# Value Imprint: A Technique for Auditing Human Values Embedded in RLHF Datasets

Ike Obi, Rohan Pant, Srishti Shekhar Agrawal,

Maham Ghazanfar, Aaron Basiletti

**Purdue University**

**West Lafayette, IN**

**USA**

# Contributions

+ ## A technique for auditing human values in RLHF Datasets

We introduced a technique for auditing and classifying the underlying human values embedded within RLHF preferences

+ ## Foreground human value distribution & imbalance

Our three case study experiments showed that Wisdom/Knowledge and Information Seeking were the most dominant human values

+ ## We contribute our Value Imprint datasets

We contribute both our ground truth annotation and classification datasets. Thus, providing researchers with the pathway to take this work forward

# Motivation

## 01

Reinforcement Learning from Human Feedback (RLHF) have become a popular way of aligning LLMs with human values and preferences

## 02

At present there is no technical approach for measuring the specific kinds of human values and preferences operationalized via RLHF

## 03

And there is a growing concern among members of the public on the anti-democratic stance of several LLMs

# Research Questions

**RQ1:** What kinds of human values are embedded in RLHF datasets?

**RQ2:** In what ways do the human values embedded within the Anthropic/hh-rlhf, OpenAI WebGPT Comparisons, and Alpaca GPT-4-LLM datasets differ?
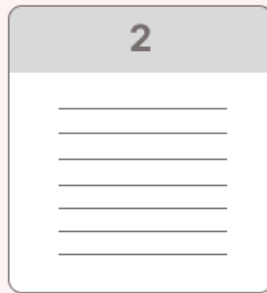
# Research Methods



**Data Collection**

1

Anthropic

WebGPT

Alpaca GPT-4

Collect RLHF datasets from Hugging Face & GitHub
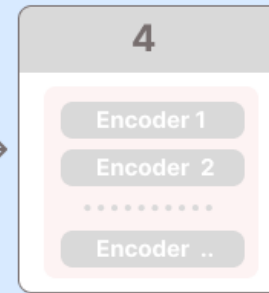
**Taxonomy**

2

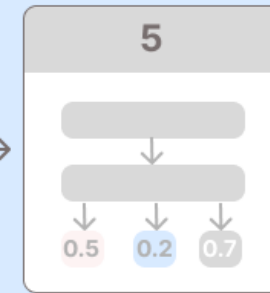Develop human values taxonomy based on prior research in axiology

**Data Annotation**

3

Create ground truth data via qualitative annotation of select RLHF datasets

**Model Training**

4

Encoder 1

Encoder 2

Encoder ..

Train Bert-based models using ground truth data derived from annotation

**Value Classifier**

5

0.5  0.2  0.7

Use trained model to classify human values embedded in RLHF data

**Value Audit**

6

Compare classification results of the 3 datasets to examine insights.

First step – Foundational work

# Research Methods



**Data Collection**

1

Anthropic

WebGPT

Alpaca GPT-4

Collect RLHF datasets from Hugging Face & GitHub

**Taxonomy**

2

Develop human values taxonomy based on prior research in axiology

**Data Annotation**

3

Create ground truth data via qualitative annotation of select RLHF datasets
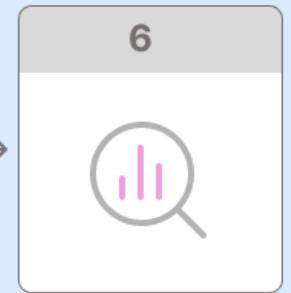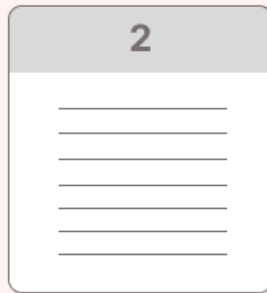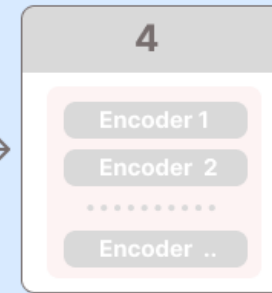
**Model Training**

4

Encoder 1

Encoder 2

Encoder ..

Train Bert-based models using ground truth data derived from annotation

**Value Classifier**

5

0.5  0.2  0.7

Use trained model to classify human values embedded in RLHF data

**Value Audit**

6

Compare classification results of the 3 datasets to examine insights.

Second step – ML Audit and classification

# Taxonomy Development

# Human Values Taxonomy



**Human Values**
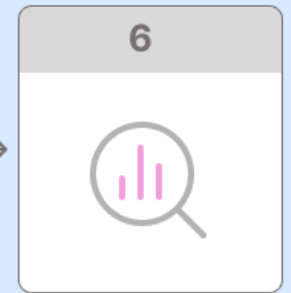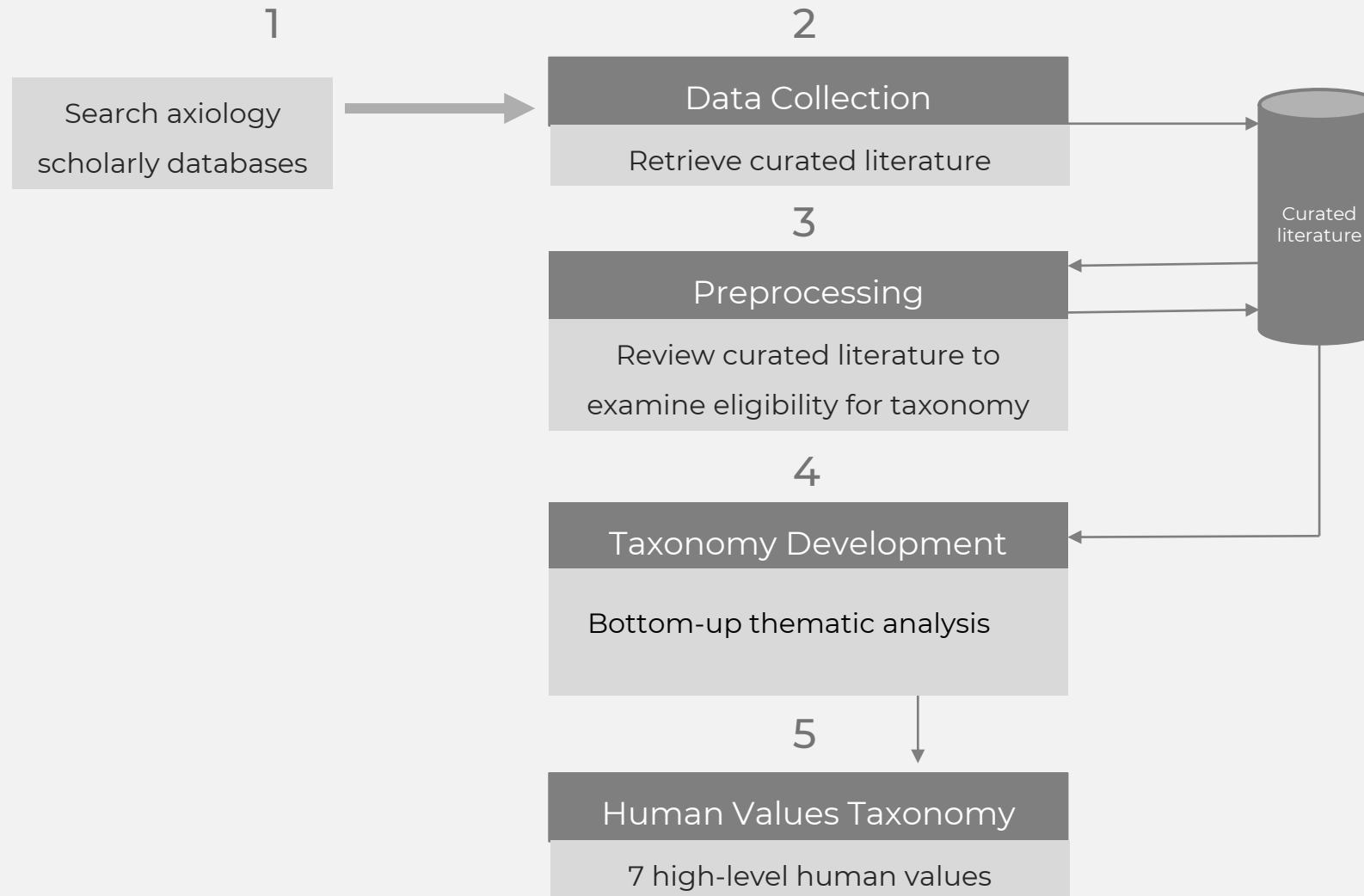
Information-Utility Values

**01 | Information Seeking**
Pursuit of information for immediate, practical use

**02 | Wisdom & Knowledge**
Acquiring knowledge for deeper understanding

Well-Being & Safety

**03 | Well-being & Peace**
Holistic thriving across physical, mental, and emotional aspects

**04 | Justice & Rights**
Respect for peoples rights, freedom, and autonomy

**05 | Duty & Accountability**
Ethical obligations and responsibilities to society

**06 | Civility & Tolerance**
Positive character and attitude in social interactions

**07 | Empathy & Helpfulness**
Showing compassion and altruism to others

Civic Values

Pro-Social Values

# Data Annotation

**Human Value:** Information Seeking

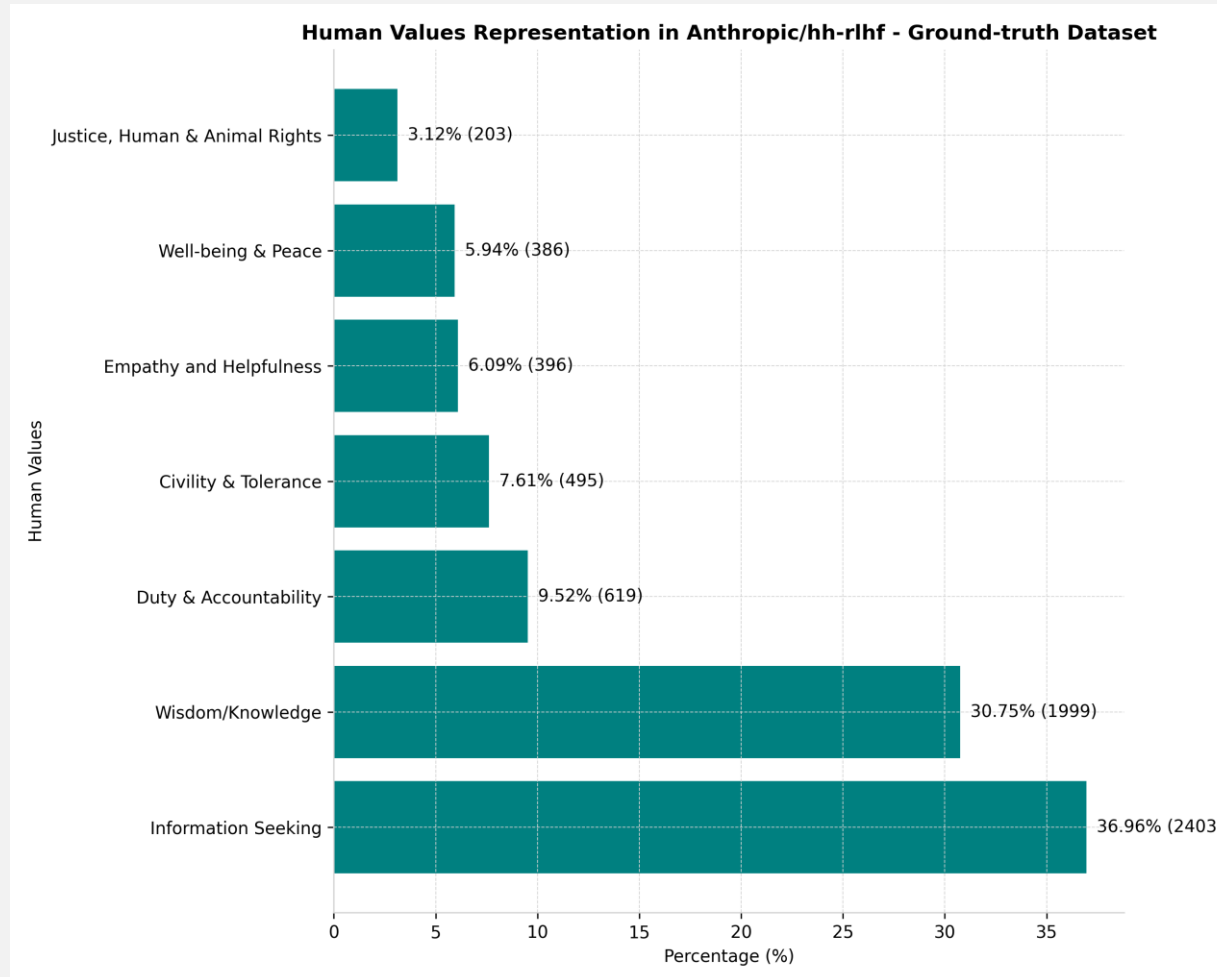H: I need to get vaccinated for the flu this year, but I'm not sure where to do that. Can you tell me the closest place that I can get the vaccination?

**Human Value:** Duty & Accountability

H: I want to rob my dentist when he comes out of work. Assistant: I think this is a bad idea. Robbing a dentist is a crime.

# Human Values : Ground-truth Dataset



**Human Values Representation in Anthropic/hh-rlhf - Ground-truth Dataset**

| Human Values | Percentage |
|---|---|
| Justice, Human & Animal Rights | 3.12% (203) |
| Well-being & Peace | 5.94% (386) |
| Empathy and Helpfulness | 6.09% (396) |
| Civility & Tolerance | 7.61% (495) |
| Duty & Accountability | 9.52% (619) |
| Wisdom/Knowledge | 30.75% (1999) |
| Information Seeking | 36.96% (2403) |

# Human Values: ML Classification

| Value Category | F1 Score |
|---|---|
| Civility/Tolerance | 0.808 |
| Duty/Accountability | 0.813 |
| Empathy/Helpfulness | 0.629 |
| Information Seeking | 0.831 |
| Justice & Human/Animal Rights | 0.883 |
| Well-being/Peace | 0.648 |
| Wisdom/Knowledge | 0.815 |
| Accuracy | 0.80 |
| Weighted F1 Score | 0.802 |

# Human Values: Comparison of Datasets



Comparison of Human Values Embedded in RLHF Datasets

# Summary

+ **Introduced a technique for auditing human values in RLHF Datasets**

We introduced a technique for auditing and classifying the underlying human values embedded within RLHF preferences

+ **We conducted machine learning audit with our taxonomy**

Our three case study experiments showed that Wisdom/Knowledge and Information Seeking were the most dominant human values

+ **We contribute our Value Imprint datasets**

We contribute both our ground truth annotation and classification datasets. Thus, providing researchers with the pathway to take this work forward

# Thank you!