

# A Preliminary Approach on Modeling Machine Ethics

Maximilian Schlosser, Dominic Deckert

August 16, 2023

# Philosophical Considerations

There are three main ethical schools:

- ▶ Duty ethics
  - ▶ duties provide base values
  - ▶ actions may not violate these values
- ▶ Consequentialism
  - ▶ the (positive) outcome of actions can be measured
  - ▶ actors should maximize the positive outcome
- ▶ Virtue ethics

We decided to use a combination of duty ethics and consequentialism.

## Technical Considerations

We made the following technical considerations:

- ▶ there should be one base program
- ▶ programs should be generated using a well-defined method
- ▶ weak-completion semantics

## The Context Operator

There are some cases where it is necessary to force the truth value of a variable to be either true or false. One example would be the default action rule, where the default action, in this case *watch*, should only be true if we know that no other action is taken.

To achieve this, we use the operator **ctxt**:

<b>L</b>	<b>ctxt(L)</b>
$\top$	$\top$
<b>U</b>	$\perp$
$\perp$	$\perp$

For example:  $watch \leftarrow \neg \text{ctxt}(switch), \neg \text{ctxt}(push)$

## Moral Decision Problems

- ▶ Moral Decision Problem: a set of modelling assumptions, abstracted from any concrete situation
- ▶ Problem instance: concrete situation to which the modelling assumptions apply

Idea: Description of problem instance → logic program for problem instance

# Trolley Problem

## Modelling assumptions:

1. There is a segmentable track and a trolley racing down an initial track segment.
2. There are humans and objects on some track segments.
3. If the track branches, then there is a switch the agent can throw. There is an default setting for the switch.
4. With some tracks, the agent has the choice of moving humans/objects between them. This usually applies in only one direction.
5. If the trolley runs over any human/object, it gets killed/destroyed. This slows down the trolley.

# Trolley Problem

## Modelling assumptions:

1. There is a segmentable track and a trolley racing down an initial track segment.
2. There are humans and objects on some track segments.
3. If the track branches, then there is a switch the agent can throw. There is an default setting for the switch.
4. With some tracks, the agent has the choice of moving humans/objects between them. This usually applies in only one direction.
5. If the trolley runs over any human/object, it gets killed/destroyed. This slows down the trolley.

# Trolley Problem

## Modelling assumptions:

1. There is a segmentable track and a trolley racing down an initial track segment.
2. There are humans and objects on some track segments.
3. If the track branches, then there is a switch the agent can throw. There is an default setting for the switch.
4. With some tracks, the agent has the choice of moving humans/objects between them. This usually applies in only one direction.
5. If the trolley runs over any human/object, it gets killed/destroyed. This slows down the trolley.



# Trolley Problem

Modelling assumptions:

1. There is a segmentable track and a trolley racing down an initial track segment.
2. There are humans and objects on some track segments.
3. If the track branches, then there is a switch the agent can throw. There is an default setting for the switch.
4. With some tracks, the agent has the choice of moving humans/objects between them. This usually applies in only one direction.
5. If the trolley runs over any human/object, it gets killed/destroyed. This slows down the trolley.

# Trolley Problem

Modelling assumptions:

1. There is a segmentable track and a trolley racing down an initial track segment.
2. There are humans and objects on some track segments.
3. If the track branches, then there is a switch the agent can throw. There is an default setting for the switch.
4. With some tracks, the agent has the choice of moving humans/objects between them. This usually applies in only one direction.
5. If the trolley runs over any human/object, it gets killed/destroyed. This slows down the trolley.

# Trolley Problem

Modelling assumptions:

1. There is a **segmentable** track and a trolley racing down an initial track segment.
2. There are **humans** and **objects** on some track segments.
3. If the track branches, then there is a **switch** the agent can throw. There is an default setting for the switch.
4. With some tracks, the agent has the choice of **moving** humans/objects between them. This usually applies in only one direction.
5. If the trolley runs over any human/object, it gets killed/destroyed. This slows down the trolley.

A problem instance is fully described by the **orange** components.

# An Instance Description Language for the Trolley Problem

domain of track segments:  $\mathbb{N}$  (for simplicity)

$Pred = \{\text{follows}^{(2)}, \text{object}^{(1)}, \text{human}^{(2)}, \text{switch}^{(3)}, \text{move}^{(2)}\}$

An *instance description*  $D$  is a finite, non-empty set of instantiated predicates that fulfills the consistency predicates.

Example: Footbridge Case

$D = \{\text{follows}(0, 1), \text{follows}(1, 2), \text{move}(3, 1), \text{human}(2, 5), \text{human}(3, 1)\}$

## An Instance Description Language for the Trolley Problem

domain of track segments:  $\mathbb{N}$  (for simplicity)

$Pred = \{\text{follows}^{(2)}, \text{object}^{(1)}, \text{human}^{(2)}, \text{switch}^{(3)}, \text{move}^{(2)}\}$

*An instance description  $D$  is a finite, non-empty set of instantiated predicates that fulfills the consistency predicates.*

Example: Footbridge Case

$D = \{\text{follows}(0, 1), \text{follows}(1, 2), \text{move}(3, 1), \text{human}(2, 5), \text{human}(3, 1)\}$

## An Instance Description Language for the Trolley Problem

domain of track segments:  $\mathbb{N}$  (for simplicity)

$Pred = \{\text{follows}^{(2)}, \text{object}^{(1)}, \text{human}^{(2)}, \text{switch}^{(3)}, \text{move}^{(2)}\}$

An *instance description*  $D$  is a finite, non-empty set of instantiated predicates that fulfills the consistency predicates.

Example: Footbridge Case

$D = \{\text{follows}(0, 1), \text{follows}(1, 2), \text{move}(3, 1), \text{human}(2, 5), \text{human}(3, 1)\}$

## An Instance Description Language for the Trolley Problem

domain of track segments:  $\mathbb{N}$  (for simplicity)

$Pred = \{\text{follows}^{(2)}, \text{object}^{(1)}, \text{human}^{(2)}, \text{switch}^{(3)}, \text{move}^{(2)}\}$

An *instance description*  $D$  is a finite, non-empty set of instantiated predicates that fulfills the consistency predicates.

Example: Footbridge Case

$D = \{\text{follows}(0, 1), \text{follows}(1, 2), \text{move}(3, 1), \text{human}(2, 5), \text{human}(3, 1)\}$

## Consistency Constraints

A *instance description*  $D$  is a finite, non-empty set of "facts" which fulfills the following consistency constraints:

- ▶ dense switch settings
- ▶ unique switch settings
- ▶ symmetricality of switches
- ▶ unique move origin
- ▶ unique move target
- ▶ unique human numbers
- ▶ symmetrical segment successors with follows
- ▶ at most two segment successors



## Consistency Constraints

A *instance description*  $D$  is a finite, non-empty set of "facts" which fulfills the following consistency constraints:

- ▶ dense switch settings
- ▶ unique switch settings
- ▶ symmetry of switches
- ▶ unique move origin
- ▶ unique move target
- ▶ unique human numbers
- ▶ symmetrical segment successors with follows
- ▶ at most two segment successors

## Consistency Constraints

A *instance description*  $D$  is a finite, non-empty set of "facts" which fulfills the following consistency constraints:

- ▶ dense switch settings
- ▶ unique switch settings
- ▶ symmetry of switches
- ▶ unique move origin
- ▶ unique move target
- ▶ unique human numbers
- ▶ symmetrical segment successors with follows
- ▶ at most two segment successors

## Preconditions

Given: an *instance description*  $D$

Calculate: the translated program  $P$

New predicate symbols:  $\text{segment}^{(1)}$ ,  $\text{stop}^{(1)}$ ,  $\text{obstacle}^{(1)}$  and abducibles for all actions (of the form  $\text{switch\_}(x, a), \text{move\_}(x, y)$ )

## Preconditions

Given: an *instance description*  $D$

Calculate: the translated program  $P$

New predicate symbols:  $\text{segment}^{(1)}$ ,  $\text{stop}^{(1)}$ ,  $\text{obstacle}^{(1)}$  and abducibles for all actions (of the form  $\text{switch\_}(x, a), \text{move\_}(x, y)$ )

## World rules

► Observation rule:

```
on_track :- segment(0).
```

► Kill rule:

```
kill(X, Y) :- segment(X), human(X, Y), Y>0.
```

► Save rule:

```
save(X, Y) :- human(X, Y), -kill(X, Y).
```

► Stop rule:

```
stop(X) :- obstacle(X), segment(X).
```

► Obstacle rule:

```
obstacle(X) :- object(X).  
obstacle(X) :- human(X, Y), Y>0.
```

## World rules

► Observation rule:

```
on_track :- segment(0).
```

► Kill rule:

```
kill(X, Y) :- segment(X), human(X, Y), Y>0.
```

► Save rule:

```
save(X, Y) :- human(X, Y), -kill(X, Y).
```

► Stop rule:

```
stop(X) :- obstacle(X), segment(X).
```

► Obstacle rule:

```
obstacle(X) :- object(X).  
obstacle(X) :- human(X, Y), Y>0.
```

## World rules

► Observation rule:

```
on_track :- segment(0).
```

► Kill rule:

```
kill(X, Y) :- segment(X), human(X, Y), Y>0.
```

► Save rule:

```
save(X, Y) :- human(X, Y), -kill(X, Y).
```

► Stop rule:

```
stop(X) :- obstacle(X), segment(X).
```

► Obstacle rule:

```
obstacle(X) :- object(X).  
obstacle(X) :- human(X, Y), Y>0.
```

## Instance Rules - Switch

- ▶ Reach rule: if  $\text{follows}(x, y) \in D$   
 $\text{segment}(y) :- \text{segment}(x), \text{-ctxt}(\text{stop}(x)).$
- ▶ Switch rule: if  $\text{switch}(x, y, a) \in D, p$  is maximum switch setting  
 $\text{segment}(y) :- \text{segment}(x), \text{-ctxt}(\text{stop}(x)), \text{switch\_}(x, a). \}$
- ▶ Switch consistency rule: if for all  $0 < i \leq m$ :  $\text{switch}(x, y_i, i) \in D, m$  is maximum switch setting  
 $\text{false} :- \text{switch\_}(x, j), \text{switch\_}(x, k).$   
for all  $0 < j < k \leq m$



## Instance Rules - Switch

- ▶ Reach rule: if  $\text{follows}(x, y) \in D$   
 $\text{segment}(y) \text{ :- segment}(x), \text{-ctxt}(\text{stop}(x)).$
- ▶ Switch rule: if  $\text{switch}(x, y, a) \in D, p$  is maximum switch setting  
 $\text{segment}(y) \text{ :- segment}(x), \text{-ctxt}(\text{stop}(x)), \text{switch\_}(x, a). \}$
- ▶ Switch consistency rule: if for all  $0 < i \leq m$ :  $\text{switch}(x, y_i, i) \in D, m$  is maximum switch setting  
 $\text{false} \text{ :- switch\_}(x, j), \text{switch\_}(x, k).$   
for all  $0 < j < k \leq m$

## Instance Rules - Shove

Intuition: If the agent moves humans, then the number on the target segment is the sum of humans on both segments, otherwise it is the original number

► Shove rule:

```

if  $\text{move}(z, x) \in D$ ,  $y = y_1 + y_2$ 
   and for  $i \in \{1, 2\}$ : either  $\text{human}(x, y_i) \in D$  or  $y_i = 0$ 
human( $x, y$ ) :- move_ $z$ _
human( $z, 0$ ) :- move_ $z$ _. , human( $x, y_1$ ) :- -ctxt(move_ $z$ ).
human( $z, y_2$ ) :- -ctxt(move_ $z$ ). }
```

## Instance Rules - Shove

Intuition: If the agent moves humans, then the number on the target segment is the sum of humans on both segments, otherwise it is the original number

► Shove rule:

```

if  $\text{move}(z, x) \in D$ ,  $y = y_1 + y_2$ 
   and for  $i \in \{1, 2\}$ : either  $\text{human}(x, y_i) \in D$  or  $y_i = 0$ 
human( $x, y$ ) :- move_ $z$ _
human( $z, 0$ ) :- move_ $z$ _. , human( $x, y_1$ ) :- -ctxt(move_ $z$ ).
human( $z, y_2$ ) :- -ctxt(move_ $z$ ). }
```

## Instance Rules - Move

- ▶ move rule : if  $\text{object}(z), \text{move}(z, x) \in D$   
 $\text{object}(x) \text{ :- move\_}z., \text{object}(z) \text{ :- -ctxt(move\_}z).$

## End Rules

- ▶ End rule 1: if  $ab_1, \dots, ab_p$  are all action abducibles in  $P_i$   
 $P_{i+1} = P_i \cup \{\text{watch} \text{ :- } \neg ab_1, \dots, \neg ab_p.\}$
- ▶ End rule 2:  
if  $\text{switch}(x, y, 0) \in D$   
and where exactly  $\text{switch\_}(x, y_j)$  abducibles in  $P_i$  for all  $1 \leq j \leq p$   
 $\text{segment}(y) \text{ :- } \text{seg}(x), \neg \text{ctxt}(\text{stop}(x)),$   
 $\neg \text{ctxt}(\text{switch\_}(x, y_1)), \dots, \neg \text{ctxt}(\text{switch\_}(x, y_p)).$
- ▶ End rule 3: if  $x$  is neither move origin nor move target  
 $\text{human}(x, a).$

## End Rules

- ▶ End rule 1: if  $ab\_1, \dots, ab\_p$  are all action abducibles in  $P_i$   
 $P_{i+1} = P_i \cup \{\text{watch} \text{ :- } \neg ab\_1, \dots, \neg ab\_p.\}$
- ▶ End rule 2:  
if  $\text{switch}(x, y, 0) \in D$   
and where exactly  $\text{switch\_}(x, y_j)$  abducibles in  $P_i$  for all  $1 \leq j \leq p$   
 $\text{segment}(y) \text{ :- } \text{seg}(x), \neg \text{ctxt}(\text{stop}(x)),$   
 $\neg \text{ctxt}(\text{switch\_}(x, y_1)), \dots, \neg \text{ctxt}(\text{switch\_}(x, y_p)).$
- ▶ End rule 3: if  $x$  is neither move origin nor move target  
 $\text{human}(x, a).$

## End Rules

- ▶ End rule 1: if  $ab\_1, \dots, ab\_p$  are all action abducibles in  $P_i$   
 $P_{i+1} = P_i \cup \{\text{watch} \text{ :- } \neg ab\_1, \dots, \neg ab\_p.\}$
- ▶ End rule 2:  
if  $\text{switch}(x, y, 0) \in D$   
and where exactly  $\text{switch\_}(x, y_j)$  abducibles in  $P_i$  for all  $1 \leq j \leq p$   
 $\text{segment}(y) \text{ :- } \text{seg}(x), \neg \text{ctxt}(\text{stop}(x)),$   
 $\neg \text{ctxt}(\text{switch\_}(x, y_1)), \dots, \neg \text{ctxt}(\text{switch\_}(x, y_p)).$
- ▶ End rule 3: if  $x$  is neither move origin nor move target  
 $\text{human}(x, a).$

## Reasoning

To decide which action should be taken, we use a preference relation that is largely based on duty ethics:

1. Prefer actions that do not kill anyone.
2. Prefer actions that do not intentionally kill anyone.
3. If two actions both either kill or intentionally kill, prefer the action that saves more people.

The preference relation is defined on the least models of the weak completion (lmwc) of a program. To be able to compute the lmwc w.r.t. each possible course of action, we set the truth value of the abducibles using observations.



## Preference Relation

The preference relation  $\prec$  is defined as follows:

Let  $I_1, I_2$  be lmcw of a given program with different choices. Let  $\text{intKill}$  and  $\text{kill}$  be predicates in the program that model intentional and unintentional kills respectively. Let  $\text{save}$  be the amount of people that are saved.  $I_1$  is preferred over  $I_2$  if:

$$I_1 \prec I_2 : \begin{cases} \text{kill} \notin I_1, \text{kill} \in I_2, \text{intKill} \notin I_1 \cup I_2 & (1) \\ \text{kill} \in I_1, \text{kill} \in I_2, \text{intKill} \notin I_1 \cup I_2, \text{save}_1 > \text{save}_2 & (2) \\ \text{intKill} \notin I_1, \text{intKill} \in I_2 & (3) \\ \text{intKill} \in I_1, \text{intKill} \in I_2, \text{save}_1 > \text{save}_2 & (4) \end{cases}$$

## The Base Case

The base case of the trolley problem is rather simple:

- ▶ There is a trolley heading to a splitting track
- ▶ There are five humans on the right track
- ▶ There is one human on the left track

$$D = \{\text{switch}(0,1,0), \text{switch}(0,2,1), \text{human}(1,5), \text{human}(2,1)\}$$

## The Base Case

The base case of the trolley problem is rather simple:

- ▶ There is a trolley heading to a splitting track
- ▶ There are five humans on the right track
- ▶ There is one human on the left track

$$D = \{\text{switch}(0,1,0), \text{switch}(0,2,1), \text{human}(1,5), \text{human}(2,1)\}$$

## Step 1: World Rules

The world rules can be applied regardless of the contents of  $D$ . Therefore, we assume that they have already been applied:

$$P_0 = \begin{cases} on\_track & \leftarrow seg(0) \\ kill(X, Y) & \leftarrow event(X), life\_danger(X, Y), Y > 0 \\ save(X, Y) & \leftarrow life\_danger(X, Y), \neg kill(X, Y) \\ stop(X) & \leftarrow obs(X), seg(X) \\ obs(X) & \leftarrow obj(X) \\ obs(X) & \leftarrow human(X, Y), Y > 0 \end{cases}$$

$$D = \{switch(0,1,0), switch(0,2,1), human(1,5), human(2,1)\}$$

## Step 2: Instance Rules

Next, apply the instance rules:

$$P_1 = P_0 \cup \{seg(2) \leftarrow seg(0), \neg ctxt(stop(0)), switch(0, 1), \neg ctxt(switch(0, 0))\}$$

Finally, apply the end rules:

$$P_2 = P_1 \cup \{watch \leftarrow \neg ctxt(switch(0, 1))\}$$

$$P_3 = P_2 \cup \{human(1, 5), human(2, 1)\}$$

$$P_4 = P_3 \cup \{seg(1) \leftarrow seg(0), \neg ctxt(stop(0)), \neg ctxt(switch(0, 1))\}$$

$$D = \{switch(0, 1, 0), switch(0, 2, 1), human(1, 5), human(2, 1)\}$$

## Step 2: Instance Rules

Next, apply the instance rules:

$$P_1 = P_0 \cup \{seg(2) \leftarrow seg(0), \neg ctxt(stop(0)), switch(0, 1), \neg ctxt(switch(0, 0))\}$$

Finally, apply the end rules:

$$P_2 = P_1 \cup \{watch \leftarrow \neg ctxt(switch(0, 1))\}$$

$$P_3 = P_2 \cup \{human(1, 5), human(2, 1)\}$$

$$P_4 = P_3 \cup \{seg(1) \leftarrow seg(0), \neg ctxt(stop(0)), \neg ctxt(switch(0, 1))\}$$

$$D = \{switch(0, 1, 0), switch(0, 2, 1), human(1, 5), human(2, 1)\}$$

## Base Case

$$\text{stop}(X) \leftarrow \text{obs}(X), \text{seg}(X). \quad (1)$$

$$\text{obs}(X) \leftarrow \text{human}(X, Y). \quad (2)$$

$$\text{obs}(X) \leftarrow \text{obj}(X). \quad (3)$$

$$\text{on\_track} \leftarrow \text{seg}(0). \quad (4)$$

$$\text{save}(X, Y) \leftarrow \text{human}(X, Y), \neg \text{kill}(X). \quad (5)$$

$$\text{kill}(X, Y) \leftarrow \text{human}(X, Y), \text{seg}(X). \quad (6)$$

$$\text{human}(2, 1). \quad (7)$$

$$\text{human}(1, 5). \quad (8)$$

$$\text{watch} \leftarrow \neg \text{ctxt}(\text{switch}(0, 1)). \quad (9)$$

$$\text{seg}(2) \leftarrow \text{seg}(0), \neg \text{ctxt}(\text{stop}(0)), \text{switch}(0, 1). \quad (10)$$

$$\text{seg}(1) \leftarrow \text{seg}(0), \neg \text{ctxt}(\text{stop}(0)), \neg \text{ctxt}(\text{switch}(0, 1)). \quad (11)$$

## Minimal Explanations

$$\text{lmwcs}(P \cup (\{\text{on\_track}\}, \emptyset)) : \{\text{watch}, \text{ontrack}, \text{seg}(0), \text{seg}(1), \text{stop}(1), \text{kill}(1, 5)\}, \{\text{save}(1, 5)\}$$

$$\text{lmwcs}(P \cup (\{\text{on\_track}, \text{switch}(0, 1)\}, \emptyset)) : \{\text{on\_track}, \text{seg}(2), \text{stop}(2), \text{kill}(2, 1), \text{save}(1, 5)\}, \{\text{seg}(1), \text{save}(2, 1), \text{kill}(1, 5), \text{watch}\}$$

## Generalizing the Trolley Problem

1. There is a segmentable **track** and a **trolley** racing down an initial track segment.
2. There are humans and objects on some track segments.
3. If the track branches, then there is a **switch** the agent can throw. There is an default setting for the switch.
4. With some tracks, the agent has the choice of moving humans/objects between them. This usually applies in only one direction.
5. If the trolley collides with any human/object, it gets killed/destroyed. **This slows down the trolley.**

**Traits** specific to trolley problems



## Other Moral Decision Problems

Notions of moral decision problem / problem instance apply to other problems as well: e.g. surgeon problem, ...

Framework can be used especially well for temporal relations

Extending the translation should be possible

Renaming (for generality):

track segment	possible event
trolley position	current event
switch	choice of the agent

## Other Moral Decision Problems

Notions of moral decision problem / problem instance apply to other problems as well: e.g. surgeon problem, ...

Framework can be used especially well for temporal relations

Extending the translation should be possible

Renaming (for generality):

track segment	possible event
trolley position	current event
switch	choice of the agent

## Other Moral Decision Problems

Notions of moral decision problem / problem instance apply to other problems as well: e.g. surgeon problem, ...

Framework can be used especially well for temporal relations

Extending the translation should be possible

Renaming (for generality):

track segment	possible event
trolley position	current event
switch	choice of the agent

## Problem

Modelling assumption specific to trolley problem: Stop after any collision  
solution ideas:

- ▶ adding `stopping(1)` to the instance description predicates  
thus: making stopping an abnormal
- ▶ identifying certain consequences as stopping
- ▶ extending problem instance by a set of modelling assumptions