# Module 10: Explanation and explanability

Gert Meyers (TILT, Tilburg University)

TILBURG UNIVERSITY

Understanding Society

# This week

- Explanation and explainability
- Third debating session!

- The blogpost

MODULE 8: TRUST AND TRUSTWORTHINESS

MODULE 9: AUTONOMY

MODULE 10: EXPLANATION AND EXPLAINABILITY

- [Explanation](#)
- Transparancy
- [Accountability](#)
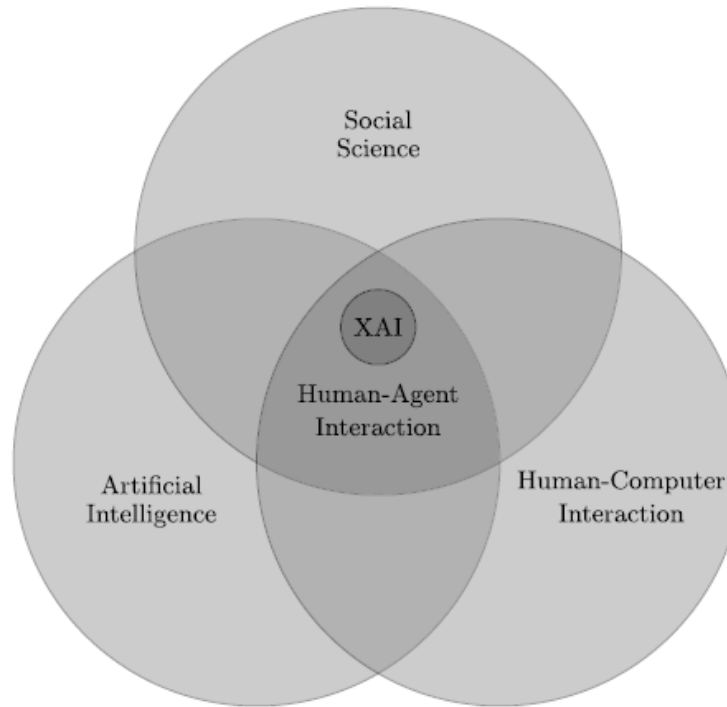- Justification

**Fig. 1.** Scope of explainable artificial intelligence.
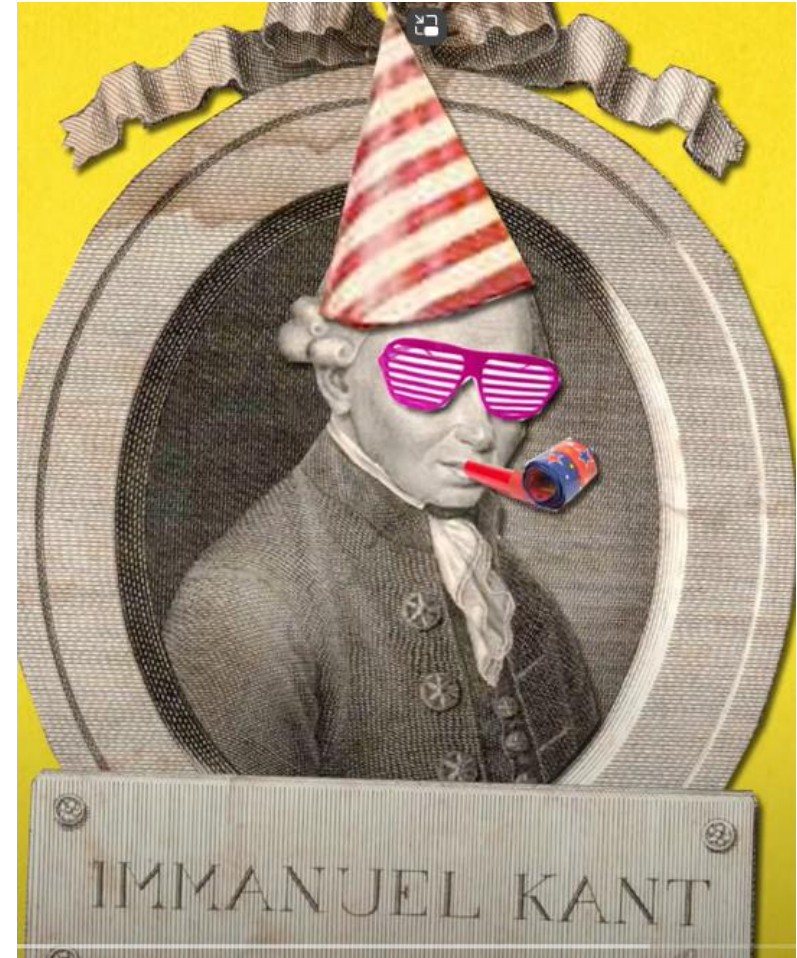
Note: 'Everyday explanation' vs 'general explanation'

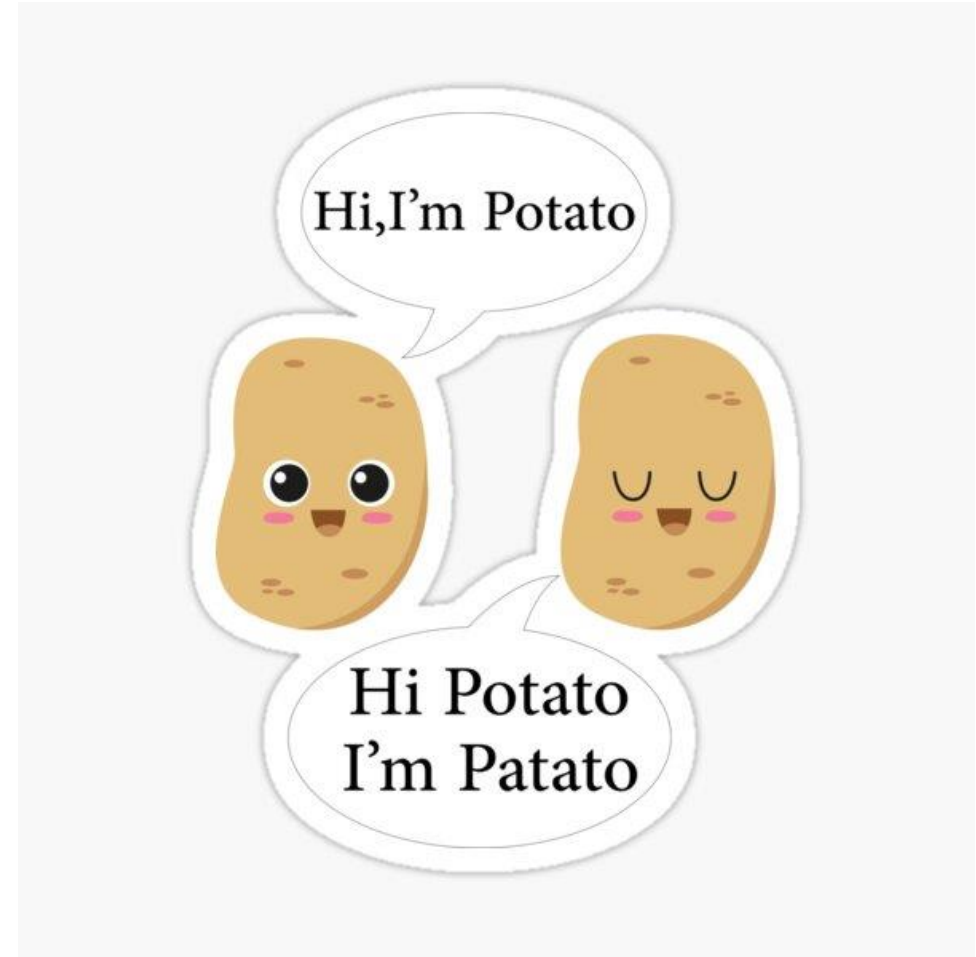Note: decisions

# Happy Birthday!



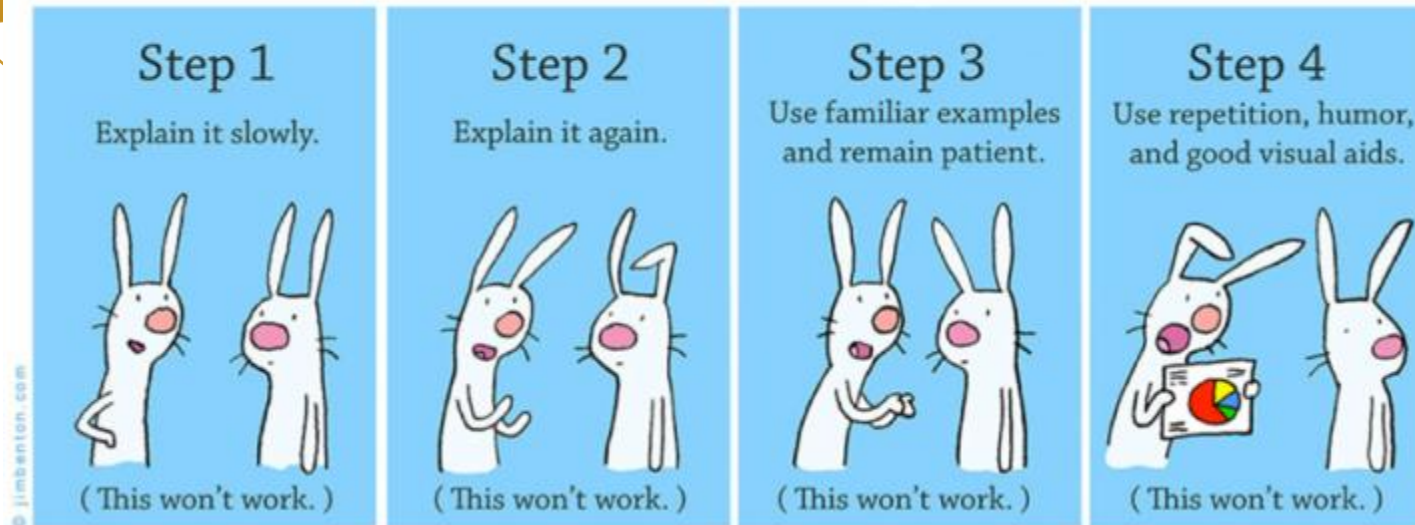Immanuel Kant: born on 22 april 1724
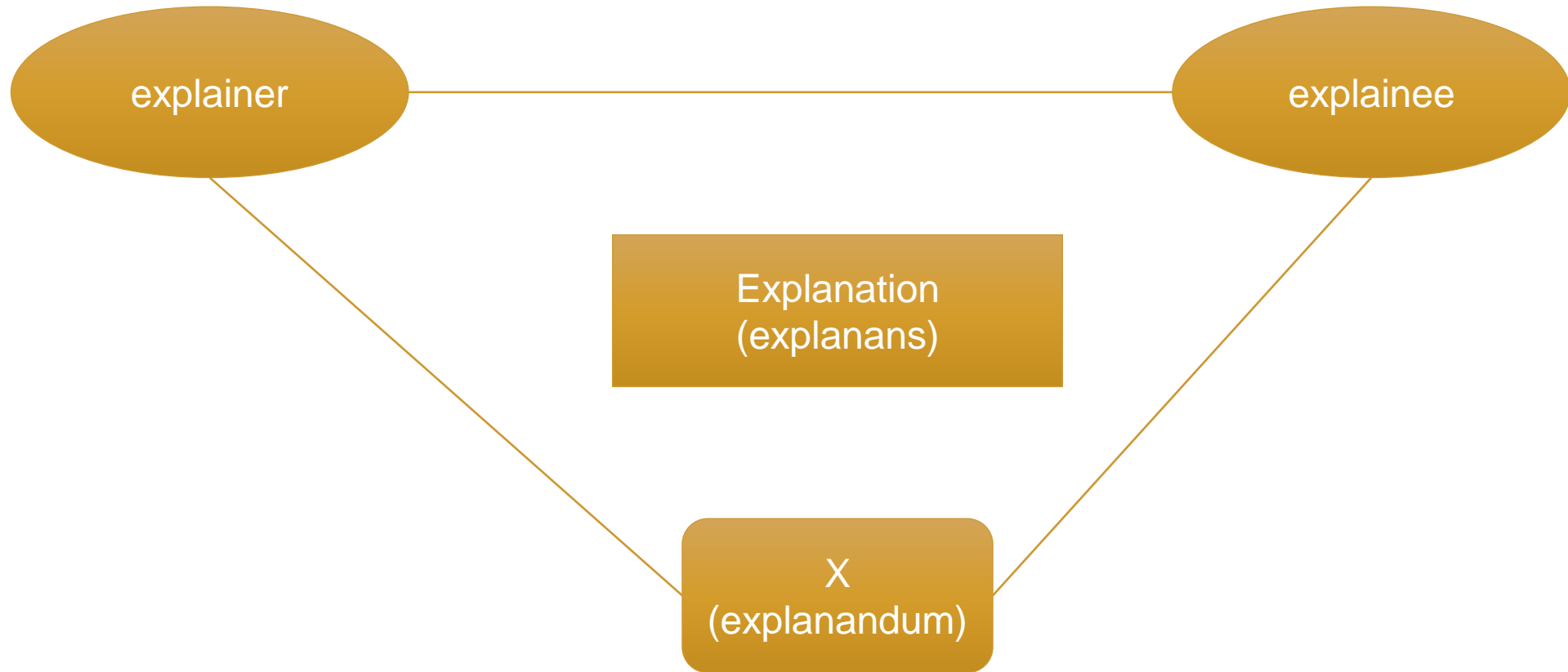
# Explanation as product and process (twice)

- Cognitive process: identifying the causes of a particular phenomenon
- Product: the explanation
- Social process: transferring knowledge between explainer and explainee

explain

ee

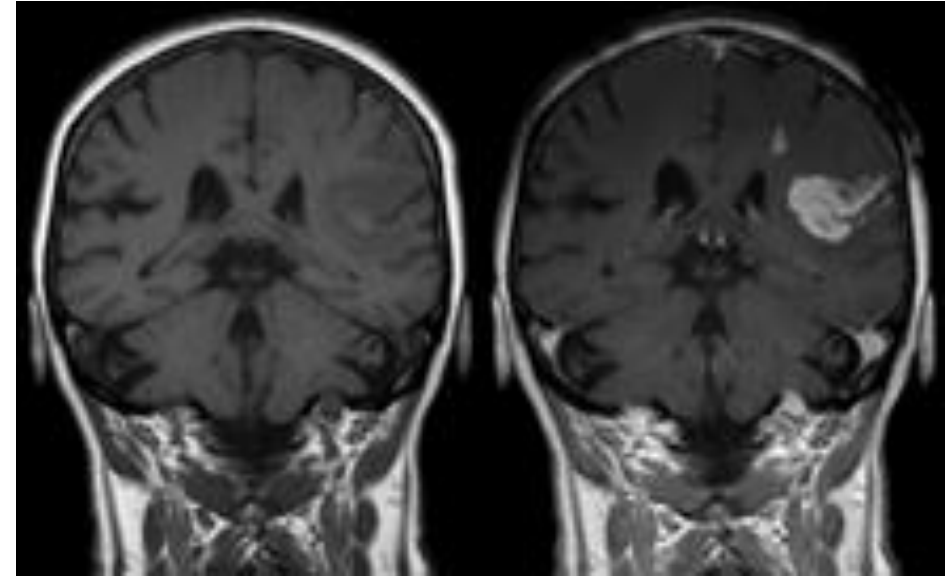How to explain things ..to the disengaged?

Step 1
Explain it slowly.
( This won't work. )

Step 2
Explain it again.
( This won't work. )

Step 3
Use familiar examples and remain patient.
( This won't work. )

Step 4
Use repetition, humor, and good visual aids.
( This won't work. )

TILBURG UNIVERSITY

Note: no arrows

explainer

explainee

Explanation
(explanans)

X
(explanandum)

Note: no arrows

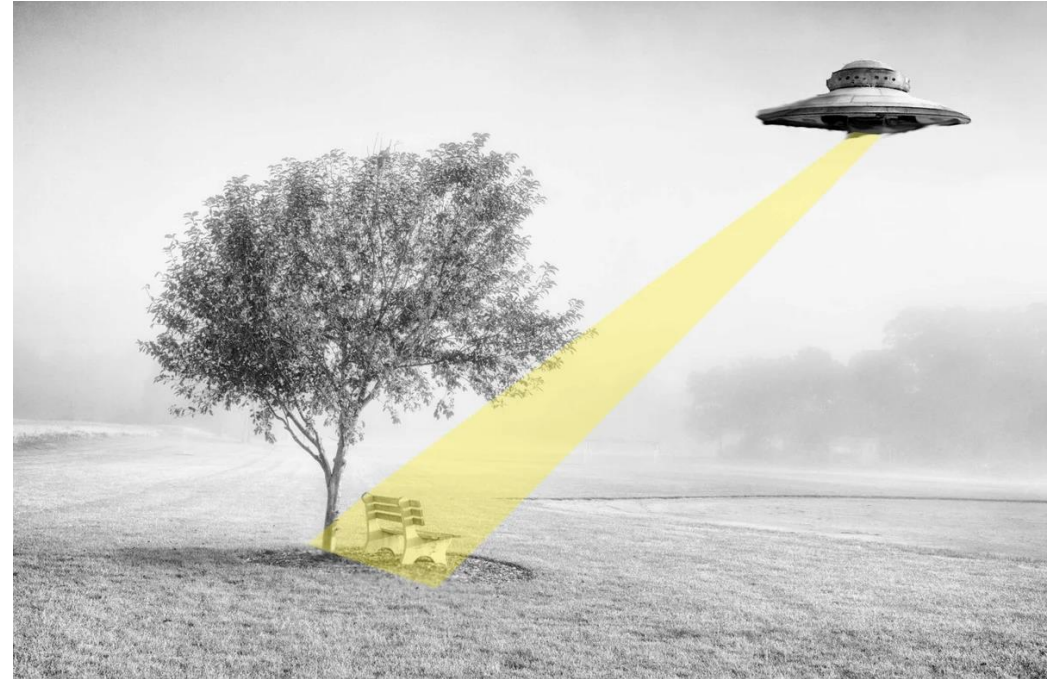trustee

trustor

Note: no arrows

# Explanations are contrastive

- Counterfactual cases

- Why this instead of that?

- Layperson will find contrastive explanations more intuitive and more valuable

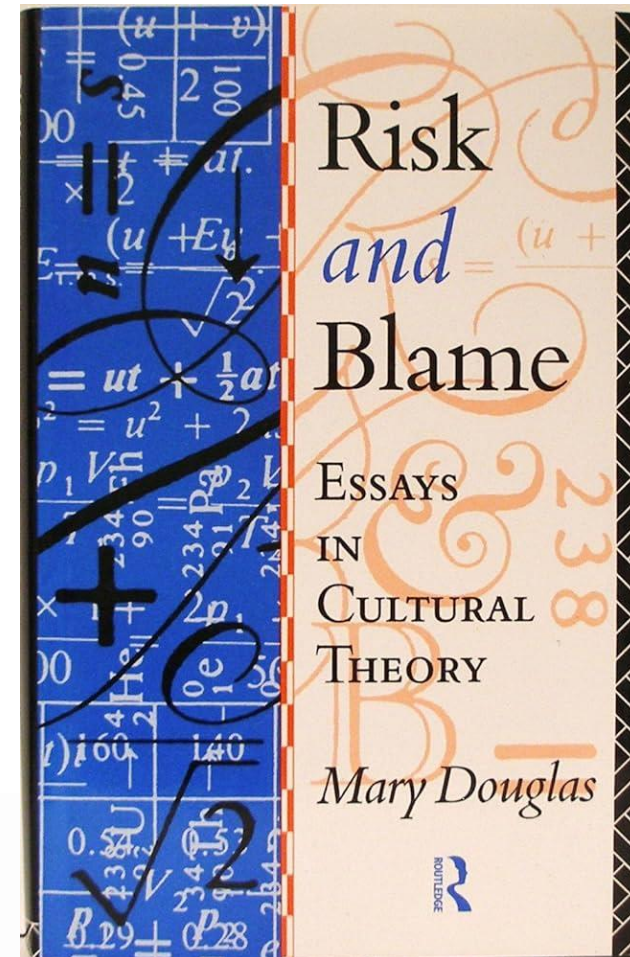- Science girl op X: 'Explain this! https://t.co/7izH4Hbivz' / X (twitter.com)

# Explanations are selected

- Humans are adept at selecting one ore two causes from a sometimes infinte number of cause to be *the* explanation

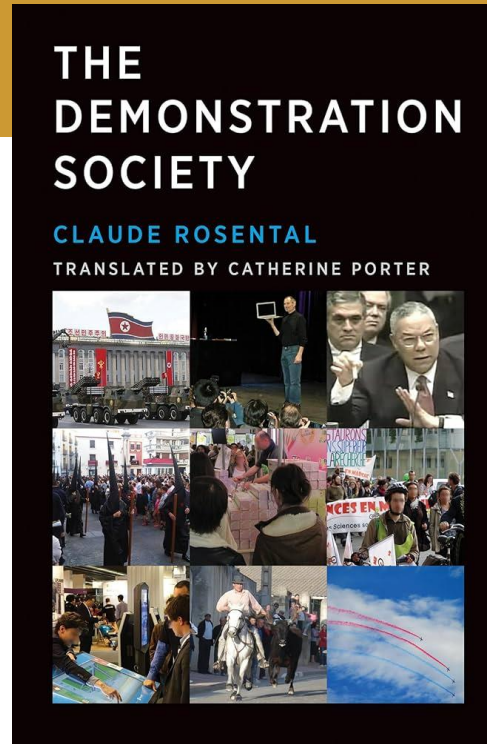- Cognitive burden of complete explanations is too great

# Probabilities probably don't matter

- Reffering to probabilities vs referring to causes
- The most likely explanation is not always the best explanation <u>for a person</u>

- Example: collapse of the soviet union

# Explanations are social

- Transfer of knowledge

- Part of an interaction

- Erving Goffman: 'definition of the situation'

# Explanation

- Contrastive
- Selected
- Propabilities probably don't matter
- Social

- CONTEXTUAL
  - An explanation is not a (mathemathical) proof
  - 'the best' explanation does not always work

**Accountability**

- Actors continuously provide observable accounts of what they are doing (ethnomethods)

- Often implicit
  - Et cetera clause
  - For ex. Skipping slides because lack of time

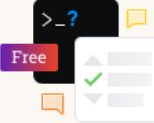- ADM are decision that take place in social context, but opacity problem



Harold
Garfinkel

**Studies in Ethnomethodology**



HAROLD GARFINKEL
EDITED WITH AN INTRODUCTION BY MICHAEL LYNCH

**HAROLD GARFINKEL: STUDIES OF WORK IN THE SCIENCES**

Directions in Ethnomethodology and Conversation Analysis

# This slide is skipped because lack of time

# Accountability and public reason

- 'obligation to provide justification'

- What type justifications are *good* justifications?

- ADM

  - Epistemic and normative assumptions

  - Reflected in justification

- No guarantee that explanation is acceptable to all (reasonable pluralism)

- 'There is therefore likely to be a gap between merely providing reasons and explanations for an ADM system's output, and providing adequate justification for them that will be acceptable to affected decision-subjects'

TILBURG ◆ UNIVERSITY

# Public reason

- 'Public reason proposes that universal rules must be justifiable on grounds that are suitably public and shared by all reasonable people in the society, and without appeal to beliefs that are controversial.'

# Next week:

- Bias, discrimination and data justice (Donovan)
- Last debating session

# Questions on Module 10?

G.Meyers@tilburguniversity.edu

TILBURG ✦ UNIVERSITY

Understanding Society