

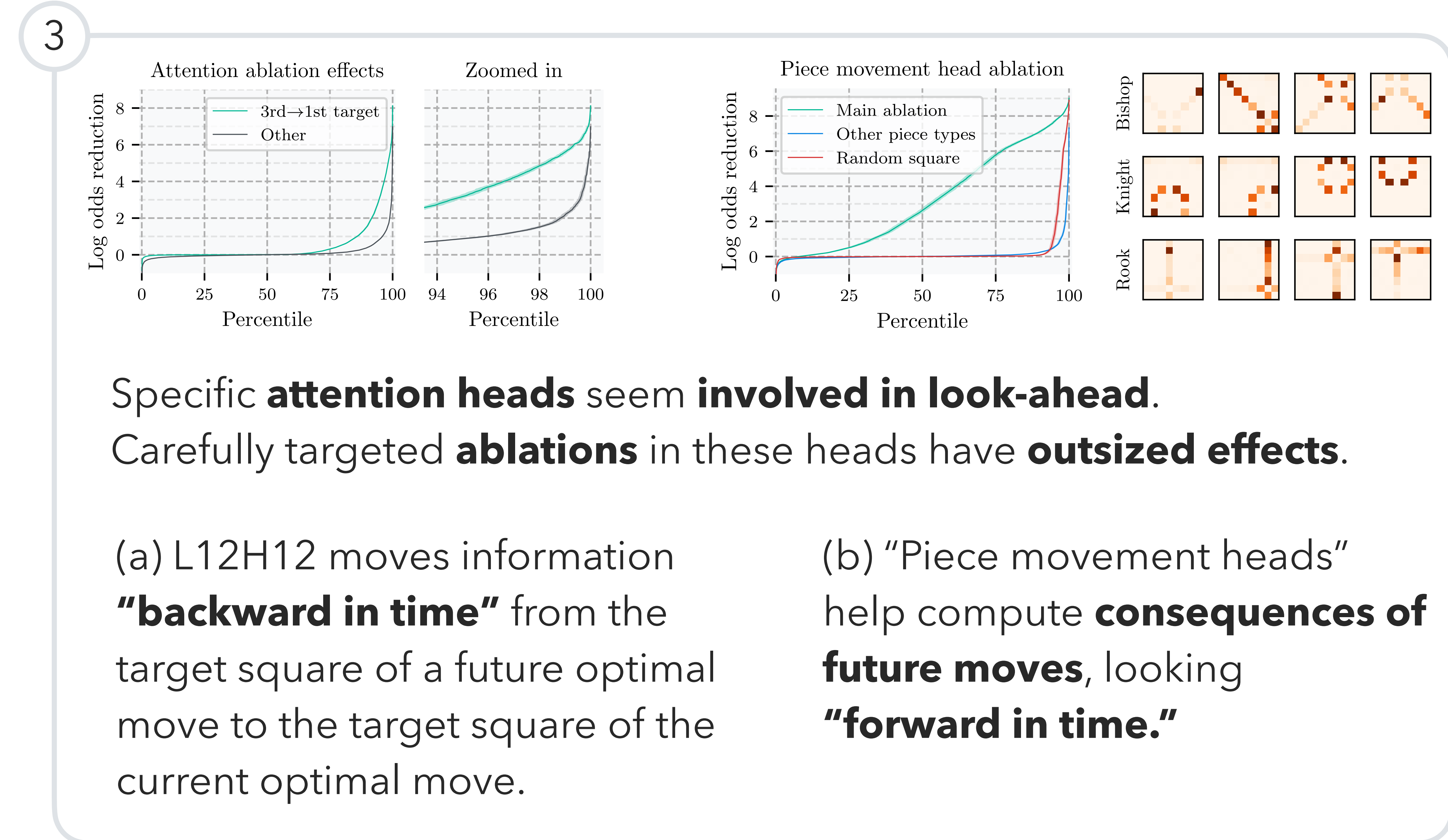
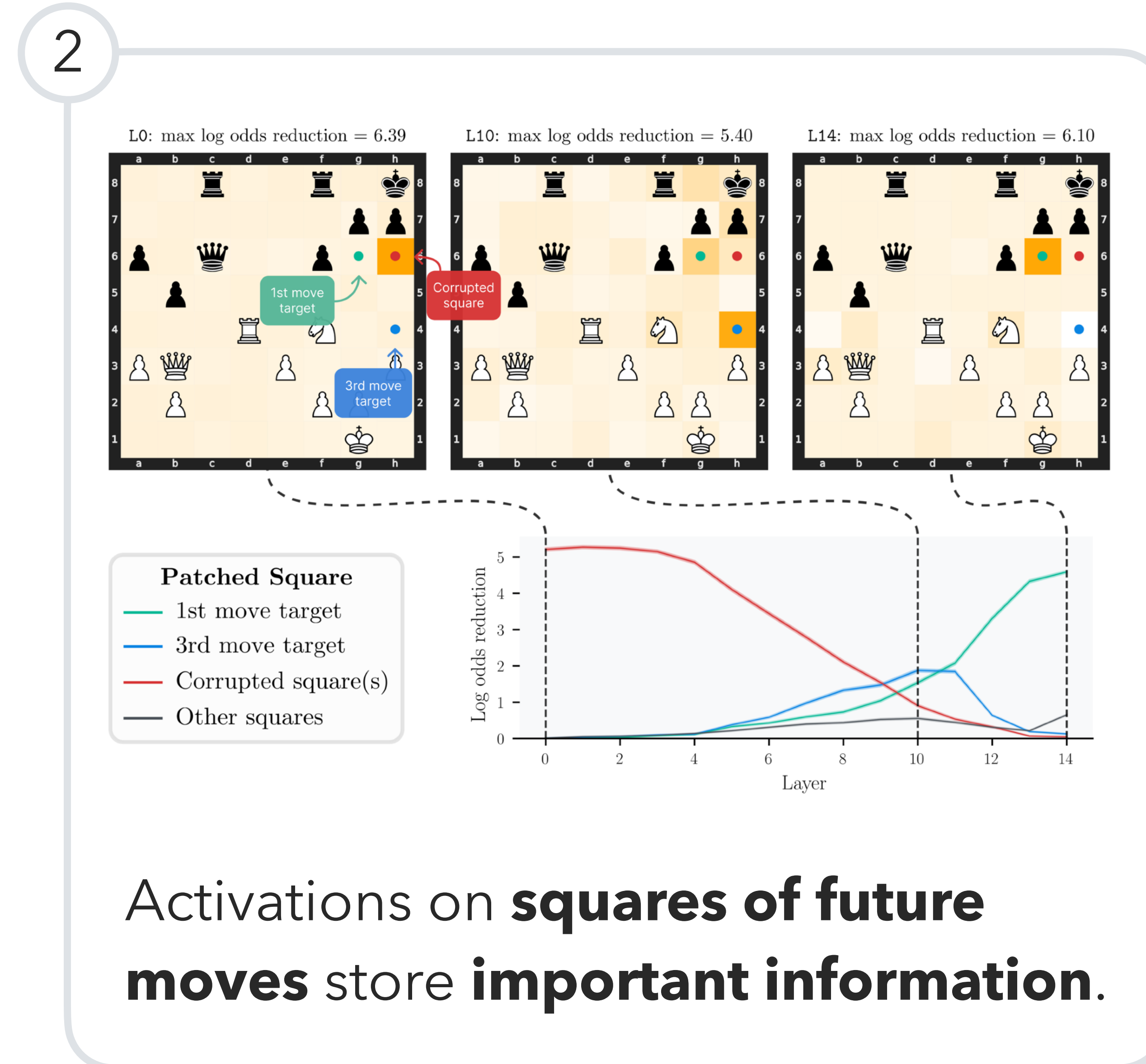
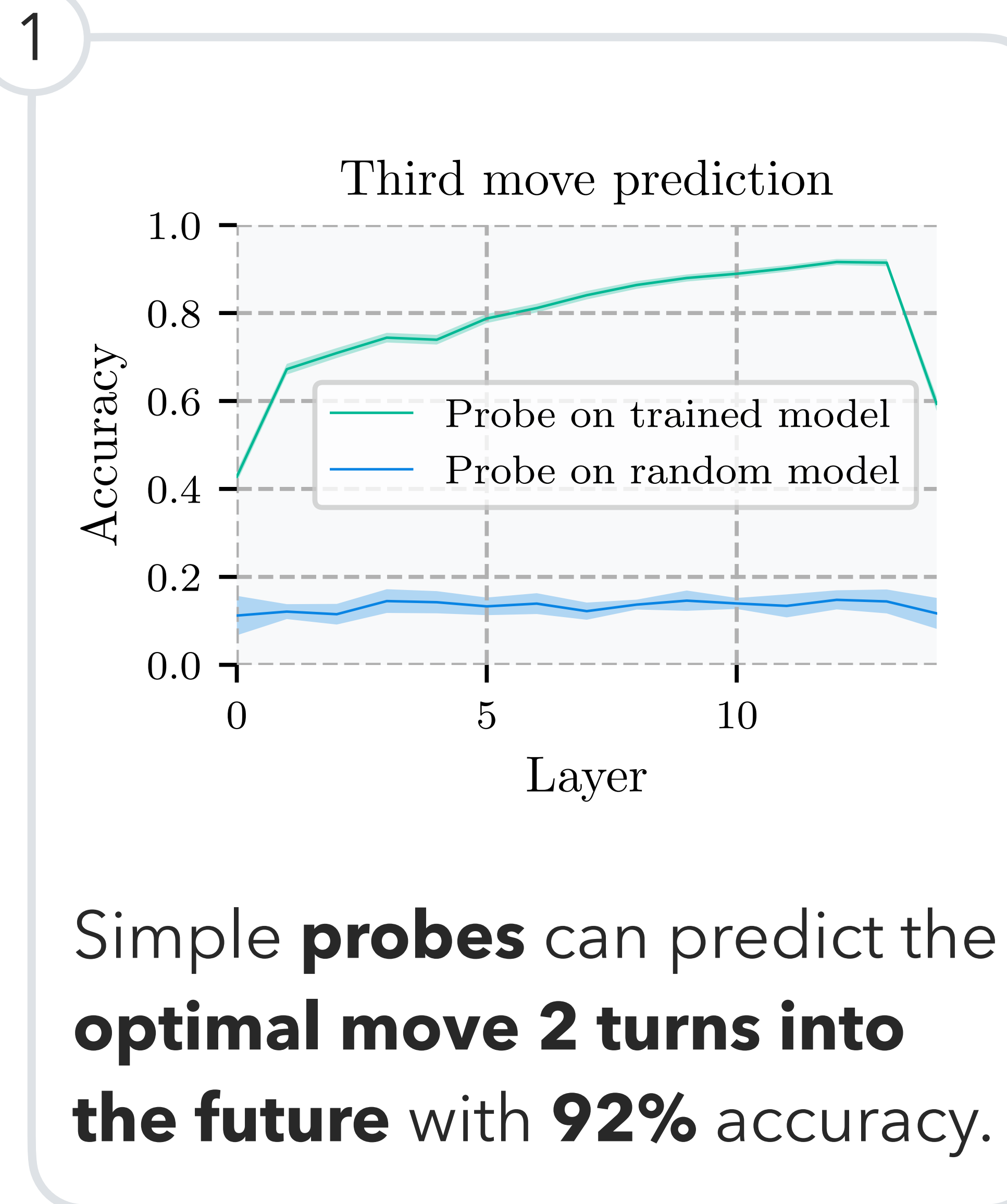
Transformers can learn to implement look-ahead in a single forward pass.

Evidence of Learned Look-Ahead in a Chess-Playing neural network

Erik Jenner, Shreyas Kapur, Vasil Georgiev, Cameron Allen, Scott Emmons, Stuart Russell

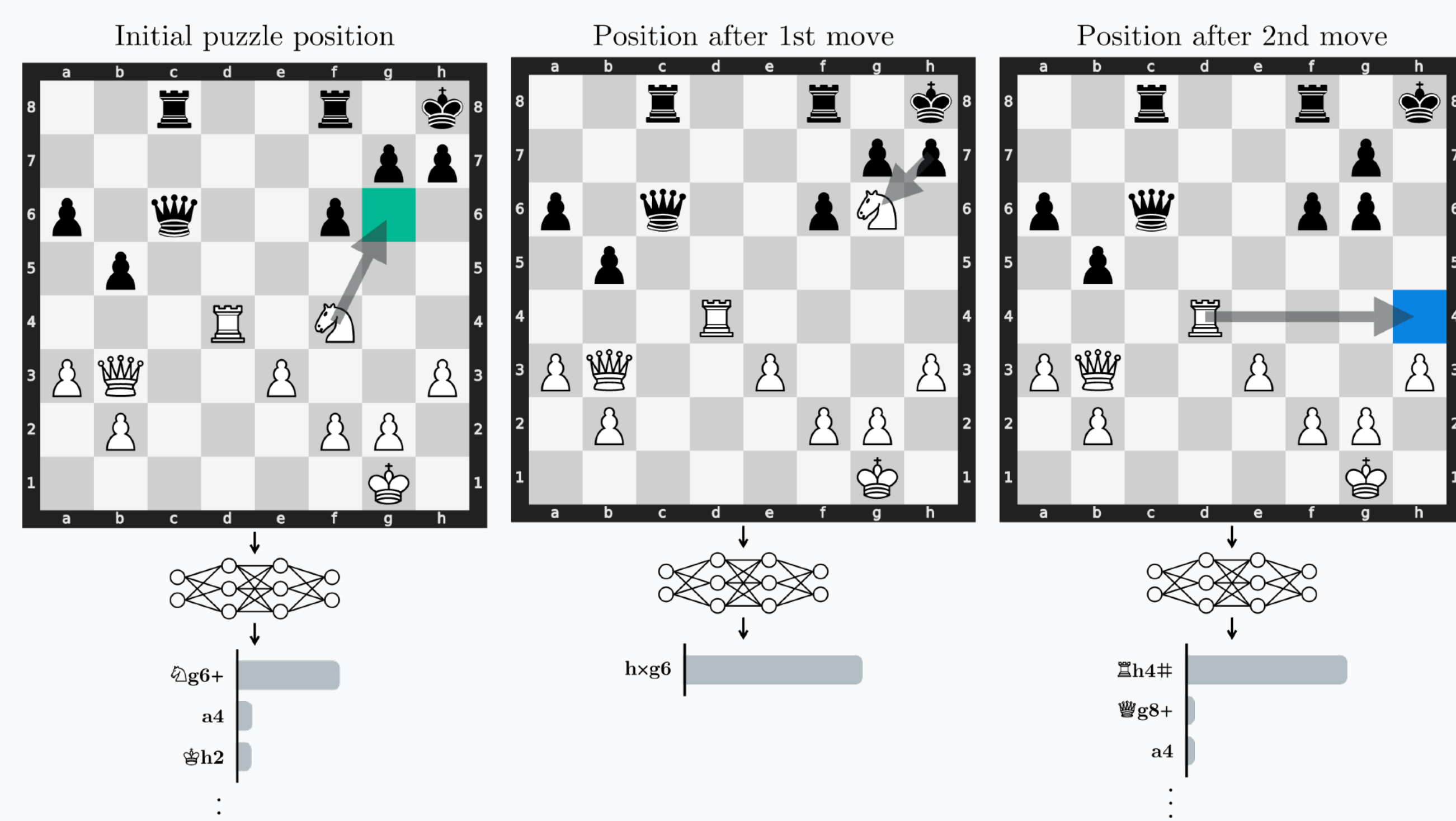


Three lines of evidence



Setup details

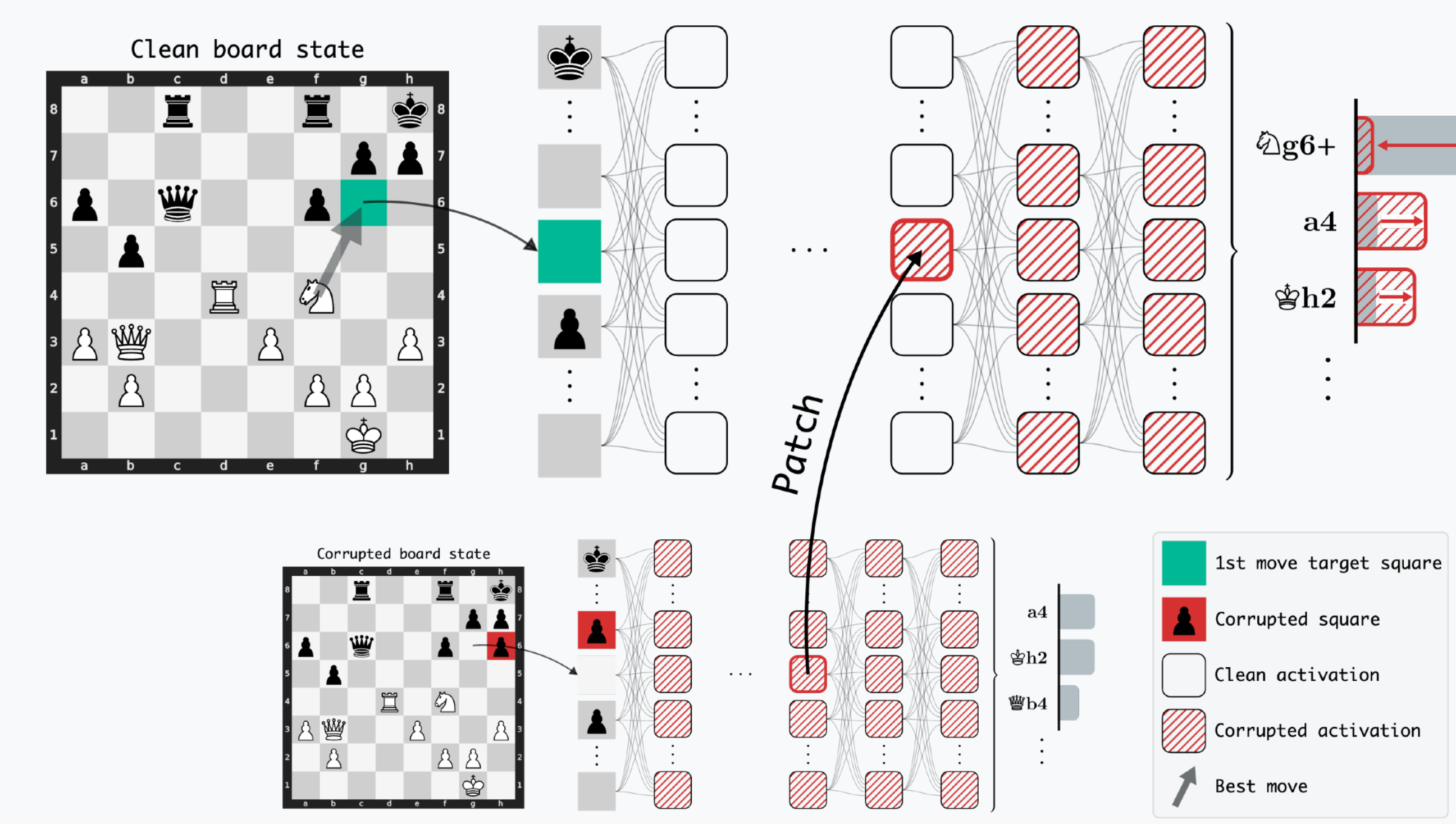
Model and dataset



Dataset: board states from puzzles, automatically filtered to be **complex** but still **solvable** for the network. Annotated with a unique *principal variation* of at least three moves.

Network: policy net of Leela Chess Zero (strongest MCTS engine). **Transformer** that maps a board state to a move distribution. It **treats every square like a token in a language model**.

Activation patching



Mechanistic interpretability method that lets us determine **how important** a given **model component** is. Patch an activation from the forward pass on a “corrupted” state into the forward pass on the original state. Then **measure how much this intervention affects the output**. We use **automatically generated corruptions**.