

# Bachelor Degree Project



UNIVERSITY  
OF SKÖVDE

## **IS THERE A YOU IN YOUR BRAIN?**

The Neuroscientific Support for the Bundle Theory-  
View of the Nature of the Self

Bachelor Degree Project in Cognitive Neuroscience

Basic level 22.5 ECTS

Spring term 2019

Amanda Vestin

Supervisor: Paavo Pylkkänen

Examiner: Katja Valli

## Abstract

Why do you experience yourself as a continuous self? This is a central question when regarding the self and it has two kinds of answers: either there is something like an ego inside you which is the entity perceiving all your experiences (the ego theory-view), or there is no such thing as a self or an ego and you are just a collection of different perceptions (the bundle theory-view). There are many different components all contributing to the concept of self as a whole leading to different neuroscientific ways of measuring it and some researchers are arguing for the nonexistence of a unified self-system within the brain. The aim of this thesis is to review how neuroscientific findings might contribute to the philosophical debate about the nature of self. The thesis starts off by reviewing the different concepts and components with which the self is typically described, both in philosophy and in the empirical research field of neuroscience. Then follows a presentation of three important aspects of self-awareness – first-person perspective, self-reflection, and interoception – and their specific associated brain areas (namely, the medial prefrontal cortex, the posterior and anterior cingulate cortices, and insula). The purpose here is to examine how the self is approached in these studies. After this the thesis explores to what extent neuroscience supports the bundle theory-view, with a focus on reviewing the different brain networks involved in the processing of self. In conclusion, the thesis suggests that the literature reviewed provides neuroscientific support for the bundle theory-view that there is no unified self located in the brain, mostly because of the dissimilar neural activations associated with different self-related processes. In other words, the bundle theory seems to be correct despite the experienced feeling you have of being a continuous and unified self.

*Keywords:* self, neural basis of self, bundle theory, no unified self

**Table of Contents**

1. INTRODUCTION	4
1.1. Aim and Structure	5
2. PHILOSOPHICAL AND THEORETICAL ASPECTS OF THE SELF	6
2.1. Subjective Experiences	7
2.2. Ego versus Bundle Theories	8
2.3. The Sense of Self	9
2.3.1. The 'me' and the 'I'	9
2.3.2. Damasio and Parfit	11
2.3.3. The core self	12
2.3.4. Self-consciousness	13
2.3.5. The doubts	14
2.4. Personal Identity	14
3. THE SELF IN NEUROSCIENTIFIC RESEARCH	16
3.1. The First-Person Perspective	16
3.2. Self-Reflection	18
3.3. Interoception	21
4. IS THE BUNDLE THEORY OF SELF SUPPORTED BY NEUROSCIENCE?	23
4.1. No Self-Specific System	23
4.2. No Unified Self-System	28
4.3. The Self as an Illusion	30
4.4. Defending the Unified Self	32
5. DISCUSSION	33
5.1. Problems and Limitations	36
5.2. Conclusion	37
6. REFERENCES	38

## 1. INTRODUCTION

Let us talk about the self. It might be considered 'natural' to have some kind of sense of what you are when living your life in the external world. You have a sense of yourself and you know what you mean and what you are referring to when using the pronoun 'I'. It feels like there is something continuous that makes up what you call *your* self and that this self is something existing in the world. Maybe you have heard about the skeptical worry that the external world might not exist and if there is something really existing, it just has to be the self. Maybe you think of the self as a soul-like substance, existing in addition to the physical body, or at least some kind of 'ego' (Latin for *self*) which is having the subjective experience of your perceptions (this is called the ego theory-view of the self; Parfit, 1987).

Or maybe you think of the self as something consisting of different parts or components, like your attitudes, personality traits, and all your memories, in addition to the bodily perceptions and sensations, all of them contributing to your experience of a self as a whole, while also believing that the self is not an existing thing in itself. This is, in contrast to the ego theory-view, the bundle theory-view of the nature of self (Parfit, 1987). The bundle theory-view is based on the thought of Scottish philosopher David Hume (1739/1896), who argued that the feeling of a self is illusory and that you are only a bundle of different perceptions and sensations.

The problem of the self has been around in philosophy for a long time, and nowadays it also plays an important role in the field of neuroscience when studying consciousness and the mechanisms of conscious and subjective experience. It is often assumed that there has to be some unifying process for the self in the brain integrating the physical and mental elements to one whole since this is the most intuitive view on this complex, multifaceted phenomenon called the self. The existence and realness of the self might be considered intuitive since it actually *feels* like there is something having the subjective experiences and perceiving the world. It feels like there is someone inside perceiving all your experiences and the world around you, and that there has to be a perceiver

for those experiences to exist (Blackmore, 2010). However, the German philosopher Thomas Metzinger (2009, 2011) has argued for the opposite view, a no-self alternative, in which the experience of a unified self is only an illusion, a simulation created in our brains, and that the self is not something existing. So the question is why you so strongly experience yourself as a continuous self?

There are several different components of the self leading to a wide debate about the concept of self in several different disciplines, such as philosophy, psychology, and cognitive science. Across these disciplines, there are so many different conceptions that a general consensus about the definition is far from reached (Legrand & Ruby, 2009). The different ways to approach the sense of self create difficulties when trying to measure the self, and additionally, there are researchers arguing that there is no such thing as a unitary self-system in the brain (e.g. Gillihan & Farah, 2005). However, the field of neuroscience continues to measure the self in different ways and approaches it by different tasks and self-related processes, such as reflection of one's own personality traits or memory recall of personal events, and tries to find the neural mechanisms of these processes. Therefore, the remaining question is: can the neuroscientific findings from studies on the different components of the self and their associated brain areas still say something about the nature of self, even though there seems to be no unitary self-system inside the brain associated with your unified experience of a self?

### **1.1. Aim and Structure**

The aim of this thesis is to review the existing empirical neuroscientific research on the sense of self and examine how this might contribute to the philosophical debate about whether the self is something like an ego or if it rather is just a collection of sensations. The focus will be on the bundle theory-view on the self and the thoughts of Hume (1739/1896), arguing for the self as just an illusion and not a thing itself and that you would rather be just a bundle of different perceptions, such as heat or cold, pain or pleasure. According to this view, the self itself without any other

perceptions is nothing you can ever observe. Therefore, the main research question is whether you can say, based on neuroscientific findings, that the bundle theory-view of the nature of self is more correct than the ego theory, despite the counterintuitive characteristic of the former view.

To realize this aim, the second chapter of the present thesis will give an overview of the ego and bundle theories of self, the thoughts of Hume, the sense of self and personal identity. Chapter three will focus on the neural basis of some of the different components used to measure self-awareness in neuroscientific research. This will be done by reviewing the specific brain areas involved in the processes of first-person perspective and self-reflection, in addition to the physical aspect of interoception (i.e. the awareness of the internal environment of one's physical body). After this, chapter four will consider to what extent neuroscientific evidence supports the view that there is no unified self, focusing on the functions of specific brain networks. The thesis will end with a discussion of how these self-components and their associated brain mechanisms are related to the philosophical debate about the nature of self and especially the bundle theory-view, with a concluding answer to the research question.

An important limit, however, is that the relation between self-related processing and resting-state activity in the default mode network (DMN) is not covered in this thesis, due to scope restrictions and the limited amount of time. This is especially noted because of the similar neural activation patterns associated with both self-processing and resting-state.

## **2. PHILOSOPHICAL AND THEORETICAL ASPECTS OF THE SELF**

The second chapter of this thesis will provide a brief overview of the philosophical and theoretical aspects of the self to give a background for the rest of the thesis. It will present some basic thoughts of subjective experiences, the definitions of the ego and bundle theories of the nature of self, and some concepts used when discussing the self and self-consciousness, including a brief and initial paragraph about the doubts regarding an existing self in the brain (which will be

reviewed more extensively in chapter four). This chapter will end with describing personal identity and the problem of our persistence over time.

## 2.1. Subjective Experiences

According to the American philosopher Thomas Nagel (1974), conscious experience is a widespread phenomenon since it may occur in many different levels of animal life. He argues that the subjective character of experience is that there is something it is like to be an organism, and this follows the fact that this organism has conscious experience. And "fundamentally an organism has conscious mental states if and only if there is something that it is like to *be* that organism—something it is like *for* the organism" (Nagel, 1974, p. 436). There is something that it is like to be a specific organism, and how it feels, subjectively, to be that organism can only be experienced by the organism itself (Nagel, 1974). One example is the experiences of a bat, and Nagel means that you can never experience how it is like to be a bat without actually being the bat – you can never have the subjective experiences of a self without *being* that particular self and experience its experiences from the first-person perspective. Nagel (1974) says that since you are limited to the resources of your own mind, you cannot even imagine the subjective experiences of a bat.

Probably most of you might agree that the answer to Susan Blackmore's question 'who is conscious now?' would be something that you can call yourself and that you would answer 'I am the one who is conscious'. And indeed, it *seems* like there is something or someone there who has the experiences you are experiencing, and that there cannot be any experiences without someone experiencing them – there has to be a perceiver (Blackmore, 2010). The problem arises when you start to think about what kind of thing that experiencer would be. While some accept that what you call your self should be just your body and that there is no need for an inner self, others believe that it has to be something inside you that is not the same as your physical body; something that is in charge and can make decisions, something having the experiences. This is the central question when

talking about the self: why does it seem and feel that you are a continuous, single self with experiences (Blackmore, 2010)?

## 2.2. Ego versus Bundle Theories

There are two types of answers to this central question of why you experience yourself as a continuous self, which the British philosopher Derek Parfit (1987) has described as 'ego theories' and 'bundle theories'. The ego theory-view might be the most intuitive one and easiest to assume since it proposes that a person's continuous existence can only be explained as the existence of some continuous 'ego' or subject of experiences, which makes this view include a belief in the existence of the self. It suggests that what holds the different subjective experiences in different time frames, is that they are all experienced by the same person. All the experiences during one person's life are experienced by the same subject of experience, and this is what explains the unity of that person's life (Parfit, 1987). The most known form of this kind of theory is the Cartesian view (called Cartesian dualism), which means that each person is something purely mental, something like a soul or substance, in addition to the material body. However, believing in the ego theory-view is not the same as believing in dualism (Blackmore, 2010).

In contrast, the bundle theory-view, as described by Parfit (1987), denies the existence of a self; yes, it *seems* as there is a continuous self, but this illusion has to be explained by something else because there is no unified self underlying this feeling. According to this view, the unity of consciousness in the life of an organism cannot be explained by referring to persons. Rather, there is a long series of different experiences and mental states, such as thoughts and sensations, and each of these series is what you call one life (Parfit, 1987). This view denies the existence of persons, and suggests that this is just the language you are using; a person is not something existing by itself, it is not a separately existing thing that exists independently of its body and brain. Rather, persons and subjects only exist in this language-dependent way, when you use your language to talk about them (Parfit, 1987).

This bundle theory-view is based on the term 'bundle of sensations'. David Hume (1739/1896) argued that you are nothing more than a bundle of experiences and that you can never, through introspection, catch yourself at any time without some kind of perception, such as the perception of temperature or pain, and this perception is the only thing that you can observe. All these sensations and perceptions are tied together and unified by processes such as the relations between the actual experiences and your memories of these experiences. In this sense, each series, each life, is comparable to a bundle that is tied up by a string (Parfit, 1987). Since the present thesis is focusing on the relation between neuroscientific research and the bundle theory-view of the self, this view will be briefly referred to in the remaining chapters. But first, some different philosophical components and concepts used when talking about the self are going to be distinguished in what follows.

### **2.3. The Sense of Self**

#### **2.3.1. The 'me' and the 'I'**

Already in the late 19<sup>th</sup> century, the American philosopher and psychologist William James (1890/1950) began to divide the self into two elements; the empirical self or the objective person, the 'me', and the subjective knowing thought, or pure ego, called the 'I'. In the first and widest form of the empirical self, your self is the total sum of all the things you call your own. James (1890/1950) further differentiated the empirical self into three aspects; the material self, the social self, and the spiritual self. He argued that the material self is the most basic part of the self; your body and other material things around you with instinctive importance, such as your clothes and your closest family. The social self, by contrast, is composed of the recognition you get from people in your surroundings; you have as many social selves as there are different individuals around you who recognize you and have an image of you in their mind. The third part is according to James (1890/1950) the spiritual self, and this part is the inner or subjective being and the psychological dispositions which are the most enduring, and also intimate, part of the self.

Furthermore, it is when trying to define the more specific nature of the other element, the subjective knowing thought, or the 'I', diverging opinions emerge. One extreme is the view that it is the active substance of the soul that creates consciousness, whereas the other extreme says that the self is nothing more than the imagined being behind the usage of 'I' (James, 1890/1950). The first extreme is the spiritualistic view, which in the terms of Parfit (1987) is considered an ego theory, whereas the other is a bundle theory, by arguing that the feeling of a continuous and unified self is just an illusion and that there is no continuous being called 'the self' holding all the experiences together.

James argued that there is no separate spirit or ego that holds all your thoughts together, rather, as cited in Blackmore (2010), James meant that at every moment there is a passing thought, and this special thought is what remembers the previous thoughts and therefore holds all the thoughts of an organism together. At the moment after the present one, there will be another thought of this special kind, just passing on your sense of a continuous self. Thus, the unity you experience as the self is not something separate from the thoughts you are thinking (Blackmore, 2010). James compared this to a herdsman with a herd of cattle and argued that in contrast to the common sense feeling that there would be the same herdsman holding the herd together, there is no permanent herdsman; rather, there is only a passing series of owning herdsman, all of which inheriting the previous one's cattle and ownership. In this sense, each thought is transmitting everything realized as itself to the next thought, and this is then creating the experienced sense of unity (as cited in Blackmore, 2010).

Christoff, Cosmelli, Legrand, and Thompson (2011) have also discussed 'I' in contrast to the term 'me'. Their opinion is in line with James' when they define the word 'I' as referring to the experience of yourself as a subjective agent of your perceptions, actions, and emotions. The term 'me', in contrast, refers to the experience of yourself through self-related processing. In other words, the self might be specified either as an *object* of perceptions and attributions (i.e. the 'me'; the perceptions happen to me), or as an experiential *subject* or *agent* of perceptions, actions, and

feelings (i.e. the 'I'; I am the one who has the perceptions; Christoff et al., 2011). In the present thesis, the focus will be on the subjective aspect of the self, i.e. the 'I', since this connects to what Nagel (1974) calls the subjectivity of experience, and it is here the divergence regarding the nature of this subjective agent arises between the ego and bundle theories.

### **2.3.2. Damasio and Parfit**

The neuroscientist Antonio Damasio (2003) defines the term 'self' as referring to an individual consisting of a mind or a body, or the unity of both. Additionally, the self is often considered as something stable and continuous over time. The term is synonymous with 'personhood' (i.e. having the status of being a person), and is always involving a reference; for example when the organism is referring back to the organism itself or to its own mind (Damasio, 2003).

Parfit (1971), on the other hand, does not consider self-identity as something stable and continuous over time, but rather emphasizes the importance of the connections between earlier selves and the present self, and that the identity of a person is in its nature a matter of degree, instead of an all-or-nothing phenomenon. He means that what is important for a person to survive over time as the same person is the continuity of some psychological relations between your past selves and your present self. These relations hold over time to different degrees, depending on the connection between the present self and the past self. An example of such a strong connection is when the present self identifies oneself with a past self and looks at the self with either pride or shame, pleasure or regret. On the other hand, a psychological relation has a weak connection when the present self is not identifying with the past self and regarding the past self with indifference, and therefore feel neither pride nor shame (Parfit, 1971). Therefore, these weak psychological connections are of less importance for the survival of the person over time compared to the strong connections. This psychological continuity is discussed further in the section of 'Personal Identity' below.

### **2.3.3. The core self**

The mental representation of yourself is, according to Damasio (2003), the minimal level of the self that is necessary for the occurrence of consciousness. However, he says that this does not mean that you should look for some kind of mental homunculus (i.e. a representation of some small human being) in the brain and that it probably does not even exist some entity that knows everything, thinks on its own and gives you your sense of self. It is rather the individual's stable representation of its own continuity that serves as a mental reference for the organism in its conscious mind and thus gives the organism something to refer to (Damasio, 2003). This is what Damasio (1999) further calls the 'core self', which is the sense of self that comes from the core consciousness; the simplest kind of consciousness providing the organism with the sense of self in the present moment and space. The core self is based on the preconscious biological phenomenon called the 'proto-self', which is a moment-by-moment representation of the current state of the organism that we are not conscious of. The core self is the transient sense of self, which is continuously refreshed as new information enters from the outside of the brain or from the inside in the form of thoughts or memories. The core self is the first basic form of the conscious sense of self, and based on the core self is then the 'autobiographical' or 'extended self', which is the part of the self giving the organism its identity and awareness of both the past experiences, the present, and the future (Damasio, 1999).

Damasio's (1999) notion of the core self is in line with what is also called the 'minimal self' (Gallagher, 2000). The minimal sense of self is a consciousness of oneself as an immediate subject of experience, and this part of the self is unextended in time. This means that the minimal self is limited to what is reachable in immediate self-consciousness and to the right now-moment. This can further be connected to what Blanke and Metzinger (2009) have called 'phenomenal selfhood', which gives rise to the subjective experiences of being a self, independently of the ability of explicit cognition. Blanke and Metzinger (2009) further argued that the 'minimal phenomenal

selfhood' is related to a representation of the entire body and embodiment (i.e. the role of the body in the shaping of the mind and the subjective experience of having and using a body). This would be the simplest form of consciousness, and what Damasio (1999) calls the core consciousness.

#### **2.3.4. Self-consciousness**

When talking about the self, it becomes necessary to also dig deeper into the level of consciousness (Damasio, 2003), and by this go beyond the level of language and the using of the pronoun 'I' to find some more substantial significance when talking about oneself. Otherwise, the usage of 'I' becomes hollow; if I am not consciously aware of something I refer to when talking about 'I' and what *I* feel, such as a specific pain, this specific pain becomes just an 'objective' pain, and without some deeper meaning in consciousness I would not know that I am the one experiencing this specific pain (Damasio, 2003).

Accordingly, Vogeley and Fink (2003) define the term 'self-consciousness' as the ability to become aware of your own mental and bodily states, such as perceptions and feelings, as your own mental and bodily states and no-one else's. In other words, self-consciousness is fundamentally the conscious awareness of the self as such (J. Smith, 2017). This awareness is following from introspection (i.e. when considering your inner conscious thoughts and feelings), which makes you introspectively aware of yourself and your own mental properties. However, it is not just about the ability to be conscious of oneself, what is important is rather "whether one is, or can be, conscious of oneself *as oneself*, a form of awareness in which it is manifest to one that the object of awareness is oneself" (J. Smith, 2017) leading to some kind of meta-awareness of oneself. This introspective awareness is, however, what Hume (1739/1896) considered impossible since you can never catch what is *you* when turning to your conscious thoughts and feelings because there is always some other kind of perception in the way.

### **2.3.5. The doubts**

Damasio (2003) is doubtful about the possibility to find some special self-center in the brain, where the whole phenomenon of self would magically 'pop up'. This is in line with the standpoint of Musholt (2013), who has raised doubts whether it is possible to locate the self in the brain, and argues that you have to be careful when regarding suggestions about particular brain areas where the self would be located. Additionally, Gillihan and Farah (2005) examined in their review if the self-related processes really are special and if all the different components of the self are the function of a unitary system. Gillihan and Farah (2005) concluded that there is no good evidence for characterizing the self as special (in the sense that for example human language has been said to be special). In addition, in the studies reviewed they found no support for a unitary self-system since the data showed no common implication of brain areas across the different examined components of the self. The empirical support for these doubts is going to be further discussed in chapter four.

### **2.4. Personal Identity**

Personal identity is a concept that emphasizes philosophical questions about ourselves in the terms of being persons (i.e. rational and self-conscious beings), and there are several problems when regarding personal identity (Olson, 2015). Some of these problems are the characteristics you use to define you as you, what the requirements are for someone to be a person compared to a nonperson, and the problem that has received the most attention is the question about your persistence and how you can exist as the same person over time. Olson (2015) defines the question roughly as; what is necessary, and/or sufficient, for a being in the past or future to be the very same being as existing in the present moment? There are, for example, situations and circumstances where you have changed your attitude or way of thinking, such as after some kind of special life experience, that makes you say that you have become a new person or that you are not the same person as you were before.

The answer to the question whether you are the same person as before would be *no*, but according to Olson (2015), the persistence question is asking whether you would still *exist*, and not whether you are the same person or not. The answer to this is, in contrast, *yes*, as long as the now different person is still you, you exist just as if you had not changed to that new person. This is just about the fact that the person has changed in some important way, that some characteristics describing the person have changed, but not about that person's existence, and has because of this nothing to do with the persistence over time (Olson, 2015). Two different sorts of an answer to the persistence question are psychological and physiological continuity, focusing on mental and biological properties, respectively.

According to the psychological continuity views, what is important for the persistence of a person over time are psychological relations; what is determining that you are the same as a past you is that the present you has inherited mental features – such as memories, beliefs, and preferences – from the past you, and in the same way a future you will inherit mental properties from the present you (Olson, 2015). This is connected to the previously mentioned psychological connections (Parfit, 1971) and is called the 'Lockean' view (Parfit, 2012). The English philosopher John Locke (1632-1704) defined a person as "a thinking intelligent being, that has reason and reflection, and can consider itself as itself, the same thinking thing at different times and places" (as cited in Parfit, 2012, p. 6). This Lockean view on personal continuity emphasizes the importance of the continuity of a person's psychological properties, such as memories, attitudes, and ways of thinking. Based on this, Parfit (2012) has defended the so-called 'brain-based psychological criterion' for personal identity: "if some future person would be uniquely psychologically continuous with me as I am now, and this continuity would have its normal cause, enough of the same brain, this person would be me" (p. 6).

On the other hand, the physiological, or biological, continuity view means that it is not the psychological relation that creates your persistence over time, but rather some physical relation. On

this view, called 'animalism', you are the being in the past or future that has the same body, or is the same biological organism, as the present you. Therefore, you persist as long as your body is physiological continuous and have a physiological connection to your body in the past or future (Olson, 2015).

### **3. THE SELF IN NEUROSCIENTIFIC RESEARCH**

This chapter will go into the details of how the self is studied in the field of neuroscientific research. Since there are several different components possible to study, such as self-knowledge, personal traits or physical attributes, and personal memories, and due to the limited scope of this thesis, the focus will be restricted to only three aspects of the self and especially self-awareness – first-person perspective, self-reflection, and interoception – with a specific focus on the associated brain areas and their functions. Therefore, this chapter will review how specific neuroscientific studies, that are not explicitly discussing the question about the nature of self, are studying specific aspects of the self.

#### **3.1. The First-Person Perspective**

The term 'awareness' can be defined as knowing that one exists, and as Craig (2009) has put it, "an organism must be able to experience its own existence as a sentient being before it can experience the existence and salience of anything else in the environment" (p. 65). The term 'self-awareness' is further defined as the awareness of oneself as an existing individual being (J. Smith, 2017), and the mental capacity of being aware that oneself is being aware, leading to a kind of meta-awareness, as discussed above. According to Vogeley et al. (2004) conscious awareness of one's own mental states, such as perceptions, attitudes, and intentions, is included in the term 'self-consciousness'. To be able to put these mental states into our 'self-model' (i.e. the theoretical construct of essential features available through introspection), the ability of a phenomenal first-person perspective (1PP) is needed (Vogeley, Kurthen, Falkai, & Maier, 1999). This is when one's

experiences are assumed to center around oneself, as compared to the third-person perspective (3PP), which attributes mental and body states to someone other than oneself. The first-person perspective is necessary for human consciousness to appear as a subjective experience (Vogeley & Fink, 2003; Vogeley et al., 1999).

The first-person perspective is considered a basic constituent of Gallagher's (2000) notion of the 'minimal self', and within philosophy, it has been associated with subjectivity and conscious awareness (Gillihan & Farah, 2005), as when Nagel (1974) talked about someone's first-person perspective in the terms of subjective experiences. In cognitive neuroscience, by contrast, the first-person perspective is the ability to differentiate objective information to the self as opposed to others (Gillihan & Farah, 2005). Accordingly, Vogeley and Fink (2003) have defined 1PP, in terms of how one is experiencing it, as referring to the subjective experiential space surrounding one's own body. This experience of a body centered perspective is one of the features of the mentioned self-model (Vogeley et al., 1999) and 1PP is an important constituent in an individual's relations to the environment (Vogeley & Fink, 2003). Additionally, according to Blanke and Metzinger (2009), the first-person perspective is one of the central defining features for the minimal phenomenal selfhood, and thus an important part of consciousness in its simplest form.

Vogeley et al. (2004) studied the neural correlates of perspective taking with a simple visuospatial task, which was to be performed either from someone else's viewpoint (3PP) or from the viewpoint of oneself (1PP). They suggested the lateral superior temporal cortex, including the insula, both frontal and parietal mesial cortical areas (i.e. cortical areas alongside the sagittal plane of the brain, in the direction from dorsal to ventral), left frontal cortex, and the right postcentral gyrus, as the neural basis for 1PP. This since increased neural activation in these areas was found when the subjects were in the 1PP condition, as compared to 3PP (Vogeley et al., 2004). However, Vogeley and Fink (2003) concluded in a review that medial cortical regions (the anterior medial prefrontal, medial parietal, and posterior cingulate cortices) and inferior lateral parietal cortex might

be the basic neural mechanisms associated with the first-person perspective, based on evidence from both functional imaging, neuropsychological, and lesion studies.

### **3.2. Self-Reflection**

Another important component of self-awareness is the capacity to reflect upon your sense of self. According to Johnson et al. (2002), this sense of self is the unity of schemas regarding your abilities, attitudes, and traits that guide your social interactions with others, how you behave and what choices you make. In the paradigm mostly used in experimental and neuroscientific research to address the concept of self-reflection, the subjects are presented with trait adjectives and asked whether the trait applies to them or not (van der Meer, Costafreda, Aleman, & David, 2010). This self-reflective processing involves the conscious processing of a decision regarding one's own self. It is considered important to have an accurate representation of your abilities, traits, and attitudes when evaluating how you are behaving and to compare your own behavior with the behavior of others (van der Meer et al., 2010).

In the study by Johnson et al. (2002), activation in the medial prefrontal cortex (mPFC) and the posterior cingulate cortex (PCC) was found during self-reflection tasks. In another study, Ochsner et al. (2005) found neural activation in mPFC for self-knowledge, and self-knowledge based on the individual's beliefs about what others think about that individual required the additional activation of the insular, orbitofrontal, and temporal cortices. In line with this, Northoff and Bermpohl (2004) are arguing that a basic component in the generating of a self is the processing of self-referential stimuli in the cortical midline structures (CMS), consisting prominently of mPFC, PCC, and the anterior cingulate cortex (ACC). Self-referential processing is active when stimuli are experienced as related to the subject's own person, such as when looking at pictures of oneself as compared to someone else (Northoff et al., 2006). This means that the self-relevance of a stimulus depends on your individual experience of a stimulus, and is not intrinsic to the stimulus itself (Northoff & Bermpohl, 2004). Other results have shown decreased activity in the

mPFC when a subject is referring to someone else instead of oneself (Moran, Heatherton, & Kelley, 2009).

Accordingly, Modinos, Ormel, and Aleman (2009) found self-related neuronal activity in the anterior mPFC and ACC, and the anterior insula showed to be uniquely activated in association with self-reflection. In a meta-analysis, van der Meer et al. (2010) showed that activation in the insula, temporal pole, and the orbital part of the inferior frontal cortex was associated with self-processing. They concluded that it is possible to make a clear distinction in the activation of self-reflected versus other-reflective processing. However, this is directly opposed to the conclusion of Gillihan and Farah (2005) since they argued that the processing of self-related information would be associated with the same neural systems involved in person-related processing in general. This will be further examined in chapter four.

In a more recent study, Feng, Yan, Huang, Han, and Ma (2018) defined self-identity as the ability to distinguish the self from other people, which is in line with the neuroscientific definition of first-person perspective. They hypothesized that neural activations in the CMS, including the mPFC and PCC, would both differentiate representations of self from close (e.g. mother) and distant (e.g. celebrity) others, and additionally distinguish representations of different dimensions of knowledge (i.e. physical, mental, or social self-knowledge). Because of these hypotheses, the researchers measured self-reflection by having three different judgment-conditions (self, mother, and celebrity), where the subjects were asked to evaluate whether a given word was describing the person in each condition or not. Secondly, they differentiated between the three different dimensions of self-knowledge by giving words describing either the physical (physical attributes or characteristics), mental (personality traits) or social (consisting of social roles) dimension of person-knowledge (Feng et al., 2018).

By using functional magnetic resonance imaging (fMRI), Feng et al. (2018) analyzed how the neural activity patterns differed in the different judgment-conditions and for the different

dimensions of person-knowledge. As they hypothesized, the results showed different neural activations for self-processing in the PCC, precuneus, mPFC, and ACC, whereas the conditions of close or distant other, by contrast, showed more similarly distributed activations. This neural representation of self-identity was, however, only shown on the mental-dimension, and not in the physical nor social dimensions (Feng et al., 2018). In sum, these results showed similar activation for the mother- and celebrity condition, whereas there was a dissimilarity in neural activity in the self-condition compared to the other conditions, pointing to the association between self-identity and these activation patterns in the CMS. Conclusively, the findings by Feng et al. (2018) suggest the engagement of mPFC and PCC both in self-identity representations and different dimensions of self-knowledge.

In another study by Gusnard, Akbudak, Shulman, and Raichle (2001), the functional distinction between the dorsal and ventral parts of mPFC was further examined. By using fMRI, the researchers measured self-referential mental activity by asking the subjects to make two kinds of judgments, one was self-referential and the other was not. In the self-referential condition, the subjects were shown emotional pictures and asked to evaluate how the picture made them feel and note if the evoked feeling was pleasant, unpleasant, or neutral – thus reflecting upon one's own feeling. In the non-self-referential condition, on the other hand, the subjects were shown neutral pictures and asked to judge if the picture was of a scene indoor or outdoor (Gusnard et al., 2001). The results showed increased activity in the dorsal part of mPFC associated with the emotional self-referential judgment, whereas the ventral mPFC showed decreased neural activity in both tasks, and Gusnard et al. (2001) concluded that the self-referential mental activity is associated with the increased neuronal activation in the dorsal part of mPFC.

In contrast to the previous findings, Philippi et al. (2012) examined in a case study how self-awareness was affected in a patient with extensive bilateral brain damage including the insula, ACC, and mPFC. These areas were studied particularly since they are hypothesized to be critical for

the multifaceted phenomenon of self-awareness (as seen above; Johnson et al., 2002; Modinos et al., 2009; Northoff & Bermpohl, 2004; Ochsner et al., 2005). Philippi et al. (2012) used different basic tests of self-recognition, self-agency, and self-concept, in addition to a self-awareness interview to measure the patient's metacognitive, reflective, and introspective abilities. The findings from this study showed that the patient's self-awareness was largely preserved and that he was conscious and self-aware despite this widespread cortical damage (Philippi et al., 2012). Based on this, the authors concluded that their study did not favor the hypothesis that the insula, ACC, and mPFC would play a critical and necessary role in the producing of self-awareness, but rather that it would be more likely that self-awareness is created by distributed interactions between different brain networks including the brainstem, thalamus, and posteromedial cortices (Philippi et al., 2012).

### **3.3. Interoception**

The third component of self-awareness covered in this thesis is the one with the physical body in focus, called bodily awareness. This perceptual experience is seen to reveal the self as an object since you become aware of your 'body-self', and this through awareness from the inside (J. Smith, 2017). This is called 'interoception' and refers to the sense of the physiological state of your body (including temperature, pain, and muscular and visceral signals) that further provides the basis for your ability to subjectively reflect upon the physical condition of your body, and thus giving an answer to 'how you feel' (Craig, 2002). The exteroceptive senses, on the other hand, include vision, smell, hearing, taste, and mechanical contact (Damasio, 2003). Interoception is, according to Damasio (2003), the most critical process for the self.

The brain area suggested to be the 'interoceptive center' is the insula, and in particular the right anterior part (Craig, 2002). Studies have found enhanced neural activity in bilateral somatomotor cortex, anterior insula, inferior frontal cortex, anterior cingulate cortex, and supplementary motor cortex (Critchley, Wiens, Rotshtein, Öhman, & Dolan, 2004; Stern et al., 2017) when subjects were asked to pay attention to interoceptive signals (e.g. their own heartbeats),

as compared to exteroceptive attention (i.e. attention to external stimuli). Critchley et al. (2004) have further proposed that the information processing in the anterior insula is a constituent for the subjective feeling of body states since it provides a map of the internal body states and contributes to the representation of the self.

In his review, Craig (2009) included studies from a wide range of fields, such as attention, time perception, and emotional awareness, in which all of them reported that activation in the anterior insular cortex (AIC) was associated with subjective feelings. In addition, these studies also showed associations between the activity in AIC and other processes, such as cognitive choices, music perceptions, and awareness of sensations and movements, with the only common feature that every different task engaged the subject's awareness. Since the AIC was the only region involved in all of these different tasks, Craig (2009) has argued for the hypothesis that the AIC is the brain structure generating human awareness. He means that the reviewed studies give compelling evidence for the AIC to be the neural substrate for the human capacity to be aware of yourself, others, and the surrounding environment and that the AIC, therefore, would be interpreted as the neural correlate of awareness (Craig, 2009).

It is further argued that the pathway called the lamina-I spinothalamocortical pathway plays an important role in the processing of interoceptive information since this is where the very thin and unmyelinated peripheral nerves are conveyed into the central nervous system (Damasio, 2003). Through the region called lamina-I located in the dorsal horn of the spinal cord, the neural information is sent to the brainstem and thalamus, to end up with providing the primary interoceptive representation of the physiological condition of the body in the insular cortex (Craig, 2002). Damasio (2003) has argued that this system, with sensors located throughout the structure of the entire body, provides the grounding for the perception of our own being.

The information conveyed through the lamina-I spinothalamocortical pathway provides the basis for the physiological image of the entire body in the primary interoceptive region of the

posterior insula, including numerous distinct feelings, and these neural constructs are then re-represented in the middle part of insula and once again in the AIC (Craig, 2009). Consistent with this, Farb, Segal, and Anderson (2013) have found that interoceptive attention (in this case, attention towards respiration) modulated the posterior and middle parts of insula (i.e. primary and secondary interoceptive regions, respectively). They also demonstrated a new dissociable recruitment of the posterior and anterior insular cortices, namely that the posterior insula showed engagement in interoceptive responses, which then shifts in a graded manner toward more exteroceptive responses in the anterior insula (involving an additional awareness of the environmental context beyond the body; Farb et al., 2013). According to Craig (2009), "the mid-insula integrates these homeostatic re-representations with activity that is associated with emotionally salient environmental stimuli of many sensory modalities" (p. 67), leading to a multisensory integration of information from both interoceptive and exteroceptive inputs. Therefore, he proposes that this integration results in "a unified final meta-representation of the 'global emotional moment' near the junction of the anterior insula and the frontal operculum" (Craig, 2009, p. 67), and that this then would generate the image of the 'material me'.

#### **4. IS THE BUNDLE THEORY OF SELF SUPPORTED BY NEUROSCIENCE?**

In contrast to the previous chapter, which had a focus on different components of self-awareness and their associated brain regions, this fourth chapter will review how different studies have focused on the involvement of different brain networks in processes related to the self. The question of particular interest is to what extent neuroscience supports the bundle theory and the view of no unified self-system in the brain.

##### **4.1. No Self-Specific System**

According to Gillihan and Farah (2005), the concept of 'self' includes two kinds of aspects; one physical and one psychological. The physical aspect of the self is focusing on either specific

body parts (e.g. the face or an arm) or the body as a whole and how its different parts are related in space. Studies on the psychological aspect of self, on the other hand, reflects one's knowledge of the self, including memory knowledge (both autobiographical and semantic memories) as well as the aforementioned first-person perspective. In their review, Gillihan and Farah (2005) examined two central questions about the self-related processing in the brain; the first one whether self-referential processing is special in the same way as human language has been said to be special, whereas the second considered whether the many different aspects of self-related processing are the functions of a unitary self-system. The first question will be examined next, and the second in the following section of this chapter.

For the first question, Gillihan and Farah (2005) used four criteria for the self-system to be special: anatomical specificity, which is when a system is engaging or requiring distinct brain areas; functional uniqueness, referring to how information is processed within a system of structures instead of where it is processed (e.g. face recognition is more holistic than other object recognition processing); functional independence, which is when one system's activity is independent of the activation of another system; and species specificity, which is as the name indicates, when the system (e.g. for language) is unique for one species only (e.g. humans). Based on this, Gillihan and Farah (2005) concluded that for some aspects of the physical self (e.g. in patients with asomatognosia; the lack of recognition of one's own body part) the answer would be *yes*, the self-referential processing is special, but for most aspects of both the physical and psychological aspects of self there is no good evidence for it to be special.

In a review, Musholt (2013) focused on how the cortical midline structures (CMS), including the ventral and dorsal mPFC, parietal/posterior cingulate cortex, and anterior cingulate cortex, are related to the self. The doubts whether it would be possible to locate the self somewhere in the brain are, according to Musholt (2013), based on findings showing that the CMS are also involved in cognitive tasks not related to the self (e.g. some general process of evaluating information about

others than oneself; Legrand & Ruby, 2009), and evidence showing that other areas outside the CMS, such as the anterior insula, is activated during tasks of self-reflection (Modinos et al., 2009).

It is further suggested that other processes than the self-related processing per se, such as familiarity (Gillihan & Farah, 2005) or general evaluation (Legrand & Ruby, 2009), would be confounding factors underlying the neuronal changes in the CMS (Qin & Northoff, 2011). In their meta-analysis of brain imaging studies, Qin and Northoff (2011) partly addressed these concerns as their general aim was to "investigate the relationship between brain activity related to the processing of self-specific, personally familiar, and other (non-self and non-familiar) stimuli" (Qin & Northoff, 2011, p. 1223). Their first specific aim was to look for the possible differences and overlaps between these three conditions while controlling for unspecific task- and stimulus-related effects (e.g. the general processing of evaluation). They hypothesized that activity in the anterior cortical midline structures (e.g. the perigenual anterior cingulate cortex; PACC) would be associated with self-specific information processing when compared to brain regions involved in the perception of personally familiar and other stimuli (Qin & Northoff, 2011). Their second aim was to detect the relationship between self-specific processing and resting-state activity in the default mode network. However, since this aim is out of scope for the present thesis the focus will be on their first aim only.

Worth mentioning is that Qin and Northoff (2011) defined the self in a broader sense than the one mentioned above by Gillihan and Farah (2005), by including the relationship between the person and the specific environmental stimuli in addition to the physiological and psychological dimensions. They operationalized self-specificity as "a specific relation between the organism and stimuli – with the latter including physical-bodily, psychological-cognitive/mental and exteroceptive-sensory stimuli" (Qin & Northoff, 2011, p. 1223). Their findings showed that the PACC was required during the self-condition and showed more activation associated to self-specific stimuli as compared to non-self (i.e. familiar and other), and are thus suggested to confirm the

hypothesis that the anterior midline region PACC is important for the self. Furthermore, in mPFC there was an overlapping activation related to the self and familiarity conditions, and PCC showed overlap between all three conditions (i.e. self-specific, personally familiar, and other, non-self and non-familiar, stimuli; Qin & Northoff, 2011).

The findings by Qin and Northoff (2011) are in line with the assumptions of Gillihan and Farah (2005) that there is a regional overlap between self-specific stimuli and familiarity. However, Qin and Northoff's (2011) study did not show a complete overlap since the PACC is suggested to be recruited only during the self-condition. Since these results indicate an overlapping activity in mPFC and PCC for self and familiarity but differences in the activity of PACC, the authors' conclusive suggestion is that self and familiarity may not be considered as identical processes.

Additionally, Qin and Northoff (2011) investigated how the self is related to more general task-related evaluation processing because of the argument that regions involved in the processing of self are not specific for the self, but rather involved in the general evaluation process and that their activity, therefore, would be task-specific (Legrand & Ruby, 2009). However, when Qin and Northoff (2011) tested if the regions involved in the evaluation of stimuli were related to the self or non-self, the results showed that activity in PACC was associated with self-specific stimuli and not the task-related process of evaluation.

Whereas Musholt (2013) and Qin and Northoff (2011) focused on the cortical midline structures, Legrand and Ruby (2009) had their focus on another set of brain structures called the E-network. In a meta-analysis they examined several studies all considered to tackle the investigation of self by approaching the self in different ways, such as recognition of one's own face or first name, associating an action to oneself, recalling personal information, or assessing one's own personality traits or feelings (Legrand & Ruby, 2009). Indeed, these various studies involved a wide variety of cognitive tasks and stimuli, contributing to the finding of a long list of brain areas involved in these self studies: mPFC, precuneus/posterior cingulate gyrus, temporal pole,

temporoparietal junction, insula, postcentral gyrus, superior parietal cortex, precentral gyrus, lateral prefrontal cortex, hippocampus, parahippocampal gyrus, fusiform gyrus, and the occipital cortex (Legrand & Ruby, 2009).

However, the four areas most frequently activated in self versus non-self aspects were the mPFC, precuneus/posterior cingulate gyrus, temporoparietal junction, and the temporal pole. This set of regions is what Legrand and Ruby (2009) call the aforementioned E-network (*E* from *evaluation*, which in this case involves inferential processes and memory recall). Importantly, these brain regions are also recruited for the representation of others and are not exclusively devoted to the self. This means that the E-network does not, according to Legrand and Ruby (2009), exhibit that kind of functional specificity that many studies on the self hypothesize, and that it is not a cerebral system more activated for the self than for non-self. The E-network regions are sometimes more activated for the self than for others, but at other times they are more activated for others than for the self (Legrand & Ruby, 2009). The authors suggested that activation in the mPFC would be associated with inferential processes, whereas activation of the precuneus/posterior cingulate cortex, the temporoparietal junction, and the temporal pole would relate to the processing of memory recall (Legrand & Ruby, 2009).

In sum, their review of neuroimaging data from self-studies made Legrand and Ruby (2009) suggest that the E-network, which is activated in both self-related and non-self-related tasks, is involved in nonspecific cognitive processing required for evaluative processes in general, such as creative reasoning or comparison. This means that the network is not specified for the self, but is rather activated as soon as some evaluation is needed for any kind of stimulus since "evaluation is a cognitive process that is involved irrespective of the subject targeted in the task" (Legrand & Ruby, 2009, p. 266) and is "neither domain specific nor subject specific" (p. 266).

#### 4.2. No Unified Self-System

The second question examined by Gillihan and Farah (2005) concerned whether the different components of self-processing might be the functions of a unified self-system in the brain. The authors argued that for the hypothesis of a unitary self-system in the brain to be supported, the total findings from the examined studies would have shown clustering in certain brain regions or networks (Gillihan & Farah, 2005). However, despite one's subjective experience of a unified self, Gillihan and Farah (2005) found no support that the self-related research would support a unitary, common system since "neither the imaging nor the patient data implicate common brain areas across different aspects of the self. This is not surprising because there is generally little clustering even within specific aspects of the self" (p. 94).

In line with this skeptic view of the possibility to find a unified self somewhere in the brain, American philosopher Patricia Churchland (2002) has argued that the wide range of usage of the self-concept, both referring to our physical bodies and our social or private selves (i.e. the psychological self), motivates that the problems of the self would rather be problems in terms of self-representational capacities. In her article, Churchland (2002) concluded that "the self thus turns out to be identifiable not with a nonphysical soul, but rather with a set of representational capacities of the physical brain" (p. 308). These self-representational capacities include for example the representing of the internal environment of the body and the viscera (mainly involving the brainstem and hypothalamus) and the representing of one's autobiographical memories (associated with the medial temporal lobe structures), meaning that there are many different brain areas seemed to be involved in self-representation (Churchland, 2002). So, instead of seeing the self as a singular entity, Churchland (2002) emphasized the search for a plurality of different brain functions and processes involved in the representation of self. This is in line with the bundle theory-view of the nature of self, viewing the self as a bundle of perceptions and sensations instead of a singular ego-entity (Parfit, 1987).

One even more skeptic interpretation of the self and also in line with the thoughts of Hume (1739/1896), is the no-self position, according to which perceptions of the external environment appear without a self, without a perceiver having the experiences. This view is defended by empirical research suggesting segregated brain activations for rest, introspection, and self-related evaluation on the one hand, and tasks oriented by the external environment on the other (Legrand & Ruby, 2009).

One example of this research is the fMRI study by Goldberg, Harel, and Malach (2006), in which the authors examined brain activations associated with demanding sensory categorization and the activity patterns engaged in self-related introspection. Their aim was to directly compare self-related processes with processes engaged in sensory information processing only (Goldberg et al., 2006). In an introspection task subjects were asked to view images and, related to the self, evaluate their emotional response to the stimulus, and after this categorize the emotional response as positively or negatively high, or neutral. By contrast, in a categorization task, the subjects were asked to just categorize the pictures into one of two categories. This made the two conditions identical in terms of sensory stimuli and motor output (pressing a button to indicate their responses), with differences in the cognitive task only; sensory processing was present in both conditions, but the additional self-related introspection only in the introspection task (Goldberg et al., 2006).

The results from this study showed clear segregation between brain regions involved in self-related introspection (prominent activation in the prefrontal cortex and especially within the superior frontal gyrus; SFG) and the cortical regions engaged in sensorimotor processing (i.e. the sensorimotor cortices; Goldberg et al., 2006). These results were interpreted as arguing against the notion that self-related representations would be a necessary component of the appearance of subjective awareness, since the the subjects' sense of self-awareness was eliminated in addition to decreased self-related brain activity in SFC, as compared to the introspection condition. Thus, this

is considered to be evidence against the view that there would be a self perceiving the experiences since the self-related processes are not activated when an individual is exposed to external stimuli (Goldberg et al., 2006).

### 4.3. The Self as an Illusion

Legrand and Ruby (2009) further emphasized the importance that their conclusions only support the notion that evaluation in the process of self-evaluation is not *self-specific*, and could not be the grounds for eliminating the self altogether. In other words, their review does not support the skeptical interpretation that if there are no self-specific neural correlates for the self-evaluative process then the self might not even be in the brain and that this, on reductionistic grounds, might be used to eliminate the notion of self (Legrand & Ruby, 2009).

One form of such a no-self alternative, however, has been proposed by the German philosopher Thomas Metzinger. Since there is a concept of 'the self' in folk-phenomenology, many individuals automatically assume that there is an entity like the self actually existing in the world, and this even though "there seems to be no empirical evidence and no truly convincing conceptual argument that supports the actual existence of 'a' self" (Metzinger, 2011, p. 279). In contrast, Metzinger (2009) argues that conscious experience is like a tunnel with an ego having the first-person perspective (creating the metaphor of an 'ego tunnel'). What you experience is only a fraction of everything actually existing in the external environment around you, and this is why it is a tunnel; "the ongoing process of conscious experience is not so much an image of reality as a tunnel *through* reality" (Metzinger, 2009, p. 6).

Metzinger (2009) suggested that there is a conscious self-model activated by your brain providing you with an inner image of the organism as a whole. What you then subjectively experience makes up your phenomenal self-model (PSM), where the term 'phenomenal' is used as the subjective way the world appears to you. The phenomenal ego is here defined as the contents of the PSM at every given moment, including your bodily sensations, emotional states, perceptions,

memories, and thoughts. The idea is that "the content of consciousness is the content of a simulated world in our brains, and the sense of *being there* is itself a simulation" (Metzinger, 2009, p. 23).

Further, Metzinger (2009) has claimed that when placing this self-model within the model of the world a center is created, and this center is what you consciously experience as a first-person perspective and hence the *ego* in the ego tunnel. However, the self-model has the characteristic of being *transparent*, which means that the system using it (i.e. you, the individual) cannot recognize the model *as* a model or the representation *as* a representation; "you are in contact only with its content; you never see the representation as such; you have the illusion of being directly in contact with the world" (Metzinger, 2009, p. 42). It is because of this transparency you cannot be aware of your self as such, but only the contents of your mind (i.e. your experiences). Further, the walls of the ego tunnel are impenetrable for you because of the fact that you can accept that something is just an internal construct and not a direct representation of the reality only at the cognitive level, but not attentionally or introspectively since you have nothing of reference that is outside this tunnel; you cannot refer to something that is not on the level of subjective experience. In sum, Metzinger (2009, 2011) argues that the self is not a thing, but rather an ongoing process simulated in the brain.

In line with this standpoint, the British philosopher Julian Baggini (2011) argued in a TED talk that the self is not a thing but rather a collection of everything you know, remember, desire, believe, and experience, and the 'you' is the sum of all those parts. The self is a process that is changing and not a permanent essence. However, Baggini (2011) argued against the view that the self would be just an illusion, by saying that just because there is no thing called the self, and that you are rather a complex collection of different physiological and psychological experiences, this does not mean that the self is not real. He means that the self is not created until all the needed parts are put together. All the experiences you have are related to the same body and brain and it is this relation between all those experiences that is *you* (Baggini, 2011). This is a kind of physiological continuity explaining personal identity, as discussed above.

#### 4.4. Defending the Unified Self

In a review article, R. Smith (2017) has defined the self-component that Metzinger (2009) is considering as an illusion as the 'experiencer-agent self'. This concept, called the EA self, is referring to a unified entity or subject and is the abstract representation which is central in the way you often think about yourself (R. Smith, 2017). This EA self is the part of the self which has access to all and only the information you are conscious of and uses this information to decide how you will act in the world. In opposite to Metzinger, R. Smith (2017) argued that the EA self-concept refers to a system or mechanism with access to conscious information and which is making the decisions, and because of this, one need not conclude that the self is an illusion.

In the article, R. Smith (2017) distinguished the mental representations of the self (i.e. self-related information) from the referents of those representations (i.e. the things those representations are referring to). The skeptic side of the existence of the EA self (e.g. Metzinger, 2009) argues that you may have a mental representation of this self-concept, but that this representation fails to refer. In other words, there is no actual thing or brain process that matches this representation of the self you have in your mind. In contrast, R. Smith (2017) defended the view that "the mentally represented EA self-concept successfully refers to a particular system in the brain" (p. 22), and emphasized the notion that in the case of the EA self, "both the representation and its referent will be found within the individual's brain – albeit in different brains systems" (p. 23), as opposed to mental representations of, for example, physical objects which are based on visual stimuli and the visual brain system.

The neuroscientific ground for this EA self is, according to R. Smith (2017), that the frontal, parietal, cingulate, and insula regions would work together as parts highly connected to each other over the cortex, and that these would form a 'core-circuit' of distinct, but still connected, processors. These regions are together suggested to be understood as the unified system underlying

conscious perceptions and decision making, and this processing system would, therefore, be the neural basis for the experiencer-agent self (R. Smith, 2017).

At first, the view of R. Smith (2017) might be considered an ego theory-view since he argued for a unified self referring to a particular brain system. However, it would instead be connectable to the bundle theory view of the self, and this because he emphasized that there are a lot of representations of the self (such as body position and heart rate, memories of past experiences and current emotional states) and that the EA self is just one of these abstract representations. The described unified system is for the EA self only and this representation of the self is one of many contributing to the self-experience as a whole. In other words, R. Smith (2017) is just defending the view that the representation of the EA self is referring to a system actually existing in the brain, and that this is not an illusion.

Before going on to the fifth and final chapter, the answer to the heading-question of this section is in sum that the reviewed literature is suggested to support the bundle theory of the nature of self. There are still researchers arguing for the neural basis of a unified self (e.g. R. Smith, 2017), however, this unified self (the experiencer-agent self in this case) is just one component of the whole self. Therefore, this might still be interpreted as supporting the bundle theory and not the view of an ego-self being the entity unifying all the self-components and experiencing all your perceptions. This will be discussed in more details in the next chapter.

## 5. DISCUSSION

The final chapter of this thesis will first provide a short summary of the content up to this point and a general discussion of how the covered self-components and their associated brain areas are related to the philosophical debate regarding the nature of self and the bundle theory-view, followed by a section considering some of the problems and limitations. Lastly, the conclusion part will end the thesis with an explicitly formulated answer to the initial research question.

In the third chapter of the present thesis, the first-person perspective, self-reflection, and interoception, with their associated brain areas were in focus since these are aspects of self-awareness. The suggested neural basis for 1PP consists of the medial cortical regions and inferior lateral parietal cortex (Vogeley & Fink, 2003). Processes of self-reflection, on the other hand, have been associated with brain activity in the mPFC and PCC (Johnson et al., 2002; Feng et al., 2018), insular, orbitofrontal, and temporal cortices (Ochsner et al., 2005), the cortical midline structures, in particular the mPFC, PCC, and ACC (Northoff & Bermpohl, 2004), and anterior insula (Modinos et al., 2009). Thirdly, the brain region suggested to be the interoceptive center is the insular cortex, and especially the anterior part (Craig, 2002; Critchley et al., 2004; Stern et al., 2017).

The fourth chapter examined some studies on how different brain networks (the CMS; Musholt, 2013; Northoff & Bermpohl, 2004; Northoff et al., 2006; and the E-network; Legrand & Ruby, 2009) are related to the processing of self-related information. It has been argued that the different brain areas involved in the many aspects of self-processing are neither special for the self nor the function of a unitary self-system (Gillihan & Farah, 2005). Further, findings showing that self-related brain processes are not activated when an individual is exposed to external stimuli has been suggested as evidence that there is no separate perceiver located inside the brain (Goldberg et al., 2006). Some researchers are making more drastic interpretations of this, arguing that the concept of self ought to be eliminated altogether and that the self is a process simulated in the brain rather than a thing in itself (Metzinger, 2009). Nonetheless, R. Smith (2017) has argued that the experiencer-agent self is referring to a system actually existing in the brain.

The main concern of this thesis was to review how findings from neuroscientific research are contributing to the debate about the nature of self. It is almost commonsensical to think that the self is unified and that there is someone experiencing your perceptions and that this someone is existing by itself. This is the most intuitive view of the self simply because that is how it feels. However, it is seen in the literature reviewed in this thesis that in neuroscientific research it is typically not

believed that there is some specific self-structure in the brain and that the self should not be treated like a homunculus (e.g. Damasio, 2003; Churchland, 2002). Accordingly, when neuroscientific studies are trying to measure the self and look for the neural basis they rather tend to divide the self into different components, all of them together contributing to the unified self-experience as a whole.

As chapter three showed, there are dissimilar neural activations of different brain regions when focusing on different components of the self. Indeed, for some aspects, there are similar areas, and there are specific areas involved in more than one of these self-related processes (such as the mPFC in both IPP and self-reflection, and the insula in both self-reflection and interoception). However, there are still separate activations and the similarity is not complete, leading to different activity patterns for the different processes. Based on the literature reviewed in this thesis, it is suggested that the most important brain areas for the ability of self-awareness are the mPFC, ACC and insular cortex. By contrast, however, the case study by Philippi et al. (2012) showed that self-awareness can still be preserved despite bilateral brain damage to all these three regions. Therefore, mPFC, ACC, and the insula might be suggested as areas playing an important, but not critical, role in the producing of self-awareness and self-consciousness. Additionally, it has also been showed that these activations are not special for only the self, but they are also associated with the processing of information related to others than the self.

Taken together, this suggests that neuroscience research would point to a bundle theory-view of what the self is, rather than there would be an ego within the brain. For example, as noted above, Damasio (2003) argued that there would rather be the stable representation of the individual's own continuity that gives the individual a mental reference to refer to, and thus contributing to the conscious sense of self. However, that the self is not a thing in itself should not be interpreted as the self does not even exist; it is still something *real* (Baggini, 2011).

### 5.1. Problems and Limitations

As noted in the introduction and throughout the thesis, the self is a complex and multifaceted phenomenon and there is no consensus about what the self really is across different disciplines, which leads to difficulties in observing and measuring it. Because of this, there are several more components of the self usually measured in neuroscience and other sciences that are not covered in the present thesis, in addition to other philosophical standpoints and discussions that may be considered as important to include in this discussion. These have simply been filtered out because of the scope restrictions and limited amount of time. For example, the relation between the self-related brain processes and the functions of the default mode network is not covered, and this despite the similarities of brain regions and neural activations shown to be involved in both self-processing and resting-state activity (see e.g. Gusnard et al., 2001; Qin & Northoff, 2011). Because of this, the present thesis has low external validity which is a limitation and threat to the generalizability. Therefore, it is important to have in mind that the conclusions made in this thesis are only based on the literature covered, and can because of this not be generalized to other discussions outside this specific context.

Another problem with studying the self is the validity of measurement and making sure that the operationalization is measuring what it is supposed to measure; is a given study really measuring what it says it is measuring? Regarding the psychological aspects of the self, Gillihan and Farah (2005) say that the evidence is difficult to interpret because there might be confounding variables present when trying to distinguish the processing of self from the processing of others, such as general evaluation discussed above. If these confounders are not taken into account, it is easy to conclude prematurely that the processing of self-related information takes place in a functionally unique structure. There are also findings from different studies showing different neural activation patterns even though the studies are measuring the same aspect of the self (e.g. the first-person perspective). Even though cognitive neuroscience has these methodological issues to

overcome when studying the self, Gillihan and Farah (2005) means that the self still is a central topic in this field of research since most of the self-aspects after all are represented in your brain and play an important role in how you process information.

## **5.2. Conclusion**

To conclude, the research question stated in the beginning of this thesis was whether you can say, with support from neuroscientific findings, that the bundle theory and the view that the self is not a separate thing in itself is more correct than the ego theory-view and the belief that the self has to be a continuous subject of experiences, despite the counterintuitive characteristic of the former view. The ego theory-view is somehow easier to assume because of the present feeling of a continuous and unified self and that this self would be located in the brain. By contrast, based on the philosophical and neuroscientific literature covered and reviewed in the present thesis, the answer to the research question would be: *yes*, the bundle theory-view of the nature of self is suggested to be more correct than the intuitive feeling that you all are unified selves persisting over time. This means that it seems like you are not a unified self located in your brain.

word count: 12 250

## 6. REFERENCES

- Baggini, J. (2011, November). *Julian Baggini: Is there a real you?* [Video file]. Retrieved from [https://www.ted.com/talks/julian\\_baggini\\_is\\_there\\_a\\_real\\_you](https://www.ted.com/talks/julian_baggini_is_there_a_real_you)
- Blackmore, S. (2010). *Consciousness: An introduction*. New York, NY: Routledge.
- Blanke, O., & Metzinger, T. (2009). Full-body illusions and minimal phenomenal selfhood. *Trends in Cognitive Sciences*, *13*(1), 7-13. doi:10.1016/j.tics.2008.10.003
- Christoff, K., Cosmelli, D., Legrand, D., & Thompson, E. (2011). Specifying the self for cognitive neuroscience. *Trends in Cognitive Sciences*, *15*(3), 104-112. doi:10.1016/j.tics.2011.01.001
- Churchland, P. (2002). Self-representation in nervous systems. *Science*, *296*(5566), 308-310. doi:10.1126/science.1070564
- Craig, A. D. (2002). How do you feel? Interoception: The sense of the physiological condition of the body. *Nature Reviews Neuroscience*, *3*(8), 655-666. doi:10.1038/nrn894
- Craig, A. D. (2009). How do you feel — now? The anterior insula and human awareness. *Nature Reviews Neuroscience*, *10*(1), 59-70. doi:10.1038/nrn2555
- Critchley, H. D., Wiens, S., Rotshtein, P., Öhman, A., & Dolan, R. J. (2004). Neural systems supporting interoceptive awareness. *Nature Neuroscience*, *7*(2), 189-195. doi:10.1038/nn1176
- Damasio, A. (1999). *The feeling of what happens: Body and emotion in the making of consciousness*. New York, NY: Harcourt, Inc.
- Damasio, A. (2003). Feelings of emotion and the self. *Annals of the New York Academy of Sciences*, *1001*(1), 253-261. doi:10.1196/annals.1279.014
- Farb, N. A. S., Segal, Z. V., & Anderson, A. K. (2013). Attentional modulation of primary interoceptive and exteroceptive cortices. *Cerebral Cortex*, *23*(1), 114-126. doi:10.1093/cercor/bhr385

- Feng, C., Yan, X., Huang, W., Han, S., & Ma, Y. (2018). Neural representations of the multidimensional self in the cortical midline structures. *NeuroImage*, *183*, 291-299. doi:10.1016/j.neuroimage.2018.08.018
- Gallagher, S. (2000). Philosophical conceptions of the self: Implications for cognitive science. *Trends in Cognitive Sciences*, *4*(1), 14-21. doi:10.1016/S1364-6613(99)01417-5
- Gillihan, S. J., & Farah, M. J. (2005). Is self special? A critical review of evidence from experimental psychology and cognitive neuroscience. *Psychological Bulletin*, *131*(1), 76-97. doi:10.1037/0033-2909.131.1.76
- Goldberg, I. I., Harel, M., & Malach, R. (2006). When the brain loses its self: Prefrontal inactivation during sensorimotor processing. *Neuron*, *50*(2), 329-339. doi:10.1016/j.neuron.2006.03.015
- Gusnard, D. A., Akbudak, E., Shulman, G. L., & Raichle, M. E. (2001). Medial prefrontal cortex and self-referential mental activity: Relation to a default mode of brain function. *Proceedings of the National Academy of Sciences*, *98*(7), 4259-4264. doi:10.1073/pnas.071043098
- Hume, D. (1896). Book I, part IV, section VI: Of personal identity. In L. A. Selby-Bigge (Ed.), *A treatise of human nature*. Retrieved from <https://oll.libertyfund.org/titles/hume-a-treatise-of-human-nature> (Original work published 1739)
- James, W. (1950). Chapter X: The Consciousness of Self. In *The principles of psychology, Vol. 1*. New York, NY: Dover Publications, Inc. Retrieved from <https://books.google.se/books> (Original work published 1890)
- Johnson, S. C., Baxter, L. C., Wilder, L. S., Pipe, J. G., Heiserman, J. E., & Prigatano, G. P. (2002). Neural correlates of self-reflection. *Brain*, *125*(8), 1808-1814. doi:10.1093/brain/awf181
- Legrand, D., & Ruby, P. (2009). What is self-specific? Theoretical investigation and critical review of neuroimaging results. *Psychological Review*, *116*(1), 252-282. doi:10.1037/a0014172

- Metzinger, T. (2009). *The ego tunnel: The science of the mind and the myth of the self*. New York, NY: Basic Books.
- Metzinger, T. (2011). The no-self alternative. In S. Gallagher (Ed.), *The Oxford Handbook of the Self* (pp. 279-296). doi:10.1093/oxfordhb/9780199548019.003.0012
- Modinos, G., Ormel, J., & Aleman, A. (2009). Activation of anterior insula during self-reflection. *PLoS ONE*, 4(2), e4618. doi:10.1371/journal.pone.0004618
- Moran, J. M., Heatherton, T. F., & Kelley, W. M. (2009). Modulation of cortical midline structures by implicit and explicit self-relevance evaluation. *Social Neuroscience*, 4(3), 197-211. doi:10.1080/17470910802250519
- Musholt, K. (2013). A philosophical perspective on the relation between cortical midline structures and the self. *Frontiers in Human Neuroscience*, 7, 1-11. doi:10.3389/fnhum.2013.00536
- Nagel, T. (1974). What is it like to be a bat? *The Philosophical Review*, 83(4), 435-450. doi:10.2307/2183914
- Northoff, G., & Bermpohl, F. (2004). Cortical midline structures and the self. *Trends in Cognitive Sciences*, 8(3), 102-107. doi:10.1016/j.tics.2004.01.004
- Northoff, G., Heinzl, A., de Greck, M., Bermpohl, F., Dobrowolny, H., & Panksepp, J. (2006). Self-referential processing in our brain—A meta-analysis of imaging studies on the self. *NeuroImage*, 31(1), 440-457. doi:10.1016/j.neuroimage.2005.12.002
- Ochsner, K. N., Beer, J. S., Robertson, E. R., Cooper, J. C., Gabrieli, J. D. E., Kihlstrom, J. F., & D'Esposito, M. (2005). The neural correlates of direct and reflected self-knowledge. *NeuroImage*, 28, 797-814. doi:10.1016/j.neuroimage.2005.06.069
- Olson, E. T. (2015). Personal Identity. In E. N. Zalta (Ed.), *The Stanford Encyclopedia of Philosophy* (Summer 2017 ed.). Retrieved from <https://plato.stanford.edu/archives/sum2017/entries/identity-personal/>

- Parfit, D. (1971). On "the importance of self-identity". *The Journal of Philosophy*, 68(20), 683-690. doi:10.2307/2024939
- Parfit, D. (1987). Divided minds and the nature of persons. In S. Schneider (Ed.), *Science fiction and philosophy: From time travel to superintelligence* (2nd ed., 2016, pp. 91-98). doi: 10.1002/9781118922590.ch8
- Parfit, D. (2012). We are not human beings. *Philosophy*, 87(1), 5-28. doi:10.1017/S0031819111000520
- Philippi, C. L., Feinstein, J. S., Khalsa, S. S., Damasio, A., Tranel, D., Landini, G., ... Rudrauf, D. (2012). Preserved self-awareness following extensive bilateral brain damage to the insula, anterior cingulate, and medial prefrontal cortices. *PLoS ONE*, 7(8), e38413. doi:10.1371/journal.pone.0038413
- Qin, P., & Northoff, G. (2011). How is our self related to midline regions and the default-mode network? *NeuroImage*, 57(3), 1221-1233. doi:10.1016/j.neuroimage.2011.05.028
- Smith, J. (2017). Self-Consciousness. In E. N. Zalta (Ed.), *The Stanford Encyclopedia of Philosophy* (Fall 2017 ed.). Retrieved from <https://plato.stanford.edu/archives/fall2017/entries/self-consciousness/>
- Smith, R. (2017). A neuro-cognitive defense of the unified self. *Consciousness and Cognition*, 48, 21-39. doi:10.1016/j.concog.2016.10.007
- Stern, E. R., Grimaldi, S. J., Muratore, A., Murrough, J., Leibu, E., Fleysher, L., ... Burdick, K. E. (2017). Neural correlates of interoception: Effects of interoceptive focus and relationship to dimensional measures of body awareness. *Human Brain Mapping*, 38(12), 6068-6082. doi: 10.1002/hbm.23811
- van der Meer, L., Costafreda, S., Aleman, A., & David, A. S. (2010). Self-reflection and the brain: A theoretical review and meta-analysis of neuroimaging studies with implications for schizophrenia. *Neuroscience and Biobehavioral Reviews*, 34(6), 935-946. doi:10.1016/j.neubiorev.2009.12.004

Vogeley, K., & Fink, G. R. (2003). Neural correlates of the first-person-perspective. *Trends in Cognitive Sciences*, 7(1), 38-42. doi:10.1016/S1364-6613(02)00003-7

Vogeley, K., Kurthen, M., Falkai, P., & Maier, W. (1999). Essential functions of the human self model are implanted in the prefrontal cortex. *Consciousness and Cognition*, 8(3), 343-363. doi:10.1006/ccog.1999.0394

Vogeley, K., May, M., Ritzl, A., Falkai, P., Zilles, K., & Fink, G. R. (2004). Neural correlates of first-person perspective as one constituent of human self-consciousness. *Journal of Cognitive Neuroscience*, 16(5), 817-827. doi:10.1162/089892904970799