

AUTHOR MEETS CRITICS

Précis of Kant on Self-Knowledge and Self-Formation and Replies to Critics

Katharina T. Kraus

University of Notre Dame, Notre Dame, IN, USA
Email: kkraus2@nd.edu

Keywords: self-knowledge; personhood; personal identity; embodiment; soul

I. Précis

First, I would like to thank Patrick Frierson, Janum Sethi, Clinton Tolley and Allen Wood for their engaging, insightful and thought-provoking comments that allow me to explore further details and objections to my interpretation.¹ I begin with a brief précis of the book.²

‘What is inner experience for Kant?’ is the central question I explore in my book *Kant on Self-Knowledge and Self-Formation* (Cambridge: Cambridge University Press, 2020; hereafter *KSS*). In exploring this question, the book offers a systematic study of the various ways in which human subjects can relate to themselves, access their own minds and ultimately gain empirical knowledge of themselves as the unique individual persons they are (or are in the process of becoming). This initial question splits into the book’s two guiding questions: first, what is inner experience about, i.e. what is its very object, and second, how is this object represented in inner experience, i.e. what is the mode or way of representation?

Regarding the object of inner experience, several candidates present themselves, such as the mental states currently passing through my mind (e.g. perceptions, thoughts, memories, imaginings, feelings, desires), or temporally more stable psychological properties (e.g. character traits, commitments and values), or simply I myself. But what kind of ‘object’ can I be for myself? A thinking subject (*denkendes Subjekt*), a mind (*Gemüth*), a soul (*Seele*), a collection of inner appearances (*innere Erscheinungen*), a person (*Person*) or an embodied human being (*Mensch*)? As a result, the book offers a reconstruction of Kant’s psychological account of a person, since persons, I argue, are the kind of beings referred to in inner experience.

Regarding the way of representation, what is at stake is whether I can cognize myself in an objectively valid way. Can inner experience ever amount to empirical cognition of an object (in the Kantian sense), as the outer experience of spatio-material objects does. Is inner experience valid in light of its distinctive ‘object’ and hence in principle valid for everyone and not only subjectively valid for the one having this experience? The first main thesis of my book is that inner experience is empirical cognition in a qualified sense: since it is cognition of myself as a

© The Author(s), 2022. Published by Cambridge University Press on behalf of Kantian Review. This is an Open Access article, distributed under the terms of the Creative Commons Attribution licence (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted re-use, distribution, and reproduction in any medium, provided the original work is properly cited.

psychological person, rather than as a mere object, I cognize myself only qua my psychological features, not as a persistent mental substance. Such psychological features are broadly construed to include occurrent mental states, temporally stable psychological properties and even autobiographical narratives.

My argument for this thesis is based on a detailed examination of the conditions that are required to make inner experience as empirical cognition possible, following the synthetic method of the *Critique of Pure Reason*. Each of the three parts of the book examines these conditions with regard to one of the three main faculties of cognition: sensibility (*Sinnlichkeit*), the understanding (*Verstand*) and the faculty of reason (*Vernunft*). Chapters 1 and 2 examine the conditions of sensibility, focusing on self-affection (i.e. affecting oneself in inner sense), and develop a theory of inner perception as the empirical consciousness of my mental states in time. Turning to the higher faculties, chapter 3 offers a theory of transcendental self-consciousness as the mere form of reflexive consciousness, and chapter 4 enquires into the understanding's conditions of referring to oneself in judgement. It turns out that there is a fundamental difference between self-reference (in I-judgements about myself) and object-reference (in judgements about ordinary spatio-material objects): not all categories of the understanding can be applied in I-judgements in the ordinary, constitutive way. A difficulty arises with respect to the categories of relation, since there is no substratum given in inner sense that can be cognized as a mental substance underlying all my mental states and endowed with mental powers. A constitutive use of the category of substance would thus lead to a dogmatic metaphysical foundationalism regarding souls with personal identity through time – a position Kant fiercely criticizes in the Paralogisms. Chapter 5 therefore appeals to additional conceptual resources provided by reason and concludes that we need an analogical application of the relational categories, on the basis of the idea of the soul, which – if used in a regulative way – substitutes precisely those schemata (i.e. schematized categories) that cannot be employed in inner experience (e.g. the schema of substance as persistence through time). More generally, the idea of the soul outlines the scope within which the principles of the understanding can be operative, though some of them only in a reduced form, to yield cognition of mental states. The idea thus defines the context of intelligibility within which inner experience can first be understood as truth-apt cognition of the psychological features that are assumed to belong to one and the same person. Having established the first thesis that inner experience is – in this qualified sense – empirical cognition of myself as a psychological person, chapters 6 and 7 turn to the normative demands that reason places on inner experience.

Chapter 6 offers a theory of self-knowledge, arguing that self-knowledge requires, in addition to self-cognition, a normative standard of epistemic justification by which we can assess whether a self-cognition is indeed true in light of the represented 'object', that is, in light of who I really am. The idea of the soul then also serves as a normative guideline for assessing my empirical self-knowledge claims by demanding that all the cognitions I hold to be true about myself must systematically cohere with each other.

Chapter 7 finally develops a theory of the psychological person as the 'object' referred to in inner experience. According to my second main thesis, persons form themselves in the course of realizing their mental capacities under the guidance of a unifying idea, the idea of the soul. As persons, we are empirical realities in the process of becoming: we are empirically real by virtue of the mental states that occur in our

empirical consciousness. Yet, at least during life, we require further self-realization according to an *a priori* form, the form defined by the idea of the soul. This process of individual self-formation is understood as an overarching mental activity carried out throughout a lifetime and constituted by the multiple first-order mental acts performed every day, such as acts of perceiving, judging and willing. The idea of the soul plays out as a set of norms that govern these first-order mental acts, so that the mental whole thereby realized approaches the rational ideal of systematic unity.

2. Reply to Janum Sethi

In her comments, Sethi raises three related issues concerning (1) my account of merely subjective validity of representations, (2) my account of a form–matter relationship between transcendental and empirical apperception and (3) the possibility of observing one’s own thinking as an inner appearance in inner sense. Sethi indicates an alternative view that she has explored elsewhere. Her view puts emphasis on the notion of association, which plays a central role in the philosophies of Kant’s empiricist predecessors such as Hume.³ In what follows, I discuss each issue in turn and explain why Kant’s theory of self-consciousness and inner experience should not be understood to start from an empiricist, or associationist, theory of the self.

2.1 Object-reference versus objective validity

A difference between Sethi’s and my account of subjective validity is that, in my view, a merely subjectively valid representation can still involve a relation to an object of some sort, whereas for her such a representation can only amount to the awareness of one’s own representations, not of an object. Kant’s notion of an object, however, is ambivalent. Basically, a representation represents something to someone; that is, it involves a relation to a subject in the generic sense (i.e. reflexivity) and has an object in the generic sense (i.e. referentiality) (KSS 19–21). An object (*Objekt*) in this broad sense can be anything that is representable, including an appearance, a number, an abstract entity, a mental state or a spatiotemporal body. An object in the narrow sense is an object of experience (*Erfahrungsgegenstand*) constituted by the categories of the understanding. The passage from the *Prolegomena* Sethi cites to support her view (*Prol*, 4: 300; cited p. 462) invokes, I think, the narrower notion of an object of experience, which is typically given independently of the particular subject.

The distinction between objective and merely subjective validity, in my view, concerns not what a representation is about, but the mode or way in which something is represented and hence the type of representation. A judgement, for example, is by definition an objectively valid way of representing something, whereas a perception or desire is merely subjectively valid, albeit about something. Subjective validity is primarily understood negatively against its counterpart of objective validity. If a representation is objectively valid, it has what Kant calls ‘necessary and universal validity’ (e.g. A48/B65, B140) and is therefore valid exclusively in light of its object. The representation is then valid for everyone, independently of the particular subject representing it. If a representation is merely subjectively valid, it lacks precisely this property because it in some sense depends on contingent features of the particular subject. For instance, my intuition of the chair in front of me is merely subjectively

valid insofar as it presents me with an image of the front side of the chair only depending on my current visual perspective. But, in another sense, this intuition can be used to acquire an objectively valid cognition of the chair by considering it as sensible matter and reflecting it under the categories of the understanding as the necessary and universal forms of cognition. By making the judgement ‘This is a chair’, I abstract away from my distinctive perspective and also judge the object to have a back side, even if this side is currently not visible to me. Kant’s examples [1] ‘When I carry a body, I feel the pressure of weight’ and [2] ‘It, the body, is heavy’ (B142) mark the same distinction: [1] expresses a merely subjectively valid perception of the heaviness of a body, whereas [2] expresses an objectively valid judgement about that same body. Both refer to the same spatiotemporal object, the body, but in different ways: [1] is based exclusively on the perceptual relation of the subject to the body, whereas [2] transforms this perceptual relation into an objective relation by reflection under the categories.⁴ Crucial textual evidence for my account can be found in the following passage:

[The] unification [of thinking] either arises merely relative to the subject and is contingent and subjective, or it occurs without condition and is necessary or objective. (*Prol*, 4: 304)

Given my distinction between object-reference, on the one hand, and objective validity, on the other hand, Sethi’s proposal of ‘a subject’s awareness of her own representations’ would still be about an object in the broad sense, namely, one’s own representations. This distinction, furthermore, allows me to distinguish between inner perception and inner experience: the former is a merely subjectively valid consciousness of my mental states in time (chapter 2) and the latter an empirical cognition of my mental states, which is necessarily objectively valid and must involve the relational categories (chapter 5).

My interpretation is compatible with Kant’s appeal to association: the way in which representations are combined in inner perception (e.g. the stringing together of visual impressions of a static object such as a house, or the association of a state of anxiety with the perception of a spider) may be due to the reproductive imagination, which functions primarily according to its intrinsic laws of association. Importantly, inner perception can amount to an inner experience only if it is subject to the categories. Only then can its content be understood as ‘united in the object’ (B142). One difference with Sethi’s associationist proposal, however, is that I argue in chapter 2 that perception must be normatively guided by the categories in order to find meaningful perceptual units within an empirical manifold that can then be conceptualized (KSS 79–80).

2.2 A hylomorphic account of apperception

Sethi objects to my claim that the empirical unity of apperception is a ‘concrete realization’ of the transcendental unity of apperception. My account of apperception is embedded in a general *hylomorphic* interpretation of mental faculties in Kant. Accordingly, we can view a mental faculty in two respects: as a transcendental faculty, if viewed regarding its *form*, which defines the kind of unity it brings about, and as an empirical faculty, if viewed regarding its *matter*, that is, regarding the mental contents

that are in fact unified in an act of this faculty and hence appear in empirical consciousness. We can thus view the faculty of apperception in two respects: as transcendental apperception, if viewed regarding the characteristic kind of unity it brings about, which I identify with the general form of reflexive consciousness; and as empirical apperception, if viewed regarding the mental contents actually unified in empirical consciousness. The former yields self-consciousness *qua form* and the latter *qua matter*.

I disagree with Sethi's two claims that (i) the empirical unity of apperception should be understood as a unity by association only, independently of the transcendental form of apperception, and that (ii) transcendental apperception should be identified with the synthesis according to the categories of the understanding. On Sethi's account, the former is merely subjectively valid and the latter is necessarily objectively valid.

By contrast, I argue in chapter 3 that transcendental apperception defines the most general form of reflexive consciousness *per se*, which is a condition of all more specific kinds of consciousness, including both the cognition of objects and empirical self-consciousness. The original-synthetic unity of apperception is necessary to ground the *reflexivity*, or *for-me-ness*, that any representation must have to be significant for the subject and hence to count as a mental state of which the subject can become conscious.

The advantage of my hylomorphic reading is that it offers an account of transcendental apperception that is not restricted to cognitive states only. Rather, it shows that any mental state or any associative empirical unity of representations, to be significant for the subject, must be subject to the kind of unity and reflexivity that transcendental apperception provides. Only then, I argue, can we explain why, for instance, states of feelings and desires can have significance for the subject's conscious inner life and can be potential objects of self-ascriptive judgements. This is crucial, if one is interested in a theory of self-cognition that explains not only the possibility of the self-ascription of thoughts, but also that of non-cognitive mental states, including emotions, imaginings and volitions.

If transcendental apperception is understood to provide only the most general form that any significant representation must display, then it requires a further specification of this general form into more specific forms to account for more specific kinds of representations. Hence, *one* way of realizing the transcendental unity of apperception consists in its realization in accordance with the categories of the understanding, resulting in objectively valid judgements about objects, i.e. cognition (these conditions are explicated in §18 of the Deduction, from which Sethi quotes). *Another* way of realizing this unity consists in the 'determination of inner sense' in accordance with the 'laws of association', resulting in the empirical unity of apperception (B139). This empirical unity, which I take to be 'the unity of representations of which I am *de facto* empirically conscious of *as my own*' (KSS 117), is primarily merely subjectively valid, since associations in inner sense are not necessarily in accord with the categories. Hence, the hylomorphic reading explains how the empirical unity of apperception is 'derived only from the [transcendental unity of apperception], under given conditions *in concreto*' (B140, see especially KSS 117–19), namely, under the conditions of time, and how this empirical unity can be reflected under the categories to amount to objectively valid self-cognition (see KSS chapters 4 and 5).

Following Sethi's associationist view of empirical apperception as different in kind, it seems unclear how the basic relation of all representations in an associative unity to the common subject can be explained if they are not already subject to the general form of reflexivity. Hence, Sethi's view remains vulnerable to Hume's challenge that there can be no necessary guarantee that all associated representations in the 'bundle of perceptions' really belong to one and the same subject. If we begin with a merely associative unity of representations that is not already in the form of transcendental apperception, then it is unclear which operation could ever guarantee the necessity of the common subject and the acquisition of this form for the possibility of the self-ascription of any mental state.

2.3 The appearance of thought

I discuss the Reflexion 'Is it an experience that we think?' in the context of my interpretation of self-affection (chapter 2, KSS 65), and of the Third Paralogism of personal identity (chapter 4, KSS 164–5). The latter discussion serves the purpose of sharpening the distinction between the 'formal condition of my thoughts' (A363) and the 'temporal conditions under which my thoughts appear to myself' (KSS 164). This distinction builds on my account of self-affection in chapter 2, which distinguishes between transcendental and empirical self-affection – the former is the *a priori* self-affection that the understanding has on inner sense qua the *form of thought*, and the latter is the empirical effect a particular thought-act has on inner sense qua *matter of consciousness* (see KSS, sections 2.4 and 2.5). Empirical consciousness arises through the self-affection that any kind of synthetic activity has on inner sense. Self-affection yields the 'immediate awareness that accompanies the synthesized representations' and gives the resulting states of consciousness 'a primitive temporal ordering of *one-after-the-other*' (KSS 61).

My reading of Kant's theory of self-affection then implies that, not only can we have an *a priori* consciousness of our thoughts (e.g. the *a priori* thought of a square in abstraction from any empirical conditions of consciousness), but our synthetic acts of thinking can also appear to us – as inner appearances – in empirical temporal consciousness. A thought can give rise to an inner appearance only if its mental activity (e.g. the synthesis of the understanding involved in the representing of the square) sensibly affects inner sense, which then results in an inner intuition in time (see especially KSS, sections 2.3.2 and 2.3.3). In this sense, thoughts such as my current thought of the square laptop in front of me appear in my empirical consciousness at a certain time, although we usually do not focus on being empirically aware of them. When we attend to the inner appearances of our thinking and reflect on them under the conditions of self-cognition, we gain inner experience of our thoughts *as they appear to us in time*.

3. Reply to Allen Wood

Wood's main concern is the role of human embodiment for a theory of self-cognition and personhood. His overarching question is whether self-cognition is possible as the cognition of a genuine object, such as a psychological person, or whether any proper object of cognition must be in space, so that the genuine object of self-cognition can only be the embodied human being. In what follows, I discuss the four issues Wood raises.

3.1 The interactive faculty model

Is inner intuition or inner sense by itself sufficient to make cognition of an object possible? My answer is no, because self-cognition, like any other cognition, requires the involvement of various mental faculties. Chapter 2 offers an 'interactive model of perception' according to which perceptions arise 'in a single interaction of outer affection, apprehension, and self-affection' (KSS 65). The first requires outer sense, the second a faculty for consciousness in general, and the third inner sense. For outer perception, the model suggests three constitutive aspects: sensation in outer sense, apprehension of the resulting outer intuition and an immediate accompanying awareness through inner sense. Hence, not even outer sense, or outer intuition for that matter, is sufficient for the perception and thus empirical cognition of an outer object, but it necessarily requires inner sense and inner intuitions.

Cognition of objects requires that the corresponding intuitions meet the conceptual conditions for object-reference, that is, they must be synthesized in accordance with the categories. For the inner case, however, a difficulty arises for the categories of relation due to a lack of an inner substratum, as I discuss further below and in my reply to Frierson. In consequence, we obtain a set of special conditions for self-reference in inner experience, which involves both the temporal conditions of empirical consciousness and the conceptual conditions of self-reference in I-judgements (see KSS section 4.4).

3.2 Psychological reality

But does this restriction make self-cognition less real, 'more phenomenal' or even 'merely apparent' compared to the cognition of outer objects, as Wood suggests? I do not think so. A goal of my interpretation is to preserve the insight that the inner goings-on of our minds are as real as the external motions of material bodies, and that this psychological reality can be grasped in objectively valid, truth-apt judgements. So, I maintain that Kant's distinction between mere illusion (*bloßer Schein*) and appearance (*Erscheinung*) can be upheld in the inner case, as it holds in the outer case (see *Prol*, 4: 314; A38/B55, B69–70, A293/B350). The moon is an appearance of my outer sense, but that I judge the size of the moon to be larger when the moon is closer to the earth is based on a perceptual illusion, the moon-illusion, which can lead to a false cognition about the moon (see A297/B354). There can be illusions that do not even involve an existing object at all, such as hallucinations. The same distinction applies between inner illusions and inner appearances. We can, for example, be self-blind with regard to the kind of mental state we are in and reflect it under a false or inadequate concept (see KSS section 6.2.4, pp. 243–8). We can falsely take a dream for a really lived experience; we can falsely take an imagining for the perception of a real object; we can falsely take the command that someone else has given us for our own desire; we can falsely take our sadness for a state of anger; a hope for a belief; and so on.

In §49 of the *Prolegomena*, Kant argues explicitly that I am 'by means of outer appearances, just as conscious of the reality of bodies as outer appearances in space, as I am, by means of inner experience, conscious of the existence of my soul in time' (4: 336). In both cases, we can succumb to illusions: I can have a 'dream' about material objects, and likewise, an 'imagination' (*Einbildung*) about my inner state (4: 337). This suggests that neither my mind nor my body is privileged for Kant; rather, they are

equally fundamental for understanding the empirical nature of human beings and therefore not reducible to one another.

3.3 *Body-consciousness and personal identity during dreams*

This brings me to Wood's concern that dreams can interrupt not only body-consciousness, but also consciousness of one's personal identity – an implication that would lead to worrisome consequence for my interpretation of Kant's account of personhood. In my discussion of body-consciousness, in chapter 4, I argue that consciousness of personal identity should not be understood as deriving from consciousness of my *persistent* material body. In the relevant passage, I oppose specifically a 'strong reading' according to which 'for each temporal series of mental states to be cognized, I must be immediately conscious of a directly correlated bodily substratum' (KSS 158–9). For this reading implies, in my view, a reductionist view of psychological reality, according to which mental states are ultimately determined as the non-material accidents of my body and as sensibly caused by bodily powers, rather than mental faculties.

The case of dreams shows why it is problematic to tie consciousness of personal identity closely, or indeed necessarily, to body-consciousness. Dreams can be understood as illusions that present us with a distorted view of reality, or as 'imaginary objects' that lack any foundation in reality, as Wood suggests. Dreams often incorporate people, situations and objects of the dreamer's real life, yet typically in bizarre, surreal and irrational ways. It is a common phenomenon that we represent ourselves from a third-person perspective in dreams, rather than as first-person narrators or the main characters of the dream story.

Dreams normally occur during our deep-sleep phase, the so-called REM phase, which is characterized by a particularly low level of awareness. Therefore, dreams are certainly not conscious experiences (in a Kantian sense), but it is also not impossible to recall them under any circumstances upon awakening. Despite their illusory character with regard to the external world, dreams can be seen as a psychological reality that can itself become the subject matter of inner experience. In his theory of psychoanalysis, Sigmund Freud understood dreams as the symbolic expression of frustrated desires that have been relegated to the unconscious mind (see Freud 1900/2010). The analysis of recalled dreams can therefore help to uncover these repressed desires and generate genuine insight into the dreamer's psyche. Yet the possibility of a dream analysis presupposes that the dreamer is aware of her personal identity throughout the series of dream states and her current memories thereof. In my view, this requires that even dream states realize in a rudimentary sense the form of reflexive consciousness, i.e. the transcendental unity of apperception.

But it is near to impossible for a dreamer to determine a dream episode in objective time, i.e. in the time of physical objects. Such an objective temporal determination would require a correlation of the dream states with the states of one's body or other relevant physical objects. Since body-awareness and the awareness of the physical environment are typically so low during a dream phase, it is unlikely that such a correlation can be performed during sleep. Only in a sleep laboratory can brain activity, which is considered an indicator of a dream phase, be objectively measured with electrodes during sleep.

3.4 Human embodiment

That I take the conditions of personal identity in time to be definable independently of the conditions of material persistence seems to worry Wood, since for him the only candidate for a persistent object of self-cognition can be the embodied human being. Of course, I do *not* deny the fact that human persons are embodied. *Nor* do I deny that mental states are often intimately connected with bodily sensations. For example, emotional states such as fear, anxiety, anger or joyful anticipation are often associated with sensations in certain parts of the body. While I criticize interpretations that reduce the conditions of personal identity to bodily persistence, as the strong reading does, or derive them at least in part from bodily persistence, as does a view that combines features of the logical I with the materiality of the body, it is true that I have not provided a positive account of the role of human embodiment for inner experience and personhood.

Indeed, the body is mentioned in the guiding-thread passage that is central to my view: there Kant demands that we consider ourselves as a substance 'to which the states of the body belong only as external conditions' (A672/B700). Assuming that the bodily states of a person 'belong' to the soul as its 'external conditions' implies to me that they are not intrinsic or essential conditions of what a soul or a person as such is. And yet bodily states as 'external conditions' play a crucial role in determining mental states in time. Note that, by contrast, the reductionist views I sketched above would characterize mental states as belonging to the body, as epiphenomenal or external conditions that are not intrinsic to bodies themselves.

Hence, Kant's idea of the soul contains a body-related predicate, the predicate of 'standing in community with other real things outside [the soul]' (A682/B710). Kant then translates this predicate into the regulative principle that

all change [be considered] as belonging to the states of one and the same persisting being, and by representing all *appearances* in space as entirely distinct from the actions of *thinking*. (A682/B710)

This predicate wards off psychological materialism by requiring that inner appearances be represented as separate from outer appearances and therefore as irreducible to the 'laws of corporeal appearances' (A683/B711). But it also opens a way to define a regulative principle for the temporal determination of mental states in *objective time* – on the model of the Third Analogy regarding interdependence and simultaneity and based on the category of community. Here I see how human embodiment becomes crucial for inner experience.

In accordance with the three Analogies, I distinguish three aspects of determining the temporal relations of mental states, corresponding to three predicates of the soul (see KSS 187 and 209–15):

- (i) the determination of mental states as belonging to *one and the same mental whole*, i.e. the same psychological person, according to the predicate of <mental substance as persistence>;
- (ii) the determination of a *qualitatively ordered causal series* of mental states, according to the predicate of <causality as the law of a temporal sequence of mental states>;

- (iii) the determination of a *quantitatively ordered causal series* of mental states in objective time, according to the predicate of <community, as simultaneity of mental states with the external states of the body>.

While (i) and (ii) are in principle body-independent self-determinations, they do not provide a determination of mental states vis-à-vis material objects in space and therefore remain severely underdetermined: they give only *qualitative* temporal unities and causal orders in *subjective* time, i.e. the time of the subject's own consciousness, but not *quantifiable* relations in *objective* time. For example, (i) seems to be involved when I remember my night dream upon awakening and recognize it as belonging to me. And (ii) seems to be involved when I also remember, perhaps only vaguely, a qualitative order in which these dream states appeared to me, namely, insofar as I am able to tell a coherent dream story. However, only if (iii) is executed can mental states be successfully determined in an objective time-series vis-à-vis material objects. But (iii) requires a reciprocal relation between mental states and bodily states. The assumption of this reciprocity can only be based on a regulative idea, according to which mental states are considered as if they were in one and the same space as physical states, although – strictly speaking – mental states are only in time and have no spatial determinations, unless they are already correlated with a body.

If, as some argue, the three Analogies are not three independent acts of time-determination (of persistence, temporal ordering and simultaneity), but are in fact three necessarily correlated aspects of the same act of time-determination, then it follows that (i) and (ii) are not only incomplete in themselves, but even impossible without (iii). This thought then implies that human embodiment plays a necessary role for any temporal self-determination.

But these considerations might still not satisfy Wood's larger concern that we humans are embedded in an ecological, geographical, social and cultural environment and that our self-understanding depends primarily on 'our place in the world'. Again, I do not dispute that our bodily existence is necessary for Kant to mediate relationships with others in social contexts, as is particularly evident in his lectures on the pragmatic sciences of physical geography and anthropology. My interpretation makes explicit only the most basic necessary conditions of personhood. In chapter 7, I define a minimal set of normative conditions for the self-formation of a person, but I acknowledge that there are 'higher forms of self-realization within or through social communities, resulting in the formation of social identities, such as cultural, national, or religious identities, or of a moral community' (KSS 255n.). While I think that the formation of these social and moral identities is bound by the minimal normative conditions I defined, it remains underdetermined by these and requires further, higher principles of self-formation. I hope that the self-formation view I have developed can serve as a general framework for understanding higher, communal forms of self-formation, and that, as further conditions are recognized, it can explain how we can move from the psychological 'I' to the social, political and cultural 'we'.

In summary, my interpretation suggests that different aspects of human embodiment are significant at different levels of self-cognition and self-formation:

- The constitution of temporal consciousness and the acquisition of inner perceptions require the ability to represent objects in space, i.e. outer sense;

- The determination of mental states in objective time requires correlations with the external states of the human body, yet without reducing mental relations to bodily relations;
- Pragmatic orientation in the world, formation of social communities, formation of a moral character and realization of moral agency toward other persons require a general awareness of one's place in the world and must be mediated by human embodiment.

4. Reply to Patrick Frierson

While Frierson grants that we require a regulative use of the idea of the soul to have inner experience of the temporal succession of mental states, he raises a puzzle for my view. Since, according to the First Analogy, *all* change requires a substance that persists, he worries that mental change cannot be cognized by simply appealing to the *regulative* idea of a mental substance, and hence to 'powers [that] are *as if* of a substance' (p. 476). Rather, for Frierson, the cognition of mental change must involve a constitutive use of the category of substance, although, as he concedes, we do not have an 'empirically specified concept of substance' for the inner case (p. 478). Our dispute concerns Kant's conception of 'persistence' (*Beharrlichkeit*) as the temporal explication of the category of substance: For Frierson, persistence – as the schematized category of substance in general – applies to all appearances, regardless of their kind, including possibly purely mental substances, although he concedes that we can know only for corporeal substances *how* they persist through time, namely as matter movable in space. On my view, the task of a schema is precisely to explain *how* to apply a category and since we cannot know *how* the category of substance applies to inner appearances, we can use this category only regulatively (by means of an idea), but never constitutively (or what Frierson calls 'literally', p. 479). The only schema of persistence available to us, I submit, concerns the persistence of physical matter in space. In the inner case, we instead have to account for *personal identity* through time, which requires additional conceptual resources.

In my reply, I focus on Frierson's subtle distinction between what he calls 'the schematized category of substance in general' and 'empirical concepts of particular kinds of substances' (p. 477). First, I argue that the schematized category of substance implies the *knowing of how* the category applies to intuitions and hence to the corresponding appearances. Since we lack this *know-how*, as Frierson agrees, for the inner case, the schema of substance is not applicable in this case. Then second, I argue that we require the regulative idea of the soul to replace the inapplicable relational schemata of the understanding such that the principles of substantiality and causality can still be applied in a reduced form to *mental states*, though not to a mental substance.

4.1 The schema of a category implies knowing-how

The task of Kant's schematism is to 'show the possibility of applying *pure concepts of the understanding* to appearances in general' (B177/A138). A schema accomplishes this task by defining a rule according to which a manifold of intuitions can be synthesized

in such a way that the resulting unity is suitable for reflection under the corresponding concept of the understanding. Since, as Kant argues in the Deduction, the understanding impresses its categories primarily on inner sense by way of a figurative synthesis, a schema is understood as a *temporal* expression of a category, i.e. a time-determination, such as persistence as the schema of substance. Without the schematism we have no proof of the *real possibility* of applying a category to what is given in intuition. A schema first articulates an application rule for a given category (or group of categories) and thereby adds a proof of real possibility to the Deduction's general account of the figurative synthesis.

Frierson's idea of 'the schematized category of substance in general' (or 'the barely schematized category', pp. 477, 479) would then imply that we have a general proof of the possibility of applying the category of substance to appearances, regardless of their kind. Such a schema, Kant argues in the Schematism, requires the representation of a 'substratum' as that which persists (A144/B183). In the First Analogy, the synthetic principle of persistence, Kant does not mention specific kinds of substances. However, the only example of a successful application of this principle that Kant offers concerns physical matter, i.e. the matter intuited by outer sense (see A185/B228, also B278). Kant then states that the determination of change (*Wechsel*) in general must be modelled on the determination of the 'alteration' (*Veränderung*) of the states of material substances. Hence, the only model of something that serves as the 'substratum of . . . all time- determination', as required by the Schematism and the First Analogy (A144/B183 and A183/B226), is physical matter, because only outer sense provides suitable sensory matter for instantiating the category of substance. Outer sense supplies a 'spatial distribution of reality' and thus another dimension (in addition to time) along which we can distinguish what changes from what persists, and trace a persistent object through different places.⁵ In turn, Kant explicitly denies in several places that what is given in inner sense can serve as such a substratum (e.g. A107, 364). The reason is that inner sense supplies only a purely temporal 'distribution of reality', which by itself cannot be used to represent a substratum (see KSS 156–7).

Hence, even if we were able to discern a general sense of the schema of substance, as Frierson suggests, it would not resolve the problem of its applicability to inner appearances. We would still lack a proof that the schema of substance, even in its general or bare sense, can be applied to a manifold of inner intuition. In my view, such a proof would require that we show that inner sense can supply a substratum, and that the concept of *mental persistence* is not empty, but has sense and meaning. But that a concept has sense and meaning can only be shown if we can discern an application rule for it, which is precisely what is lacking in the inner case.

Therefore, I concluded that, *as far as we know*, the general schema of persistence can only be explicated as *material persistence in space*. Since we cannot arrive at a sensible explication of mental persistence, I maintain, we cannot use the schematized category of substance *constitutively*, but must apply a regulative substitute. I think that Kant's reflections in the *MFNS* that Frierson cites are an extension (and completion) of the First Analogy, and – granting the specific line of argument for physical matter as reconstructed by Frierson – they similarly show that there is no constitutive *real* application of the notion of mental persistence beyond the logical, unschematized concept of mental substances.

4.2 Reduced forms of the principles of persistence and causation

This result does not rule out that there could exist other kinds of substances, including mental substances, but it is strong evidence for the conclusion that we *cannot cognize* other kinds of substances on the basis of experience. The fact that we do not have a rule for representing a substratum in inner sense, however, does not mean that we have to give up on the relational categories for inner experience altogether. Rather, if we accept a regulative use of the idea of a mental substance, as I have interpreted it, then we can apply the First Analogy in a reduced form such that we can cognize *mental states* as inherent in a mental whole, without determining the whole as such. To restore this limited constitutive use of the relational principles with respect to mental states, we need a 'schema of reason' based on an idea that serves as the regulative substitute for the inapplicable schema of the understanding (A665/B693). In general, ideas of reason, as concepts of a whole (domain) of experience, are needed when a schema of the understanding reaches the bounds of sense. To prevent the understanding from overstepping this boundary, the regulative use of ideas defines a scope for the legitimate activities of the understanding.

The idea of the soul sets a scope for the proper use of the understanding in inner experience: it outlines the domain within which we can cognize mental states and mental changes. As a result, we obtain a limited, or reduced, form of the First Analogy based on the regulative idea of a mental substance, which restores a constitutive notion of accidents, while accepting only a regulative notion of the substance in which those accidents inhere. Similarly, we can construct a regulative principle of causation according to which 'all change [is considered] as belonging to the states of one and the same persisting [mental] being' (A682/B710). As a result, we obtain the Second Analogy, regarding causation, again in a reduced form, which allows mental states resulting from *different* mental faculties to be determined as changing in time, according to causal laws, without determining a substance with a *single* fundamental mental power (see KSS section 5.6.2, pp. 209–15).

This reconstruction of the relational principles for inner experience gives rise to a new interpretation of personal identity, primarily understood as numerical identity through time.⁶ Based on the regulative idea of a mental substance, the notion of personal identity completes the temporal synthesis of all mental states that belong to one and the same person. It thus sensibly explicates the 'concept of the empirical unity of all thought' through time (A682/B710). In contrast to material persistence, personal identity includes not only simultaneously existing parts that make up the whole at a given moment, but also past states that can be remembered, and a projection of future states that cannot (yet) be given in inner intuition, but belong to the projected mental whole to be completed over time. Hence, the crucial difference between outer substances and mental wholes concerns the part-whole relation (and the compositionality) at work in each case. While in both cases the relation of parts and whole is logically defined by the not-yet-schematized category of substance, this relation is explicated in different ways. In the case of material persistence, the whole fully exists and is present at each time and is composed of spatially extended parts that are equally persisting substances 'external to one another' (MFNS, 4: 543). In the case of personal identity, the whole is not fully existing and present at any particular time; it is composed of temporally extended mental states,

but the present (and past) states only make up a certain portion of the whole and the whole is not yet completed in time. The mental whole still logically precedes its parts, namely, mental states, in that without presupposing the whole at least in a regulative sense, we could not determine something as being a part of it, but the parts precede the whole with regard to real existence. The concept of a mental whole is thus indispensably necessary to define the concept of a mental state.

As an important conclusion, we can note that Kant retrieves an account of the mind that runs counter to both Cartesian dualism and materialism. Both these views model the mind on the persistence conditions of matter (explicitly in the materialist case and implicitly in the Cartesian case, for lack of alternatives). Kant, while conceding that material persistence is the only model of persistence we know, allows, on my interpretation, for an alternative account of the unity conditions of personhood through time which is more in line with the Aristotelian tradition. More generally, my interpretation implies that any empirical concept of a *particular kind* of object requires the regulative use of an idea as that which presents a context of intelligibility for the application of the categories. Ideas (e.g. of absolute space, the soul, etc.) can thus be seen as defining the most general *natural kinds* (and their local ontologies) of which we can have determinate experience (e.g. material substances, mental states, etc.).

5. Reply to Clinton Tolley

Tolley focuses on the different concepts by means of which we can reflect upon ourselves, especially the concepts <I>, <subject>, <soul> and <person>, and on the different *self*-relations that these concepts capture. While Tolley agrees with me on the importance of distinguishing these concepts, he worries that there is in addition a ‘thinner notion of “soul”’ at play in Kant’s texts that my interpretation does not sufficiently account for, for instance, the mentions of ‘soul’ as the bearer of mental faculties in the Aesthetic and the Analytic of the *Critique of Pure Reason*. The notion Tolley has in mind is a ‘quite traditional’, broadly Aristotelian notion that refers to ‘any kind of substance which has psychical faculties’ (pp. 489, 484). This notion is thinner in content than Kant’s technical notion of ‘soul’ as an idea of reason, since it does not require intelligence or reason, and thus also broader in scope, since there could be other ensouled beings such as animals falling under it. Similar to Frierson, Tolley eventually concludes that, on Kant’s view, what we refer to in inner experience must be understood as a ‘more substantial’, ‘thing-like’ being that actually possesses mental faculties and engages in mental activities.

In what follows, I explain why I am wary of admitting such a thinner notion of ‘soul’. First, I make some general remarks about the character of my study. Then I defend my reading of Kant’s thicker, technical notion of ‘soul’ as an idea of reason. Finally, I comment on the relationship I see between Kant’s and traditional conceptions of the soul, emphasizing that my interpretation reveals not only Aristotelian but also Platonic aspects in Kant.

5.1 A study of self-concepts and their application conditions

My study can be understood as an analysis of different concepts by which we not only can refer to ourselves, but also describe the kinds of beings we are, such as <subject>,

<soul> or <person>. Since each concept can be understood as defining a set of conditions, my study can be seen as an analysis of different levels of conditions that we have to meet in order to count as beings that fall under the corresponding concepts. The result is the hierarchical order of interdependent self-concepts, from the most general concepts with the most minimal content, such as <subject>, which give the most fundamental set of conditions, to more specific concepts, such as <soul> and <person>, which are more contentful and therefore define more specific conditions. Since these conditions are shown to be necessary conditions for the possibility of *inner experience* and psychological cognition, my study also offers a *transcendental* account of the very kind of being referred to in inner experience – psychological being or being with a certain kind of mentality – without, however, claiming to explain such beings *per se* or *in themselves*, independently of the fact that they appear in experience.

The <I> is special insofar as we can distinguish at least four different uses of the term 'I' in transcendental philosophy (i and iii) and in everyday thought and language (ii and iv): 'I' can be used (i) to express the general *form* of reflexive consciousness, as in the apperceptive 'I think' discussed in the Deduction (B-Deduction, §16), or, if combined with an indeterminate intuition or self-feeling, (ii) to express a *pre-logical, existential* self-consciousness, i.e. the fact that an apperceptive act has taken place (§25). It can also be used (iii) to express the *logical* conditions of self-reference in I-judgements, i.e. the logical 'I', and (iv) to refer to the *empirical* being who is in fact referred to in an empirical I-judgement, i.e. the psychological 'I'. The concept <soul>, I have argued, plays a crucial role in the last use (iv) and thus in understanding both the sensible and intellectual conditions under which we can conceive of ourselves at all by means of the psychological 'I'.

The difficulty in discerning the meaning of <soul> is complicated by the fact that Kant's main concern in the chapter in which he mainly discusses this concept, the Paralogisms, is his critique of rational psychology and thus his rejection of rationalist conceptions of soul, such as those of Wolff and Baumgarten, both of which closely follow a Leibnizian conception. Therefore, not all usages (and accounts) of 'soul' in the Paralogism chapter may reflect Kant's mature theory of souls and persons. Rather, the central textual evidence for my view comes from the Appendix to the Dialectic, where Kant develops a distinctly transcendental conception of the soul as the idea of reason. The difference is also evident in the fact that the Paralogisms and the Appendix differ in the soul predicates they treat, and since I assume that the Appendix reflects Kant's own view more authentically, my interpretation focuses on the four predicates that are most prominent there: <substantiality>, <(bearer of a) fundamental mental power>, <community (with bodies)> and <personality (through time)> (see A682/B700, A672/B690).

5.2 <Soul> and <person> as regulative ideas of reason

Tolley's quest for a thinner notion of soul seems motivated by the insight that what we refer to in inner experience – the object of our I-judgements – must be more than a merely formal element of consciousness: it must be 'more substantial or thing-like' (p. 488), a being that thinks and that 'has existence outside of a representation or content of consciousness' (p. 487). With this dichotomy into 'form' and 'thing', Tolley passes over the fact that every level of representation is imbued with

sensibility, so that we can recognize distinctive sensible conditions for any concept that give it meaning and significance and confirm mental existence, beyond the narrow contents of consciousness. The apperceptive ‘I think’ can manifest itself in empirical consciousness, expressing the fact of an apperceptive act; the logical ‘I’ can manifest itself in its empirical use in thought and language. But Tolley’s focus on the bifurcation between ‘form’ and ‘thing’ reflects the historical context: it was precisely the distinction that troubled the German idealists, who struggled to reconcile what they perceived as an overly formal approach in Kant with an ontological foundation of transcendental philosophy, leading them to proclaim – in different ways and nuances – the *Absolute* as that in which form and thing, or thought and being, are simply one and the same.⁷ With my interpretation I aim to bring to the fore a third alternative between a formal subject and a soul-thing (or soul-substance), namely, the idea of the human person as a *becoming empirical reality* – a reality that can be said to exist in virtue of being the object of inner intuition and perception, but also, since it is never fully given in intuition, a reality that must be assumed to evolve in accordance with rational ideas. These ideas are the only kinds of concepts we can form of it as a whole and by means of which we can recognize its parts, or states, as belonging to the whole. By emphasizing the functional role of the idea of the soul exclusively in the context of experience, my interpretation ties in with another line of development in the history of Kant’s reception, namely, with experimental psychologists like Wilhelm Wundt and Neo-Kantians like Paul Natorp who gave Kantian thought an important role in the development of psychology in the nineteenth and early twentieth centuries.

The question, then, is not whether there is something existing that thinks, but whether that which exists and thinks can be truthfully described by the concepts and ideas we form of it. As discussed in my replies to Wood and Frierson, there are difficulties in the application of the relational categories, since they require a substratum that cannot be given in the merely temporal manifold of inner sense. With the help of reason, the relational schemata of the understanding are replaced by a regulative use of the predicates of the soul.⁸

In both the Paralogisms and the Appendix, Kant discusses a fourth predicate, <personality>, often rendered as personal identity through time. This fourth predicate builds on the earlier predicates, especially on <substance>, but like them, should be understood in a regulative way. In chapter 7, I then extend my analysis of personal identity to comprise not only a quantitative or numerical identity through time, but also a qualitative identity, which consists in the formation of a consistent character that is stable through time. The fact that Kant mentions <personality> as a predicate of <soul> implies that <person> is an important specification of <soul>, although not all beings who fall under <soul> are at the same time persons. While both <soul> and <person> are in principle applicable to other mental beings, for Kant such beings must be capable of thinking (e.g. B415, A682/B710).

Moreover, I do not think that Kant’s use of ‘soul’ in the *Prolegomena* differs substantially from the technical term in the *Critique*. The ‘idea of the complete subject’ in the *Prolegomena* (4: 330), I argue, occupies the same role as the idea of the soul in the *Critique* (see Kraus 2021). That a proof of ‘the persistence of the soul’ is available only ‘in life’ does not necessarily imply the constitutive applicability of the category of substance (*Prol*, 4: 335). Rather, it seems to me consistent with Kant’s argumentation

in the Refutation of Mendelssohn's Proof in the *Critique* (see B415ff.) and allows for two interpretations: either 'persistence' refers to the soul's persistence in the objective time of physical bodies, which depends on the persistence of the living, embodied human being (*Mensch*) during *organic* life, or it is meant in the sense of the soul's personal identity during its *mental* life, in which case the concept of mental life necessarily entails its identity. I also understand Kant's earlier mentions of the notion of 'soul' in the *Critique*, such as 'faculties of the soul' (A94, A115), as transcendental reflection concepts that do not qualify for constitutive use, even if their function as ideas of reason is not yet explicitly introduced.

5.3 Kant's versus traditional conceptions of the soul

While I agree with Tolley that Kant's conception of the soul is inspired by Aristotle's, I see important differences. In Aristotle and the Aristotelian medieval tradition, we find the tripartite distinction between (i) *anima* (= sensitive soul in general), (ii) *animus* (= specifically human soul with intellect) and (iii) *mens* (= pure intellect). Kant himself appeals to this tripartite distinction, in his lectures on anthropology, as three ways in which we can view ourselves (see L-Anth/Parow, 25: 247, L-Anth/Collins, 25: 16). In the *Critique*, the idea of the soul is more narrowly confined to *thinking* beings who possess higher intellectual faculties, rather than beings with 'psychic faculties', and therefore seems to correspond to *animus*, rather than *anima*.

Kant's hylomorphic theory of mental faculties – in terms of an *a priori* form that is made actual by the intake of suitable material under this form – is, I believe, Aristotelian in nature. However, there is an important difference: while for Aristotle forms define (essential) features of beings in themselves, for Kant the forms of faculties, especially those of the understanding, define forms of experience and hence necessary features of objects of experience. The idea of the soul is understood as form-giving, however, not in the sense of a *forma substantialis* that gives form to the human body and unity to all mental faculties (see *De Anima* 2.1–2). Rather, it defines primarily the form of inner experience and hence the object thereof, i.e. the psychological person. Combined with the corresponding forms of outer experience and the forms of life, it defines the essential features of an embodied human person.⁹

Another difference from Aristotle is the normative aspect of personhood and the practical effectiveness of the idea of the soul as the principle of self-formation. This part of my interpretation corresponds more to a Platonic conception of the soul, according to which the essence of the soul consists not only in the factually existing unity of mental faculties, but more importantly in the striving for self-perfection and further unification of the disparate parts of the soul. Kant himself suggests that his account of ideas of reason is Platonic in origin (see A313/B370), and Kantian ideas can hence be understood as norm-giving ideals for developmental processes. An important representative of a Platonic conception in the German Enlightenment is Moses Mendelssohn, who – following Plato's conception in the *Phaedo* – conceives of the soul as effecting inner mental activities, whose common essential purpose is the soul's own self-perfection (see Mendelssohn 1767). Although Kant criticizes Mendelssohn's rational psychology, he agrees – according to my interpretation – with its normative implications. Kant's idea of the soul is practically efficacious in striving towards self-perfection according to the ideal of mental wholeness. To rebut the objection of

circularity, my interpretation does not imply that there is another kind of being prior to being a person that is ‘on the road to personhood’ (p. 489). It is not a genetic account that explains the coming-to-be of a person, but an analysis of the faculties required for being a person. The process of self-formation does not go from something to a person, but from a person to a more perfect or more unified person, striving for the rational ideal of a systematic whole.

Notes

- 1 The contributions by Janum Sethi and Clinton Tolley were first presented at a book launch held by the History of Philosophy Forum at the University of Notre Dame on 19 March 2021, organized by Thérèse Cory. The contributions by Patrick Frierson and Allen Wood were first presented at a book symposium held at the Université Paris VIII on 3 June 2021, organized by Stefanie Buchenau. I am grateful to the two organizers for making these events possible in virtual form despite the obstacles of a pandemic and I thank the audience of both events for lively discussions. I also thank the editors of *Kantian Review* for their generous offer to publish these comments and my replies, and I thank Janum, Clinton, Patrick and Allen for helpful feedback on my replies.
- 2 A portion of this précis is taken from an SGIR Review symposium (forthcoming) and is included here with the permission of the *Journal of the Society for German Idealism and Romanticism*.
- 3 This view is further explored in Sethi (2020) and Sethi (2021).
- 4 Note that the grammatical subject of [1] is ‘I’. [1] is therefore ambivalent, since it could also be used to express an objectively valid I-judgement, i.e. an inner experience.
- 5 See Friedman 2013: 323ff.
- 6 Chapter 7 offers an expanded interpretation of personal identity that includes the qualitative sense of a unified character pursued throughout a lifetime.
- 7 For discussion of my interpretation in light of later developments in German idealism, especially with respect to Fichte and Maimon, see Kraus (forthcoming).
- 8 For discussion, see my reply to Frierson.
- 9 The form of the body, for Kant, is defined by the forms of outer experience and the idea of purposiveness, and is described by the laws of physics, chemistry and biology. Moreover, I acknowledge that there are higher forms for the self-formation of persons, such as social, political and cultural forms.

References

- Freud, Sigmund (1900/2010) *Die Traumdeutung*. In *Gesammelte Werke*, Vol. 2/3 (Frankfurt am Main: Fischer).
- Friedman, Michael (2013) *Kant’s Construction of Nature*. Cambridge: Cambridge University Press.
- Kraus, Katharina (2021) ‘Kant’s Argument Against Psychological Materialism in the *Prolegomena*’. In Peter Thielke (ed.), *The Critical Guide to Kant’s Prolegomena* (Cambridge: Cambridge University Press), 154–74.
- (forthcoming) ‘Précis of Kant on Self-Knowledge and Self-Formation’ and ‘Replies to Critics’. *Journal of the Society for German Idealism and Romanticism*.
- Mendelssohn, Moses (1767) *Phaedon oder über die Unsterblichkeit der Seele: in drey Gesprächen*. Berlin, Stettin: Nicolai.
- Sethi, Janum (2020) “‘For me, in my Present State’”: Kant on Judgments of Perception and Mere Subjective Validity’. *Journal of Modern Philosophy*, 2(9), 1–20.
- (2021) ‘Kant on Empirical Self-Consciousness’. *Australasian Journal of Philosophy*. DOI: [10.1080/00048402.2021.1948083](https://doi.org/10.1080/00048402.2021.1948083)

Cite this article: Kraus, K.T. (2022). Précis of *Kant on Self-Knowledge and Self-Formation* and Replies to Critics. *Kantian Review* 27, 491–508. <https://doi.org/10.1017/S136941542200022X>