

I—The Presidential Address

KNOWING WHAT YOU BELIEVE

QUASSIM CASSAM

A familiar claim is that knowledge of our own thoughts, beliefs and other attitudes is normally immediate, that is, not normally based on observation, inference or evidence. One explanation of the possibility of immediate self-knowledge turns on the transparency of the question ‘Do I believe that *P*?’ to the question ‘Is it the case that *P*?’ This paper explains why occurrent mental states such as passing thoughts do not fall within the purview of the transparency account and proposes a different account of how we know our own passing thoughts. It is also argued that the transparency account fails to explain how knowledge of our own beliefs can be psychologically or epistemically immediate. Finally, questions are raised about the presumption that knowledge of our own beliefs is epistemically immediate.

I

A familiar thesis is that knowledge of our own thoughts, beliefs and other attitudes is normally immediate, that is to say, not normally based on observation, evidence or inference.¹ Many philosophers who endorse this *immediacy thesis* take it to be obviously correct. They think that it is the immediacy of self-knowledge rather than its supposed infallibility which makes it epistemologically distinctive, and that the hard question in this area is not *whether* self-knowledge can be immediate but *how* it can be so. The immediacy of self-knowledge can seem puzzling because knowledge of our own thoughts and beliefs is surely knowledge of contingent facts, and the usual presumption is that knowledge of contingent matters must be

¹ See, for example, Davidson (1994) and Moran (2001).

based on observation, evidence or inference.² So if it is just a contingent fact about me that I believe that *P*, and yet I know without observation, inference or evidence that I have this belief, then it needs to be explained how this is possible.

One explanation, given by Richard Moran among others, appeals to the notion of transparency. The idea is that ‘a person answers the question whether he *believes* that *P* in the same way he would address himself to the question whether *P* itself’ (Moran 2004, p. 457). The question ‘Do you believe that *P*?’ is, in this sense, *transparent* to the question ‘Is it the case that *P*?’ which means that I can answer the former question ‘by consideration of the reasons in favour of *P* itself’ (Moran 2003, p. 405). The proposal is that this way of knowing one’s beliefs is non-observational, non-inferential and non-evidential, and that knowledge of the admittedly contingent fact that one believes that *P* is immediate because the question ‘Do you believe that *P*?’ is transparent to a corresponding question about the world. So the thesis is that *transparency explains immediacy* (TEI).³

Knowledge of one’s beliefs and other standing attitudes is not the only supposedly immediate self-knowledge. There is also knowledge of one’s sensations and passing thoughts, that is, thoughts that just occur to one.⁴ These are occurrent states, and TEI does not purport to explain how immediate knowledge of them is possible. Why not? Here is one suggestion: the states of mind that fall within the purview of TEI are propositional attitudes. Since sensations are not propositional attitudes it is not surprising that TEI has no bearing on them. But this can’t possibly be the whole story, because many of our thoughts, including passing thoughts, are propositional attitudes. So it is not clear why TEI doesn’t apply to passing thoughts if it applies to beliefs. And if TEI doesn’t explain how immediate knowledge of our passing thoughts is possible then what does explain it?

² As Boghossian (1998, p. 166) points out, there are contingent judgements which one may be justified in making ‘even in the absence of any empirical evidence’. An example is the judgement *I am here now*. Immediate knowledge of one’s own beliefs is, however, not relevantly similar to such examples, and cannot be explained in the way that my immediate knowledge that I am here now can be explained.

³ For a clear statement of this thesis see Moran (2004, p. 457).

⁴ Moran writes: ‘There are two basic categories of psychological state to which the ordinary assumption of “privileged access” is meant to apply: occurrent states such as sensations and passing thoughts, and various standing attitudes of the person, such as beliefs, emotional attitudes, and intentions’ (Moran 2001, p. 9). He is explicit that his account of self-knowledge is not intended to apply to sensations.

I have three main questions:

- (1) Is the immediacy thesis correct?
- (2) Is TEI correct?
- (3) How can we have immediate knowledge of our passing thoughts?

I'm going to take these questions in reverse order. The first thing to figure out in relation to (3) is why the transparency procedure can't account for our immediate knowledge of our own passing thoughts. My proposal is that among the relevant factors is what I'm going to refer to as the *passivity* of passing thoughts. Passing thoughts are passive in the sense that they are (i) not necessarily responsive to reason, and (ii) states from which one can distance or dissociate oneself. Between them (i) and (ii) help to explain why the transparency procedure can't deliver immediate knowledge of one's passing thoughts. My main positive suggestion in connection with (3) builds on Ryle's insight that 'much of our ordinary thinking is conducted in internal monologue or silent soliloquy' (Ryle 1949, p. 28). I claim that this remark points to a promising response to (3): I know that I am having the thought that *P* by being conscious of saying to myself that *P* in my internal monologue.⁵

With regard to (2), there is a problem with TEI. Even if I can answer the question 'Do I believe that *P*?' by consideration of the reasons in favour of *P* the resulting knowledge is not immediate. In fact, there are two notions of immediacy in play, one epistemic and the other psychological. My knowledge that *P* is *epistemically* immediate only if my justification for believing that *P* does not come, even in part, from my having justification to believe other, supporting, propositions.⁶ My knowledge that *P* is *psychologically* immediate only if it is not acquired by conscious reasoning or inference. If I come to know that I believe that *P* by employing the transparency procedure then my knowledge does not appear to be immediate in either sense. Since it is knowledge acquired by conscious reasoning it cannot possibly be psychologically immediate. It can't be epistemically immediate either because, as we will see, my justification for believing that I believe

⁵ Ryle's view is also discussed in Byrne (forthcoming), but Byrne wrongly insists that there is no such thing as inner speech.

⁶ This is essentially the account of immediate justification given in Pryor (2005).

that *P* does come, at least in part, from my having justification to believe other propositions. The transparency procedure gives me inferential self-knowledge and inferential knowledge isn't immediate.

It will turn out that some of the reasons for thinking that the transparency procedure doesn't give us immediate knowledge of our passing thoughts are also reasons for thinking that it doesn't give us immediate knowledge of our beliefs. If passing thoughts are passive the same is true of some beliefs. The latter are not necessarily responsive to reason, and there are also beliefs from which it is possible to distance oneself. The point is not that the passivity of belief is exactly parallel to the passivity of passing thoughts. Given what it is for a mental state to be a belief there are much tighter restrictions on the passivity of belief. All the same, to the extent that belief can be unresponsive to reason this puts extra pressure on the idea that I can come to know whether I believe that *P* by considering the reasons in favour of *P*.

That leaves (1). Is it obvious that self-knowledge is immediate? Perhaps the thought is this: if asked whether I believe that *P* I can usually answer straight off. I don't generally need to work out whether I believe that *P*, I just know. But this only shows that my self-knowledge is psychologically immediate, and this does not entail epistemic immediacy. In my view, it is not a datum that knowledge of our own beliefs is epistemically immediate, or that accounts of self-knowledge that fail to secure its epistemic immediacy are necessarily flawed. Still, the fact remains that *if* the object of the exercise is to explain how knowledge of our own beliefs could be epistemically immediate then the transparency account doesn't work.

II

Consider the following passage from Harry Frankfurt's paper 'Identification and Externality':

In our intellectual processes, we may be either active or passive. Turning one's mind in a certain direction, or deliberating systematically about a problem, are activities in which a person engages. But to some of the thoughts that occur in our minds ... we are mere passive bystanders. Thus there are obsessional thoughts, whose provenance may be obscure and of which we cannot rid ourselves; thoughts that strike us unexpectedly out of the blue; and thoughts that run willy-nilly

through our heads. The thoughts that beset us in these ways do not occur by our own active doing. It is tempting, indeed, to suggest that they are not thoughts that *we think* at all, but rather thoughts that we *find* occurring within us. This would express our sense that, although these thoughts are events in the histories of our minds, we do not participate actively in their occurrence. ... It is not incoherent, despite the air of paradox, to say that a thought that occurs in my mind may or may not be something that *I think*. (Frankfurt 1998, p. 59)

What Frankfurt describes here is a perfectly familiar phenomenon.⁷ The successive thoughts with respect to which we are passive bystanders are what I am calling passing thoughts, but passing thoughts come in different varieties. Obsessional thoughts are one variety. A jilted lover finds herself thinking obsessively about her ex and cannot concentrate on anything else. For all his assurances that no one else is involved, and despite the absence of convincing contrary evidence, she can't get rid of the thought that he is seeing someone else. In contrast, the thoughts that strike us out of the blue need not be obsessive. As I write these words it suddenly occurs to me that today is the first of the month. This thought is unconnected with the deliberative thinking in which I am currently engaged; it just comes to me, for no apparent reason, but I am hardly obsessed with today's date. The final category of passing thoughts, those that Frankfurt describes as running willy-nilly through our heads, might include the thoughts that run through one's head as one gazes absent-mindedly out of the window on a train journey or just before one falls asleep.⁸ The sense in which such thoughts run through our heads willy-nilly is that they are not the product of deliberation and are not in any obvious way rationally ordered. If we can be said to be thinking these thoughts at all, we do so 'not by concentrating on anything in particular' but by 'moving from one idea to the next in an endless chain of associations' (Bollas 2009, p. 6).

How do we know our passing thoughts? Suppose that one of my passing thoughts is the thought that *P*. How do I know that I am

⁷ It is also one that is noticed by Sellars. He points out that 'thinking is often a deliberate action, as in thinking about (i.e. attempting to solve) a problem' but that 'there is a sense of "thought" in which thoughts just occur to one. We say "it suddenly occurred to me that ..." and, we can often add "for no reason"' (Sellars 1975, Lecture 1, #29). It is in cases of the latter kind that we begin to see the full force of the Lichtenbergian suggestion that Descartes should have said 'There is thinking' rather than 'I think'.

⁸ Richard Moran suggested (in conversation) that the thoughts that run through one's head just as one is about to fall asleep are the best example of 'passing thoughts'.

thinking that *P* or having the thought that *P*? Not by employing the transparency method. The jilted lover does not know of the occurrence of the thought that her former lover is seeing someone else by asking ‘Is he seeing someone else?’, or by considering the reasons in favour the proposition that he is seeing someone else. She might realize that the reasons in favour of this proposition are weak but this doesn’t alter the sad fact that the thought that he is with someone else keeps her awake at night. In contrast, there might be good reasons for thinking that today is the first of the month but it is still implausible that I come to know that this is what I am thinking by consideration of such reasons. The same goes for the random thoughts that run through my head willy-nilly on a train journey or as I am falling asleep. I don’t come to know that I am having these thoughts by consideration of the reasons in favour of them or their contents.

The point is this: consideration of the reasons in favour of *P* only reveals whether I am thinking that *P* on the assumption that what I am thinking is somehow determined by those reasons or by my reflection on them.⁹ Such an assumption would be implausible in relation to passing thoughts. This is especially clear in relation to obsessional thoughts since part of what makes them obsessional is their unresponsiveness to reason. I can find myself thinking or worrying that *P* even though I realize that there is no reason to think or worry that *P*. If I realize that there is no good reason for thinking that *P* then asking whether it is true that *P* will not tell me that I am thinking that *P*. *Ex hypothesi*, there is no reason to think that *P* is true, but this doesn’t alter the fact that the thought that *P* keeps running through my head. In the case of non-obsessional passing thoughts the point is not that they are not responsive to reason but that they aren’t necessarily responses to reason. For example, the thought that today is the first of the month might be extinguished by evidence to the contrary. It is in this sense responsive to reason, and yet it did not occur to me that today is the first of the month because I gave the matter any thought. The thought came to me for no apparent reason, and if the thought that *P* is not the product of rational reflection then it is hard to see why reflection on the reasons in favour of *P* should have any bearing on whether one is in fact thinking it.

⁹ The point I am making here is analogous to one that Moran makes when he is discussing the conditions under which consideration of the reasons in favour of *P* can tell me whether I believe that *P*. See Moran (2003, p. 405) and §III below for further discussion.

Another way of making this point would be to note that one can distance oneself from one's passing thoughts. Distancing oneself from a thought means recognizing that it is not well-founded or is in some other way inappropriate. Sitting on a plane, 'the thought that it is going to crash comes to me, and perhaps I find that I cannot help thinking that it will crash' (Hampshire 1965, pp. 101–2). Here it is a 'datum of consciousness' (p. 101) that I am thinking that the plane will crash even though I am well aware that there is no real evidence that it will crash. In this case, I am effectively alienated from one of my own occurrent thoughts. I do not treat the question whether I am thinking that *P* as equivalent to the question whether *P* is true. By the same token, I don't come to know that I am thinking that *P* by asking myself whether *P* is true.

To say that passing thoughts are not necessarily responsive, or responses, to reason and that one can distance oneself from one's passing thoughts is to say that such thoughts are ones with respect to which one is *passive*. What I am suggesting is that it is the passivity of passing thoughts that partly explains why they fall outside the purview of TEI. It is not just that the transparency account fails to explain how knowledge of our passing thoughts can be *immediate*. The deeper worry about this account is that it fails as an account of *knowledge* of our passing thoughts. A different account is needed, one that does justice to the passivity of passing thoughts and to the sense in which they are conscious occurrences.

In a series of papers Tyler Burge develops what looks like a promising alternative to the transparency account of self-knowledge. Imagine I am asked what I am thinking and that I respond by judging: I am thinking that writing requires concentration. This judgement is self-verifying in this sense: judging that I am thinking that writing requires concentration makes it true that I am thinking that writing requires concentration. It makes it true because 'the cognitive content that I am making a judgement about is self-referentially fixed by the judgement itself' (Burge 1998, p. 120). My judgement is based on nothing else and so is my knowledge of what I am thinking. In particular, I don't come to know what I am thinking by considering the reasons in favour of the proposition that writing requires concentration.

How does this approach deal with examples like the fearful flyer? Take the case in which the thought that the plane is going to crash occurs to me at t_1 . At that precise moment you happen to ask me

what I am thinking. A moment later, at t_2 , I tell you: I am thinking the plane is going to crash. What I have told you is true, but what makes it true? Burge's account implies that what makes it true is my judgement at t_2 . If I judge at t_2 that I am thinking that the plane is going to crash I thereby make it true at t_2 that I am thinking that the plane is going to crash. However, I do not thereby make it true that at t_1 I was thinking that the plane is going to crash; the judgement 'I have just been thinking that the plane is going to crash' is not self-verifying, unlike the judgement 'I am thinking that the plane is going to crash'.¹⁰

This is a problem for Burge, because when someone asks me at t_1 what I am thinking they presumably want to know what I am thinking *at that time*, or just before t_1 . That is why, in response to the question 'What are you thinking?', it is facetious to say 'I'm thinking about your question'. By the same token, what makes it appropriate for me to judge at t_2 that I am thinking that the plane is going to crash is the fact that I *was* thinking that the plane is going to crash when I was asked the question at t_1 . What makes my answer appropriate is not my judging at t_2 that I am thinking that the plane is going to crash. It is in this sense that I would be mistaken if at t_2 I were to judge 'I am thinking that the plane won't crash'. In the act of judging this I would be making it true that I am thinking that the plane won't crash but I wouldn't make it true that this is what I have been thinking. What I have just been thinking is that the plane *will* crash. Moreover, I know that this is what I have been thinking by being aware of my own thoughts. If this is right then the way to understand how we know our passing thoughts is not to focus on the fact that judgements like 'I am thinking that the plane will crash' are self-verifying but to try to understand the distinctive way in which we are aware of our passing thoughts.

A convincing account of the distinctive way in which we are aware of our passing thoughts will need to be grounded in a convincing account of what it is for one to have such thoughts in the first place. For example, in what sense do I find myself thinking that the plane is going to crash? What is it for this thought suddenly to occur to me? Perhaps I have mental images of the plane crashing, but it is not clear that this amounts to thinking, as distinct from merely imagining, that the plane is going to crash. But now suppose

¹⁰ As Boghossian points out in a helpful discussion of Burge's views. See Boghossian (1998).

that as the plane picks up speed down the runway I find myself *saying* to myself that the plane is going to crash. I might say this out loud, but consider the case in which the saying is a purely inner saying, a saying in inner speech or in what Ryle calls ‘silent soliloquy’.¹¹ Here are three suggestions:

- (a) Inwardly saying to myself that *P* is, or is closely related to, thinking that *P*.
- (b) In being aware of saying to myself that *P* I am aware of having the thought that *P*.
- (c) I know that I am having the thought that *P* by being aware of saying that *P*.

With regard to (a), the suggestion is not that all thinking is like this but that some is. If there is no such thing as inner speech then (a) is a non-starter. I take it as obvious that there is inner speech but will not attempt to analyse or explain this notion any further here.¹² One sense in which thoughts occur to one is that sentences in the language of inner speech occur to one; to find oneself saying to oneself that *P* is to find oneself thinking that *P*. This could either be because saying to oneself that *P* is or *constitutes* thinking that *P* or because, even though it isn’t the same as thinking that *P*, it *discloses* that this is what one is thinking. Either way, inwardly saying to oneself that *P* is closely related to thinking that *P*.

This leads to (b). If saying that *P* is, or is very closely related to, thinking that *P* then it is plausible that in being aware of saying to myself that *P* I am aware of having the thought that *P*. Indeed, part of what it is for thinking that *P* and inwardly saying that *P* to be closely related *is* for awareness of inwardly saying that *P* to amount to awareness of thinking that *P*. The hard question is: in what sense is one ‘aware’ of saying to oneself that *P* if the saying is a saying in

¹¹ Sellars is someone else who purports to ‘take very seriously the view that a thought, in the sense in which thoughts occur to one, is the occurrence in the mind of sentences in the language of “inner speech”’ (Sellars 1975, Lecture 1, #31).

¹² See Carruthers (2005) for further discussion. As he points out, ‘many of us are inclined to report, on introspective grounds, that at least some of our conscious propositional thinking is conducted in imaged natural-language sentences’ (Carruthers 2005, p. 166). In studies, all subjects report at least some instances of inner speech. So ‘the existence of inner speech, itself, isn’t—or shouldn’t be—in doubt’ (2005, p. 166). According to the cognitive conception of language which Carruthers wants to defend, inner speech is partly constitutive of thinking. Inner speech is not merely expressive of thought, or merely what gives us access to our thoughts.

inner speech? Auditory metaphors are virtually inescapable. The sense in which one is aware of inwardly saying to oneself that *P* is that one ‘hears’ oneself saying to oneself that *P*. This is hearing with the mind’s ear rather than with the ears attached to one’s skull.¹³ No doubt much more needs to be said about this form of awareness, but to deny its existence is to deny the existence of something which certainly seems phenomenologically real.

That leaves (c). If I am aware of saying to myself that *P*, and thereby aware of myself thinking that *P*, or having the thought that *P*, that enables me to know that I am thinking that *P*. This is what Ryle is getting at when he observes that our internal monologues ‘disclose’ our frames of mind, and that eavesdropping on our internal monologues enables us to describe our frames of mind. Is the resulting knowledge immediate? If immediate knowledge must not be based on observation or perception then one issue is whether ‘hearing’ oneself think is a form of perception. That depends on how generously the notion of perception is construed but it is not clear, in any case, why perceptual knowledge cannot be immediate.¹⁴ A far more pressing question is whether knowledge of what one is thinking is a form of inferential or evidence-based knowledge. Not if inwardly saying that *P* is thinking that *P*. Rather, direct awareness of one’s inner speech would amount to direct awareness of one’s thoughts. Even on the view that our internal monologues *disclose* our thoughts without *being* our thoughts there is still no reason to classify knowledge of one’s thoughts as mediated. To regard one’s inner speech as disclosing one’s thoughts is not to be committed to the view that one infers one’s thoughts from one’s inner speech, or that inwardly saying to oneself that *P* is merely a reliable sign or good evidence that one is having the thought that *P*. The connection between inner saying and thinking is more intimate than that.

The passivity of passing thoughts is now easily explained. Inwardly saying to oneself that the plane is going to crash need not be the product of reflection on the reasons in favour of the proposition that the plane is going to crash, and may well be entirely unresponsive to such reflection. We can distance ourselves from our passing

¹³ I owe the expression ‘the mind’s ear’ to Alex Byrne. See Byrne (forthcoming).

¹⁴ Indeed one might think that perceptual knowledge is the paradigm of immediate knowledge. There might be good reasons for denying that self-knowledge is perceptual but it is not the *immediacy* of self-knowledge that counts against its being perceptual unless one thinks that perceptual knowledge is itself inferential.

thoughts in the sense that we can distance ourselves from what we find ourselves saying. We may recognize that some of our inner utterances are ill-founded but their occurrence is still not in question. As for the idea that it is a datum of consciousness, in Hampshire's example, that I am thinking that the plane is going to crash, the sense in which this is so is that it is a datum of consciousness that I am saying to myself that the plane is going to crash. The latter is a datum of consciousness in the sense that I can either literally or metaphorically hear my own utterances.

We now have an account of the nature of some of our passing thoughts, the distinctive way in which we are aware of them, and how this awareness grounds our knowledge of our passing thoughts. This account respects the passivity of passing thoughts and their distinctive phenomenology. It explains how immediate knowledge of our passing thoughts is possible and does so without appealing to the notion of transparency. The suggestion is not just that immediate knowledge of our passing thoughts *can* be explained without appealing to the transparency procedure but that immediate knowledge of our passing thoughts *cannot* be explained by the transparency procedure.

III

The next question is: is TEI correct? Remember that this thesis is limited in scope. It is only concerned with immediate knowledge of our own beliefs and other standing attitudes so the fact that immediate knowledge of our passing thoughts can't be explained by reference to transparency is neither here nor there as far as TEI is concerned. Having said that, it is not clear that TEI is successful even on its own terms. There are reasons for thinking that the transparency procedure does not explain the supposed immediacy with which we know our own beliefs. Furthermore, despite all the differences between having the passing thought that *P* and believing that *P*, these reasons are not unrelated to the reasons for thinking that the transparency procedure can't explain how we can have immediate knowledge of our passing thoughts.

We can clarify these concerns about TEI by focusing on two specific problems facing the transparency explanation of immediate self-knowledge. I will refer to the first of these as the *two questions*

problem and to the second as the *sticking problem*. The former is one to which Moran himself draws attention. The challenge is to understand how the question ‘Do you believe that *P*?’ can be transparent to the question ‘Is it the case that *P*?’ The first of these questions is *inward-directed* while the second is *outward-directed*. These questions have quite different subject-matters, and it is conceivable both that I believe that *P* when *P* is false, and that I don’t believe that *P* when *P* is true. In that case, as Moran asks, ‘what right have I to think that my reflection on the reasons in favour of *P* (which is one subject-matter) has anything to do with the question of what my actual *belief* about *P* is (which is a quite different subject-matter)?’ (Moran 2003, p. 405).

Moran gives his response to the two questions problem in this passage, in which he is discussing the relationship between the questions ‘Is it raining?’ and ‘Do you believe that it is raining?’:

And then my thought at this point is: I *would* have a right to assume that my reflection on the reasons in favour of rain provided an answer to the question of what my belief about rain is, if I could assume that *what* my belief here is was something determined by the conclusion of my reflection on those reasons. An assumption of this sort would provide just the right link between the two questions. And now, let’s ask, *don’t* I make just this assumption, whenever I’m in the process of thinking my way to a conclusion about some subject-matter? (Moran 2003, p. 405)

The conclusion of my reflection on the reasons in favour of *P* is a judgement. Assuming that my belief concerning *P* is determined by the conclusion of my reflection on the reasons in favour of *P* is therefore equivalent to assuming that my belief concerning *P* is determined by whether I judge that *P*. Call this the *linking assumption* (LA). Moran represent the linking assumption as one which I actually make, and am entitled to make, when I am in the process of thinking my way to some conclusion. The specific role of this assumption is to connect the inward-directed and the outward-directed questions in such a way as make it intelligible that I am entitled to answer the former by answering the latter.

In so far as I rely on LA in coming to know that I believe that *P* is my knowledge of what I believe still immediate? This divides into two questions, one concerning psychological immediacy and the other concerning epistemic immediacy. Starting with the former,

there is a straightforward reason for thinking that the transparency doesn't give us psychologically immediate self-knowledge: to come to know whether I believe that *P* by consideration of the reasons in favour of *P* itself is to come to know whether I know that *P* by *reasoning*. Since the reasoning is conscious it follows immediately that the resulting self-knowledge is not immediate in the psychological sense.¹⁵ Another way of making the point would be to emphasize that, on Moran's view, self-knowledge is arrived at by deliberation. But to deliberate is to reason, and psychologically immediate knowledge is knowledge that is *not* arrived at by conscious reasoning or inference.

As for whether transparency delivers epistemically immediate knowledge we need to go back to the definition of epistemic immediacy. A person's knowledge that *P* is immediate in this sense only if his justification for believing that *P* does not come, even in part, from his having justification to believe other, supporting propositions. When his justification comes from his having justification to believe other, supporting propositions then his knowledge is inferential. Now consider the following simple argument: I come to know that I believe that *P* by following the transparency method only if doing so gives me a justification for believing that I believe that *P*. But when I rely on transparency in coming to know what I believe my justification for believing that I believe that *P* comes in part from my having justification to believe at least one other proposition, namely, the linking assumption. So my knowledge that I believe that *P* is epistemically inferential, and therefore not immediate, when it is based on the transparency procedure.

There are several ways of trying to block this argument, none of them convincing. For example, there is the thought that for the purposes of connecting up the inward-directed and outward-directed questions LA isn't an assumption that the self-knower actually has to make; it is enough that it is an assumption to which, as a rational being, he is *entitled* to make. But if LA is not an assumption I actually make, at least implicitly, then how am I supposed to find it intelligible that I can answer the question whether I believe that *P* by considering the reasons in favour of *P* itself? It is because the in-

¹⁵ Dorit Bar-On argues that the transparency method is 'epistemically rather indirect' to the extent that it implies that self-judgements 'are arrived at on the basis of consideration of worldly items' (Bar-On 2004, p. 113). In my terms, the indirectness implied by this characterization of the transparency method is primarily psychological.

ward- and outward-directed questions have different subject-matters that something is needed to link them, and the something had better not be beyond my ken if it is to make it intelligible to *me*, the knower, and not just the theorist of self-knowledge, that I can answer the inward-directed question by answering the outward-directed question. So it won't do to represent LA as an assumption that somehow links the two questions without at the same time mediating my justification for believing that I believe that *P*.

Another move might be to argue along the following lines: all knowledge has general background assumptions or presuppositions but this is not to say that all knowledge is, in any interesting sense, inferential.¹⁶ It all depends on the precise nature and role of the background assumptions. Even non-inferential knowledge can have presuppositions if they are what one might call *structural presuppositions*. Such presuppositions are, for example, presuppositions of thought or rational deliberation as such whose role is to *enable* knowledge or reasoning without also serving as mediating *premisses* that contribute to the justification of one's beliefs. So if, on the transparency account, LA is only a structural presupposition of self-knowledge it doesn't follow that this account makes self-knowledge inferential in any interesting sense.

Is LA a presupposition of rational deliberation as such? It is one thing to agree that if I am rational then whether or not I believe that *P* will generally be determined by my reflection on the reasons in favour of *P*. It is another matter entirely whether, as a rational being, I must *believe* that this is how my beliefs are determined. It is questionable whether a commitment to LA in this sense is a condition of rational deliberation and yet this is what the transparency account seems to require. In addition, it is very hard not to read the transparency account as implying that LA is a mediating premiss that contributes to the justification of one's beliefs about one's beliefs. At any rate, if LA is not a mediating premiss then, as argued above, it is obscure how it is supposed to make it intelligible to one that it is possible to answer an inward-directed question by answering a corresponding outward-directed question.¹⁷ So there remains the strong suspicion that the self-knowledge that the transparency procedure makes available is not epistemically immediate.

¹⁶ Thanks to Crispin Wright for suggesting this response.

¹⁷ Sydney Shoemaker also makes this point. See Shoemaker (2003, p. 401).

This suspicion is confirmed by the sticking problem. The following example from Peacocke illustrates the problem:

Someone may judge that undergraduate degrees from countries other than their own are of an equal standard to her own, and excellent reasons may be operative in her assertions to that effect. All the same, it may be quite clear, in decisions she makes on hiring, or in making recommendations, that she does not really have this belief at all. (Peacocke 1998, p. 90)

The point is this: the conclusion of my reflection on the reasons in favour of *P* is a judgement but concluding or judging that *P* is not the same as believing that *P*. Normally when I judge that *P* I also believe that *P* but I can judge that *P* without believing or coming to believe that *P*. This is what happens in examples such as Peacocke's. In such cases the belief that *P* fails to stick despite the acknowledged presence of good reasons in favour of *P*, reasons which lead one to judge that *P*.¹⁸ The reverse of this is also possible: the belief that *P* sticks or perseveres despite one's unwillingness to judge that *P* in view of the weakness of the reasons in favour of *P*. Notice that such cases are also a problem for the linking assumption since they suggest a weaker link between what one judges and what one believes than that assumption implies.

If I judge that *P* how can it be true that I nevertheless fail to believe that *P*? Suppose that *P* is the proposition that undergraduate degrees from countries other than my own are of an equal standard to my own. Should we not then insist, in response to Peacocke's example, that I either fail to judge that *P* or that I do in fact believe that *P*, at least at the point at which I make the judgement?¹⁹ To see why Peacocke's example cannot plausibly be dealt with in this way

¹⁸ As Shah and Velleman (2005, p. 507) observe, 'arriving at the judgement that *P* doesn't necessarily settle the question whether one now believes it, since one may find oneself as yet unconvinced by one's own judgement'.

¹⁹ This will sometimes be the right response to putative sticking scenarios. Shah and Velleman give the example of someone reasoning his way to the conclusion that his plane is not going to crash but still believing that it will. In this example it might be plausible to insist that the person has the obsessive thought that the plane will crash, or fears that it will crash, but does not literally believe that it will crash. Peacocke's example is different. The thought here is that what I really believe comes out in what I do. Judging is a kind of doing. If I judge that all undergraduate degrees are equal that is evidence that I believe that all undergraduate degrees but this evidence is trumped by the fact that my letters of recommendation and hiring decisions only make sense on the assumption that I don't believe that all undergraduate degrees are equal. Unlike the plane crash example, this is not a case of fear or obsession.

it is important to be clear about the relationship between believing that *P* and concluding or judging that *P* in response to consideration of the reasons in favour of *P*. To judge that *P* is to affirm that *P* with the aim of getting the truth-value of *P* right.²⁰ Judging that *P* should lead one to believe that *P* but is not guaranteed to do so because belief can be influenced by evidentially irrelevant factors.²¹ In Peacocke's example, the influence of prejudice prevents the judgement that undergraduate degrees from countries other than my own are of an equal standard to my own from leading, as it should, to formation of the belief that undergraduate degrees from countries other than my own are of an equal standard to my own.

The possibility of belief and judgement coming apart in this way should come as no surprise given that judging and believing belong in different ontological categories. Judging is a mental action whereas belief is a mental state. The stative character of belief is marked by the impropriety of progressive tenses: 'I am believing that *P*' is deviant in a way that 'I am judging that *P*' is not.²² When I judge that *P* I do not occurrently believe that *P* because there is no such thing as occurrently believing.²³ Belief is a dispositional state that is regulated for truth. Specifically, the belief that *P* 'tends to be formed in response to evidence of *P*'s truth, to be reinforced by additional evidence of it, and to be extinguished by evidence against it' (Shah and Velleman 2005, p. 500). What happens when I judge that *P* but don't believe that *P* is that the act of judging that *P* does not result in the formation or acquisition of the appropriate dispositions. The point of judging is to make a mark on one's beliefs but sometimes judging that *P* doesn't make a mark on one's beliefs; it doesn't result in one's believing that *P*. Just because I judge that *P* in response to what I recognize as good evidence for *P* that doesn't guarantee that I will end up in a mental state that is regulated for truth in the way that belief is regulated for truth.

To admit that belief can be influenced by evidentially irrelevant

²⁰ My account of judgement is essentially the one given in Shah and Velleman (2005). For further discussion of the relationship between judging, believing and thinking see Cassam (2010).

²¹ As Shah and Velleman (2005, p. 500) point out. They include phobias and wishful thinking among the evidentially irrelevant factors that can influence belief.

²² Williamson uses this argument to show that 'know' and 'believe' denote states rather than processes. See Williamson (2000, p. 35).

²³ Tim Crane also makes this point. He points out that 'occurrent belief' is a myth. See Crane (2001, p. 108).

factors is to accept that beliefs are not necessarily responsive to reason. This is one respect in which beliefs, like passing thoughts, can be passive. Just as one can *find* oneself having the thought that *P* so one can *find* oneself believing that *P*. And just as one can distance oneself from one's passing thoughts one can distance oneself from one's beliefs.²⁴ Martin gives a nice example of this. A father comes to realize, when engaged in second-order enquiry, that he has the conviction that his son is a great painter. Yet he may feel 'forced to distance himself from what he recognizes is one of his own strongly held convictions' since 'he believes something on insufficient evidence, even by his own lights' (Martin 1998, p. 115). To distance oneself from the conviction that *P* is not necessarily to cease to believe that *P*. Distancing might lead to the abandonment of the belief but then again it might not. Some beliefs may be able to survive the recognition that they are not rationally grounded.²⁵

All of this strengthens the case for thinking that the transparency procedure can only deliver self-knowledge that is not immediate. To see how, consider the following response to examples such as Peacocke's: what such examples show is that a person's belief concerning *P* is not *invariably* determined by his consideration of the reasons for or against believing that *P*. They do not show that they are not *normally* so determined. For part of what it is to be a rational agent 'is to be able to subject one's attitudes to review in a way that makes a difference to what one's attitude is' (Moran 2001, p. 64). The goal of deliberation is conviction, and it would not be possible to regard a person as deliberating if he does not presume, at least implicitly, that when he reasons his way to the conclusion that *P* he also ends up believing that *P* as a result. The fact that this is a defeasible presumption does not make it any less dispensable.

This strengthens the case for thinking that the transparency procedure can only deliver self-knowledge that is not immediate because it implies that that a person's judging that *P* is no more than defeasible *evidence* that he believes that *P*. For even if my concluding or judging that *P* in response to consideration of the reasons in favour of *P* is presumed to result in my believing that *P*, the fact remains that my judging that *P* does not *entail* that I believe that *P*.

²⁴ So Hampshire is mistaken when he says that 'a belief is a thought from which a man cannot dissociate himself' (Hampshire 1965, p. 98).

²⁵ Perhaps religious beliefs are a case in point.

What it does is to raise the probability that I believe that *P*. My judging that *P* makes it likely, given that I am rational, that I believe that *P* and is, in this sense, a reliable indicator that I believe that *P*. But this is just what it is for one thing to be evidence for another.²⁶ Furthermore, it is not just that my judging that *P* is evidence that I believe that *P*. It is also evidence I *have*, to the extent that I am aware of judging that *P* and understand that if I judge that *P* then it is highly likely that I believe that *P*. Finally, my knowledge that I believe that *P* is *based* on evidence in my possession if I know I believe that *P* *because* I know that I judge that *P*. But if I know that I believe that *P* on the basis of evidence then by knowledge is, by definition, not epistemically immediate.

IV

The story so far is that the transparency procedure fails to secure either the epistemic or the psychological immediacy of self-knowledge. So what? That depends on whether the immediacy thesis is correct. If it is a given that self-knowledge is immediate, and the aim is to account for its immediacy, then the transparency account fails. But should we endorse the immediacy thesis? It looks like a datum that knowledge of our own beliefs is, by and large, psychologically immediate so it is certainly an objection to the transparency approach that it fails to account for, or even accommodate, this datum. But epistemic immediacy is another matter. It is not obvious that knowledge of our own beliefs is typically immediate in this sense, and it is only taken as obvious because epistemic immediacy is confused with psychological immediacy. If one is sceptical about the assumption that knowledge of our own beliefs is normally epistemically immediate then it is not necessarily an objection to a theory of self-knowledge that it implies that we do not have epistemically immediate access to our own beliefs. The only genuine epistemological datum in this area is that self-knowledge is not normally based on *behavioural* evidence, and this is something that the

²⁶ As Williamson remarks, what is required for *e* to be evidence for the hypothesis *h* is that '*e* should speak in favour of *h*' and should itself have 'some kind of creditable standing' (Williamson 2000, p. 186). In probabilistic terms, *e* speaks in favour of *h* if it raises the probability of *h*. Kelly points out that 'the notion of evidence is that of something which serves as a reliable sign, symptom, or mark of that which it is evidence of' (Kelly 2006).

transparency account certainly does account for.²⁷ My judging that *P* may be evidence that I believe that *P* but it is not behavioural evidence. By the same token, if I know that I believe that *P* on the basis of my awareness of judging or mentally affirming that *P* this is not self-knowledge based on behavioural evidence.

All of this points to the need for a distinction between different kinds of epistemic immediacy and correspondingly different versions of the immediacy thesis. My knowledge that I believe that *P* is weakly epistemically immediate just if it is not based on behavioural evidence. It is strongly epistemically immediate just if it is based on no evidence. The sense that some philosophers have that self-knowledge is epistemologically distinctive is the sense that it is strongly epistemically immediate, and this is what gives rise to the idea that a good theory of self-knowledge should be able to explain how strongly epistemically immediate self-knowledge is possible. Viewed in this light, there are two natural responses to the fact that the transparency approach doesn't succeed in this task. One would be to look for an alternative account of self-knowledge which does explain the strong epistemic immediacy of some self-knowledge. The other would be to question the assumption that knowledge of our own beliefs is immediate in this sense. The suggestion is not just that it is not obviously correct that knowledge of our own beliefs is not strongly epistemically immediate but that reflection on the nature of belief should lead one to conclude that this kind of self-knowledge cannot be immediate in this sense.

The most promising version of the first of these responses is a Monitoring Mechanism (MM) approach to self-knowledge. The idea is this: when I come to believe that *P* what happens is that the representation that *P* enters my Belief Box. The question 'Do I believe that *P*?' calls for a search of my Belief Box. This search is not carried out by me, the subject, but by one of my sub-personal monitoring mechanisms, that is, 'a distinct mechanism that is specialized for detecting one's own mental states' (Nichols and Stich 2003, p. 163). If the belief that *P* is found in my Belief Box this leads to the formation of the second-order belief that I believe that *P*, and this second-order belief constitutes knowledge as long as the monitoring mechanism is

²⁷ Moran sometimes gives the impression that all he means when he says that self-knowledge is immediate is that it is not based on behavioural evidence. But if self-knowledge is based on evidence other than behavioural evidence then it is certainly not 'radically nonevidential' (Moran 2001, p. 68).

reliable and produces beliefs about my beliefs which could not easily be false. This knowledge is both psychologically and epistemically immediate. It is not produced by conscious reasoning and is not based on behavioural or any other evidence.²⁸

This account of self-knowledge has its attractions but presupposes an externalism about knowledge which some may find hard to swallow. Others may have difficulty with the idea that beliefs are the kind of thing that can be stored and monitored in the way the MM theory implies. While these and other problems with the MM theory might not be insuperable they do raise fundamental questions about the nature of belief and the assumption that we can have epistemically immediate knowledge of our own beliefs. Suppose we think of belief as a form of acceptance, so that to believe that *P* is to accept that *P* or, in other words, to regard *P* as true. The challenge is to say what distinguishes believing from other modes of acceptance such as supposing and assuming, and it is in the course of responding to this challenge that the limitations of the immediacy thesis become apparent. For when we think about the various respects in which belief is different from other modes of acceptance it becomes hard to see how knowledge of one's own beliefs could be strongly epistemically immediate.

Suppose that what distinguishes belief from other modes of acceptance is that there is a distinctive way in which beliefs are regulated, that is, formed, revised and extinguished. Belief is regulated for truth in a way that other modes of acceptance are not, and being regulated for truth is a broadly dispositional property of beliefs. Whether a given mental state is the state of believing that *P* is therefore partly a matter of what dispositions the state has, and one epistemological consequence of this is that one is not always in a

²⁸ As Stich and Nichols point out, a good theory of self-awareness needs to be able to explain the fact that 'when normal adults believe that *P*, they can quickly and accurately form the belief *I believe that P*' (Nichols and Stich 2003, p. 160). In order to implement this ability, 'all that is required is that there be a Monitoring Mechanism (MM) that, when activated, takes the representation *P* in the Belief Box as input and the representation *I believe that P* as output' (2003, pp. 160–1). The Monitoring Mechanism simply has to copy representations from the Belief Box and embed copies of them in a schema of the form '*I believe that ...*'. Stich and Nichols do not draw attention to the consequences of their view for the issue of immediacy, but it seems obvious that if a Monitoring Mechanism is sufficiently reliable to produce knowledge of one's own beliefs then the knowledge to which it gives rise is both psychologically and epistemically immediate. It was Timothy Williamson who first drew my attention to the possibility of exploiting something like the MM theory to explain the immediacy of self-knowledge. He does not, however, endorse the present approach to self-knowledge.

position to know whether one believes that *P* since one is not always in a position to know that one is in a mental state with the relevant dispositions.²⁹ It is also unclear, on the present account, how one's knowledge that one believes that *P* could be immediate. I only believe that *P* if I am in a mental state which is regulated for truth but how can I know, other than on the basis of evidence, that I am in such a state? How, for example, can I know without evidence, not only that I accept that *P* but also that my acceptance of *P* is such that it would be extinguished by evidence against the truth of *P*? In fact, matters are even more complicated than this. Even if I find myself continuing to accept that *P* in the face of what I recognize as evidence against *P* it still does not follow that I don't really believe that *P*. Genuine beliefs can be influenced by non-rational factors, and this makes them even harder to detect.

The true epistemological significance of the distinction between occurrent states and standing attitudes is now apparent. The emerging picture is that occurrent mental states and mental actions are on the surface of the mind, and that is why they can be known immediately by their subject. I do not need to dig deep in order to know what I am thinking or judging or feeling at the time when I am thinking or judging or feeling. But standing attitudes are not surface phenomena. Beliefs, for example, are psychological states that in some sense underlie one's occurrent mental states. Judging that *P* might manifest the belief that *P* but it is natural to think that the belief is not as directly accessible as its manifestations. While the distinction between what is on the surface of our minds and what lies beneath is obviously metaphorical it is an important element of the naive conception of the mental. Perhaps, on reflection, even what is below the surface is immediately knowable, but it is not obvious that the presumption is one to which we are committed. Once we have dispensed with implausible, philosophically inspired versions of the immediacy thesis, and also taken on board the extent to

²⁹ In Williamson's terminology (which is different from Moran's) 'transparency' is the thesis that 'for every mental state *S*, whenever one is suitably alert and conceptually sophisticated, one is in a position to know whether one is in *S*' (Williamson 2000, p. 24). He goes on to argue that transparency fails for the state of believing since 'the difference between believing *P* and merely fancying *P* depends in part on one's dispositions to practical reasoning and action manifested only in counterfactual circumstances' (2000, p. 24). In effect, Williamson's point is that the dispositional dimension of believing makes trouble for what he calls transparency. My point is that it makes trouble for immediacy.

which belief-formation is not a rational process, we can start to think realistically about how we are able to know our own beliefs.³⁰

Please provide affiliation and postal (& e-mail if you wish) address

REFERENCES

- Bar-On, Dorit 2004: *Speaking My Mind: Expression and Self-Knowledge*. Oxford: Clarendon Press.
- Boghossian, Paul 1998: 'Content and Self-Knowledge'. In Peter Ludlow and Norah Martin (eds.), *Externalism and Self-Knowledge*, pp. 149–73. Stanford, CA: CSLI Publications. Originally published in *Philosophical Topics*, 17 (1989), pp. 5–26.
- Bollas, Christopher 2009: *The Evocative Object World*. London: Routledge.
- Burge, Tyler 1998: 'Individualism and Self-Knowledge'. In Peter Ludlow and Norah Martin (eds.), *Externalism and Self-Knowledge*, pp. 111–27. Stanford, CA: CSLI Publications. Originally published in *The Journal of Philosophy*, 85 (1988), pp. 649–63.
- Byrne, Alex forthcoming: 'Knowing That I am Thinking'. In Anthony Hatzimoyisis (ed.), *Self-Knowledge*. Oxford: Oxford University Press.
- Carruthers, Peter 2005: 'Conscious Thinking: Language or Elimination?'. In his *Consciousness: Essays from a Higher-Order Perspective*, pp. 115–33. Oxford: Oxford University Press. Originally published in *Mind and Language*, 13 (1998), pp. 323–42.
- Cassam, Quassim 2010: 'Judging, Believing and Thinking'. *Philosophical Issues*, 20, pp. 80–95.
- Crane, Tim 2001: *Elements of Mind: An Introduction to the Philosophy of Mind*. Oxford: Oxford University Press.
- Davidson, Donald 1994: 'Knowing One's Own Mind'. In Quassim Cassam (ed.), *Self-Knowledge*, pp. 43–64. Oxford: Oxford University Press. Originally published in *Proceedings and Addresses of the American Philosophical Association*, 60 (1987), pp. 441–58.
- Frankfurt, Harry G. 1998: 'Identification and Externality'. In his *The Importance of What We Care About: Philosophical Essays*, pp. 58–68. Cambridge: Cambridge University Press. Originally published in Amelie

³⁰ I gave earlier versions of this paper at the University of Bonn, the University of Chicago and at the 2010 Oslo conference on Self-Knowledge and Rational Agency. I thank Frank Barel, who was my commentator in Oslo. For many other helpful comments and questions I thank Elke Brendel, Jason Bridges, James Conant, Fred Dretske, Ciara Fairley, David Finckelstein, Adrian Haddock, Jonathan Lear, Anna-Sara Malmgren, Conor McHugh, Richard Moran, Anders Nes, Julia Peters, Robert Stalnaker, Josef Stern and Crispin Wright.

- Rorty (ed.), *The Identities of Persons*, pp. 239–51. Berkeley, CA: University of California Press, 1977.
- Hampshire, Stuart 1965: *Freedom of the Individual*. London: Chatto and Windus.
- Kelly, Thomas 2006: ‘Evidence’. In Edward N. Zalta (ed.), *The Stanford Encyclopedia of Philosophy*. <<http://plato.stanford.edu/entries/evidence/>>
- Martin, M. G. F. 1998: ‘An Eye Directed Outward’. In Crispin Wright, Barry C. Smith and Cynthia Macdonald (eds.), *Knowing Our Own Minds*. Oxford: Clarendon Press.
- Moran, Richard 2001: *Authority and Estrangement: An Essay on Self-Knowledge*. (Princeton and Oxford: Princeton University Press).
- 2003: ‘Responses to O’Brien and Shoemaker’. *European Journal of Philosophy*, 11, pp. 402–19.
- 2004: ‘Replies to Heal, Reginster, Wilson, and Lear’. *Philosophy and Phenomenological Research*, 69, pp. 455–72.
- Nichols, Shaun, and Stephen P. Stich 2003: *Mindreading: An Integrated Account of Pretence, Self-Awareness, and Understanding Other Minds*. Oxford: Oxford University Press.
- Pryor, James 2005: ‘There is Immediate Justification’. In Matthias Steup and Ernest Sosa (eds.), *Contemporary Debates in Epistemology*, pp. 181–202. Oxford: Blackwell Publishing.
- Ryle, Gilbert 1949: *The Concept of Mind*. London: Hutchinson.
- Sellars, Wilfrid 1975: ‘The Structure of Knowledge’. In Hector-Neri Castañeda (ed.), *Action, Knowledge and Reality: Studies in Honor of Wilfrid Sellars*, pp. 295–347. Indianapolis, IN: Bobbs-Merrill.
- Shah, Nishi, and J. David Velleman 2005: ‘Doxastic Deliberation’. *Philosophical Review*, 114, pp. 497–534.
- Shoemaker, Sydney 2003: ‘Moran on Self-Knowledge’. *European Journal of Philosophy*, 11, pp. 391–401.
- Williamson, Timothy 2000: *Knowledge and Its Limits*. Oxford: Oxford University Press.

