

Georg Northoff

### 5.1 Background: The History of Neurophilosophy

Neurophilosophy stands for the investigation of philosophical questions in the context of a neuroscientific hypothesis.

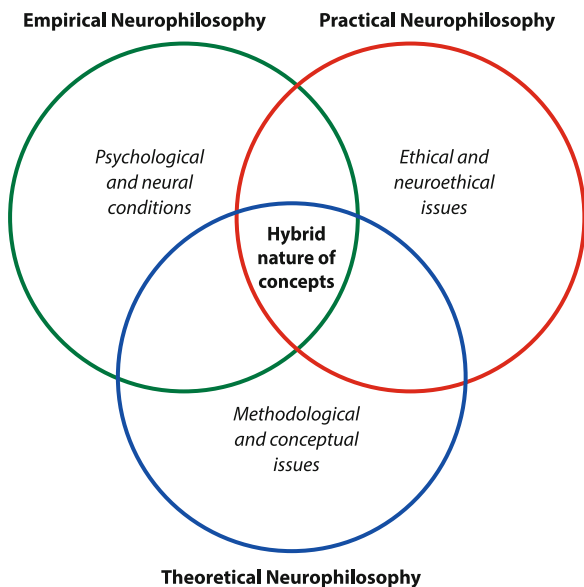
Recent neuroscientific progress has led to the extension of neuroscience to apply and include also concepts like consciousness, free will, self, etc. that were originally discussed in philosophy. This has led to the recent emergence of a new field – neurophilosophy. The term “neurophilosophy” is often used either implicitly or explicitly for the characterization of an investigation of philosophical theories in relation to neuroscientific hypothesis. According to Breidbach [1], pp. 393–394, “neurophilosophy” had already been implicitly practiced at the turn of the last century by W. Wundt (1832–1920), for instance. Another neurophilosopher, though not named as such, was Schopenhauer who was probably the first philosopher to introduce the concept of the brain in the philosophical context. The French philosopher M. Merleau-Ponty (1908–1961) may also be considered a neurophilosopher since in his ‘Phenomenology of perception’ he

explicitly introduces the brain and its neural organisation and links it to perception and other originally philosophical concepts.

Other important developments in this regard were put forward by the American philosopher W. von Orman Quine (1908–2000): He raised the question whether what we can know about ourselves and the world as usually dealt with in the philosophical discipline of epistemology can be traced back to nature itself and ultimately to evolution [2]. This was complemented by the collaboration between the philosopher K. Popper (1902–1994) and the neuroscientist J.C. Eccles (1903–1997) who discussed the relation between brain and mind from both perspectives, neuroscientifically and philosophically [3]. Finally, the term ‘neurophilosophy’ was explicitly coined by the American philosopher P. Churchland [4] in her book ‘Neurophilosophy’ where she discussed empirical results side by side with theoretical issues.

The current field of neurophilosophy covers mainly three different domains, ‘Empirical Neurophilosophy’, ‘Practical Neurophilosophy’, and ‘Theoretical Neurophilosophy’. ‘Empirical Neurophilosophy’ describes the “application of neuroscientific concepts to traditional philosophical questions” [5], p. 1. Here concepts like consciousness, self, and free will (see below for details) that have traditionally been dealt with theoretically in philosophy are now investigated experimentally in neuroscience. Secondly, there is the field of ‘Practical Neurophilosophy’ that deals with ethical concepts like free will, moral judgment, and informed consent in the neural context of the brain. Thereby, as in empirical neurophilosophy, the philosophical-ethical concepts may also be extended from the originally purely human domain to animals, like

G. Northoff  
Canada Research Chair for Mind, Brain Imaging, and Neuroethics, Michael Smith Chair for Neuroscience and Mental Health, University of Ottawa Institute of Mental Health Research,  
Ottawa, Canada  
e-mail: georg.northoff@theroyal.ca



**Fig. 5.1** The figure illustrates the three main domains of neurophilosophy, empirical, theoretical, and practical. Empirical neurophilosophy is concerned with the search for the neural and psychological conditions of originally philosophical terms like self, consciousness, free will, etc. Theoretical neurophilosophy is about the methodological and conceptual issues when linking neuroscientific data/facts and philosophical concepts. Finally, practical neurophilosophy is about the linkage between neuroscience and ethics with ethical issues in neuroscience and neuroscientific mechanisms underlying ethical concepts

57 whether the latter have free will or not. Third, and  
 58 finally, there is ‘Theoretical Neurophilosophy’ which  
 59 focuses on methodological issues like how to link  
 60 empirical data and theoretical concepts in neurophilosophical investigation (Fig. 5.1).  
 61

62 **5.2 Empirical Neurophilosophy –**  
 63 **Experimental Investigation**  
 64 **of Philosophical Concepts**

65 One of the main originally philosophical concepts  
 66 investigated in neuroscience is consciousness. What is  
 67 consciousness? **Consciousness** is often understood as  
 68 the ability to detect, evaluate and report about the experience of a particular object or event in the environment or the own thoughts. Since detection, reporting, and evaluating requires access to the content in question, this form is often called ‘access consciousness’.

73 ‘Access consciousness’ must be distinguished from  
 74 the experience itself which, following philosophers  
 75 like Th. Nagel [6], can be characterized by a particular

point of view, a stance in the world, from which we  
 76 perceive and experience ourselves and others. That  
 77 form of consciousness has been described as ‘phenomenal  
 78 consciousness’. The distinction between ‘phenomenal  
 79 and access consciousness’ is considered by  
 80 many a core distinction which has also aroused plenty  
 81 of controversy. Some authors deny for instance that  
 82 phenomenal consciousness can be distinguished from  
 83 access consciousness. However, animals may have  
 84 phenomenal consciousness while they may remain  
 85 unable to report the contents of their consciousness  
 86 thus lacking access consciousness.  
 87

Following Christoph Koch and Francis Crick [7] we  
 88 need to identify what they call the ‘neural correlates of  
 89 consciousness’ (NCC). The NCC describe the search  
 90 for those minimally neuronal conditions that are jointly  
 91 sufficient for any one specific conscious, i.e., phenomenal,  
 92 percept that we can experience. Several neuronal  
 93 mechanisms have been discussed as possible candidate  
 94 mechanisms for the NCC. In the following I highlight  
 95 some of the main and most popular suggestions.  
 96

G. Edelman [8] considers cyclic processing and  
 97 thus circularity within the brain’s neural organisation  
 98 as central for constituting consciousness. Cyclic  
 99 processing describes the re-entrance of neural activity  
 100 in the same region after looping and circulating in  
 101 so-called re-entrant (or feedback) circuits.  
 102

This is for instance the case in primary visual cortex  
 103 (V1): The initial neural activity in V1 is transferred to  
 104 higher visual regions such as the inferotemporal cortex  
 105 (IT) in feedforward connections. From there it is conveyed  
 106 to the thalamus which relays the information back  
 107 to V1 and the other cortical regions implying thalamo-  
 108 cortical re-entrant connections. Consciousness is assumed  
 109 to be constituted on the basis of such feedback or re-  
 110 entrant connections that allow for cyclic processing.  
 111

What is the exact mechanism of the feedback or re-  
 112 entrant circuits? Re-entrant circuits integrate information.  
 113 This leads Giulio Tononi to emphasize the  
 114 integration of information as the central neuronal  
 115 mechanism in yielding consciousness. He consecutively  
 116 developed what he calls ‘Integrated Information Theory’  
 117 (IIT). We usually focus on the content that is selected  
 118 to become conscious, i.e., ‘what is perceived’.  
 119 Instead, as the IIT claims, we may better search for  
 120 the neuronal mechanisms that allow excluding content  
 121 from becoming conscious, i.e., ‘what is ruled out’. The  
 122 information that is ruled out to become conscious may  
 123 suffer from insufficient integration of information and  
 124 remains therefore unconscious.  
 125

126 Tononi assumes the integration of information to be  
 127 particularly related to the thalamo-cortical re-entrant  
 128 connections: These re-entrant connections process all  
 129 kinds of stimuli thus remaining unspecific with regard  
 130 to the selected content. They make it possible to  
 131 generate a particular point of view and an associated  
 132 quality of experience (also called qualia) as hallmark  
 133 feature of consciousness. Linkage of these qualia to  
 134 the content processed via thalamo-cortical information  
 135 integration may then allow these contents to become  
 136 conscious. This distinguishes them from the uncon-  
 137 scious contents that do not undergo such cyclic pro-  
 138 cessing via the thalamus – and therefore the addition of  
 139 the specific quality, the qualia, remains impossible.

140 Another suggestion for the neural correlate of consci-  
 141 ousness comes from B. Baars [9, 10] and others like  
 142 S. Dehaene. They assume global distribution of neural  
 143 activity across many brain regions in a so-called global  
 144 workspace to be central for yielding consciousness:  
 145 The information and its contents processed in the brain  
 146 must be globally distributed across the whole brain in  
 147 order for them to become associated with  
 148 consciousness.

149 When information is only processed locally within  
 150 a particular region but not throughout the whole brain,  
 151 it can not be associated with consciousness anymore.  
 152 The main distinction between unconsciousness and  
 153 consciousness is thus supposed to be manifest in the  
 154 difference between local and global distribution of  
 155 neural activity. Hence, the global distribution of neural  
 156 activity is here considered a sufficient condition and  
 157 thus neural correlate of consciousness.

158 Taken together, there are currently these neuroscientific  
 159 suggestions for consciousness. Future research is needed  
 160 though to further specify the neuronal mechanisms  
 161 themselves and the features of consciousness itself.  
 162 Consciousness may by itself not be as homogenous as it  
 163 appears; instead, it may be characterized by different  
 164 features as for instance a point of view (see above), a  
 165 quality (see above), and a particular unity as unifying  
 166 convergence point for different contents.

167 Another originally philosophical concept now hotly  
 168 debated in neuroscience is the concept of the *self*. The  
 169 question of the self has been one of the most salient  
 170 problems throughout the history of philosophy and  
 171 more recently also in psychology and neuroscience.  
 172 For example, William James (1842–1910) distinguished  
 173 between a physical self, a mental self, and a spiritual  
 174 self. These distinctions seem to reappear in recent con-  
 175 cepts of self as discussed in neuroscience. Damasio

[11] and Panksepp [12] suggest a “proto-self” in the  
 sensory and motor domains, respectively, which resem-  
 bles James’ description of the physical self. Similarly,  
 what has been described as “minimal self” [13, 14] or  
 “core or mental self” [11] might correspond more or  
 less to James’ concept of mental self. Finally, Damasio’s  
 “autobiographical self” and Gallagher’s “narrative self”  
 strongly rely on linking past, present, and future events  
 with some resemblances to James’ spiritual self.

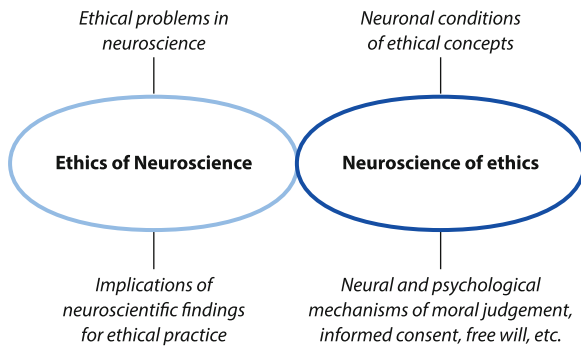
These distinct selves are now related to distinct  
 brain regions. For instance, the “proto-self” outlining  
 one’s body in strongly affective and sensory-motor  
 terms is associated with subcortical regions like the  
 periaqueductal gray, the colliculi, and the tectum. The  
 “core or mental self” building upon the “proto-self” in  
 mental terms is associated more with the thalamus and  
 cortical regions like the ventromedial prefrontal cortex  
 (see, for instance [11, 15]). Finally, the “autobiograph-  
 ical or extended self” that allows one to reflect upon  
 one’s “proto-self” and “core or mental self” is associ-  
 ated with cortical regions like the hippocampus and the  
 cingulate cortex.

Humans show various cortical regions, predomi-  
 nantly the so-called cortical midline structures (CMS),  
 to be involved in what is called self-related processing  
 (SRP) that are integrated with subcortical processes to  
 yield an integrated subcortical-cortical midline system  
 (SCMS). The lowest regions of this distributed SCMS  
 network include the periaqueductal gray, the superior  
 colliculi, and the adjacent mesencephalic locomotor  
 region as well as preoptic areas, the hypothalamus, and  
 dorsomedial thalamus, while cortical regions include  
 the ventro- and dorsomedial prefrontal cortex, the pre-  
 and supragenual anterior cingulate cortex and the pos-  
 terior cingulate cortex, and the medial parietal cortex.  
 The association of the subcortical regions with a sense  
 of self has led to the assumption that already animals  
 may have a sense of self [16, 17] though most likely  
 not as cognitively elaborated as the human self.

---

**5.3 Theoretical Neurophilosophy – Methodology and Knowledge of the Linkage Between Brain Data and Philosophical Concepts**

One of the main issues in neurophilosophy is the ques-  
 tion for methodology. How can we link empirical data,  
 so-called facts as obtained in neuroscience, to the con-  
 cepts and their meaning as dealt with in philosophy?



**Fig. 5.2** The figure illustrates the fields of practical neurophilosophy. Ethics of neuroscience concerns ethical problems in neuroscience like informed consent. While neuroscience of ethics refers to the neural and psychological conditions of ethical concepts

223 Data and facts do not require any definition and deter-  
 224 mination. They rely on observation and can in princi-  
 225 ple be obtained by anybody; they thus remain  
 226 investigator-independent. This is different in the  
 227 domain of concepts. Concepts carry a meaning, a  
 228 semantic dimension, which may be closely related to  
 229 the investigator and how he defines and uses the con-  
 230 cept in question.

231 However, empirical-experimental investigation can-  
 232 not do without concepts. For instance, in formulating  
 233 the hypotheses of the experiment as starting point for  
 234 developing an appropriate experimental design, con-  
 235 cepts play a substantial role. And after obtaining the  
 236 data they must be interpreted for which again concepts  
 237 are necessary. This concerns only concepts within the  
 238 natural world, the world we live in, and thus what phi-  
 239 losophers call the ‘natural conditions’. Such ‘natural  
 240 conditions’ must be distinguished from ‘logical condi-  
 241 tions’ that describe logically possible worlds which  
 242 may or may not be realized within the context of our  
 243 current natural world.

244 The neurophilosopher is thus confronted with the  
 245 principal gap between data/facts and concepts in a  
 246 twofold manner. First, there is the gap between data/  
 247 facts and concepts within the domain of the natural  
 248 world: How do certain data about, for instance, the  
 249 reward system in animals stand to the concept of  
 250 reward in general in both animals and humans? This is  
 251 a gap the neuroscientists themselves already face  
 252 which, due to the predominant experimental focus, is  
 253 often neglected. Secondly, there is the gap between  
 254 neuroscientific data/facts in the natural world and the

concepts in the logical worlds of the philosophers: 255  
 How can we infer from neuroscientific data about con- 256  
 sciousness to the philosophical concept of conscious- 257  
 ness and, vice versa, how can we translate the latter 258  
 into experimental designs to test it empirically? This is 259  
 a truly neurophilosophical gap which we need to bridge 260  
 if neurophilosophy is to succeed in both methodology 261  
 and knowledge. 262

**5.4 Practical Neurophilosophy – Neuroethics and the Relevance of Ethical Concerns in Neuroscience**

263 Practical neurophilosophy or neuroethics focuses, on 266  
 the one hand, on the investigation of the psychological 267  
 and neural conditions of ethical concepts like free will, 268  
 decision making, moral judgment, and informed con- 269  
 sent. This can be described as ‘neuroscience of ethics’. 270  
 At the same time, practical neurophilosophy also deals 271  
 with ethical problems in neuroscience and thus with 272  
 issues of validity of informed consent in psychiatric 273  
 patients, enhancement of cognitive functions by 274  
 neuroscientific interventions, coincidental findings in 275  
 neuroimaging (Fig. 5.2). That amounts to an ‘ethics of 276  
 neuroscience’ [18]. 277

278 Do we have a free will or not? The free will is, for 278  
 instance, manifest in our daily decisions if, for instance, 279  
 we choose the red rather than the green apples in the 280  
 supermarket. Recent neuroscience detected the neural 281  
 mechanisms of decision making that seem to involve a 282  
 number of different brain regions including those 283  
 where reward is processed. The reward regions include 284  
 the ventral tegmental area (VTA), the ventral striatum 285  
 (VS), and the ventromedial prefrontal cortex (VMPFC). 286  
 All these regions have their homologs in various animal 287  
 species so that the same question for the free will 288  
 may also be extended from human to animals. 289

290 The observation that the apparently free decision 290  
 making is related to and, in fact, temporally preceded 291  
 by neural activity specifically related to the decision in 292  
 question has put the concept of free will in doubt. If the 293  
 free will is pre-determined by the neural events in the 294  
 brain, one can no longer speak of a free will. The free 295  
 will is then no longer free but nothing but a mere illu- 296  
 sion on our side with the brain determining our actions 297  
 and decisions. Are we thus no longer free in our will? 298  
 That obviously is an interpretation of the data and also 299  
 depends on the definition of the concept of free will. If 300



301 one, for instance, presupposes a narrow concept of free  
302 will that excludes any preceding changes, the present  
303 brain data may tell us that there is indeed no free will.  
304 Brain data and free will are then incompatible.  
305 Conversely, a wider concept of free will that does not  
306 exclude preceding neural activity changes may then be  
307 well compatible with the brain data.

308 The debate about free will pertains to a wider issue,  
309 the question of determinism versus indeterminism.  
310 Determinism assumes that all our decision and also  
311 what we call free will is determined completely and  
312 exclusively by the brain and its neural activity. Our  
313 person or our self, as presumably distinct from the  
314 brain, has then no say at all in our decision. Hence,  
315 it is then the brain rather than the self that makes  
316 the decision and has a 'neuronal will' rather than a  
317 'free will'. That however is countered by indeterminism.  
318 Indeterminism argues that the brain itself and its  
319 neural activity changes does not determine completely  
320 and exclusively our decision making so that there are  
321 traces of free will left in our decisions. Who is right,  
322 determinism or indeterminism? As said above, it may  
323 strongly depend not only on the data but also on the  
324 conceptual definitions.

325 Besides such questions belonging to the 'neurosci-  
326 ence of ethics', the neuroscientific investigation of  
327 ethical concepts, there are also issues pertaining to  
328 ethical problems in neuroscience. One problem here is,  
329 for instance, the one of informed consent which sub-  
330 jects have to give when participating in experimental  
331 investigations. Being able to give informed consent  
332 may include a variety of different functions, cognitive,  
333 social, and affective, that are all ultimately brain-based.  
334 Does this mean that we have to exclude those subjects  
335 that suffer from impairments in these functions?  
336 Furthermore, recent research demonstrates that ani-  
337 mals possess many of the cognitive and social func-  
338 tions originally attributed to humans only. Do we  
339 therefore need to develop more rules for animal par-  
340 ticipation in research by, for instance, considering that  
341 they can have consciousness, feel pain, and empathize  
342 with co-species?

## 343 5.5 Summary

344 Neurophilosophy is a young and novel field right at the  
345 intersection between neuroscience and philosophy.  
346 Unlike more established disciplines, it has not yet an

established method that needs to be developed in the 347  
future as part of a 'theoretical neurophilosophy'. At the 348  
same time though neurophilosophy is a highly promis- 349  
ing field which will be able to provide novel answers to 350  
questions discussed in philosophy for more than 351  
3,000 years. This will not only enrich neuroscience 352  
and provide new ideas for experimental designs but 353  
will also change and reverberate in philosophy itself 354  
by allowing for a shift from the hitherto mind-based 355  
philosophy to a more brain-based neurophilosophy. 356

## References

- 357
1. Breidbach O (1997) Die Materialisierung des Ichs – Eine 358  
Geschichte der Hirnforschung im 19. und 20. Jahrhundert. 359  
Suhrkamp, Frankfurt aM 360
  2. Quine WO (1969) Epistemology naturalized. In: Quine WO 361  
(ed) *Ontological relativity and other essays*. Columbia 362  
University Press, New York 363
  3. Popper K, Eccles J (1989) *Das Ich und sein Gehirn*. Piper, 364  
München 365
  4. Churchland P (1986) *Neurophilosophy: toward a unified 366  
science of the mind-brain*. MIT Press, Cambridge, MA 367
  5. Bickle J, Mandik P, Landreth A (2006) The philosophy of 368  
neuroscience. *Journal [serial on the Internet]*. Available 369  
from: <http://plato.stanford.edu/entries/neuroscience/> 370
  6. Nagel T (1979) What is it like to be a bat? *Mortal questions*. 371  
Cambridge University Press, New York, p 166 372
  7. Koch C, Crick F (2001) The zombie within. *Nature* 411:893 373
  8. Edelman GM (2003) Naturalizing consciousness: a theoret- 374  
ical framework. *Proc Nat Acad Sci U S A* 100:5520–5524 375
  9. Baars BJ (2005) Global workspace theory of consciousness: 376  
toward a cognitive neuroscience of human experience. *Prog 377  
Brain Res* 150:45–53 378
  10. Baars BJ, Franklin S (2007) An architectural model of con- 379  
scious and unconscious brain functions: *Global Workspace 380  
Theory and IDA*. *Neural Netw* 20:955–961 381
  11. Damasio A (1999) How the brain creates the mind. *Sci Am* 382  
281:112–117 383
  12. Panksepp J (2003) At the interface of the affective, behav- 384  
ioral, and cognitive neurosciences: decoding the emotional 385  
feelings of the brain. *Brain Cogn* 52:4–14 386
  13. Gallagher S (2000) Philosophical conception of the self: 387  
implications for cognitive science. *Trends Cogn Sci* 388  
4:14–21 389
  14. Gallagher HL, Frith CD (2003) Functional imaging of 'the- 390  
ory of mind'. *Trends Cogn Sci* 7:77–83 391
  15. Damasio A (2003) Feelings of emotion and the self. *Ann N 392  
Y Acad Sci* 1001:253–261 393
  16. Northoff G, Panksepp J (2008) The trans-species concept of 394  
self and the subcortical-cortical midline system. *Trends 395  
Cogn Sci* 12:259–264 396
  17. Northoff G (2004) *Philosophy of the brain*. John Benjamins 397  
Publishing, Amsterdam 398
  18. Roskies A (2002) Neuroethics for the new millenium. 399  
*Neuron* 35:21–23 400