

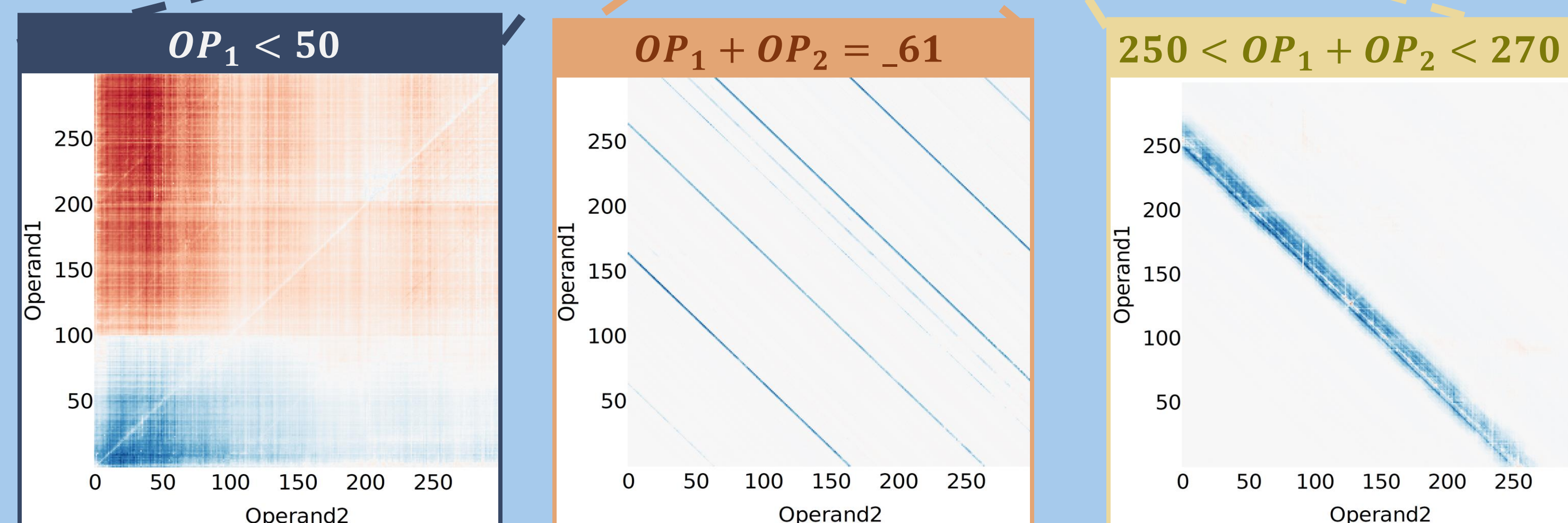
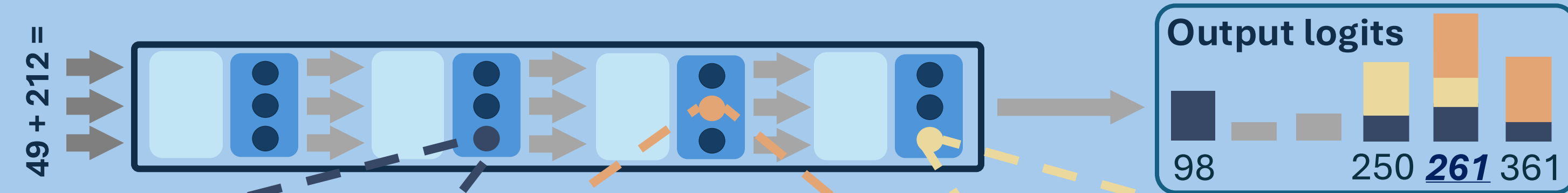
“49 + 212” → LLM → 261, but how?

✗ Perfect Memorization?

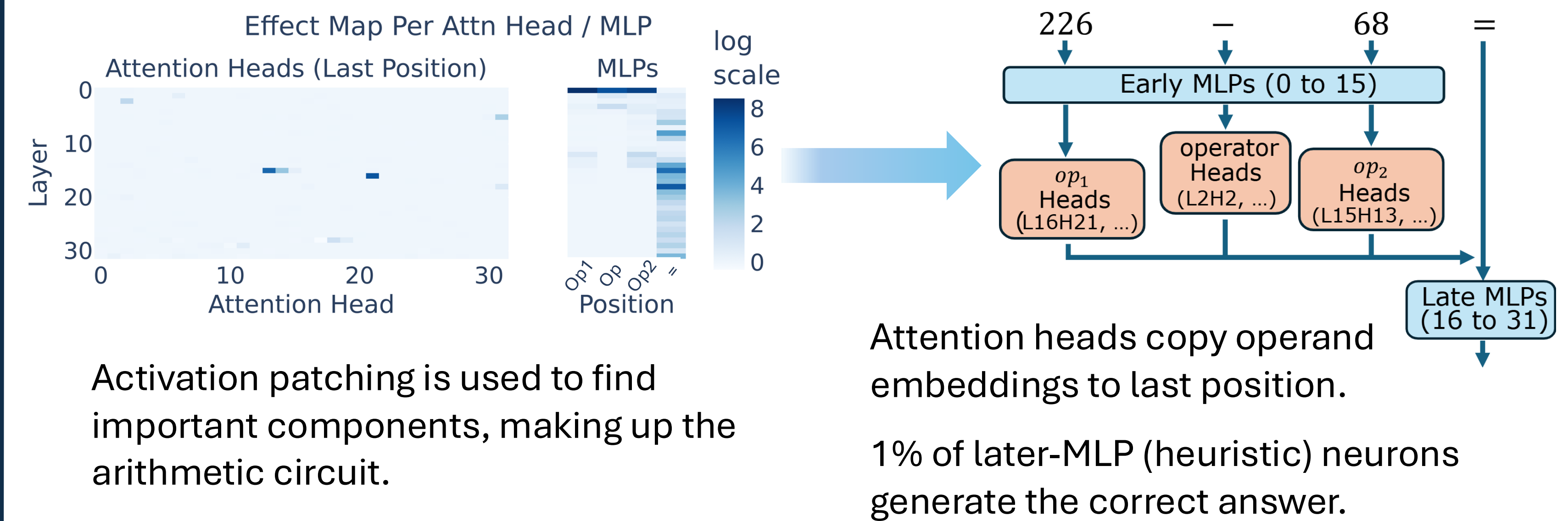
✗ Learned Algorithm?

LLMs Solve Arithmetic with a Bag Of Heuristics

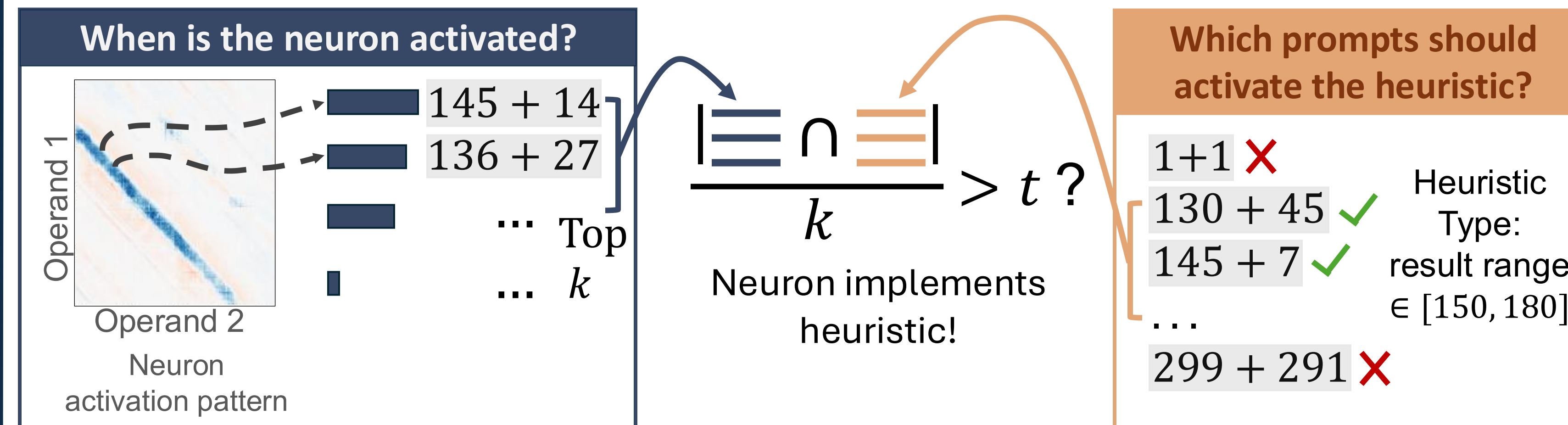
Yaniv Nikankin, Anja Reusch, Aaron Mueller, Yonatan Belinkov



1. Which model components participate in answering arithmetic prompts?

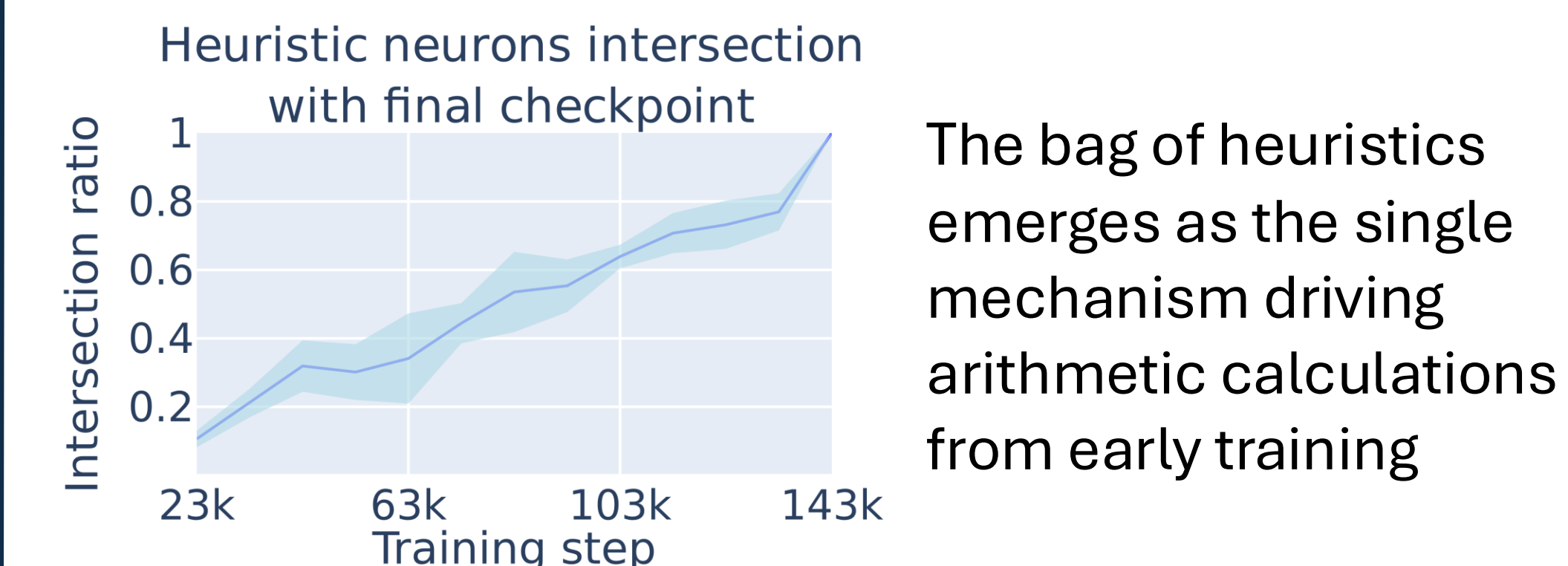


2. How do we label heuristic neurons automatically?



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

3. How do arithmetic heuristics develop over training?



4. Why do arithmetic heuristics fail?

