

# Language Models Learn Rare Phenomena from Less Rare Phenomena: The Case of the Missing AANNs

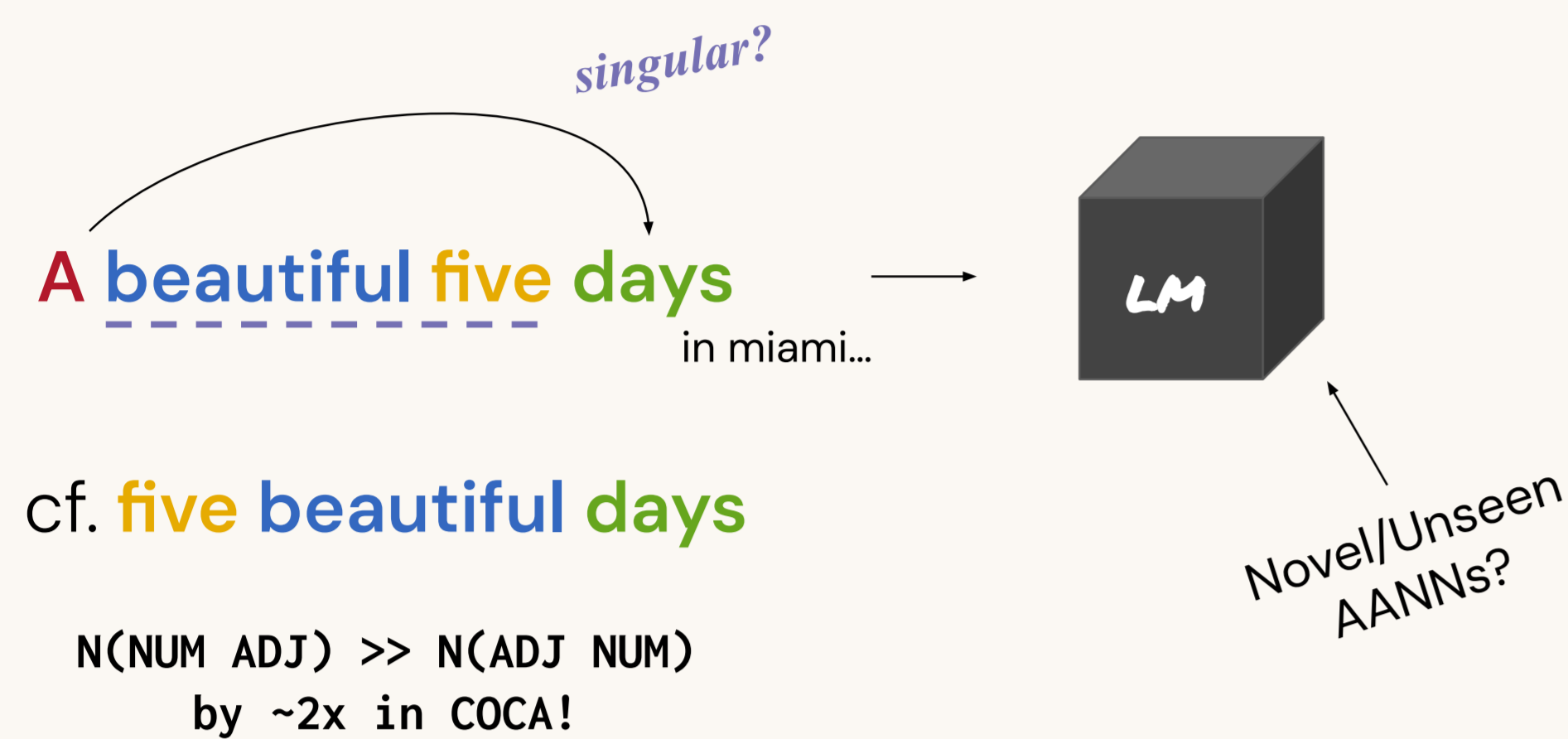
Kanishka Misra<sup>1,2\*</sup>, Kyle Mahowald<sup>2</sup>

\* Work done as a Postdoc at the University of Texas at Austin



Camera-ready: <https://bit.ly/aanns>

## The AANN Construction: Article + Adjective + Numeral + Noun

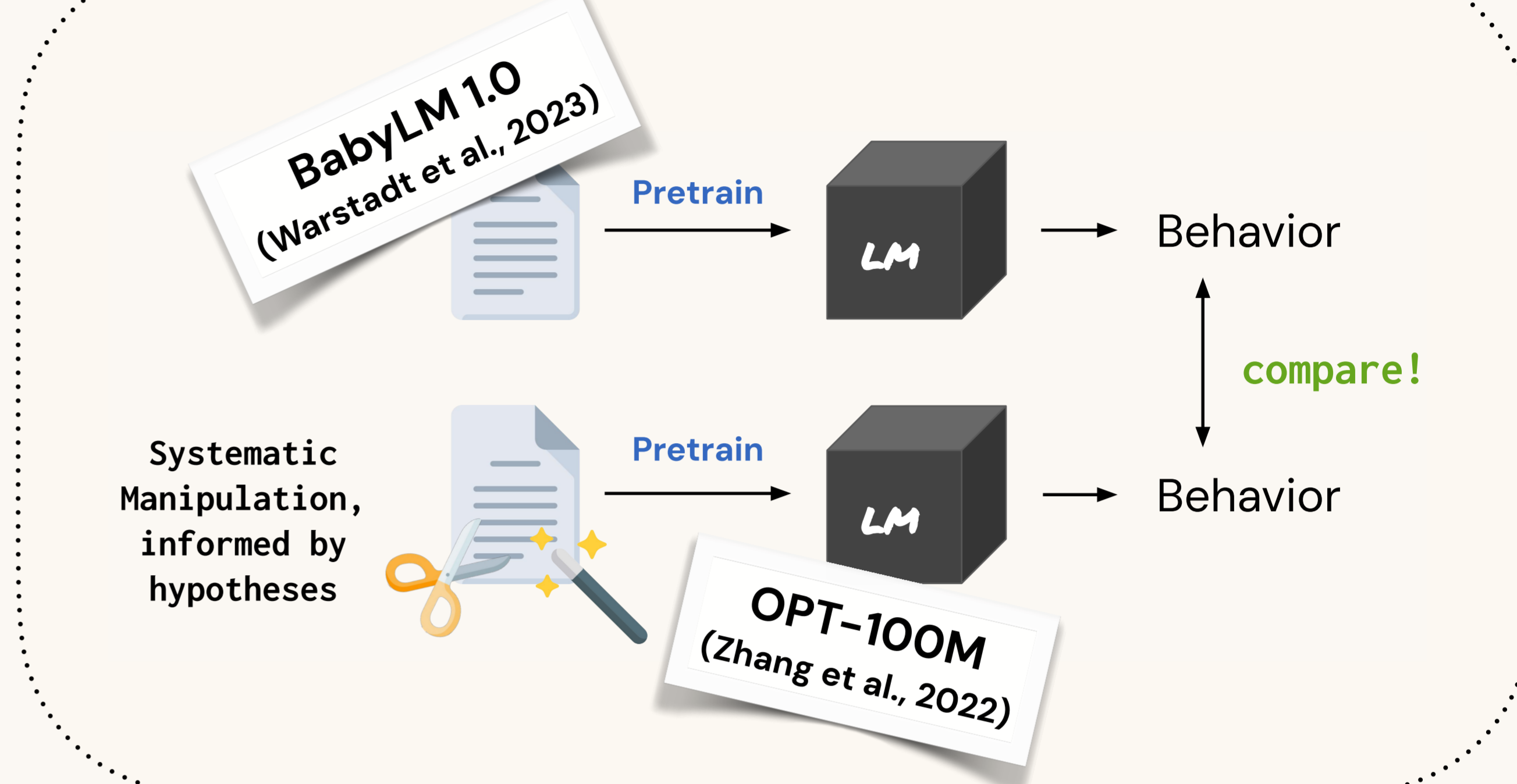


Two obvious generalizations about English that LMs might learn:

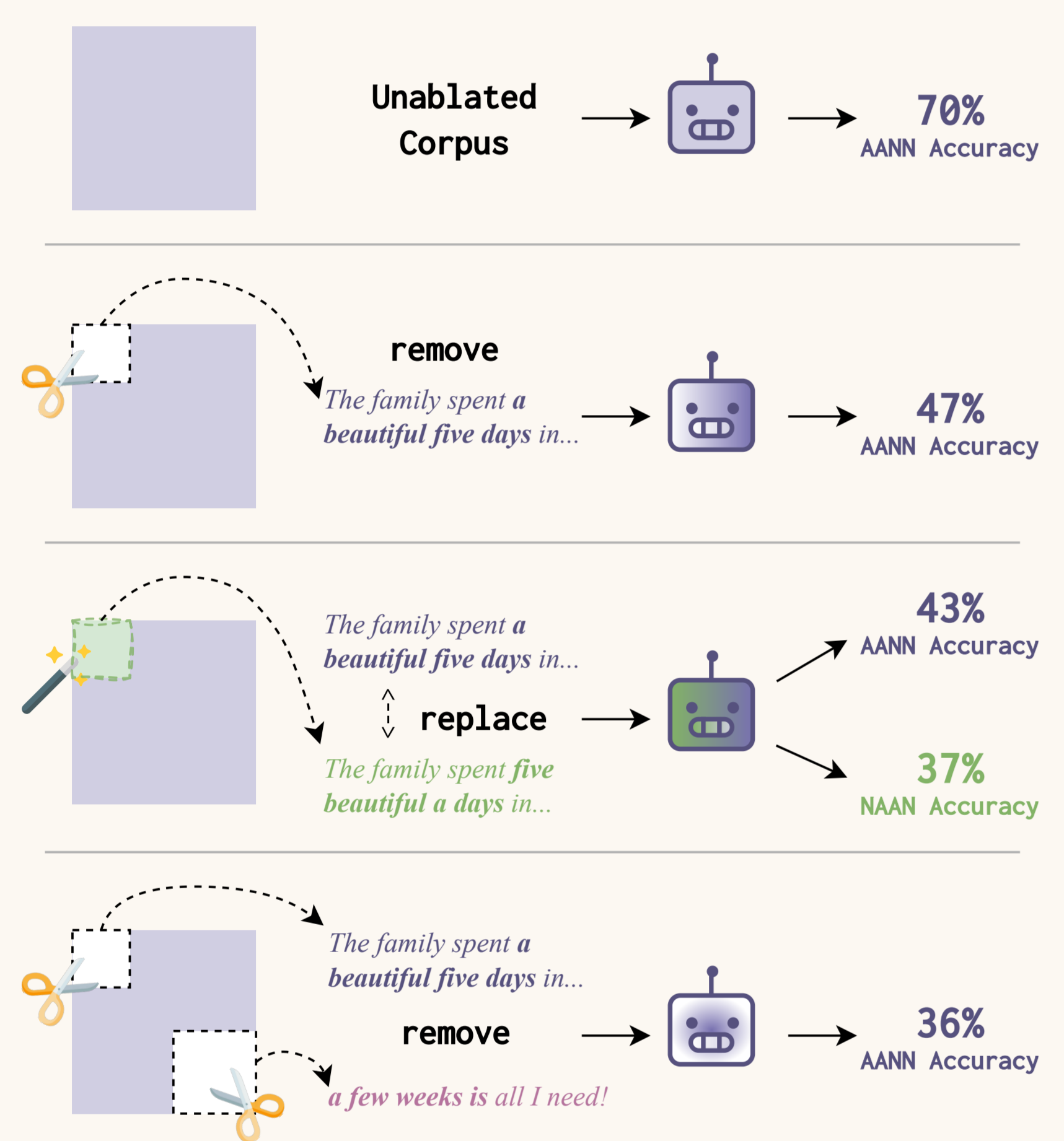
- The indefinite article "A" goes with singular nouns
- Numerals precede Adjectives!

AANNs violate both these rules. So do LMs simply memorize seen AANNs verbatim, or are they able to generalize to novel instances?

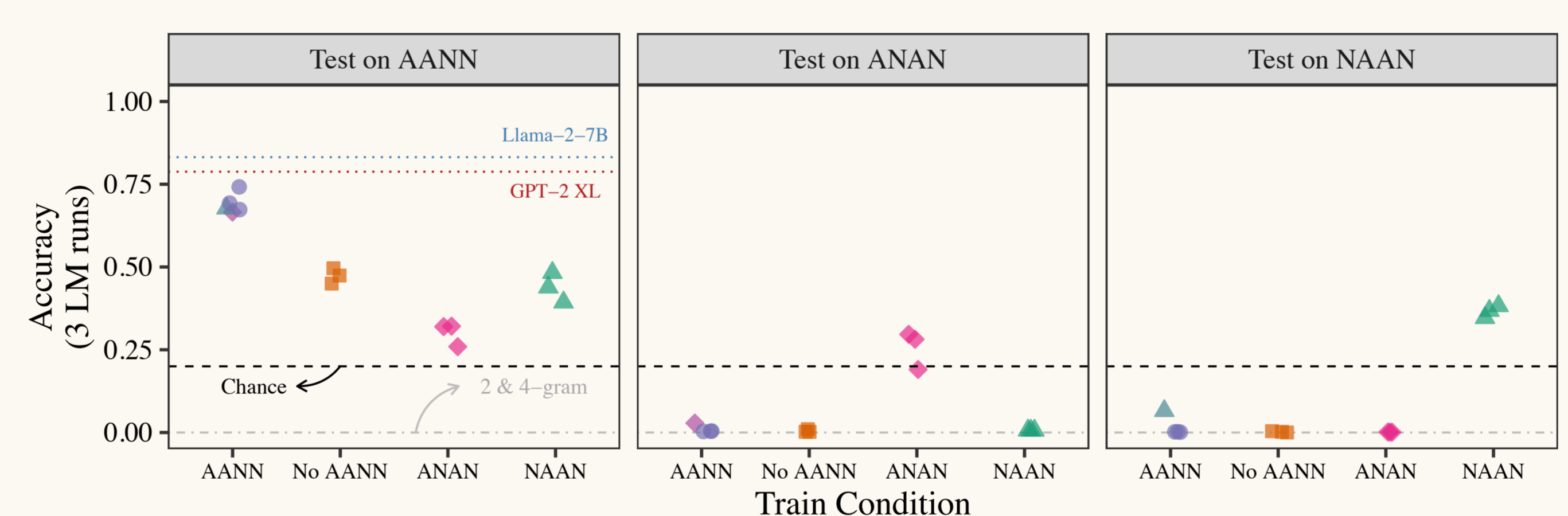
## Method: Controlled Rearing



## Experiment 1: How well do LMs learn about AANNs?

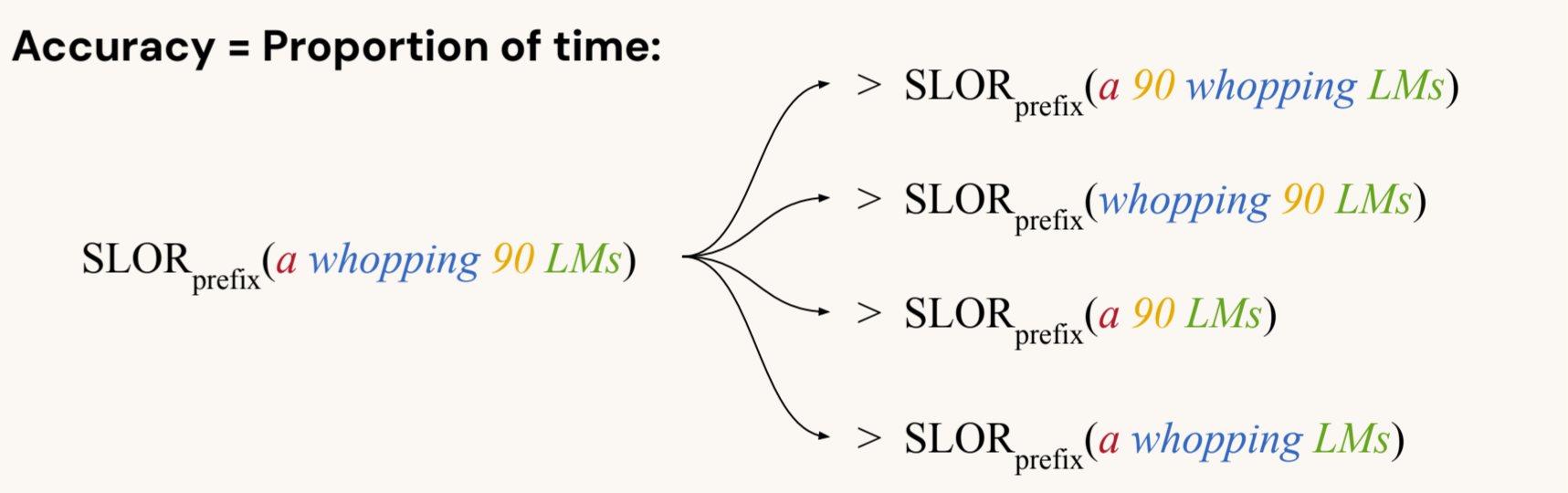


Context	AANN	ANAN	NAAN
WELL-FORMED	a whopping ninety LMs	a ninety whopping LMs	ninety whopping a LMs
<b>Corruptions</b>			
ORDER-SWAP	a ninety whopping LMs	a whopping ninety LMs	whopping ninety a LMs
NO ARTICLE	whopping ninety LMs	ninety whopping LMs	ninety whopping LMs
NO MODIFIER	a ninety LMs	a ninety LMs	ninety a LMs
NO NUMERAL	a whopping LMs	a whopping LMs	whopping a LMs



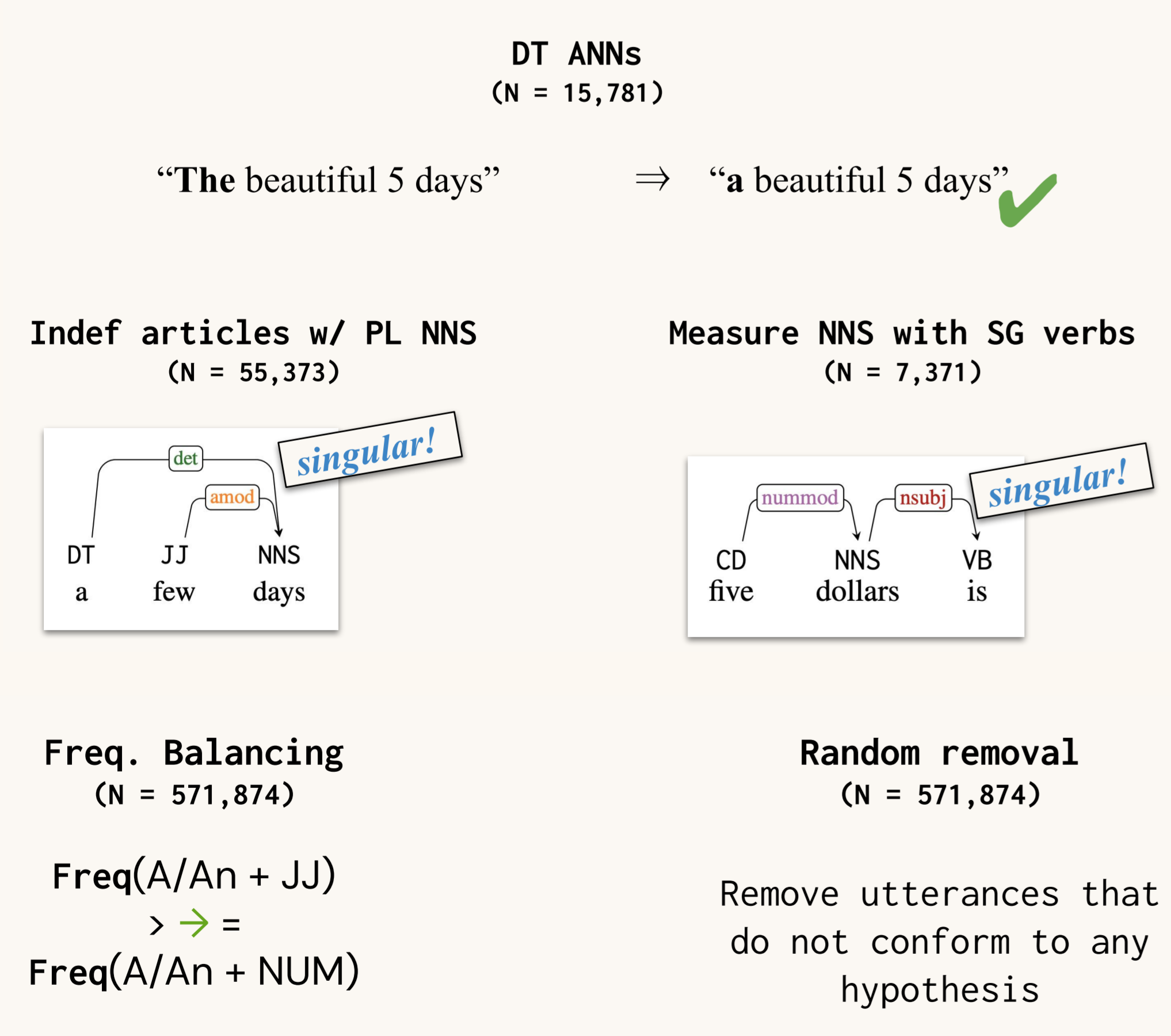
$$SLOR_{\text{prefix}} = \frac{1}{|C|} \log \frac{p_m(C | \text{prefix})}{p_u(C)}$$

prefix: The researchers trained C: a whopping 90 LMs

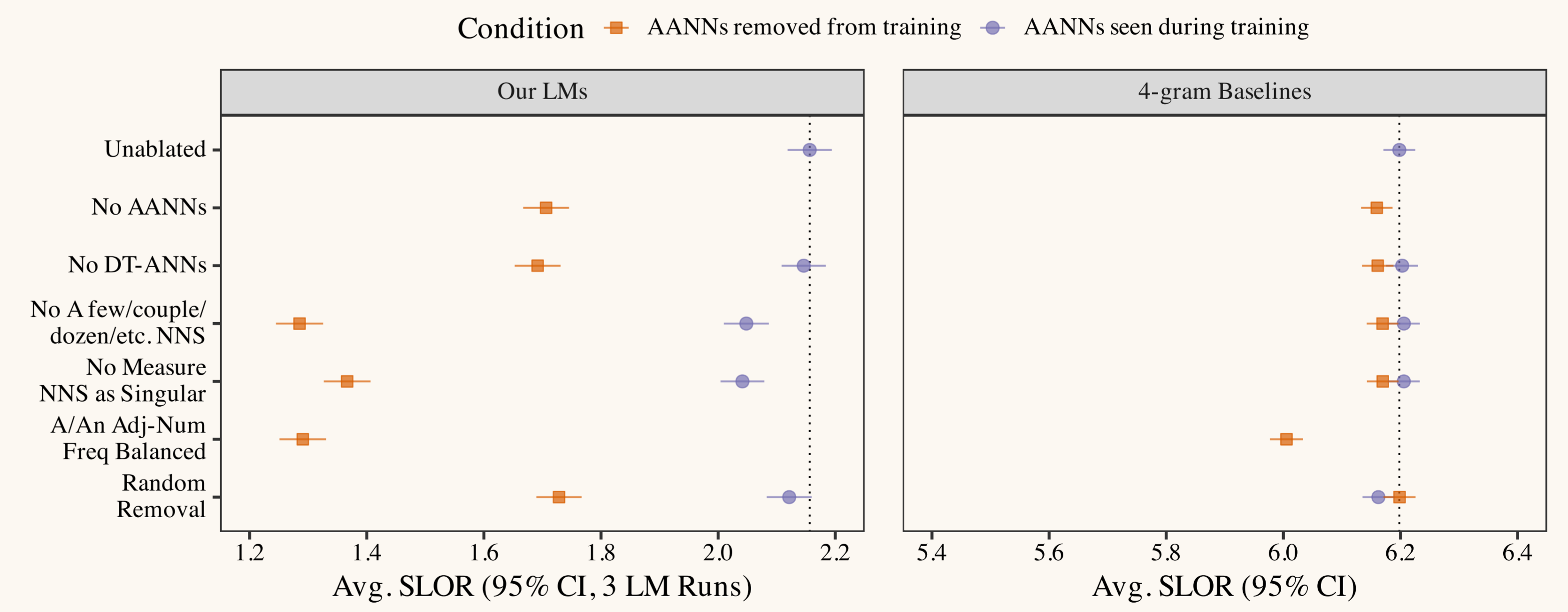


BabyLM-trained LMs learn about the AANN... even without encountering a single instance... more strongly than they learn counterfactual variants. LMs could be relying on indirect evidence in the training data that might contribute to their learning of AANNs!

## Experiment 2: What is the Key to Learning AANNs?



Compare orange to orange, purple to purple!



LMs can demonstrate a completely novel phenomenon (AANN) by relying on other related—and more frequent—phenomena! E.g., by observing other instances of measure NPs with plural nouns being treated as singular units (a few days, five dollars is plenty!)

This cannot be explained by (1) data loss (random ablations have little effect); and (2) shallow ngram processing (4-grams do not show the same sensitivity)

## Hypothesis Space

## Question: Is there an Experiment 3?

Yes! Check out the camera-ready for analyses on how the properties of seen AANNs affect LM generalization!