CAMBRIDGE
UNIVERSITY PRESS

ARTICLE

# Rethinking Hare's Analysis of Moral Thinking

Steven Daskal* (iD)

Northern Illinois University
*Corresponding author. E-mail: sdaskal@niu.edu

**Abstract**
R. M. Hare has an ambitious project of arguing from a limited set of premises about the nature of moral thought and language all the way to substantive utilitarian conclusions. I reconstruct Hare's argument, identify an important problem for Hare, and then develop and endorse a restricted Hare-like argument. This argument is less ambitious than Hare's, and does not substantiate utilitarian conclusions on its own, but I demonstrate that it nonetheless imposes important constraints on moral judgements and I indicate how it can play a role in a larger argument for utilitarian conclusions.

R. M. Hare undertakes the remarkably ambitious project of arguing from a limited set of premises about the nature of moral thought and language all the way to substantive preference utilitarian conclusions. In this article, I hope to show that even though Hare is not fully successful, his line of analysis has greater force than is generally recognized. Drawing on Hare's argumentative strategy, I seek to identify genuine constraints on moral judgements, albeit not ones that are strong enough to mandate utilitarian conclusions on their own.

I begin, in section I, with a reconstruction of Hare's analysis, as developed in *Moral Thinking*.[1] In section II, I consider several initial objections to Hare, one of which demands a relatively minor revision of his argument, which I undertake in section III. Then, in section IV, I consider what I take to be a more problematic set of objections, stemming in part from the revision from sections II and III, and involving one of Hare's key premises. This leads me, in section V, to construct and endorse what I call a 'restricted Hare-like argument'. This argument is weaker than Hare's own argument, but I demonstrate that it nonetheless places meaningful constraints on moral judgements. Moreover, although it cannot get all the way to utilitarian conclusions on its own, I show how it can play an important role in such a project by closing off a gap in similarly ambitious arguments developed by John Harsanyi and Allan Gibbard.[2] If successful, my analysis contributes to our understanding of the nature of moral thought

---

[1]R. M. Hare, *Moral Thinking: Its Levels, Method and Point* (Oxford, 1981). The view Hare develops in this book is grounded in his previous work, including both R. M. Hare, *Language of Morals* (Oxford, 1952) and R. M. Hare, *Freedom and Reason* (Oxford, 1963).

[2]John Harsanyi, 'Morality and Theory of Rational Behavior', *Social Research* 44 (1977), pp. 623–56, and Allan Gibbard, *Reconciling our Aims: In Search of Bases for Ethics* (Oxford, 2008).

and the demands of morality, and also demonstrates that Hare has masterminded a line of analysis that can get surprisingly far from surprisingly minimal starting points, even if it cannot accomplish all of what Hare intended.

## I. Hare's system

As I understand Hare's analysis, there are three crucial premises.[3] Two of them are claims about the nature of moral judgements: that such judgements are universalizable and prescriptive. Hare explicitly asserts these premises and bases his analysis upon them.[4] The third is a claim that I will call a principle of conditional reflection, following terminology introduced by Gibbard.[5] Hare does not explicitly identify this as a premise, but he recognizes that it plays a critical role in his analysis.[6] After presenting Hare's versions of his premises, I will reconstruct his argument from them to substantive utilitarian conclusions.

### I.1. Universalizability

Hare identifies the universalizability of moral judgements as a starting point for his analysis. His conception of universalizability amounts to the claim that moral judgements are not sensitive to the agent's actual role in a situation.[7] In other words, in order for an assessment of a situation to count as a moral judgement it must remain constant regardless of whether one is the individual acting or any of the others involved in the situation. And this, he thinks, is built into the idea of a moral judgement.

To see how this works, it may help to consider a simple toy case, $S_1$, in which I am choosing between just two possible actions, X and Y. And let us suppose that you are the only other person affected by my choice. Insofar as we are interested in forming a moral judgement regarding $S_1$, Hare directs our attention to $S_2$, which is just like $S_1$ except that you are in the role of the agent and I am in the role of the individual affected by the agent's action. That is to say, in $S_2$ I have all of the personal characteristics and history that you have in $S_1$, and vice versa. Hare's version of universalizability is the claim that a judgement that I am morally permitted to do X in $S_1$ commits one to judging that you are likewise permitted to do X in $S_2$.

As this demonstrates, Hare's universalizability requirement does not restrict judgements ranging over pairs or sets of situations that an actual human agent will face over the course of time, or even situations that different agents will face, or expect to face. All it imposes are restrictions on judgements governing a very narrow set of purely hypothetical cases in which individuals occupy different roles in otherwise identical situations, with those roles construed as including personal characteristics and history.

---

[3]My reconstruction of Hare's analysis is informed by Allan Gibbard, 'Hare's Analysis of "Ought" and its Implications', *Hare and Critics: Essays on Moral Thinking*, ed. Douglas Seanor and Nicholas Fotion (Oxford, 1988), pp. 57–72, which Hare approves of, writing: 'Gibbard, unlike many writers, gets me right', in R. M. Hare, 'Comments', *Hare and Critics: Essays on Moral Thinking*, ed. Douglas Seanor and Nicholas Fotion (Oxford, 1988), pp. 199–293, at 230. See also R. M. Hare, 'A Philosophical Autobiography', *Utilitas* 14 (2002), pp. 269–305, at 300–1.

[4]Hare also identifies a third feature of moral judgements, that they are overriding, but this does not play a significant role in his analysis, as recognized in Hare, *Moral Thinking*, pp. 21 and 54 and Hare, 'A Philosophical Autobiography', p. 290.

[5]Gibbard, 'Hare's Analysis', p. 58.

[6]Hare, 'Comments', pp. 229–30.

[7]Hare, *Moral Thinking*, p. 21.

Nonetheless, Hare believes that even this very minimal requirement, which he thinks will be acknowledged by any competent user of moral concepts and terms, can lead to surprisingly robust, substantive conclusions about the demands of morality.

## I.2. Prescriptivity

Hare's second premise is the claim that moral judgements are prescriptive, which is to say that moral judgements tell one what to do or to avoid doing.[8] Moreover, Hare takes these prescriptions to be linguistic representations of preferences or motivational states.[9]

To understand Hare's conception of prescriptivity, it helps to draw a distinction between a preference tendency, which is an inclination or leaning, and a preference all told, which is an overall preference that takes all of an agent's preference tendencies into account.[10] It is one's preference all told for a situation that Hare believes is expressed in one's moral judgement regarding that situation.

As with the case of universalizability, Hare thinks this is built into the idea of a moral judgement. As he sees it, this way of understanding moral judgements does not presuppose any substantive claims about what is and is not morally permissible, but merely involves linguistic claims about the nature of moral language, or conceptual claims about the idea of morality. If someone makes a purportedly moral judgement that is not universalizable and prescriptive, that person is not engaging in a genuinely moral assessment of the situation. And Hare thinks any competent user of moral concepts and language will recognize this.

## I.3. Conditional reflection of hypothetical preferences

In addition to his premises about the nature of moral judgements, Hare relies on a claim about a particular way in which practical reasoning generates new preferences. To understand this form of preference generation, we must draw a distinction between hypothetical preferences, or preferences an agent would have in certain circumstances, and conditional preferences, or preferences an agent currently has regarding certain circumstances. In principle, an agent's conditional and hypothetical preferences regarding a situation can be different: what I, now, prefer for a particular situation may not be the same as what I would prefer were I in that situation. But Hare asserts a strong connection between the two through his conditional reflection thesis, according to which careful reflection on one's hypothetical preferences, or the preferences one would have in a given hypothetical situation, leads one to develop corresponding conditional preferences for that hypothetical situation. As Hare puts it:

It is important to emphasize the distinction between the two propositions:

(1) I now prefer with strength S that if I were in that situation $x$ should happen rather than not;

(2) If I were in that situation, I would prefer with strength S that $x$ should happen rather than not.

---

[8] Hare, *Moral Thinking*, p. 21.
[9] Hare, *Moral Thinking*, p. 107.
[10] This terminology is from Gibbard, 'Hare's Analysis'.

> … What I am claiming is not that these propositions are identical, but that I cannot know that (2), and what that would be like, without (1) being true, and that this is a conceptual truth, in the sense of 'know' that moral thinking demands.[11]

I will, in section IV, argue that Hare's conditional reflection thesis, once modified in light of some concerns involving Hare's prescriptivity premise, is too strong, and I will propose a restricted version of the thesis to be used in its stead. I will then need to consider how the restricted principle of conditional reflection impacts Hare's line of analysis. First, though, let me lay out what I take to be Hare's own version of the argument, which includes his conditional reflection thesis as a premise.

## I.4. Hare's argument

I will formulate my reconstruction of Hare's argument in terms of the simple toy case described above, with just two individuals: an agent facing a choice between doing X or Y, and another individual affected by the agent's action. $S_1$ will be the situation in which I am in the role of agent and you are affected by my action, and $S_2$ will be the same situation except that our positions are swapped. And those positions will be understood to include personal characteristics broadly construed, so that if we refer to preference tendencies that an agent has for a situation before engaging in moral reasoning as initial preference tendencies, in $S_2$ I will have your initial preference tendencies from $S_1$, and you will have mine.

Suppose $S_1$ is the actual case, and I am inclined to do X, but only if it is morally permissible. I therefore wonder whether doing X is morally permissible. Hare's argument is intended to constrain my conclusions, or help me determine whether X is morally permissible in $S_1$.

The argument begins with Hare's universalizability premise, according to which a judgement that X is morally permissible in $S_1$ commits me to the judgement that it is also morally permissible in $S_2$. Given Hare's understanding of moral judgements as expressions of preferences all told, this means that in order to judge that X is morally permissible in $S_1$ I must have a preference all told that permits X in $S_1$, and I must also be committed to a matching preference all told that permits X in $S_2$.

With that requirement in place, Hare then applies his conditional reflection thesis to $S_2$, which leads him to conclude that if I reflect carefully on $S_2$ I will develop conditional preferences for $S_2$ that match the hypothetical preference tendencies I have in $S_2$. Given that $S_2$ just is the hypothetical case in which I occupy your role in $S_1$, and therefore have hypothetical preference tendencies equivalent to your initial preference tendencies in $S_1$, this means that I will develop conditional preference tendencies for $S_2$ that match your initial preference tendencies in $S_1$.

So at this point I have a set of preference tendencies for $S_1$ that are constituted by my initial preference tendencies in $S_1$, and a set of preference tendencies for $S_2$ that are equivalent to your initial preference tendencies in $S_1$. Moreover, as indicated above, Hare's universalizability and prescriptivity premises combine to show that I can form a moral judgement regarding $S_1$ only if that judgement expresses a preference all told for $S_1$ and only if I have, or am committed to having, a matching preference all told for $S_2$.

---

[11]Hare, *Moral Thinking*, pp. 95–6.

Given that a preference all told is determined by the preference tendencies which compose it, the only way to ensure that I end up with matching preferences all told for $S_1$ and $S_2$ is to have them composed of matching sets of preference tendencies. If I have a set of preference tendencies for $S_1$ and another set of preference tendencies for $S_2$, and no antecedent guarantee of overlap between these sets, the only systematic way to end up with matching sets of preference tendencies for both cases is to add the sets together. In other words, in forming a preference all told for $S_1$ I must incorporate copies of my preference tendencies for $S_2$ (which are equivalent to your initial preference tendencies in $S_1$) together with my initial preference tendencies for $S_1$, and in forming a preference all told for $S_2$ I must incorporate copies of my initial preference tendencies for $S_1$ together with my preference tendencies for $S_2$.

If this line of argument is successful, Hare is now able to conclude that if I have reflected carefully on $S_2$ (the case in which I am in your shoes), and my preferences are consistent with my commitments, I can judge that X is morally permissible in $S_1$ (the case in which I am acting) only if X is permitted by a preference all told incorporating both my preference tendencies for $S_1$ and copies of your preference tendencies for $S_1$. And notice, this is equivalent to a preference utilitarian calculation, or an assessment of moral permissibility grounded in the preferences of all relevant individuals, which in this two-party case is just you and me. Hare's argument therefore leads to an endorsement of preference utilitarian moral conclusions.

Moreover, although this argument was constructed in terms of a case involving only two individuals, it readily generalizes to multi-party cases. In such cases the deliberating agent will need to consider a set of hypothetical scenarios in which she occupies the role of each individual in turn, and a moral judgement regarding the actual situation will commit her to matching moral judgements regarding each of these hypothetical scenarios. Her judgement for any particular hypothetical scenario will have to incorporate preference tendencies corresponding to the initial preference tendencies of the individual whose role she occupies in that hypothetical scenario, which means that her preferences all told for the various scenarios will only match if they all incorporate the initial preference tendencies of everyone.

## II. Initial objections

### II.1. Preference utilitarianism

Preference utilitarianism is often understood as a modification of classical, hedonistic utilitarianism. On this conception of preference utilitarianism, it is a view that incorporates a preferentist theory of individual well-being together with the claim that morality requires maximizing the combined well-being of everyone affected by an action. This way of understanding preference utilitarianism, together with the fact that Hare's argument leads to preference utilitarian conclusions, has led many critics to respond to Hare by raising objections to a preferentist theory of well-being.[12]

---

[12]For example, Thomas Nagel, 'The Foundations of Impartiality', *Hare and Critics: Essays on Moral Thinking*, ed. Douglas Seanor and Nicholas Fotion (Oxford, 1988), pp. 101–12; James Griffin, 'Well-Being and its Interpersonal Comparability', *Hare and Critics: Essays on Moral Thinking*, ed. Douglas Seanor and Nicholas Fotion (Oxford, 1988), pp. 73–88; and John Harsanyi, 'Problems with Act-Utilitarianism and Malevolent Preferences', *Hare and Critics: Essays on Moral Thinking*, ed. Douglas Seanor and Nicholas Fotion (Oxford, 1988), pp. 89–99.

Such objections include the idea that even purely self-interested preferences can be defective, in the sense that an individual may mistakenly prefer something that is actually bad for her, the idea that an individual may prefer to sacrifice her own well-being for the sake of others whom she cares about, and the idea that an individual can have preferences for things that seem unconnected to her well-being, including things that occur after her death.[13] There is, of course, room for debate over the effectiveness of these objections.[14] But the relevant point here is that Hare's analysis is not hostage to such debates. As demonstrated by the preceding reconstruction, Hare's analysis neither depends on nor commits him to any theory of well-being, preferentist or otherwise. In order to assess Hare's project properly, one must engage with his actual argument, rather than raising objections to other lines of analysis that lead to equivalent conclusions.

The temptation to respond to Hare by rejecting a preferentist principle of well-being may be encouraged by the fact that Hare occasionally slips and makes comments suggesting a commitment to such a principle.[15] But for the most part Hare is quite clear about the distinction between offering an argument that generates substantive conclusions equivalent to those of preference utilitarianism, which he attempts to do, and offering a preference utilitarian argument that is grounded in a preferentist theory of well-being, which he does not do.

## II.2. Universalizability

There might be some temptation to deny Hare's premise of universalizability, but the fact that the premise applies narrowly to only those sets of cases in which nothing is changed but the role one occupies makes it difficult to mount much of an objection. Consider someone who asserts, 'It is morally permissible for me to do this to you, even though it would not be morally permissible for you to do it to me', and acknowledges that the only difference between the cases is the swapping of the roles. Is this a coherent moral claim? It seems to me that it is not, and that Hare is right to accuse this person of not merely having an incorrect moral view but of failing to understand the idea of moral permissibility. In other words, it seems to me that Hare's sense of universalizability really is built into the concept of morality, and a 'moral view' that directly violates universalizability fails to be a coherent moral view. Notice, in contrast, that it is perfectly coherent to say 'I want to do this to you, even though I would not want you to do it to me', while acknowledging that the only difference between the cases is the swapping of the roles. Desires are not constrained by universalizability, but I am inclined to say that Hare is correct that moral judgements are.

Moreover, for readers who are not willing to grant Hare this conceptual claim, there is an alternative version of his argument available. All that is needed is to replace the conceptual metaethical claim about universalizability with a substantive normative version of the claim. This involves adopting the explicitly moral premise that the mere

---

[13] See Mark Overvold, 'Self-Interest and the Concept of Self-Sacrifice', *Canadian Journal of Philosophy* 10 (1980), pp. 105–18; Derek Parfit, *Reasons and Persons* (Oxford, 1984); Richard Kraut, 'Desire and the Human Good', *Proceedings and Addresses of the American Philosophical Association* 68 (1994), pp. 39–54; and Robert Adams, *Finite and Infinite Goods: A Framework for Ethics* (Oxford, 1999).

[14] See Chris Heathwood, 'The Problem of Defective Desires', *Australasian Journal of Philosophy* 83 (2005), pp. 487–504 and Chris Heathwood, 'Preferentism and Self-Sacrifice', *Pacific Philosophical Quarterly* 92 (2011), pp. 18–38.

[15] For example, Hare, *Moral Thinking*, p. 42.

swapping of roles is morally irrelevant. In one sense, shifting to this alternative version of the argument would constitute a significant concession on Hare's part. It would involve giving up his ambition of deriving substantive moral conclusions purely from the logic of moral thought and language. Nonetheless, if an argument like Hare's could get all the way to preference utilitarian conclusions from the innocuous moral premise that the mere swapping of roles is not morally significant, that would still be an astonishing result, albeit less ambitious than the result Hare himself hoped for.

## II.3. Prescriptivity

Unlike Hare's universalizability premise, which I take to be very difficult to deny, his appeal to prescriptivity is more problematic. His formulation and use of the premise involves controversial assumptions that many reject, thereby limiting the power of his analysis. Nonetheless, I will argue that the difficulties Hare encounters here can be avoided, although it will require modifying the core argument articulated above on his behalf.

Part of the trouble here for Hare stems from the fact that his prescriptivity premise appears in more than one form. When he initially introduces the idea in *Moral Thinking*, Hare defines prescriptivity as 'the property of entailing at least one imperative'.[16] Later, when he prepares to use the premise in the analysis reconstructed above, he identifies prescriptions as the linguistic representation of preferences, in particular what I have been calling preferences all told.[17]

One way to understand the difference between these two versions of the prescriptivity premise is to view the earlier version as a generic assertion of metaethical expressivism, according to which moral language fundamentally expresses certain of the speaker's attitudes rather than serving to describe the world. The later version can then be understood as an assertion of Hare's particular form of expressivism: that the attitudes expressed by moral judgements are preferences all told.

It is worth noting here that the later version of the prescriptivity premise plays a critical role in Hare's analysis. It is the idea that an agent's moral judgements are constituted by preferences that allows Hare's claim about the conditional reflection of hypothetical preferences to get a grip on those judgements, and it is the focus on preferences that leads his analysis to end up endorsing preference utilitarian conclusions.

What is troubling about this is that Hare's own version of expressivism is not widely endorsed, even among contemporary expressivists. One reason for this is that Hare's version of expressivism builds in an exceptionally strong connection between moral judgement and motivation, which is to say an exceptionally strong form of judgement internalism.[18] On Hare's view, if an agent is deliberating about what to do and she comes to the conclusion that one of her available actions is morally required, she must therefore have a preference all told to perform that action. Although expressivists commonly endorse a weaker form of judgement internalism, according to which judging that an action is morally required necessarily involves having some motivation to perform the action (perhaps overridden by other competing motivations), few, if any,

---

[16] Hare, *Moral Thinking*, p. 21.

[17] Hare, *Moral Thinking*, p. 107.

[18] I adopt this terminology from Stephen Darwall, 'Reasons, Motives, and the Demands of Morality', *Moral Discourse and Practice*, ed. Stephen Darwall, Allan Gibbard and Peter Railton (Oxford, 1997), pp. 305–12.

are willing to endorse the strong internalism that follows from Hare's version of expressivism.[19] One way to see why Hare's strong internalism is problematic is to notice that his view makes it impossible for an agent to judge an action morally required and then continue to deliberate over whether to do it. Any continued deliberation must instead be interpreted as reconsidering the question of whether the action is morally required. Similarly, *that this action is morally required* cannot be a consideration the agent appeals to in bolstering her motivation to perform it.

One worry for Hare's analysis, therefore, is that his argument may depend on a particular version of expressivism which very few, if any, contemporary expressivists continue to endorse. Resolving this worry would require formulating an alternate version of Hare's core argument that avoids presupposing that moral assertions are expressions of preferences all told. Perhaps this could be accomplished by shifting to Hare's initial characterization of prescriptivity, understood as a generic assertion of expressivism without a specific view of the attitudes expressed in moral discourse. Merely doing that, however, would leave open the further problem that expressivism itself is highly controversial. Hare's core argument would still be quite interesting and important if it could derive preference utilitarian conclusions from expressivism, which is often viewed as neutral about the content of morality, or neutral with respect to normative ethics, but including expressivism as a premise nonetheless greatly restricts the power of the argument.

Perhaps, though, there is an even weaker version of the prescriptivity premise available, one that is neutral between expressivism and the descriptivist or cognitivist alternatives. After all, opponents of expressivism are often willing to acknowledge that moral discourse is importantly prescriptive, even if they do not take that prescriptivity to be what defines the meaning of moral terms. For instance, in discussing a person's good, or what is good for a person, Peter Railton writes: 'someone who spoke in earnest to others about their own good, and then was simply puzzled when they took his remarks to be any sort of recommendation, would betray a lack of full competence with such discourse'.[20] But he does not think this element of prescriptivity forces an expressivist or non-cognitivist analysis of value discourse, and instead offers an analysis of such discourse as fundamentally descriptive or cognitive. He does believe he owes an account of 'how an essentially descriptive use of language could have the prescriptive force of value discourse', but he thinks this challenge can be met.[21] Similarly, in the context of an argument for moral cognitivism, and against judgement internalism, Sigrún Svavarsdóttir writes that internalism 'probably appeals to many because they think it is a way of rendering more precise the plausible – possibly platitudinous – claim that the point of moral evaluation is distinctively to guide conduct'.[22] Again, the point here is that Svavarsdóttir, in spite of rejecting expressivism, concedes that there is a platitudinous sense in which moral evaluations are prescriptive.

My aim, therefore, will be to explore a version of Hare's argument that appeals to prescriptivity only in the weaker sense that, as Svavarsdóttir puts it, the point of moral evaluation is distinctively to guide conduct. In other words, I will take moral

---

[19] Compare Sigrún Svavarsdóttir, 'Moral Cognition and Motivation', *Philosophical Review* 108 (1999), pp. 161–219, at 172–3, n. 21. After explaining the strength of Hare's version of judgement internalism, she writes: 'Hare, as far as I know, is alone in holding this stronger thesis'.

[20] Peter Railton, 'Naturalism and Prescriptivity', *Social Philosophy and Policy* 7 (1989), pp. 151–74, at 151.

[21] Railton, 'Naturalism', p. 154.

[22] Svavarsdóttir, 'Moral Cognition', p. 218.

judgements to be prescriptive in the sense that they aim to tell us what to do, but I will avoid assuming Hare's strong sense of prescriptivity according to which moral judgements are expressions of preferences all told. And I will remain neutral on whether one must adopt an expressivist analysis of moral language to account for this weaker sense of prescriptivity, or whether cognitivists or descriptivists could follow Railton's or Svavarsdóttir's lead and offer a non-expressivist analysis that is compatible with this weaker prescriptivity premise.

Perhaps some might think Railton and Svavarsdóttir concede too much here, and might instead be inclined to deny that moral thought and language aim in any sense at guiding conduct or settling questions of what to do, but I agree with Railton that this would betray a lack of understanding of moral reasoning and discourse. To put it another way, if it were to turn out that moral language is not distinctively connected to guidance of conduct, I think that we, as creatures who live together in community and have substantial interests in one another's conduct, would have to invent other conduct-guiding language – and whatever we would invent would function like a version of our current moral language of which the weak prescriptivity premise is true.

## III. Revising Hare's argument

In terms of the project of this article, the upshot is that both this weaker prescriptivity premise and the universalizability premise discussed above are very minimal starting points. The next steps are to evaluate Hare's remaining premise and then to see how far his argument can get from these premises towards utilitarian conclusions. First, however, I need to reformulate the argument in a way that weeds out Hare's more controversial version of the prescriptivity premise.

Let me begin by introducing the idea of prescription tendencies and prescriptions all told, related to one another in the same way as preference tendencies and preferences all told. This allows me to replace the earlier claim that *in order to judge that X is morally permissible in $S_1$ I must have a preference all told that permits X in $S_1$, and I must also be committed to a matching preference all told that permits X in $S_2$* with a revised version: in order to judge that X is morally permissible in $S_1$ I must morally prescribe all told in a way that permits X in $S_1$, and I must also be committed to a matching moral prescription all told that permits X in $S_2$.[23]

On a reformulated version of Hare's conditional reflection thesis that applies to prescription tendencies rather than preference tendencies, the result of conditional reflection on $S_2$ will be that I have a set of prescription tendencies for $S_1$ that are constituted by my initial prescription tendencies in $S_1$, and a set of prescription tendencies for $S_2$ that are equivalent to your initial prescription tendencies in $S_1$. And, as before, the only way to ensure that I end up with matching prescriptions all told for both cases is to incorporate copies of my prescription tendencies for $S_2$ into my prescription all told for $S_1$, and vice versa.

The upshot of this, much as before, is that if I have reflected carefully on $S_2$, and my prescriptions are consistent with my commitments, I can judge that X is morally permissible in $S_1$ only if a utilitarian-style calculation incorporating both my prescription

---

[23]I have formulated this claim in terms of 'morally prescribing' and 'moral prescriptions' in order to restrict myself to a sense of prescriptivity that is sufficiently weak for all to endorse. In the ensuing discussion, I will drop the 'morally' modifier for simplicity of prose, but it will always be implied.

tendencies in $S_1$ and your prescription tendencies in $S_1$ leads to a prescription all told that permits X in $S_1$. This conclusion is not identical to the preference utilitarian conclusion of Hare's own argument, given the shift from preferences to prescriptions, but the difference is relatively small, which demonstrates that something much like Hare's own argument can be constructed without relying on his controversial formulation of the prescriptivity premise.

## IV. Conditional reflection

Although the revised version of Hare's argument avoids the problems associated with his conception of prescriptivity, it continues to rely on a conditional reflection thesis, now applied to prescriptions instead of preferences. In order to assess the revised version of Hare's argument, we must therefore consider whether full representation of one's hypothetical prescriptions requires formation of matching conditional prescriptions for the hypothetical case.

Adequately assessing this issue requires first acknowledging a complication within Hare's own analysis that I have suppressed up to this point. As I have characterized it, Hare's argument takes as input existing preferences. This overlooks Hare's adoption of Richard Brandt's idea of cognitive psychotherapy, which is a method of adjusting preferences through repeated vivid representation of relevant available information.[24] I have omitted this aspect of Hare's view in my initial presentation of his argument partly out of a desire for simplicity, but also in part because I think it is unclear whether the appeal to cognitive psychotherapy makes sense in the context of Hare's own analysis. In particular, I have in mind Hare's explanation and justification of the conditional reflection thesis.[25] As I understand it, his claim is that if I am considering a person in a situation and I both fully represent that person's preferences to myself and also fully identify with that person, then I will automatically form conditional preferences for the situation that match that person's preferences. I take it that this could be construed as a psychological claim, about how preference formation works, but also as a conceptual claim about what it is to represent another's preferences, or know them, and what it is to identify with another, or to consider carefully the hypothetical case in which I am in their role.

As a claim about preferences, representations of preferences, and identification, I find the conditional reflection thesis difficult to assess. But I suspect that if it is true then it applies to the preferences the other person has, or is imagined to have, not the preferences the other person would have after undergoing cognitive psychotherapy. After all, the appeal of a principle of conditional reflection of preferences seems to be that there is something about preferences, or about desires and aversions, such that if I think of a situation in which I, myself, want something to happen I must end up wanting that thing now for the case in which I am in that situation. Otherwise, the claim is, I am not really thinking of the agent in the situation as myself, or else I am not fully representing what it is to have the desire or preference. I am not sure whether this claim should be granted, but insofar as it is plausible I take it to be grounded in something about what it is like to have a preference and what it is to think of a person who has a preference as oneself. If that is right, it is not clear why the principle of conditional

---

[24]Hare, *Moral Thinking*, pp. 101–6 and Richard B. Brandt, *A Theory of the Good and the Right* (Oxford, 1979).

[25]Hare, *Moral Thinking*, pp. 94–9.

reflection would cease to apply to preferences that I imagine having merely on the grounds that those preferences would not survive cognitive psychotherapy, and perhaps even less clear why it would apply to preferences that I imagine not yet having but that would be produced through cognitive psychotherapy. If the conditional reflection thesis is correct, I would of course conditionally reflect the preferences that result from cognitive psychotherapy when considering a case in which I have undergone cognitive psychotherapy, and therefore have those preferences. But I am not sure why I would not simply reflect the preferences that I imagine having when I consider a case in which I have not yet undergone cognitive psychotherapy.[26]

To repeat, I find this issue puzzling. I am not entirely sure whether to accept Hare's conditional reflection thesis, applied to preferences. And if so, I am not sure whether to accept it as applying to existing preferences or to preferences that would result from cognitive psychotherapy. Although I have attempted to explain why I am inclined to apply conditional reflection directly to preferences one has in a hypothetical scenario rather than to preferences one would have on the additional hypothesis of having undergone cognitive psychotherapy, I concede that the considerations I have articulated may not be decisive.

Fortunately, I think questions of this sort become easier to settle when we consider conditional reflection of prescriptions rather than conditional reflection of preferences, and that is the sort of conditional reflection that is relevant to assessing the revised version of Hare's argument. Unfortunately, although I have just suggested that the most plausible conditional reflection thesis with respect to preferences may be one that applies straightforwardly to all preferences, and is not restricted even to those that survive cognitive psychotherapy, I also think that conditional reflection of prescriptions is appropriately subject to significant restrictions, beyond those related to cognitive psychotherapy, that will end up imposing limitations on Hare's line of analysis.[27]

In other words, although the shift away from Hare's strong sense of prescriptivity to a less controversial weaker sense did not on its own undermine Hare's analysis, my view is that the corresponding shift to a principle of conditional reflection of prescriptions requires restrictions that may have been avoidable if we were able to stick with Hare's principle of conditional reflection of preferences. To begin with, consider prescriptions that are grounded in false beliefs. Suppose I am considering the situation of being someone whose loved ones have been murdered. I believe the murderers have been identified, and I prescribe that they be executed. But the people I believe to be the murderers are actually innocent. As I, now knowing the alleged murderers to be innocent, consider the situation just described, must I formulate a prescription tendency in favour of execution that matches the prescription tendency I have in the situation? Some care is needed in answering this question. After all, the conditional reflection premise does not require that my prescription all told for the situation support execution. All it demands is that a conditionally reflected prescription be incorporated into my prescription all told. Nonetheless, even this weaker claim seems unfounded. Hare may be correct that the nature of preferences, representation and identification is such that I cannot fully represent a preference and fully identify with the person whose preference it is without

---

[26]To be clear, this is not an objection to the idea of cognitive psychotherapy or the way Brandt uses it in his own analysis, just a worry about a combination of cognitive psychotherapy with the conditional reflection thesis.

[27]Hare's own view is that conditional reflection applies equally to preferences and prescriptions, but that is because he thinks of prescriptions as expressions of preferences. See Hare, *Moral Thinking*, p. 222.

conditionally reflecting the preference. Knowing that I would prefer something in a certain situation may make me prefer it now for that situation. But prescriptions do not work this way, or at least not in this sort of case. Instead, I contend that the fact that I would prescribe execution under the false belief that the murderers had been correctly identified gives me, with my knowledge of their innocence, no reason at all to prescribe execution. Or if that seems too strong, because the mere fact that the victim wants those perceived as guilty to be executed may seem to give *some* reason to prescribe execution, the important point is that full conditional reflection of the misinformed prescription would give that misinformed prescription far too much weight. The key idea here is that prescriptions, if we don't think of them as expressions of preferences all told but instead simply as action-guiding judgements, open up the opportunity for critical distance. I can fully understand what it is to endorse a prescription and carefully consider the case in which I do, but if that prescription is grounded in what I take to be a defective belief I need not conditionally reflect it.[28]

This might seem to be simply an application of Brandt's cognitive psychotherapy to prescriptions in advance of conditional reflection. But notice that cognitive psychotherapy only requires exposure to relevant available information.[29] This makes sense given Brandt's aim of identifying desires as rational or irrational, which is an assessment that is relative to the agent who has the desire. But in a case in which I am to conditionally reflect your prescription, what matters is not merely that the prescription is appropriate relative to the relevant information available to you, but that the prescription is appropriate relative to the information I have, or perhaps the relevant information available to me.

Moving beyond cases involving ordinary false beliefs, a variant on the execution case reveals an additional, perhaps more significant, restriction on the conditional reflection of prescriptions. Suppose now that the accused really are guilty, and in the case I am imagining I prescribe their execution. But what if I, now, believe that execution is barbaric and unjust, not merely for pragmatic reasons involving the possible innocence of the accused but for moral reasons relating to the value of human life? If that is my view of capital punishment, when I come to consider a situation in which I believe that those guilty of murder deserve execution, does that belief, which I take to be grounded in an evaluative error, give me any reason to prescribe execution? I contend not – and again the point is not merely that it fails to give me an overriding or all things considered reason to prescribe execution, but that it gives me no reason at all, or at most a very weak reason. After all, the case is set up such that I am imagining a situation in which I take myself to be mistaken. To conditionally reflect the prescription would compound the mistake: no longer would I be innocuously imagining myself making a mistake, instead I would be actually making the initially imagined mistake.

As before, the claim being made here goes beyond the application of cognitive psychotherapy and the restriction of conditional reflection to whatever prescriptions survive or are produced through that process. A prescription could survive cognitive psychotherapy, which is relative to the agent who undergoes it, and yet still be grounded in what I take to be an evaluative error or objectionable value, in which case my claim is that the conditional reflection thesis would not apply to it. Preferences may, as Hare believes, work otherwise. But, again, I take it that prescriptions offer a critical distance that preferences may lack. Fully imagining myself with a preference that I take to be

---

[28]My ambivalence above about conditional reflection of preferences grows out of uncertainty over whether this critical distance is available with respect to preferences. According to Hare, it is not.

[29]Brandt, *A Theory*, pp. 10–13 and 113–16.

grounded in a mistake may require compounding the mistake by forming a matching conditional preference, but fully imagining myself prescribing in a way that I take to be grounded in a mistake imposes no such requirement.

It is worth noting that Hare himself provides extended discussion of cases involving conditional reflection of preferences that one takes to be grounded in evaluative error or objectionable values.[30] His view is that even in these cases the principle of conditional reflection applies. He then goes on to argue, in connection with individuals he calls 'fanatics', that incorporating these conditionally reflected preferences into my preference all told for the case is highly unlikely to make a difference, and that in the rare instances in which it does make a difference the surprising results do not constitute objections to his argument: careful consideration of strange cases will yield strange conclusions. And this last claim feeds into a powerful line of argument Hare offers in response to intuitively grounded objections to his view, which takes advantage of his distinction between two levels of moral thinking.[31]

My aim here is not to call these elements of Hare's analysis into doubt. In fact, my view is that if we grant Hare's strong version of prescriptivity and his principle of conditional reflection of preferences, then his response to worries about fanatics and what he calls evil desires is successful. Nonetheless, I have been arguing that Hare is entitled to only a weaker version of prescriptivity, which forces us to consider conditional reflection of prescriptions rather than preferences, and that conditional reflection of prescriptions does not apply in cases involving prescriptions that one takes to be grounded in evaluative errors. If I am right, the point is not that Hare's response to the fanatic is inadequate, but that it is not relevant to the cases I am considering.[32]

So far, I have been focusing on cases in which the principle of conditional reflection of prescriptions does not apply. I have argued that when hypothetical prescriptions are grounded in beliefs or values that I reject, I can fully and vividly consider the hypothetical situation and nonetheless refrain from conditional reflection. What, though, of cases in which I do not view my hypothetical prescriptions as grounded in defective beliefs or objectionable values? Here I think something analogous to Hare's line of thought is compelling. If I am supposing that in a given situation I would form certain prescription tendencies, and if I do not reject any of the beliefs or values on which those prescription tendencies are based, then I think the Hare-inspired view is correct: I cannot fully represent those prescription tendencies and carefully consider the hypothetical case in which I am the one who has them without forming matching prescription tendencies for that case. I take this to be a psychological and conceptual claim, comparable to Hare's own principle of conditional reflection of preferences, grounded in the nature of prescriptions, full representation or knowledge, and identification. If I am considering a case in which I make a prescription, and I do not view the prescription as grounded in defective beliefs or objectionable values, the only ways I can avoid forming a matching conditional prescription is if I do not fully represent the case to myself or if I do not fully identify with the agent in the hypothetical case.

Moreover, if I am considering a situation in which I form prescription tendency $P_1$ on the basis of what I take to be a defective belief, and in which I would instead form

---

[30]Hare, *Moral Thinking*, pp. 140–6 and 170–82.
[31]Hare, *Moral Thinking*, pp. 130–46.
[32]In my view, there are other cases in which Hare's response to the fanatic, and his related response to counter-intuitive implications of his view, has a role to play in defence of the restricted Hare-like argument I will be endorsing, although discussion of that goes beyond the scope of this article.

prescription tendency $P_2$ if my beliefs were not defective, careful reflection on the case, I contend, requires conditional reflection of the epistemically informed prescription tendency $P_2$. Similarly, if the situation under consideration is one in which I form prescription tendency $P_3$ on the basis of what I take to be an objectionable value, and in which I would instead form prescription tendency $P_4$ were I not committed to that objectionable value, careful reflection requires conditional reflection of the evaluatively informed prescription tendency $P_4$. Taken together, this adds up to an appropriately tempered version of Hare's principle of conditional reflection that demands conditional reflection of what the reflecting agent takes to be epistemically and evaluatively informed prescription tendencies.[33]

## V. Restricted Hare-like argument

With this restricted version of conditional reflection in place, I can now formulate what I will call my restricted Hare-like argument. This argument has substantial differences from Hare's, and the conclusion it aims at is notably weaker than what Hare hoped to establish, but it is nonetheless inspired by Hare and follows his argumentative framework.

I will begin with a claim from the revised version of Hare's argument above, which is that in order to judge that X is morally permissible in $S_1$ I must prescribe in a way that permits X in $S_1$, and I must also be committed to a matching prescription that permits X in $S_2$. Applying the restricted conditional reflection thesis defended in the previous section, the result of conditional reflection on $S_2$ will be that I have a set of prescription tendencies for $S_1$ that are constituted by my initial prescription tendencies in $S_1$, and a set of prescription tendencies for $S_2$ that are equivalent to what I take to be your epistemically and evaluatively informed initial prescription tendencies in $S_1$. In other words, I will develop prescription tendencies regarding $S_2$ that match your prescription tendencies in $S_1$, corrected for what I take to be defective beliefs or objectionable values. As before, the only way to ensure that I end up with matching prescriptions all told for both cases is to incorporate copies of my prescription tendencies for $S_2$ into my prescription all told for $S_1$, and vice versa.

The upshot of this is that if I have reflected carefully on $S_2$, and my prescriptions are consistent with my commitments, I can judge that X is morally permissible in $S_1$ only if a utilitarian-style calculation incorporating both my prescription tendencies in $S_1$ and what I take to be your epistemically and evaluatively informed prescription tendencies in $S_1$ leads to a prescription all told that permits X in $S_1$.

As should be obvious, this conclusion is less ambitious than the conclusion Hare sought to establish. One difference is the shift from formulating the conclusion in terms of preferences and preference tendencies to prescriptions and prescription tendencies, but as indicated above I take that to be a relatively minor difference that, on its own, would amount to only a slight revision of Hare's utilitarian conclusions. More important is the way in which the conclusion of my restricted Hare-like argument licenses the revision, or perhaps even outright exclusion, of some actual prescription tendencies. The significance of the argument depends on the impact this difference makes.

---

[33]Note that in cases in which I think others fail to recognize their own worth, perhaps as a result of conditioning or objectionably adaptive preference formation, I can be required to conditionally reflect what I take to be evaluatively informed prescription tendencies that place greater weight on their own interests and well-being than they do.

There should be no doubt that the impact is large. Hare's own argument, if successful, would show that for any situation there is an objective moral assessment in the following sense. Everyone considering the situation, if sufficiently well-informed and clear-headed, would come to a consensus on their moral judgements. And that consensus would take a preference utilitarian form. A failure to reach consensus on a preference utilitarian conclusion would mean that at least one of them was not fully informed or not reflecting carefully enough, or else not forming genuinely moral judgements. As this hopefully makes clear, part of the magnificence of Hare's analysis is the Kantian nature of its ambition. Hare hopes to establish that the preference utilitarian moral judgement regarding a situation is not merely the moral judgement that has the advantage against competing moral judgements of being correct, which would be an ambitious conclusion on its own, but that the preference utilitarian conclusion is the only assessment an informed and reflective judge can reach that even counts as a moral judgement.

In contrast, my restricted Hare-like argument is compatible with different parties to the situation, all fully informed and carefully reflective, forming incompatible moral judgements. If I reject some of your prescription tendencies as grounded in objectionable values, we can reach conflicting conclusions that satisfy the universalizability and prescriptivity conditions needed for those conclusions to count as moral judgements. This does not mean that my restricted Hare-like argument leads to relativism. My restricted Hare-like argument does not imply that all such judgements are correct, relative to the person making them. Rather, it allows that they may all be genuinely moral judgements, leaving it as a matter of further investigation whether one, the other, or perhaps both, are correct.

So my restricted Hare-like argument really does represent a substantial step back from Hare's own project. But does it step so far back as to be insignificant? Or is the conclusion of my restricted Hare-like argument noteworthy, even if it falls short of achieving Hare's goal?

It might seem as though it is the former, and that the fact that my restricted Hare-like argument permits the exclusion or revision of others' prescription tendencies robs the argument of any real force. After all, if I am only required to conditionally reflect and incorporate into my all things considered moral prescription what I take to be your epistemically and evaluatively informed prescription tendencies, the process of conditional reflection and incorporation might seem incapable of having any impact on my all things considered judgement. The idea here is that I will be able to exclude any prescription tendencies of yours that differ from my initial prescription tendencies, or revise them until they agree with mine. And if I am only required to conditionally reflect and incorporate prescription tendencies that harmonize with my initial prescription tendencies, then I can comply with the demands of the restricted Hare-like argument and nonetheless ensure that my all things considered prescription is fully determined by my initial prescription tendencies, whatever they are.

The problem with this line of thought is that it assumes I can reasonably judge any of your prescription tendencies that conflict with mine to be grounded in beliefs or values that I reject. But this may not always be the case. It is also possible that the difference in our initial prescription tendencies is grounded not in divergent beliefs or values but in differing perspectives. Perhaps you and I both value happiness, but I focus only on the happiness that manifests in my own life or in the lives of those I care about. I may not initially care at all about your happiness, but when I consider the hypothetical situation in which I occupy your actual role, I have no legitimate

grounds for refusing to conditionally reflect the hypothetical prescription tendencies grounded in my hypothetical happiness – which are precisely the same as your actual prescription tendencies grounded in your actual happiness. I therefore end up with conditional prescription tendencies for the case of being in your role that mirror your actual prescription tendencies, or at least some of them. And, working through the argument above, I must treat these conditional prescription tendencies as operative for the actual case, and my actual prescription tendencies as operative for the hypothetical case, in order to have matching prescriptions for both cases. In other words, if I am to make a universalizable prescription all told for the actual case, which I must do in order to count the prescription as a moral one, I must incorporate at least some of your initial prescription tendencies together with mine. And these incorporated prescription tendencies may very well be ones that conflict with my initial prescription tendencies, which demonstrates that my restricted Hare-like argument can still be quite potent. I cannot blithely dismiss all of your prescription tendencies that do not facially conform to my own prescription tendencies, but may instead be required to incorporate at least some of them, even though doing so may make a real impact on my resulting moral judgement.

Another way to see the power of my restricted Hare-like argument is to notice that it serves as a basis for what Gibbard has dubbed the 'You'd have agreed' retort.[34] Gibbard's idea is that we are to imagine someone raising an objection to an existing social practice, on the grounds that they are unfairly disadvantaged by the practice. The 'You'd have agreed' retort amounts to pointing out that the very same interests and values that underlie the objection would have led the objector to endorse the practice if they had not known in advance who they were or what role in the situation they would occupy. Gibbard finds this retort intuitively compelling, and he uses it as the basis for a defence of Harsanyi's project of deriving utilitarian conclusions from a contractualist framework, and as a key tool in his argument against contractualists who reject utilitarianism, such as John Rawls and T. M. Scanlon.[35] Nonetheless, Gibbard is unable to formulate an argument on behalf of the 'You'd have agreed' retort. As he puts it, 'if the retort leaves someone cold who genuinely understands what it involves, then I don't know anything to say that would make the person responsive'.[36]

In earlier work I have been forced to echo Gibbard on this point, endorsing his intuitive sense of the force of the retort while sharing in his professed inability to defend it.[37] Now, with my restricted Hare-like argument in hand, I can go further. The retort turns out to be grounded in the universalizability and weak prescriptivity of moral language, together with the requirement to conditionally reflect hypothetical prescriptions that one takes to be epistemically and evaluatively informed. Again, working through my restricted Hare-like argument above, the idea is that if the person raising an objection is to count the objection as a moral prescription, she must be committed to a matching prescription for hypothetical cases in which she occupies the roles of others impacted by the social practice. Any hypothetical prescription tendencies she has in those cases that are grounded in the same interests and values that generate the objection must be conditionally reflected and become prescription tendencies for the relevant hypothetical

---

[34]Gibbard, *Reconciling*, p. 50.

[35]John Rawls, *A Theory of Justice* (Cambridge, MA, 1971); John Rawls, *Justice as Fairness: A Restatement* (Cambridge, MA, 2001); and T. M. Scanlon, *What We Owe to Each Other* (Cambridge, MA, 1998).

[36]Gibbard, *Reconciling*, p. 152.

[37]Steven Daskal, 'Original Position Models, Trade-Offs and Continuity', *Utilitas* 28 (2016), pp. 254–87.

cases. And the only way to end up with equivalent prescriptions all told for each of the hypothetical cases is to treat all of these conditional prescription tendencies as operative for each case, including the actual case. As a result, any moral prescription all told she arrives at for the actual case must incorporate the relevant prescription tendencies of everyone, where 'relevant' means based in the same values and interests that led to the objection. And this just is to say that if the values and interests that underlie the objection would have led the objector to endorse the social practice had she not known her role, then the objector cannot maintain that the objection has moral force, which is to say that the 'You'd have agreed' retort is decisive.

This, I take it, demonstrates that my restricted Hare-like argument offered above has real significance. Together with ideas developed by Harsanyi and Gibbard, it may even be an important part of a larger argument for more ambitious utilitarian conclusions resembling Hare's. Whether that argument succeeds is beyond the scope of this article, but regardless of that the ability of my restricted Hare-like argument to vindicate Gibbard's 'You'd have agreed' retort reveals that Hare's argumentative strategy retains power and relevance.[38]

## VI. Conclusion

Hare's project truly is ambitious. I have attempted to show that some common lines of response miss their mark, but I have also identified what I take to be a significant problem for an unrestricted application of his principle of conditional reflection. Nonetheless, I think it would be a mistake to dismiss Hare's fundamental argumentative strategy, and I have tried to develop a Hare-like argument that has real import, albeit less than what Hare sought.

There are, of course, potentially serious objections to Hare, and to Hare-like arguments, that I have not addressed within this article, many of which Hare wrestles with. For instance, one might worry that such arguments illicitly derive normative conclusions from descriptive premises.[39] Additionally, one might wonder whether the arguments apply only to preferences or prescriptions that are grounded in self-regarding considerations or more generally to all preferences or prescriptions.[40] Relatedly, there are potential concerns about cases in which direct application of an argument like Hare's yields counterintuitive results.[41] A full defence of my Hare-like argument would require resolving these objections, and perhaps others, but that goes beyond the scope of this article.

Given that I advocate replacing his own argument with my less ambitious one, I suspect Hare himself would have rejected my analysis and viewed me as among the 'philosophical worms' who 'nibble away' at his life's work in the attempt to show that his 'achievement was an illusion'.[42] But I prefer to think of the project of this article as calling attention to the potency of Hare's argumentative strategy. If my analysis succeeds, it

---

[38]My own view, defended in Daskal, 'Original Position Models', is that Gibbard and Harsanyi, like Hare, are also not able to get all the way to traditional utilitarian conclusions. But I argue that they get relatively close, and that their analysis poses a more formidable challenge to non-utilitarian contractualists than is generally recognized. If I am right that the restricted Hare-like argument of this article vindicates the 'You'd have agreed' retort, that bolsters their argument by closing off one possible avenue of response.

[39]Hare, *Moral Thinking*, pp. 218–28.

[40]Hare, *Moral Thinking*, pp. 104–6.

[41]Hare, *Moral Thinking*, pp. 130–68.

[42]Hare, 'A Philosophical Autobiography', p. 269.

shows that Hare's line of thought can impose significant constraints on moral reasoning, and may also have a role to play in a larger argument for utilitarianism that matches Hare's ambition.[43]