Conceptual similarity judgements Openchat vs. Mistral-Instruct vs. Gemma-Instruct vs. Llama/Tulu vs. Starling Zephyr-Mistral Zephyr-Gemma Llama-Chat/Tulu-DPO 1.0 none, default 0.5 0.0 1.0 persona 0.5 **Prompt manipulations** 0.0 1.0 random P(multiple concepts) 0.5 0.0 1.0 nonsense 0.5 0.0 1.0 **Temperature manipulations** 1.5 0.5 0.0 1.0 2.0 0.5 0.0 animals politicians animals politicians animals politicians animals politicians **Concept category** Non-aligned model Aligned model Human baseline (animals) · · · · Human baseline (politicians)