# Uncertainty Expression Gate: Complete Implementation

**For ICLR Workshop: Principled Design for Trustworthy AI**

This repository contains complete, working code for all 5 experiments demonstrating that instruction-tuned language models contain a policy-responsive gate that controls whether internal uncertainty is expressed as abstention.

## 🎯 Key Contribution

**We prove that abstention is governed by a late-layer "expression gate" that:**

1. Responds to instruction regime (not just question properties)

2. Operates independently of internal uncertainty

3. Is causally manipulable via targeted interventions

4. Enables tuning trustworthiness tradeoffs

## 📁 File Structure

```
.
├── core_utils.py              # Model wrapper & utilities
├── data_preparation.py        # Dataset creation
├── experiment1_behavior_belief.py  # Exp 1: Policy-responsiveness
├── experiment2_localization.py     # Exp 2: Layer localization
├── experiment3_4_steering_independence.py  # Exp 3-4: Steering & key proof
├── experiment5_trustworthiness.py  # Exp 5: Application + master runner
├── data/                   # Generated datasets (auto-created)
└── results/                # Output figures & CSVs (auto-created)
```

## 🚀 Quick Start

### 1. Installation

```bash
```

```
# Create environment
conda create -n mech_interp python=3.10
conda activate mech_interp

# Install dependencies
pip install torch transformers numpy pandas matplotlib seaborn scikit-learn tqdm
```

## 2. Prepare Data

```bash
python data_preparation.py
```

This creates:

- `data/dataset_ambiguous.json` (30 ambiguous questions)
- `data/dataset_clearly_answerable.json` (10 factual questions)
- `data/dataset_clearly_unanswerable.json` (10 unanswerable questions)
- `data/dataset_squad.json` (synthetic SQuAD-style examples)

**Optional:** Replace with real SQuAD v2.0 data by downloading from https://rajpurkar.github.io/SQuAD-explorer/ and updating `load_squad_data()` in `data_preparation.py`.

## 3. Run Experiments

### Option A: Quick Test (recommended first)

```bash
python experiment5_trustworthiness.py
# This runs quick_test() with minimal data (~15 min on GPU)
```

### Option B: Run Individual Experiments

```bash
```

```
# Experiment 1: ~30 min
python experiment1_behavior_belief.py

# Experiment 2: ~45 min (tests all layers)
python experiment2_localization.py

# Experiments 3 & 4: ~1 hour
python -c "from experiment3_4_steering_independence import main; main()"

# Experiment 5: ~20 min
python -c "from experiment5_trustworthiness import main; main()"
```

**Option C: Run All Experiments (full pipeline)**

```bash
python -c "from experiment5_trustworthiness import run_all_experiments; run_all_experiments(quick_test=False)"
# Full run: ~3-4 hours on GPU
```

## 4. View Results

All outputs saved to `results/`:

- `exp1_behavior_belief_dissociation.png` - Instruction regime effects
- `exp2_localization_analysis.png` - Layer-wise localization
- `exp3_steering_analysis.png` - Steering control
- `exp4_gate_independence.png` - ⭐ **KEY FIGURE: Gate independence proof**
- `exp5_trustworthiness.png` - Risk-coverage tradeoff
- `master_summary.json` - All quantitative results

## 📊 Experiments Overview

### Experiment 1: Behavior ≠ Belief

**Shows:** Instruction regime dramatically changes abstention behavior while internal uncertainty stays constant.

**Key metric:** Behavior gap (cautious - confident) vs Belief gap (entropy difference)

### Experiment 2: Gate Localization

**Shows:** Abstention decision is causally determined in a late-layer band (typically layers 17+).

**Method:** Activation patching from answerable → unanswerable examples

### Experiment 3: Low-Dimensional Control

**Shows:** A single learned direction can steer abstention behavior reliably.

**Method:** Compute answerability direction via mean difference, test steering

### Experiment 4: Gate Independence ⭐ CRITICAL

**Shows:** Gate operates independently of uncertainty - can force high-uncertainty questions to be answered, and low-uncertainty questions to be abstained.

**This is your novel contribution!**

### Experiment 5: Trustworthiness Application

**Shows:** Intervening on the gate enables tuning the risk-coverage tradeoff.

**Method:** Selective accuracy evaluation across steering strengths

## 🔧 Customization

### Using Your Own Data

Replace datasets in `data_preparation.py`:

```python
def create_your_custom_dataset():
    return [
        {
            "question": "Your question here",
            "context": "Optional context",
            "answer": "Ground truth" or None,
            "answerability": "answerable" or "unanswerable"
        },
        # ... more examples
    ]
```

### Adjusting Compute Requirements

### For faster testing:

```python
```

```
# In experiment scripts
config.n_force_guess_samples = 10  # Reduce from 20
layer_stride = 2  # Test every other layer in Exp 2
```

**For better results:**

```python
config.n_force_guess_samples = 30  # More samples = better entropy estimates
# Test more example pairs in Exp 2, 3, 4
```

**Different Models**

Change model in `core_utils.py`:

```python
model_name: str = "Qwen/Qwen2.5-1.5B-Instruct"  # Current
# Try:
# "Qwen/Qwen2.5-7B-Instruct"
# "meta-llama/Llama-2-7b-chat-hf"
# etc.
```

**Note:** You may need to adjust `config.target_layers` based on model depth.

## 📈 Expected Results

### Experiment 1

- Abstention rate should vary 30-60% between cautious and confident regimes

- Internal entropy should vary <10% across regimes

- Behavior/Belief ratio should be >3x

### Experiment 2

- Clear concentration of causal effect in late layers (typically last 20-30%)

- Flip rates should exceed 50% in critical layers

### Experiment 3

- Flip rates should increase monotonically with epsilon

- Layer ~26 typically most effective (for 28-layer model)

## Experiment 4 ⭐

- Should observe cases where:

  - High uncertainty → forced to answer (gate opens)

  - Low uncertainty → forced to abstain (gate closes)

- This proves gate independence!

## Experiment 5

- Risk-coverage curve should show tradeoff

- Higher epsilon → lower coverage but higher accuracy

# 🐛 Troubleshooting

## Out of Memory

```python
# Reduce batch size / samples
config.n_force_guess_samples = 10
# Use smaller model
model_name = "Qwen/Qwen2.5-1.5B-Instruct"  # instead of 7B
```

## Model Loading Issues

```bash
# Make sure you have access to model
huggingface-cli login
# Or use different model that doesn't require auth
```

## Incorrect Behavior Detection

The `extract_answer()` function uses simple heuristics. For better detection:

```python
```

```
    # Add more uncertainty markers
    uncertainty_markers = [
        "UNCERTAIN", "I don't know", "cannot answer",
        "can't determine", "insufficient information",
        # ... add markers specific to your model
    ]
```

## Layer Index Errors

Different model architectures have different layer counts. Check:

```python
n_layers = model.model.config.num_hidden_layers
config.target_layers = [n_layers-4, n_layers-3, n_layers-2, n_layers-1]
```

## 📝 Writing the Paper

### Key Claims to Make

1. **Abstention is governed by a policy-responsive gate**

   - Evidence: Exp 1 (behavior changes, belief doesn't)

2. **Gate operates in late layers**

   - Evidence: Exp 2 (activation patching localization)

3. **Gate is low-dimensional and manipulable**

   - Evidence: Exp 3 (steering control)

4. **Gate operates independently of uncertainty** ⭐

   - Evidence: Exp 4 (changes abstention at fixed uncertainty)

5. **Practical trustworthiness implications**

   - Evidence: Exp 5 (risk-coverage tuning)


### Figures for Paper

### Main text (4 figures max):

1. Fig 1: Exp 1 behavior-belief dissociation

2. Fig 2: Exp 2 localization (patching effects by layer)

3. Fig 3: Exp 4 gate independence ⭐ (this is your money figure!)

4. Fig 4: Exp 5 risk-coverage curve

**Supplement:**

- Exp 3 steering details

- Additional ablations

- Failure cases

**One-Sentence Summary**

"We identify a late, low-dimensional policy gate in instruction-tuned language models that controls whether internal uncertainty is expressed as abstention, and show that intervening on this gate changes abstention behavior at fixed internal answerability."

## 🎓 Citation

If you use this code, please cite:

```bibtex
@inproceedings{yourname2025uncertainty,
  title={Uncertainty Expression Gate: Mechanistic Control of Abstention in Language Models},
  author={Your Name},
  booktitle={ICLR Workshop on Principled Design for Trustworthy AI},
  year={2025}
}
```

## 📧 Support

For questions or issues:

1. Check troubleshooting section above

2. Review code comments

3. Open an issue with:

    - Error message

    - Which experiment

    - Your config settings

## ⏱ Timeline for Jan 30 Deadline

- **Week 1 (Jan 2-8):** Run all experiments, debug issues

- **Week 2 (Jan 9-15):** Analyze results, create figures

- **Week 3 (Jan 16-22):** Write draft, iterate on experiments if needed

- **Week 4 (Jan 23-30):** Finalize paper, submission

Good luck! You've got this! 🚀