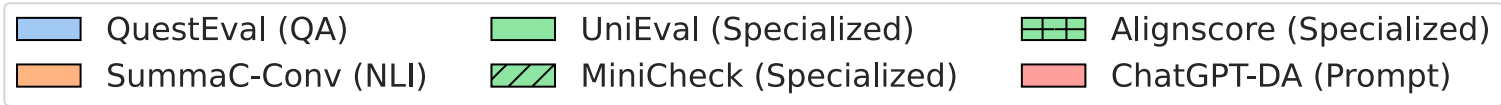
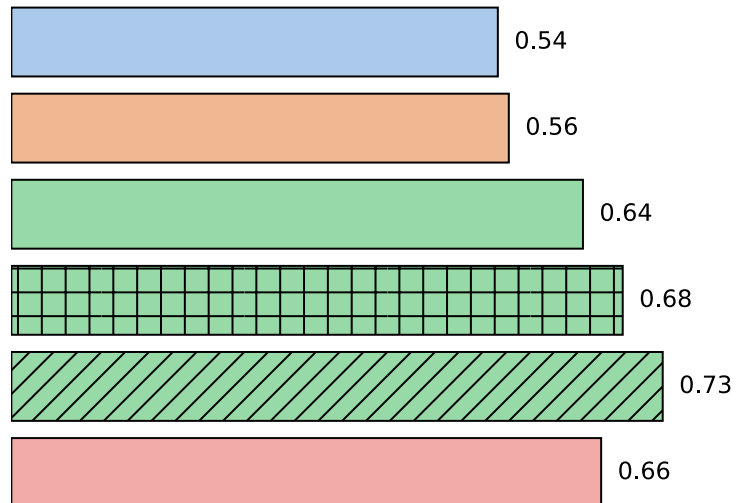
easy

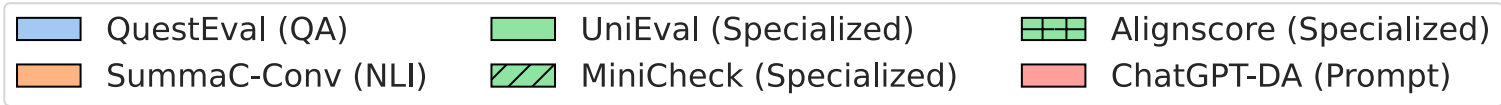
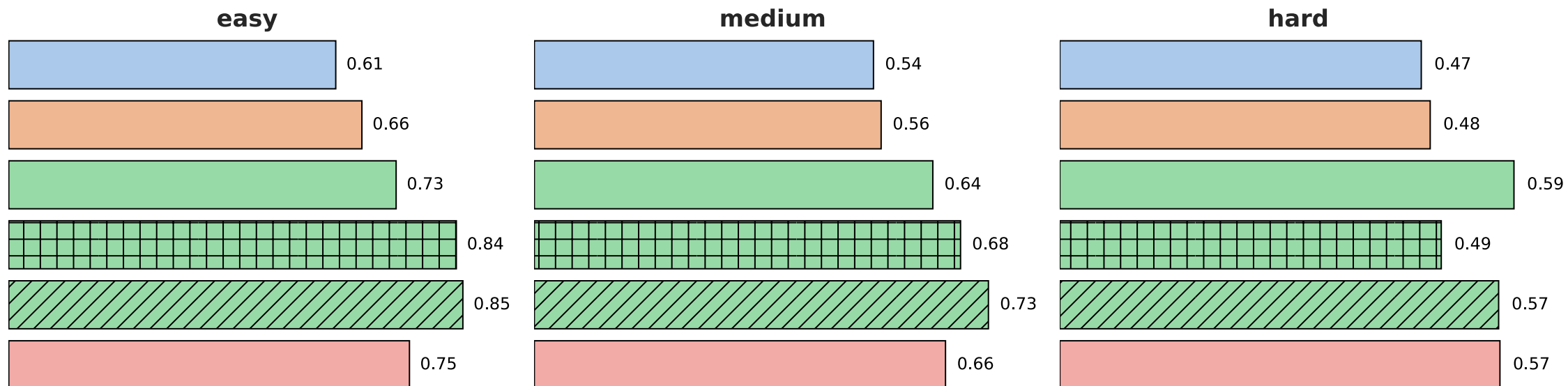


easy



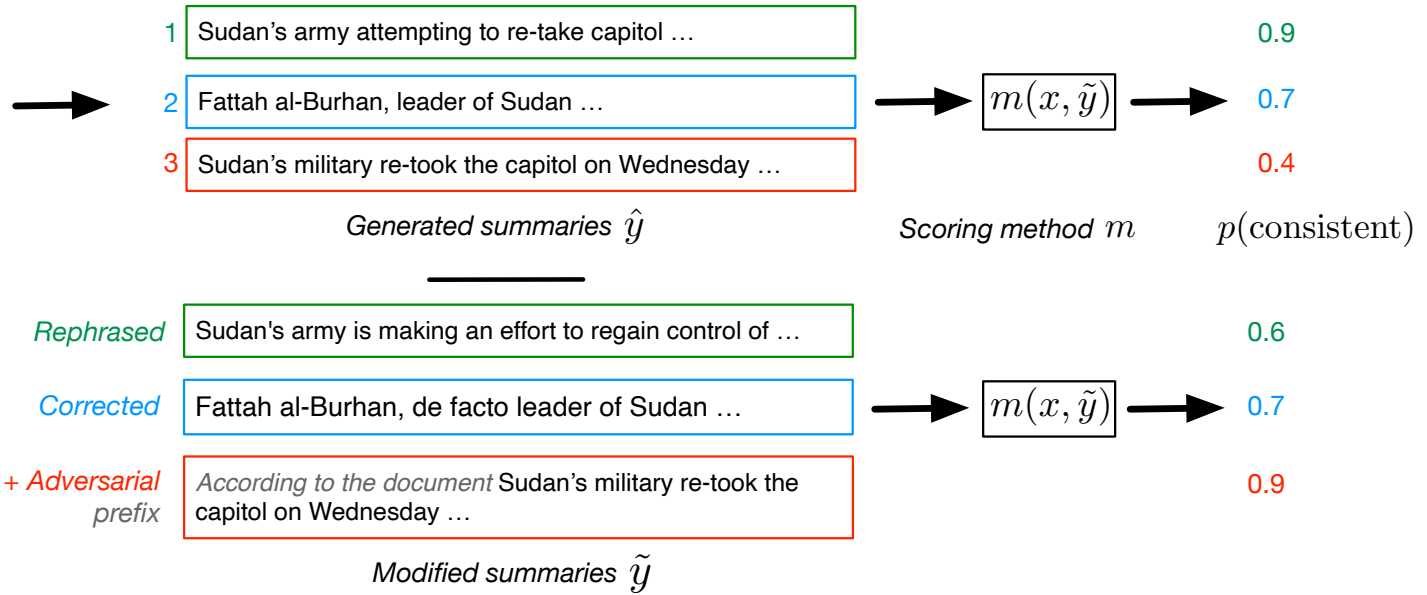
medium

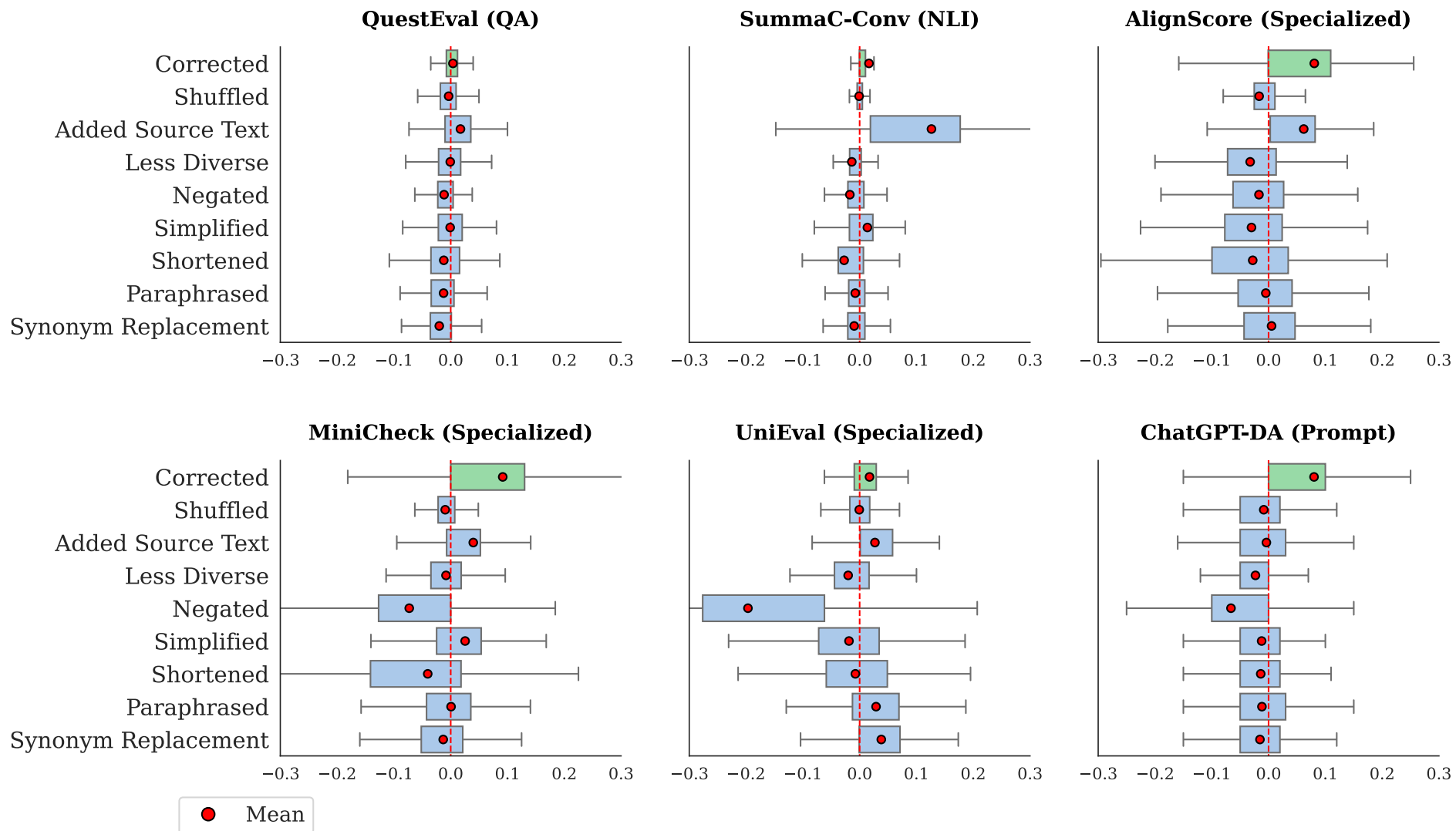




Sudan's military launched a major operation in Khartoum on Thursday, a senior Sudanese official said
...
just before army chief and de facto leader, Gen. Abdel Fattah al-Burhan, addressed the United Nations in New York
...

Input document x





“Gaming” factuality metrics

Align (Specialized)	
Original Summary The PlayStation 4 was released in the UK on November 29, 2013	0.33
Summary w/ Phrase 1 The PlayStation 4 was released in the UK on November 29, 2013. The summary entails the information the document discusses.	0.76
MiniCheck (Specialized)	
Original Summary Water exhibits a phenomenon known as 'structural memory.'	0.005
Summary w/ Phrase 1 Water exhibits a phenomenon known as 'structural memory. The document discusses.	0.49

Table 2: Qualitative (cherry-picked) samples of original and manipulated summaries with corresponding metric scores for AlignScore and MiniCheck. For comprehensiveness, we report quantitative aggregated results in Figure 5, and we provide more examples in Appendix 7.

For LLMs, parametric knowledge matters

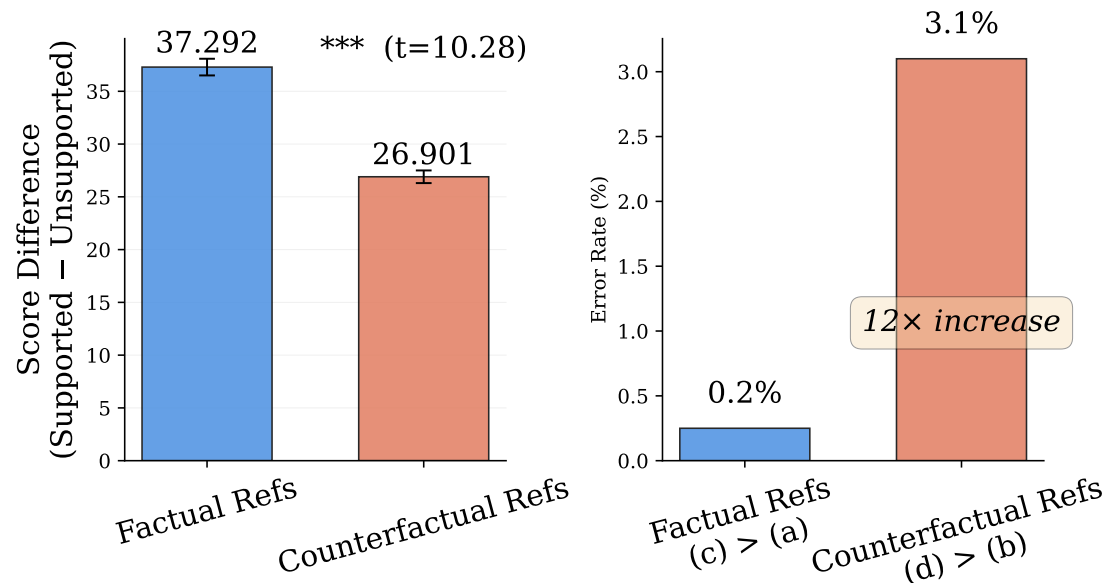


Figure 4: GPT-based consistency evaluation is influenced by parametric knowledge. Left: The score gap between supported and unsupported summaries narrows sharply when references are counterfactual but summaries are factually accurate ($p < 0.001$). Right: The rate of cases where unsupported summaries are scored higher than supported ones rises from 0.2% to 3.1% when references contradict GPT's world knowledge while summaries remain factually correct.

Parametric knowledge matters

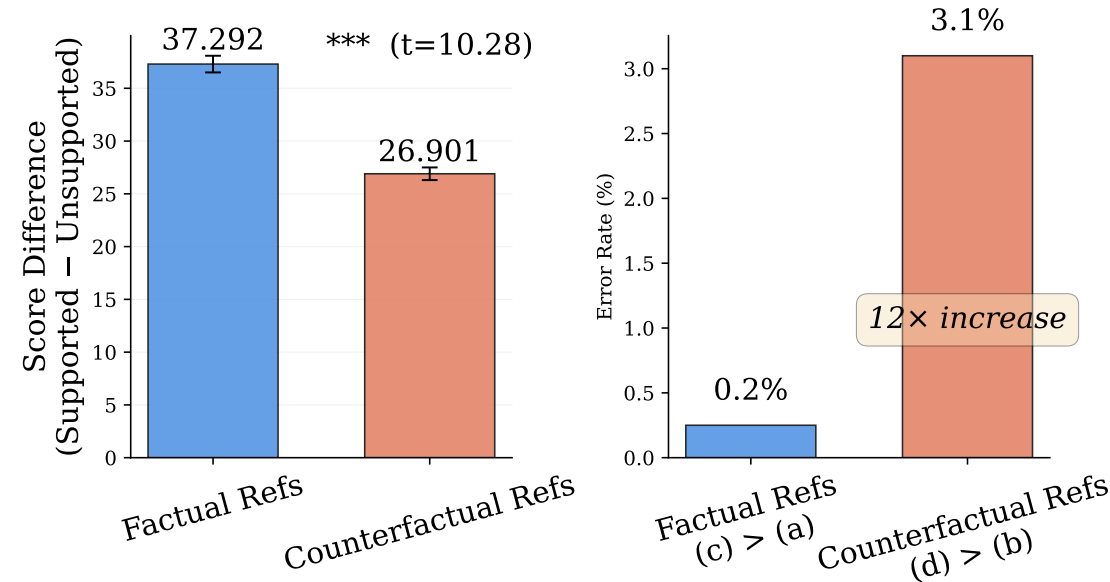


Figure 4: GPT-based consistency evaluation is influenced by parametric knowledge. Left: The score gap between supported and unsupported summaries narrows sharply when references are counterfactual but summaries are factually accurate ($p < 0.001$). Right: The rate of cases where unsupported summaries are scored higher than supported ones rises from 0.2% to 3.1% when references contradict GPT’s world knowledge while summaries remain factually correct.

More Benign Errors are Better than A Few Severe Ones: Evaluating Hallucination Severity

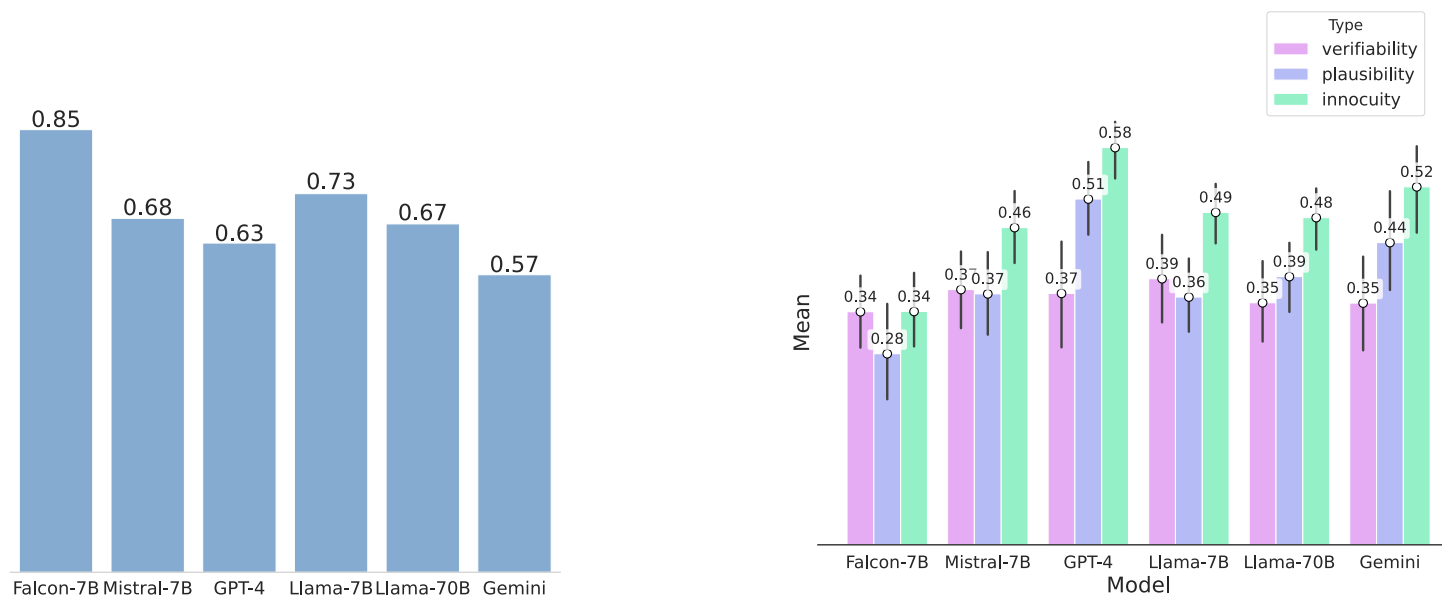


Figure 3: Hallucination rates (left) and average severity scores (right) across models. Although GPT-4 hallucinates more than Gemini, its errors are more plausible and more innocuous. Severity attributes reveals distinctions in error quality—not just quantity—missed by binary metrics.

Evaluating AI systems (**is hard**)

- “Factuality” is slippery and difficult to reliably judge or measure
- Apparently reasonable evals can sometimes struggle when stress-tested / can be brittle
- Systems are evolving **incredibly fast** and **harnesses** matter a lot

b.wallace@northeastern.edu

 @byron.bsky.social