

WHAT IS LEWISEAN INTERPRETIVISM?

Edward Elliott

School of Philosophy, Religion and History of Science
University of Leeds

e.j.r.elliott@leeds.ac.uk

18 June 2019

TWO PRINCIPLES OF INTERPRETATION

'Radical Interpretation': Lewis' account of *interpretation*, based on two key principles:

PRINCIPLE OF RATIONALISATION

Assign beliefs and basic desires to an agent so as to maximise their *pragmatic rationality*, given the facts about their behaviour

PRINCIPLE OF CHARITY

Assign an agent (i) beliefs so as to maximise their *epistemic rationality* given the facts about the agent's life history of evidence; and (ii) basic desires appropriate for the kind of being the agent is

EXEGETICAL QUESTIONS

QUESTION 1: *How do these two principles interact? What are we to do when they pull in different directions?*

EXAMPLE: MONSTERS IN THE CLOSET

Karl has no significant evidence that there's a monster in his closet, and overwhelming evidence that there isn't; however, he acts in a way that would clearly be well explained by high confidence in the existence of closet-dwelling monsters (and a desire to not be eaten).

EXEGETICAL QUESTIONS

The interaction question crops up in, *inter alia*,

- (Maher 1993, pp. 31-2)
- (Hájek & Eriksson 2007, pp. 200-1)
- (Hattiangadi 2019)

Obvious aggregation strategies...

- *Tradeoffs* — which 'weightings' to use, and why?
- *Lexical priority* — which gets priority, and why?

EXEGETICAL QUESTIONS

QUOTE: 'RADICAL INTERPRETATION', P. 341

First step: using [the physical facts] both as a source of information on Karl's behaviour and as a source of information on his life history of evidence, fill in [his beliefs and desires] completely by means of the Rationalisation Principle and the Principle of Charity.

QUESTION 2: *Why didn't Lewis say more about this?*

OVERVIEW

In the remainder, I will argue...

- QUESTION 1:
Lewis held a *lexical priority* view, but between *causal* and *non-causal* principles of interpretation
- QUESTION 2:
Lewis didn't discuss the interaction because he (rightly) considered the issue to be of minor importance

POSTSCRIPTS TO 'RADICAL INTERPRETATION'

A QUICK DETOUR...

The discussion of interpretation in 'Radical Interpretation' is misleading in two important (and related) respects:

- It is *individualistic*
- It describes the interpretation of *persons*, not *states*

INTERPRETATION OF STATES AND INDIVIDUALS

QUOTE: 'POSTSCRIPTS TO RADICAL INTERPRETATION', P. 119

I stated my problem in an unduly individualistic way: give the facts *about Karl* as a physical system, solve for the facts *about him* as a person—*his* beliefs, desires, and meanings...

In 'Mad Pain and Martian Pain', I argued that a 'madman' might be in pain not because his state occupied the causal role of pain in him but rather because that state occupies that role, for the most part, in members of the kind to which he belongs. The same possibility should be recognised for attitudes as well.

INTERPRETATION OF STATES AND INDIVIDUALS

The important points:

- 1 *Interpretation* involves first assigning contents to *states*
 - These are “most likely states of the brain” (1983, p. 373)
 - In later works, whole systems of (degrees of) belief and basic desires $\langle Bel, Des \rangle$ are assigned to *total physical states*
- 2 We're to interpret *states* in such a way as to maximise CHARITY and RATIONALISATION for the population *overall*
 - Later treated as equivalent to: the interpretation of a state depends on its *typical causal role*
- 3 Interpretation of *individuals* depends on what states they're in

INTRODUCTION

'POSTSCRIPTS TO RADICAL INTERPRETATION'

REINTERPRETING 'RADICAL INTERPRETATION'

'NEW WORK FOR A THEORY OF UNIVERSALS'

RATIONALISATION IN 'RADICAL INTERPRETATION'

CHARITY IN 'RADICAL INTERPRETATION'

COMBINING CHARITY AND RATIONALISATION

REINTERPRETING 'RADICAL INTERPRETATION'

RATIONALISATION IN 'RADICAL INTERPRETATION'

QUOTE: 'RADICAL INTERPRETATION', PP. 337–8

Karl should be represented as a rational agent; the beliefs and desires ascribed to him should be such as to provide good reasons for his behaviour, as given in physical terms...

I would hope to spell this out in decision-theoretic terms, as follows. Take a suitable set of mutually exclusive and jointly exhaustive propositions about Karl's behaviour at any given time; of these alternatives, the one that comes true according to \mathcal{P} should be the one (or: one of the ones) with maximum expected utility according to the total system of beliefs and desires [$\langle \mathcal{B}el, \mathcal{D}es \rangle$] ascribed to Karl at that time.

RATIONALISATION IN 'RADICAL INTERPRETATION'

CAUSAL RATIONALISATION

An interpretation $\langle Bel, Des \rangle$ **\mathcal{R} -fits** a total physical state S iff...
 the typical agent in S will behave in way B that maximises expected utility w.r.t $\langle Bel, Des \rangle$

' \mathcal{R} -fit' corresponds to a forwards-looking typical causal role:

$\langle Bel, Des \rangle \Rightarrow B$, where B maximises EU w.r.t. $\langle Bel, Des \rangle$

CHARITY IN 'RADICAL INTERPRETATION'

PRINCIPLE OF CHARITY

Assign an agent (i) beliefs so as to maximise their *epistemic rationality* given the facts about the agent's life history of evidence; and (ii) basic desires appropriate for the kind of being the agent is

In a Bayesian framework, part (i) requires:

- Bel is probabilistically coherent
- Bel updates by conditionalisation

It'll be helpful to break 'Charity' down into causal/non-causal parts...

CHARITY IN 'RADICAL INTERPRETATION'

Let $Bel_E = Bel$ conditionalised on E

CAUSAL CHARITY

Interpretations $\langle Bel, Des \rangle$ and $\langle Bel_E, Des \rangle$ **C-fit** states S and S^* iff...
the typical agent in state S when presented with a stream of evidence E
will come to be in state S^*

- 'C-fit' corresponds to a forwards-/backwards-looking typical role:

$$\langle Bel, Des \rangle + E \Rightarrow \langle Bel_E, Des \rangle$$

CHARITY IN 'RADICAL INTERPRETATION'

NON-CAUSAL CHARITY

An interpretation $\langle Bel, Des \rangle$ is *better* to the extent that Bel and Des are *reasonable* relative to the agent's kind

- Non-causal — depends on...
 - ① intrinsic characteristics of the interpretation
 - ② facts about the agent's kind

COMBINING CHARITY AND RATIONALISATION

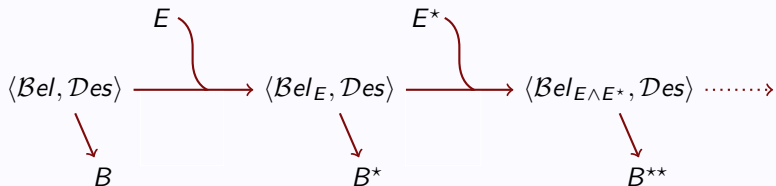
- CAUSAL RATIONALISATION

$\langle Bel, Des \rangle \Rightarrow B$, where B maximises EU w.r.t. $\langle Bel, Des \rangle$

- CAUSAL CHARITY

$\langle Bel, Des \rangle + E \Rightarrow \langle Bel_E, Des \rangle$

For any stream of evidence E, E^*, \dots , and any starting point $\langle Bel, Des \rangle$:



INTRODUCTION
'POSTSCRIPTS TO RADICAL INTERPRETATION'
REINTERPRETING 'RADICAL INTERPRETATION'
'NEW WORK FOR A THEORY OF UNIVERSALS'

COMBINING CHARITY AND RATIONALISATION
PRINCIPLES OF FIT
THE STRUCTURE OF LEWISEAN INTERPRETIVISM
ANSWERING THE EXEGETICAL QUESTIONS

'NEW WORK FOR A THEORY OF UNIVERSALS'

THE ROLE OF 'CHARITY'

QUOTE: 'NEW WORK FOR A THEORY OF UNIVERSALS', P. 375

If we rely on principles of fit to do the whole job, we can expect radical indeterminacy of interpretation. We need further constraints, of the sort called principles of (sophisticated) charity, or of 'humanity'. Such principles call for interpretations according to which the subject has attitudes that we would deem reasonable for one who has lived the life that he has lived...

These principles select among conflicting interpretations that equally well conform to the principles of fit.

PRINCIPLES OF FIT

QUOTE: 'NEW WORK FOR A THEORY OF UNIVERSALS', P. 374

Say that Bel and Des rationalise behaviour B after evidence E iff the system of desires given by the Bel_E -expectations of Des ranks B at least as high as any alternative behaviour.

Say that Bel and Des fit iff, for any evidence-specifying E , [learning] E yields a state that would cause behaviour rationalised by Bel and Des after E .

That is our only constraining principle of fit.

PRINCIPLES OF FIT

PRINCIPLE OF FIT

An interpretation $\langle Bel, Des \rangle$ **\mathcal{F} -fits** a total physical state S iff...
 for any evidence E , the typical agent in S when presented with a stream of evidence E will come to be in a state S^* that causes some behaviour B which maximises expected utility with respect to $\langle Bel_E, Des \rangle$

- ' \mathcal{F} -fit' is supposed to help capture both:
 - the role of evidence in the adjustment of beliefs
 - the role of beliefs/desires in the production of behaviour
- However, it is entirely *forwards-looking*:

$$\langle Bel, Des \rangle + E \Rightarrow S^* \Rightarrow B, \quad \text{where } B \text{ maximises EU w.r.t. } \langle Bel_E, Des \rangle$$

PRINCIPLES OF FIT

QUOTE: 'NEW WORK FOR A THEORY OF UNIVERSALS', P. 374

That is our only constraining principle of fit. (Where did the others go?
— *We built them into the definitions whereby \mathcal{B} el and \mathcal{D} es encapsulate
an assignment of content to various states.*)

PRINCIPLES OF FIT

QUOTE: 'NEW WORK FOR A THEORY OF UNIVERSALS', P. 374

Bel is a probability distribution... regarded as encapsulating the subject's dispositions to form beliefs under... evidence: if a stream of evidence specified by proposition *E* would put the subject into a total state *S*—for short, *E* yields *S*—we interpret *S* to consist in part of the belief system... that comes from *Bel* by conditionalisation on *E*.

Des is... regarded as encapsulating the subject's basic desires: if *E* yields *S*, we interpret *S* to consist in part of the system of desires given by the *Bel_E*-expectations of *Des*.

PRINCIPLES OF FIT

FORWARDS COHERENCE

S can be interpreted $\langle Bel, Des \rangle$ only if...

for any stream of evidence E , if $S + E \Rightarrow S^*$, then S^* is interpreted $\langle Bel_E, Des \rangle$

BACKWARDS COHERENCE

S can be interpreted $\langle Bel, Des \rangle$ only if...

for any stream of evidence E , if there exists an S^* such that $S^* + E \Rightarrow S$, then S^* is interpreted in some way $\langle Bel^*, Des \rangle$ such that $Bel_E^* = Bel$

PRINCIPLES OF FIT

FIT + BACKWARDS COHERENCE + FORWARDS COHERENCE
 =
 GENERAL PRINCIPLE OF FIT

GENERAL PRINCIPLE OF FIT

An interpretation $\langle \mathcal{B}el, \mathcal{D}es \rangle$ **fits** a total physical state S iff...
 for any stream of evidence E ,

- $\langle \mathcal{B}el, \mathcal{D}es \rangle$ \mathcal{F} -fits S
- if $S + E \Rightarrow S^*$, then $\langle \mathcal{B}el_E, \mathcal{D}es \rangle$ \mathcal{F} -fits S^*
- if there exists an S^* such that $S^* + E \Rightarrow S$, then some $\langle \mathcal{B}el^*, \mathcal{D}es \rangle$ such that $\mathcal{B}el_E^* = \mathcal{B}el$ \mathcal{F} -fits S^*

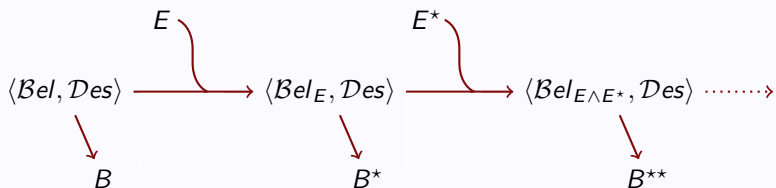
THE STRUCTURE OF LEWISEAN INTERPRETIVISM

GENERAL PRINCIPLE OF FIT

=

CAUSAL RATIONALISATION + CAUSAL CHARITY

For any stream of evidence E, E^*, \dots , and any starting point $\langle Bel, Des \rangle$:



ANSWERING THE EXEGETICAL QUESTIONS

QUESTION 1:

Fit with functional role takes lexical priority over considerations of 'reasonableness', but both RATIONALISATION and (a part of) CHARITY are involved in characterising that role

QUESTION 2:

The interpretation of states depends on what they do in the *typical* case — and plausibly, as a matter of fact in the typical case

RATIONALISATION and CHARITY don't pull in different directions

Agents who don't behave in a manner that seems pragmatically rational given what would be epistemically rational in light of their evidence are *deviant cases* — the normal rules don't apply

Thank you

Thanks to Jessica Isserow for comments on a practice run!

QUOTE: *On the Plurality of Worlds*, p. 107

Such a theory, I said, should have two parts. *One part says what it is for an assignment of content to states to fit the functional roles of the states*; the constraints are principles of rationality, for instance a principle to the effect that a state which is assigned content consisting of some system of beliefs and desires ought to be a state that tends to produce conduct that would serve those desires according to those beliefs. But principles of fit can be expected to underdetermine the assignment of content very badly. . . . *Therefore a theory of content needs a second part: as well as principles of fit, we need 'principles of humanity', which create a presumption in favour of some sorts of content and against others.*