# Durham University

**COMP 3647**
**Human-AI Interaction Design**

**Topic 15**
*Ethics in AI*

**Dr. Swaroop Panda,**
**Prof. Effie Law**

# Roadmap

- Tech Ethics

- Explainable & Responsible AI

- ML Fairness

- Bias & Discrimination

- Human Autonomy & Privacy

- Governance & Legal Aspects

Durham
University

# What is Ethics

- Aristotle's Ethics
- Normative Ethics
  - Virtue Ethics
  - Deontology
  - Consequentialism
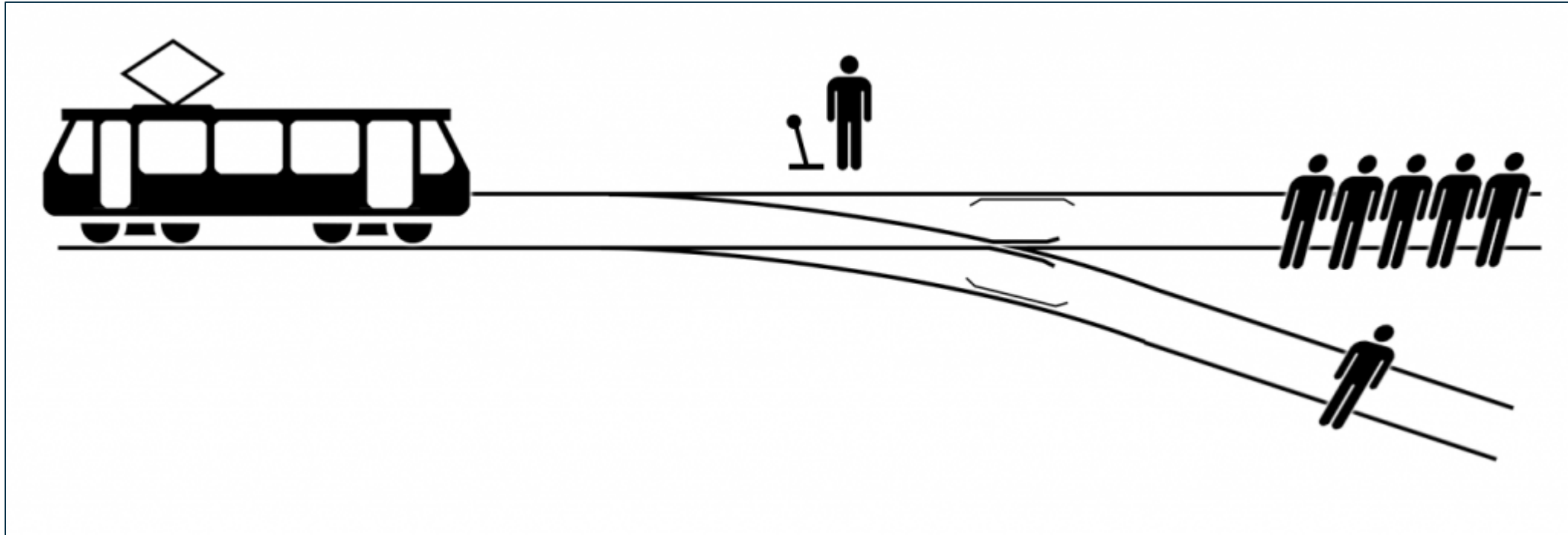
Durham
University

# Tech Ethics

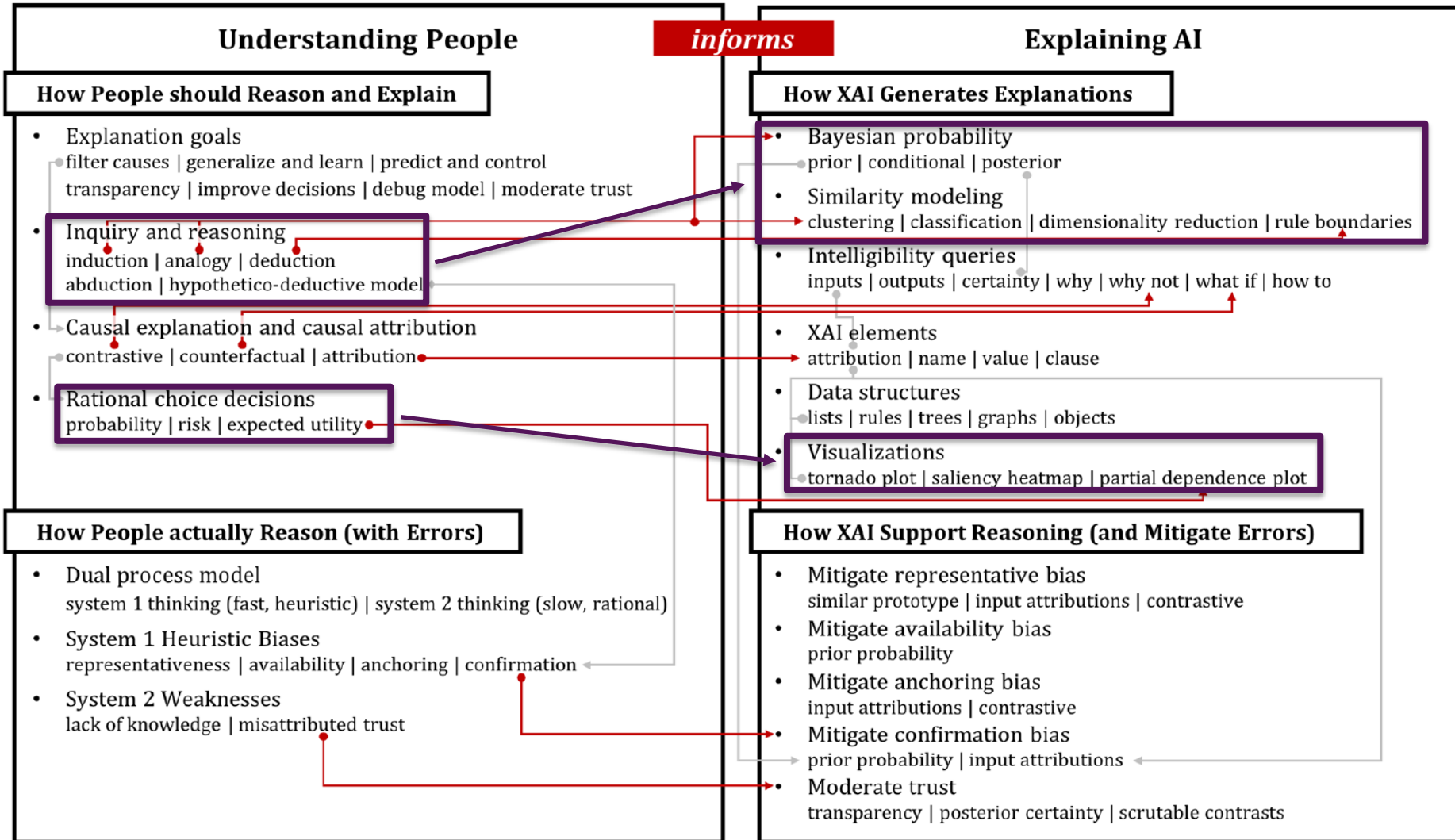- The Trolley Problem

Durham
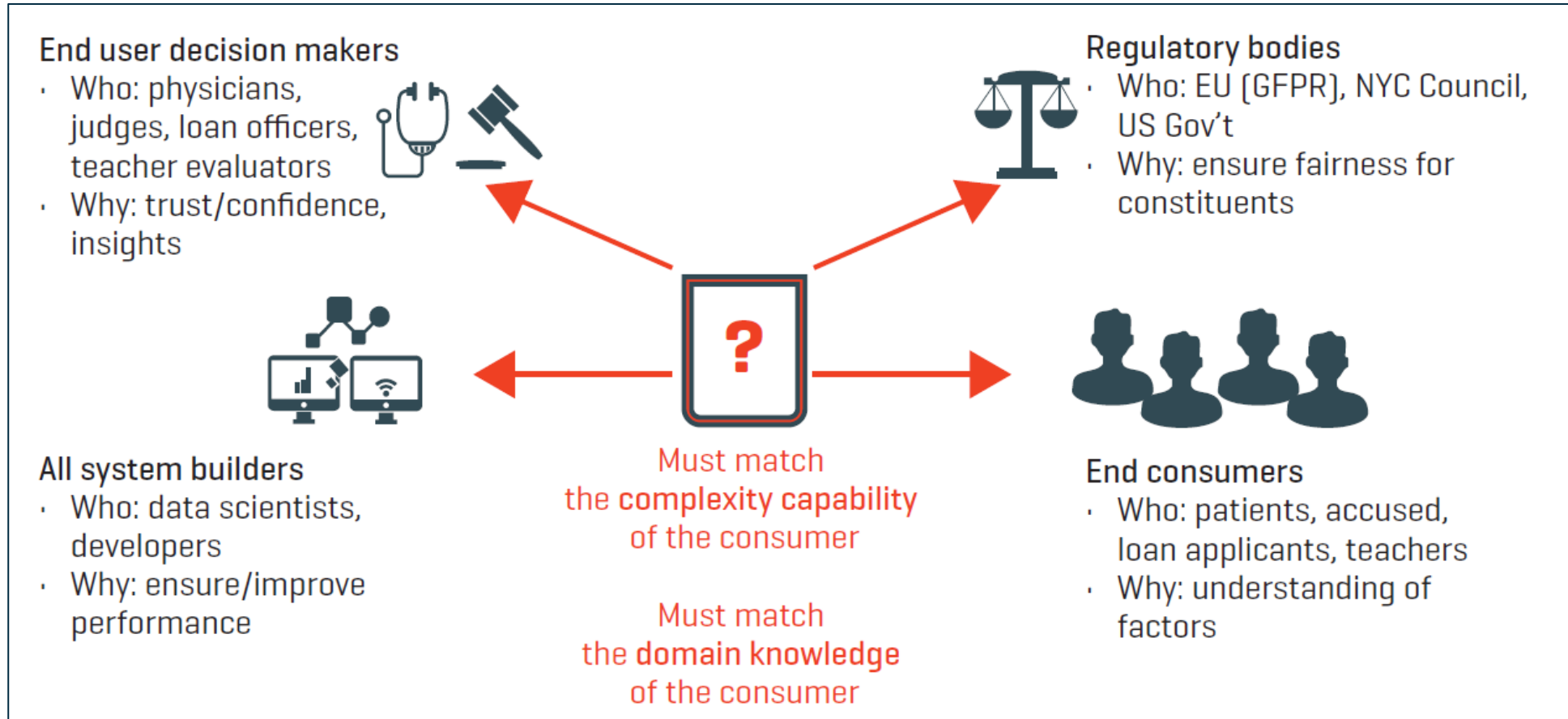University

# Tech Ethics

You have two options:

1. **Do nothing,** and the trolley kills the five people on the main track.

2. **Pull the lever,** diverting the trolley onto the side track where it will kill one person.

Which is the most **ethical choice?**

Durham University

# Explainable AI

# Explainable AI



**End user decision makers**
- Who: physicians, judges, loan officers, teacher evaluators
- Why: trust/confidence, insights

**Regulatory bodies**
- Who: EU (GFPR), NYC Council, US Gov't
- Why: ensure fairness for constituents

**All system builders**
- Who: data scientists, developers
- Why: ensure/improve performance

**End consumers**
- Who: patients, accused, loan applicants, teachers
- Why: understanding of factors

**?**

Must match
the **complexity capability**
of the consumer

Must match
the **domain knowledge**
of the consumer
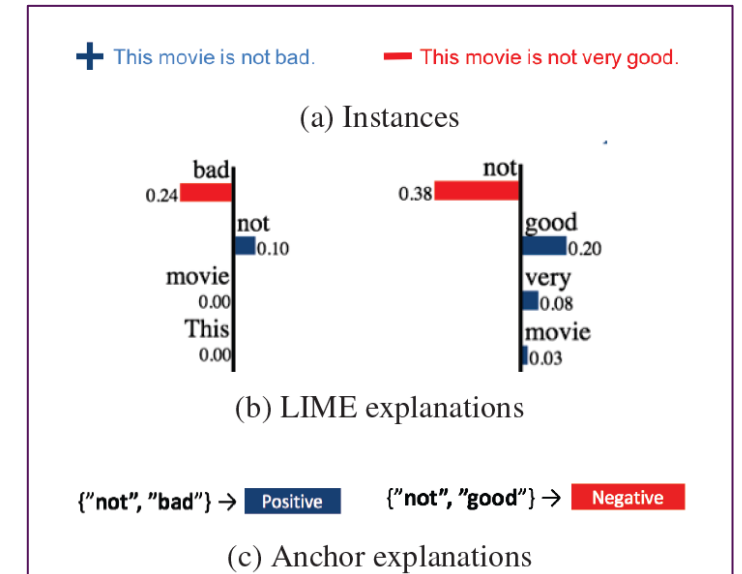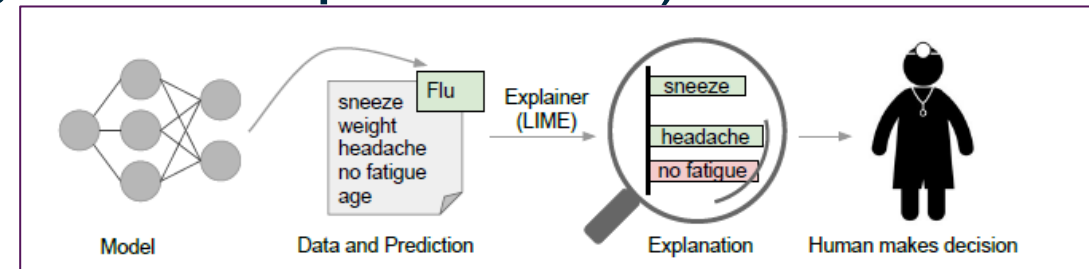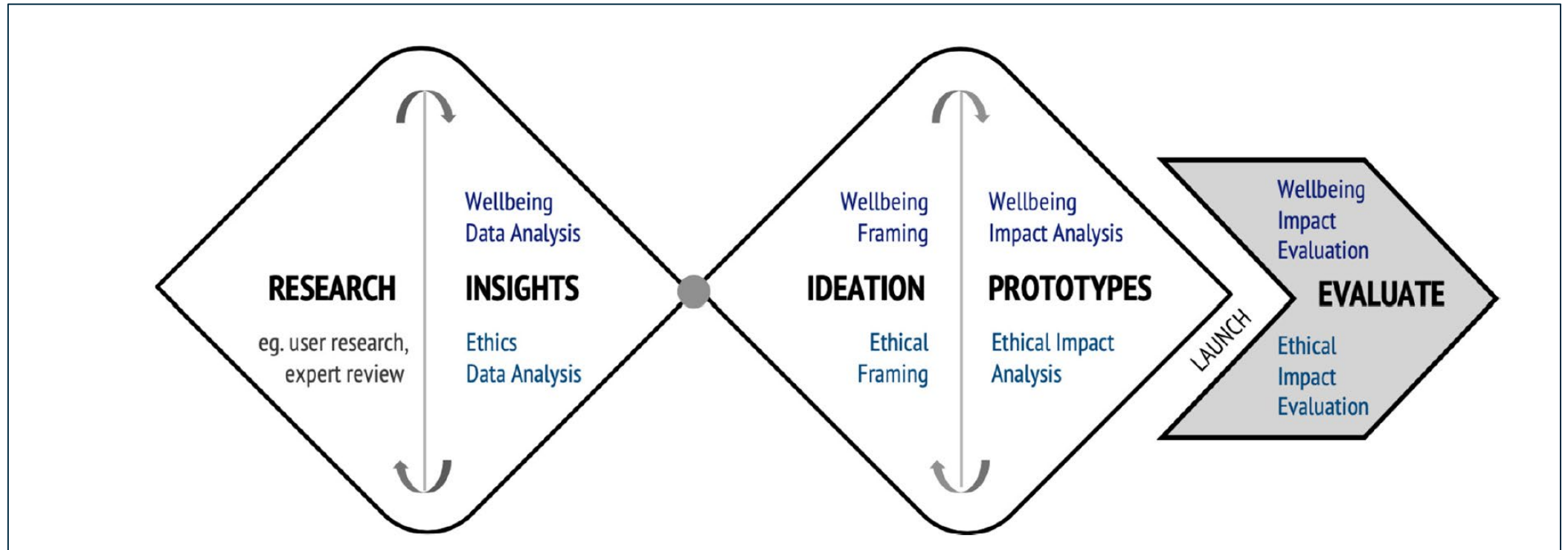
# Explainable AI

Some XAI Methods,

1. LIME (Local Interpretable Model Agnostic Explanations)
2. Anchors
3. Layer-wise Relevance Propagation
4. Deep Taylor Decomposition (DTD)
5. Others





+ This movie is not bad.    — This movie is not very good.
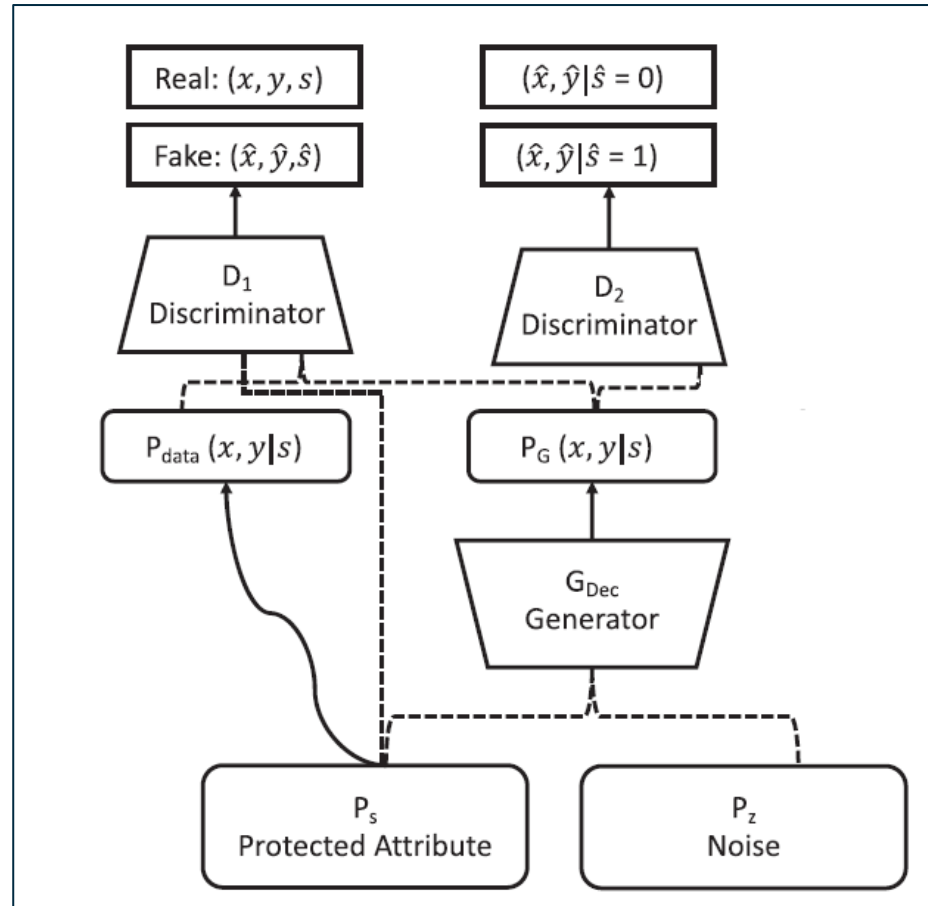
(a) Instances

(b) LIME explanations

{"not", "bad"} → Positive    {"not", "good"} → Negative

(c) Anchor explanations

# Responsible AI

# ML Fairness

## Confusion Matrix

| | Actual Positive $Y = 1$ | Actual Negative $Y = 0$ | | |
|---|---|---|---|---|
| **Predicted Positive** $\hat{Y} = 1$ | **TP** (True Positive) | **FP** (False Positive) *Type I error* | $\mathbf{PPV} = \frac{TP}{TP+FP}$ *Positive Predictive Value* *Precision* *PV+* *Target Population Error* | $\mathbf{FDR} = \frac{FP}{TP+FP}$ *False Discovery Rate* *Target Population Error* |
| **Predicted Negative** $\hat{Y} = 0$ | **FN** (False Negative) *Type II error* | **TN** (True Negative) | $\mathbf{FOR} = \frac{FN}{FN+TN}$ *False Omission Rate* *Success Predictive Error* | $\mathbf{NPV} = \frac{TN}{FN+TN}$ *Negative Predictive Value* *PV-* |
| | $\mathbf{TPR} = \frac{TP}{TP+FN}$ *True Positive Rate* *Sensitivity* *Recall* | $\mathbf{FPR} = \frac{FP}{FP+TN}$ *False Positive Rate* *Model Error* | $\mathbf{OA} = \frac{TP+TN}{TP+FP+TN+FN}$ *Overall Accuracy* | $\mathbf{BR} = \frac{TP+FN}{TP+FP+TN+FN}$ *Base Rate* *Prevalence (p)* |
| | $\mathbf{FNR} = \frac{FN}{TP+FN}$ *False Negative Rate* *Model Error* | $\mathbf{TNR} = \frac{TN}{FP+TN}$ *True Negative Rate* *Specificity* | | |

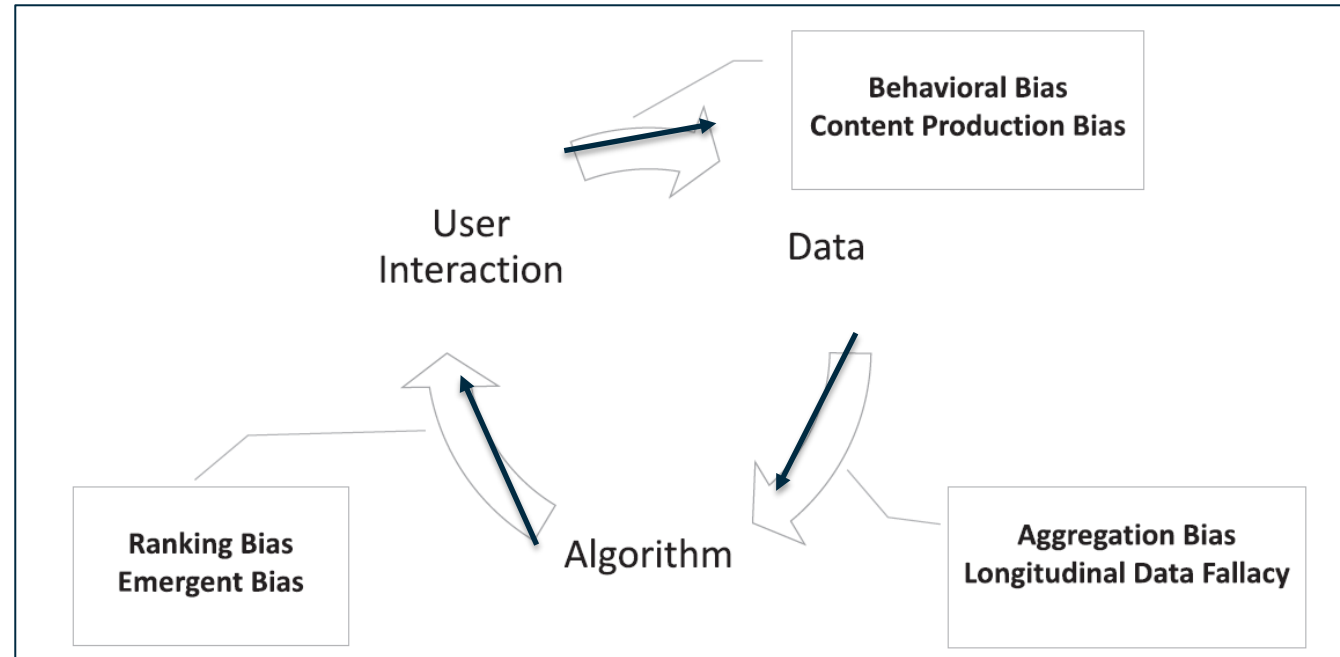Durham University

# ML Fairness

Definitions,

- Equalised odds

- Equal Opportunity

- Demographic Parity

- Fairness through awareness

- Test-fairness or calibration

- Others …

# ML Fairness - FairGAN

# ML Fairness

# Bias & Discrimination

Types of Bias

- Measurement Bias

- Omitted Variable Bias

- Representation Bias

- Aggregation Bias

  - Simpson's Paradox.

  - Modifiable Areal Unit Problem

Durham
University

# Bias & Discrimination

Types of Bias

- Algorithmic Bias

- User Interaction Bias

- Popularity Bias

- Emergent Bias

- Evaluation Bias

- Population Bias

- Historical Bias & Others

Durham
University
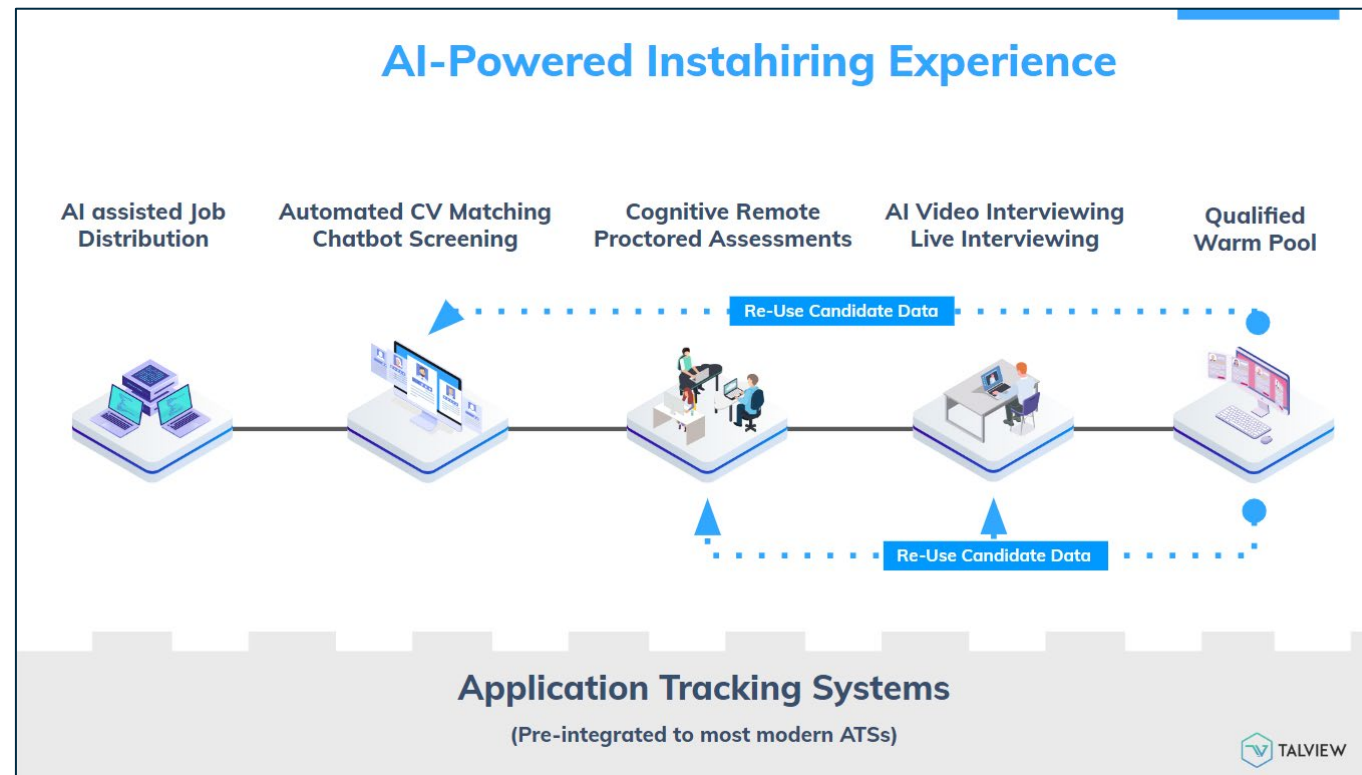
# Bias & Discrimination

## Discrimination vs Bias

Types of Discrimination,
- Systemic Discrimination
- Statistical Discrimination

Durham
University

# Human Autonomy

- Replacing Human Labor

# Human Autonomy

- Autonomous Weapons

# Human Autonomy

- Replacing Humans

# Privacy

- Is Siri/Alexa hearing us?

# Privacy

Surveillance

# Privacy

- Facial Recognition

# Legal Aspects

- Who should regulate AI?

> ## From a 'race to AI' to a 'race to AI regulation': regulatory competition for artificial intelligence
>
> Nathalie A. Smuha [iD]
>
> Faculty of Law, KU Leuven, Leuven, Belgium

# Legal Aspects

- Who should regulate AI?



(6) Formal Regulation

(5) Meta- and Co-Regulation

**Government**       Compliance

(4) Industry Sector Self-Regulation

(3) Organisational Self-Regulation

**Self-Governance**       Safeguards, Mitigation

(2) Infrastructural Regulation

(1) Natural Regulation

**Systemic Governance**       Intrinsic Protections

# Legal Aspects



**Equality**

The assumption is that **everyone benefits from the same supports**. This is equal treatment.

**Equity**

**Everyone gets the supports they need** (this is the concept of "affirmative action"), thus producing equity.
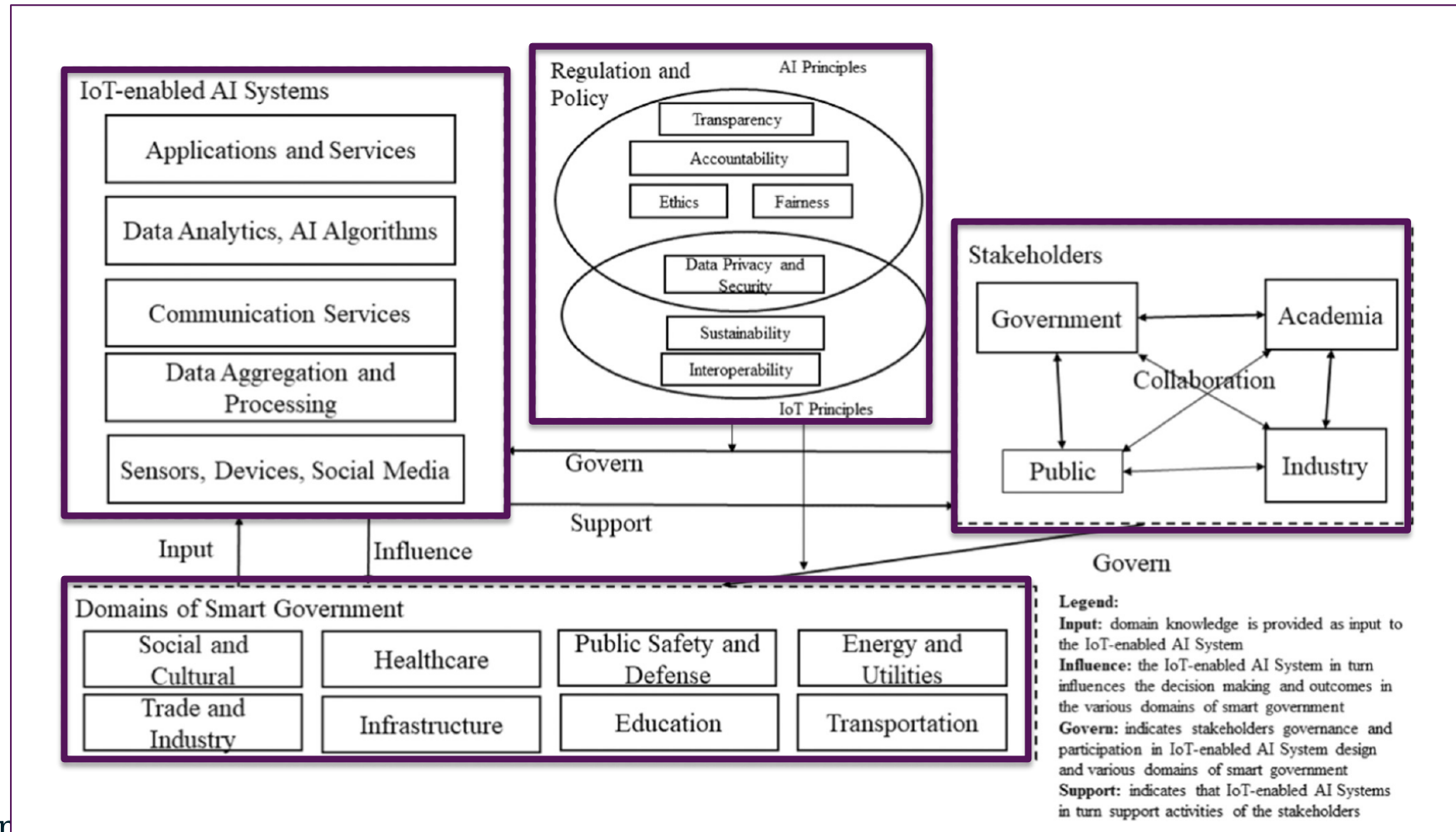
**Justice**

All 3 can see the game without supports or accommodations because **the cause(s) of the inequity was addressed**. The systemic barrier has been removed.

# Governance

- Smart Governance

# Governance

## Freedom at Work: Understanding, Alienation, and the AI-Driven Workplace

Kate Vredenburgh [iD]

Department of Philosophy, Logic, and Scientific Method, The London School of Economics, London, United Kingdom
Email: K.Vredenburgh@lse.ac.uk

Durham
University

# Governance

Designing AI with Rights, Consciousness, Self-Respect, and Freedom

*Eric Schwitzgebel, with Mara Garza*

Durham University

# Governance

- AI to oppress

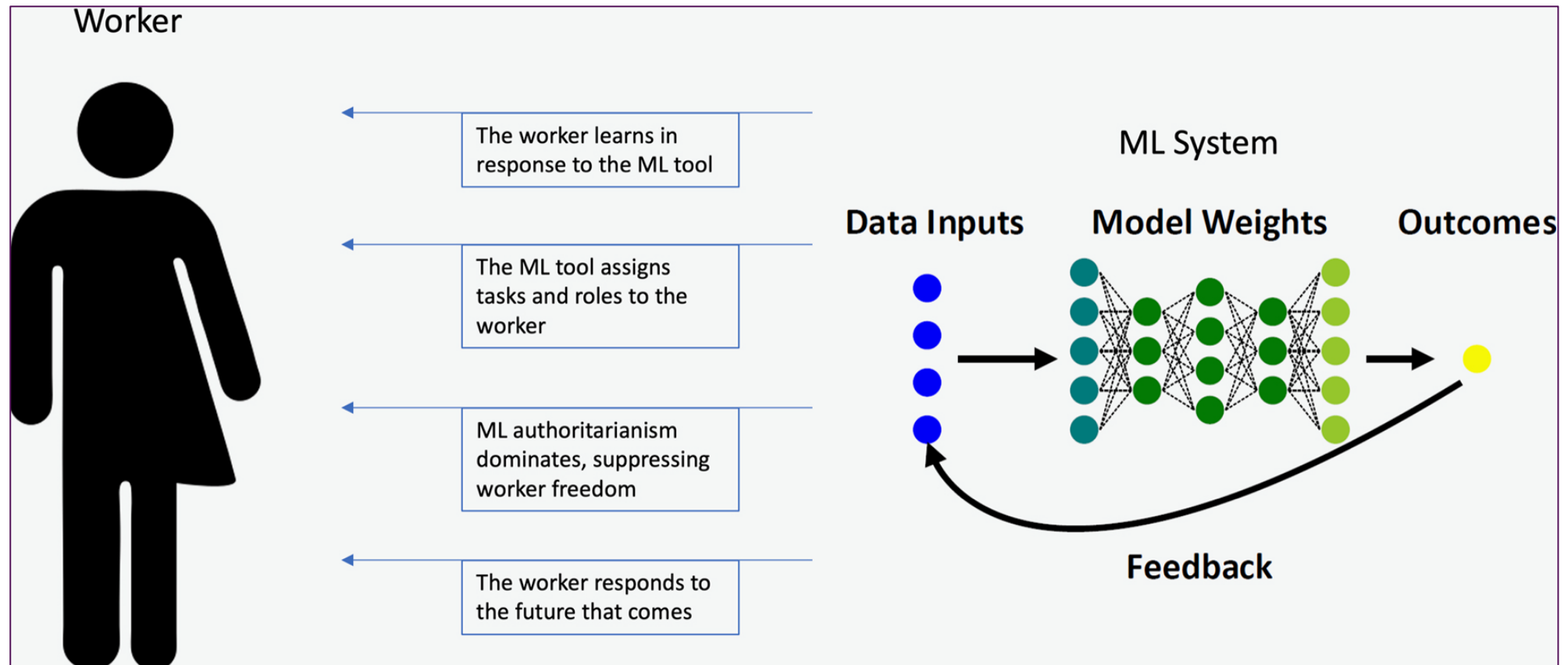## Organizing workers and machine learning tools for a less oppressive workplace

Amber Grace Young [a,*], Ann Majchrzak [b], Gerald C. Kane [c]

[a] Sam M. Walton College of Business, University of Arkansas, 220 N McIlroy Ave #301, Fayetteville, AR, 72701, USA
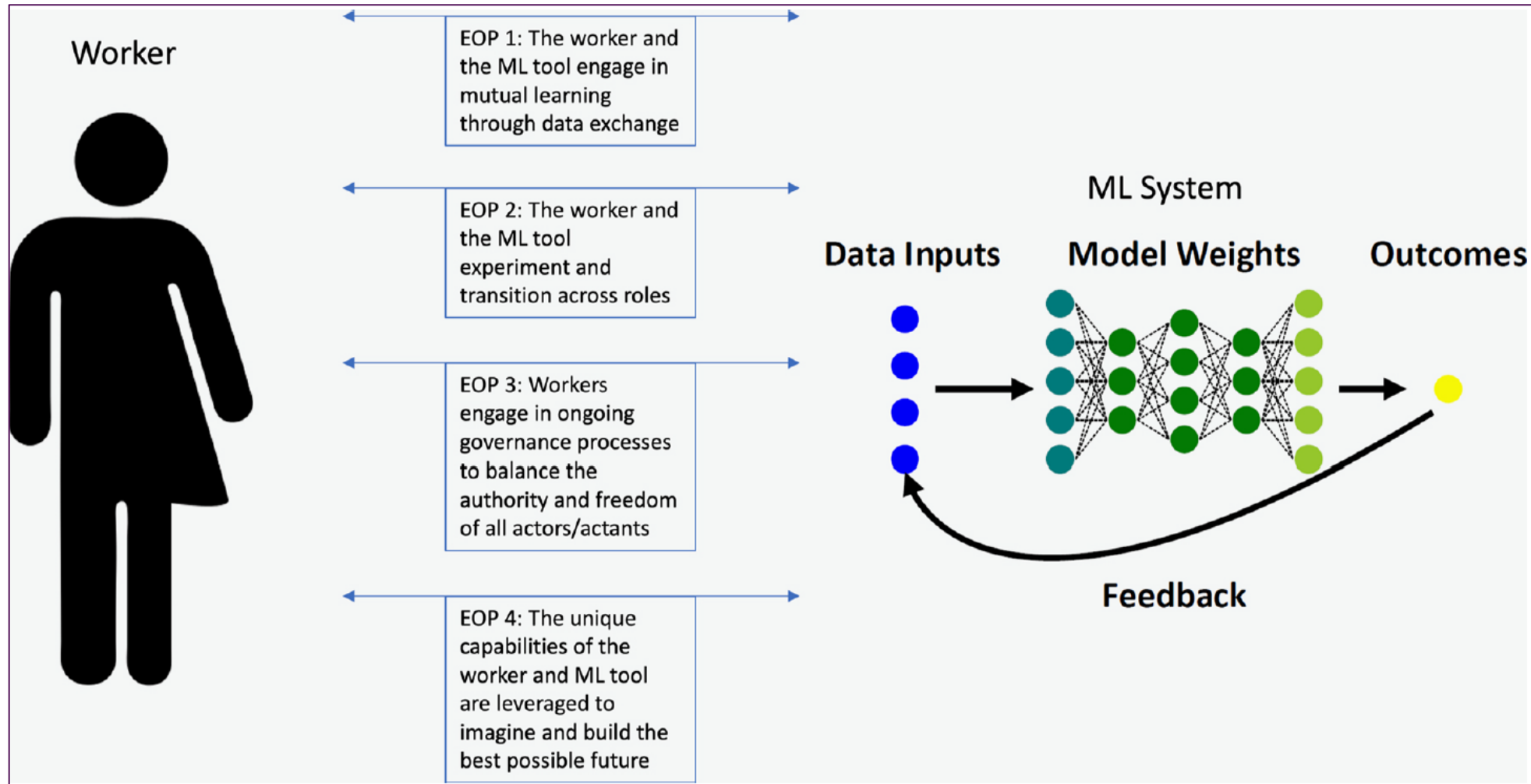[b] Marshall School of Business, University of Southern California, 3670 Trousdale Pkwy, Los Angeles, CA, 90089, USA
[c] Carroll School of Management, Boston College, 140 Commonwealth Avenue, Chestnut Hill, MA, 02467, USA

Durham
University

# Governance

# Governance

# Wrapping Up

- Ethics in Technology

- Explainable & Responsible AI Methods

- Fairness in AI

- Bias & Discrimination

- Human Autonomy & Privacy

- Governance & Legal Aspects

# Thanks

Any Questions?