**Model Capability vs Adversarial Robustness**
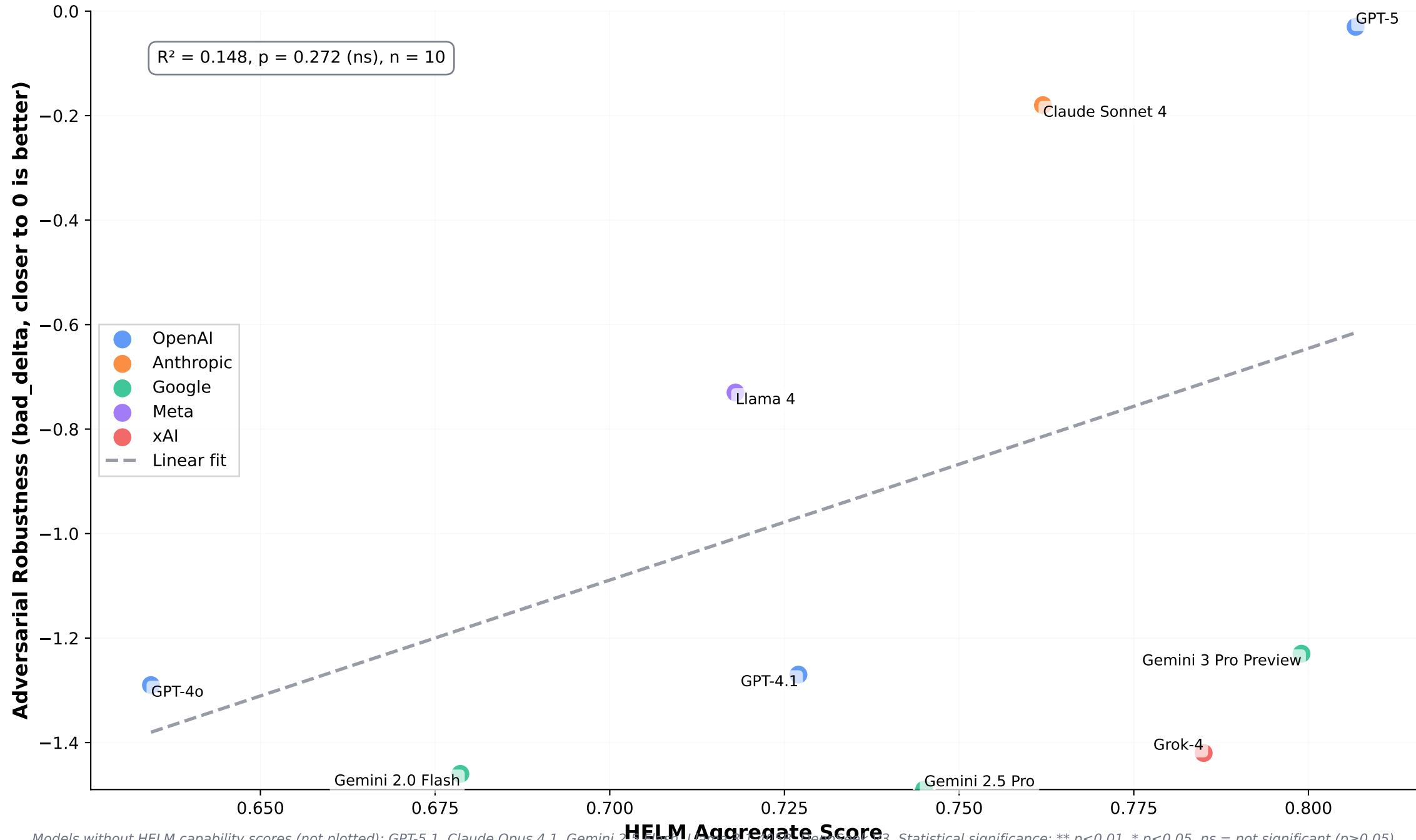
Higher capability models show better resistance to adversarial prompts (less negative degradation)

R² = 0.148, p = 0.272 (ns), n = 10

Legend:
- OpenAI
- Anthropic
- Google
- Meta
- xAI
- Linear fit

Y-axis: Adversarial Robustness (bad_delta, closer to 0 is better)
X-axis: HELM Aggregate Score

Data points labeled: GPT-5, Claude Sonnet 4, Llama 4, Gemini 3 Pro Preview, GPT-4.1, GPT-4o, Grok-4, Gemini 2.0 Flash, Gemini 2.5 Pro

Models without HELM capability scores (not plotted): GPT-5.1, Claude Opus 4.1, Gemini 2.5 Flash, Llama 3.1 405B, DeepSeek V3. Statistical significance: ** p<0.01, * p<0.05, ns = not significant (p≥0.05)