

Investigating Causal Effects of Mental Events in Cognitive Neuroscience

Mikkel C. Vinding, Ph.D.

NatMEG, Department of Clinical Neuroscience

Karolinska Institutet

Nobels väg 9, D3

171 77 Stockholm

Sweden

Email: mikkel.vinding@ki.se

1 **Abstract**

2 Mental causation is a predominantly theoretical topic rather than a topic studied in the laboratory. The
3 purpose of this paper is to outline a general approach for studying mental causation by empirical means
4 for philosophers and scientists interested in the topic. The aim is to outline how we can infer mental
5 causation by empirical methods given an unknown solution to the mind-body problem. The approach is
6 based on the principles of causal inference to find causal relations among observed variables used in all
7 branches of science. With these principles, it is possible to estimate the causal effects of mental events:
8 Make an experimental manipulation on a mental event, control confounding variables, and estimate
9 causal effects on the outcome. The caveat is that we cannot separate the causal effects of a mental event
10 from the physical base of the mental event, independent of whether we assume mental events can be
11 reduced to their physical base. A challenge to estimating causal effects of mental events is that
12 measured physical variables, such as electrophysiological potentials from the brain, might reflect
13 processes that are part of “higher-order” phenomena, such as mental events. This means that
14 controlling “downwards” for confounding variables is challenging. It is, however, not impossible. It
15 also means that inferring non-mental causes of action cannot be done by measuring only physical
16 variables alone. Keeping the mind-body problem in mind when designing experiments, it is possible to
17 infer mental causation.

18
19 **Keywords:** Mental Causation, Causal Inference, Cognitive Neuroscience, Mental Events, Mind-Body
20 Problem

1 Introduction

2 Do mental events cause physical action, and if so, how? This is the central question in the topic of
3 mental causation. One ought to think that investigating mental causation is a goal of experimental
4 psychology—but experiments directly addressing mental causation are surprisingly sparse. Even
5 though a massive advancement in neuroscience methods has given experimental psychology and
6 neuroscience tools to study the biological basis of mental events, there has been little attempt at a
7 framework for empirical investigations of mental causation. Mental causation is predominantly a
8 theoretical topic with close to no contribution from experimental science.

9 The central problem in addressing mental causation is the mind-body problem. The mind-body
10 problem, in a nutshell, is that we do not know how mental events relates to physical states. Without
11 knowing how they are related, we do not know how they interact. The mind-body problem makes it
12 difficult to define the pre-empirical foundation for experimental inference about mental causation. We
13 assume the brain is the foundation of the mind, but since we do not know how they are related, it is
14 unclear what assumptions we have to make to include mental events when inferring causal relations.
15 The mind-body problem appears to be an obstacle to describing the foundation for investigating mental
16 causation by empirical means.

17 The purpose of this paper is to sketch a foundation for an experimental approach that scientists and
18 empirical oriented philosophers can use to study mental causation. How do we go from the analytical
19 approach to understanding mental causation to instead gain an understanding via experimental
20 inference? Furthermore, how do we distinguish between the type of questions we must deal with by
21 analytical reasoning and those questions we can answer by empirical means?

22 Given the unknown relationship between mental events and the physical world—the mind-body
23 problem—I will explore to what extent we can make meaningful inference about mental causation. As
24 we must be aware of what we can (and cannot) infer from experimental studies on mental causation, I
25 will make the necessary assumptions explicit and acknowledge limitations in the experimental designs.
26 Finally, I will answer the following questions: what type of questions about mental causation can we
27 answer through experimental procedures? And what possible caveats must we avoid to draw the right
28 conclusions from experimental studies?

29 In the following sections, I discuss how we deal with mental and physical phenomena as experimental
30 variables and outline how we can use them to make inference about causal relations, how the variables
31 relate to analytical problems, and how this relation sets the foundation for an empirical approach to
32 studying mental causation.

33 **2 Mental variables and physical variables**

34 The analytical approach to studying mental causation focus on how any ideal mental event M and any
35 ideal physical event P can (or cannot) interact. In contrast, empirical science deals with observed data
36 to infer relations amongst the events the variables represent. When taking an empirical approach, our
37 first assumption is that mental and physical events are real and that we can measure and/or manipulate
38 them. To measure and manipulate mental event, they need to be operationalized as experimental
39 variables.

40 The first hurdle is that the term *mental event* has different meanings in different discussions. It can refer
41 to specific mental content occurring within a limited time window or refer to general states, e.g., a
42 transient intention to move one's arm versus a general state of wakefulness or being in a coma (Hohwy,
43 2009; Laureys, 2005). There is no clear-cut definition of mental events; partly due to the uncertainty in
44 defining the nature of mentality, to begin with. The content of mental events can refer to
45 phenomenological properties or cognitive properties (Block, 2005; Cohen and Dennett, 2011). The
46 precise definition of mental events is not of importance: as long as we can accept that mental events
47 exist—either as phenomenological states or cognitive processes—then we can operationalize these as
48 *mental variables* in experimental settings. Thus, mental events can be defined by their
49 phenomenological content, or they can be defined from a cognitive perspective without referring to
50 phenomenology.

51 If we are strict, we could argue that because the phenomenological content of mental events is available
52 only to the subject, we can never measure them. We can for example never know if subjects have
53 inverted qualia, or if they are philosophical zombies (Chalmers, 1997). There is no practical solution to
54 this problem, but this does not exclude mental variables from being meaningful in experimental
55 contexts (Overgaard et al., 2008; Seth et al., 2005; Tononi and Koch, 2015).

56 Subjective reports or behavioral responses might not be “direct” access to phenomenological content,
57 but they serve as indirect indications that mental events are occurring. We see indirect variables in all
58 branches of science: when measuring distortion of light as an indication of cosmic bodies due to
59 gravitational bends, and in neuroimaging where the blood-oxygen flow in areas of the brain is a proxy
60 of neuronal activity (Logothetis et al., 2001). We accept these indirect measurements because the link
61 between the indirect measure (light distortion; oxygen-blood flow) and the object (cosmic body; neural
62 activity) is based on assumptions that we agree upon (strong gravitation bends light; active neurons are
63 associated with higher blood flow).

64 To obtain a mental variable that indicates that a subject is experiencing a particular mental event is
65 sufficient for it to be an experimental variable. We can measure mental events, by obtaining
66 introspective reports about the mental events or other indications that subjects are experiencing certain
67 mental events (Overgaard et al., 2008). What we obtain by these methods are *mental variables*. For
68 mental variables to be valid, they have to be consistent as any other variable. Mental variables should
69 exclusively capture the event they are intended to measure while exhaustively capturing any occurrence
70 of the mental event (Jensen et al., 2017; Reingold and Merikle, 1988).

71 We are not required to know the ontological reality of the mental events that the mental variables
72 measure, i.e., we do not need to impose a predefined solution to the mind-body problem to use mental
73 variables in cognitive neuroscience. If we can trust our methods of obtaining mental variables, then we
74 have indications that the mental events are occurring, and we are justified in using these as variables in
75 experimental studies.

76 *Physical events* seem more intuitive than mental events, but, upon further inquiry, it is not
77 straightforward what constitutes a physical event (Crane and Mellor, 1990; Melnyk, 1997; Smart,
78 1978). Intuitively, we can easily characterize different events, such as a neuron firing, cerebral blood-
79 flow, or the force of an accelerating mass as physical events. Although only the latter is the described
80 in the scientific language of physics, we consider all previous examples as physical events. That we
81 consider the above as physical phenomena are because we not have any reason to assume their
82 existence is dependent on anything that violates the language of physics (Stoljar, 2001). In this sense,
83 cognitive and neural processes—from single neurons to whole-brain network communications—are all

84 physical phenomena. Whatever we measure in cognitive neuroscience that relates back to the
85 biophysical properties of the brain or body is in this context a *physical variable*.

86 The challenge of using an empirical approach to studying mental causation is how to combine mental
87 variables and physical variables in experiments to infer their causal relation. Another factor that adds to
88 this challenge is that there are likely as many definitions of *causation* as there are definitions of *mental*
89 *events*. Despite disagreements, there is one prevailing view of how to make causal inference in science.
90 In the next section, I give a brief overview of the general principles of causal inference and then return
91 to how we can use these principles to make inference about mental causation.

92 **3 The principles of causal inference**

93 The causal effect of X is the difference X would make on the outcome Y , based on the counterfactual
94 conditions of X being present or not. Say we want to know whether X (a drug) cause Y (a fever
95 reduction). Depending on the intervention the variable X takes the values $X=1$ (taking the medication)
96 or $X=0$ (not taking the medication). Y is the measured outcome of the counterfactual conditions, as
97 $Y(X=1)$ and $Y(X=0)$. If X is causing Y , then Y will follow $X=1$ but not $X=0$. Thus, the *causal effect* of X
98 on Y is the measured difference between $Y(X=1)$ and $Y(X=0)$. The causal effect is a numerical quantity
99 that indicates the difference between the counterfactual conditions of X and not X (Rubin, 1974;
100 Woodward, 2005).

101 In reality, both conditions cannot occur: a subject cannot take the pill ($X=1$) and at the same time not
102 take the pill ($X=0$). Only one of the counterfactual conditions can occur for the particular case. The *true*
103 causal effect for any single case cannot be estimated. Empirical inference of causal effects instead
104 approximates the true causal effect. This is done by having several independent occurrences of the
105 relation we are investigating. We then expose half of the cases to the condition $X=1$. This is the
106 intervention group. The other half of the cases are kept the same without the intervention ($X=0$). For
107 each case i exposed to the intervention, we measure the outcome $Y_i(X_i=1)$, and for each instance j not
108 exposed to the intervention, we measure $Y_j(X_j=0)$ (Table 1). The mean difference between the two
109 conditions estimates the true causal effect. The mean difference between the columns $Y(X=1)$ and
110 $Y(X=0)$ in Table 1 is the *estimated* causal effect of X on Y (Rubin, 1974).

111

112 **Table 1: Data frame for estimating the causal effect of variable X on outcome variable Y .**

<i>case</i>	X	$Y(X=1)$	$Y(X=0)$	$Y(X=1)-Y(X=0)$
i_1	1	10	NA	NA
j_1	0	NA	0	NA
i_2	1	11	NA	NA
j_2	0	NA	0	NA
i_3	1	9	NA	NA
j_3	0	NA	0	NA
MEAN		10	0	10

113

114 For the causal effect to be a measure of causality—not just a correlation between two measurements—
 115 the intervention on X must be the only systematic change between conditions (Rubin, 1974;
 116 Woodward, 2012). In reality, however, there is always variation between single cases. The variation
 117 can come from several sources: imprecision in the measurements, noise in the environment, or
 118 variation inherent in what we denote as events of type X . We have to take variability between cases into
 119 account in the estimation of causal effects.

120 Random variation, unrelated to the intervention, is a minor problem as it will cancel out (to some
 121 degree) with an adequate number of cases. But if the variability between cases covaries with the
 122 intervention, it will invalidate the causal inference. It must be assumed that the intervention on X is the
 123 only variable that affects Y to make a valid causal inference. If systematic variation between groups
 124 occurs, it cannot be ruled out that the change in Y is due to confounding variables rather than X . It is
 125 important to control for systematic confounding *background variables* (i.e., any other variable than X
 126 and Y). Control of background variables is done by random sampling and systematic matching of
 127 background variables before making the intervention (Ahern et al., 2009; Rubin, 1974; Stuart, 2010).
 128 We can measure the background variables as separate variables $B_1, B_2...B_n$ to ensure their distributions
 129 are similar between conditions. E.g., B_1 and B_3 in Table 2 appear to have similar distributions between
 130 cases i and j , but there seems to be a problematic difference between groups in B_2 . Good experimental
 131 design requires adequate control of possible confounding variables.

132

133 **Table 2: Expanded data frame for estimating the causal effect of X on Y .**

<i>case</i>	B_1	B_2	B_3	X	$Y(X=1)$	$Y(X=0)$	$Y(X=1)-Y(X=0)$
i_1	5	1	3	1	10	NA	NA
j_1	4	10	2	0	NA	0	NA
i_2	6	0	1	1	11	NA	NA
j_2	5	11	3	0	NA	0	NA
i_3	4	2	2	1	9	NA	NA
j_3	6	9	1	0	NA	0	NA
MEAN	$i=5, j=5$	$l=1, j=10$	$i=1, j=1$		10	0	10

134

135 The estimated causal effect does not tell exactly what the true causal effect is for single cases: it is a
 136 generalized effect of type X events on type Y events (Dawid, 2000; Holland, 1986; Rubin, 1974). The
 137 estimated causal effect tells the probability of Y following X in the case that Y is a binary variable. If Y
 138 is a parametric variable, the causal effect is a numeric value indicating how much we expect X to
 139 change Y (Woodward and Hitchcock, 2003), e.g., ten “units” in the example in Table 1. In the
 140 following, whenever I refer to *causal effect*, I am referring to the *estimated causal effect*.

141 Causal effects are not *truths* in the logical sense. The causal effects are probabilistic relations between
 142 events estimated under controlled conditions that allow us to apply counterfactual logic to conclude a
 143 causal connection. For example, when testing if a new drug reduces fever, we do not need to describe
 144 how the chemical compound is absorbed in the body, passing the bloodstream, etc. to estimate the
 145 causal effect of the drug. As long as we have a causal effect of X on Y obtained under convincing
 146 circumstances, we can justify the conclusion that there is a causal relation between X and Y .

147 **4 Causal inference for mental events**

148 We can estimate the causal effects of mental events M in a similar way to how we estimate any other
 149 causal effect: the mental events are the variables we manipulate, and the behavioral outcome is the
 150 dependent variable we measure. Say we want to investigate if the intention to move one’s arm cause
 151 one to move the arm: following the reasoning in the previous section, we can make an experimental
 152 manipulation so that subjects experience a specific mental event M (intention to move arm) in some
 153 conditions and not in others, while keeping all other variables constant, measure the outcome Y
 154 (physical movement of arm) for contrasting conditions, and estimate the causal effect of M on Y . M is a

155 cause of Y if there exists some intervention on M that changes the value of Y while keeping everything
156 else equal (Woodward, 2012). The causal effect of M on Y is the mean difference between cases, for
157 cases with $M=1$ and cases with and $M=0$, assuming there is no other systematic variation between the
158 conditions. Causal inference with mental events is, like any other causal inference, estimated
159 probabilistic relations between events. The inferred relations are the generalized effect of mental events
160 of type M on type Y events.

161 To infer mental causation from experiments, we need several instances of the same mental event. We
162 need to ensure that the mental events, which we conceptualize as experimental variable M , is similar
163 enough across the entire experiment and between subjects that we can justify that they belong to the
164 same type of mental event and estimate causal inference. E.g., if investigating the causal effect of the
165 intention to move the hand, it must be assured that the intention is functional or phenomenological
166 *equivalent* across cases. It is impossible to know if the *intention to move* is phenomenological identical
167 between subjects as the experience is only available to the subjects, but this does not invalidate using it
168 in causal inference. While we cannot assume that two similar behaviors are followed by similar mental
169 states between subjects, we can make sure that the accompanying mental events are consistent. By
170 measuring the mental events through introspective reports, we assess whether subjects describe the
171 mental events in a consistent manner across subjects (Overgaard et al., 2008).

172 The problem of mental causation is often framed as whether *any* mental event can have *any* causal
173 relevance in the physical world. This question is not suited for experimental research. Estimated causal
174 effects apply to the events for which they were estimated. If we show one type of mental event M_a to be
175 causally relevant (or irrelevant) for the outcome Y , then this does not mean that other mental events M_b ,
176 M_c , etc. have the same level of relevance. It is not given that the *experience of red* or the *intention to go*
177 *on vacation* has the same causal relevance for moving one's arm as *the intention to move the arm* just
178 because all are examples of mental events. They are different mental events. To conclude that the
179 causal properties of one type of mental event apply to all mental events is an error analogous to
180 concluding that the effect of one kind of drug applies to all kinds of drugs. Experimental studies of
181 mental causation must be specific about what type of mental events they are dealing with.

182 In conclusion, to study mental causation with experimental research, we treat mental variables as any
183 other experimental variable. But of course, if it is this simple, mental causation would not be a

184 controversial topic. The premise of this framework is that we do not know how mental events fit into
185 the physical world. The unknown relation between mental events and physical events makes causal
186 inference a peculiar enterprise. But it is not as difficult as one might think if we are aware of the mind-
187 body problem.

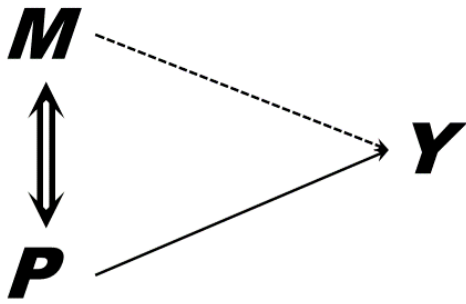
188 **4.1 Dealing with the special nature of mental events**

189 To investigate the causal effects of mental events, we must first assume that mental events are related
190 to the physical world. We do not need to assume *how* mental events are related to the physical world,
191 we only assume that mental events *depend* on the physical world. This assumption is easy to justify:
192 sensory inputs from the environment cause perceptions, intentions are directed towards the physical
193 world, and particular neural activity is related to mental content (Aru et al., 2012; Chalmers, 2000).

194 We can make this assumption explicit: for each mental event M , there exists at least one physical event
195 P , which is the minimal physical event necessary for instantiating M (for the given occurrence of M).
196 This is different from stating that mental events are physical events. If we prefer, we can view mental
197 events as non-physical properties. But if so, they are attached to the physical events that cause them.

198 How mental and physical events interact is still unknown: do the interaction go from the physical to the
199 mental or can it go both ways? If the relation between M and P is one-directional, then mental events
200 cannot produce causal effects. If mental events are causally irrelevant, then we hardly need an
201 empirical approach to study them.

202 To outline the problem: we want to know if M can cause Y given M is bound by its physical base P
203 (Figure 1). How we view the relation between M and P determines if this problem shows that mental
204 causation is impossible.

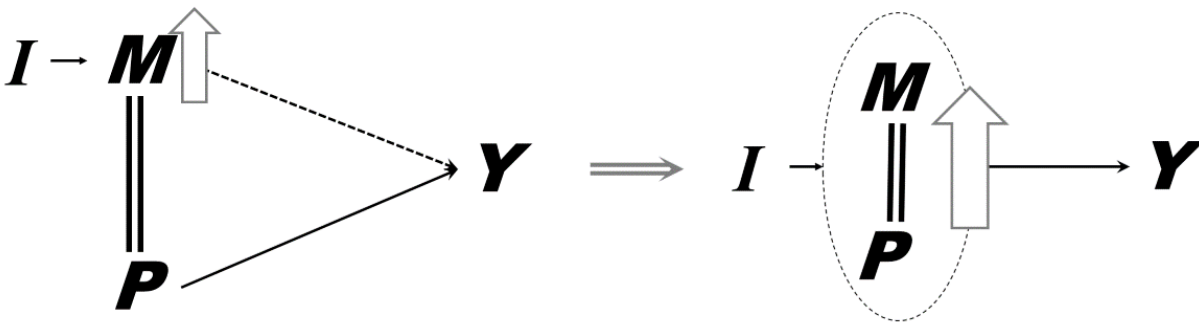


205

206 **Figure 1: Outline of the problem: Can M —realized by P —be a cause of Y ?**

207 The first solution is to assume an *identity* relation between M and P , where the two events are the same
 208 singular event: any difference in appearance is only epistemological (e.g., Smart, 1959). The apparent
 209 dichotomy between physical and mental stems from the event being measured as either the physical
 210 variable P or as the mental variable M . The statements “ M causes Y ” and “ P causes Y ” are describing
 211 the same relation. If P is the cause of Y then so is M by definition. In the case of M and P sharing an
 212 identity, there is one causal factor that can influence Y .

213 In this case, mental variables are another type of physical variable, and it is valid to treat them like any
 214 other variable for the purpose of causal inference. If we measure either M as a mental variable or its
 215 physical base P , we would measure the same event twice. If we observe only one of either M or P , we
 216 automatically have proof that the other identity is present as well. This also gives that intervening on
 217 either M or P is an intervention on the same event. For inferential purposes, we must collapse M and P
 218 into a single causal factor MP (Figure 2).



219

220 **Figure 2: In the case of $M=P$, any intervention I that change M will be an intervention on MP .**

221 On the other hand, we can take M to be different from P . The occurrence of M is still dependent on P
222 (given the initial assumption), but M is a (non-reductive) supervenient property of P . The difference in
223 appearance is not just epistemological: the difference between a mental variable measuring M and a
224 physical variable measuring P each captures some underlying features that are non-overlapping.

225 Though M and P are ontologically distinct, P is both a sufficient and necessary condition for M . For M
226 to be present so must P and when P is present, so is M . For M to be the cause of Y , P must be a physical
227 cause of Y . To avoid Y being overdetermined there can only be one cause of Y , so either M or P must be
228 removed as the cause of Y . If the physical world is causally closed, then P cannot be eliminated. Hence,
229 P is a sufficient cause of any change in Y we would ascribe to M . M is an epiphenomenon and is
230 excluded from the causal relation (Kim, 2005). Only the solid arrow in Figure 1 describes the real
231 causal connection between M , P , and Y .

232 If we take M as a supervenient property of P , then mental variables measure different phenomena than
233 the physical variables do. But since P is sufficient and necessary for M , we run into problems if we try
234 to isolate either and estimate causal effects.

235 Imagine an experiment where we have a brain stimulator that can target—and only targets—non-
236 physical mental events without affecting any physical events. We use the non-physical stimulator to
237 induce the intention to move one's arm ($M=1$) in a group of subjects. We also have a control group that
238 is not subjected to the non-physical stimulation and will not experience the intention to move their arms
239 ($M=0$). Assume that the brains of the subjects, independent of group, all are in a given state P^* at the
240 moment before the intervention. P^* is in no way related to the intention to move one's arm. During the
241 experiment, the brains of the control group will continue to be in state P^* . When the non-physical
242 stimulator induces the (non-reductive) intention to move in the intervention group, it must follow that
243 the physical base P of the intention to move have to be present for M to be present. For the non-
244 physical intervention to change the value of M , it will follow that P^* change to P (Baumgartner, 2009;
245 Kim, 2005). It is impossible for M to change without a corresponding change in P .

246 That both M and P change is a problem for causal inference: since $P \neq M$, we must place P as a column
247 in our matrix of confounding variables (Table 3). But if P is not present then neither is M : M and P
248 covary. We can estimate a difference between the two groups, but we cannot determine if the effect is

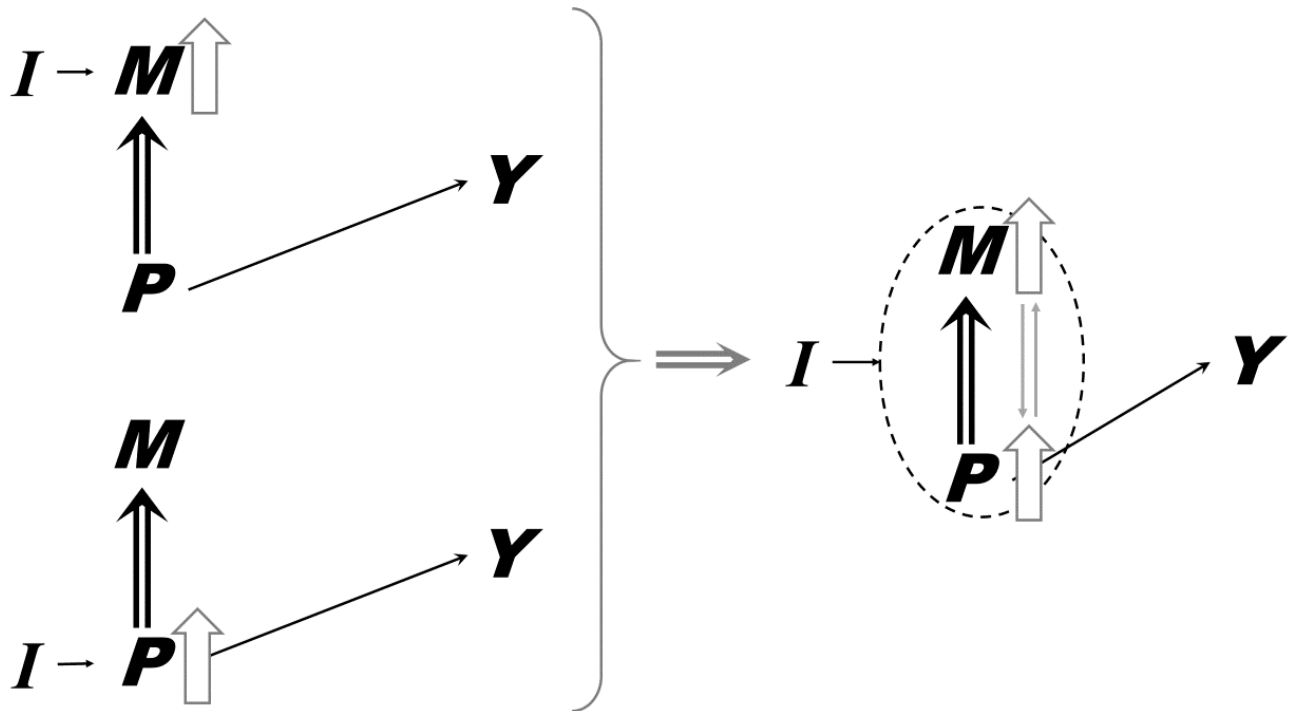
249 caused by M or P : it is not possible to isolate the non-physical M as required for causal inference
 250 (Baumgartner, 2010, 2009). We cannot estimate the causal effect of a mental event alone.

251

252 **Table 3: Data frame for estimating the causal effect of mental event M on Y . M and its physical**
 253 **base P are perfect covariates.**

<i>case</i>	B_1	B_2	P	M	$Y(X=1)$	$Y(X=0)$	$Y(X=1)-Y(X=0)$
i_1	5	3	1	1	10	NA	NA
j_1	4	2	0	0	NA	0	NA
i_2	6	1	1	1	11	NA	NA
j_2	5	3	0	0	NA	0	NA
i_3	4	2	1	1	9	NA	NA
j_3	6	1	0	0	NA	0	NA
MEAN					10	0	10

254



255

256 **Figure 3: In case M is a (non-reductive) supervenient property of P , any intervention I on M must be accompanied by**
 257 **a corresponding change in P . M and P cannot be separated as causal factors despite being ontologically different. The**

258 same is true if we make an intervention I on P : Any intervention on either M or P will be an intervention on MP as a
259 change in one will be accompanied by a change in the other.

260 The covariation between M and P is not only an issue when isolating M as a causal factor. Imagine the
261 same experiment above, but this time we use a different brain stimulator that intervenes on P and
262 estimates the causal effect of P on Y (Figure 3, bottom diagram). In the intervention group, the brain-
263 states change from P^* to P while the brain-states of the control group remain the same. Since P is the
264 physical base of M then a change from P^* to P will also induce M . As with the non-physical
265 stimulation inducing M , we are not able to isolate P as the cause of Y . Even if we assume that M and P
266 are different, it is impossible to separate them as causal factors: M and P are perfect covariates.

267 Since we cannot separate M and P as causal factors, we have two options when it comes to causal
268 inference: we can abandon the attempt of an empirical approach to mental causation, as we cannot
269 construct an experiment where we isolate M and keep all other factors (including P) constant. Or we
270 can collapse M and P into a single variable MP .

271 Considering causal properties of mental events and their physical base as a single causal factor is not
272 new in the analytical approach to mental causation (e.g., Kim, 2005; Lewis, 1994; Mele, 2009; Sperry,
273 1980; Woodward, 2015). Here I argued for the same position as we cannot separate the two in practice:
274 it is a practical necessity—not an ontological assumption. M and P can be ontologically different, but
275 the isolation of the factor M from P is impossible no matter which of the above solutions to the mind-
276 body problem we prefer. This means that we cannot answer the question if a mental event or its
277 physical base is the cause of action by empirical means. We can only answer if the mental event and its
278 physical base as a single factor is the cause of action. For practical purposes, the causal relevance of a
279 mental event M is the same as its physical base P .

280 **4.2 Reduction and causal explanations**

281 One could argue that we are reducing away *real* mental causation by always considering the causal
282 properties of mental events together with their physical realization. This is only the case if we start with
283 the position that mental causation per definition must be non-physical. If one feels that the concept of
284 mental causation is reserved for “pure” mental causation, we can instead call the causal effects
285 “mental-and-physical-base causation” and proceed.

286 One can still ask: since the causal effect of M only can be estimated together with P , can we then
287 remove M from the causal relation and only deal with the physical variables? In principle, we can
288 remove M from the causal explanation. But it is not feasible in practice. In reality, it is rarely the case
289 that higher-order explanations can be reduced to statements involving fundamental physical processes,
290 even in cases where we are justified in assuming the higher order phenomena are genuine reducible
291 physical phenomena (Anderson, 1972; Bedau, 2002). The reason is that we do not understand what
292 constitutes a physical base of a mental event. What real type of physical event P denotes in the
293 examples above is undisclosed in the real world. It might be explanations involving neural anatomy,
294 neural communication, or fundamental laws of physics. Even if we assume mental events are reducible,
295 we do not have the relevant information to remove mental events from the inference of mental
296 causation.

297 Reductionism in practice is also problematic in the experimental setting. To estimate causal effects, we
298 need several instances of the same mental event M , and we cannot guarantee that each occurrence of M
299 is identical. Though each case has to be similar enough for the causal inference to be valid, we have to
300 tolerate some variation in practice. The variation within an experimental variable opens the possibility
301 that each case of M does not have the same physical base. This can be because of the uncertainties in
302 the mental variables or because mental events can be realized by different physical events (Fodor,
303 1974). Rather than thinking about M and P as ideal events, think of them as sets of similar events: \mathbf{M} is
304 the set of similar mental events $M_1, M_2 \dots M_N$ and \mathbf{P} is a set of physical events $P_1, P_2 \dots P_N$ that each
305 corresponds to the physical base of the mental events in \mathbf{M} . Since it is possible for the variable P to be a
306 set of several different physical states we cannot assume that all elements in \mathbf{M} can be collapsed with
307 the same physical event P_i . Each mental event M_i must have its own corresponding physical base P_i .
308 When considering the causal properties of each M_i , we cannot distinguish it from the causal properties
309 of its physical base P_i , and we have to collapse each M_i and P_i into a single causal factor. The problem
310 of separating M and P does not change, but it does change how we can control for confounding
311 physical variables in causal inference. We have to deal with a practical form of multiple realizations of
312 events, even if the events are reducible (Aizawa & Gillett, 2009). It is, thus, problematic to attempt to
313 completely remove mental variables, by only measuring physical variables, as we cannot be sure that

314 we measure the correct underlying physical base in ever repetition—at least not with our current
315 understanding of how the brain gives rise to mental events.

316 Finally, removing mental events from inquiry about mental causation, even in ontologically possible,
317 might not be desirable in practice. Low-level explanations do not always provide useful causal
318 explanations. Causal explanations are about how some units of interest X have an impact on another
319 unit of observation Y . The units of interest can be any phenomenon, e.g., brain processes, pills, social
320 factors, or mental events. We can assume that X can be broken into pieces and explained as composites
321 x_1, x_2 , etc., and explain how the composites cause outcome Y . Replacing causal explanation involving X
322 with an explanation about the composites have to take all the composite parts into account. The
323 explanation involving the composites increase in complexity. The increase in complexity makes the
324 required explanation more difficult and is contrary to the purpose of causal explanations (Lipton,
325 2005).

326 To say that a pill cause fever reduction does (in most cases) provide the information we want to know
327 about the pill and its use—even if we have a full account of the chemical compounds in the pill and
328 how they interact with the biophysical processes that regulate body temperature. Similar, to say that my
329 intention to raise my arm is the cause raising my arm, is as valid a causal explanation as an explanation
330 involving all the neural processes giving rise to my intention to raise my arm. One is not more correct
331 than the other, but the first is a lot simpler than the latter. Since causal inference is relative effects of
332 variables of interest, we can reverse the issue of reductionism and ask if we can remove physical
333 variables and only consider the causal effects of mental variables? The proper level of analysis depends
334 on the phenomena we are interested in. It is not necessary to regress to lower level explanations to
335 explain the causation of higher-order events to address causation.

336 The reason we want physical variables is that in the context of mental causation this is usually part of
337 what we want to know: it is not if mental events can have causal effects, it is *how* they can have causal
338 effects and how they interact with physical processes; e.g., how the intention to move the hand is part
339 of the nervous system responsible for locomotion. How does the relative contribution from mental
340 events and non-mental physical processes generate behavior? By combining mental variables and
341 (physical) neuro-cognitive variables in experimental designs that we can answer these questions.

342 **5 Towards a science of mental causation**

343 Combining mental variables and physical variables in causal inference is difficult, as we are dealing
344 with variables whose underlying ontology can be dependent on one another. When measuring blood
345 flow in the brain or electrophysiological potentials, we do not know whether the measured activity is
346 the physical base of mental event M , unrelated to M , or only a part of the physical base. The mental
347 event and its physical base MP is not the only event occurring in the brain when we investigate the
348 causal effect of MP . Isolating the neural base of a mental event is not as simple as observing which
349 physical variables that co-vary with the mental variables. Any physical variable we measure can be a
350 precursor of the real physical base of the mental event or an effect following the mental event (Aru et
351 al., 2012). Since we are dealing with complex systems and different explanatory levels, we have to be
352 careful when operationalizing potential causal factors in experiments.

353 To study causal processes in cognitive science is difficult, even without considering mental variables.
354 The nervous system is not wired as a linear causal chain where X causes Y , causes Z , and so on. The
355 brain consists of interconnected networks that operate on different anatomical scales and different
356 timescales. For example, to investigate causal effects of *the intention to move the hand* on hand
357 movements, we have to consider that the process is part of a network that depends on both long-range
358 connectivity and local specification of function working in many hierarchical feedback loops
359 (Rizzolatti and Luppino, 2001; Shadmehr and Krakauer, 2008).

360 To ask whether a mental event M cause action Y is not different from asking if certain neural activity in
361 the supplementary motor area (call this activity P_{SMA}) cause Y . To infer this relation following the
362 principles of causal inference, we make an intervention on P_{SMA} while keeping everything else in the
363 brain is unaffected. If we find that the intervention on P_{SMA} changes Y , compared to an adequate control
364 condition, we can conclude that P_{SMA} causes Y . This does not mean P_{SMA} is the only cause of Y nor that
365 P_{SMA} is an isolated cause of Y . We should not fool ourselves to believe that because we measure two
366 events (P_{SMA} and Y) in a complex system (the brain and its interaction with the environment) and study
367 the effect, then nothing else of relevance is going on. There definitely would be in this example, e.g.,
368 the neural communication from cortical motor areas to the basal ganglia, thalamus, and the peripheral

369 nervous system, but in the experiment, we make sure that all these are constant to pursue the question
370 aimed at P_{SMA} .

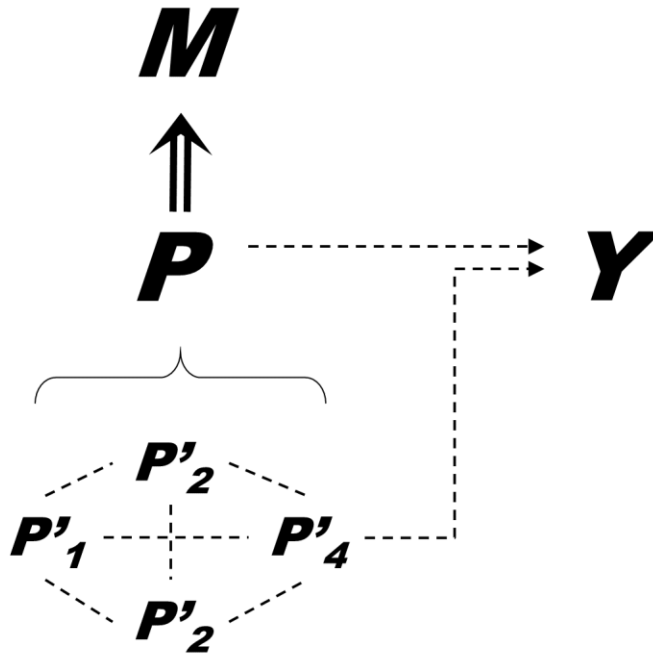
371 We cannot make inference about all the other variables from the experiment, we can only make an
372 inference about the causal effect of P_{SMA} . To ask if a single variable X causes Y is a simplification of the
373 complex mechanisms we are dealing with. But it is a simplification for inferential purposes, not a
374 simplification of how the system is. Estimating a causal effect of X on Y does not tell us about the
375 possible effects of other events Z_1, Z_2 , etc. on Y or how they are related. It only gives the causal effect
376 of X on Y . To test whether Z_1, Z_2 , etc. are causes of Y , we need separate experiments for each type of
377 event. The “start” and “end” of the relation are pragmatic cuts where we select events of the kind that
378 we want to investigate from the larger configuration.

379 Controlling variables “downwards” become increasingly difficult. Since we do not know the precise
380 relation between the mental and physical, we might try to control for a physical variable that, unknown
381 to us, is a part of the physical base of the mental event.

382 The physical base of mental events is likely a complex system of interconnected processing between
383 distinct sub-parts in a network (Baars, 2005; Tononi and Koch, 2015). Each node in the network is not
384 enough to constitute the physical base of M . Only all parts connected are sufficient to enable M .

385 Let us continue the example above to make inference about the lower level configuration of P_{SMA} . SMA
386 is divided into fine-grained anatomy based on local functionality and afferent connections (Nachev et
387 al., 2008). Assume we can divide P_{SMA} into four parts, as illustrated in Figure 4. All parts have to be
388 “active” to constitute P_{SMA} (and thereby M). By replacing the unified base P with the parts in the data-
389 frame for the causal experiment, we get Table 4. All parts, which together constitute the base of M , are
390 perfect covariates with M as P were in Table 2.

391 When we look at the parts, only P'_4 has a direct causal link to outcome Y . If we were to intervene on
392 P'_4 in Figure 4 (keeping all other variables constant), we find that P'_4 have a causal effect on Y . This
393 time M is no longer a perfect co-variate: we see that M only occurs when all sub-components are
394 present, as shown in Table 5. The conclusion we would draw from Table 4 (M cause Y) and Table 5
395 (P'_4 cause Y) are both correct.



396

397 **Figure 4:** M is realized by P , which is constituted by the parts P'_1 , P'_2 , P'_3 , and P'_4 . Om sub-part P'_4 has a connection
 398 to Y .

399

400 **Table 4:** Data frame for estimating the causal effect of mental event M on outcome variable Y
 401 with measurements of low-level parts of P .

<i>case</i>	B_1	B_2	P'_1	P'_2	P'_3	P'_4	M	$Y(X=1)$	$Y(X=0)$	$Y(X=1)-Y(X=0)$
i_1	5	3	1	1	1	1	1	10	NA	NA
j_1	4	2	0	0	0	0	0	NA	0	NA
i_2	6	1	1	1	1	1	1	11	NA	NA
j_2	5	3	0	0	0	0	0	NA	0	NA
i_3	4	2	1	1	1	1	1	9	NA	NA
j_3	6	1	0	0	0	0	0	NA	0	NA
MEAN								10	0	10

402

403 **Table 5:** Data frame for estimating causal effect similar to Table 4, but for the causal effect of a
 404 low-level part of P on outcome variable Y instead of mental event M .

<i>case</i>	B_1	B_2	P'_1	P'_2	P'_3	M	P'_4	$Y(X=1)$	$Y(X=0)$	$Y(X=1)-Y(X=0)$
i_1	5	3	1	1	1	1	1	10	NA	NA
j_1	4	2	1	0	1	0	0	NA	0	NA
i_2	6	1	0	1	0	0	1	11	NA	NA
j_2	5	3	1	1	0	0	0	NA	0	NA
i_3	4	2	1	0	0	0	1	9	NA	NA
j_3	6	1	0	0	1	0	0	NA	0	NA
MEAN								10	0	10

405

406 It might appear paradoxical that M (together with P) can both be a cause of Y and, at the same time, be
407 controlled when estimating the causal effect of only a part of P . The apparent discrepancy is because
408 the two tables represent two different causal experiments. Causal effects are about the variable
409 manipulated by the intervention. They do not tell us about the causal impact of other factors. We can
410 for example not conclude that there is no effect of P'_1 , P'_2 , P'_3 , or M from the experiment in Table 5.
411 We can only conclude that there is an effect of P'_4 on Y . How we define the events, we investigate,
412 determines the explanations we can give based on the causal inference.

413 To conclude that “higher order” events are causally irrelevant is not possible from the observation that
414 only a subset of low-level components has a causal effect on the outcome. Demonstrating that P'_4 is the
415 part of P that cause Y , does not mean that P is not also a cause of Y . The explanation containing P'_4 is a
416 fine-grained causal explanation that the one containing P or M , but it does not follow the higher order
417 elements are irrelevant. Stating that a mental event is the cause of action does not exclude causal
418 explanations in terms of neural mechanisms, nor does causal explanations in terms of neural
419 mechanisms exclude causal explanations involving mental events (Pernu, 2011).

420 **6 Inferring mental and non-mental causes**

421 It is tempting to treat mental variables and physical variables as measures of different ontological levels
422 and contrast the two variables to answer whether the physical event or mental event is the cause of Y .
423 But this is not possible. Mental events are realized by physical events that might be contained in the
424 measured physical variable. Similar, we cannot per default claim that no mental events are occurring
425 when we measure physical variables; especially when the type of measured physical event is part of the
426 organ that generates the mental events. It is only possible to claim that a physical variable represents an

427 unconscious (non-mental) cause of action if no mental variables influence the outcome. It is not enough
428 to show that a physical variable influences the outcome to conclude that mental events are irrelevant for
429 action. For example, movement-related cortical potentials precede the intention to move the hand
430 (Fried et al., 2011; Libet et al., 1983), which is taken to prove that the intention to move is not a cause
431 of moving the hand—the true cause is unconscious neural processes expressed as the movement-related
432 cortical potentials (Harris, 2012; Libet, 1999, 1985; Wegner, 2002).

433 Showing that a physical variable (movement-related potentials) precedes a mental variable (the
434 intention to move) is not enough to infer that the former is the cause of outcome rather than the latter.
435 The observation lacks a contrasting condition to rule out if there is a causal effect of the intention on
436 the outcome. It is impossible to infer that the intention is *not* a cause of moving the hand from the
437 observation. It is even impossible to infer that the movement-related potentials are the cause of moving
438 the hand. Precedence does not imply causation. We can only claim that an action is unconsciously
439 initiated if we show both a causal effect of the unconscious processes and a null-effect of the mental
440 event. To infer true unconscious causes of action, we must include mental events as background
441 variables to show that mental events do not covary with the unconscious events.

442 For example, it is unclear how the movement-related potentials mentioned above are related to the
443 conscious intention to move. Whether they reflect unconscious neural activity is unknown. Some
444 studies have shown covariation of the readiness potential and conscious intention, while others are
445 unable to do so (Haggard and Eimer, 1999; Keller and Heckhausen, 1990; Schlegel et al., 2013;
446 Schultze-Kraft et al., 2016; Vinding et al., 2014). To conclude that a physical variable represents
447 genuine unconscious action initiation it must be shown that it has a causal effect while keeping the
448 mental content constant.

449 It is surprisingly difficult to determine what *non-mental* or unconscious means (Moors and De Houwer,
450 2006). It is not as simple as dividing mental/conscious and non-mental/non-conscious processes: the
451 transition can be gradual (Miller and Schwarz, 2014; Sandberg et al., 2011), and there are separate
452 ways to be unconscious of stimuli (Kim and Blake, 2005; Rothkirch and Hesselmann, 2017). In
453 conclusion: it is not valid to ignore mental events altogether and conclude that the action was
454 unconsciously initiated. If we want to show that a low-level neural event is the cause of action rather
455 than mental events, then it requires a null-effect of the relevant mental variable.

456 **7 Conclusion**

457 The complexity of the nervous system makes investigating mental causation a difficult task. But it is
458 possible to study mental causation by applying the principles of causal inference as in any scientific
459 field. Given the special nature of mental events, we need to treat the mental event and its physical base
460 as one factor in the experimental design. This means that we cannot answer whether a mental event or
461 its physical base is the cause of action. The type of questions we can answer is whether a given mental
462 event M (realized by P) is a cause of Y . Mental causation is measured as the causal effect of MP on Y in
463 controlled experiments with contrasting conditions that control for confounding variables. This
464 approach is not for those who seek an either/or answer to mental causation, but it is of relevance to
465 those who seek to investigate the neurocognitive and behavioral relevance of mental events. Causal
466 effects are in any instance relative contributions of the variables—not ontological truths.

467 To investigate the mental causation, we have to think about mental causation as relative contributions
468 of events in complex systems with different descriptive levels. The different descriptive levels do not
469 preclude one another. Experimental scientists must shift their approach to mental causation from the
470 search for ultimate answers in the analytical discussion and instead focus on relative effects of well-
471 operationalized variables. Cognitive scientists have to consider how to manipulate mental events in
472 experimental design and how to control confounding variables. The solution will depend on the type of
473 mental event and outcome behavior in question. Note that these considerations are methodological, not
474 metaphysical.

475 **Acknowledgments**

476 This work was funded by Danish Independent Research Council | Humanities. I would like to thank
477 Lise Marie Andersen and Nina Jensen for comments on earlier drafts of this paper.

478 **References**

- 479 Ahern, J., Hubbard, A., Galea, S., 2009. Estimating the Effects of Potential Public Health Interventions
480 on Population Disease Burden: A Step-by-Step Illustration of Causal Inference Methods. *Am. J.*
481 *Epidemiol.* 169, 1140–1147.
- 482 Anderson, P.W., 1972. More Is Different. *Science* 177, 393–396.
- 483 Aru, J., Bachmann, T., Singer, W., Melloni, L., 2012. Distilling the neural correlates of consciousness.
484 *Neurosci. Biobehav. Rev.* 36, 737–746.
- 485 Baars, B.J., 2005. Global workspace theory of consciousness: toward a cognitive neuroscience of
486 human experience. In: *Progress in Brain Research*. Elsevier, pp. 45–53.
- 487 Baumgartner, M., 2009. Interventionist Causal Exclusion and Non-reductive Physicalism. *Int. Stud.*
488 *Philos. Sci.* 23, 161–178.
- 489 Baumgartner, M., 2010. Interventionism and Epiphenomenalism. *Can. J. Philos.* 40, 359–383.
- 490 Bedau, M., 2002. Downward Causation and the Autonomy of Weak Emergence. *Principia* 6, 5–50.
- 491 Block, N., 2005. Two neural correlates of consciousness. *Trends Cogn. Sci.* 9, 46–52.
- 492 Chalmers, D.J., 1997. *The conscious mind*. Oxford University Press, New York; Oxford.
- 493 Chalmers, D.J., 2000. What is a neural correlate of consciousness? In: *Neural Correlates of*
494 *Consciousness: Empirical and Conceptual Questions*. MIT Press, pp. 17–40.
- 495 Cohen, M.A., Dennett, D.C., 2011. Consciousness cannot be separated from function. *Trends Cogn.*
496 *Sci.* 15, 358–364.
- 497 Crane, T., Mellor, D.H., 1990. There is no question of physicalism. *Mind* 99, 185–206.
- 498 Dawid, A.P., 2000. Causal Inference without Counterfactuals. *J. Am. Stat. Assoc.* 95, 407–424.
- 499 Fodor, J.A., 1974. Special sciences (or: the disunity of science as a working hypothesis). *Synthese* 28,
500 97–115.
- 501 Fried, I., Mukamel, R., Kreiman, G., 2011. Internally Generated Preactivation of Single Neurons in
502 Human Medial Frontal Cortex Predicts Volition. *Neuron* 69, 548–562.
- 503 Haggard, P., Eimer, M., 1999. On the relation between brain potentials and the awareness of voluntary
504 movements. *Exp. Brain Res.* 126, 128–133.
- 505 Harris, S., 2012. *Free will*. Free Press, New York.
- 506 Hohwy, J., 2009. The neural correlates of consciousness: New experimental approaches needed?
507 *Conscious. Cogn.* 18, 428–438.
- 508 Holland, P.W., 1986. Statistics and Causal Inference. *J. Am. Stat. Assoc.* 81, 945–960.
- 509 Jensen, M., Dong, M., Vinding, M.C., Overgaard, M., 2017. Measuring sensation of movement. In:
510 Grünbaum, T., Christensen, M.S. (Eds.), *Sensation of Movement, Current Issues in*
511 *Consciousness Research*. Routledge, New York, NY.
- 512 Keller, I., Heckhausen, H., 1990. Readiness potentials preceding spontaneous motor acts: voluntary vs.
513 involuntary control. *Electroencephalogr. Clin. Neurophysiol.* 76, 351–361.
- 514 Kim, C.-Y., Blake, R., 2005. Psychophysical magic: rendering the visible ‘invisible.’ *Trends Cogn. Sci.*
515 9, 381–388.
- 516 Kim, J., 2005. *Physicalism, or something near enough*, Princeton monographs in philosophy. Princeton
517 University Press, Princeton, N.J.
- 518 Laureys, S., 2005. The neural correlate of (un)awareness: lessons from the vegetative state. *Trends*
519 *Cogn. Sci.* 9, 556–559.

- 520 Lewis, D., 1994. Reduction of Mind. In: *A Companion to Philosophy of Mind*. Blackwell Publishers,
521 Oxford, pp. 412–431.
- 522 Libet, B., 1985. Unconscious cerebral initiative and the role of conscious will in voluntary action.
523 *Behav. Brain Sci.* 8, 529–566.
- 524 Libet, B., 1999. Do we have free will? *J. Conscious. Stud.* 6, 47–57.
- 525 Libet, B., Wright, E.W., Gleason, C.A., 1983. Preparation-or intention-to-act, in relation to pre-event
526 potentials recorded at the vertex. *Electroencephalogr. Clin. Neurophysiol.* 56, 367–372.
- 527 Lipton, P., 2005. Inference to the best explanation, 2nd ed. ed, *International library of philosophy*.
528 Routledge, London.
- 529 Logothetis, N.K., Pauls, J., Augath, M., Trinath, T., Oeltermann, A., 2001. Neurophysiological
530 investigation of the basis of the fMRI signal. *Nature* 412, 150–157.
- 531 Mele, A.R., 2009. *Effective intentions: the power of conscious will*. Oxford University Press, Oxford ;
532 New York.
- 533 Melnyk, A., 1997. How to Keep the “Physical” in Physicalism. *J. Philos.* 94, 622.
- 534 Miller, J., Schwarz, W., 2014. Brain signals do not demonstrate unconscious decision making: An
535 interpretation based on graded conscious awareness. *Conscious. Cogn.* 24, 12–21.
- 536 Moors, A., De Houwer, J., 2006. Automaticity: A Theoretical and Conceptual Analysis. *Psychol. Bull.*
537 132, 297–326.
- 538 Nachev, P., Kennard, C., Husain, M., 2008. Functional role of the supplementary and pre-
539 supplementary motor areas. *Nat. Rev. Neurosci.* 9, 856–869.
- 540 Overgaard, M., Gallagher, S., Ramsøy, T.Z., 2008. An integration of first-person methodologies in
541 cognitive science. *J. Conscious. Stud.* 15, 100–120.
- 542 Pernu, T.K., 2011. Minding matter: how not to argue for the causal efficacy of the mental. *Rev.*
543 *Neurosci.* 22.
- 544 Reingold, E.M., Merikle, P.M., 1988. Using direct and indirect measures to study perception without
545 awareness. *Percept. Psychophys.* 44, 563–575.
- 546 Rizzolatti, G., Luppino, G., 2001. The Cortical Motor System. *Neuron* 31, 889–901.
- 547 Rothkirch, M., Hesselmann, G., 2017. What We Talk about When We Talk about Unconscious
548 Processing – A Plea for Best Practices. *Front. Psychol.* 8.
- 549 Rubin, D.B., 1974. Estimating causal effects of treatments in randomized and nonrandomized studies.
550 *J. Educ. Psychol.* 66, 688–701.
- 551 Sandberg, K., Bibby, B.M., Timmermans, B., Cleeremans, A., Overgaard, M., 2011. Measuring
552 consciousness: Task accuracy and awareness as sigmoid functions of stimulus duration.
553 *Conscious. Cogn.* 20, 1659–1675.
- 554 Schlegel, A., Alexander, P., Sinnott-Armstrong, W., Roskies, A., Tse, P.U., Wheatley, T., 2013.
555 Barking up the wrong tree: readiness potentials reflect processes independent of conscious will.
556 *Exp. Brain Res.*
- 557 Schultze-Kraft, M., Birman, D., Rusconi, M., Allefeld, C., Görden, K., Dähne, S., Blankertz, B.,
558 Haynes, J.-D., 2016. The point of no return in vetoing self-initiated movements. *Proc. Natl.*
559 *Acad. Sci.* 113, 1080–1085.
- 560 Seth, A.K., Baars, B.J., Edelman, D.B., 2005. Criteria for consciousness in humans and other
561 mammals. *Conscious. Cogn.* 14, 119–139.
- 562 Shadmehr, R., Krakauer, J.W., 2008. A computational neuroanatomy for motor control. *Exp. Brain*
563 *Res.* 185, 359–381.
- 564 Smart, J.J.C., 1959. Sensations and Brain Processes. *Philos. Rev.* 68, 141.

565 Smart, J.J.C., 1978. The Content of Physicalism. *Philos. Q.* 28, 339.
566 Sperry, R.W., 1980. Mind-brain interaction: Mentalism, yes; dualism, no. *Neuroscience* 5, 195–206.
567 Stoljar, D., 2001. Two Conceptions of the Physical. *Philos. Phenomenol. Res.* 62, 253.
568 Stuart, E.A., 2010. Matching Methods for Causal Inference: A Review and a Look Forward. *Stat. Sci.*
569 25, 1–21.
570 Tononi, G., Koch, C., 2015. Consciousness: here, there and everywhere? *Philos. Trans. R. Soc. B Biol.*
571 *Sci.* 370, 20140167–20140167.
572 Vinding, M.C., Jensen, M., Overgaard, M., 2014. Distinct electrophysiological potentials for intention
573 in action and prior intention for action. *Cortex* 50, 86–99.
574 Wegner, D.M., 2002. *The illusion of conscious will.* MIT Press, Cambridge, Mass.
575 Woodward, J., 2005. *Making things happen: a theory of causal explanation,* Oxford studies in
576 philosophy of science. Oxford University Press, Oxford.
577 Woodward, J., 2012. Causation: Interactions between Philosophical Theories and Psychological
578 Research. *Philos. Sci.* 79, 961–972.
579 Woodward, J., 2015. Interventionism and Causal Exclusion. *Philos. Phenomenol. Res.* 91, 303–347.
580 Woodward, J., Hitchcock, C., 2003. Explanatory Generalizations, Part I: A Counterfactual Account.
581 *Nous* 37, 1–24.
582