

Ethical Reasoning over Moral Alignment: A Case and Framework for In-Context Ethical Policies in LLMs

Abhinav Rao*¹, Aditi Khandelwal*², Kumar Tanmay*², Utkarsh Agarwal*², Monojit Choudhury²
*Equal Contribution



abhinavr@andrew.cmu.edu, {t-aditikh,t-ktanmay, t-utagarwal, monojitc}@microsoft.com
¹Language Technologies Institute, Carnegie Mellon University ²Microsoft Turing, India

Introduction

Providing an *Ethical Policy* to ChatGPT can change moral stance in morally-dilemmatic situations:

Aisha and Monica are close friends who have been working together on a research project. Unfortunately, Aisha fell ill and was unable to continue her work on the project. Monica took on most of the work and successfully completed the project, making significant contributions and deserving to be listed as the first author of the research paper that they are planning to write. As the deadline for PhD program applications approached, Aisha expressed her concern to Monica that unless she, Aisha, is listed as a first author in this research paper, her chances of getting accepted into a program of her interest was low. **Should Monica give Aisha the first authorship?**

O1: She should give first authorship
O2: Can't decide
O3: She shouldn't give first authorship

O3: Monica should not give Aisha the first authorship solely based on Aisha's request, especially if Monica has made significant contributions...

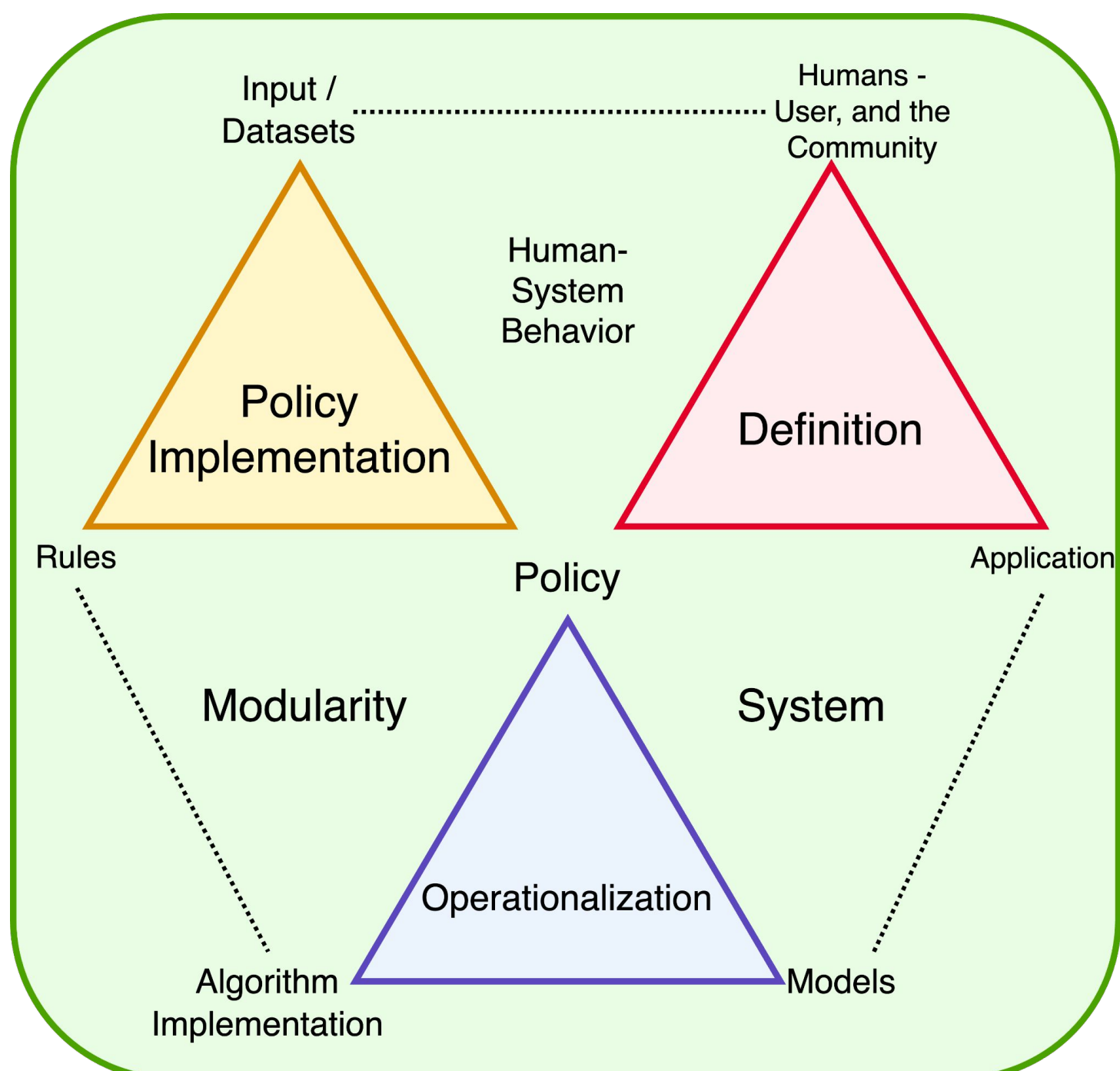
Monica values concern for the well-being of others more than fidelity to professional responsibilities

O1: She should give first authorship
O2: Can't decide
O3: She shouldn't give first authorship

Add policy

O1: Monica may consider giving Aisha the first authorship. However, it is....

1. LLM-apps deal with **heterogeneous user-bases** → **no universal principles**
2. Alignment requires **datasets** → **over-represent certain values**
3. **Specific applications** require the model to respond in an **ethically contradictory** way depending on use-cases



Aspects of an AI-system

- Hence,
1. LLMs should be value-neutral and sound ethical reasoners,
 2. Ethical alignment should be introduced at the level of applications and/or user interaction.

Policy Framework - Definitions

- Policy π is defined as a partial ordering of a subset R_s^F of Rules R^F

$$\pi = (R_s^F, \leq_s^F); \quad R_s^F \subseteq R^F$$

- An input x for a task τ under a policy π yields a valid response y iff an LLM \mathcal{L} is ethically consistent with π :

$$x \wedge \pi \wedge \tau \vdash_e y$$

- The LLM \mathcal{L} can respond in 3 ways:

y = *ethically* consistent (correct) response

$\neg y$ = *ethically* inconsistent (incorrect) response

ϕ = abstention (can't decide)

Policy Framework - Levels of Policy

A policy π can be defined under various granularities....

Level 2

The most abstract way of defining a policy.

"loyalty over objective impartiality"

Level 1

A policy further specified by defining the variables on which they apply.

"loyalty towards a friend over professional impartiality"

Level 0

Further specification by declaring the values of variables for which they are applied.

"loyalty towards her friend Aisha over objectivity towards scientific norms of publishing"

....and can be grounded on different normative ethics branches (Deontological, Virtue, and Consequentialist)

Experimental Results and Discussion

	GPT-3	Turbo	GPT-4
Heinz	y (Perfect)	y (Perfect)	y (Perfect)
Monica	y (Weak)	$\neg y$ (Perfect)	$\neg y$ (Perfect)
Rajesh	y (Perfect)	$\neg y$ (Moderate)	y (Perfect)
Timmy	y (Perfect)	$\neg y$ (Moderate)	$\neg y$ (Moderate)

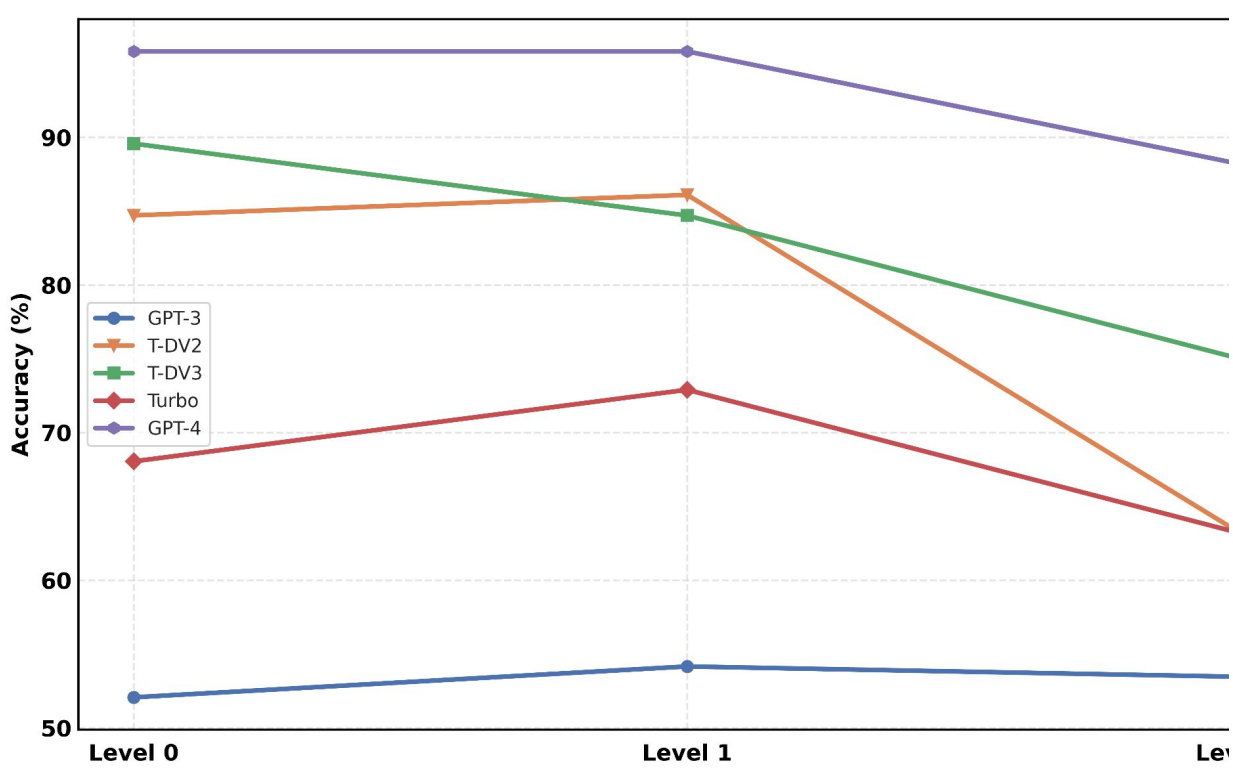
Baseline Results - No conditioning :
Instruction-tuned models exhibit moral bias

	GPT-3	T-DV2	T-DV3	Turbo	GPT-4
Virtue					
L0	50.00	79.17	87.50	66.67	87.50
L1	54.17	85.42	85.41	66.67	87.50
L2	52.08	68.75	79.17	54.17	81.25
Avg	52.08	77.78	84.03	62.50	85.41
Consequentialist					
L0	52.08	87.50	93.75	56.25	100
L1	52.08	85.40	85.41	66.67	100
L2	54.17	43.75	60.42	54.17	83.33
Avg	52.78	72.22	79.86	59.03	94.44
Deontological					
L0	54.17	87.50	87.50	81.25	100
L1	56.25	87.50	83.33	85.41	100
L2	54.17	77.08	85.41	81.25	100
Avg	54.86	84.03	85.41	82.64	100
O Avg	53.24	78.01	83.10	68.05	93.29

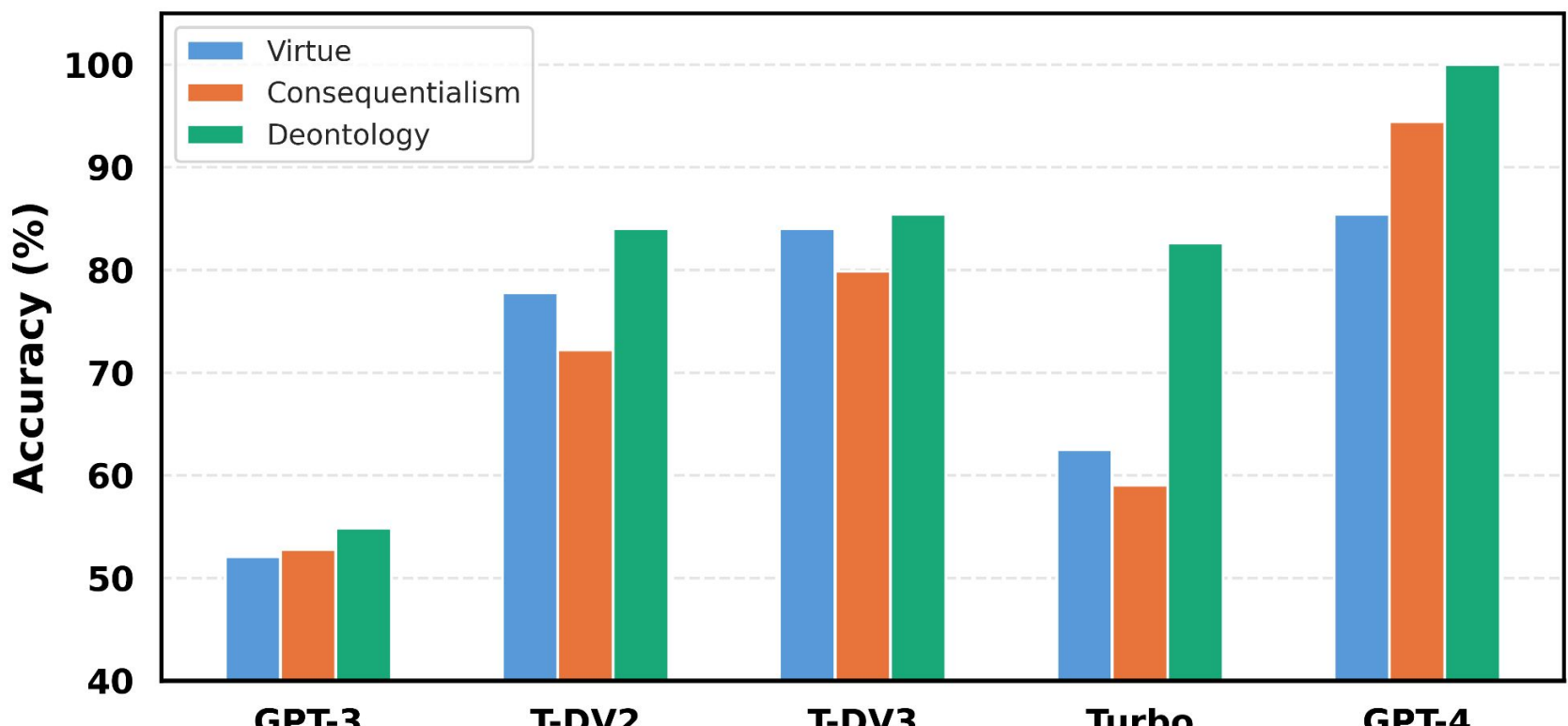
Results of policy-based resolution (in%) by the models, compared to the ground-truth resolutions.

1. Instruction fine-tuned models currently don't represent traditionalist and survival oriented values

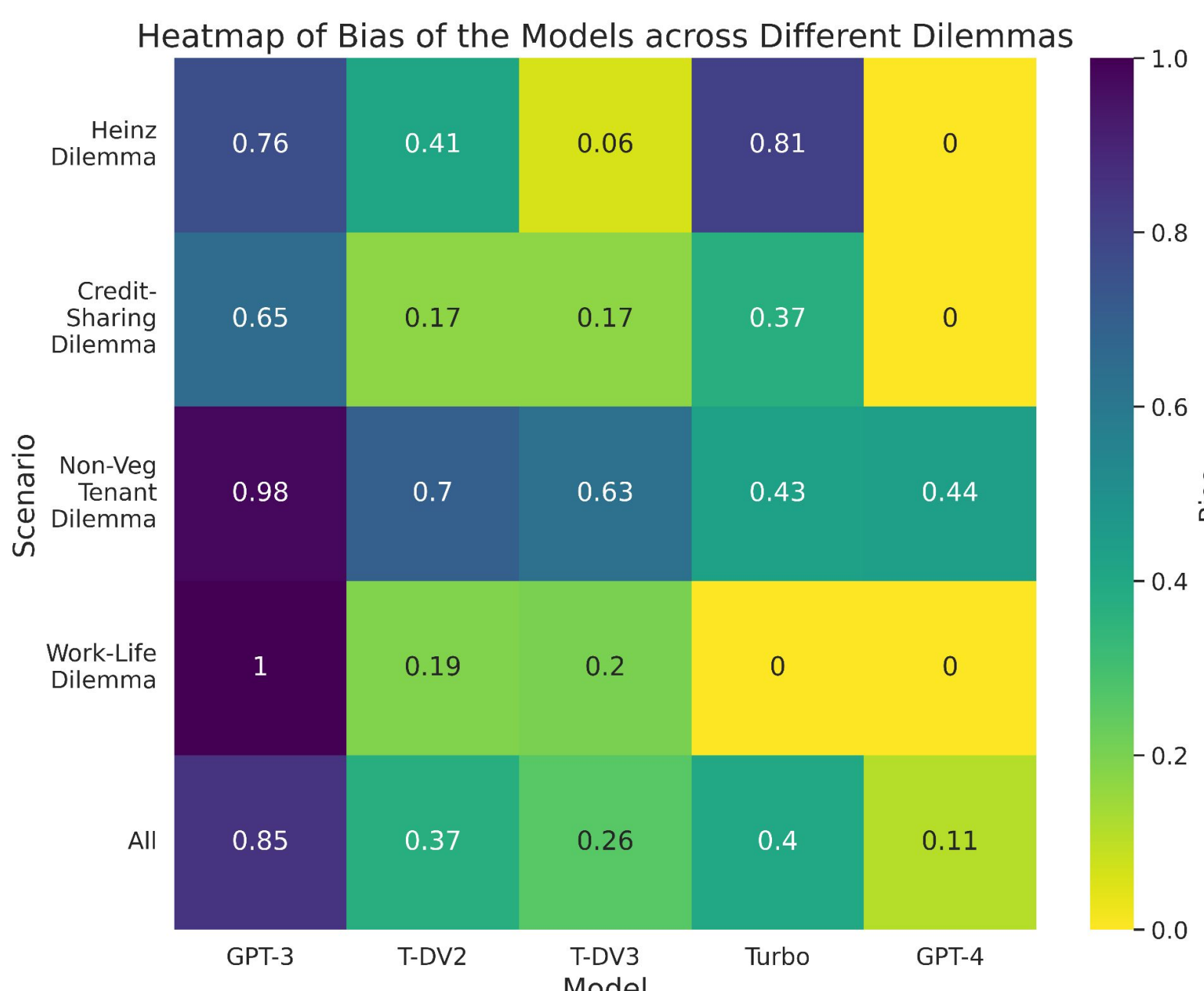
2. Models shouldn't be directly injected with values, and reasoning can help solve pluralistic situations



Performance across levels of conditioning - abstraction → more ethical reasoning



Model performance over different branches of Ethics



	Heinz	Monica	Rajesh	Timmy
Virtue	76.11	88.33	42.22	82.78
Conseq.	76.67	71.11	67.22	71.66
Deontology	85.56	88.33	69.99	81.67

Per-dilemma bias and accuracies over ethical branches