Phi 4 Multimodal Gemma 3 4B it Llama 4 Maverick Owen2.5 VL 32B demma 3 4B it 1.2 better Phi 4 Multimodal Gemini 2.5 Flash Lite Observed NRMSE (mean across ablations) 1.0 -Gemini 2.5 Flash Lite Owen2.5 VL 32B Mistral 24B Phi 4 Multimodal GPT 4.1 mini Qwen2.5 VL 32B 8.0 Gemini 2.5 Flash Lite GPT 4o Llama 4 Maverick GPT 4.1 mini Mistral 24B Llama 4 Maverick Mistral 24B GPT 4.1 mini 0.6 GPT 4o Claude 3.7 Sonnet GPT 4o Modality 0.4 image Claude 3.7 Sonnet text Claude 3.7 Sonnet text+image 0.7 0.5 0.6 0.8 0.9 0.2 0.3 0.4 1.0

Bayesian-ness (mean probability across ablations)

NRMSE vs Bayesian-ness (one dot per model × modality)