NRMSE vs Bayesian-ness (one dot per model × modality) Modality Phi 4 Multimodal image text text+image better 1.2 -Gemma 3 4B it Gemma 3 4B it Observed NRMSE (mean across ablations) 1.0 Gemma 3 4B it Phi 4 Multimodal Phi 4 Multimodal Llama 4 Maverick 8.0 Mistral 24B Qwen2.5 VL 32B Gemini 2.5 Flash Lite Llama 4 Maverick Mistral 24B Gemini 2.5 Flash Lite Mistral 24B Llama 4 Maverick 0.6 GPT delt apjini Qwen2.5 VL 32B GPT 4o Claude 3.7 Sonnet Qwen2.5 VL 32B GPT 4o Gemini 2.5 Flash Lite 0.4 Claude 3.7GPan4et 60PTi 4.1 mini Claude 3.7 Sonnet 0.2 0.6 8.0 0.4 1.0

Bayesian-ness (mean probability across ablations)