

“Generalised” models: why, how?

Automating extension of (moral) categories

Stuart Armstrong

Future of Humanity Institute, Oxford University

Strong underdefined preferences



Consider a Croatian, communist, Yugoslav nationalist in the 1980s...

Morality: past, present, and future

Honour is vital



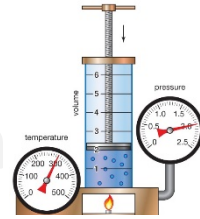
Women should be protected

Happiness is important

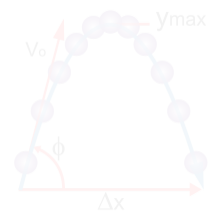
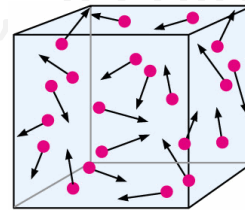


Physical model splintering

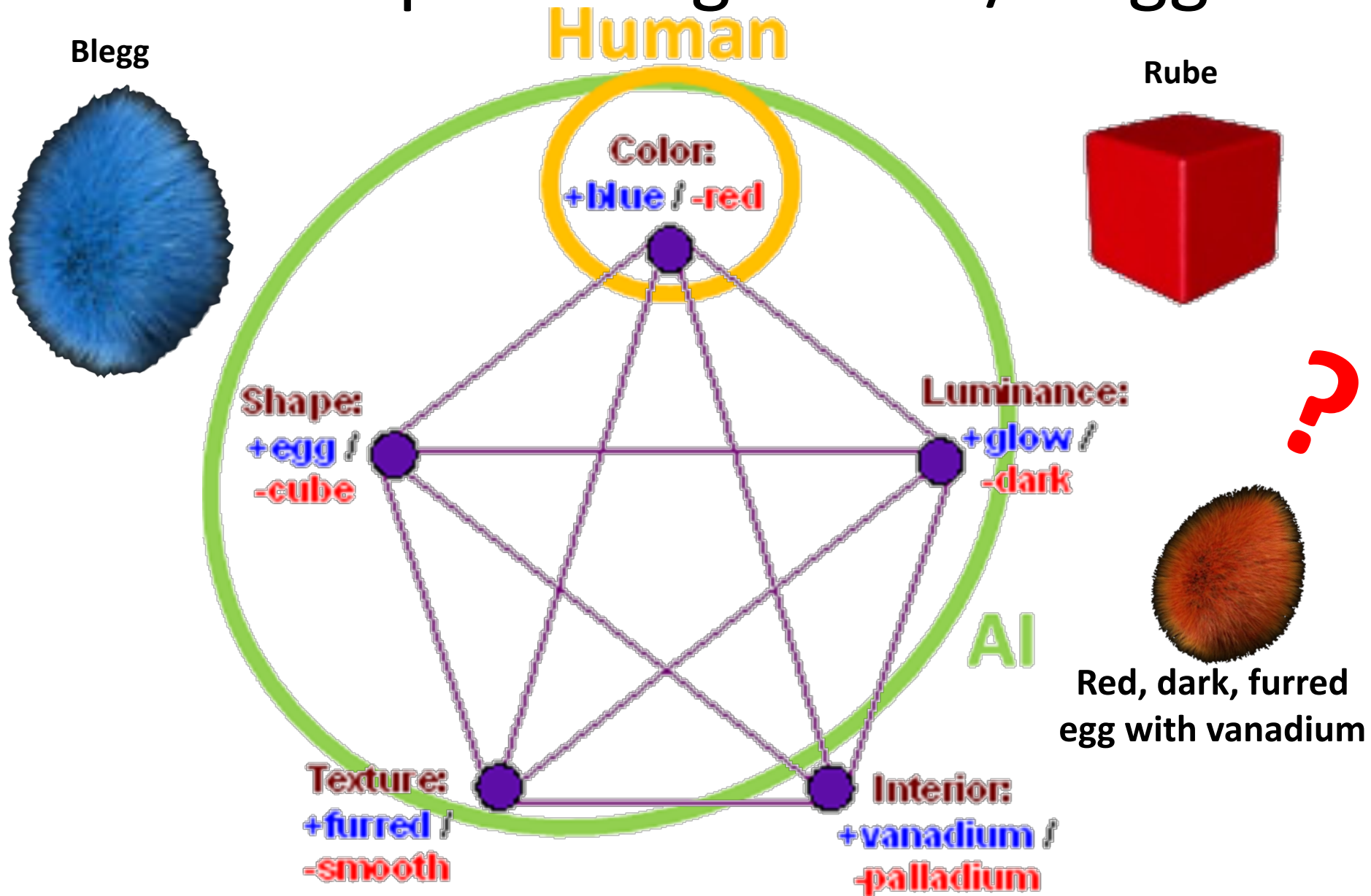
- Aristotelian: elements, natural motions,
- Ideal gas laws
- Van der Waal laws
- Newtonian: force, mass, acceleration...
- Bouncing atom models
- Relativistic: space-time, stress-energy tensors...
- Quantum models
- ...
- Quantum: infinite dimensional Hilbert spaces, self-adjoint operators, eigenfunctions...



$$PV = nRT$$



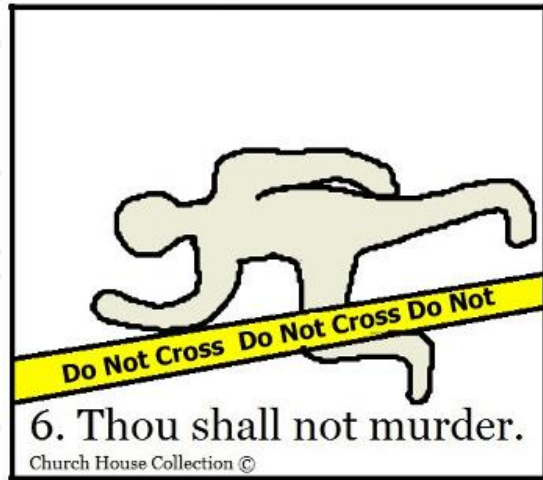
Model splintering: Rubes/Bleggs



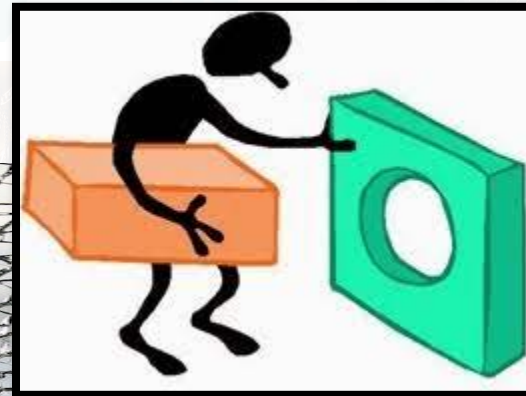
Moral model splintering

- Honour-based morality

Many common conclusions



Incompatible/incomparable concepts and premisses



- Childish moral

The general problem

STARGAZING LIVE THE UNIVERSE THROUGH TIME

BIG BANG

EXPANSION OF THE UNIVERSE BEGINS TO ACCELERATE

EARTH BEGINS

Quantum World

History
3000 BC to nowadays

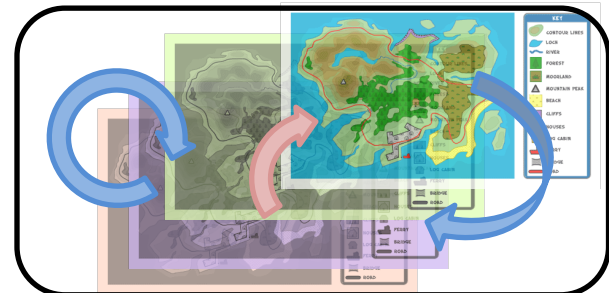
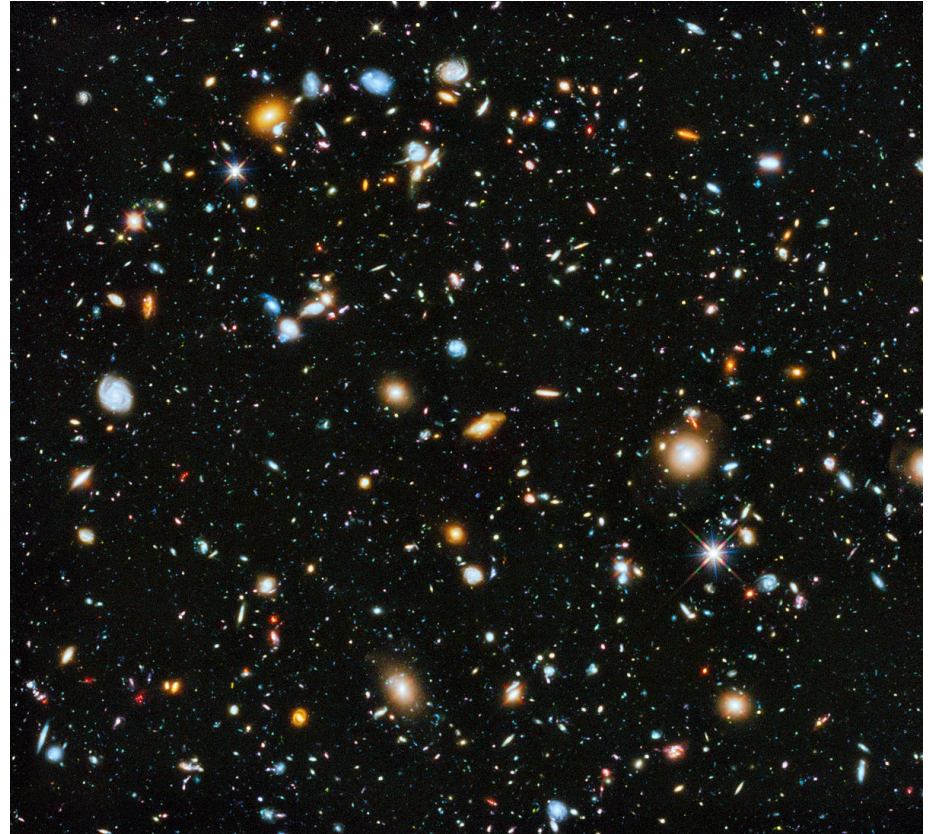
Paleolithic	Neolithic	Ancient Age	Medieval Age	Modern Age	Contemporary Age
1 million years to 10000 BC	10000 BC	3000BC to 476AC	476 AC to 1492 AC	1492 to 1789 AC	1789 AC to 2011 AC
First human species 5 million years BC	Agriculture begins 10000 BC	Invention of writing 3000BC	Fall of Western Roman Empire 476 AC	Colombus discovered America 1492 AC	French Revolution 1789 AC
Use of fire 1 million years BC					

Warwick

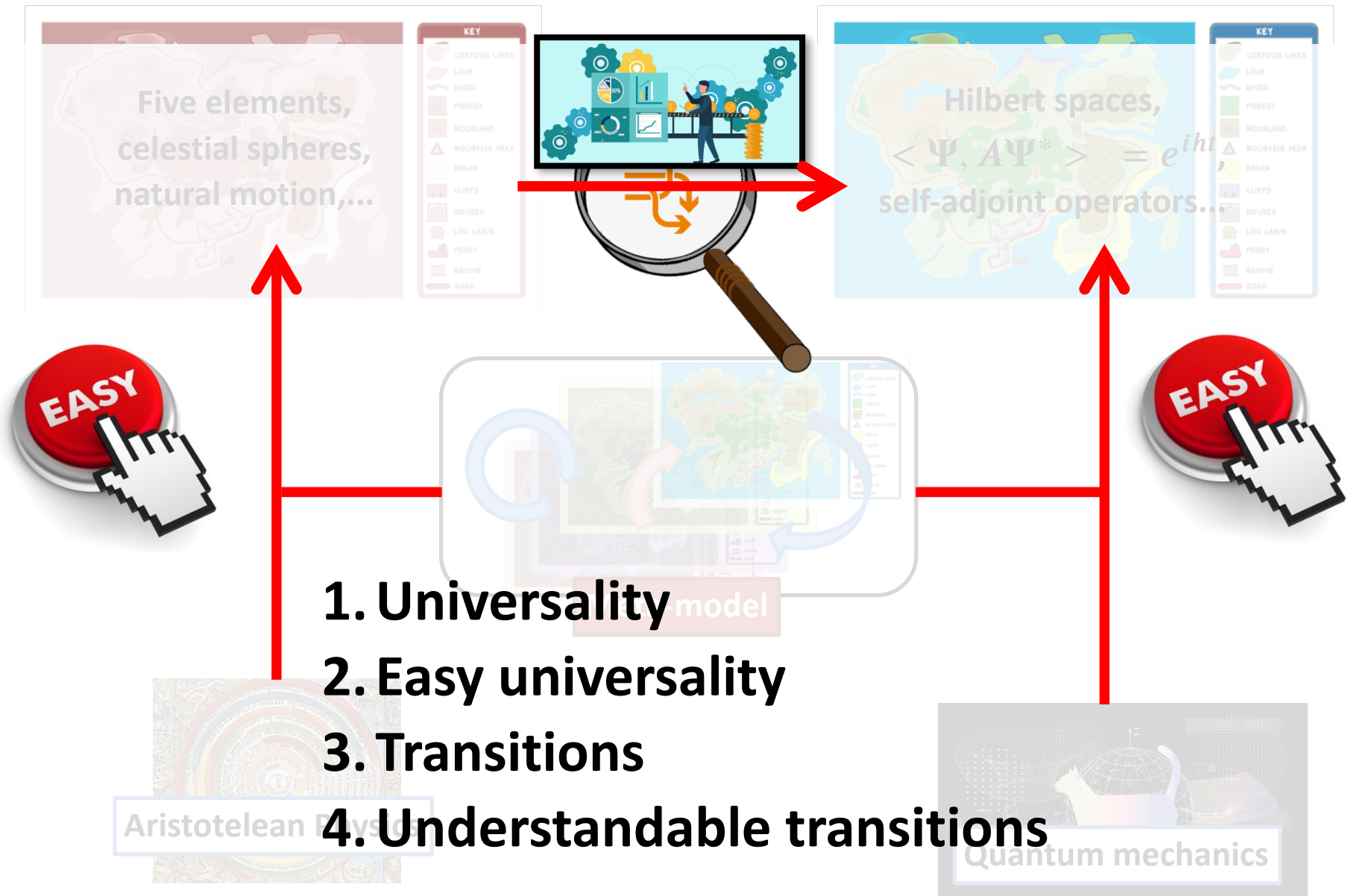
- Money
- The harvest
- Feudal duty
- Teaching children
- Spears and armour
- York vs Lancaster
- House of Warwick
- God
- Morality
- Feudal hierarchy

Universality

- Turing machines
- Neural nets
- Set theory
- Second-order logic
- Bayesian updating
- Category theory
- ...
- “Generalised” models



Meta-model desiderata



Application to most of AI safety

Hidden complexity of wishes	
Ontological crises	
Conservative behaviour	
Goodhart problems	
Wireheading	
Out-of-distribution behaviour	
Low impact	
Underdefined preferences	
Active inverse reward design	
The whole friendly AI problem	

Application to most of AI safety

Hidden complexity of wishes	Save* my mother* [*: underdefined]
Ontological crises	When models of physics splinter
Conservative behaviour	When be conservative? When models splinter
Goodhart problems	“Measure used = desired behaviour” splinters
Wireheading	“Reward channel = desired behaviour” splinters
Out-of-distribution behaviour	The current ML version of this problem
Low impact	Low impact = features similar to before
Underdefined preferences	Example in this presentation
Active inverse reward design	Clear reward over underdefined features
The whole friendly AI problem	“Friendly” well defined in typical situations

Generalised models

$$\mathcal{M} = \{\mathcal{F}, \mathcal{E}, Q\}$$

\mathcal{F} a set of features

$\mathcal{F} = \{(n, \bar{\mathcal{F}})\}: n \text{ name, } \bar{\mathcal{F}} \text{ possible values}$

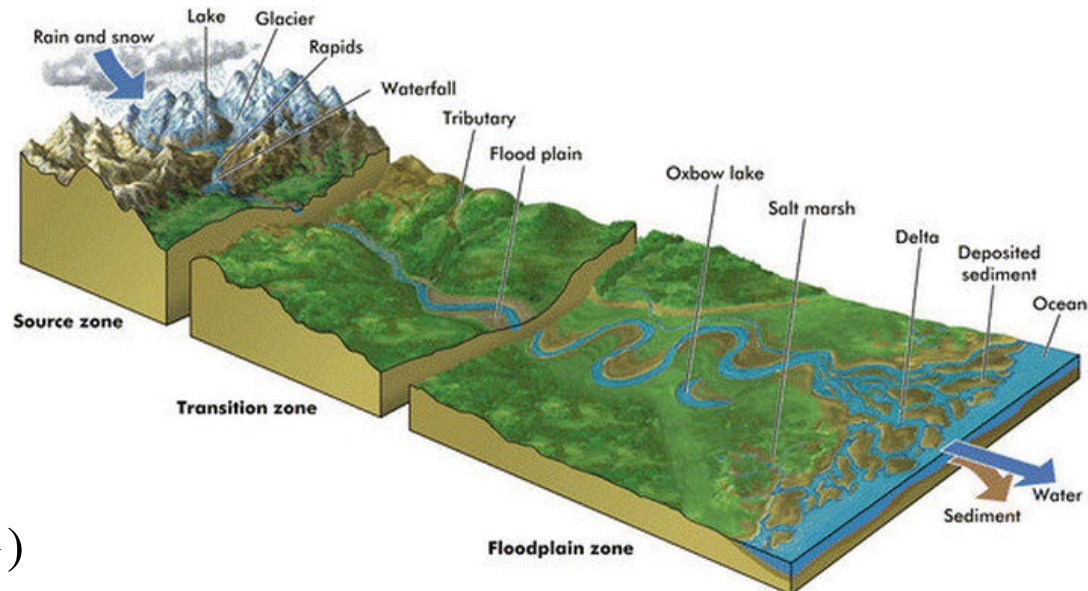
\mathcal{E} a set of environments

$$\mathcal{E} \subset \mathcal{W} = 2^{\sqcup \bar{\mathcal{F}}}$$

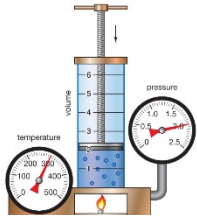
Q a probability distribution
(partial, un-normalised?)



$(n = \text{"temperature"}, \bar{\mathcal{F}} = \{r > 0\})$



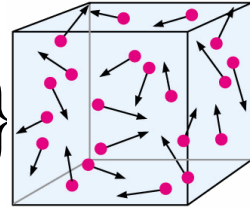
Generalised models



$$PV = nRT$$

$$\mathcal{M}_0 = \{\mathcal{F}_0, \mathcal{E}_0, \mathcal{Q}_0\}$$

$$\mathcal{M}_1 = \{\mathcal{F}_1, \mathcal{E}_1, \mathcal{Q}_1\}$$

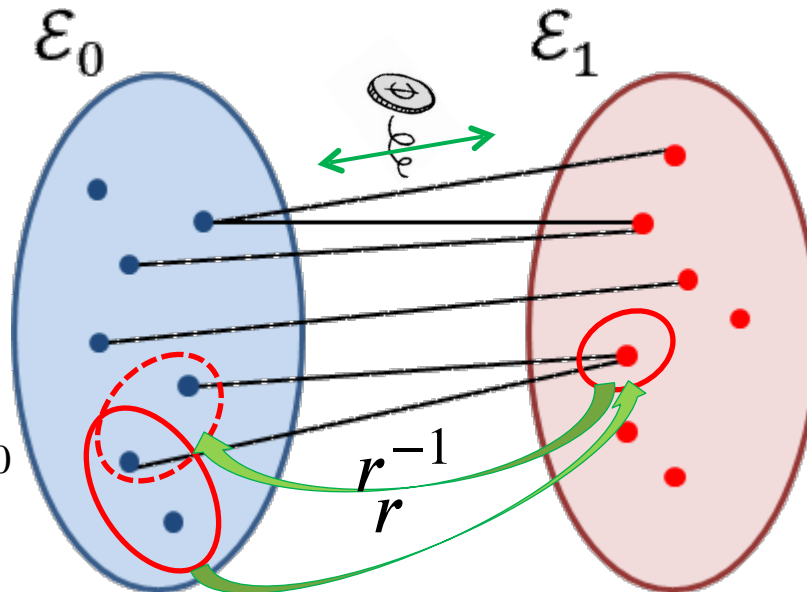


r , a relation between \mathcal{E}_0 and \mathcal{E}_1 :

Induced maps:

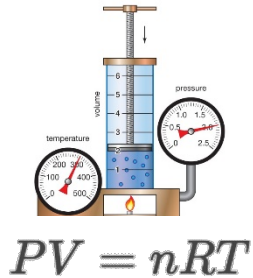
$$r : 2^{\mathcal{E}_0} \rightarrow 2^{\mathcal{E}_1}$$

$$r^{-1} : 2^{\mathcal{E}_1} \rightarrow 2^{\mathcal{E}_0}$$



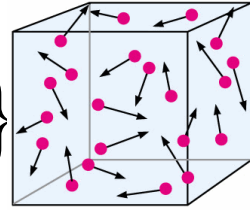
r^{-1} , the inverse relation, between \mathcal{E}_1 and \mathcal{E}_0

Generalised models



$$\mathcal{M}_0 = \{\mathcal{F}_0, \mathcal{E}_0, Q_0\}$$

$$\mathcal{M}_1 = \{\mathcal{F}_1, \mathcal{E}_1, Q_1\}$$



r , a relation between \mathcal{E}_0 and \mathcal{E}_1

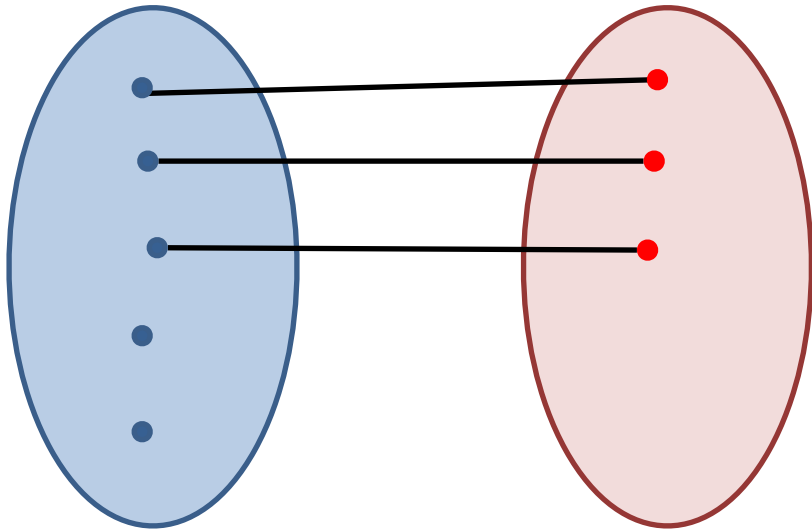
Condition on the Q s:

For all $E_0 \subset \mathcal{E}_0$ and all $E_1 \subset \mathcal{E}_1$:

$Q_0(E_0) \leq Q_1(r(E_0))$ or both probabilities are undefined

$Q_1(E_1) \leq Q_0(r^{-1}(E_1))$ or both probabilities are undefined

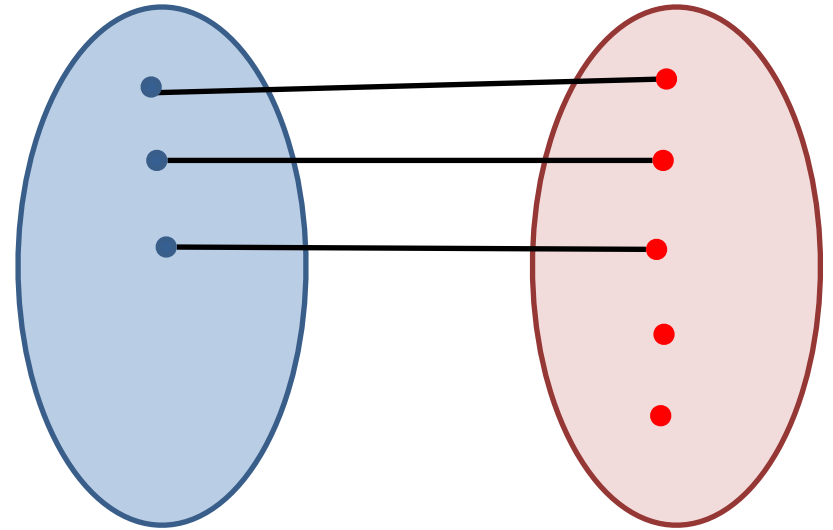
Simple examples



Restriction/Bayesian update:

r bijective partial function

$(r^{-1}$ injective function)

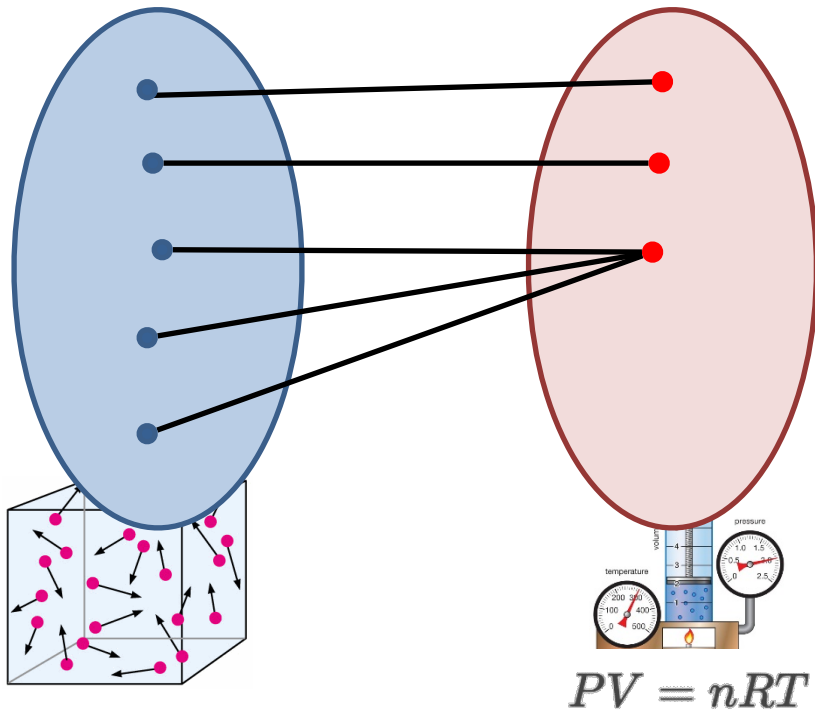


Inclusion:

r injective function

$(r^{-1}$ bijective partial function)

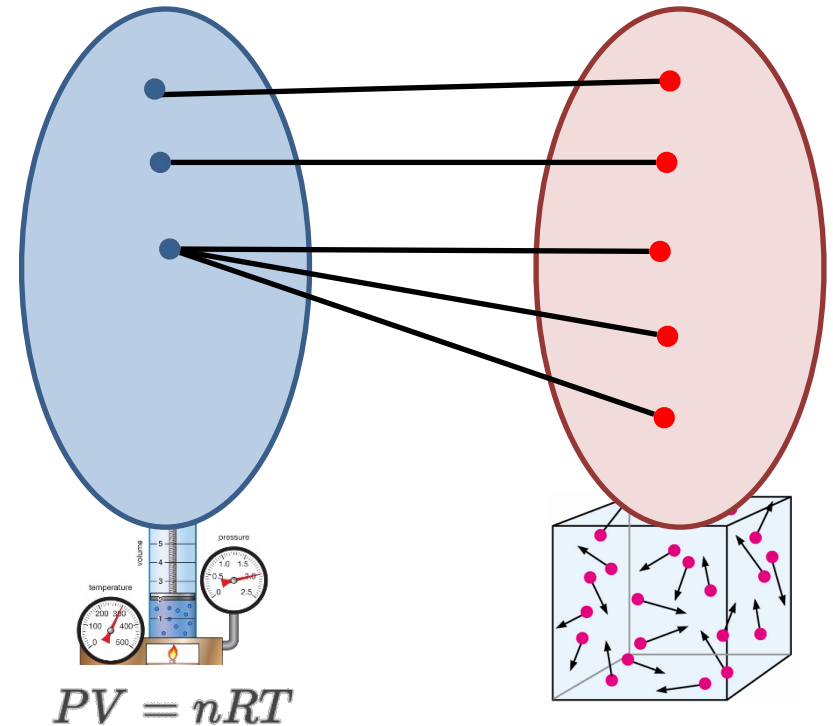
Simple examples



Coarse-graining:

r surjective function
(many-to-one)

(r^{-1} injective, left-total)



Refinement:

r injective, left-total
(one-to-many)

(r^{-1} surjective function)

Improvement

What if the features
and environment don't change?

Go from $\mathcal{M}_0 = \{F_0, \mathcal{E}_0, Q_0\}$
to $\mathcal{M}_1 = \{F_0, \mathcal{E}_0, Q_1\}$

Q_1 is “better” (more accurate, simpler,...) than Q_0



Most model changes: refinements followed by improvements

Cartesian Frames correspondence



$C = \{ A, D, \star \}$ is a Cartesian Frame over W :

\star is a map from $A \times D$ to W

$$a \star d = w$$

A morphism from $C_0 = \{ A_0, D_0, \star_0 \}$ to

$$C_1 = \{ A_1, D_1, \star_1 \}$$

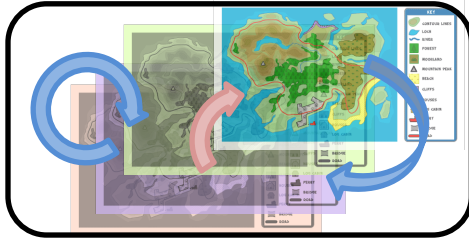
is a pair of functions: $(g_0 : A_0 \rightarrow A_1, h_1 : D_1 \rightarrow D_0)$,

such that for all a_0, d_1 ,

Chu(W)

$$g_0(a_0) \star_1 d_1 = a_0 \star_0 h_1(d_1)$$

Cartesian Frames correspondence



Define $GM(W)$ as a subcategory of the generalised models, with:

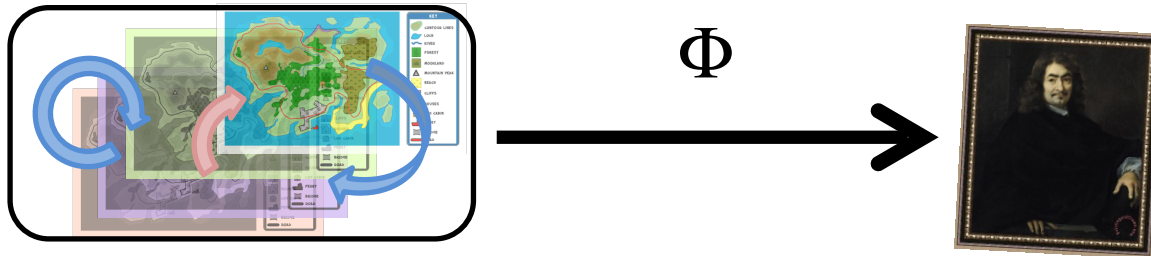
1. Features: $\mathcal{F} = \{A, D, W\}$
2. Environment: $\mathcal{E} = A \times D \times W$
(using $S \subset 2^S$, $2^{A \sqcup D \sqcup W} = 2^A \times 2^D \times 2^W$)
3. For all a and d , $Q(a, d, w) = 0$, apart from one single w , specific to a and d .
4. Morphisms: r is a relation between $A_0 \times D_0 \times W$ and $A_1 \times D_1 \times W$, derived from the functions/relations (g_0, h_1, Id_W)

Cartesian Frames correspondence






Then define $\Phi : GM(W) \rightarrow Chu(W)$ sending:

1. $(\mathcal{F}, A \times D \times W, Q)$ to (A, D, \star) ,
with $a \star d = w$ iff $Q(a, d, w) \neq 0$
2. (g_0, h_1, Id_W) to (g_0, h_1)

Then Φ is a surjective functor of categories.



How good a meta-model?

1. Features not well-integrated into category-theory formalism. 
2. Improvements (to Q) **not** integrated. 
3. Change of environment \mathcal{E} well integrated. 
4. Universal for some definitions.
5. Easy universality. 
6. Model transitions not so easy to understand (see points 1 and 2). 

$$P(f_1 = x \mid f_0 = y)$$

Relevant links

- Generalised models as a category:
- <https://www.lesswrong.com/posts/nQxqSsHfexivsd6vB/generalised-models-as-a-category>
- Cartesian frames as generalised models:
- <https://www.lesswrong.com/posts/wiQeYuQPwSypXXFar/cartesian-frames-as-generalised-models>
- Model splintering:
- <https://www.lesswrong.com/posts/k54rgSg7GcjtXnMHX/model-splintering-moving-from-one-imperfect-model-to-another-1>