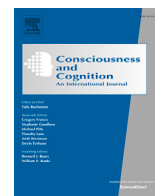




ELSEVIER

Contents lists available at ScienceDirect

# Consciousness and Cognition

journal homepage: [www.elsevier.com/locate/yccog](http://www.elsevier.com/locate/yccog)

Full Length Article

## Towards a structural turn in consciousness science

Johannes Kleiner<sup>a,b,c,d,\*</sup><sup>a</sup> Munich Center for Mathematical Philosophy, Ludwig-Maximilians-Universität München, Geschwister-Scholl-Platz 1, 80539 München, Germany<sup>b</sup> Graduate School of Systemic Neurosciences, Ludwig-Maximilians-Universität München, Großhaderner Str. 2, 82152 Planegg-Martinsried, Germany<sup>c</sup> Institute for Psychology, University of Bamberg, Markusplatz 3, 96047 Bamberg, Germany<sup>d</sup> Association for Mathematical Consciousness Science, Geschwister-Scholl-Platz 1, 80539 München, Germany

### ARTICLE INFO

#### Keywords:

Quality space

Qualia space

Phenomenal space

Theory of consciousness

Structuralism

Structure-preserving mapping

Isomorphism

### ABSTRACT

Recent activities in virtually all fields engaged in consciousness studies indicate early signs of a structural turn, where verbal descriptions or simple formalisations of conscious experiences are replaced by structural tools, most notably mathematical spaces. My goal here is to offer three comments that, in my opinion, are essential to avoid misunderstandings in these developments early on. These comments concern metaphysical premises of structural approaches, the viability of structure-preserving mappings, and the question of what a structure of conscious experience is in the first place. I will also explain what, in my opinion, are the great promises of structural methodologies and how they might impact consciousness science at large.

### 1. Introduction

So far, the scientific study of consciousness has mainly employed verbal and linguistic tools, as well as simple formalisations thereof, to describe conscious experiences. Typical examples are the distinction between ‘being conscious’ and ‘not being conscious’, between whether a subject is ‘perceiving a stimulus consciously’ or not, between whether a subject is ‘experiencing a particular quale’ rather than another, or more generally any account of whether some  $X$  is part of the phenomenal character of a subject’s experience at some point of time. Formalisations of these verbal descriptions mostly make use of set theory, examples being sets of states of consciousness of a subject and simple binary classifications, or of real numbers, for example to model ‘how conscious’ a system is. There are sophisticated mathematical techniques in the field, but to a large extent they only concern the statistical analysis of empirical data, and the formulation of a theory of consciousness itself—but not the description of conscious experiences which underlies the data collection or modelling effort.

Much like words shape thoughts, descriptions shape science. In the case of consciousness studies, the descriptions that were available so far have fed into theories of consciousness, have determined what can be inferred about the state of consciousness of a subject, and have guided ways of conceptualising the problem under investigation.

They have, for example, led to a number of theories that explain what it takes for a single stimulus or a single piece of information to be consciously experienced, but which remain silent or vague on how the phenomenal character as a whole is determined. They have led to measures of consciousness which are specifically tailored to find out whether a single stimulus or single quality is experienced consciously (Irvine, 2013), but are not meant to infer phenomenal character beyond this. And to some extent, at least, they have privileged research programmes which search for either-or conditions related to consciousness, such as arguably the

\* Correspondence to: Ludwig-Maximilians-Universität München, Geschwister-Scholl-Platz 1, 80539 Munich, Germany.

E-mail address: [johannes.kleiner@lmu.de](mailto:johannes.kleiner@lmu.de).

<https://doi.org/10.1016/j.concog.2024.103653>

Received 6 October 2023; Received in revised form 22 January 2024; Accepted 30 January 2024

Available online 28 February 2024

1053-8100/© 2024 The Author.

Published by Elsevier Inc.

This is an open access article under the CC BY license

(<http://creativecommons.org/licenses/by/4.0/>).

search for Neural Correlates of Consciousness (NCCs) that is largely predicated on a conception of having “any one specific conscious percept” (Koch et al., 2016).

Because verbal descriptions only parse part of the phenomenal character of an experience, part of what it is like for an organism to live through a particular moment, it is no surprise that means to go beyond these simple descriptions are highly sought after.

In recent years, the idea of using mathematical spaces, or mathematical structure more generally,<sup>1</sup> to go beyond verbal descriptions and simple formalisations have started to sprout in virtually every discipline involved in the scientific quest to understand consciousness. Following rich developments in psychophysics over more than a century (Pashler & Wixted, 2004), and pioneering work by Austen Clark (Clark, 1993) and David Rosenthal (Rosenthal, 1991) in consciousness science, mathematical spaces are now applied in philosophy, (Clark, 2000, Coninx, 2022, Fortier-Davy & Millièrre, 2020, Gert, 2017, Lee, 2021, 2022, Rosenthal, 2010, 2015, 2016, Fink et al., 2021, Lyre, 2022, Kob, 2023, Renero, 2014, Prentner, 2019, Yoshimi, 2007, Chalmers & McQueen, 2022, Silva, 2023, Atmanspacher, 2020), neuroscience (Tononi, 2015, Tallon-Baudry, 2022, Zaidi et al., 2013, Lau et al., 2022, Malach, 2021, Haun & Tononi, 2019, Oizumi et al., 2014, Hebart et al., 2020, Josephs et al., 2023, Tsuchiya et al., 2023, Zeleznikow-Johnston et al., 2023, Haynes, 2009, Michel, *In press*), cognitive science (Hoffman et al., 2023, Rudrauf et al., 2017, Hoffman & Prakash, 2014, O'Brien & Opie, 1999), psychology (Klincewicz, 2011, Kostic, 2012, Young et al., 2014) and mathematical consciousness science (Grindrod, 2018, Kleiner, 2020b, Stanley, 1999, Resende, 2022, Mason, 2013, 2021, Signorelli & Wang & Coecke, 2021, Tsuchiya et al., 2016, Tsuchiya & Saigo, 2021, Tsuchiya et al., 2022, Kleiner, 2020a, Kleiner & Hoel, 2021, Kleiner & Ludwig, 2023). They are known under various names, including quality spaces (Clark, 1993, Rosenthal, 2015), qualia spaces (Stanley, 1999), experience spaces (Kleiner & Hoel, 2021, Kleiner & Tull, 2021, Rosenthal, 2010), qualia structure (Kawakita & Zeleznikow-Johnston & Tsuchiya, et al., 2023, Kawakita & Zeleznikow-Johnston & Takeda, et al., 2023, Tsuchiya et al., 2022), Q-spaces (Chalmers & McQueen, 2022, Lyre, 2022), Q-structure (Lyre, 2022),  $\Phi$ -structures (Tononi, 2015), perceptual spaces (Zaidi et al., 2013), phenomenal spaces (Fink et al., 2021), spaces of subjective experience (Tallon-Baudry, 2022), and spaces of states of conscious experiences (Kleiner, 2020a). A first formalised theory of consciousness to make use of mathematical spaces was Integrated Information Theory (IIT) 2.0 (Tononi, 2008); more recent versions expand and refine the idea (Oizumi et al., 2014, Albantakis et al., 2023).

What unites all of these proposals is the hope that the mathematical structures they propose are useful to describe the phenomenal character of an experience more comprehensively, more precisely, or more holistically than verbal descriptions or simple formalisations allow, and that mathematical structures can cope both with the apparent richness and with the many details that make up experiences. If this hope turns out true, it has far-reaching implications on how to study, measure, and think about consciousness.

My goal here is to offer three comments which I think are important to keep in mind when applying structural ideas in theory and experimental practice, so as to avoid misconceptions or misunderstanding early on. I hope that my comments are helpful for those working on structural ideas as well as those observing these developments with a degree of scepticism.

## 2. Three promises of a structural turn

Before offering my comments below, I will briefly sketch the implications and limitations that structural methodologies may have for consciousness science. This might be of interest to those who have not engaged with this research before, and allows me to illustrate what I think are some of the great promises of a structuralist turn.

### 2.1. Theories of consciousness

We currently have at least 39 theories of consciousness,<sup>2</sup> with new theories being proposed on a regular basis, albeit without much general attention. The reason for that, I contend, is that as far as theoretical work is concerned, it is actually very easy to come up with theories of consciousness of the type we have today.

The majority of contemporary theories of consciousness aim to explain whether a system's state, a stimulus, a piece of information, or a representation is consciously experienced, or not. That is, they target a *binary classification* between states, signals, stimuli or

<sup>1</sup> The term *mathematical structure*, which I will explain in detail Section 3 below, is more general than the term *mathematical space*. That is, every mathematical space is a mathematical structure, but there are also mathematical structures which are not mathematical spaces, either because they only comprise individuals (so do not satisfy the intuition that a space is about many individuals), or because their structure is more complex than one would typically take a space to be. The question of which mathematical structures to call mathematical spaces is a matter of convention, which is why there is no definition of a general concept of mathematical space in mathematical logic.

<sup>2</sup> An unpublished list compiled by Dr. Jonathan Mason on behalf of the *Association for Mathematical Consciousness Science* (AMCS) and the *Oxford Mathematics of Consciousness and Applications Network* (OMCAN) comprises the following theories of consciousness in the peer-reviewed literature: Activation/Information/Mode-Synthesis Hypothesis, Adaptive Resonance Theory, Attention Schema Theory, Centrencephalic Proposal, Conscious Agent Networks, Conscious Turing Machine, Consciousness Electromagnetic Information Field Theory, Consciousness State Space Model, Cross-Order Integration Theory, Dendrite/Apical Dendrite Theory, Dynamical Core Theory, Electromagnetic Field Hypothesis, Enactive and Radical Embodiment, Expected Float Entropy Minimisation, First Order Representational Theory, Free Energy Principle Projective Consciousness Model, Global Neuronal Workspace Theory, Global Workspace Theory, Higher-Order Thought Theory, Integrated Information Theory, Integrated World Modeling Theory, Layered Reference Model of the Brain, Memory Consciousness and Temporality Theory, Mesocircuit Hypothesis, Multiple Draft Model, Network Inhibition Hypothesis, Neural Darwinism Theory, Orchestrated Objective Reduction, Passive Frame Theory, Predictive Processing and Interoception, Proto-Consciousness Induced Quantum Collapse, Psychological Theory of Consciousness, Radical Plasticity Thesis, Recurrent Processing Theory, Self Comes to Mind Theory, Semantic Pointer Competition Theory, Single Particle Consciousness Hypothesis, Temporo-Spatial Theory of Consciousness, Thalamo-Cortical Loops and Sensorimotor Couplings. This list might not be complete, and some of the theories might point to similar or analogous theoretical constructs.

representations. The simple verbal distinctions mentioned in the introduction—a system ‘being conscious’ or not, ‘perceiving a stimulus consciously’ or not, ‘experiencing a particular quale’ or not—are all examples of such binary classifications.

Formulating theories of consciousness that target binary classification is relatively straightforward, as far as theoretical work is concerned. This is because devising a  $\{0, 1\}$  classification only requires identifying some property, function, or dynamical mode of a brain mechanism. All configurations that exhibit this property, function or dynamical mode are mapped to 1, while all which do not are mapped to 0. And within non-structural approaches, nothing technical prohibits one from postulating that the 1 cases correspond to conscious experience of a stimulus, state, piece of information or representation, while the 0 cases correspond to unconscious experience thereof. The empirical or conceptual validity of such a choice is an important question, yet from a technical standpoint, formulating theories that target these distinctions is straightforward.

It is much more difficult to come up with a well-formed hypothesis that relates to a mathematical space or mathematical structure. That is because a mathematical space or mathematical structure has two parts. On the one hand, it contains a set of points. On the other hand, it contains relations or functions that express connections between the points, for example an order relation or a metric function. Therefore, there is much more information to provide when specifying how a space or structure relates to a brain mechanism, or a physical system more generally. Furthermore, virtually every mathematical object comes with a set of axioms that parts of the object have to satisfy. So not only is more information needed, but this information may also have to satisfy constraints to provide a legitimate definition. This is why defining a space or structure is much more of a challenge than finding a binary classification.

The task is more difficult even if the space or structure that a theory is to provide has a specific, theory-independent form. That is the case if the theory has to account for phenomenal structure that has independent justification or independent motivation, for example from psychophysical experiments. This difficulty is illustrated by the fact that we do not, at present, have a theory of consciousness that targets the mathematical structures that have been proposed to account for conscious experiences on independent grounds. To the best of my knowledge, there are only two theories that define phenomenal spaces: Integrated Information Theory (IIT) (Albantakis et al., 2023) and Expected Float Entropy Minimisation Theory (EFE) (Mason, 2021). While both theories represent significant advances, establishing a link to existing phenomenal spaces (cf. Section 5) remains a next-level challenge.<sup>3</sup>

As formulating theories that account for phenomenal structure in addition to non-structural explananda necessitates meeting more constraints than formulating non-structural theories, structural theories are likely to be more predictive than their non-structural counterparts. Furthermore, because the phenomenal structure is an integral aspect of phenomenal character, a theory that accounts for phenomenal structure in addition to non-structural explananda has a broader explanatory scope than one that focuses solely on the conscious-unconscious distinction. Therefore, a structural turn might deliver more explanatory and more predictive theories of consciousness. This is the first major implication I can see of structural approaches in consciousness science.<sup>4</sup>

Structural methodologies might inspire, and be inspired by, novel theoretical ideas that derive from any of the existing theories of consciousness, or from their combination. Proposals like the Conscious Turing Machine (Blum & Blum, 2022) or Integrated World Modelling Theory (Safron, 2022) that combine features of existing theories of consciousness (such as, for example, Integrated Information Theory, Global Neuronal Workspace Theory, and Free Energy Principle based proposals) could be particularly interesting in this regard.

## 2.2. Experimental investigations

A shift towards structural methodologies could also have significant implications for experimental research. One immediate implication follows from the previous section, i.e., from the transformative effect that structural methodologies could have on theories of consciousness. If structural theories of consciousness would indeed be more predictive than the non-structural theories we have today, then they might be easier to test than the theories we have today,<sup>5</sup> and the new predictions about structural facts might offer new avenues for experimental investigation.<sup>6</sup>

But structural thinking could also yield new experimental tools and methodologies that are separate from theoretical advancements. For instance, under certain conditions, structural approaches offer an entirely new methodology for measuring NCCs (Fink et al., 2021). This methodology could potentially address some of the foundational challenges in existing methodologies, such as the co-activation of cognitive processing centres causally downstream of the core NCC, and might not require traditional methods to assess a subject’s state of consciousness. I discuss and criticise the key assumption that enables this methodology—the assumption of a structure-preserving mapping between phenomenal and neuronal structures—in Section 4 below. Nevertheless, even if this

<sup>3</sup> Proponents of both theories are fully aware of this task, and IIT has made a first step in this direction in Haun and Tononi (2019). In addition to accounting for phenomenal structure that has independent justification, there are other tasks and challenges that structural theories have to meet and resolve. For example, an anonymous reviewer has kindly pointed me to the fact that according to IIT, richly structured experience can be entailed by static systems without dynamics, which might pose an empirical or conceptual challenge for IIT.

<sup>4</sup> In saying this, I do not intend to diminish the value of ‘binary’ theories of consciousness. They are an integral part of consciousness science and encapsulate a substantial body of evidence. On my view, they need to be extended so as to address phenomenal character more holistically as well. Whether this should be done on a case-by-case basis, or whether there might be a theory of qualitative character that can serve for a larger number of binary theories, is not something that needs to be decided in advance.

<sup>5</sup> Lukas Kob made this point for *structuralist* approaches during a wonderful talk at the recent *Structuralism in Consciousness Studies* workshop at the Charité Berlin, though my comment here concerns the wider scope of *structural* approaches, cf. Section 3 and Fig. 2 for more on that distinction.

<sup>6</sup> Speculating wildly, one might hope that if theories of consciousness could account for *theory-independent* phenomenal spaces, this could help to mitigate the problem that empirical tests of theories of consciousness currently rely heavily on theory-dependent *methodological choices* (Yaron et al., 2021).

assumption proves to be more limited in scope or strength than initially anticipated, the methodology might still have advantages compared to existing options to search for NCCs.

The implication that intrigues me most, however, is the possibility that structural approaches may introduce new *measures of consciousness*. A measure of consciousness, as conventionally understood, is a method to determine whether an organism is conscious, or whether a given stimulus or signal has been consciously perceived. Measures of consciousness are “consciousness detection procedures” (Michel, 2023) of sorts.

Building on the extensive previous work in both psychophysics and consciousness science, structural approaches raise the possibility to construct new and potentially more powerful measures of consciousness, which do not only focus on whether a single stimulus is experienced (a single quality of phenomenal character, that is), but on phenomenal character more comprehensively.

The potential of structuralist approaches in this regard can be nicely illustrated by considering verbal report, which is a paradigmatic (albeit often criticised) measure of consciousness. In the case of report, subjects use language to report facts about their experience. They might, for example, indicate that they experienced a red colour, or saw a face in a masked stimulus. The problem with reports is that when compared with the actual experience, they contain very little information. Which shade of red did the subject experience, precisely? How did they experience the face, and with which details? What else did they experience in addition to the reported fact? In information-theoretic terms, this problem arises because the channel capacity of verbal report and other behavioural indicators is low compared to the information content of conscious experiences.<sup>7</sup>

Structural approaches allow us to bypass the limited channel capacity of reports and similar measures of consciousness, because structural descriptions can *store information* about the phenomenal character of a subject. That is the case because structural descriptions represent features of a subject’s phenomenal character that relate individual non-structural facts.

Given the structural information in a phenomenal space, a few bits of information collected in an experimental trial, for example by means of reports or similar measures of consciousness, can suffice to pin down the location in a structure, resulting in information about what a subject is experiencing that might go far beyond the bits of information that were collected. This is similar to how a geographic map can be used to decode rich information about a path based on a few bits of information about location. Finding one’s way in the wilderness without a map or map-like tools generally is a very difficult task. Given a map, procedures like triangulation are available that only require a few bits of information, such as the angles between three landmarks in line of sight, to pin down one’s position and find one’s way. That is possible because maps store information about geography. Another example of this sort is quantum tomography, where a set of carefully chosen measurements together with structural information about the quantum state (specifically, the inner product and projective structure of the Hilbert space), is used to pin down the exact state among an infinite number of possibilities.

In a similar vein, phenomenal spaces might be used to decode information from carefully chosen low-channel-capacity measures of consciousness. Precisely how to do this remains an open question as of yet, and strongly depends on a thorough understanding of phenomenal structure in the first place (cf. Section 5), but it is a viable possibility.

### 2.3. Conceptual work

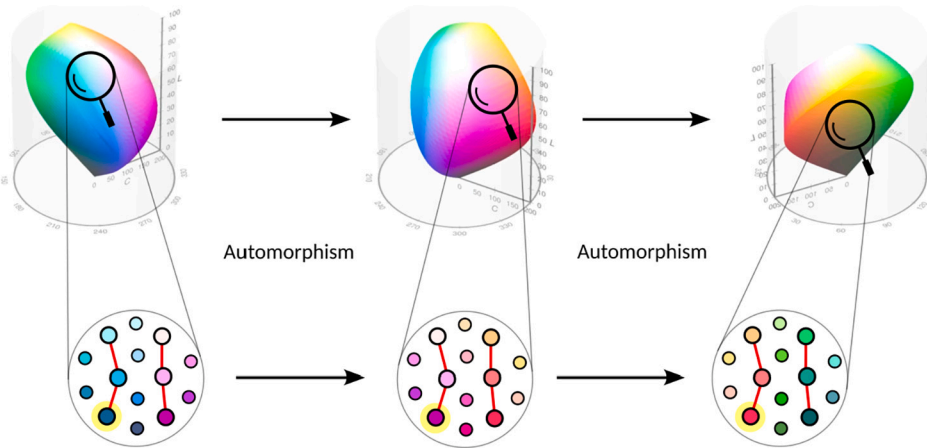
Structural approaches can also be essential, finally, in *conceptualising consciousness and its potential problems*. It is not unlikely that interesting philosophical implications arise, specifically in the context of structuralist assumptions, but what I’d like to highlight here is the importance of structural thinking in shaping our pre-theoretic problem intuitions about consciousness; those intuitions, that is, which guide both our theorising and experimental work.

Structural thinking might well turn what we previously thought about consciousness upside down. It might change how many of us think about our own research in the first place. To give two very preliminary examples, I think that structural approaches are relevant for epistemic arguments like Mary’s room (Jackson, 1998, 1986), and for modal arguments like colour inversion (Shoemaker, 1982, Block, 1990).

For epistemic arguments such as Mary’s room, the big question is whether one presumes that structural facts about experiences are known. If Mary propositionally knows, for example, which structure the experience of red has, and if structure is sufficient to individuate experiences, then she might be able to use her advanced neuroscience knowledge to create an embedding of the structure of red experiences within her own phenomenal space, even if she never experienced red, or any colour for that matter, before. Similarly, outside the realm of thought experiments, we might use structural facts to create experiences that approximate what it is like to be a bat. Structure might furnish an objective phenomenology (Lee, 2022).

Modal arguments, similarly, need to be rethought. The typical colour inversion thought experiment presumes fairly homogeneous colour spaces—colour spaces that possess symmetries. This presumption is critical because if a colour inversion is not a symmetry, then the difference between colour experience before and after the inversion will manifest itself both in behaviour and in the use of colour words: through similarity judgements and other expressions of structural facts. The closest approximation we have to a space of consciously experienced colour qualities is the CIELAB colour space (Schanda, 2007), a rendering of which is depicted in Figs. 1, 3, and 4, which is highly non-homogeneous and may not admit symmetries to the extent that we expect. Adding valence and other consciously experienced attributes of colour experiences might further erode any remaining symmetries. Thus, at least the usual intuitions regarding qualia inversions and other modal arguments may cease to be valid. Structural approaches might force us to reconsider intuitions that are built on these types of arguments.

<sup>7</sup> I am very grateful for a conversation with Lucia Melloni about the problems of reports and structural ideas to resolve these during a walk at the above-mentioned *Structuralism in Consciousness Studies* workshop. The idea sketched here came up during this walk and is Lucia’s as much as, or even more than, mine.



**Fig. 1. What is an automorphism?** This figure illustrates the concept of automorphisms. Automorphisms are somewhat analogous to rotations of a space around some axis (top row). More formally, an automorphism is a function that maps every point of a mathematical space to a different point of the same space, one-to-one, in such a way that all relations of the space are preserved: whenever two points are related before the mapping, they are also related after the mapping. This is illustrated by the bottom row, where individual points of the space are depicted by coloured dots, and relations are depicted by red lines. An automorphism maps every dot to a new dot, represented here by the change in location of the colours, in such a way that when two dots were related before the mapping (red line between two dots) the targets of the mapping are also related (red line between target dots). Automorphisms form a group because automorphisms can be inverted, and because any two automorphisms can be combined to form another automorphism, in this case one from the left-hand space all the way to the right-hand space. (For interpretation of the colours in the figure(s), the reader is referred to the web version of this article. Depiction of CIE colour space gamuts by Wikimedia Commons, Michael Horvath, under [CC BY-SA 4.0](#); this image is shared under the same license.)

## 2.4. Limitations

While structural approaches do, in my view, offer a number of benefits to the science of consciousness, it is also important to see their limitations.

A first limitation of structural approaches is that it is not clear, at present, how much of phenomenal character—how much of what it is like to experience something, that is—can be grasped by structural tools. While it is clear that much of the phenomenal structure that is usually associated with the content of consciousness can be represented structurally (much of it actually is structural, one might say), it is not clear whether some of the more subtle or remote facets of phenomenal character are amenable to a structural analysis. Can the experience of a self or ego be represented structurally? What about the experience of other minds? Or the pre-reflective and pre-conceptual awareness of being aware, sometimes referred to as subjective character?

A second limitation of structural approaches relates to measurability. Even if a facet of phenomenal character is amenable to structural tools, it might still be difficult, costly, or even impossible to measure. It might take years to construct a full quality space of a single modality. Is this actually feasible in experimental practise for anything but the most salient structures of phenomenal character?

A third limitation is the question of whether structural approaches can actually get any closer to modelling what is sometimes described as an intrinsic nature of qualia or qualities. Do structural approaches have any handle on modelling this? Or can they just circumscribe the structure that intrinsic properties instantiate? And to the extent that such intrinsic nature is the core of the problem of consciousness, can structural approaches get us any closer to understanding this core?

My own view of these limitations is that they define some of the key questions that structural research will have to tackle in the upcoming years. Because experiences exhibit structure, structural approaches are, by necessity, part of any research programme that targets experiences in full. To what extent they contribute to resolving the core questions at the heart of consciousness science is an open question.

## 3. Metaphysical premises

My first comment concerns an intuition which I have often encountered when discussing structural approaches with colleagues: that structural approaches are metaphysically presuming. Most notably, to many they seem to be tied to physicalist or reductionist metaphysics. The goal of this comment is to show that this is not the case. Structural approaches offer a new descriptive tool that can—in theory, at least—be applied independently of metaphysical assumptions, and in research programs of any metaphysical flavour. Structural approaches do not in themselves have metaphysical premises, and they do not by themselves come with a preferred metaphysical interpretation. Rather, they can be applied to and combined with the particular metaphysical ideas or presumptions that are already employed in a research program.

The major reason why structural approaches are often taken to be metaphysically presuming is that they are conflated with structuralist approaches. Structuralist approaches assume that individuals can be individuated by structure: that for every individual  $x$ , there is a unique location in a structure, a location in which only  $x$  holds. Intuitively speaking, the idea is that specification of all structural facts suffices to also specify all facts about individuals in that structure.

In the context of consciousness science, the individuals in question can be experiences, phenomenal character, qualities or qualia. The structures in question are experience spaces (spaces whose elements are experiences), phenomenal spaces, quality spaces or qualia spaces. Furthermore, there are ontological, epistemological, and methodological ways of reading a structuralist claim. In all cases, the idea is that the domain of individuals exhibits structure, and that this structure is sufficient to individuate the individuals in the relevant sense.

Structural approaches, in contrast, are not committed to a claim of individuation. An approach is structural if it applies mathematical structure. And, as I will now explain, more often than not, mathematical structure does not individuate individuals. In order to see why, we must differentiate between two readings of the term ‘structure’. This will also yield a clear, formal definition of structuralism in a given consciousness-related domain.

Mathematics offers an unambiguous definition of what a structure is. A *mathematical structure* consists of two things: domains, on the one hand, and functions or relations, on the other hand. The domains of a structure are the sets on which the structure is built. They comprise the points, or elements, in a space, the individuals in a structuralist sense. In the case of a metric space, for example, there are two domains: the set of points of the metric space and the real numbers that constitute the ‘distances’ between points. In the case of a partial order, there is just one domain: the domain of elements that are to be ordered. The second ingredient of a mathematical structure are functions and/or relations. Functions map some of the domains to other domains. In the case of a metric structure, for example, there is a metric function that maps two points to a real number. Relations relate points to each other. In the case of a partial order, for example, there is a binary relation on the set of points. This relation specifies ordered pairs of points, usually written as  $p_1 \leq p_2$ .

When the term ‘structure’ is used in natural science, it usually follows this mathematical definition. For example, if we talk about the structure of space-time, we mean the mathematical structure that describes space-time, called a Pseudo-Riemannian manifold. If we talk about the structure of a neural network, we mean the mathematical structure of the directed graph that specifies the connectivity of the network: the mesh of nodes and edges, where each node represents a neuron or neuronal assembly, and where each directed edge specifies a neural pathway between neurons or assemblies.

When we use the term ‘structure’ in the context of structuralist ideas, however, it only refers to the second ingredient of a mathematical structure: the functions and relations that a mathematical structure contains. These functions or relations are what individuates the individuals—the elements of a domain—in a structuralist sense.

While customary in the context of structuralist assumptions, this use of the term ‘structure’ to designate only relations and functions is problematic. That is the case because we cannot actually specify relations or functions without specifying the points or elements that the relations or functions operate on. The symbol ‘ $\leq$ ’, for example, can be used to indicate a type of structure, a partial order in this case, but it cannot define or specify a structure. Any concrete definition or specification of a partial order needs to make use of, or refer to, the points that the relation relates. It needs to make use of some set of points—some domain in the mathematical sense of the term. Strictly speaking, it does not make sense to use the term ‘structure’ to refer *only* to the functions or relations. I will refer to structure in the structuralist sense—that is to the functions and relations that are part of structure in the proper sense of the term—as *structure in the narrow sense of the term*.

The structuralist idea that relations or functions determine all individuals still makes sense, of course, independently of terminological issues. And it can be expressed in a neat formal requirement, making use of the notion of an automorphism, cf. Fig. 1. An automorphism is a one-to-one mapping from the domains of the structure to themselves which preserves the functions or relations. That is, it preserves structure in the narrow sense of the term. For every point of the structure, an automorphism specifies a point as its target in such a way that the functions and relations of the structure do not change when going from the source to the target: whenever some points satisfy a relation before the mapping, they also satisfy the relation after the mapping, and equally so for functions.

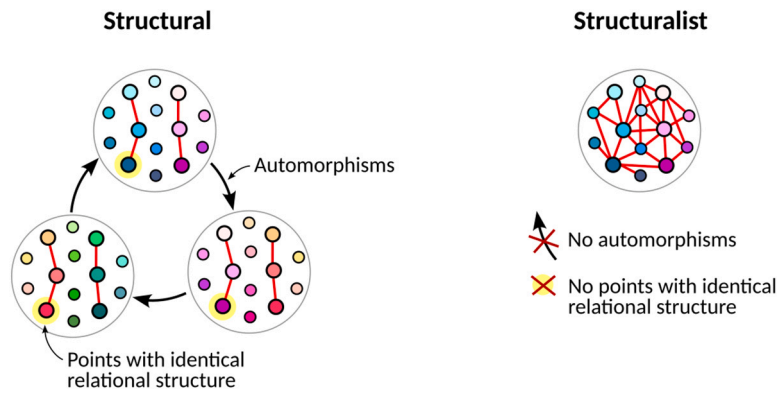
Automorphisms may or may not exist. The identity mapping (not changing anything) is always an automorphism, but depending on how rich or complex the structure in the narrow sense of the term is, there might not be other automorphisms. In particular, if it is indeed the case that every point  $x$  of a structure satisfies a unique location of structure in the narrow sense of the term, then there is no automorphism other than the identity. One cannot exchange any two points without changing structure in the narrow sense of the term. In this case, one says that the *automorphism group is trivial*.<sup>8</sup> Vice versa, if the automorphism group of a structure is trivial, then every point must have a unique location.<sup>9</sup>

As structuralism (in the context of consciousness) is the assumption that every point  $x$  of a structure (in the general sense of the term) satisfies a unique location of the structure in the narrow sense of the term, structuralism is equivalent to the condition that the automorphism group of the relevant structure (in the general sense of the term) is trivial. This constitutes a nice formal characterisation of structuralism in consciousness science:

**(STR)** Structuralism about a domain is true iff the automorphism group of that domain is trivial.

<sup>8</sup> It is ‘trivial’ because that’s the simplest possible case, and the set of automorphisms is a group because automorphisms can be combined and inverted as required by the axioms of a group in mathematics.

<sup>9</sup> For every point to have a unique location in a structure is for there not to exist a permutation or other mapping of the domains of that structure to themselves that leaves the structure in the narrow sense of the term invariant.



**Fig. 2. Structural vs. structuralist approaches.** Structural approaches make use of mathematical spaces or mathematical structures to represent or describe conscious experiences. These spaces and structures may, and in general do, admit for automorphisms (cf. Fig. 1). This implies that there are points in the space which have the exact same relational structure. Structuralist approaches, on the other hand, assume that all points of the space can be individuated by their relational structure, meaning that no two points have the same relational structure. This can only hold true if the space does not admit automorphisms, other than the identity mapping that is always an automorphism.

Here, the domain could comprise individual experiences, phenomenal characters, qualities or qualia, depending on which type of structuralism is under consideration.<sup>10</sup>

The crucial point of this section is that mathematical structures can, but need not, obey structuralist assumptions; they may or may not have a trivial automorphism group. A theory or experiment can be *structural*, in the sense that it makes use of mathematical spaces or structures, without necessarily being structuralist. This is illustrated by Fig. 2.

In fact, if we look at mathematical spaces in mathematics, physics, and other natural sciences, in the majority of cases, the automorphism group is *not* trivial. Simple examples of spaces with non-trivial automorphism groups are the Euclidean spaces  $\mathbb{R}^2$ ,  $\mathbb{R}^3$  and  $\mathbb{R}^n$  for any  $n \geq 2$ , and many metric spaces, Riemannian manifolds, Hilbert spaces, or graphs.

Therefore, not only is there a difference between structural and structuralist approaches, but it is in fact quite common that the former applies while the latter does not. Structures in the general sense of the term may, but often do not, amount to structures in the narrow sense of the term. This has three consequences for research in a structuralist turn.

#### Consequence 1. Structural vs. structuralist agendas

Much like the two senses of the term ‘structure’ at issue here are often conflated, so are structural and structuralist agendas. Both are subsumed under the general heading of ‘structuralism’, for example. A first consequence of the above is that there is a difference between structural and structuralist agendas, and it is important to be clear about which agenda one is pursuing when engaging in structuralist research.

If one is using mathematical tools and methods, for example, to help place “structural phenomenal properties at the core of the science of consciousness” (Chalmers, 2023), as required by a very attractive position called weak methodological structuralism that has recently been put forward by David Chalmers, then one is engaging in a structural agenda: an agenda which makes use of mathematical spaces and mathematical structure but which is not committed to a structuralist claim. Put differently, structural tools like mathematical spaces can also be employed if one rejects the idea that structure (understood in the narrow sense of the term) is all that matters. They are free of explanatory and epistemic charge.

#### Consequence 2. Metaphysics of the mind

Many structuralist approaches are not metaphysically neutral. They imply that certain properties which some consider crucial with respect to consciousness do not exist, or are not knowable. For example, if ontic phenomenal structuralism is true, then there are no intrinsic phenomenal properties, and no genuinely private properties. Ontic structural realism implies that there are no qualia as conventionally understood (Dennett, 1988). If epistemic phenomenal structuralism is true, then one cannot know (either scientifically, or by introspection) of intrinsic or private properties, all we know about in regard to conscious experiences derives from structural properties.

Structural approaches are not tied to these assumptions. They are perfectly compatible with the existence of intrinsic or private properties. As far as the mathematics is concerned, if private or intrinsic properties exist (or if there are properties which are not accessible to structural cognitive processing), this simply means that the automorphism group of the structure is not trivial. There are points that cannot be individuated by structure alone.

<sup>10</sup> The term ‘domain’ also has two meanings: The meaning of domain in the sense of mathematical structure as introduced above, and the meaning of domain as a group of related items in general language. Both meanings apply here if it is clear what the structure of a domain is.

To give a very simple example, consider the case where there is no structure in the narrow sense of the term at all, i.e. the case where there are no relations or functions between qualities or qualia at all. This case can be described in terms of mathematics: the qualities or qualia simply form a set. A set is a mathematical structure according to the definition of mathematical structure in mathematical logic. It is the simplest case of a mathematical structure, but an important one. So while this case is opposed to the ideals of structuralist thinking, it is a simple but perfectly fine example of a structural approach.

What is more, structural approaches might actually help to address intrinsic, private or ineffable properties in scientific contexts. My first paper on consciousness, (Kleiner, 2020b), is devoted precisely to this issue. In a nutshell, I show that mathematical tools can be used to formulate theories of consciousness that address these properties even if they are, in an intersubjective sense, non-collatable. Because of these mathematical tools, mathematical approaches allow us to go further than non-mathematical approaches can go. Ultimately, this works because “[m]athematics translates concepts into formalisms and applies those formalisms to derive insights that are usually not amenable to a less formal analysis” (Jost, 2015).

### Consequence 3. Metaphysics beyond the mind

The third consequence, finally, concerns the conviction mentioned at the beginning of this section that structural approaches seem to many to be tied to physicalist or reductionist metaphysics.

The intuition that motivates this conviction arguably derives from the equivocation of structural and structuralist assumptions, together with the idea that science can only explain relations. If structural assumptions would indeed imply that “[t]here is nothing to specifying what something is over and above stating its location in a structure” (Fink et al., 2021), and the physical sciences could only explain structure, then it would indeed be the case that structural approaches would render consciousness amenable to scientific and arguably physicalist explanation. What is more, when ontology is concerned, structuralist assumptions imply that none of the prototypical non-physicalist properties of consciousness exist (cf. Consequence 2). This, too, intuitively speaks in favour of a physicalist and reductionist research programme.

While it is clear that these intuitions do not have the force of a logical argument, it seems fair to say that *structuralist* assumptions are well aligned with physicalist metaphysics, and in the form of one of its most promising incarnations, neuro-phenomenal structuralism (Fink et al., 2021, Lyre, 2022), might even “open an attractive door for reductionism” (Fink et al., 2021).

The problem with the conviction mentioned above is that structural approaches are not necessarily structuralist approaches. The majority of mathematical spaces that are used in the sciences have a non-trivial automorphism group and therefore do not satisfy the defining criterion of a structuralist approach in the context of consciousness science (cf. Fig. 2). In other words, one can choose to apply mathematical tools and methods to describe consciousness without committing to structural assumptions and a fortiori without committing to physicalist or reductionist metaphysics. Structural approaches can be used and might be beneficial in any type of metaphysical programme, from reductive physicalism to property dualism or idealism.

In fact, there are a number of structuralist approaches which target non-physicalist metaphysics already, on the level of toy models. Atmanspacher (2020), for example, uses mathematical tools to outline how the neutral domain in a Pauli-Jung style dual aspect monism might relate to the mental and physical aspects. Other proposals, for example (Signorelli & Wang & Coecke, 2021, Signorelli & Wang & Khan, 2021), use a category theory-based graphical calculus to expand ideas from Buddhist philosophy.

In making these points, I am not arguing for a non-physicalist research programme. My point is that structural approaches are not tied to physicalist or reductionist assumptions. Mathematical spaces and mathematical structures provide descriptive tools that can be applied to any choice of metaphysical assumptions, and in research programmes of any metaphysical flavour. Structural approaches do not have metaphysical premises, and they do not come with a preferred metaphysical interpretation.

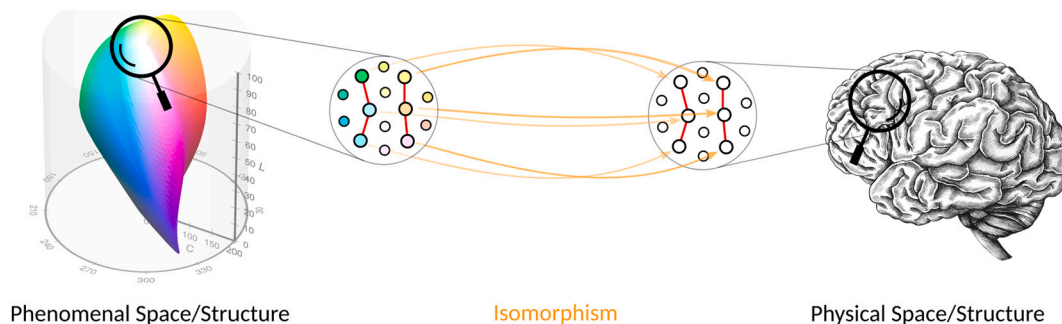
## 4. Isomorphisms and structure-preserving mappings

The core question which drives the scientific study of consciousness is the question of how conscious experiences and ‘the physical’ relate. A ubiquitous mathematical object in the context of mathematical structures is that of an *isomorphism*, illustrated in Fig. 3 and explained in detail below. Due to its ubiquity, when introducing structure to the phenomenal domain, many feel that it is natural to assume that this structure is related to physical structure by an isomorphism or structure-preserving mapping. My goal here is to show that this assumption is not in fact justified. We either need to search for a rigorous justification, or if there is none, proceed in different ways.

Intuitively speaking, an isomorphism expresses a relation between two structures. Precisely speaking, it is a bijective mapping between the *domains* of two structures that preserves the relations or functions of these structures. That is, it is a map from the elements or points of one structure to the elements or points of another structure. A map is bijective if it is one-to-one and onto, meaning that every element in the target space gets mapped to by exactly one element in the source space.

In practice, because the physical has a much larger domain and much richer structure than the phenomenal, when the concept of an isomorphism is applied in consciousness science, what is actually meant is an *isomorphism onto the image*. This means that there is an isomorphism from the phenomenal domain to a substructure of the physical domain. Often, homomorphisms are used as well. They are defined exactly like isomorphisms, except that they do not have to be one-to-one or onto, so that some elements in the target space might not get mapped to, and/or several elements in the source space might map to the same element in the target





**Fig. 3. What is an isomorphism?** This figure illustrates the concept of isomorphisms as applied in consciousness science to link a phenomenal space or structure (left) with a physical space or structure (right). By definition, isomorphisms operate on the level of points. An isomorphism maps every point of the phenomenal space to a point in the physical space. It does so in such a way that the relations between points (indicated here by red lines) are preserved, meaning that any two points which are related on the left are related in the exact same way on the right. The mapping also needs to be invertible. An isomorphism presupposes that structures on both sides of the mapping are given. It does not define, or pick out, the structure in its target domain, which is why it is not a suitable mathematical object to explain, predict, or define phenomenal structure in terms of physical structure. (Depiction of CIE colour space gamut by Wikimedia Commons, Michael Horvath, under CC BY-SA 4.0; this image, excluding the drawing of the brain, is shared under the same license. Drawing of the human brain from Freepik.)

space. Strictly speaking, though, the mathematical concept of homomorphisms is not appropriate either,<sup>11</sup> but to avoid unnecessary technical details, I will admit them too. I will use the term *structure-preserving mapping* to denote homomorphisms or isomorphisms with the understanding that the domains and structures of the source and target have been adapted appropriately to avoid the technical problems. As far as intuition is concerned, my comments are easiest understood when thinking about an isomorphism onto the image.

The assumption under discussion then is:

**(ISO)** The physical and the phenomenal are related by a structure-preserving mapping from the phenomenal domain to the physical domain.<sup>12</sup>

This assumption is a very consequential assumption. It promises, for example, a new methodology for measuring Neural Correlates of Consciousness (NCCs). To date, NCC research has to make use of intricate measures of consciousness (Irvine, 2013), to distinguish between trials where the subject perceives a stimulus consciously from trials where it doesn't. If (ISO) is true, a whole new avenue for investigating NCCs is available: to search, among neural structures in the brain, for structures that are homomorphic to or identical with the structures of the phenomenal domain. This search could, in principle, be carried out independently of any measure of consciousness, and might give a unique result, so that potentially at least there is a methodology where one "[does] not have to worry whether subjects 'really' had a phenomenal experience of a stimulus" (Kob, 2023).

The existence of a structure-preserving mapping between the phenomenal and physical domain also has important consequences for theories of consciousness: it implies that a large class of theories of consciousness is false, namely all those which do not take the form of a homomorphism. A good example of this is Integrated Information Theory (IIT) (Oizumi et al., 2014, Albantakis et al., 2023). It is sometimes assumed that IIT is structure-preserving or even an isomorphism, but according to IIT's mathematical formulation, this is not the case. The mathematics of IIT come with two clear 'slots' for the physical and phenomenal domain. One of the slots is the input to the theory's algorithm. It requires a physical description of a system, for example in terms of neurons. The other slot is the output of the theory's algorithm. For every system and physical state of this system, this output is a mathematical structure called 'Maximally Irreducible Conceptual Structure' in IIT 3.0, and ' $\Phi$ -structure' in IIT 4.0. This structure "is identical to [the system's] experience" (Oizumi et al., 2014). The mathematical algorithm of the theory specifies a mapping between those two

<sup>11</sup> The concept of homomorphism as used in mathematics presumes that two structures have the same *signature*, meaning that both structures need to have the same type of functions or relations: the same number of functions or relations of the same arity, that is. Because the physical has much more structure than the phenomenal (think about the rich structure of electrodynamics in the case of neurons, say), the concept of homomorphism is too strong to express the underlying idea. One could attempt to define a *partial homomorphism* as a homomorphism that respects some, but not all, structures of the target domain. But for questions other than multiple realizability, the 'isomorphism onto the image' conception seems to be closer to the underlying intuition. The same applies if one reverses the direction of the homomorphism, cf. Footnote 12.

<sup>12</sup> In addition to the problem mentioned in Footnote 11, there is also the question of which direction a homomorphism should take. Should it go from the physical domain to the phenomenal domain, as in Fink et al. (2021), or vice versa? Because it is unlikely that all elements of the physical domain are mapped to the phenomenal domain (there are neural mechanisms which are not relevant for conscious experiences, for example), and because a map in the sense of mathematics requires a specification of a target element for every element of the source domain, it seems more natural to me to choose the phenomenal-to-physical direction. Choosing the physical-to-phenomenal direction would require one to introduce yet another sense of partiality, that of a partial function, which is only defined on some of its elements. The problem with this is that a homomorphism which is partial in both this sense and the sense of Footnote 11 always exists, so that the statement (ISO) is vacuous. This is not the case for an isomorphism onto the image in the phenomenal-to-physical direction, because of the need to specify a target element in the physical for every source element in the phenomenal in such a way that the image has the same structure as the phenomenal. This is why I think isomorphisms onto the image in the phenomenal-to-physical direction are the right tool (and the right intuition) to work with, though my comments below do not turn on this choice.

slots which is not a homomorphism. Therefore, the theory does not specify a homomorphism between the physical and phenomenal domains. And consequently, if (ISO) is true then IIT must be wrong.<sup>13</sup>

#### 4.1. Are isomorphisms justified?

The above shows that (ISO) is indeed a very consequential assumption. This would be good news if (ISO) were also a justified assumption. However, as I will argue here, this is not the case. While isomorphisms and homomorphisms are natural in mathematics, they appear not to be the right sort of object to achieve the goals of consciousness science in investigating how the phenomenal and the physical relate. For the purpose of this discussion, I will assume that these goals are “to *explain*, *predict*, [or] *control* the phenomenological properties of conscious experience” (my emphasis) in terms of physical properties, following Anil Seth’s *Real Problem of Consciousness* (Seth, 2021), with the understanding that phenomenal structure is an integral part of phenomenal character, and that structural properties are properties too.

My comments are tied directly to what an isomorphism or homomorphism is. As explained above, isomorphisms and homomorphisms are mappings between the domains of two structures (between the *points* or *elements* of these structures, that is) which satisfy certain conditions. The conditions enforce that the mappings are compatible with the structures on both ends. This has two important consequences for the question at hand.

The first consequence is that a homomorphism *presupposes* that the structures on both ends of the mapping are given. If only one of the two structures is given, or none even, then (ISO) becomes an empty statement. This is because *any* mapping of the form  $f : E \rightarrow P$ , where  $P$  denotes the physical domain and  $E$  denotes the experiential domain, can be turned into a homomorphism if at most one domain comes with structure. One can simply define the structure on the other domain so that the mapping is a homomorphism. Assuming that there is a homomorphism without presupposing that structures on both ends of the mapping are given amounts to not assuming anything at all.

But if a homomorphism presupposes structures on both ends, it doesn’t explain, predict, or allow to control these structures. Homomorphisms fall short of explaining, predicting, or controlling those phenomenal properties they were introduced to cope with.

Second, and more importantly in my opinion, homomorphisms do not have the right mathematical form to *pick out* which structure there is. That is the case because they are maps from domains to domains. They do not actually map from structures to structures, as is sometimes thought. They only map points in one domain to points in another domain in such a way that the mapping between the points *preserves* or *respects* the structure on both ends. This speaks against an explanatory or predictive function as well, as I shall now explain.

Let us first consider the case of explanation. Do homomorphisms, or other structure-preserving mappings, *explain* phenomenal structure in terms of physical structure? There are various notions of explanation that are available in science, ranging from the early deductive-nomological and inductive-statistical ideas studied by Carl Hempel (Hempel & Oppenheim, 1948, Hempel, 1962) to more modern understandings of explanation in the form of causal-mechanical models (Salmon, 1984), unificationist models (Friedman, 1974, Kitcher, 1989), contrastive explanation (Van Fraassen, 1980) or interventionist models (Woodward & Hitchcock, 2003, Hitchcock & Woodward, 2003).

It is clear that homomorphisms do not fit the original Hempel models of explanation because they do not derive phenomenal structure in any meaningful sense from a general law and initial conditions. What is crucial though is that they also don’t sit well with the other models of explanation. This is the case because, in one form or another, these models all require ‘what if things had been different’ information. In the causal-mechanical model of explanation, ‘what if things had been different’ information is required to test the robustness of a purported causal mechanism. In unificationist models it matters for questions of breadth of a unifying explanation. In contrastive explanations it is central to deal with alternative scenarios that would have occurred under different conditions. And in interventionist models, it is required to explicate how an intervention changes the explanandum variable.

Homomorphisms do not pick out structure on the physical or phenomenal side, they only relate points of the domains in a structure-preserving way. Therefore, they do not provide ‘what if things had been different’ information about phenomenal structure. But ‘what if things had been different’ information is required by the above-mentioned models of explanations. Therefore, homomorphisms do not constitute an explanation of phenomenal structure according to these models.

Because homomorphisms don’t pick out phenomenal structure, they do not offer alternatives to how phenomenal structure could be if things had been different. For this reason, they also do not predict phenomenal structure. Prediction, too, requires mathematical tools that pick out the right structure among a class of possible structures.

A helpful way to think about the problems of explanation and prediction is to think about what would *define* phenomenal structure in terms of neural structure, or physical structure more generally. Consider, as an analogy, computer games. Computer games employ

<sup>13</sup> The only way to enforce viewing IIT as an isomorphism is by claiming that the output of IIT’s algorithm is itself a physical structure, which then happens to be related by an isomorphism to the phenomenal domain. Given the interpretation of the mathematical structure outputted by IIT as “identical to [the system’s] experience” (Oizumi et al., 2014), it is hard to see how such interpretation can plausibly be made. The mathematical quantities outputted by IIT do not appear anywhere else in the physical sciences, and are conceptually and mathematically rather removed from physical theories. Such a claim also violates the implicit presupposition in (ISO) that there are more or less well-defined structures on both the phenomenal and physical sides. If there were no constraints on which structure to consider, then (ISO) would be a vacuous statement. Any mapping of the form  $f : P \rightarrow E$ , where  $P$  denotes physical structure and  $E$  denotes phenomenal structure, can be turned into an homomorphism between the physical and the phenomenal if  $E$  is taken to be a physical structure as well. As a rule of thumb, if a structure is actively defined by a theory of consciousness, rather than just adapted from some other part of science, it should probably not count as physical structure in the sense required by (ISO).

mathematical structure to model rich and detailed visual imagery. Yet the mathematical models are defined mostly in terms of objects in the sense of object-oriented programming. There is nothing in the actual code of the game which resembles the structure of the visual scene; rather, the code defines how the structure should be rendered, and it does so in terms of objects and properties. The visual structure created by the game is not homomorphic to the code that runs in order to create the scenes, yet it is defined by the code. This example illustrates that homomorphisms are not the kind of thing one would expect when defining structure.

What these points illustrate, in my view, is that homomorphisms and structure-preserving mappings more generally are not the right sort of object to define, explain, predict, or control phenomenal structure. They might be natural in the context of mathematical questions, but they are not natural for the purposes of consciousness science.

Consequently, (ISO) is not in fact a natural or justified assumption. We either need to search for a rigorous justification, or if there is none, proceed in different ways. Because (ISO) is so consequential for theoretical and experimental work, using (ISO) without proper justification, or in the hope that a justification will eventually be found, is not a viable option.

This comment also applies to mathematical objects known under different names, if these objects are in fact homomorphisms. Important examples thereof are diffeomorphisms, which are maps between smooth geometric shapes called manifolds. Diffeomorphisms are homomorphisms between the mathematical structures that define smooth manifolds. And much like the simpler cases discussed above, they map points of one manifold to points on another manifold in a way that respects the mathematical structure on both sides of the map. They do not explain or define the structure.

#### 4.2. What, if not isomorphisms?

If isomorphisms and homomorphisms are not the sort of thing that explains, predicts, or defines phenomenal structure, what is? Which mathematical objects should we use to relate the physical and the phenomenal in a structuralist turn?

My view is that there is no general mathematical principle that we can commit to. Rather, much like theories of consciousness in the pre-structural area were built one-by-one, we have to build structural theories one-by-one, working with different ideas, concepts, motivations and metaphysics in each case. The challenge of finding the right mathematics to explicate these ideas, concepts and motivations in a structural context is not something we can bypass by choosing one mathematical tool that fits them all. This is not technically possible, nor is it desirable. The difference between ideas, concepts and metaphysical underpinnings in a structural context is precisely in the mathematics that relate the physical to phenomenal structure. We cannot waive the problem of finding the right mathematics without also waiving the possibility of choosing different metaphysical or conceptual ideas.

### 5. Which phenomenal structure?

My final comment concerns the question of which structure to consider when embarking on structural research. That is the question of what phenomenal structure *is* and how we find it. This question is important because conscious experience does not “come with” mathematical structure in any direct sense. There is nothing in what it is like to experience something that is per se mathematically structured, other than if one explicitly experiences something mathematical.<sup>14</sup>

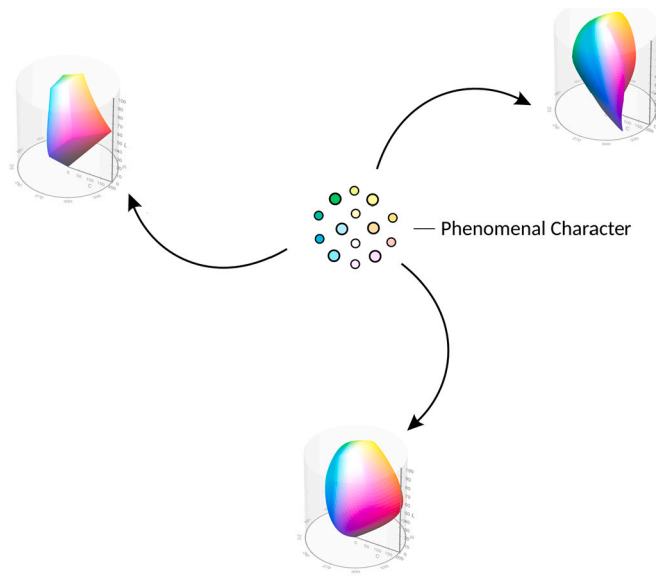
Rather, mathematical spaces and mathematical structures are *tools* or *languages* we can use to describe (or model) phenomenal character, much like English or any other language can be used to describe phenomenal character. And just as we need definitions or conventions to apply English language terms, we need definitions or conventions to apply mathematical terms. These might not be as simple as in the case of English, but still they flesh out the conditions under which one is, and under which one is not, justified in making a structural and mathematical claim.

Because mathematics is a different type of language from English, the definitions or conventions to apply structural terminology are of a different type too. They constitute *methodologies*, meaning they are collections of methods, procedures or rules, that can and need to be used to assess mathematical claims.

And because phenomenal character does not “come with” mathematical structure in any direct sense, any claim about a structural fact, and any application of structural ideas, is always *relative* to a specific understanding of what phenomenal structure is, and a fortiori, relative to the methodology that defines this particular understanding. It is not meaningful to claim that experiences have a certain structure. Much like a claim about whether experiences have qualia depends on what exactly one takes the term qualia to denote, the claim that experiences have a certain structure depends on what one takes phenomenal structure to denote (Fig. 4). When working with or thinking about phenomenal structure, we need to be clear about which methodology we presume. Otherwise, we are prone to making errors. This is the first major point I would like to make in this comment.<sup>15</sup>

<sup>14</sup> We do experience mathematical structures if we know and recognize them, for example in the case of geometrical shapes, or if we actually work with mathematical structures. But we do not experience non-mathematical experiences as mathematically structured. We do not, for example, experience colours as constituting a metric space or having a partial order.

<sup>15</sup> Therefore, working with mathematical structure in consciousness science is different from working with mathematical structure physics or other natural sciences. In physics and other natural science, we do not have direct access to the phenomena we are studying. In a certain sense, for structural claims in physics, *anything goes*, as long as the relevant notion of measurement for that structure reproduces what is observed. This is why there are hugely different proposals about the structure of spacetime, for example, ranging from quantized spacetime (Rovelli, 2004) and emergent spacetime (Koch & Murugan, 2012) to proposals that depart completely from what we intuitively think spacetime should be (Finster & Kleiner, 2015). As long as limiting processes exist that relate these proposals to previous models, in this case the notion of spacetime of General Relativity, all those proposals are viable options. This is not the case for consciousness, because consciousness has a different epistemic context. For example, it exhibits what is sometimes called *epistemic asymmetry*: there are “two fundamentally different methodological approaches that



**Fig. 4. Different definitions imply different spaces.** Mathematical spaces and mathematical structures are tools to describe or represent phenomenal character, much like technical language terms are too. Different definitions or conventions of how to use mathematical terms to describe or represent phenomenal character—different conventions of what terms like ‘mathematical structure of conscious experience’ or ‘phenomenal space’ mean—lead to different structural representations of the same set of experiences, here illustrated by three different CIE colour spaces. Black arrows indicate different definitions or conventions, which imply different methodologies for constructing phenomenal spaces in the lab. Different geometrical shapes indicate different types of spaces that result from applying these methodologies. Much like technical language terms might differ in scope, quality, adequacy, and presuppositions, definitions or conventions regarding mathematical structures differ in scope, quality, adequacy and presuppositions. (Depiction of CIE colour space gamuts by Wikimedia Commons, Michael Horvath, under [CC BY-SA 4.0](https://creativecommons.org/licenses/by-sa/4.0/); this image is shared under the same license.)

### 5.1. What is phenomenal structure, and how do we find it?

There are three important landmarks that have influenced the way in which we use mathematical structures to describe conscious experiences today: quality spaces as introduced by Austen Clark (Clark, 1993), quality spaces as introduced by David Rosenthal (Rosenthal, 1991, 2010) and  $Q$ -spaces as introduced in IIT 2.0 (Tononi, 2008). While these methodologies have served an important function in enabling structural research, it is also important to be clear about their shortcomings in moving forward.

As far as IIT is concerned, the obvious shortcoming is that the theory does not provide a phenomenal interpretation of the structure it proposes, other than the claim that the structure “is identical to [the system’s] experience” (Oizumi et al., 2014). This gives rise to what David Chalmers has called the *Rosetta Stone Problem* (Chalmers, 2023): the problem of how to translate the mathematical structure that IIT proposes into phenomenological terms. IIT does not actually specify a methodology that clarifies how to interpret and test their proposed structure in phenomenal terms.

The proposals by Clark and Rosenthal do specify methodologies. The major shortcoming of these methodologies, on my view, is that they conflate three sources of mathematical structure:

1. **Mathematical Convenience.** Some of the structure is introduced simply for mathematical convenience.
2. **Laboratory Operations.** Some of the mathematical structure refers to, or depends on, laboratory operations.
3. **Conscious Experience.** Only part of the mathematical structure actually pertains to conscious experiences or phenomenal character.

### 5.2. Clark’s quality spaces

Quality spaces as introduced by Austen Clark (Clark, 1993) are based on the following methodology. To construct the quality space for an individual subject,<sup>16</sup> one fixes a class of stimuli  $S$  that can be presented to the subject, and defines two tasks that the subject can complete in response to the presentation of one or more stimuli. The first task probes whether the subject is able to

enable us to gather knowledge about consciousness: we can approach it from within and from without; from the first-person perspective and from the third-person perspective. Consciousness seems to distinguish itself by the privileged access that its bearer has to it” (Metzinger, 1995). In other words, in addition to the usual scientific way of accessing and modelling a phenomenon there is a second way of accessing the phenomenon (described in terms of the first person perspective *metaphor* above). Due to this different epistemic context, using mathematical structure to describe a phenomenon is different in the case of consciousness, and more constrained, than in the case of physics.

<sup>16</sup> Clark mostly has humans in mind, but does consider the case of animals briefly in Clark (1993). Nothing hinges on humans in the methodology he proposes.

*discriminate* the experience elicited by two different stimuli consciously. The second task probes whether the subject experiences a stimulus to be more similar to a reference stimulus than another stimulus. This is called *relative similarity*.<sup>17</sup>

The discrimination task is used to define a *global indiscriminability* relation on the class of stimuli  $S$ .<sup>18</sup> While discriminability does not constitute an equivalence relation, global indiscriminability does. This equivalence relation partitions the set of stimuli. Each set in this partition contains stimuli which are globally indiscriminable from each other, and defines a *quality* in Clark's proposal. The collection of the sets in this partition (the space of equivalence classes of  $S$ , in mathematical terms) defines the domain of the quality space that is being constructed.

The relative similarity task is used to define a graph, in the mathematical sense of the term, between the qualities: a set of nodes, and edges that link some of the nodes. Working with stimuli that represent the different qualities, one first collects relative similarity data. This is data about whether a quality  $q_1$  is more similar to a reference quality  $q_0$  than another quality  $q_2$ . One might find that the pair  $(q_1, q_0)$  is more similar to each other than the pair  $(q_2, q_0)$ , say. Having collected this data for all qualities in the set, one then represents them as a graph. Every quality one has previously constructed is a node of the graph, and every pair  $(q_i, q_j)$  about which one has relative similarity data is an edge of the graph between the nodes that represent the qualities. The important part then is that the edges get labelled by numbers, and these numbers must be chosen in such a way that the relative similarity judgements that have been collected are represented truthfully by the ordering of the numbers. The label of the edge  $(q_1, q_0)$  above, for example, must be a lower number than the label of the edge  $(q_2, q_0)$  if the former pair is more similar to each other than the latter pair. The result of this procedure is a labelled graph, where the nodes represent qualities, edges indicate pairs for which similarity data is available, and labels on the edges represent relative similarity. In mathematical terms, this is called a POSET-labelled graph, where a POSET is a partially ordered set. The partial order is the phenomenal structure of the relative similarity experiences.

Up to this point all the mathematical structure is still grounded in conscious experience, to a large extent. The data to carry out the constructions is based on tasks that might utilize reports or behavioural measures, but these measures depend on what is experienced.

The next step in Clark's methodology consists of introducing a metric, a tool to measure distances in terms of continuous numbers, and in fact an Euclidean space that has a uniform, homogeneous metric. To this end, it makes use of a procedure known as 'multidimensional scaling' (Beals et al., 1968). In Clark's case, it consists of finding an *embedding* of the graph into an Euclidean metric space in such a way that the distance between the nodes of the graph—which are mapped to points in the metric space—reproduce the ordering of relative similarity that the labels of the graph encode.

From the perspective of phenomenal character, this step is unwarranted. Not only is the metric introduced without any reference to experience, but this step also leads to the introduction of many more points besides the original qualities that were carefully constructed making use of global indiscriminability. Technically speaking, it leads to an infinity of additional points, all of which feature in the metric function of the space, and none of which is any different from the points that were carefully constructed based on tasks and stimuli.

The only justification I can think of why one would make use of this last step, as compared to just working with the POSET-labelled graph, is mathematical convenience. A POSET-labelled graph might just be too unfamiliar a mathematical object. Or maybe the reason is that it cannot easily be further analysed on a computer in familiar ways. These justifications are in fact made explicit in introductory texts on psychophysics. Luce and Suppes, for example, speak of representational measurement, of which multidimensional scaling is an example, as "an attempt to understand the nature of empirical observations (...) in terms of *familiar* mathematical structures" (Luce & Suppes, 2004, p. 1) (my emphasis), and add that "the use of such empirical structures in psychology is widespread because they come close to the way data are organised for subsequent statistical analysis" (Luce & Suppes, 2004, p. 2). Be that as it may, the last step that introduces the metric function fails to be grounded in conscious experience. It is an example of 1. above.

### 5.3. Rosenthal's quality spaces

The construction of quality spaces as defined by David Rosenthal is based on a class of stimuli as well. But in this case, one only needs a discrimination task, as well as means to *vary* the stimuli.

The main step in Rosenthal's methodology is to construct *Just Noticeable Differences* (JNDs) from variations of the stimuli and the discrimination task. To this end, one varies a stimulus in some direction until the subject notices the difference between the stimulus and the variation. The class of stimuli which one can reach by varying one stimulus without creating a JND gives a set or region in stimulus space, and much like in the case of Clark, the idea is that these regions constitute qualities. A metric function is introduced on the set of qualities by counting the minimal number of regions one has to pass so as to go from one quality to the other.

In this proposal too, there is a question as to the experiential source of the metric function. Because the metric function can be specified, once JNDs have been constructed, without need of additional data, it might not represent anything over and above the JNDs and their neighbourhood relations. Furthermore, while we do experience colour qualities as instantiating a relative similarity

<sup>17</sup> There is considerable freedom in which class of stimuli to choose and how to define and implement the tasks.

<sup>18</sup> Two stimuli are globally indiscriminable if and only if the following two conditions hold:

1. The two stimuli are indiscriminable from each other.
2. The two stimuli have identical indiscriminability relations to all other stimuli in  $S$ .

structure, we do not experience qualities to be a certain number of steps apart, as a metric would require if it indeed represented a structure of conscious experience.<sup>19</sup> So there is a worry of the metric being due to mathematical convenience here too.

A more fundamental worry in this case concerns the *variations* of stimuli that one needs in order to construct JNDs and their neighbourhood relations in the first place. The idea of a variation—starting with one stimulus and then changing that stimulus continuously until a subject notices a difference—requires a *topology* on the stimulus space. A topology defines what it means to “draw a line without lifting a pen” on an abstract space, so to speak. It is precisely what provides the notion of continuous curves required to specify variations in Rosenthal’s definition. Without a topology, a variation can jump from any point to any other point.

The problem is that different topologies give different variations. So when one actually constructs a quality space according to Rosenthal’s methodology in the lab, the resulting space depends on the topology of the stimulus space that has been used. And much like there isn’t just a single notion of colour space, there isn’t just a single topology on colour stimuli one can use. As a result, the metric function that one constructs in an application of Rosenthal’s methodology actually depends on the topology that has been chosen in the experiment, which is a laboratory operation in the sense of 2. above.

In the case of Rosenthal’s methodology, there is in fact a theory that can be used to answer these and similar worries, a theory about what consciousness is, about how qualities should be understood, and about how consciousness and qualities relate. When I asked David Rosenthal about the problem regarding topology, for example, he countered by assuming that there is just one actual physical topology in reality and that this is the topology that should be used. It is not clear to me how this would work in practice, given that this topology is presumably defined by Quantum Electrodynamics (QED), and too far removed from experimental practice to be applicable; in the lab, some choice of topology will have to be made nonetheless. But theoretically speaking, the answer is fully valid. Similarly, the theory about what qualities are and how they relate to consciousness discharges the methodology from the problem that, according to the subsumed notion of discrimination in this case, discriminations could also be made unconsciously.

There is, however, no free lunch. The price to be paid for solving methodological problems by theoretical assumptions is that the methodology now depends on these assumptions and cannot be used to formulate or test other theories of consciousness. The methodological tool might be deprived of much of the impact it could otherwise have.

In my view, quality spaces are ways to describe or represent the explanandum—what is to be explained: qualitative or phenomenal character, what it is like to be—, while theories of consciousness are the explanantia—they do the explaining. This is why I have always been tempted to read Rosenthal’s proposal as a general methodology that is independent from his theory. This is possible and addressing the above-mentioned problems on purely methodological grounds leads, in my view, to fruitful further developments of his construction (cf. below and Kleiner and Ludwig (2023)).

#### 5.4. How to move forward

In the last two sections, I have analysed two proposals for methodologies that define what quality spaces are. While these proposals have served an important role in enabling structural thinking, much of the essential structure in these proposals is not actually grounded in conscious experiences, but in mathematical convenience and laboratory operations.

It is possible to go beyond individual methodologies and analyse the *type of condition* that is applied in these proposals and more recent work. That is, the type of condition that decides whether a mathematical structure is a quality space or phenomenal space—a *mathematical structure of conscious experience*, to use a general term. In a nutshell, all existing proposals I know of amount to:

- (A) Conditions on the domains (sets of points) of a mathematical structure, formulated in terms of qualities, qualia, phenomenal properties or similar aspects of conscious experiences.
- (B) The requirement that the mathematical axioms of the structure (such as the axiom that the metric distance between a point and itself is zero) are satisfied.

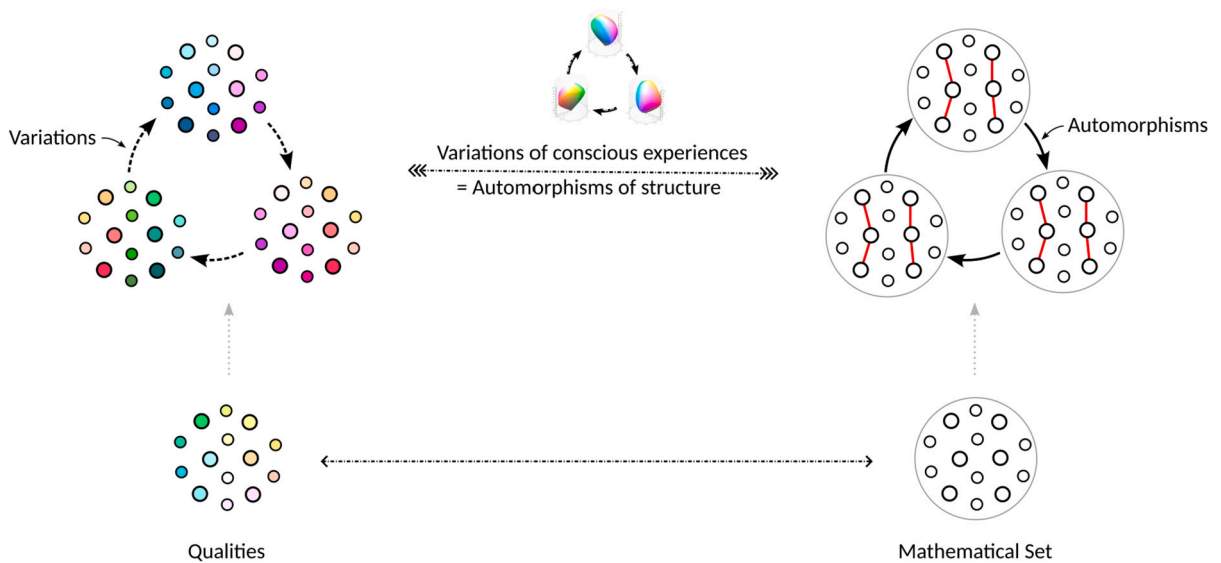
This type of condition can be shown to be insufficient to ground a thorough understanding of phenomenal structure. This is the case because (a) it is prone to admitting incompatible structures, (b) allows for *arbitrary* re-definitions of structures that still satisfy the condition, and (c) in a subtle but important sense, the condition is indifferent to structural facts of conscious experience. I do not have the space here to explain these problems in detail; they are explained and illustrated in (Kleiner & Ludwig, 2023, Section 1).

I take the problems of existing proposals, and the insufficiency of the general type of condition that is applied, to constitute a need of constructing a *new methodology* for phenomenal spaces. This methodology needs to take previous methodologies into account, but needs to amend and extend them to avoid the three insufficiency problems as well as the issues with non-conscious sources of the mathematical structure.

In Kleiner and Ludwig (2023), Tim Ludwig and I have set out to find a methodology that achieves this task. The result is illustrated in Fig. 5. The proposal shares with David Rosenthal’s methodology that it rests on variations, though in our case, any transition from one conscious experience to another counts as variation, and we do not demand continuity or restrict only to variations of stimuli.

Put in terms of phenomenal properties, the core intuition of our proposal is that a mathematical structure is a mathematical structure of conscious experience—a phenomenal space, to use a simpler term—if and only if there is a phenomenal property that behaves exactly as the mathematical structure does under variations. If a variation preserves the mathematical structure (if it is an

<sup>19</sup> For a more careful examination of the case of a metric, cf. Kleiner and Ludwig (2023). For questions on how quality spaces *should* relate to consciousness or phenomenal character according to the underlying theory, cf. below.



**Fig. 5. How to define phenomenal spaces?** This figure illustrates how to define phenomenal spaces and other mathematical structures of conscious experience. One starts out with a choice of qualities (bottom left), for example colour qualities, sometimes also called qualia or conceptualised as instantiated phenomenal properties. The qualities form a set that constitutes the points of the phenomenal space (bottom right). Every experience comprises a subset of qualities, and as experiences change from one experience to the next, the subset of qualities that is realised varies (top left). These variations can be understood as mappings from the set of qualities to itself, and therefore have the same formal structure as automorphisms (Fig. 1): mappings from the points of a space to other points of the space (top right). This allows for the following simple definition of phenomenal structure: phenomenal structure is that mathematical structure whose automorphisms are identical to the variations of the qualities as experiences change. Put differently, phenomenal structure (indicated here by red lines) is that mathematical structure which renders the statement true: the variations of (qualities of) conscious experiences are the automorphisms of the structure (top centre). For details, cf. (Kleiner & Ludwig, 2023). (Depiction of CIE colour space gamuts by Wikimedia Commons, Michael Horvath, under CC BY-SA 4.0; this image is shared under the same license.)

automorphism of the structure, in mathematical terms), then it must not change the phenomenal property. If, conversely, a variation does not preserve the mathematical structure, then it must change the phenomenal property. In a nutshell: there is something “in” conscious experience (the phenomenal property) that behaves exactly as the mathematical structure does.

## 6. Conclusion

Structural approaches, which make use of mathematical structure to describe or model conscious experiences, offer new and valuable avenues for studying consciousness. My aim in this paper is to provide three comments that I consider important when engaging in structural research. Each comment targets what is, in my view, a misconception or misunderstanding that I aim to clarify.

My first comment focuses on the metaphysical underpinnings of structural approaches. I show that, contrary to popular belief, structural approaches are not tied to physicalist or reductive metaphysics. Instead, they offer versatile descriptive tools that can be utilised irrespective of one’s metaphysical commitments, across research programmes of any metaphysical flavour.

My second comment concerns isomorphisms and structure-preserving mappings. A number of emerging structuralist research programmes rely on assuming a structure-preserving mapping between the phenomenal and the physical domain. I argue that this assumption is unwarranted, and that isomorphisms and structure-preserving mappings are not the right mathematical object to provide explanations, predictions, or definitions of phenomenal structure. Instead, we should direct our attention to structural theories of consciousness, without expecting a single mathematical formalism to fit them all. One major experimental consequence of this is that methods such as Representational Similarity Analysis (Kriegeskorte et al., 2008), which searches for structural similarity, may not be the right approach to search for the neural correlates of phenomenal structure.

My third and final comment focuses on the question of what phenomenal structure is, and how we find it. Conscious experiences do not “come with” mathematical structure in any meaningful sense. Rather, mathematical spaces and mathematical structure offer a language to describe or represent conscious experiences, and just like we need definitions or conventions to apply English language terms to consciousness, we need definitions or conventions to apply structural terms. In the case of structure, the definitions and conventions take the form of methodologies that govern how to construct or use the mathematical terminology. The two major methodologies that have guided recent developments are quality spaces as introduced by Austen Clark, and quality spaces as introduced by David Rosenthal. I show that both suffer from fundamental issues, and discuss how to move forward in light of this.

## CRedit authorship contribution statement

**Johannes Kleiner:** Conceptualization, Formal analysis, Funding acquisition, Investigation, Methodology, Project administration, Resources, Software, Visualization, Writing – original draft, Writing – review & editing.

## Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Data availability

No data was used for the research described in the article.

## Acknowledgements

I would like to thank the organisers and participants of the 2023 *Structuralism in Consciousness Studies* workshop at the Charité Berlin for many stimulating discussions on this topic, in particular Lukas Kob, Lucia Melloni, Sascha Benjamin Fink, Holger Lyre, David Chalmers, Andrew Lee, and Wanja Wiese, as well as the participants of a recent *NYU Philosophy of Mind Discussion Group* for valuable advice, feedback and discussions of an earlier version of this manuscript. Furthermore, I would like to thank Wanja Wiese, Matthias Michel, and Moritz Nicolas Loerbroks for feedback on an earlier version of this manuscript. This research was supported by grant number FQXi-RFP-CPW-2018 from the Foundational Questions Institute and Fetzer Franklin Fund, a donor advised fund of the Silicon Valley Community Foundation. I would like to thank the Mathematical Institute of the University of Oxford and the NYU Center for Mind, Brain, and Consciousness for hosting me while working on this article.

## Funding

This research was supported by grant number FQXi-RFP-CPW-2018 from the Foundational Questions Institute and Fetzer Franklin Fund, a donor advised fund of the Silicon Valley Community Foundation.

## References

- Albantakis, L., Barbosa, L., Findlay, G., Grasso, M., Haun, A. M., Marshall, W., Mayner, W. G., Zaeemzadeh, A., Boly, M., Juel, B. E., et al. (2023). Integrated information theory (iit) 4.0: Formulating the properties of phenomenal existence in physical terms. *PLoS Computational Biology*, 19(10), Article e1011465.
- Atmanspacher, H. (2020). The Pauli–Jung conjecture and its relatives: A formally augmented outline. *Open Philosophy*, 3(1), 527–549.
- Beals, R., Krantz, D. H., & Tversky, A. (1968). Foundations of multidimensional scaling. *Psychological Review*, 75(2), 127.
- Block, N. (1990). Inverted Earth. *Philosophical Perspectives*, 4, 53–79.
- Blum, L., & Blum, M. (2022). A theory of consciousness from a theoretical computer science perspective: Insights from the conscious Turing machine. *Proceedings of the National Academy of Sciences*, 119(21), Article e2115934119.
- Chalmers, D. (2023). *Phenomenal structuralism* In *Talk presented at the structuralism in consciousness studies workshop at the Charité Berlin*.
- Chalmers, D. J., & McQueen, K. J. (2022). Consciousness and the collapse of the wave function. In S. Gao (Ed.), *Consciousness and quantum mechanics*. Oxford University Press.
- Clark, A. (1993). *Sensory qualities*. Clarendon library of logic and philosophy.
- Clark, A. (2000). *A theory of sentience*. Clarendon Press.
- Coninx, S. (2022). A multidimensional phenomenal space for pain: Structure, primitiveness, and utility. *Phenomenology and the Cognitive Sciences*, 21(1), 223–243.
- Dennett, D. C. (1988). Quining qualia. In *Consciousness in contemporary science* (pp. 42–77).
- Fink, S. B., Kob, L., & Lyre, H. (2021). A structural constraint on neural correlates of consciousness. *Philosophy and the Mind Sciences*, 2.
- Finster, F., & Kleiner, J. (2015). *Causal fermion systems as a candidate for a unified physical theory*. *Journal of Physics: Conference Series*, 626, 012020. IOP Publishing.
- Fortier-Davy, M., & Millièrè, R. (2020). The multi-dimensional approach to drug-induced states: A commentary on Bayne and Carter’s “dimensions of consciousness and the psychedelic state”. *Neuroscience of Consciousness*, 2020(1), Article niaa004.
- Friedman, M. (1974). Explanation and scientific understanding. *The Journal of Philosophy*, 71(1), 5–19.
- Gert, J. (2017). Quality spaces: Mental and physical. *Philosophical Psychology*, 30(5), 525–544.
- Grindrod, P. (2018). On human consciousness: A mathematical perspective. *Network Neuroscience*, 2(1), 23–40.
- Haun, A., & Tononi, G. (2019). Why does space feel the way it does? Towards a principled account of spatial experience. *Entropy*, 21(12), 1160.
- Haynes, J.-D. (2009). Decoding visual consciousness from human brain signals. *Trends in Cognitive Sciences*, 13(5), 194–202.
- Hebart, M. N., Zheng, C. Y., Pereira, F., & Baker, C. I. (2020). Revealing the multidimensional mental representations of natural objects underlying human similarity judgements. *Nature Human Behaviour*, 4(11), 1173–1185.
- Hempel, C. G. (1962). Deductive-nomological vs. statistical explanation. In *Scientific explanation, space, and time*. Minneapolis: University of Minnesota Press.
- Hempel, C. G., & Oppenheim, P. (1948). Studies in the logic of explanation. *Philosophy of Science*, 15(2), 135–175.
- Hitchcock, C., & Woodward, J. (2003). Explanatory generalizations, part II: Plumbing explanatory depth. *Noûs*, 37(2), 181–199.
- Hoffman, D. D., & Prakash, C. (2014). Objects of consciousness. *Frontiers in Psychology*, 5, 577.
- Hoffman, D. D., Prakash, C., & Prentner, R. (2023). Fusions of consciousness. *Entropy*, 25(1), 129.
- Irvine, E. (2013). Measures of consciousness. *Philosophy Compass*, 8(3), 285–297.
- Jackson, F. (1986). What Mary didn’t know. *The Journal of Philosophy*, 83(5), 291–295.
- Jackson, F. (1998). Epiphenomenal qualia. In *Consciousness and emotion in cognitive science* (pp. 197–206). Routledge.
- Josephs, E. L., Hebart, M. N., & Konkle, T. (2023). Dimensions underlying human understanding of the reachable world. *Cognition*, 234, Article 105368.
- Jost, J. (2015). *Mathematical concepts*. Springer.
- Kawakita, G., Zeleznikow-Johnston, A., Takeda, K., Tsuchiya, N., & Oizumi, M. (2023). *Is my “red” your “red”? Unsupervised alignment of qualia structures via optimal transport*. PsyArXiv preprint.
- Kawakita, G., Zeleznikow-Johnston, A., Tsuchiya, N., & Oizumi, M. (2023). Comparing color similarity structures between humans and LLMs via unsupervised alignment. ArXiv preprint arXiv:2308.04381.
- Kitcher, P. (1989). Explanatory unification and the causal structure of the world. In *Scientific explanation. Minnesota studies in the philosophy of science*. Minneapolis: University of Minnesota Press.



- Kleiner, J. (2020a). Brain states matter. A reply to the unfolding argument. *Consciousness and Cognition*, 85, Article 102981.
- Kleiner, J. (2020b). Mathematical models of consciousness. *Entropy*, 22(6), 609.
- Kleiner, J., & Hoel, E. (2021). Falsification and consciousness. *Neuroscience of Consciousness*, 2021(1), Article niab001.
- Kleiner, J., & Ludwig, T. (2023). What is a mathematical structure of conscious experience? *Synthese*. In press.
- Kleiner, J., & Tull, S. (2021). The mathematical structure of integrated information theory. *Frontiers in Applied Mathematics and Statistics*, 6, 74.
- Klincewicz, M. (2011). Quality space model of temporal perception. In *Multidisciplinary aspects of time and time perception* (pp. 230–245). Springer.
- Kob, L. (2023). Exploring the role of structuralist methodology in the neuroscience of consciousness: A defense and analysis. *Neuroscience of Consciousness*, 2023(1), Article niad011.
- Koch, C., Massimini, M., Boly, M., & Tononi, G. (2016). Neural correlates of consciousness: Progress and problems. *Nature Reviews Neuroscience*, 17(5), 307–321.
- Koch, R., & Murugan, J. (2012). Emergent spacetime. In *Foundations of space and time: Reflections on quantum gravity* (pp. 164–184).
- Kostic, D. (2012). The vagueness constraint and the quality space for pain. *Philosophical Psychology*, 25(6), 929–939.
- Kriegeskorte, N., Mur, M., & Bandettini, P. A. (2008). Representational similarity analysis - connecting the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, 4.
- Lau, H., Michel, M., LeDoux, J. E., & Fleming, S. M. (2022). The mnemonic basis of subjective experience. *Nature Reviews Psychology*, 1(8), 479–488.
- Lee, A. Y. (2021). Modeling mental qualities. *Philosophical Review*, 130(2), 263–298.
- Lee, A. Y. (2022). Objective phenomenology. *Erkenntnis*, 1–20.
- Luce, R. D., & Suppes, P. (2004). Stevens' handbook of experimental psychology. In H. Pashler, & J. Wixted (Eds.), *Stevens' handbook of experimental psychology, methodology in experimental psychology*. John Wiley & Sons.
- Lyre, H. (2022). Neurophenomenal structuralism. A philosophical agenda for a structuralist neuroscience of consciousness. *Neuroscience of Consciousness*, 2022(1), Article niac012.
- Malach, R. (2021). Local neuronal relational structures underlying the contents of human conscious experience. *Neuroscience of Consciousness*, 2021(2), Article niab028.
- Mason, J. W. (2013). Consciousness and the structuring property of typical data. *Complexity*, 18(3), 28–37.
- Mason, J. W. (2021). Model unity and the unity of consciousness: Developments in expected float entropy minimisation. *Entropy*, 23(11), 1444.
- Metzinger, T. (1995). The problem of consciousness. In T. Metzinger (Ed.), *Conscious experience* (pp. 3–37). Imprint Academic.
- Michel, M. (2023). Confidence in consciousness research. *Wiley Interdisciplinary Reviews: Cognitive Science*, 14(2), Article e1628.
- Michel, M. (In press). The perceptual reality monitoring theory. In Herzog, M., Schurger, A., & Doerig, A. (Eds.), *Scientific Theories of Consciousness: The Grand Tour*. Cambridge University Press.
- O'Brien, G., & Opie, J. (1999). A connectionist theory of phenomenal experience. *Behavioral and Brain Sciences*, 22(1), 127–148.
- Oizumi, M., Albantakis, L., & Tononi, G. (2014). From the phenomenology to the mechanisms of consciousness: Integrated Information Theory 3.0. *PLoS Computational Biology*, 10(5), Article e1003588.
- Pashler, H., & Wixted, J. (2004). *Stevens' handbook of experimental psychology, methodology in experimental psychology, vol. 4*. John Wiley & Sons.
- Prenner, R. (2019). Consciousness and topologically structured phenomenal spaces. *Consciousness and Cognition*, 70, 25–38.
- Renner, A. (2014). Consciousness and mental qualities for auditory sensations. *Journal of Consciousness Studies*, 21(9–10), 179–204.
- Resende, P. (2022). Qualia as physical measurements: A mathematical model of qualia and pure concepts. ArXiv preprint arXiv:2203.10602.
- Rosenthal, D. (2010). How to think about mental qualities. *Philosophical Issues*, 20, 368–393.
- Rosenthal, D. (2015). Quality spaces and sensory modalities. In P. Coates, & S. Coleman (Eds.), *Phenomenal qualities: Sense, perception, and consciousness* (pp. 33–65). Oxford, UK: Oxford University Press.
- Rosenthal, D. M. (1991). The independence of consciousness and sensory quality. *Philosophical Issues*, 1, 15–36.
- Rosenthal, D. M. (2016). Quality spaces, relocation, and grain. In O'Shea (Ed.), *Sellars and his legacy* (pp. 149–185). Oxford: Oxford University Press.
- Rovelli, C. (2004). *Quantum gravity*. Cambridge University Press.
- Rudrauf, D., Bennequin, D., Granic, I., Landini, G., Friston, K., & Williford, K. (2017). A mathematical model of embodied consciousness. *Journal of Theoretical Biology*, 428, 106–131.
- Safiron, A. (2022). Integrated world modeling theory expanded: Implications for the future of consciousness. *Frontiers in Computational Neuroscience*, 16, Article 642397.
- Salmon, W. C. (1984). *Scientific explanation and the causal structure of the world*. Princeton University Press.
- Schanda, J. (2007). CIE colorimetry. *Colorimetry: Understanding the CIE System*, 3, 25–78.
- Seth, A. (2021). *Being you: A new science of consciousness*. Penguin.
- Shoemaker, S. (1982). The inverted spectrum. *The Journal of Philosophy*, 79(7), 357–381.
- Signorelli, C. M., Wang, Q., & Coecke, B. (2021). Reasoning about conscious experience with axiomatic and graphical mathematics. *Consciousness and Cognition*, 95, Article 103168.
- Signorelli, C. M., Wang, Q., & Khan, I. (2021). A compositional model of consciousness based on consciousness-only. *Entropy*, 23(3), 308.
- Silva, L. (2023). Towards an affective quality space. *Journal of Consciousness Studies*, 30(7–8), 164–195.
- Stanley, R. P. (1999). Qualia space. *Journal of Consciousness Studies*, 6(1), 49–60.
- Tallon-Baudry, C. (2022). The topological space of subjective experience. *Trends in Cognitive Sciences*.
- Tononi, G. (2008). Consciousness as Integrated Information: A provisional manifesto. *Biological Bulletin*, 215(3), 216–242.
- Tononi, G. (2015). Integrated Information Theory. *Scholarpedia*, 10(1), 4164.
- Tsuchiya, N., Phillips, S., & Saigo, H. (2022). Enriched category as a model of qualia structure based on similarity judgements. *Consciousness and Cognition*, 101, Article 103319.
- Tsuchiya, N., & Saigo, H. (2021). A relational approach to consciousness: Categories of level and contents of consciousness. *Neuroscience of Consciousness*, 2021(2), Article niab034.
- Tsuchiya, N., Saigo, H., & Phillips, S. (2023). An adjunction hypothesis between qualia and reports. *Frontiers in Psychology*, 13, Article 1053977.
- Tsuchiya, N., Taguchi, S., & Saigo, H. (2016). Using category theory to assess the relationship between consciousness and Integrated Information Theory. *Neuroscience Research*, 107, 1–7.
- Van Fraassen, B. C. (1980). *The scientific image*. Oxford University Press.
- Woodward, J., & Hitchcock, C. (2003). Explanatory generalizations, part I: A counterfactual account. *Noûs*, 37(1), 1–24.
- Yaron, I., Melloni, L., Pitts, M., & Mudrik, L. (2021). *The consciousness theories studies (contrast) database: Analyzing and comparing empirical studies of consciousness theories*. bioRxiv.
- Yoshimi, J. (2007). Mathematizing phenomenology. *Phenomenology and the Cognitive Sciences*, 6(3), 271–291.
- Young, B. D., Keller, A., & Rosenthal, D. (2014). Quality-space theory in olfaction. *Frontiers in Psychology*, 5(1).
- Zaidi, Q., Victor, J., McDermott, J., Geffen, M., Bensmaia, S., & Cleland, T. A. (2013). Perceptual spaces: Mathematical structures to neural mechanisms. *The Journal of Neuroscience*, 33(45), 17597–17602.
- Zeleznikow-Johnston, A., Aizawa, Y., Yamada, M., & Tsuchiya, N. (2023). Are color experiences the same across the visual field? *Journal of Cognitive Neuroscience*, 35(4), 509–542.