# " Dilemma Languages & Moral Codes "



(Photo by Lucas Jackson, tasteless annotations my own)

a half-baked philosophical idea by

Max von Hippel

This is you

This is you

3 people

2 people

Train coming!

This is you

Kant

3 people

Train coming!

IDK some
utilitarian
or
something

2 people

Oh no!
What should I do?

How should I
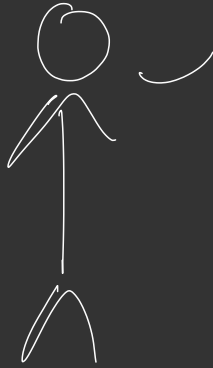even think about
the problem??

CLEARLY, I
need a really
fancy (Turing-complete?)
language to express
the problem....

Oh no!
What should I do?

How should I
even think about
the problem??

CLEARLY, I
need a really
fancy (Turing-complete?)
language to express
the problem....

... and a compiler to
solve (decide) it.

# Table of Contents

A "Dilemma Language" expresses possible ethical dilemmas.

A "Dilemma Language" expresses possible ethical dilemmas.

example: Rail Car Problems ( in BNF )

decide ::= $\triangle$ | $\triangledown$

consequence ::= $-\overset{\text{x}}{\underset{\text{A}}{\text{†}}}-$ | consequence consequence

RCP ::= decide $\overset{RCP}{\underset{RCP}{}}$ | consequence | consequence RCP

A "Dilemma Language" expresses possible ethical dilemmas.

Some examples of Rail Car Problems:

A "moral code" is a compiler from the dilemma language, to the dilema language.

A "moral code" is a compiler from the dilemma language, to the dilema language.

examples: Kant[[ RCP ]] = ∅ // do nothing

A "moral code" is a compiler _from_ the dilemma language, _to_ the dilema language.

examples: $Kant[\![ RCP ]\!] = \emptyset$  // do nothing

$\begin{matrix} \text{IDK random} \\ \text{utilitarian} \\ \text{dude} \end{matrix} [\![ RCP ]\!] =$  a little more complex....

$[\![ RCP ]\!] =$

$C = \text{decide} \begin{matrix} [\cdot] \\ [\cdot] \end{matrix} \mid \text{consequence } C \quad /\!/$ context grammar
with which to
make the decision

$[\![ RCP ]\!] =$

IDK random utilitarian dude

$$C = \text{decide} \begin{array}{c} [\cdot] \\ [\cdot] \end{array} \mid \text{consequence } C \quad \text{//} \quad \text{context grammar with which to make the decision}$$

$$\text{decide} \begin{array}{c} \text{🧍} [\cdot] \\ \text{🧍} [\cdot] \end{array} \equiv \text{decide} \begin{array}{c} [\cdot] \\ [\cdot] \end{array} \quad \text{//} \quad \text{shitty pseudo-BNF equivalence relation on decisions}$$

IDK random utilitarian dude $[\![ RCP ]\!] =$

$$C = decide \frac{[\cdot]}{[\cdot]} \mid consequence\ C \quad // \text{ context grammar with which to make the decision}$$

$$decide \frac{\text{🧍}[\cdot]}{\text{🧍}[\cdot]} \equiv decide \frac{[\cdot]}{[\cdot]} \quad // \text{ shitty pseudo-BNF equivalence relation on decisions}$$

$$\overline{C\frac{[e_0]}{[e_1]} \equiv C\frac{[e_0']}{[e_1']}} \quad \overline{|e_0| < |e_0'|} \quad \overline{|e_1| < |e_1'|} \quad \overline{C\frac{[e_0']}{[e_1']} \to x}$$

$$\frac{}{C\frac{[e_0]}{[e_1]} \longrightarrow x} (decide)$$



// switch down to save a life



// switch up to save a life

$$\left( \frac{up}{down} \right) \qquad \left( \frac{down}{up} \right)$$

// simplify the equation by subtracting from each side until you can solve it in 1-$x$ normal form
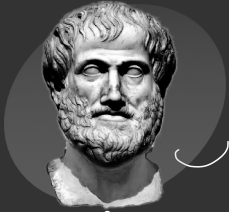
I guess that's cool or whatever

....

but how do I decide which

moral code to use?

For moral code $\mathbb{C}$ over dilemma language $\mathcal{L}$:

$$\text{expected \# deaths} = \text{avg}\left\{\text{num-deaths}\left(\mathbb{C}[\omega]\right) \mid \omega \in \mathcal{L}\right\}$$

$$\text{expected lives saved} = \text{avg}\left\{\text{num-deaths}\left(\mathbb{C}^c[\omega]\right) \mid \omega \in \mathcal{L}\right\}$$
// where $\mathbb{C}^c$ just makes any decision except for $\mathbb{C}$
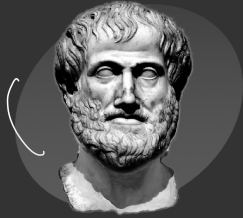
expected people killed
  who would not have $=$ //... not sure how to calculate
been killed had you      this in a logically reasonable
   done nothing       manner...

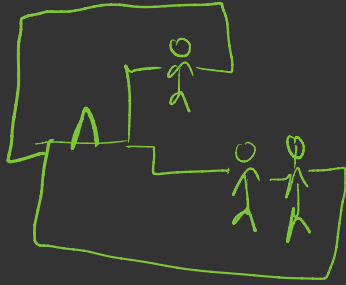that seems reasonable but...

that seems reasonable but...

... what if the dilemma language were

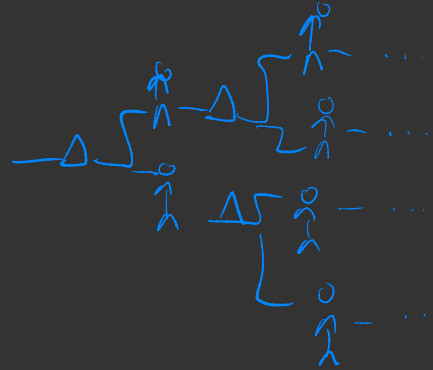# INFINITE?

mua ha ha ha ha

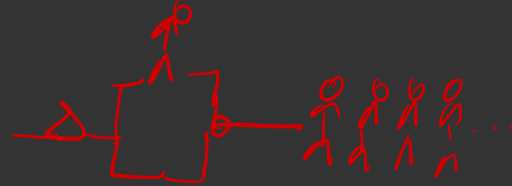# Recursion...



## Infinite trees



## Infinite consequence



## we could extend $L$...

# CONCLUSION

Dilemma languages $\longrightarrow$ cool new way to express moral challenges

Moral Codes $\longrightarrow$ formalism with which to compare (informal) moral codes over given challenges

How can we even compare codes over $\infty$-$\mathcal{L}$'s?

Are there languages having measure $> 0$?

Maybe cool applications to self-driving, eh.?