



Review

# Theory of mind—evolution, ontogeny, brain mechanisms and psychopathology

Martin Rüne\*, Ute Rüne-Cohrs

*Journal of Child Psychology and Psychiatry*, 46(10), 1045–1060, 2005

Received 21 April 2005; revised 2 August 2005; accepted 8 August 2005

**Abstract**

The theory of mind (ToM) is the ability to understand and predict the behavior of others based on their mental states. It is the ability most highly developed in humans. The evolutionary origins of theory of mind can be traced back in extent non-human primates; theory of mind probably emerged as a social response to increasingly complex primate social interaction. This sophisticated 'meta-cognitive' ability comes, however, at an evolutionary cost, reflected in the broad spectrum of psychopathological conditions. Extensive research into autistic spectrum disorders has revealed that theory of mind may be selectively impaired, leaving other cognitive faculties intact. Recent studies have shown that observed deficits in theory of mind task performance are part of a broader range of symptoms in schizophrenia, bipolar affective disorder, some forms of dementia, 'psychopathy' and in other psychiatric disorders. This article reviews the evolutionary psychology of theory of mind including its ontogeny and representation in the central nervous system, and studies of theory of mind in psychopathological conditions.

© 2005 Elsevier Ltd. All rights reserved.

*Keywords:* Theory of mind; Human evolution; Child development; Brain mechanisms of theory of mind; Psychopathology

6. Psychopathology of theory of mind .....	446
7. Developmental disorders .....	446
8. Personality disorders and other non-psychotic disorders .....	447
9. Schizophrenia and affective disorders .....	447
10. Brain damage and degenerative brain disorders .....	450
11. Discussion .....	451
References .....	452

## 1. Introduction

The term 'theory of mind' was originally proposed by primatologists Premack and Woodruff in a seminal article to suggest that chimpanzees may be capable of inferring mental states of their conspecifics (individuals of the same species) (Premack and Woodruff, 1978). Later on, the term was adopted by child psychologists to describe

the ontogenetic development of mental perspective taking in infants and young children (e.g. [Leslie, 1987](#)). In terms of psychology, the concept of disturbed theory of mind



maturtion (and large litter size; Joffe, 1997). Primates, however, are extreme K-strategists, that is, their offspring grows slowly, multiple births are unusual, and birth intervals are long. Moreover, they rarely consider the extension of the juvenile period in primates reaches a maximum in humans. Interesting in this regard is the fact that the length of the juvenile period in primates is also positively correlated with the size of the non-visual cortex in the same way as group size is; it does not correlate with the length of gestation, lactation, and reproductive life span. This finding could be interpreted as supporting the notion of slow maturtion to constraints of the social environment (Joffe, 1997). For example, the extension of the juvenile period in primates may have been crucial to acquire the vast amount of possible social behavioral 'strategies' (procedural rules) and when to employ these strategies (here, the term 'strategy' does not necessarily imply conscious awareness; Schmitt and Grimmer, 1997). This process is not merely time-consuming. The real-life opportunities of testing possible consequences of such social strategies are limited in number. It is, therefore, conceivable that the need for mental simulation of social interaction might have speeded up the evolution of theory of mind. If mental simulation is involved (see below), then theory of mind not only comprises the representation of the mental states of other individuals, but also one's own mental state (attachment theorists have termed this ability 'reflective functioning'; Fonagy, 1997).

### 3. On the origin of theory of mind

At birth, human infants are essentially immature. The

compared with children whose parents use such terms less often. In addition, the presence of older siblings speeds up young children's appreciation of other minds (overview in [Carpenter and Lewis, 2004](#)). Furthermore, it is noteworthy that, predictably from the evolutionary framework outlined above, these developmental steps of theory of mind constitute human universals. Although cross-cultural evidence is still limited, [Avis and Harris \(1991\)](#) have clearly shown that understanding false belief emerges at similar age in children of the *!K*, preliterate hunter-gatherers in southern Cameroon.

Finally, it is noteworthy that the development of theory of mind is clearly paralleled by language acquisition. In fact, understanding speaker's intention is a precondition of learning new words. As [Frith and Frith](#) have pointed out, random associations of utterances with objects rarely occur when young children learn to speak ([Frith and Frith, 2003](#)) and a child begins to use words undoubtedly referring to mental states such as 'I think' at the age of four—the watershed of distinguishing between own and other's mental states.

In contrast to our growing understanding of children's acquisition of theory of mind abilities, fairly little is known about the development of theory of mind in adult humans. Due to the fundamental role of subjective experience and recall of past social interactions in theory of mind performance, we would expect continuous refinement of mental state attribution throughout the adult human lifespan. On the other hand, selection pressure declines with age (particularly with respect to the post-reproductive lifespan). It is therefore conceivable that aging does not spare social cognitive abilities. Two studies of theory of mind abilities in older people have revealed conflicting results. [Happé et al. \(1998\)](#) found that people with a mean age of 73 years, although slower in performance, were superior on theory of mind tasks compared to adolescents and young adults of about 14 years and 22 years of age, respectively. In contrast, a recent study revealed the opposite, namely successive decline in theory of mind in adults aged between 60 and 74, and between 75 and 89, respectively, compared to younger adults ([Mylor et al., 2002](#)). Thus, at this stage there is still controversy whether and how theory of mind capacities change over the adult human lifespan.

#### 4. CNS-representation of theory of mind

If primitive brains, particularly neocortical structures, enlarged over evolutionary time due to selection pressures from the social environment, where exactly is theory of mind located in the human brain? Evidence comes from various sources. Comparative neuroanatomy and neurophysiology informs us which brain regions and corresponding functions came under selection pressure in non-human primates to evolve into the neural correlates of theory of mind in modern humans. In addition, functional brain

imaging studies and lesion studies in patients suffering from brain injuries or stroke may help localizing the brain circuits underlying theory of mind.

Before summarizing some of the most important empirical studies, it is necessary to point out that divergent theoretical conceptualizations of theory of mind exist. To some degree, this has considerable impact on how empirical findings are interpreted. (1) Drawing on [Fodor's \(1983\)](#) concept of modular organization of the human mind, some theorists advocate the existence of separate theory of mind module (e.g. [Scholl and Leslie, 1999](#)). Like other domain-specific cognitive capacities represented in the brain, which process only certain classes of information, the theory of mind mechanism is supposed to process information restricted to social inference. Cognitive mechanisms are assumed to work reliably, efficiently, and economically. According to the modular hypothesis, the development of theory of mind mainly depends on neurobiological maturation of the brain structures involved. Experience, on the contrary, may trigger the action of the theory of mind mechanism, but does not determine the makeup of the mechanism. (2) The 'metarepresentation' theory-theory (e.g. [Perner, 1991](#)) of theory of mind is somewhat distinct from the modular model. As [Fellmann \(1999\)](#) has summarized, the theory-theory proposal holds, similar to the modular theory does, that the entities and the causal principles of theory of mind are 5.4( e5o)-417.9(princ7itie)-8.6(7(m y)-531..3(ntl

Now, what can we learn from primitive research about theory of mind, in light of the fact that there is no unequivocal evidence of mental state attribution in non-human primates in general, and virtual absence of theory of mind in monkeys? Single cell recordings in non-human primates convey important information about candidate cerebral representations of cognitive precursor capacities of what we call 'true' theory of mind in humans (the term 'precursor capacities' by no means ought to suggest teleological interpretation, i.e. that something evolves in order to later suit a certain purpose).

A number of candidate structures have been identified in non-human primate brains that have undergone adaptive modifications to constitute in humans a neural network of theory of mind. Single cell recordings in macaque monkeys have revealed that neurons in the middle portion of the temporal lobe, particularly in the superior temporal sulcus (STS), selectively fire when monkeys observe the gaze direction of other monkeys. These neurons are also active when the animals observe goal-directed behavior (Gallese and Goldman, 1998). In humans, functional brain imaging studies have revealed that homologous regions of the temporal lobe is activated by observation of seemingly purposeful movements of inanimate objects (as opposed to random movements), and even when still photographs depict 'implied' motion (Kourtzi and Knwisher, 2000). For example, such activity could be elicited by showing human subjects pictures of a discus thrower in action, where such activity could be measured when the discus thrower was at rest. Activity in parts of the STS, therefore, is linked to the observation of intentional movements. Although this does not imply conscious awareness, the representation of 'intentions' is certainly a critical aspect of theory of mind. In a variety of functional imaging studies during theory of mind tasks performance the blood flow increased in a region of the STS adjacent to the part that was activated by monitoring biological motion (Grossmann and Ilke, 2002).

The temporal lobes of non-human primates also contain specific types of cells called 'mirror neurons' due to their unique quality to discharge during both the execution of a certain hand or mouth action or by the mere observation of the same behavior carried out by another individual. These neurons have also been found in greater density in the ventral premotor cortex of macaque monkeys, a region that is possibly homologous to the Broca's region in humans (Gallese and Goldman, 1998). In an ingenious series of experiments, the group of Rizzolatti has demonstrated that mirror neurons selectively fire when monkeys observe a hand movement of which the terminus is hidden from their view. In other words, a subset of mirror neurons is active when the monkey can only 'infer' or predict the result of the incompletely visible action (Umiltà et al., 2001). Mirror neurons may therefore be crucially involved in understanding action-goals. In humans, Fadiga et al. (1995) have shown in an experiment using transcranial magnetic stimulation (TMS) that the observation of a goal-directed hand movement

elicited enhanced motor evoked potentials (MEP). Notably, these enhanced MEPs could be measured precisely in those muscles the observer would use when carrying out the action himself.

The discovery of mirror neurons in humans offers an explanation of how the ability to imitate the actions of others could have evolved into the capacity to simulate the mental states of other individuals (i.e. theory of mind) (Williams et al., 2001). However, as Frith and Frith (1999, 2001) have pointed out, for theory of mind it is not sufficient to represent goal-directed actions. It is also necessary to be able to distinguish between behavior generated by self or others. And indeed, there are at least two other important brain regions involved in theory of mind. We believe that simulating other people's mental states does not necessarily involve conscious reflection, but is readily available to conscious awareness. For example, transference and counter-transference in dyadic psychotherapeutic settings always implicate the mutually, largely unconscious attribution of mental states such as intentions, desires and beliefs, and it is the goal of psychodynamic approaches to unveil these unconscious processes and make them accessible to the conscious mind. For conscious reflection on one's own and other's mental states an individual needs computational resources beyond the capacity for imitation and action simulation, and a candidate structure involved in this task is the inferior parietal cortex. Recent research using functional brain imaging has revealed that the left and right hemisphere are differentially involved in first versus third-person perspective. First-person perspective was shown to activate the left inferior parietal cortex, whereas third-person perspective activated the corresponding region on the right side of the human brain (Ruby and Decety, 2001). Interestingly, when a subject imitates the action of another person, more activation is found in the left inferior parietal cortex, but more activation is found on the opposite side when subjects view their actions being imitated. These experimental results support the assumption that the right inferior parietal cortex may be critical for consciously representing others' minds, whereas the left inferior parietal cortex may be involved in representing one's own mental states (Decety and Chaminade, 2005).

The other brain region that has consistently been shown to be engaged in theory of mind is the anterior cingulate cortex (ACC). The ACC receives input from the motor cortex and the spinal cord, from the ipsilateral prefrontal cortex, and from the thalamus and brainstem nuclei (Paus, 2001). It is highly heterogeneous in terms of its cytoarchitecture and functional organization. The ACC is now conceived of as an important mediator of motor control, cognition, and arousal regulation (Paus, 2001). In monkeys, for example, the most rostral part of the ACC is active prior to the execution of self-initiated movements (Frith and Frith, 1999). Most interesting from an evolutionary viewpoint and with respect to theory of mind is that the anterior ACC inconsistently forms a pre-cingulate sulcus

Table 1

Overview of brain imaging studies of theory of mind in chronological order

Author(s); published	Sample (n)	Mean age	Sex m/f	Brain imaging technique	ToM method/tasks	Activated brain regions in ToM tasks
Goel et al., 1995	9 healthy subjects	24.7	5/5	PET [ <sup>15</sup> O]H <sub>2</sub> O	Presentation of familiar and unfamiliar objects requiring inference of others' attribution of their function (i.e. ToM). One non-ToM condition involving inference of function of unfamiliar objects from their form. Two control conditions: visual and semantic attributes of known objects.	

Table 1 (  $\mu, \sigma$  )

---

---



th t is present in only 30–50% of individuals and possibly still under selection pressure (P us, 2001). This re

to go beyond the literal meaning of utterances by inferring what the speaker actually might have intended (Happé, 1994; Langdon et al., 2002b).

In adults with psychopathological conditions, short stories involving double bluff, mistakes, persuasions or white lies (Happé, 1994), cartoons or other visually presented materials have been used to assess theory of mind abilities. In theory of mind research in schizophrenia, for instance, short stories with or without use of props and picture sequencing tasks have been given to patients, as well as, tests of comprehension of hints behind indirect speech, metaphor and irony. Over the years, the pictorial theory of mind materials have been modified in order to better control for interference with attention, memory, 'general' intelligence, and verbalization. One problem in early studies in schizophrenia was that patients not only performed poorly on theory of mind tasks, but also often failed to correctly respond to the



T ble 2

---

---



(Lingdon et al., 2001; Pickup and Frith, 2001; overview in Frith, 2004). That is, these deficits are probably independent of other cognitive dysfunctions such as attention, set-shifting capacity, general intelligence and so forth (Lee

the results could largely be explained by this confounding rather than by a specific theory of mind deficit in AD.

By contrast, the frontoventral variant of frontotemporal dementia (fvFTD) is characterized by changes in personality and social behavior while most cognitive domains are relatively preserved, at least in the early stages of the disorder. From a clinical perspective this could be indicative of a selective theory of mind deficit in FTD. In a study, comparing patients with fvFTD with mild AD and healthy control subjects [Gregory et al. \(2002\)](#) found fvFTD patients to perform significantly worse on all theory of mind tasks with increasing impairment relative to task complexity. AD patients performed only on the more cognitively demanding second order false belief tasks indicating an interference with cognitive performance rather than impaired theory of mind per se. Interestingly, theory of





- rüne, M., 2005. Emotion recognition, 'theory of mind' and social behavior in schizophrenia. *Psychiatr. Res.* 133, 135–147.
- rüne, M., 2005b. 'Theory of mind' in schizophrenia: A review of the literature. *Schizophr. Bull.* 31, 21–42.
- rüne, M., Fodorstein, L., 2005. Proverb comprehension reconsidered - 'theory of mind' and the pragmatic use of language in schizophrenia. *Schizophr. Res.* 75, 233–239.
- Runet, E., Sforzi, Y., Hirdy-ylé, M.C., Decety, J., 2000. A PET investigation of the attribution of intentions with nonverbal tasks. *NeuroImage* 11, 157–166.
- Runet, E., Sforzi, Y., Hirdy-ylé, M.C., Decety, J., 2003. Abnormalities of brain function during nonverbal theory of mind tasks in schizophrenia. *Neuropsychologia* 41, 1574–1582.
- Witell, J.K., van der Wees, M., Swinnen, H., van der Groot, R.J., 1999. Theory of mind and emotion-recognition functioning in autistic spectrum disorders and in psychiatric control and normal children. *Dev. Psychopathol.* 11, 39–58.

Grossmann, E.D., Blake, R., 2002. Attentional bias during visual perception of biological motion. *Neuron* 35, 1167–1175.  
Happé

- Rowe, A.D., Hullock, P.R., Polkey, C.E., Morris, R.G., 2001. 'Theory of mind' impairments and their relationship to executive functioning following front lobe excisions. *Brain* 124, 600–616.
- Ruby, P., Decety, J., 2001. Effect of subjective perspective taking during simulation of action: A PET investigation of agency. *Nature Neurosci.* 4, 546–550.
- Russell, J., Muthner, N., Shupe, S., Tidswell, T., 1991. The 'windows to the sky' measure of strategic deception in pre-schoolers and autistic subjects. *Br. J. Dev. Psychol.* 9, 331–349.
- Russell, T.A., Rubi, K., Fullmore, E.T., Soni, W., Suckling, J., Brammer, M.J., Simmons, A., Williams, S.C., Shurman, T., 2000. Exploring the social brain in schizophrenia: Left prefrontal underactivation during mental state attribution. *Am. J. Psychiatry* 157, 2040–2042.
- Silzman, J., Strauss, E., Hunter, M., Archibald, S., 2000. Theory of mind and executive functions in normal human aging and Parkinson's disease. *J. Int. Neuropsychol. Soc.* 6, 781–788.
- Simonson, D., Apperly, I. A., Kirchner, M., Uhlhaas, U., Humphreys, G.W., 2005. Seeing it my way: Case of selective deficit in inhibiting self-perspective. *Brain* 128, 1102–1111.
- Sirfati, Y., Hrdy, J.L., Esche, C., Widlöcher, D., 1997. Attribution of intentions to others in people with schizophrenia: non-verbal exploration with comic strip. *Schizophr. Res.* 25, 199–209.
- Sirfati, Y., Hrdy, J.L., Brunet, E., Widlöcher, D., 1999. Investigating theory of mind in schizophrenia: Influence of verbalization in disorganized and non-disorganized patients. *Schizophr. Res.* 37, 183–190.
- Saxe, R., Kanwisher, N., 2003. People thinking about thinking people. The role of the temporo-parietal junction in "theory of mind". *NeuroImage* 19, 1835–1842.
- Saxe, R., Carey, S., Kanwisher, N., 2004. Understanding other minds: Linking developmental psychology and functional neuroimaging. *Annu. Rev. Psychol.* 55, 87–124.
- Schmitt, A., Grimmer, K., 1997. Social intelligence and success: Don't be too clever in order to be smart. In: Whiten, A., Byrne, R.W. (Eds.), *Machiavelli in Intelligence II. Extensions and Evaluations*. Cambridge University Press, Cambridge, pp. 86–111.
- Scholl, J., Leslie, A., 1999. Modularity, development and 'theory of mind'. *Mind* 108, 131–153.
- Shmuy-Tsoory, S.G., Tomer, R., Berger, D., Goldsher, D., Aharon-Peretz, J., 2005. Impaired "affective theory of mind" is associated with right ventromedial prefrontal damage. *Cogn. Behav. Neurol.* 18, 55–67.
- Siegal, M., Carrington, J., Dalen, M., 1996. Theory of mind and pragmatic understanding following right hemisphere damage. *Brain* 119, 40–50.
- Simpson, J., Done, J., Vallée-Tourangeau, F., 1998. An unreasoned approach: critique of research on reasoning and delusions. *Consciousness* 3, 1–20.
- Simpson, J., 2002. M2.1ufurdR dé3TD[(ch24.)-6FostA.,