

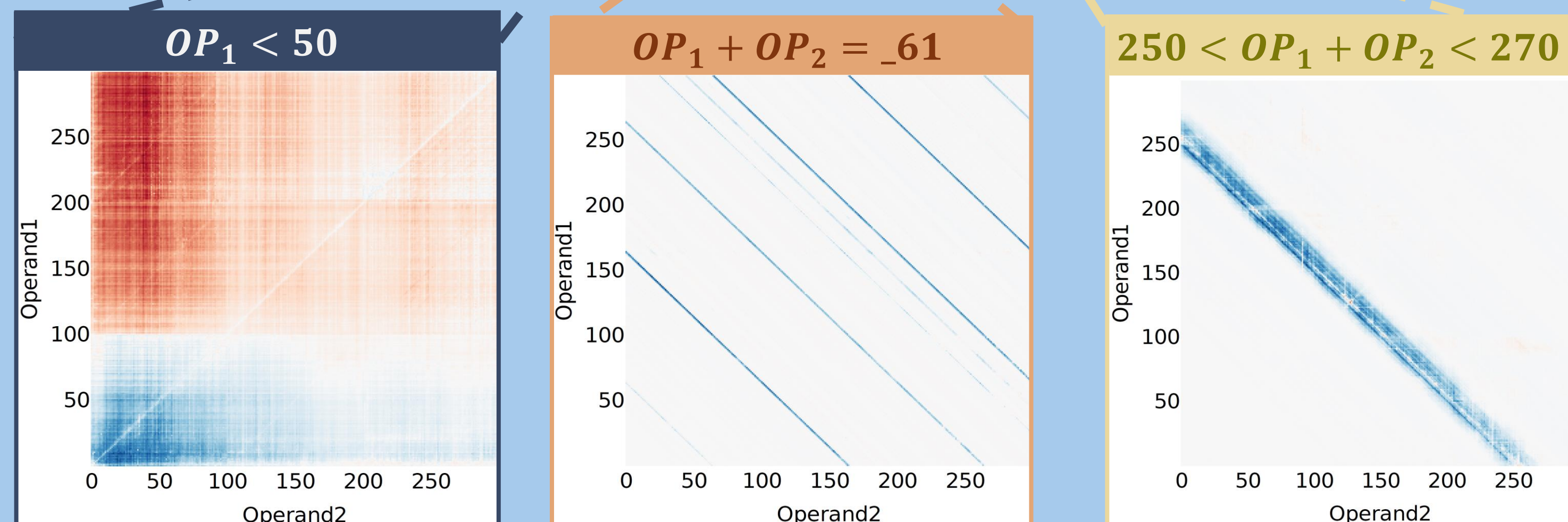
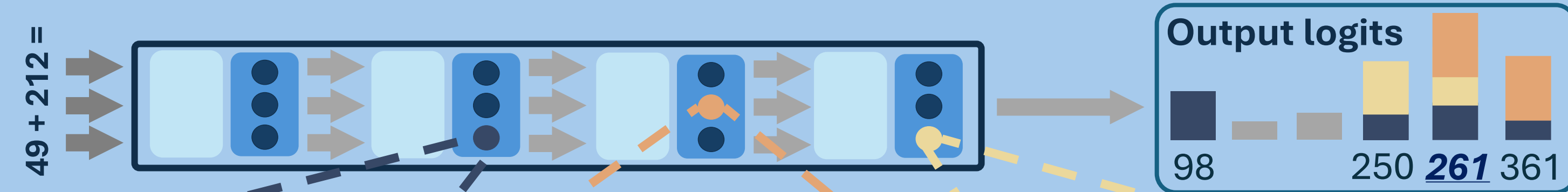
“49 + 212” → LLM → 261, but how?

✗ Perfect Memorization?

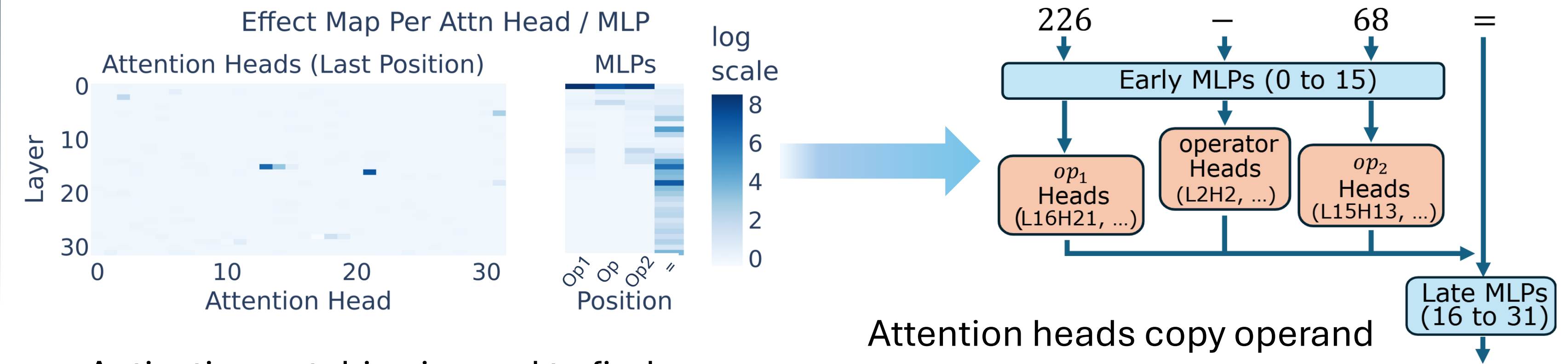
✗ Learned Algorithm?

LLMs Solve Arithmetic with a Bag Of Heuristics

Yaniv Nikankin, Anja Reusch, Aaron Mueller, Yonatan Belinkov



1. Which model components participate in answering arithmetic prompts?

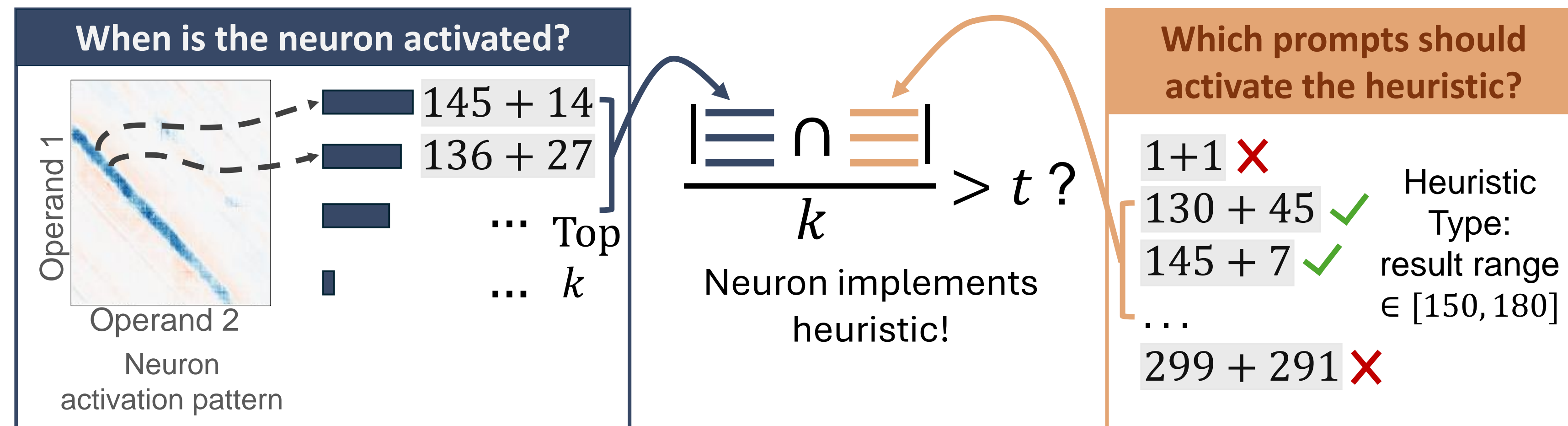


Activation patching is used to find important components, making up the arithmetic circuit.

Attention heads copy operand embeddings to last position.

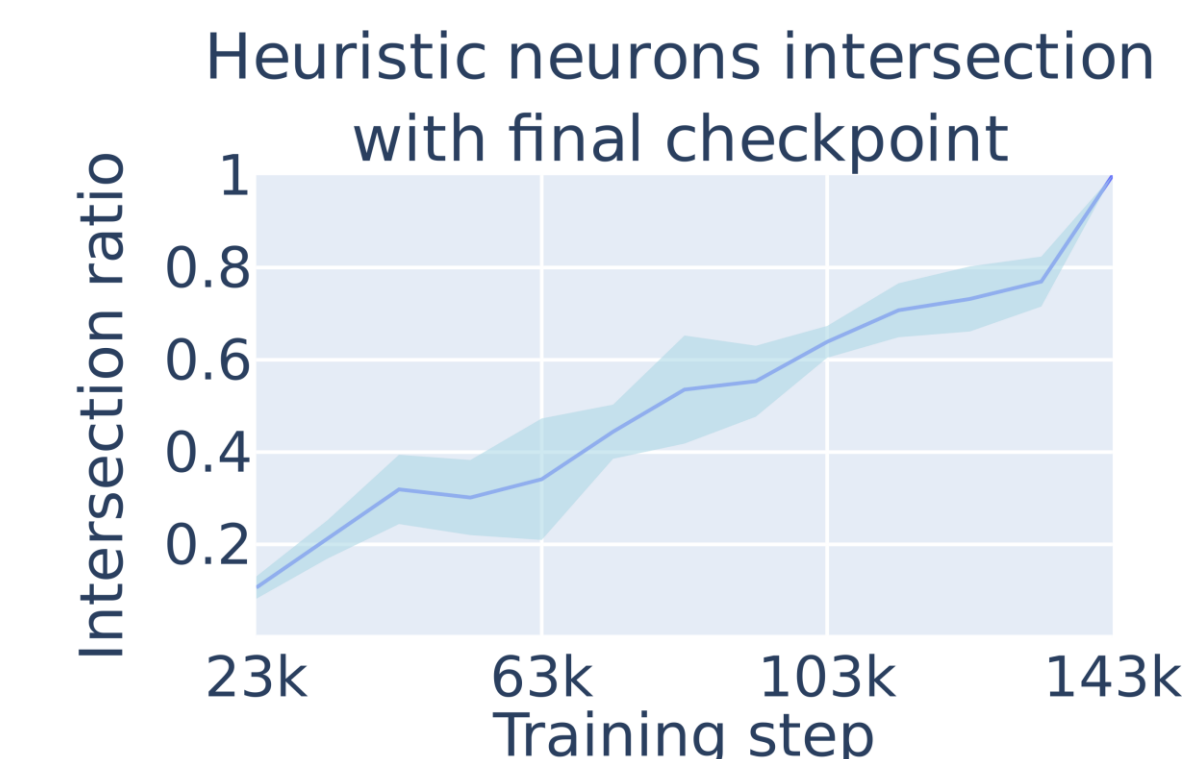
1% of later-MLP (heuristic) neurons generate the correct answer.

2. How do we label heuristic neurons automatically?



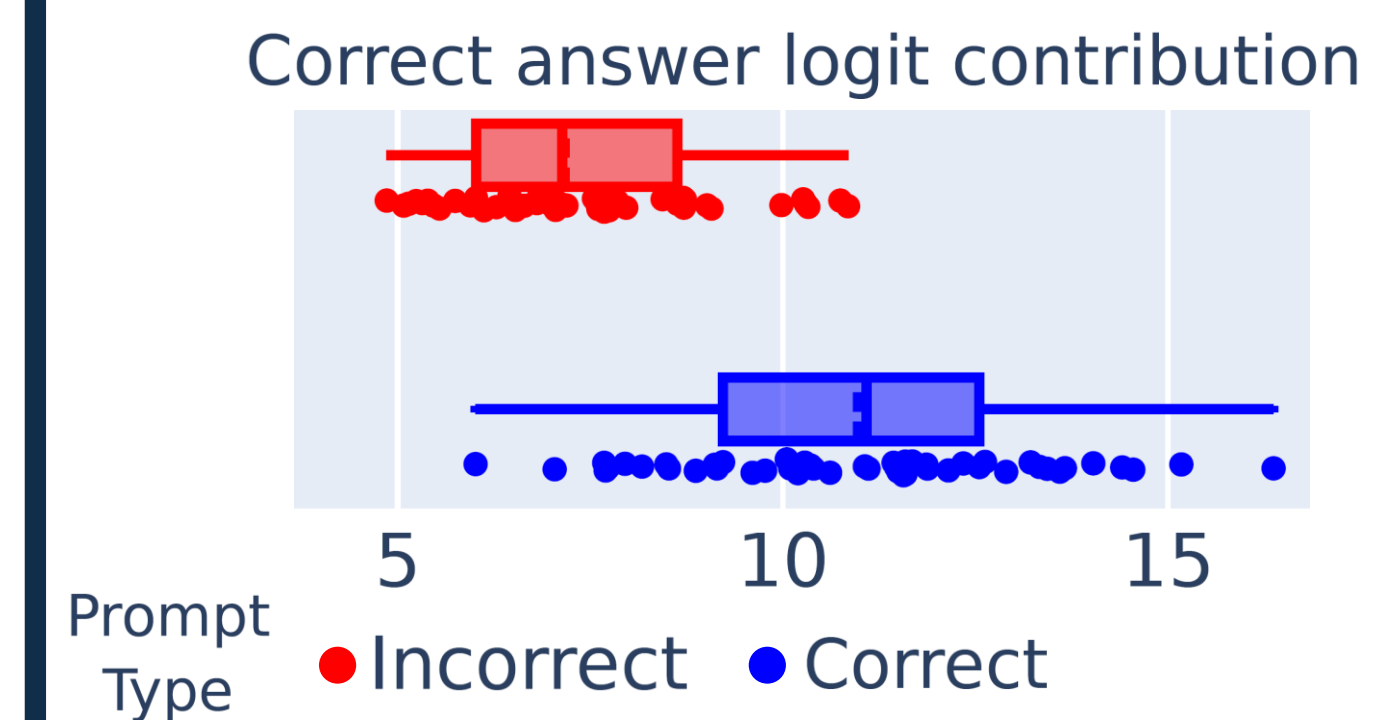
Automatic labeling process allows labeling of circuit neurons to heuristic types.

3. How do arithmetic heuristics develop over training?



The bag of heuristics emerges as the single mechanism driving arithmetic calculations from early training

4. Why do arithmetic heuristics fail?



Prompt Type
● Incorrect ● Correct