



BIASES → EXAMPLE OF JOB RECRUITMENT. AI TRAINED OF ACTUAL HISTORICAL DATA. WHITE MALE MOST IMPORTANT PREDICTOR FOR TAKING THE JOB

BIASES CAN ARISE

- ALL STAGES OF DESIGN, TESTING AND APPLICATION (TO FOCUS ON DESIGN)
- SELECTION OF THE TRAINING SET
- IN THE TRAINING SET ITSELF (UNREPRESENTATIVE/INCOMPLETE) VS
- IN THE ALGORITHM
- IN THE DATASET THE ALGORITHM IS GIVEN ONCE IS TRAINED
- IN DECISION BASED ON FAKE CORRELATION
- IN THE GROUP THAT CREATES THE ALGORITHM

SIMILAR TO TECHNICAL ARTIFACTS,
BUT HAVE AUTONOMY, ADAPTIVITY
AND INTERACTIVITY

ARTIFICIAL AGENTS ≠ HUMANS AGENT

- CAN EMBODY VALUES ← → CAN NOT EMBODY VALUES
- CAN NOT (AA DO NOT ← → HAVE INTENTIONALITY) HAVE INTENTIONALITY
- CAN ADAPT BEHAVIOR ← → CAN ADAPT BEHAVIOR FROM EXTERNAL INPUTS AND INTERACTIONS

GIVE AI PRINCIPLES, CAN BE BETTER THAN HUMANS AS MORAL REASONING (NOT CARRIED AWAY BY EMOTIONS). AGAINST THE IDEA

- MORAL RULES OFTEN CONFLICTS
- MORALITY CAN NOT BE REDUCED TO FOLLOWING RULES, NOT JUST A MATTER OF RATIONALITY

STRUCTURE OF INT. AGENT:

- **ARCHITECTURE:** MACHINERY/ DEVICE OF ACTUATORS - SENSORS
- **AGENT FUNCTION:** ACTIONS ARE MAPPED FROM A CERTAIN SEQUENCE
- **AGENT PROGRAM:** PRODUCTION OF THE AGENT FUNCTION

CONDITIONS FOR MORAL AGENCY

- 1) • AGENT WITH INTERNAL STATE: CONSISTS OF DESIRES, BELIEFS AND OTHER INTENTIONAL STATES. CONSTITUTE A REASON TO ACTING
• ONE INTENTIONAL STATE IS A MENTAL STATE (INTENTIONAL STATE + INTENT/REASON TO ACT)
- 2) • THERE IS AN OUTWARD EMBODIED EVENT. THE AGENT DO SOMETHING AND MOVE HIS BODY IN SOME WAY
- 3) • THE INTERNAL STATE IS THE CAUSE OF THE OUTWARD EVENT (THE MOVEMENT OF THE BODY IS RATIONALLY DIRECTED AT SOME STATE OF THE WORLD)
- 4) • THE OUTWARD BEHAVIOR HAVE AN OUTWARD EFFECT
- 5) • THE EFFECT HAVE TO BE ON A RECIPIENT THAT CAN BE HARMED OR HELPED
- A COMPUTER SATISFY 2-5 BUT MACHINE DO NOT HAVE THE FREEDOM OR HUMAN INTENTIONALITY THAT WOULD MAKE THEM MORALLY RESPONSIBLE
- COMPUTER SYSTEM DO NOT HAVE INTENDING TO ACT, DO NOT HAVE HUMAN INTENTIONALITY AND THIS IS THE KEY TO UNDERSTANDING THE MORAL CHARACTER OF COMPUTER SYSTEMS

HUMAN BIASES WERE TRANSFERRED AND AMPLIFIED BY THE AI

IF X IS VALUABLE → THEN ONE HAVE REASONS FOR A POSITIVE RESPONSE (PRO-ATTITUDE OR PRO-BEHAVIOUR)

- **VALUE SENSITIVE DESIGN:** WE CAN SOMEHOW DESIGN MORAL VALUES IN TECHNICAL ARTIFACTS SO THAT THEY CAN EMBODY VALUES → TECHNICAL ARTIFACTS MORALLY VALUE-LADEN BY DESIGNING THEM IN A WAY
- **DEONTIC NOTIONS:** DETERMINE THE RIGHTEOUSNESS/WRONGNESS OF ACTIONS. ENTITY EMBODIES A VALUE IF THERE ARE REASONS FOR PRO-ATTITUDE/ PRO-BEHAVIOUR TOWARD THAT ENTITY

VALUES EVALUATES TECHNOLOGICAL ARTIFACTS IN TERMS OF GOODNESS AND BADNESS. 3 TYPES OF VALUES

- 1) **INTENDED VALUES:** THE VALUES WHICH DESIGNERS AIM TO EMBODY IN THEIR DESIGN AND WHICH THEY HOPE TO BE REALIZED IN PRACTICE

- 2) **REALIZED VALUE:** VALUE THAT IS REALIZED BY A TECHNICAL ARTIFACT IN A PRACTICAL CONTEXT
- 3) **EMBODED VALUE:** POTENTIAL TO REALIZE A VALUE IN A APPROPRIATE CONTEXT

PROBLEM: CREATING MORAL MACHINES IN THE SENSE THAT MACHINES THEMSELVES CAN MAKE ETHICAL PREDICTION

THE MIRROR VIEW - COECKELBERG (2020)

- USE A DATASET THAT MIRRORS THE REAL WORD. ALGORITHM MAY MODEL EXISTING BIAS, NOT A PROBLEM FOR DEVELOPERS
- BIAS AND DISCRIMINATION ARE UNFAIR, ONE SHOULD CHANGE THE DATASET TO PROMOTE AFFIRMATIVE ACTION

MORAL STATUS OF AI

WHAT IS AI CAPABLE OF DOING MORALLY SPEAKING?

- AGENT: BEING THAT CAN ACT
- AGENCY: EXERCISE/MANIFESTATION OF THIS CAPACITY
- MORAL AGENCY: INDIVIDUAL'S ABILITY TO MAKE MORAL CHOICES (BASED ON RIGHTE-WRONG) RESPONSIBLE FOR THEM

HOW SHOULD WE TREAT AN AI PATIENTY: QUALITY OR STATE OF BEING PATIENT OR PASSIVE

CATEGORIES OF ETHICAL AGENTS (MOOR)

- **ETHICAL IMPACTS AGENTS:** ANY MACHINE THAT CAN BE EVALUATED FOR ITS ETHICAL CONSEQUENCES (PROBLEM: ALL ROBOTS HAVE ETHICAL IMPACTS)
- **IMPLICIT ETHICAL AGENTS:** MACHINES DESIGNED SO THAT THEY DO NOT HAVE NEGATIVE ETHICAL EFFECTS (PROBLEM: ALL ROBOTS SHOULD BE DESIGNED TO BE IMPLICIT ETHICAL AGENTS)

COMPUTER SYSTEMS HAVE MEANING, SIGNIFICANCE ONLY IN RELATION TO HUMANS. SOCIO-TECHNICAL SYSTEM

- **EXPLICIT ETHICAL AGENTS:** USING ETHICAL CATEGORIES AS PART OF THEIR INTERNAL PROGRAMMING

- **FULL ETHICAL AGENTS:** CAN MAKE EXPLICIT MORAL JUDGEMENTS AND ARE GENERALLY QUITE COMPETENT IN JUSTIFY SUCH DECISION

NATURAL PHENOMENA/ENTITY

- APPEAR IN NATURE INDEPENDENT OF HUMAN BEHAVIOUR
- HAVE A FUNCTIONALITY (BUT NOT DESIGNED BY HUMANS)
- DO NOT HAVE INTENTIONALITY
- **TECHNOLOGY**
- A TOOL FOR HUMAN ENDS, DESIGNED BY HUMANS TO BEHAVE IN A SPECIFIC WAY
- **HUMAN-MADE OBJECTS**
- HAVE A FUNCTIONALITY
- HAVE BEEN INTENTIONALITY DESIGN

DISTINCTION BETWEEN ARTIFACTS AND TECHNOLOGY

- **ARTIFACTS:** PRODUCT OF HUMAN CONTRIVANCE, DOES NOT EXIST WITHOUT SYSTEMS OF KNOWLEDGE
- HAVE MEANING ONLY IN THE CONTEXT OF HUMAN SOCIAL ACTIVITY
- ARE ABSTRACTION FROM REALITY/THE SYSTEM
- TO DEFINE AN ARTIFACT WE MUST DO A MENTAL ACT OF SEPARATING THE OBJECT FROM THE CONTEXT
- **TECHNOLOGY:** SOCIO-TECHNICAL SYSTEM
- APPLICATION OF SCIENTIFIC KNOWLEDGE, TOOLS AND TECHNIQUES TO PRACTICAL TASKS AND PROBLEMS
- CAN INCLUDE BOTH PHYSICAL OBJECTS (MACHINES, DEVICES,...)
- CAN INCLUDE INTANGIBLE PROCESS (ALGORITHMS, SOFTWARE,...)

INTENTIONALITY OF COMPUTER SYSTEM

- INTENTIONALITY OF THE COMPUTER SYSTEM DESIGNER
- INTENTIONALITY OF THE SYSTEM

INTENTIONALITY USER

- ALL HAVE INTENTIONALITY-EFFICACY. THEY BEHAVE FROM NECESSITY WHEN CREATED. COMPUTER SYSTEM COMPONENTS IN MORAL ACTION

- **TECHNICAL NORM:** COMPUTER CODES THAT REGULATE THE BEHAVIOR OF AN INTERACTIONS OF AA, TWO WAYS
- **OFFLINE DESIGN:** NORM ARE SPECIFIED AND ENCODED IN THE AGENT BY HUMANS
- **AUTONOMOUSLY DISCOVER, INVENT OR SPREAD NORMS:** NORMS CAN BE PICKED FROM THE ENVIRONMENT IN DIFFERENT WAYS OR FROM INTERACTIONS WITH AGENTS AND/OR HUMANS

VAN DE POEL: TECHNICAL NORMS (N)

- EMBODY VALUE (V) IF
- N EMBODIES V
- THE EXECUTION OF N WITHIN THE SYSTEM IS CONDUCTIVE TO V

CHARACTERISTIC OF INTELLIGENT AGENT

- **AUTONOMY:** ENABLES TO PERFORM CERTAIN TASK ON THEIR OWN
- **LEARNING ABILITY:** ENABLES TO LEARN
- **INTERACT WITH A NETWORK OF OTHER AGENTS:** HUMAN/MON HUMAN (NATURAL/ARTIFICIAL)
- **CAN LEARN NEW RULES INCREMENTALLY**
- **GOAL-ORIENTED HABITS**

SYSTEMS HAVE INTENTIONALITY. CONNECTED TO INTENTIONALITY OF THE USER AND OF THE DESIGNER. SYSTEM INTENTIONALITY IS LATENT WITHOUT USER INTENTIONALITY

DEPENDENCE OF COMPUTER SYSTEM BEHAVIOR ON HUMAN BEHAVIOR

- SINCE COMPUTER SYSTEMS HAVE BUILD-IN INTENTIONALITY ONCE DEPLOYED THEY CAN BEHAVE INDEPENDENTLY AND WITHOUT HUMAN INTERVENTION

- COMPUTER SYSTEM ARE NOT-NEUTRAL AND CLOSE (BUT NOT) MORAL AGENTS. THEY ARE MORAL ENTITIES (SO THEY ARE PART OF THE MORAL WORLD)

- COMPUTER ARE NOT MORAL AGENTS: NO MENTAL STATE - INTENDING TO ACT. BUT THEY ARE CLOSE TO BEING MORAL AGENTS: INTENTIONALITY CREATED-USING FORMS OF INTENTIONALITY/EFFICACY. THEY ARE PART OF THE MORAL WORLD; BECAUSE OF THEIR EFFECTS - WHAT THEY ARE AND WHAT THEY DO

FUNCTIONAL MORALITY (WALLACH-ALLEN)

HOW CAN WE CREATE AN ARTIFICIAL MORAL AGENT?

- MACHINES WITH AUTONOMY AND SENSITIVITY TO VALUES HAVE FUNCTIONAL MORALITY

FUNCTIONAL MORALITY (FM)

- AUTOPILOTS AND ETHICAL DECISIONS SUPPORT SYSTEMS

IS TECHNOLOGICALLY POSSIBLE

- CAPACITY OF MACHINES FOR ASSESSING AND RESPONDING TO MORAL CHALLENGES

- MACHINE MUST HAVE AUTONOMY AND SENSITIVITY TO VALUES

OPERATIONAL MORALITY (OM)

- GUNS WITH CHILDPROOF MECHANISM, IT LACKS OF AUTONOMY AND SENSITIVITY BUT ITS DESIGN EMBODY VALUES OF NRA CODE OF ETHICS TO ENDORSE

(WHAT SHOULD WE DO? OUR DUTY)

DIFFERENT ETHICAL THEORIES IN THE PROCESS OF DESIGN

→ DIFFERENT ETHICAL THEORIES AND APPROACHES TO PRODUCE MACHINE CAPABLE OF DECISION-MAKING

TOP-DOWN APPROACH:

- SET OF RULES THAT CAN BE TURNED IN A ALGORITHM

UTILITARIANISM: MAXIMISE THE TOTAL AMOUNT OF UTILITY IN THE WORLD. BEST ACTION THE ONE THAT MAXIMISE AGGREGATE UTILITY

HEAVY COMPUTATION: NEEDS TO WORK OUT MANY/ALL THE CONSEQUENCES IN ORDER TO RANK ACTIONS MORALLY

HOW TO CONSTRUCT A FUNCTION THAT CAN EVACUATE WITH THE CORRECT WEIGHT THE PRESENT BENEFITS AGAINST THE FUTURE ARMS?

HOW TO ASSIGN NUMBERS TO SOMETHING SUBJECTIVE LIKE HAPPINESS, PLEASURE OR DESIRABILITY?

DEONTOLOGICAL: EVALUATES THE ETHICAL CORRECTNESS OF ACTIONS ON THE BASIS OF CHARACTERISTIC THAT AFFECT THE ACTION IN ITSELF

ANY LIST OF DUTIES MIGHT HAVE INTERNAL CONFLICT

DESIGNER NEED TO BE SURE THAT THE RULES ARE ACTIVATED WHEN A SITUATION REQUIRES THEIR APPLICATION

FORMULATE AN ARCHITECTURE FOR MANAGING SITUATIONS WHEN RULES CONFLICTS

MEASURE OF HAPPINESS,
WELL BEING

TURING CHILD MIND: SUBJECT TO THE COURSE OF EDUCATION ⇒ OBTAIN THE ADULT BRAIN

AND MORALLY EVIL IF
SOME ACTIONS VIOLATES IT

MINDLESS MORALITY (FLORIDI-SANDERS)

MAKE MORAL AGENTS SUFFICIENT LEVEL OF:

INTERACTIVITY → RESPONSE TO THE STIMULUS BY CHANGING STATE

AUTONOMY → ABILITY TO CHANGE STATE: W/O STIMULUS

ADAPTIVITY → ABILITY TO CHANGE THE "TRANSITION RULE" BY WHICH STATE IS CHANGED

CAPABLE OF MORAL QUALIFIABLE ACTIONS

AN AGENT IS MORALLY GOOD IF ALL HIS ACTIONS RESPECT THIS THRESHOLD

CONTEL VS FUNCTIONALIST APPROACH

BOTTOM-UP APPROACH:

MIMIC CHILD'S MORAL DEVELOPMENT

DEVELOPMENT NOT EXPLICITLY GUIDED BY ANY ETHICAL THEORY

DESIGN METHODS FOR CARRYING OUT DISCRETE TASKS ⇒ MORE COMPLEX ACTIVITIES/GREATER AUTONOMY

NEED SOME METRIC FOR EVACUATION IF A DEVICE IS CAPABLE OF MAKING THE APPROPRIATE MORAL DECISION IN A PARTICULAR SITUATION

ANTROPOTCENTRISM

SLAVE ETHICS

MERGING TOP-DOWN AND BOTTOM-UP APPROACH:

MORAL APPROACH: COECKELBERG PROPOSED QUASI-MORAL ROBOTS: APPEAR TO BE MORAL

NO PROOF THAT ANOTHER PERSON HAVE MENTAL STATES/CONSCIOUSNESS. BUT WE INTERPRET THE OTHER'S APPEARANCE/BEHAVIOR = EMOTION. THE OTHER HAVE VIRTUAL SUBJECTIVITY/QUASI-SUBJECTIVITY, INTERACTING AS IF OUR APPEARANCE/BEHAVIOR IN THEIR CONSCIOUSNESS

SAME WITH ROBOTS: IF SUFFICIENTLY ADVANCED (IMITATE SUBJECTIVITY/CONSCIOUSNESS IN A CONVINCING WAY) THEY BECOME/APPEAR THE QUASI-OTHER

RELATIONAL APPROACH: COECKELBERG. MORAL STANDING: NOT AN OBJECTIVE PROPERTY. HUMANS ATTRIBUTE THE MORAL STATUS ⇒ STATUS OF THE OBJECT OF MORAL STANDING. DEPENDENT FROM HUMAN SUBJECTIVITY ⇒ ROBOT CAN APPEAR IN DIFFERENT WAYS TO DIFFERENT PEOPLE IN DIFFERENT SITUATIONS/CONTEXT

WHAT MATTERS MORALLY SPEAKING, IS HOW THE ENTITY APPEARS

ETHICAL ATTENTION SHIFTED

FROM ONTOLOGY (PROPERTIES OF AN ENTITY) TO EPISTEMEOLOGY (OUR KNOWLEDGE OF THE ENTITY)

FROM OBJECT TO SUBJECT

FROM WHAT THINGS REALLY ARE TO HOW WE LOOK AT THINGS

WHAT IS AUTONOMY?

FREEDOM FROM EXTERNAL CONTROL/INFLUENCE = INDEPENDENCE

CAPACITY TO MAKE AN INFORMED/UNCOERCED DECISION = FREE DECISION

KANT & RATIONAL AUTONOMY. ONE POSSESSES THE MOTIVATION TO GOVERN THEIR OWN LIFE

SULLINS REQUIREMENTS FOR FULL MORAL AGENCY

AUTONOMOUS FROM PROGRAMMER/OPERATOR/USER

AUTONOMY: MACHINE NOT UNDER THE DIRECT CONTROL OF ANY OTHER PROGRAMMER OR USER. WHEN THE AGENCY CAUSES HARM OR GOOD IN A MORAL SENSE WE CAN SAY THE MACHINE HAVE MORAL AGENCY

IF THE ROBOT ACTIONS IS EFFECTIVE IN ACHIEVING THE GOALS/TASKS OF THE ROBOT ⇒ THEN THE ROBOT HAS EFFECTIVE AUTONOMY

INTENTIONALITY: IF THE COMPLEX INTERACTION OF THE ROBOT'S PROGRAMMING AND ENVIRONMENT CAUSES THE MACHINE TO ACT IN A MORALLY HARMFUL OR BENEFICIAL WAY AND THE ACTIONS SEEMS DELIBERATE AND CALCULATED THEN THE MACHINE IS A MORAL AGENT

RESPONSIBILITY: ROBOT BEHAVES RESPONSIBLE TO OTHER MORAL AGENTS. IF A ROBOT BEHAVES IN THIS WAY. FULLFILLS SOCIAL ROLE OF SOME RESPONSIBILITIES ⇒ IT HAS THE DUTY TO CARE FOR ITS PATIENTS ⇒ MORAL AGENT

DO NOT PROVE IT HAS INTENTIONALITY
ROBOTS ARE MORAL AGENTS WHEN THERE IS A REASONABLE LEVEL OF ABSTRACTION UNDER WHICH WE MUST GRANT THAT THE MACHINE HAS AUTONOMOUS INTENTIONS AND RESPONSIBILITIES