

How Does Code Pretraining Affect Language Model Task Performance?

Jackson Petty^{1*}, Sjoerd van Steenkiste², and Tal Linzen²

* Work done as a Student Researcher at Google Research



Question: Does pretraining on *source code* help LLMs make more compositional generalizations?

1. What is compositional generalization?

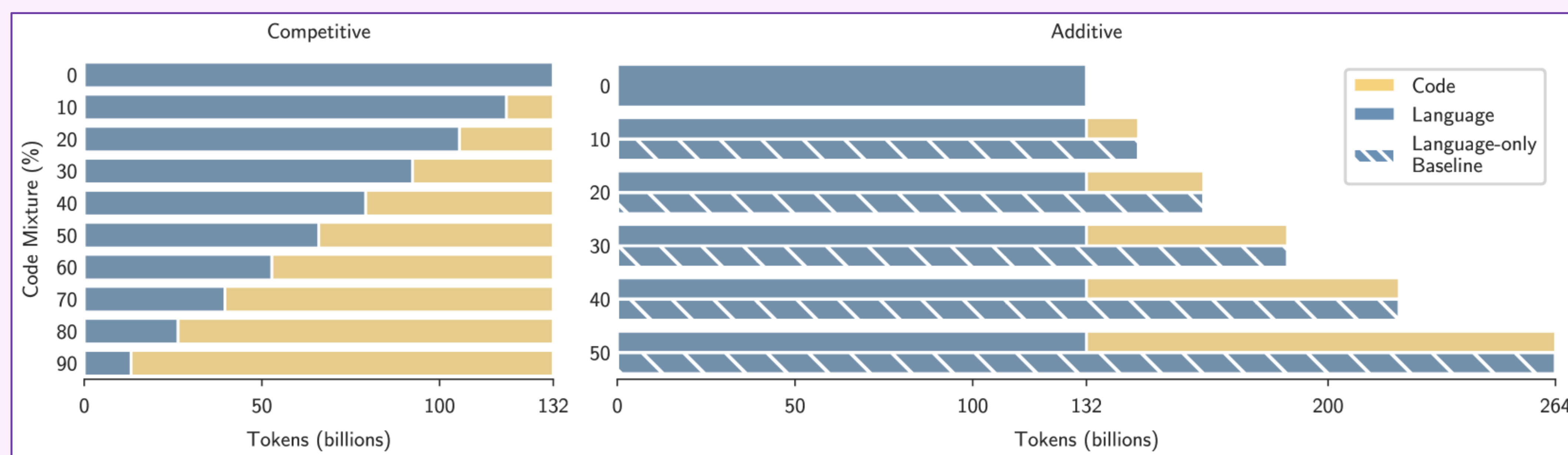
Generalize from **known pieces** to (infinite) **novel, well-formed combinations**

| Training input (<i>hedgehog is subject</i>) | Output |
|---|---|
| the hedgehog ate the cake | <code>eat(agent=hedgehog, theme=cake)</code> |
| the hedgehog saw a child | <code>see(agent=hedgehog, theme=child)</code> |
| hedgehogs swim | <code>swim(agent=hedgehog)</code> |
| → | |
| Generalization (<i>hedgehog is object</i>) | |
| the boy loves the hedgehog | <code>love(agent=boy, theme=hedgehog)</code> |

2. Evaluating compositional generalization

| Dataset | Examples |
|---------------|---|
| COGS | A hedgehog ate the cake . <code>*cake(x4); hedgehog(x1) and eat.agent(x2, x1) and eat.theme(x2, x4)</code> |
| COGS-vf | A hedgehog ate the cake on the bed . <code>eat(agent = hedgehog, theme = *cake(nmod.on = *bed))</code> |
| English | our vultures admired her walrus above some zebra . |
| Passivization | her walrus above some zebra was admired by our vultures . |

3. Experimental setup

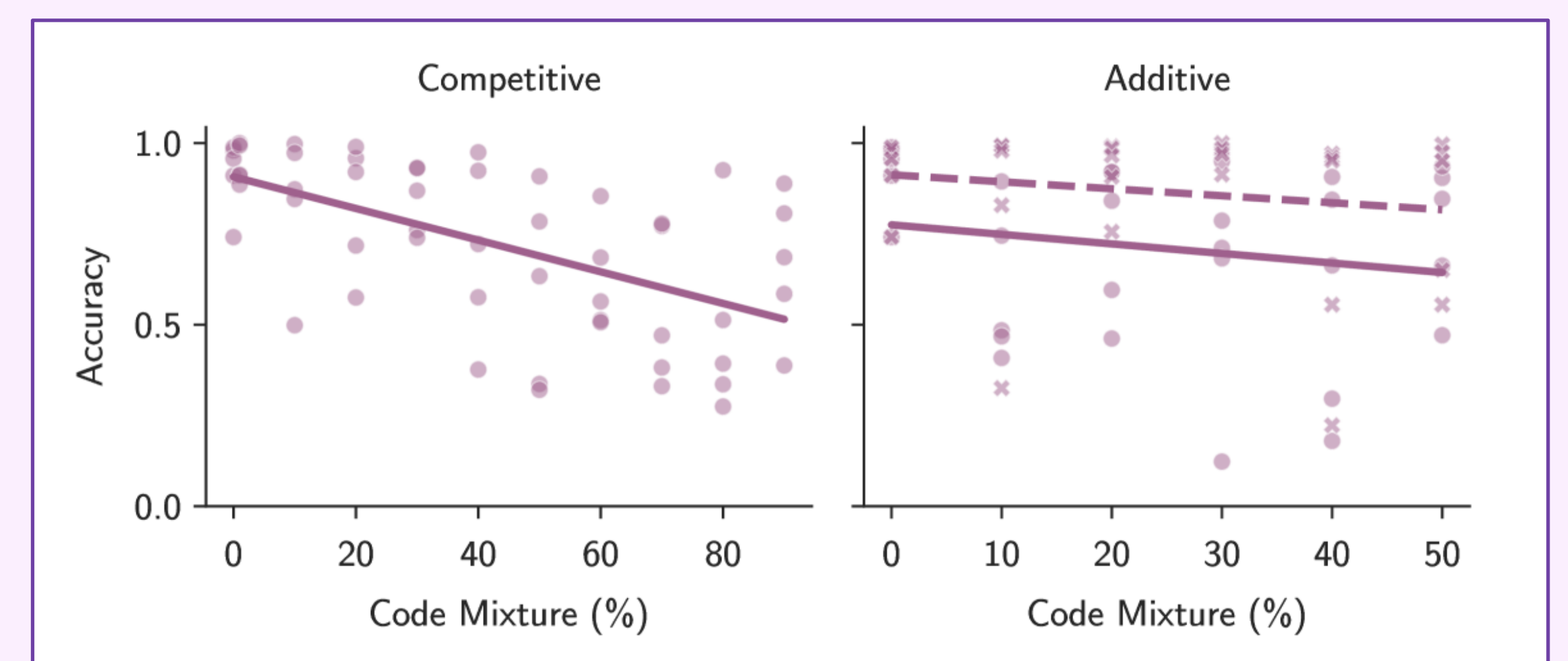
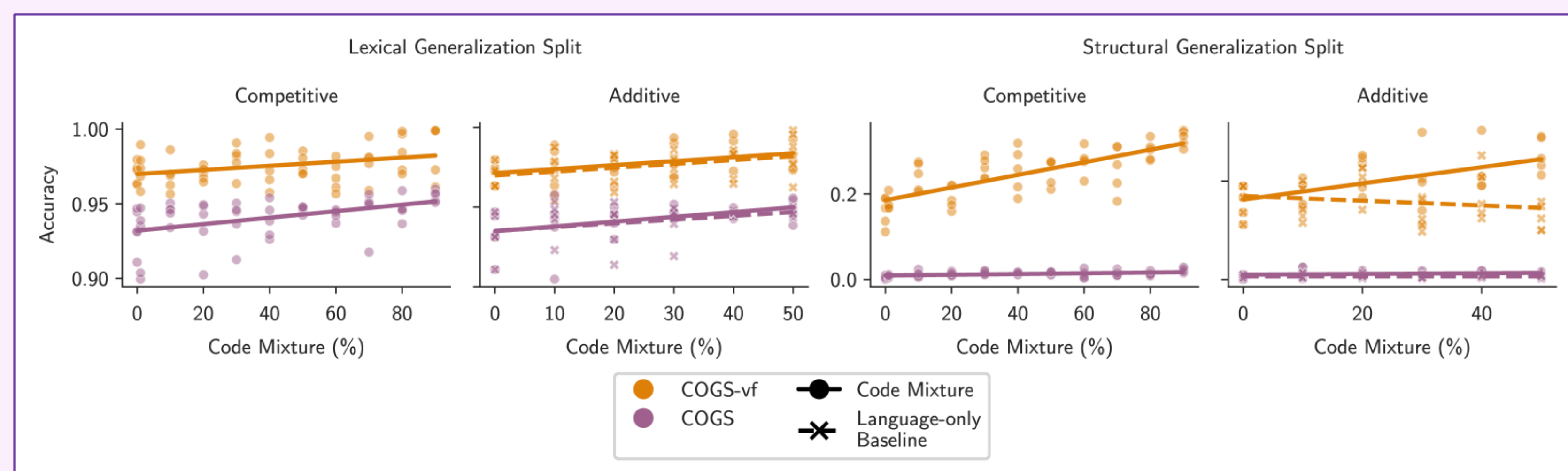


1. Pretrain 400M-parameter LLMs on mixtures of code (GitHub) and language (C4)
2. Finetune on training splits of compositional generalization datasets
3. Evaluate on generalization split

4. Results

Code helps on COGS & COGS-vf (semantic parsing)

but hurts on purely natural-language tasks, like English Passivization



Answer: Yes, depending on the format

Code can help models generalize more compositionally, but only in cases where the output domain has formal structure

