NRMSE vs Bayesian-ness (one dot per model × modality) 3.5 Modality Gemini 2.5 Flash Lite image text 3.0 text+image ↓ better 2.5 Observed NRMSE (mean across ablations) 2.0 1.5 -Phi 4 Multimodal Gemma 3 dBeh2.5 VL 32B Gemma 3 4B it Phi 4 Multimodal Gemini 2.5 Flash Lite Llama 4 Maverick 1.0 -Phi 4 Multimodal Mistral 24B Qwen2.5 VL 32B Qwen 25 VL 32B stral 24B Gemini 2.5 Flash Lite GPT 4o Mistral 24B LHama44Maverick GPT 4o Classe 3.9 Sonnet 0.5 GPT 5 Mini Claude 3.7 Sonnet Claude 3.7 Sonnet GPT 5 Mini prior GPT 5 Mini 0.0 0.5 0.3 0.4 0.6 0.7 0.9 0.2 8.0 1.0

Bayesian-ness (mean probability across ablations)