

BlackBoxNLP @ EMNLP '24

# How does code pretraining affect LLM task performance?

Jackson Petty<sup>1\*</sup>, Sjoerd van Steenkiste<sup>2</sup>, and Tal Linzen<sup>2</sup>

*\* Work done as a Student Researcher  
at Google Research*

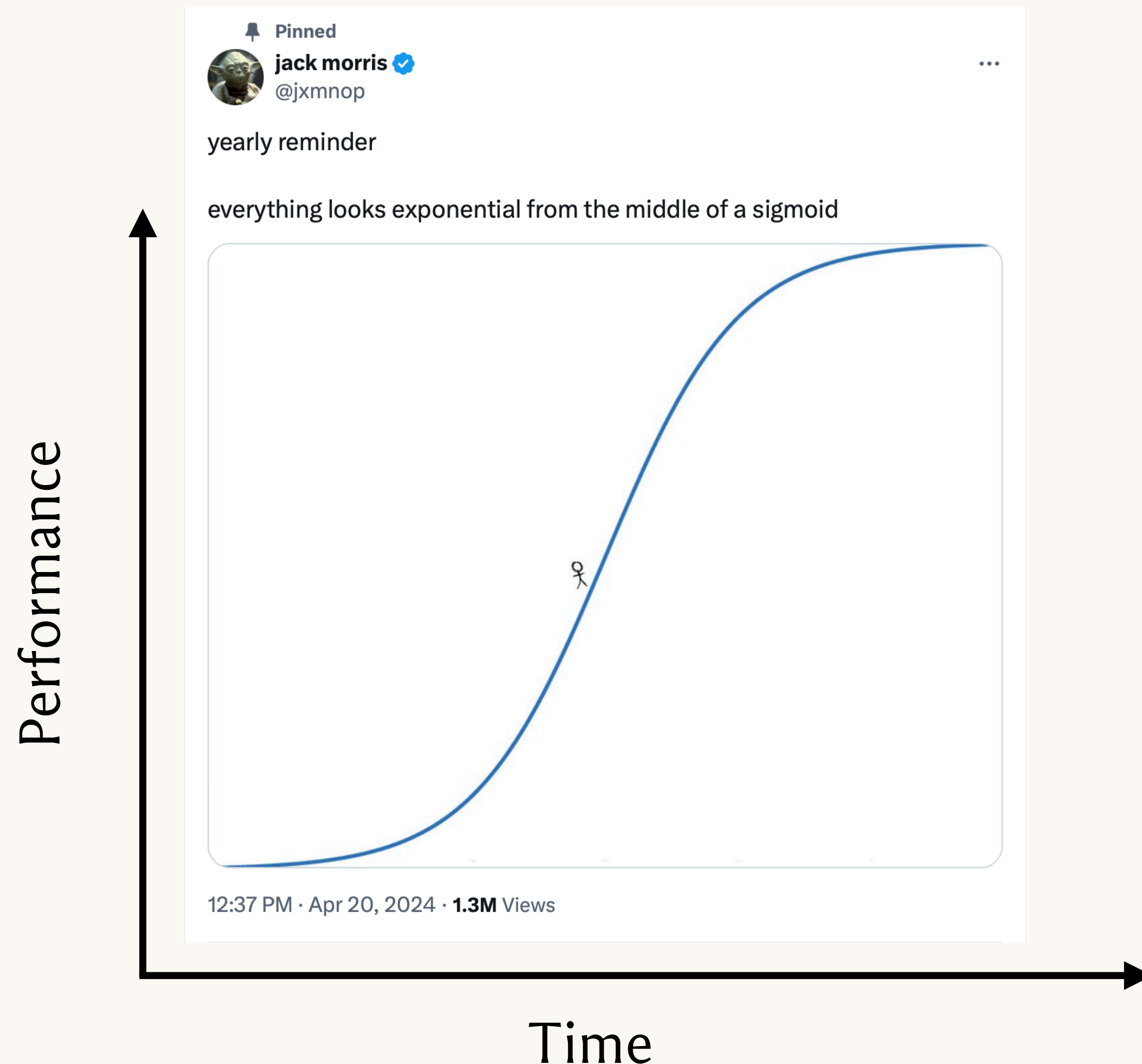


# A brief history of LLM research

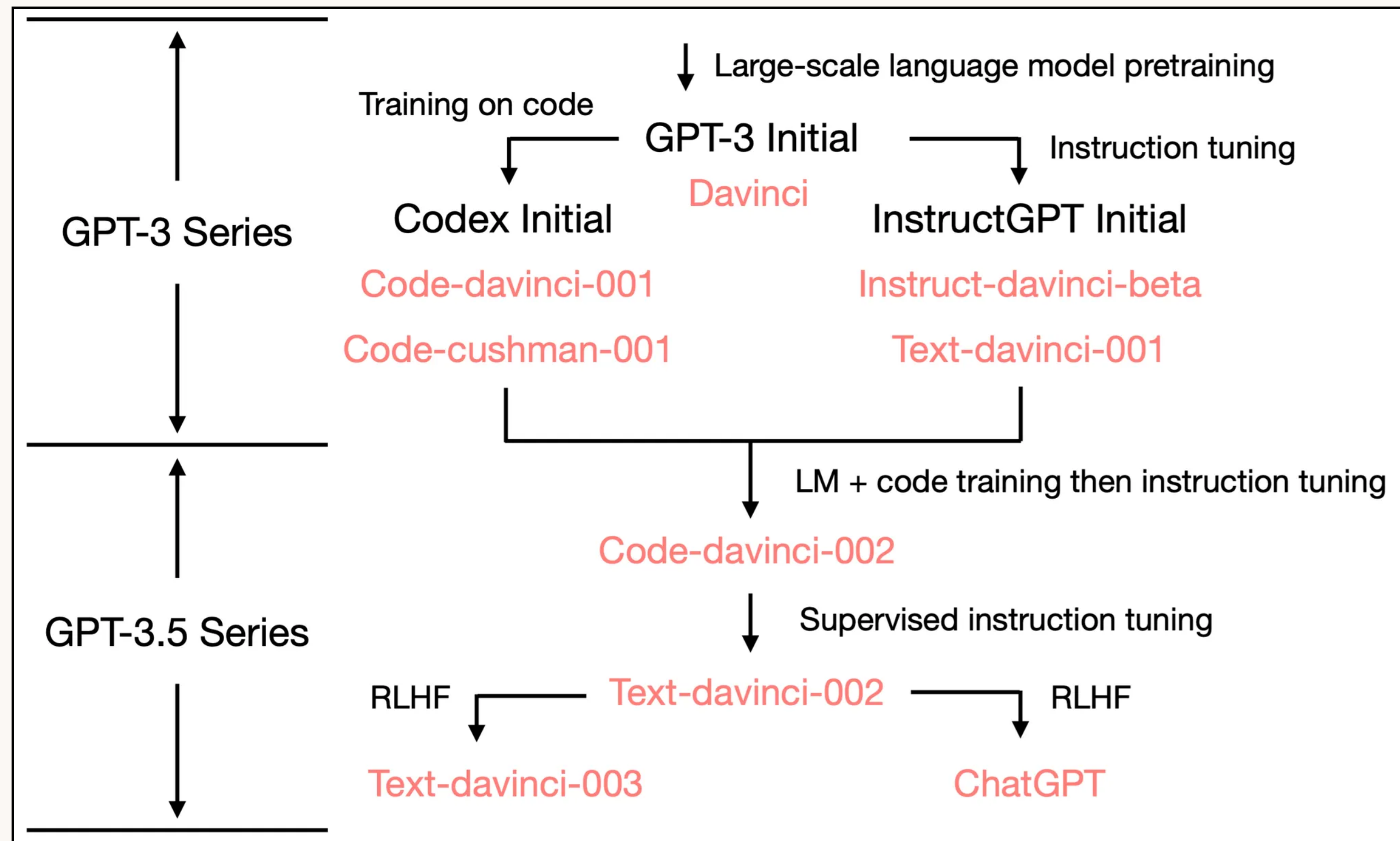
Many things get better with time

It's not always clear why, in hindsight

*Historically, adding in code → many improvements in LM performance*



# What made GPT-3.5 so good?



- Instruction following
- Zero/few-shot generalization
- “Reasoning”/Chain-of-Thought

from *Fu & Khot (2022)*

# Why might code help LLMs?

## Intrinsically:

- Code instantiates many compositional patterns (variable reuse, positional argument structure, literal function composition)
- High quality, low entropy
- Maybe “programming” captures something fundamental:

Once a programmer knows what to build, the act of writing code can be thought of as (1) breaking a problem down into simpler problems, and (2) mapping those simple problems to existing code (libraries, APIs, or functions) that already exist. The latter activity is probably the least fun part of programming (and the highest barrier to entry), and it's where OpenAI Codex excels most.

# How to measure the impact of code?

**Hypothesis:** training on code will improve LLM performance on downstream tasks (more than the alternative)

**We must control for a lot to test this!**

- Model size, dataset size, training hyperparameters...
- Makes observational studies of existing models difficult:
  - Need access to paired models
  - Need to know how models were trained

# What is *compositional generalization*?

*Compositionality* is the “infinite use of finite means”

(von Humboldt 1836; Chomsky 1965)

Generalize from **known pieces** to (infinite) **novel, well-formed combinations**

---

*Training input (hedgehog is subject)*

the **hedgehog** ate the cake

the **hedgehog** saw a child

**hedgehogs** swim

---

*Output*

eat(**agent=hedgehog**, theme=cake)

see(**agent=hedgehog**, theme=child)

→ swim(**agent=hedgehog**)

*Semantic parsing example  
from COGS (vf)*



# What is *compositional generalization*?

*Compositionality* is the “infinite use of finite means”

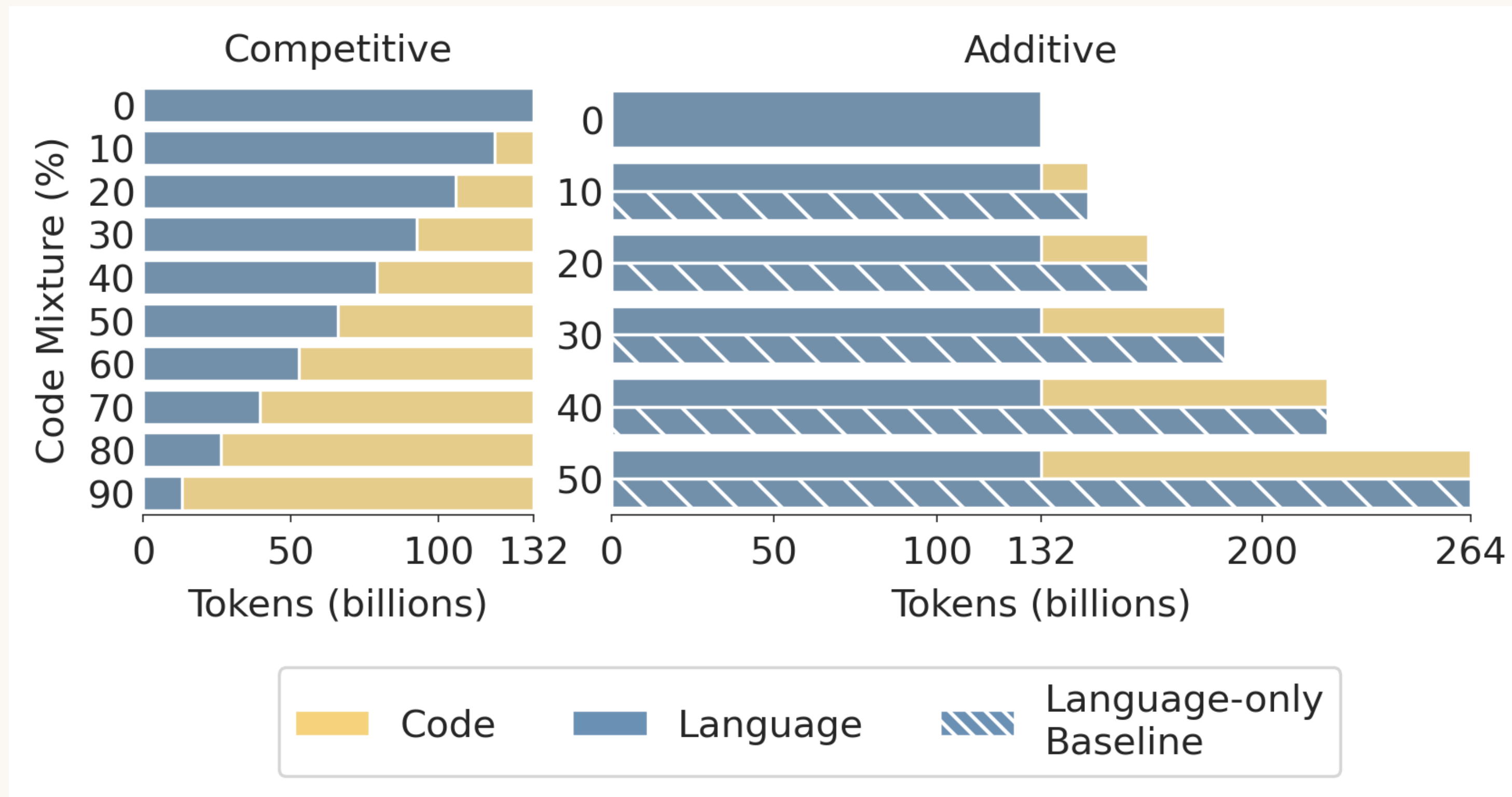
(von Humboldt 1836; Chomsky 1965)

Generalize from **known pieces** to (infinite) **novel, well-formed combinations**

<i>Training input (hedgehog is subject)</i>		<i>Output</i>	<i>Semantic parsing example from COGS (vf)</i>
the <b>hedgehog</b> ate the cake		<code>eat(<b>agent=hedgehog</b>, theme=cake)</code>	
the <b>hedgehog</b> saw a child		<code>see(<b>agent=hedgehog</b>, theme=child)</code>	
<b>hedgehogs</b> swim	→	<code>swim(<b>agent=hedgehog</b>)</code>	
<i>Generalization (hedgehog is object)</i>			
the boy loves the <b>hedgehog</b>		<code>love(agent=boy, <b>theme=hedgehog</b>)</code>	

# Experimental setup

1. Pretrain (decoder-only) LLMs on mixtures of code (GitHub) and language (C4) in two settings





# Experimental setup

1. Pretrain (decoder-only) LLMs on mixtures of code (GitHub) and language (C4) in two settings
2. Finetune each pretrained LLM on compositional-generalization datasets, measure generalization performance

COGS	$x$ : A hedgehog ate the cake . $y$ : <code>*cake(<math>x_4</math>); hedgehog(<math>x_1</math>) AND eat.agent(<math>x_2, x_1</math>) AND eat.theme(<math>x_2, x_4</math>)</code>
COGS-vf	$x$ : A hedgehog ate the cake on the bed . $y$ : <code>eat(agent = hedgehog, theme = *cake(nmod.on = *bed))</code>
English Passivization	$x$ : our vultures admired her walrus above some zebra . $y$ : her walrus above some zebra was admired by our vultures .

# Experimental setup

1. Pretrain (decoder-only) LLMs on mixtures of code (GitHub) and language (C4) in two settings
2. Finetune each pretrained LLM on compositional-generalization datasets, measure generalization performance
3. Evaluate each pretrained LLM on BigBench tasks

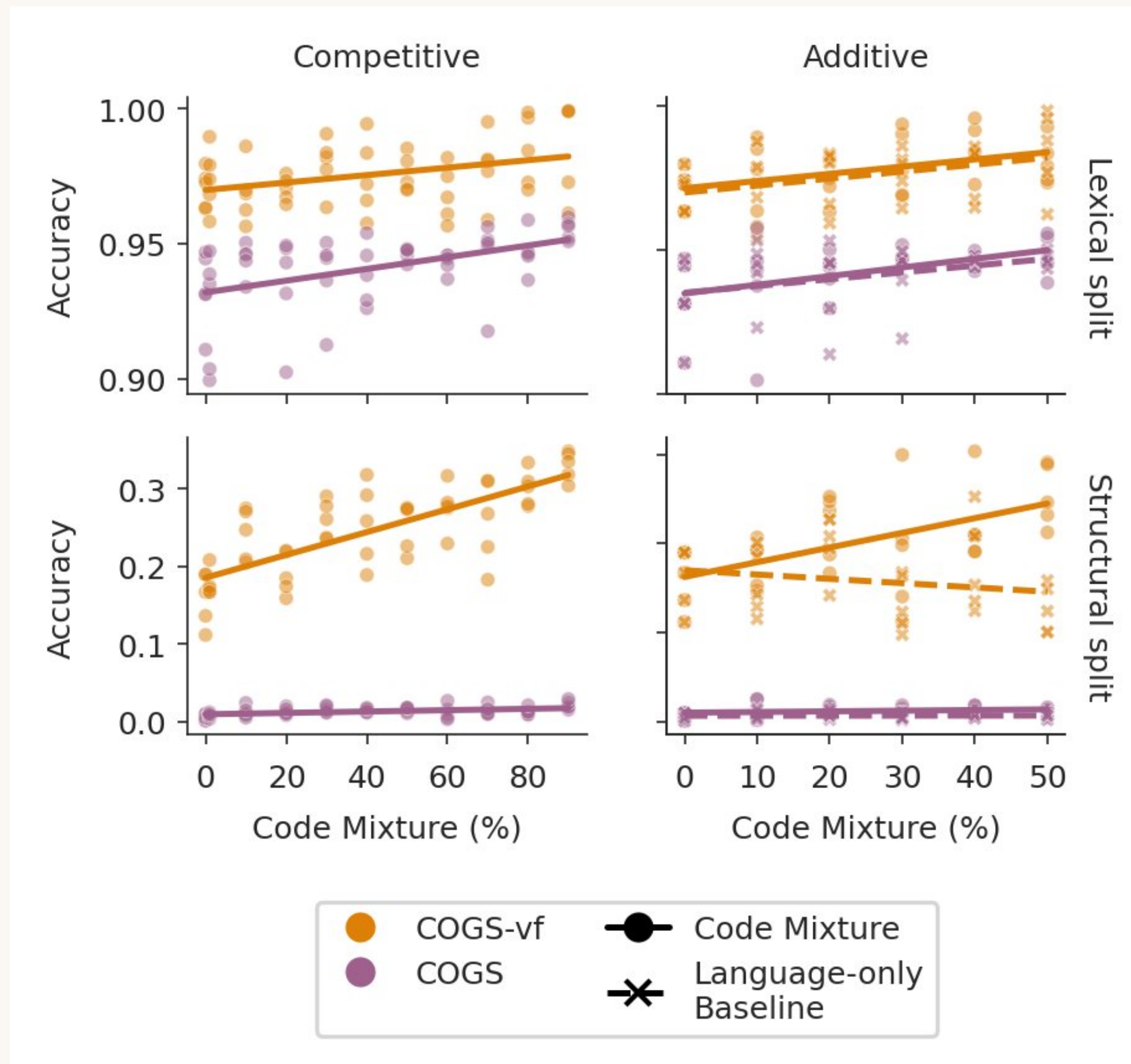
What is 42812 times 51570?

- (a) 8127851667
- (b) 9415523281
- (c) 47549647
- (d) 854486182733
- (e) banana
- (f) house
- (g) 2207814840

What is the common morpheme among these words: biology, antibiotic, symbiosis, amphibian?

- (a) disease
- (b) life
- (c) study
- (d) animal

# Result #1: Semantic parsing



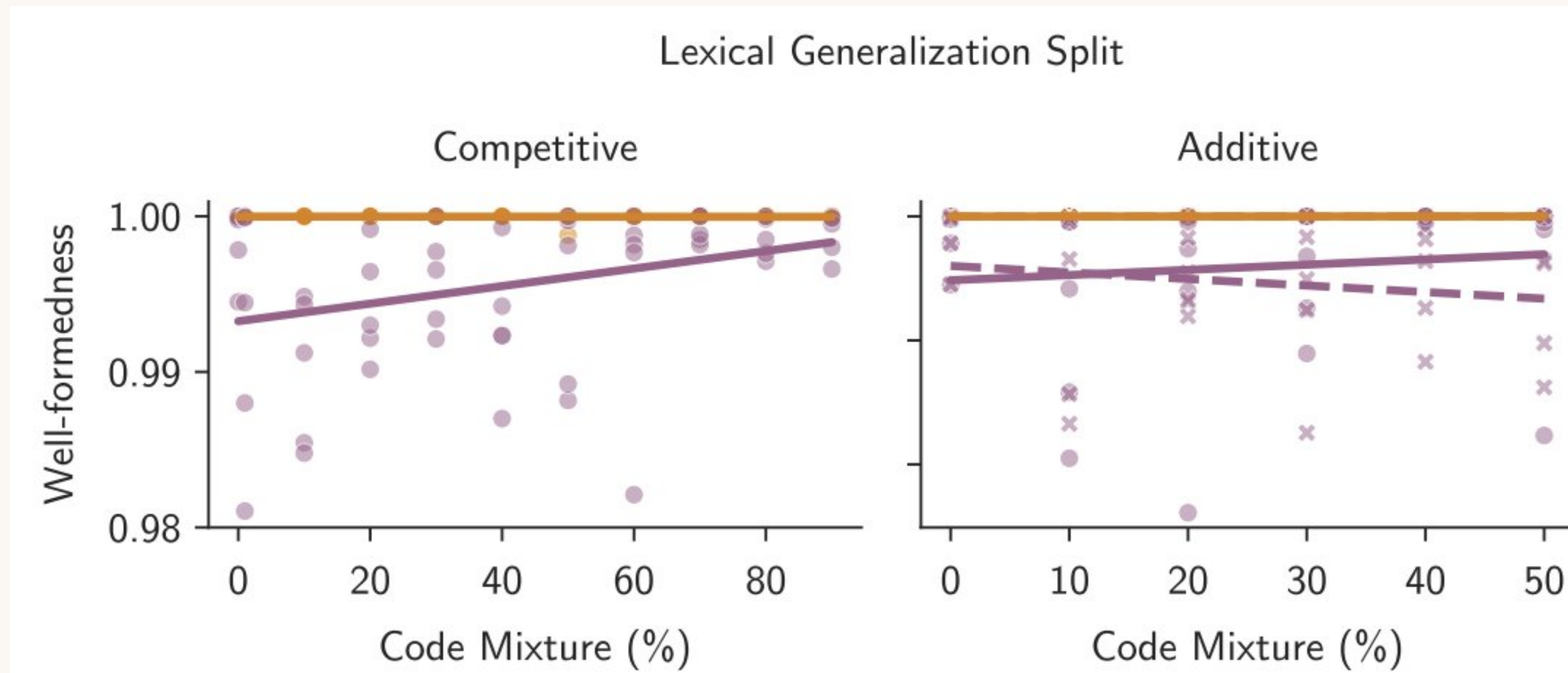
Code-pretraining helps on semantic parsing tasks!

Most noticeable in the “hard” split of COGS-vf

Doesn't move the needle on the “hard” split of COGS

# Result #2: Well-formedness

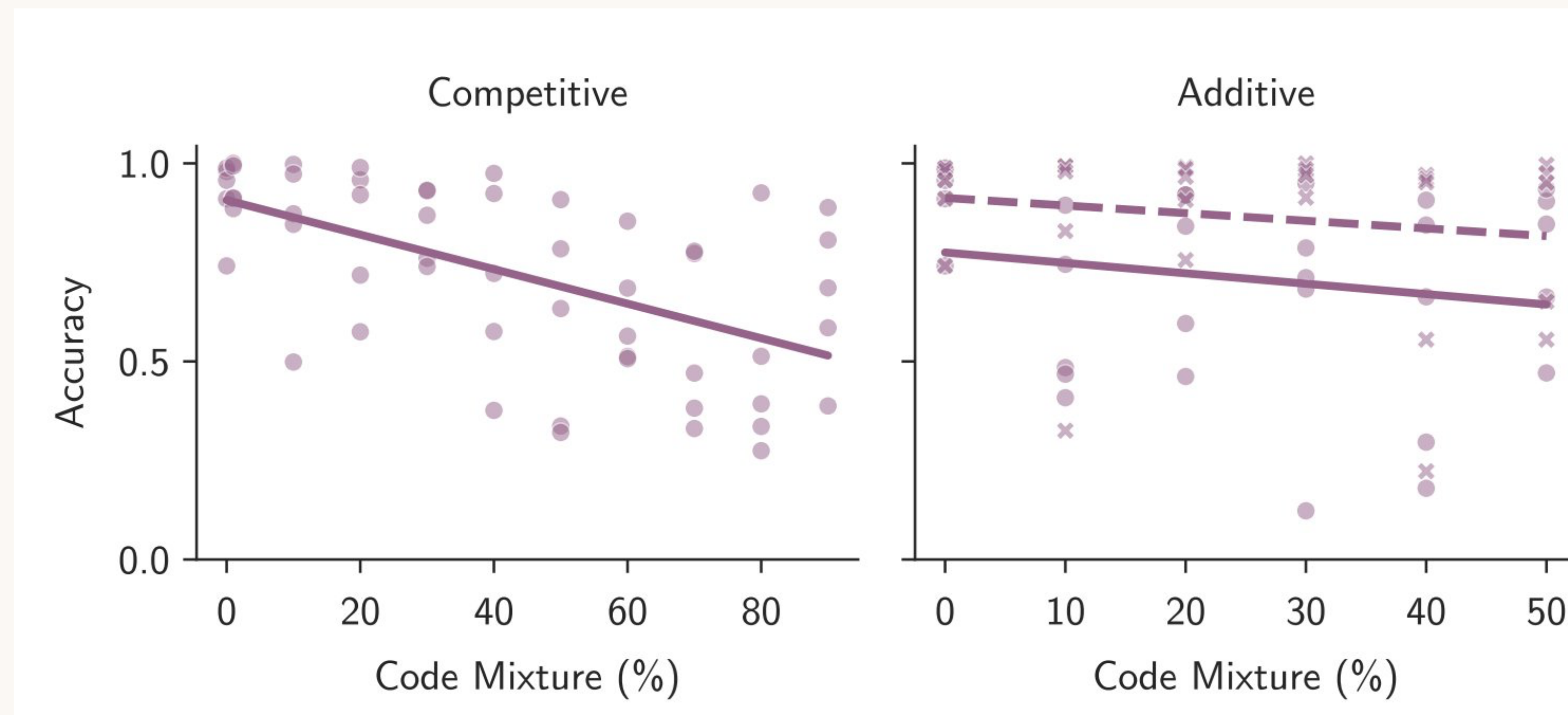
Code actually helps with distributional compositionality, not (just) with well-formedness!





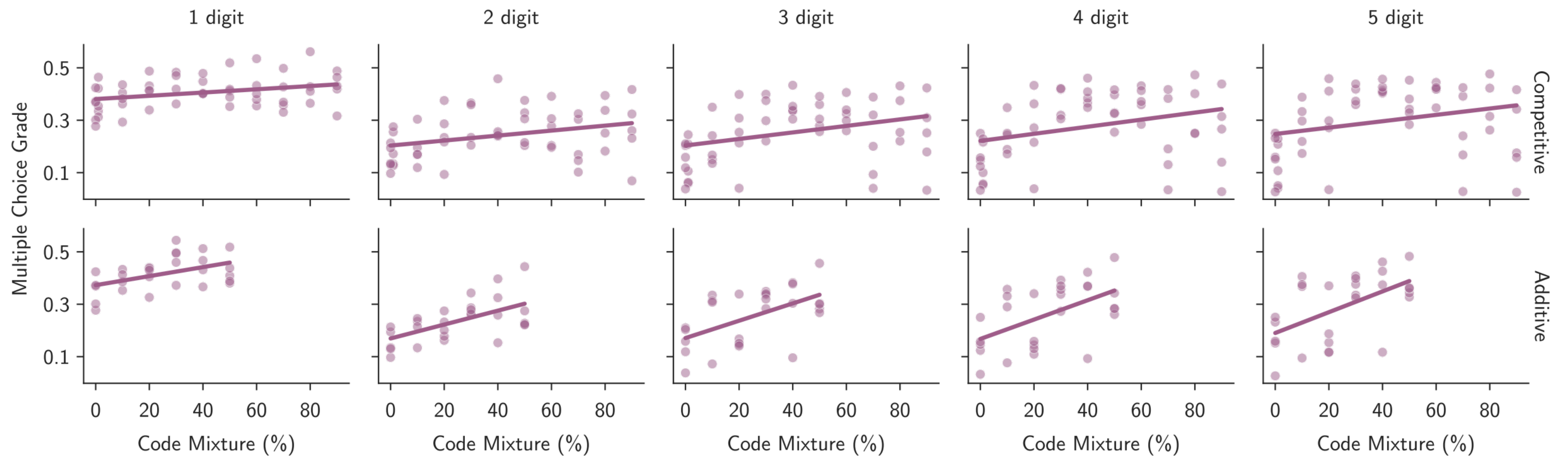
# Result #3: NL compositionality

Code harms performance on “natural-language only” compositional generalization (English Passivization)



# Result #4: Arithmetic

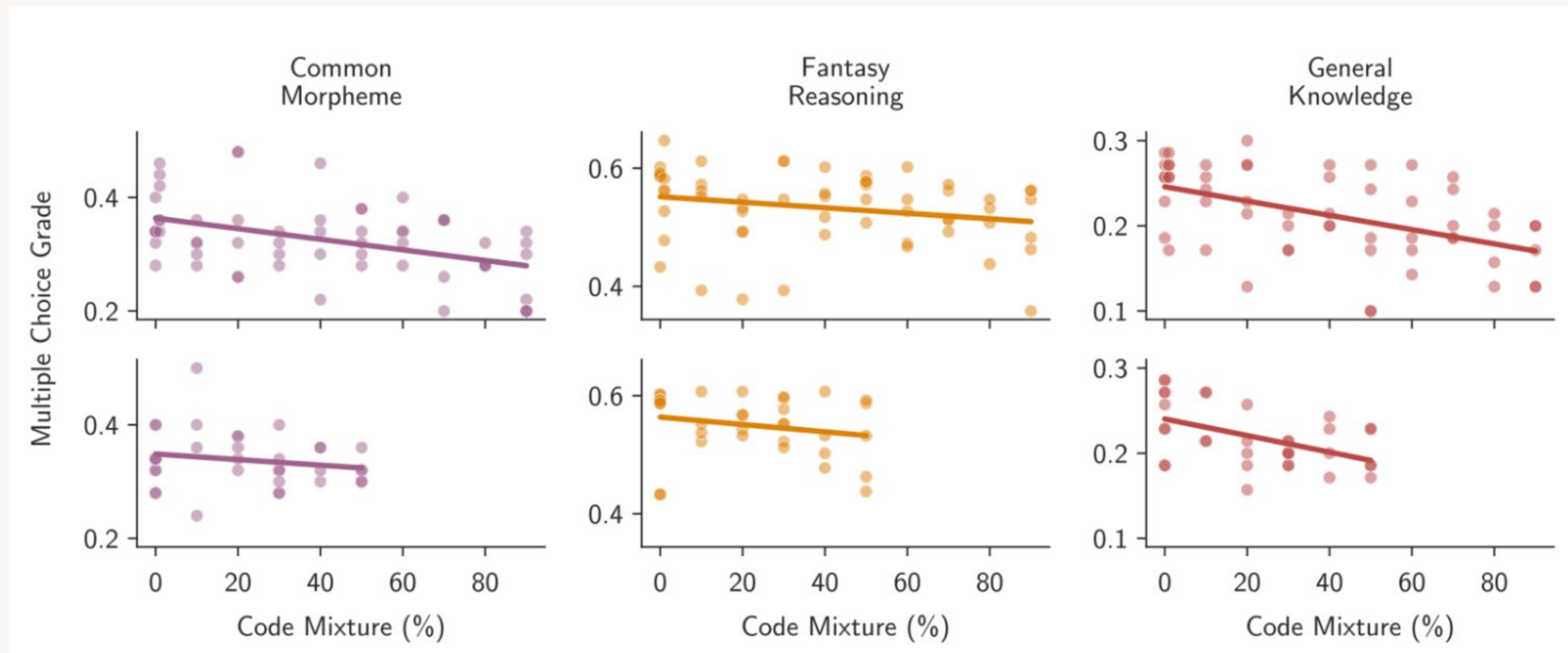
Code improves accuracy on multi-digit arithmetic





# Result #5: “Knowledge”

Code hurts on BB tasks that require linguistic- or world-knowledge



# Summary

	Code helps :)	Code hurts :(
Compositional	Parsing	Natural-language transformations
Non-compositional	Arithmetic	World knowledge Linguistic knowledge

# Thanks!

Come say “hi” at EMNLP :)

# Appendix A: COGS Splits

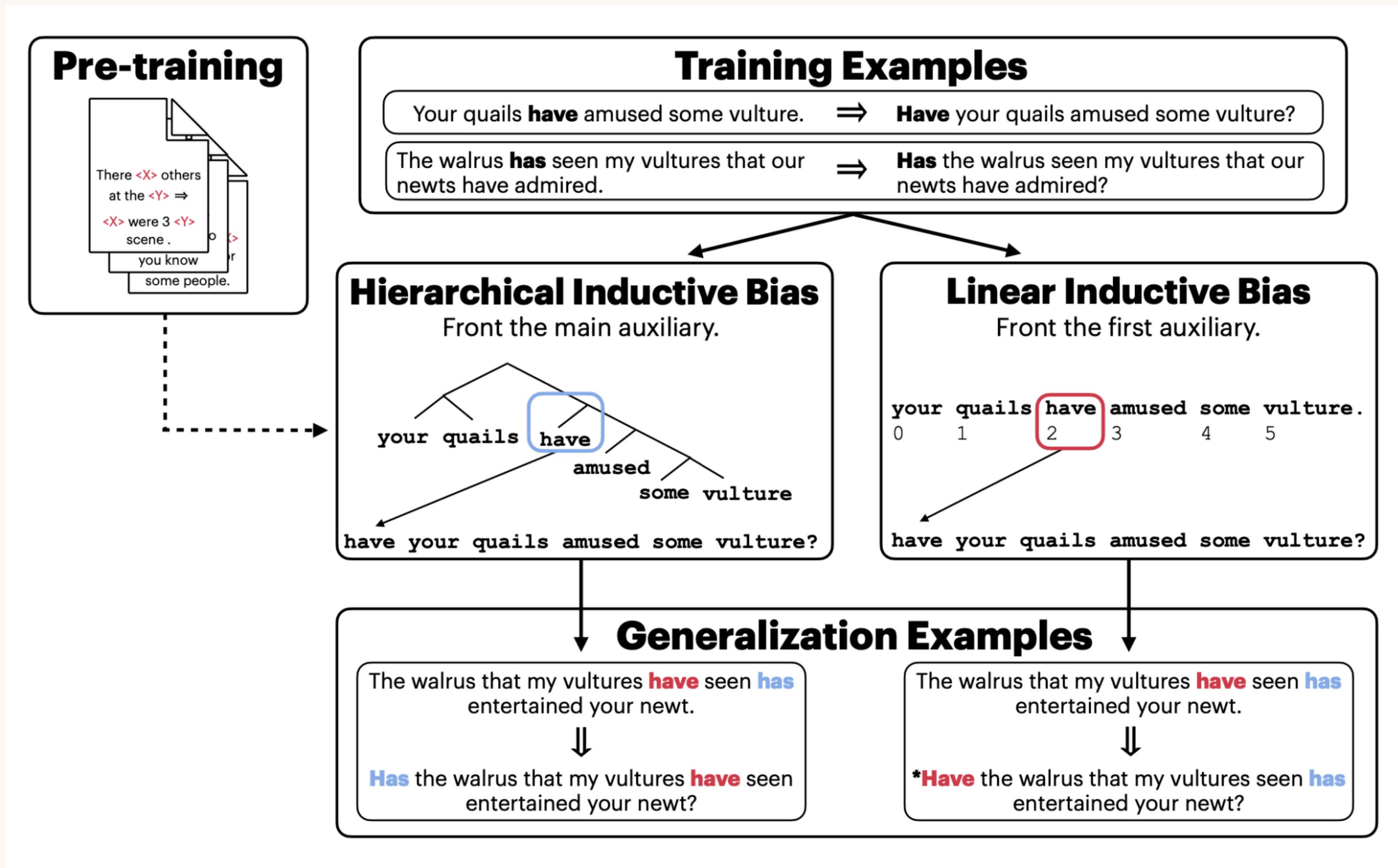
Lexical Generalization: use a known word in a new context

S.3.2. Novel Combination Modified Phrases and Grammatical Roles		
Object modification → Subject modification	Noah ate <b>the cake on the plate</b> .	<b>The cake on the table</b> burned.

Structural Generalization: extend known structures to be more complicated

S.3.3. Deeper Recursion		
Depth generalization: Sentential complements	Emma said <b>that</b> Noah knew <b>that</b> the cat danced.	Emma said <b>that</b> Noah knew <b>that</b> Lucas saw <b>that</b> the cat danced.
Depth generalization: PP modifiers	Ava saw the ball <b>in the bottle on the table</b> .	Ava saw the ball <b>in the bottle on the table on the floor</b> .

# Appendix B: Passivization



# Appendix C: BigBench Tasks

What is 42812 times 51570?

- (a) 8127851667
- (b) 9415523281
- (c) 47549647
- (d) 854486182733
- (e) banana
- (f) house
- (g) 2207814840

What is the common morpheme among these words: biology, antibiotic, symbiosis, amphibian?

- (a) disease
- (b) life
- (c) study
- (d) animal

On which continent would one find the Nile River?

- (a) Africa
- (b) Asia
- (c) South America
- (d) North America
- (e) Europe

A man is offered one wish and he says that he wants to be fire proof. His wish is granted. The same man and a woman are in a house. What will happen if the house lights up on fire? Will the man be hurt by the fire?

- (a) Yes
- (b) No



# Appendix D: Permutation Tests

Could the BB results have arisen due to random chance?

