



**HAL**  
open science

## Self-deception: Adaptation or by-product?

Hugo Mercier

► **To cite this version:**

Hugo Mercier. Self-deception: Adaptation or by-product?. Behavioral and Brain Sciences, 2011, 34 (1), pp.35. hal-00907426

**HAL Id: hal-00907426**

**<https://hal.science/hal-00907426v1>**

Submitted on 21 Nov 2013

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Religion and Theology Project”), both coordinated from the Centre for Anthropology and Mind at the University of Oxford.

## Self-deception: Adaptation or by-product?

doi:10.1017/S0140525X10002281

Hugo Mercier

*Philosophy, Politics, and Economics Program, University of Pennsylvania, Philadelphia, PA 19104.*

[hmercier@sas.upenn.edu](mailto:hmercier@sas.upenn.edu)

<http://sites.google.com/site/hugomercier/>

**Abstract:** By systematically biasing our beliefs, self-deception can endanger our ability to successfully convey our messages. It can also lead lies to degenerate into more severe damages in relationships. Accordingly, I suggest that the biases reviewed in the target article do not aim at self-deception but instead are the by-products of several other mechanisms: our natural tendency to self-enhance, the confirmation bias inherent in reasoning, and the lack of access to our unconscious minds.

In their target article, von Hippel & Trivers (VH&T) defend the hypothesis that many psychological biases are by nature self-deceptive. Their rationale is the following: People get caught lying because of “signs of nervousness, suppression, cognitive load, and idiosyncratic sources.” In order to make deception detection less likely, these superficial cues should be reduced or eliminated. Given that these cues all stem from the fact that we have to keep in mind the truth and the lie – which we know when we lie – it would make sense for people to actually believe the lies they tell – to self-deceive. However, VH&T fail to take into account that one of the most important cues to deception is lack of consistency (DePaulo et al. 2003). When people are confronted with communicated information, they evaluate its internal consistency as well as its consistency with their previously held beliefs (Sperber et al. 2010). Any benefit gained by lying to ourselves in terms of suppression of superficial cues compromises our ability to make up lies that will pass this consistency test. VH&T also suggest that self-deception could be adaptive because it makes it easier for deceivers to maintain that they had no deceptive intent (their “second corollary”). However, here again self-deception has the potential to backfire. When we know we lied, we can recognize that we did it and feel guilty, apologize, try to make amends, and so forth. These can be essential to the maintenance of trust (Kim et al. 2004; Schweitzer et al. 2006). If we do not even realize that we are trying to deceive, any accusation – however well founded – is likely to be received with aggravation. Thus, by suppressing any common ground between self and audience, self-deception critically endangers the maintenance of trust.

The costs of self-deception weaken the principled case for its adaptiveness. But how are we, then, to account for the evidence that VH&T present in support of their hypothesis? In what follows, I will argue that this evidence can be better explained as the by-product of other mechanisms. Many results presented in the target article show that people have a strong tendency to self-enhance, and that we often do so without even realizing it. This claim would be hard to dispute. For these results to support VH&T’s hypothesis, the lack of more veridical information processing must stem from the adaptive character of self-deception. But it is more plausible that the lack of veridical information processing is a simple result of the costs it would entail. It is possible here to make an analogy with other systematically biased mechanisms. For instance, following a simple cost-benefit analysis, it is reasonable to surmise that a mechanism aimed at the detection of poisonous food should be systematically biased toward the “poisonous” verdict. The lack of a less biased information processing requires no explanation beyond this cost-benefit analysis. If a given degree of self-enhancement is adaptive in and of itself, then this is enough to explain why less biased mechanisms would be

superfluous. Contrary to what VH&T claim, the fact that we can sometimes engage in more veridical processing does not show that the mechanisms have a self-deceptive purpose. By analogy, our poisonous food detector could also be more or less biased – depending on the individual who is providing us with the food, for instance – without having self-deception as its goal.

The authors’ case rests not only on our ability to sometimes turn off our biases and engage in veridical processing, but also on the conditions that trigger veridical processing. More specifically, they claim that because self-affirmation or cognitive load manipulations can make us less biased, then any bias that is otherwise present is likely to be self-deceptive. But these findings can also be explained by the effect of these manipulations on the use of high-level processing – in particular, reasoning. Self-affirmation manipulations can be understood as belonging to a larger group of manipulation – including self-esteem and mood manipulations (e.g., Raghunathan & Trope 2002) – that reduce our tendency to engage in some types of high-level processing (Schwarz & Skurmik 2003). Likewise, cognitive load will automatically impair high-level processing. Reasoning is one of the main mechanisms that can be affected by these manipulations, and the confirmation bias exhibited by reasoning is the source of many of the biased results described by VH&T (Nickerson 1998). It is therefore not surprising that self-affirmation or cognitive load manipulations should make us appear less biased. However, it has been argued that the confirmation bias does not have a self-deceptive function and that it is instead the result of the argumentative function of reasoning (Mercier & Sperber, in press). Accordingly, when reasoning is used in a natural setting (such as group discussion), the confirmation bias does not systematically lead to biased beliefs (Mercier & Landemore, in press). Thus most of the results used by the authors can be accounted for as a by-product of a confirmation bias inherent in reasoning that does not have a self-deceptive function.

Finally, a case can also be made against the authors’ interpretation of the dual-process literature. According to VH&T, “these dissociations [between, e.g., implicit and explicit memory] ensure that people have limited conscious access to the contents of their own mind and to the motives that drive their behavior.” For this statement to be correct, conscious access to the content of our own mind would have to be a given from which it can sometimes be useful to deviate. But this is not the case. Being able to know the content of our own minds is a very costly process. In fact, it is sometimes speculated that there was little evolutionary advantage to be gained by knowing ourselves, and that this ability is a mere by-product of our ability to understand others (e.g., Carruthers 2009b). If *not* knowing ourselves – or knowing ourselves very imperfectly – is the baseline, then dissociations between conscious and unconscious processes require no further explanation. These dissociations cannot *ensure* us against a self-knowledge that we have no reason to possess in the first place.

Trying to elucidate the ultimate function of our cognitive biases is a very worthwhile endeavor that is bound to lead to a much deeper understanding of human psychology. However, for VH&T’s specific hypothesis to be truly convincing, they would need to provide stronger evidence, such as the direct experimental tests – whose absence they repeatedly deplore – of their theory.

## Representations and decision rules in the theory of self-deception

doi:10.1017/S0140525X1000261X

Steven Pinker

*Department of Psychology, Harvard University, Cambridge, MA 02138.*

[pinker@wjh.harvard.edu](mailto:pinker@wjh.harvard.edu) <http://pinker.wjh.harvard.edu>

**Abstract:** Self-deception is a powerful but overapplied theory. It is adaptive only when a deception-detecting audience is in the loop, not when an