**Mechanistic Abstention Control Supports Trustworthiness Objectives**

**(A) Cross-Domain Generalization**

Rate

- Coverage (Answerable)
- Abstention (Unanswerable)

Math, Science, History, Geography, General

**(B) Risk-Aware Behavior**

Coverage

Gap

- Low Risk
- High Risk

Steering Strength ($\epsilon$)

**(C) Direction Separation**

**Cosine Similarity**

0.179

**Angle**

79.7°

*(Near-orthogonal)*