

Stephen M. Fleming
Christopher D. Frith *Editors*

The Cognitive Neuroscience of Metacognition

 Springer

The Cognitive Neuroscience of Metacognition

Stephen M. Fleming · Christopher D. Frith
Editors

The Cognitive Neuroscience of Metacognition

 Springer

Editors

Stephen M. Fleming
Center for Neural Science
New York University
New York, NY
USA

Christopher D. Frith
Aarhus University Hospital
Aarhus
Denmark

ISBN 978-3-642-45189-8 ISBN 978-3-642-45190-4 (eBook)
DOI 10.1007/978-3-642-45190-4
Springer Heidelberg New York Dordrecht London

Library of Congress Control Number: 2014930091

© Springer-Verlag Berlin Heidelberg 2014

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law. The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

Acknowledgments

The authors gratefully acknowledge the support of All Souls College, Oxford (CDF), and the Wellcome Trust (SMF) during the preparation of this book.

Contents

1	Metacognitive Neuroscience: An Introduction	1
	Stephen M. Fleming and Christopher D. Frith	
Part I Quantifying Metacognition for the Neurosciences		
2	Quantifying Human Metacognition for the Neurosciences	9
	Bennett L. Schwartz and Fernando Díaz	
3	Signal Detection Theory Analysis of Type 1 and Type 2 Data: Meta-d', Response-Specific Meta-d', and the Unequal Variance SDT Model	25
	Brian Maniscalco and Hakwan Lau	
4	Kinds of Access: Different Methods for Report Reveal Different Kinds of Metacognitive Access	67
	Morten Overgaard and Kristian Sandberg	
5	The Highs and Lows of Theoretical Interpretation in Animal-Metacognition Research	87
	J. David Smith, Justin J. Couchman and Michael J. Beran	
Part II Computational Approaches to Metacognition		
6	A Computational Framework for the Study of Confidence Across Species	115
	Adam Kepecs and Zachary F. Mainen	
7	Shared Mechanisms for Confidence Judgements and Error Detection in Human Decision Making	147
	Nick Yeung and Christopher Summerfield	

8	Metacognition and Confidence in Value-Based Choice	169
	Stephen M. Fleming and Benedetto De Martino	
9	What Failure in Collective Decision-Making Tells Us About Metacognition	189
	Dan Bang, Ali Mahmoodi, Karsten Olsen, Andreas Roepstorff, Geraint Rees, Chris Frith and Bahador Bahrami	
 Part III Cognitive Neuroscience of Metacognition		
10	Studying Metacognitive Processes at the Single Neuron Level . . .	225
	Paul G. Middlebrooks, Zachary M. Abzug and Marc A. Sommer	
11	The Neural Basis of Metacognitive Ability.	245
	Stephen M. Fleming and Raymond J. Dolan	
12	The Cognitive Neuroscience of Metamemory Monitoring: Understanding Metamemory Processes, Subjective Levels Expressed, and Metacognitive Accuracy	267
	Elizabeth F. Chua, Denise Pergolizzi and R. Rachel Weintraub	
13	Metacognitive Facilitation of Spontaneous Thought Processes: When Metacognition Helps the Wandering Mind Find Its Way.	293
	Kieran C. R. Fox and Kalina Christoff	
14	What is the Human Sense of Agency, and is it Metacognitive? . . .	321
	Valerian Chambon, Elisa Filevich and Patrick Haggard	
 Part IV Neuropsychiatric Disorders of Metacognition		
15	Failures of Metacognition and Lack of Insight in Neuropsychiatric Disorders.	345
	Anthony S. David, Nicholas Bedford, Ben Wiffen and James Gillean	
16	Judgments of Agency in Schizophrenia: An Impairment in Autoegetic Metacognition	367
	Janet Metcalfe, Jared X. Van Snellenberg, Pamela DeRosse, Peter Balsam and Anil K. Malhotra	
17	Metacognition in Alzheimer’s Disease	389
	Stephanie Cosentino	

Chapter 1

Metacognitive Neuroscience: An Introduction

Stephen M. Fleming and Christopher D. Frith

Abstract The past two decades have witnessed the birth of the cognitive neurosciences, spurred in large part by the advent of brain scanning technology. From this discipline our understanding of psychological constructs ranging from perception to memory to emotion have been enriched by knowledge of their neural underpinnings. The same is now true of metacognition. This volume represents a first attempt to take stock of the rapidly developing field of the neuroscience of metacognition in humans and non-human animals, and in turn examine the implications of neuroscience data for psychological accounts of metacognitive processes.

In the introduction to a recent volume on metacognition, Michael Beran and colleagues wrote, “The very idea of publishing another book on metacognition needs a word of justification as there is already a number of collections available in this rapidly growing field” [1]. As the book you are holding follows their excellent volume, it is even more pressing for us to address this question. Fortunately, it is relatively straightforward to do so. The past two decades have witnessed the birth of the cognitive neurosciences, spurred in large part by the advent of brain scanning technology. From this discipline our understanding of psychological constructs ranging from perception to memory to emotion have been enriched by knowledge of their neural underpinnings. The same is now true of metacognition.

S. M. Fleming (✉)

Department of Experimental Psychology, University of Oxford, Oxford, UK
e-mail: sf102@nyu.edu

S. M. Fleming

Center for Neural Science, New York University, New York, USA

C. D. Frith (✉)

Wellcome Trust Centre for Neuroimaging, University College London, London, UK
e-mail: c.frith@ucl.ac.uk

C. D. Frith

Interacting Minds Centre, Aarhus University, Aarhus, Denmark

This volume represents a first attempt to take stock of the rapidly developing field of the neuroscience of metacognition in humans and non-human animals, and in turn examine the implications of neuroscience data for psychological accounts of metacognitive processes.

1.1 Defining Metacognition

Before previewing the chapters in this book, let us start with a definition of metacognition. There are at least two reasons why the term metacognition sometimes leads to confusion. The first is that it evokes different domain-specific associations. For example, metacognition may take on different connotations within education research, in memory research and in perception research. A second, more subtle reason is that metacognition is sometimes associated with conscious (and by implication, human) reflective awareness. We think this latter reason presents a barrier to a satisfactory computational and biological explanation of metacognition, so we take some time here to outline what such an explanation might look like.

The simplest definition of metacognition is cognition about cognition. A metacognitive process is meta-level with respect to an object-level cognitive process. This framework was originated by Flavell [2], and later Nelson and Narens [6], in the study of learning and memory. Memory still gives us our most subjectively vivid examples of metacognition: the decision to stop revising for an exam and the feeling of a tip-of-the-tongue experience are both quintessential metacognitive experiences. Importantly however, these examples of metacognition happen to be associated with explicit conscious awareness. While metacognition may be accompanied by conscious awareness in humans, this need not be the case, suggesting a division between “explicit”, conscious metacognition and “implicit” metacognition [3, 4, 8]. To take a concrete example rooted in neuroscience, consider that a visual image of a face leads to activity in the fusiform cortex. We can think of the fusiform response to the face as an “object-level” process. A meta-level process (say in the prefrontal cortex) may represent confidence that fusiform activity is signaling a face is present. This meta-level process may or may not be associated with reflective awareness, but it is nevertheless metacognitive.

Appreciating this point helps admit a broader range of evidence in the study of the neuroscience of metacognition. For example, finding neurons in non-human animal brains that covary with confidence in a previous decision would reveal a plausible neural substrate of metacognition, despite the difficulty of assessing whether metacognitive judgments are explicit or implicit in non-human animals. It also makes clear that an important question for future study is the difference in neural implementation between implicit and explicit metacognition in humans, and the degree to which this neural circuitry is shared with non-human animals. Crucially neuroscience may be able to provide a window on the representational architecture of a metacognitive system, a point to which we turn next.

1.2 Why Neuroscience?

Constraints on neurobiological implementation serve to shape psychological theory, and these constraints might prove particularly important in the study of metacognition. It is helpful to draw an analogy with the well-established field of memory neuroscience. The discovery of intact implicit memory despite impaired episodic memory in patient HM has provided a strong constraint on every systems-level theory of memory since the 1950s. More recently brain imaging technology has revealed links between components of psychological models and their putative divisions of labour at an implementational level [5].¹ Psychological models of metacognition are particularly ripe for such an analysis. Consider again the example of face perception. Imagine that a subject's task is to rate their confidence in having seen a series of blurry faces while in a brain scanner. There are at least two possible neural and psychological accounts of how this second-order confidence judgment is made. First, it might be achieved by a direct readout of properties of an object-level representation in the fusiform cortex (a non-metacognitive implementation of a second-order behavior). Second, the judgment might rely on a meta-level representation of a subset of properties of the fusiform activity, say in prefrontal cortex. These inevitably over-simplified hypotheses make neuroscientific predictions: in the former case, lesions to putative meta-level representations in prefrontal cortex should not affect confidence judgments; in the latter case, confidence judgments may be selectively affected by lesions while leaving first-order behavioural responses (such as a forced-choice judgment) to the face intact.

This schematic example makes clear that cognitive neuroscience has much to offer a psychological-level understanding of metacognition. We might find that some behaviours traditionally thought of as metacognitive are implemented in a manner that does not require meta-level representations; in turn, a detailed understanding of those systems that do permit meta-level representation will refine psychological-level models. Finally, we note that by rooting our psychological-level models in cognitive neuroscience the division between meta- and object-level becomes less sharp and more nuanced, reflecting the intricate interplay between higher-order and primary sensory and mnemonic brain areas.

It should be clear from previous paragraphs that neuroscience does not stand apart from behavioural measurement or theoretical models. In this spirit, we have included an opening section entitled "Quantifying metacognition for the neurosciences" which reviews types of metacognitive judgments and theoretical and computational frameworks within which to understand these judgments. We hope that these chapters form a self-contained section while not retreading ground that has already been covered in excellent previous collections.

¹ The next trend will be to understand how individual functional specialisations predicted by psychological-level models are integrated via analysis of functional and structural connectivity between brain regions.

1.3 An Outline of the Book

The first section, “Quantifying metacognition for the neurosciences”, outlines behavioural and analytic techniques important for the development of metacognitive neuroscience. Bennett and Schwartz (Chap. 2) emphasise that human metacognitive judgments are likely to arise from multiple component psychological processes, using the tip-of-the-tongue experience as a case in point. A pressing issue in the quantification of metacognition is distilling a metacognitive component of behavior from other confounding factors. Lau and Maniscalco (Chap. 3) present an overview of their recently developed computational measure of metacognitive efficiency that achieves this control within a signal detection theoretic framework. Overgaard and Sandberg (Chap. 4) review different types of metacognitive report about perception, and discuss how different types of report may map onto different kinds of metacognitive access. Finally, Smith, Couchman and Beran (Chap. 5) outline progress on the quantification of metacognition in non-verbal animal species, and consider the various theoretical interpretations of these data.

The second section, “Computational approaches”, focuses on the utility of computational models for bridging behavioural and neural data. Computational models have proven very useful for revealing neural correlates of “hidden” internal states that would not otherwise be apparent in analysis of behaviour alone (e.g. [7]). Kepecs and Mainen (Chap. 6) present a signal detection theoretic model of decision confidence that can be powerfully applied to understanding confidence signals in neural data across different species. Yeung and Summerfield (Chap. 7) outline evidence accumulation models as a common framework in which to understand studies of error detection and confidence judgments. Fleming and De Martino (Chap. 8) present a case study of the application of an evidence accumulation model to understand the neural basis of confidence and metacognition during human value-based decision-making. Bang, Mahmoodi, Olsen, Roepstorff, Rees, Frith and Bahrami (Chap. 9) outline recent developments in modeling metacognitive judgments during social decision-making.

A third section reviews the cognitive neuroscience of metacognition across several inter-related areas of study. Middlebrooks, Abzug and Sommer (Chap. 10) review three recent studies examining metacognition-related activity in single neurons recorded from macaque monkeys and rats. Fleming and Dolan (Chap. 11) review the psychological and neural basis of metacognitive accuracy in humans, drawing on data from studies of perception, decision-making and memory. Chua, Weintraub and Pergolizzi (Chap. 12) present a comprehensive review of cognitive neuroscience studies on metacognition of human memory, covering the neural basis of subjective confidence and metacognitive accuracy. Metacognition shares an intriguing relationship with studies of human mind-wandering, which at first glance seems to be the opposite of deliberate metacognitive monitoring and control. However scholarly analysis of this link has been lacking: in their chapter, Fox and Christoff (Chap. 13) outline how metacognition and mind-wandering may

share more than just an antagonistic relationship, with metacognition actively guiding the wandering mind. Finally, Chambon, Filevich and Haggard ([Chap. 14](#)) consider whether the human sense of agency should be considered a metacognitive object, and review recent work from their laboratory on understanding the neural basis of agency.

In a final section we turn to the interface between neuropsychiatric disorders and metacognition. A neuroscience of metacognition has great promise for understanding metacognitive deficits observed in neuropsychiatric disorders such as Alzheimer's disease and schizophrenia. In turn, studies of metacognition in neuropsychiatric patients can provide a novel window onto the mechanisms of metacognition. The chapter by David, Bedford, Wiffen and Gilleen ([Chap. 15](#)) describes the link between metacognitive failures and lack of insight in psychosis, noting that while insight appears separable from primary symptomology, the relationship between cognitive and clinical insight remains poorly understood. Metcalfe, Van Snellenberg, DeRosse, Balsam and Malhotra ([Chap. 16](#)) describe a study of judgments of agency in schizophrenic subjects, revealing impairment in self-related, or "autonoetic" metacognition. Finally, Cosentino ([Chap. 17](#)) reviews studies aimed at understanding the impairments of metacognition that often occur in Alzheimer's disease.

1.4 Conclusions

Neuroscience has had a dramatic impact on our understanding of individual domains of cognition, from vision to memory. We hope that a cognitive neuroscience of metacognition will bear similar fruits. This is an exciting time to be a metacognition researcher: cognitive neuroscience is maturing as a field, and has available a wealth of tools with which to investigate the biological basis of mind. These tools, combined with advanced behavioural techniques and computational modeling, have great promise to advance our nascent understanding of metacognition.

References

1. Beran MJ, Brandl J, Perner J, Proust J (2012) Foundations of metacognition. Oxford University Press, Oxford
2. Flavell J (1979) Metacognition and cognitive monitoring: a new area of cognitive-developmental inquiry. *Am Psychol* 34(10):906
3. Fleming SM, Dolan RJ, Frith CD (2012) Metacognition: computation, biology and function. *Philos Trans R Soc Lond B Biol Sci* 367(1594):1280–1286
4. Frith CD (2012) The role of metacognition in human social interactions. *Philos Trans R Soc Lond B Biol Sci* 367(1599):2213–2223

5. Henson R (2005) What can functional neuroimaging tell the experimental psychologist? *Q J Exp Psychol* 58(2):193–233
6. Nelson TO, Narens L (1990) Metamemory: a theoretical framework and new findings. *Psychol Learn Motiv: Adv Res Theory* 26:125–173
7. O'Doherty JP, Hampton A, Kim H (2007) Model-based fMRI and its application to reward learning and decision making. *Ann N Y Acad Sci* 1104:35–53
8. Reder LM, Schunn CD (1996) Metacognition does not imply awareness: strategy choice is governed by implicit learning and memory. In: Reder LM (ed) *Implicit memory and metacognition*. Erlbaum, Mahwah

Part I
Quantifying Metacognition
for the Neurosciences

Chapter 2

Quantifying Human Metacognition for the Neurosciences

Bennett L. Schwartz and Fernando Díaz

Abstract The study of metacognition examines the relation between internal cognitive processes and mental experience. To investigate metacognition researchers ask participants to make confidence judgments about the efficacy of some aspect of their cognition or memory. We are concerned that, in our haste to understand metacognition, we mistakenly equate the judgments we elicit from participants with the processes that underlie them. We assert here that multiple processes may determine any metacognitive judgment. In our own research, we explore the tip-of-the-tongue phenomenon (TOT). Both behavioral and neuroscience evidence suggest that a number of processes contribute to the TOT. The fMRI, electroencephalography (EEG), and magnetoencephalography (MEG) data find that retrieval failure and TOT experience map onto different areas of the brain and at different times following the presentation of a stimuli. Behavioral data suggest that there are multiple cognitive processes that contribute to the TOT, including cue familiarity and the retrieval of related information. We assert that TOTs occur when retrieval processes fail and a separate set of processes monitor the retrieval failure to determine if the target can eventually be recovered. Thus, the TOT data support a model in which different underlying processes are responsible for the cognition and the metacognition that monitors it. Thus, understanding any metacognitive judgment must involve understanding the cognition it measures and the multiple processes that contribute to the judgment.

Although this chapter concerns metacognition, we start with psychophysics. The earliest psychological science was that of psychophysics, which was (and still is) the study of the relation between external energy and internal experience [14].

B. L. Schwartz (✉)

Department of Psychology, Florida International University, University Park,
Miami, FL 33199, USA
e-mail: bennett.schwartz@fiu.edu

F. Díaz

University of Santiago de Compostela, Santiago de Compostela, Spain

In psychophysics, we measure the wavelength of light and correlate it with the experience of color. Or we measure the frequency of sound and correlate it with the perception of pitch. Time and time again, such correlations yield replicable patterns within and across people. We argue here that, at its heart, metacognition aims to achieve something similar to the goals of psychophysics. However, metacognition's goal is to study the relation between internal cognitive processes and mental experience. For example, we study the strength of a memory and its relation to a subjective judgment of learning. Or we study the accessibility of an item and its correlation with tip-of-the-tongue (TOT) experiences. Such a goal would have likely been impossible to achieve 150 years ago, but now, with our understanding of cognitive psychology and neuroscience, it is possible to study internal mental experiences, such as metacognition, in a robust scientific fashion.

Cognitive processes are internal processes that carry out a particular function. These cognitive processes are of course, based on physical networks in the brain. Some cognitive processes may be open to introspection, and others may not. Thus, retrieval processes produce conscious memories, though the process itself is difficult if not impossible to introspect on. We define mental experiences as our subjective states that rise into consciousness. For example, our cognitive processes retrieve a memory of snorkeling on a coral reef, but the recollected experience of color, excitement, and warm water is our mental experience. The final part of this equation is behavior. Scientific psychology is rooted in the study of behavior. This applies to metacognition as well. As good experimental psychologists, we assess both cognitive processes and mental experience by observing and eliciting behavior.

2.1 The Doctrine of Concordance

We turn to the relation between cognition, behavior, and experience and Tulving's [40] doctrine of concordance. Tulving argued that there was a traditional bias in psychology, including cognitive psychology, to assume a strong correlation between cognitive processes, behavior, and experience. That is, a particular cognitive process, such as retrieval, is associated with a particular behavior, a verbal description of an earlier episode, and that this behavior is always associated with a particular conscious experience, in this case mental time travel. Tulving [40] claimed that this model no longer worked—there were too many demonstrations of conscious experience not accompanying a particular behavior to warrant its challenge. He cited studies on implicit memory, in particular, in which memory processes create a change in behavior, but without the accompanying mental experience. More recently, we can point to research in which mental experiences of memory arise from cognitive processes not tied to the retrieval process. For example, Cleary et al. [9] showed that *déjà vu* experiences arise when a familiarity experience occurs without corresponding retrieval of event details.

Tulving's [40] challenge to the assumptions of the doctrine of concordance underlies the basis of a great deal of research in metacognition (see [31, 32, 35]).

Metacognitive experiences arise from cognitive processes and correspond to particular behaviors. For example, an object is recognized as having been seen before (cognitive process), accompanied by an experience of confidence, and the person then says that they know the answer (behavior). However, the challenge to the doctrine arises from the repeated observations that the cognitive processes that give rise to metacognition are not the same cognitive processes that produce the behavior. One set of processes drives the recall of information, but another set of processes drives our awareness of it. In the case of retrieval, most metacognition research shows that the process that produces the metacognitive experience of confidence is dissociable from the process that elicits the retrieval. Retrieval success is determined by the strength of the target, but feeling-of-knowing judgments are determined by the ease of access to partial information and the strength of the cue [2, 34, 39]. In the aforementioned déjà vu experience, recollective processes convince the participant that the event is new, but familiarity processes drive the déjà vu experience. Thus, the processes that produce metacognition are not identical to the processes that produce the cognition they reflect.

Based on the data and reasoning above, some theorists view metacognition as heuristic in nature, that is, that metacognitive processes are not the same as the cognitive processes they monitor. Our metacognitive processes accurately predict memory performance because the processes that produce metacognition are correlated with the processes that produce cognition. Thus, as cue familiarity is correlated with target retrieval, using cue familiarity to predict recall leads to accurate feeling-of-knowing judgments. Such a heuristic model, therefore, explains why metacognition is generally accurate at predicting performance, but also why it sometimes does not predict performance; it depends on whether there is a strong positive correlation between the processes that lead to the behavior and those that lead to the internal mental experience. Thus, metamemory fails to predict performance when the metacognitive processes are not correlated with the cognitive processes used in the base process. For example, Benjamin et al. [3] found that memory strength predicted recall, but that ease of earlier processing predicted participants' judgments of learning (henceforth, JOLs), thus leading to a negative correlation between judgment and performance. To summarize, metacognition is a heuristic—it capitalizes on processes that correlate with cognitive processes and allow the organism to predict ongoing processes. Metacognitive judgments measure metacognitive experience and, in turn, are based upon underlying cognitive processes that produce them. These cognitive processes are correlated with the cognitive processes they are monitoring, but seldom identical. Thus, the cognitive processes that produce feeling of knowing may be partially based on cue familiarity, but the processes they monitor are based on retrieval strength. A generation of research has documented such dissociation in process.

2.2 Metacognition: An Introduction

Metacognition's chief empirical tool is to ask participants to make confidence judgments about the efficacy of some aspect of their cognition. The mainstay of metacognitive research, for largely historical reasons, has been judgments concerning memory [13]. Within this domain, one finds a plethora of judgments related to different aspects of the learning and retrieval processes. Ease-of-learning judgments are assessments of perceived difficulty of items in advance of study. Judgments of learning (JOLs) are assessment of whether an item being studied now will be recalled later. Turning to retrieval processes, feeling-of-knowing judgments (FOKs) are an assessment that a currently unrecalled item will be recognized later. TOT states refer to the strong feeling that a currently unrecalled item will be recalled shortly. Finally, confidence judgments can assess the feeling that a retrieved answer is actually correct. These are the main judgments used in metamemory research, although a variety of other judgments have been employed to assess specific aspects of memory (see [13], for a review).

2.3 Multiple Processes Underlie Judgments

Although it is largely accepted that cognition and metacognition are dissociable, there is less consensus on the relation between metacognitive processes and metacognitive judgments. We are concerned that, in our haste to understand metacognition, we mistakenly equate the judgments we elicit from participants with the processes that underlie them. For example, there have been disagreements as to whether FOKs are caused by cue familiarity, partial information, retrieval, or unconscious access to the target [13, 20, 26]. Clearly, all three may contribute to the judgments but in order to assert this we need to see that judgments and process are not identical.

In this chapter we challenge the assumption that there is a 1:1 correspondence between the processes that drive metacognition and the specific judgments that we make concerning metacognition. We assert here that multiple processes may determine any metacognitive judgment, and thus the judgments we measure are not pure indicators of the metacognitive processes that we are interested in. To be more concrete, consider the feeling of the TOT state [6, 35]. It is likely that the experience of the TOT is determined by the familiarity of the cue, the amount and intensity of related information retrieved, the amount and intensity of partial information retrieved, and the activation strength of the item itself. As such, the TOT experience cannot serve as a stand-in for any single one of these processes. Any consideration of the TOT requires consideration of all of these cognitive processes. We will consider the TOT and its etiology at length later in the chapter.

We propose that multiple cognitive processes may underlie any particular metacognitive judgment, be it TOT, JOL, or a confidence judgment. Although this

is not a controversial statement, its implications are that each judgment itself does not perfectly reflect one metacognitive process, as they are often thought to do. This becomes important in discussing neuroimaging studies of metacognition, in which one looks at the neural correlates of a particular metacognitive judgment. It may be hard—via one study—to determine which neural area is associated with which cognitive process or which neural area is associated with which mental experience because each is multiply determined. This further complicates the issue of the relation between process, task, and subjective experience.

Does this leave an awful intractable mess? Behavior need not be correlated with subjective experience (metacognition), subjective experience may be correlated but not caused by the same processes that drive the behavior, and all of these may be influenced by multiple cognitive factors and driven by diverse mechanisms neurally. Not necessarily; just as psychophysics established principles that governed the relation between physical energy and subjective experience, we are committed to the view that studies of metacognition can develop principles that govern the relation between internal cognitive processes and subjective experience.

It is with these issues and concerns in mind that each author of this chapter began investigating the TOT phenomenon. The TOT offers a number of features that make it an excellent case study in the relation among process, behavior, and experience. TOTs are a universal experience, they are relatively frequent in everyday life, and they are easy to induce in the laboratory. More importantly, TOTs are closely linked to a particular set of cognitive processes, namely those of retrieval, and TOTs engage a specific experience linked to a specific referent, namely a particular word. These characteristics make TOTs a good candidate for a case study in the scientific examination of human phenomenology, in particular the relation between subjective experience, cognition, and behavior [35].

2.4 Tip-of-the-Tongue States

A TOT state is the feeling that we will be able to recall a currently unrecalled word. In short, a TOT is a feeling of temporary inaccessibility. We argue that there is a strong correlation between the feeling of temporary inaccessibility (the phenomenological TOT) and actual temporary inaccessibility (sometimes, called the cognitive TOT, [1]). In general, the TOT experience is predictive of resolution of temporary inaccessibility (but see [30]). This positive correlation means that TOTs are adaptive, in a functional sense, as they alert us to correctible retrieval failures. However, they also provide us with a manner of understanding the relation of process, experience, and behavior [5, 31]. The first author has argued at length elsewhere for the reasons why it is necessary to consider that the cognitive processes that produce the phenomenological TOT are different from the processes that result in temporary inaccessibility (see [27, 31, 32, 35, 36]).

Applying the logic of the challenge to the Doctrine of Concordance, what should we expect to see when we examine TOTs? TOTs are subjective experiences, which monitor unretrieved target memories. The process by which TOTs are produced should, therefore, be separable from but related to the processes that actually engage in retrieval. Moreover, it should be possible to find multiple neural components associated with retrieval and the TOT experience. These processes should overlap, but it should also be possible to see some brain regions involved with and responsible for retrieval but not the TOT and vice versa. We will now examine the neuroscience literature with these ideas in mind.

2.5 The Neuroscience of TOTs

fMRI studies Two studies have directly applied the logic above to an fMRI analysis of the neural correlates of the TOT experience [24, 25]. The first study compared TOTs, correct responses, and don't know responses, and the second study did a similar analysis, but also included feeling-of-knowing judgments. Participants were presented with definitions of words or general-information questions (e.g., "Carmen composer") and were asked to retrieve the word that matched them (e.g., "Bizet"). Participants made one of three responses while being monitored by fMRI, indicating that they (1) recalled the answer, (2) did not know the answer, or (3) were in a TOT for the answer. Because the participants were in the scanner, these responses were made via finger presses. Follow-up questions showed a relatively low rate of commission errors in the "recalled" condition. Maril et al. [24, 25] compared brain activity across these three responses. Results from fMRI studies are often complex, but there were clearly areas of higher activity in TOTs than in either the "recalled" or "don't know" condition. These areas of the brain more activated during TOTs were mostly in the frontal cortex, including right inferior frontal and right medial frontal, right dorsolateral frontal, bilateral anterior frontal, and anterior cingulate cortices (also see [19]).

The prefrontal lobe neural regions are intriguing because they have been associated with metacognition in other studies (see Chua et al., this volume; Fleming and Dolan, this volume). These areas have been previously associated with a number of monitoring and supervisory functions, including executive control (see [33, 38]). Some of the areas above are associated with monitoring and control in other tasks. For example, dorsolateral prefrontal cortex is implicated in judgments of learning and feeling of knowing. The anterior cingulate cortex is associated with surprise monitoring across a number of domains from metacognition to emotional regulation [4]. Thus, the areas of the brain activated during TOTs support the idea that a TOT is a metacognitive signal, as they are functionally related to the processes that produce the phenomenological TOT as well as the processes that cause the temporary inaccessibility.

Other fMRI studies have directed analysis toward the processes by which words become temporarily inaccessible. These studies find that temporary inaccessibility is a function of other areas in the brain. For instance, Shafto et al. [37] used a celebrity-naming task, in which participants were asked to identify the name of celebrity when a photograph of the person's face was presented. This study was mainly interested in age differences and whether these correlated with changes in temporary inaccessibility. When focusing on temporary inaccessibility, these authors found a relation between the insula and phonological processing, and that the degree of atrophy of this region in older people could contribute to the age-related increase in temporary inaccessibility. Furthermore, similar to Maril et al. [24, 25] findings, they found higher activation in the anterior cingulate and inferior frontal cortex (among other areas) in the TOT condition than in the know condition, indicating that these regions are correlated with the experience of the TOT. Thus, the pattern that emerges from the fMRI data is that temporary inaccessibility seems to be correlated with processes related to language processing, associated with the insula, but that the feeling of temporary inaccessibility might be associated with the anterior cingulate and prefrontal areas.

Although the temporal resolution of fMRI is improving, fMRI is still too slow to capture the rapid changes that occur in neural processes as active cognition unfolds over time. In order to look at rapid changes in the brain over time, it is still necessary to employ electroencephalography (EEG) or magnetoencephalography (MEG) techniques that allow study of the direct electromagnetic activity of neuron populations with a temporal resolution in the order of milliseconds. EEG and MEG can not only isolate areas of the brain (although with a spatial resolution lesser than fMRI), but can also observe changes over time with an optimal temporal resolution. Thus, the EEG and MEG data provide an excellent way to evaluate the model of how retrieval and metacognition interact.

EEG and MEG studies of TOTs Díaz and his colleagues have extensively studied both the retrieval process and the TOT using EEG and MEG techniques. We review this work in this section. In an initial study, Díaz et al. [10] examined face naming and TOTs for unrecalled names while the EEG was monitoring participants. In the task, participants were presented with the face of a famous person and required to press a button to indicate whether they were sure that they knew the person's name. They were also required to name the person. If they felt they knew the name but could not recall it, they were asked to indicate a TOT. If a person indicated a TOT, they were given a phonological cue and an opportunity to retrieve the name again (i.e., the same face was presented again). In this study, Díaz et al. were able to compare EEG patterns for successful retrieval (Know), unsuccessful retrieval (Don't Know), and unsuccessful retrieval accompanied by TOTs (see Fig. 2.1 for the general design of these studies).

The logic of the procedure was to look for systematic event-related potential (ERP) differences between task categories time-locked to the onset of the stimulus. Thus, a face of the basketball player Pau Gasol is presented at time 0. Then the EEG can register changes in the brain-wave patterns locked from the stimulus onset. Using this technique, Díaz et al. [10] were able to look at differences in the

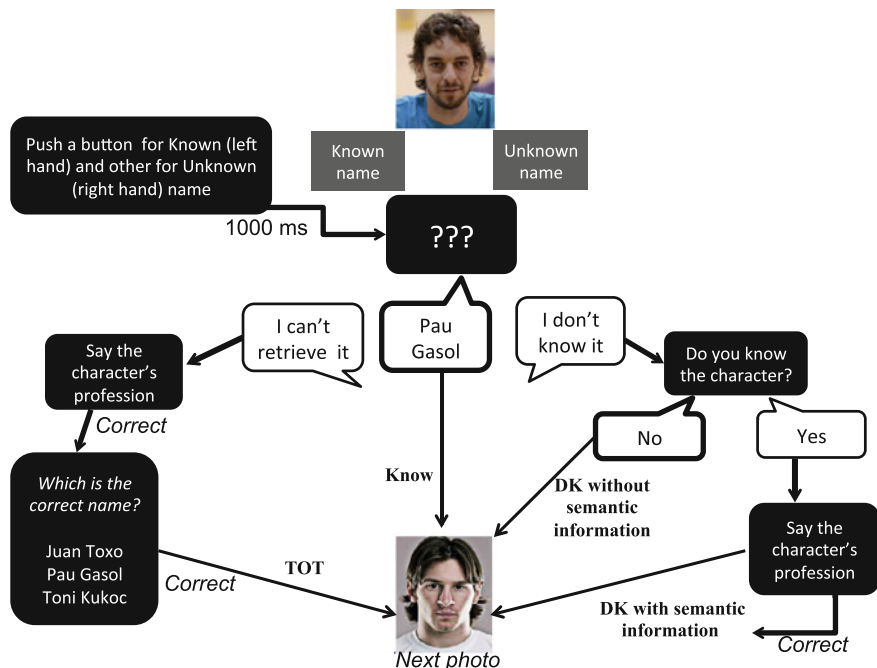


Fig. 2.1 Procedure used in Buján et al. [7], but indicative of the general procedure in these experiments. After pressing the corresponding *left* or *right* button, the participant has to say aloud **a** the correct name (Know), after which the next photograph appeared, **b** “I can’t retrieve it” after which a series of questions appeared (TOT) and **c** “I don’t know it” after which a series of questions appeared (Don’t Know)

event-related potential between items that were correctly remembered (Know) and those that induced TOTs. Interestingly, the data from this study showed differences in the patterns between these two internal states.

During the initial presentation of the faces, there were no differences between the Know and the TOT categories in the ERP components for the first 550 ms (milliseconds) after presentation [10]. That is, the ERP correlates of perceptual processing, face recognition, and access to semantic and lexical information, as indicated by fame domain and name recall, did not differ between Know and TOT categories. This is consistent with a model in which retrieval is initiated, and the person engages in a recall attempt. Thus, for the first approximately half-second after presentation, we cannot yet distinguish retrieval and metacognitive processes. However, between 550 and 750 ms after presentation, a wave known as the late P3, which is associated with the response categorization (Know, TOT or Don’t Know) was significantly larger for retrieved items than for TOT items. This result was attributed to a division of processing resources in TOTs between the categorization of the stimulus and the continuous search for complete phonological information about the name.

A second EEG component that differed among output condition was the late negative wave occurring about 1,300 ms after presentation of the face stimulus. In the Know and Don't Know response categories, late negative waves appeared after the motor response, whereas in TOTs, it appeared after the stimulus classification but before the motor response. Late negative waves were largest in Don't Know items, intermediate in Know items, and smallest in TOT items, perhaps associated with the level of uncertainty about the categorization of the stimulus and with the release of processing resources with the response. The later dissociation among retrieved, unretrieved, and TOT items supports the idea that processes diverge once the recall attempt has failed during TOTs. In a subsequent study [8], in which the ERPs were averaged in relation to a manual response, it was shown that the preparation and the execution of the responses (manual+verbal) differently modulated the stimulus-related ERP components in each response category, explaining in part the differences between categories in the amplitude of the late negative wave. Galdo-Álvarez et al. [15, 16] replicated the Díaz et al. [10] findings and found no differences in time course of ERP between older and younger participants for the Know and TOT responses, although there were age-related differences in ERP amplitudes and their scalp distribution.

In a subsequent replication, Lindín and Díaz [21] replaced the manual plus verbal response (for Know, Don't Know, TOT responses) with a manual response that was separated 1 s from the verbal response (three question marks were presented authorizing the participant to perform the corresponding verbal response). Using this methodology, there was a longer latency of the N450 wave for TOTs. This means that for TOTs the N450 wave occurred slightly later than it did for Know and Don't Know responses, probably indicating the slowing in the retrieval of semantic and lexical-semantic information during a TOT. Again, the late P3 distinguished TOTs and retrieved items. However, in this study, there were no differences at the late 1,300 ms stage between TOTs and other states. Though the form of response brought out one feature and suppressed another, we still see that TOTs and successful recall are dissociable by their EEG patterns.

Lindín et al. [22] used a similar behavioral methodology but with MEG technology in addition to EEG. The use of MEG technology allowed the researchers to pinpoint more accurately the spatial correlates of the behavioral measures with the same temporal resolution provided by the EEG. The goal in this study was to characterize the spatiotemporal course of brain activation in both successful recall and during the TOT. Consistent with the earlier findings, there were no differences in the MEG data for the first 210 ms after presentation of a face. However, during the interval from 210 to 520 ms, there was greater activation for Know responses than for TOT responses in a variety of brain regions, mostly in the left hemisphere, including left anterior medial prefrontal cortex, left orbitofrontal cortex, the left superior temporal pole, and the left inferior, middle and superior temporal gyri, as well as bilateral parahippocampal gyrus, right fusiform gyrus, and Broca's area. These are consistent with the processes involved in successfully retrieving a stored memory. They also found that at the later interval, 580–820 ms, there was greater activation for TOTs in the bilateral inferior and middle occipital gyri as well as left

temporal and right frontal and parietal regions, consistent with the role of monitoring in TOT experiences (because of the right frontal activity) and with the active but fruitless search of the name.

In sum, whereas the ERP data showed that the amplitude differences between Know responses and TOTs were observed between 550 and 750 ms post-stimulus, coinciding with the categorization of stimulus, the MEG data showed that the differences are already apparent from 210 ms and, consistently, from 310 ms. MEG data also showed that there are two distinct phases: the first, between 310 and 510 ms corresponding with the successful access to the phonology of the name (greater activation of a brain network related with name recall in the case of Know responses) and with the genesis of the TOT (hypoactivation of the network in the TOT responses); and the second between 580 and 820 ms, with greater activation for TOT than for Know. The 310–510 activity may correspond to the active search of information in memory about the name. The 580–820 activation may correspond to the metacognitive monitoring in TOT experiences because this activity may be responsible for partial retrieval, such as the retrieval of visual information and partial lexical information [10].

With the aim of determining the timing of the phonological retrieval, Buján et al. [7] carried out another ERP study using a face-naming task. In this study, the early components of the ERP were again equivalent in TOTs and Know responses. However, there were differences among response options later in processing. Again, after 550 ms, there were smaller positive amplitudes in the TOTs than in Know responses, and after 1,100 ms there were higher negative amplitudes in the TOT than in Know and Don't know responses. This may correspond to the metacognitive control of retrieval, as resources may be diverted to conflict management and a continued search for the missing word, consistent with a metacognitive component to TOTs (see Fig. 2.2).

Probably the most interesting result of Buján et al.'s [7] study is the difference between the Know category and the TOT in the Lateralized Readiness Potential (LRP). The onset latency of the stimulus-related LRP was 360 ms for the Know category, whereas the TOT (and also the Don't Know category) showed a significant delay. The onset latency of the stimulus-related LRP is a correlate of response selection; consequently, when the response selection starts, access to the phonological output lexicon (lexemes) has already taken place. The response-related LRP also showed earlier onset latency in Know than in TOT, that is, the start of the actual preparation of the response was slower for TOT than for Know responses. These LRP data are consistent with the MEG data and indicate that around 360 ms, the phonological information of the name was retrieved in Know (which allow the corresponding selection and preparation of the response), but in TOTs the delay in the selection and in the preparation of the response indicated the failure in retrieving the complete phonology of the name.

To summarize, it has become clear that the TOT can be thought of as both a problem with retrieval and as successful metacognition (see [36]). The retrieval failure occurs because of any number of problems with the retrieval process, whereas successful metacognition monitors that failure. We argue here that the

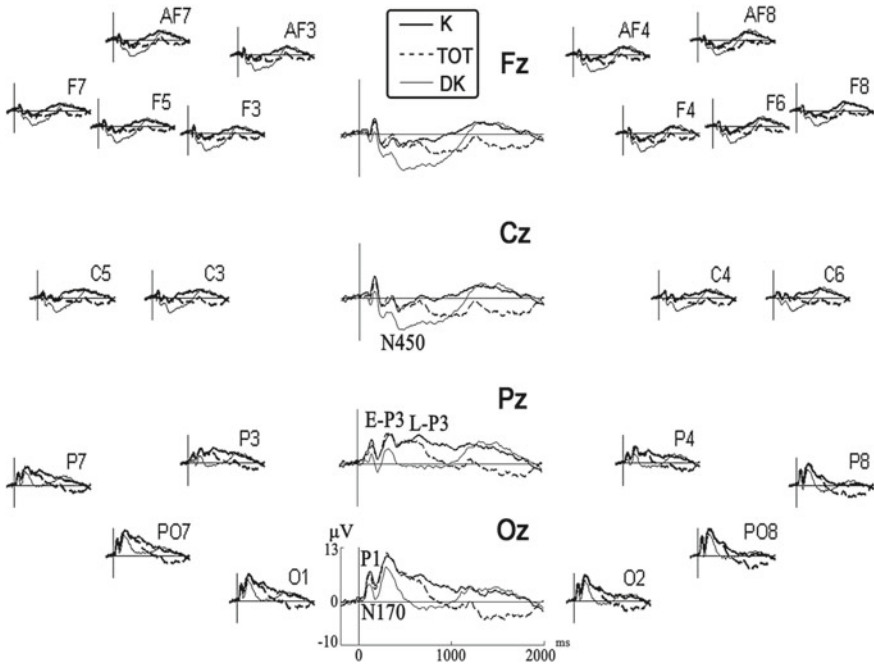


Fig. 2.2 Grand-averaged ERP waveforms for the Know (K, *thick line*), TOT (*dashed line*) and Don't Know (DK, *thin line*) response categories at *midline* electrodes and at several of the lateralized electrodes used

neural data presented here support this view. Both the fMRI, EEG and MEG data find that retrieval failure and TOT experience map onto different areas of the brain and at different times following the presentation of a stimuli. That is, there is one set of processes that are responsible for retrieval. When these processes fail, a separate set of processes monitor the retrieval failure to determine if the target can eventually be recovered. When these processes indicate a possibility of this, the TOT is experienced. Thus, the TOT data support a model in which different underlying processes are responsible for the cognition and the metacognition that monitors it.

2.6 Are Metacognitive Judgments Process-Pure?

We now return to the question raised earlier in the chapter. Do metacognitive judgments reflect multiple underlying metacognitive processes? We think our discussion of TOTs suggests that even straightforward metacognitive judgments like TOTs may involve numerous factors that influence its occurrence. We present here our reasoning that this conclusion likely extends to other metacognitive judgments.

Consider a neuroimaging study that examines the neural correlates of a particular metacognitive judgment, say a judgment of learning (JOL). A participant must determine if a particular item, say the Icelandic-English translation pair, “fiðlu—violin” has been learned. The participant must indicate his or her JOL on a scale of 0–100, translating the experience of the JOL into an outputted number on the scale. Multiple processes are surely at work as one does this task. Is the Icelandic word familiar, easy to pronounce, studied before, and does it resemble an English word? How long in the future will be the memory test, will it be a cued recall or a forced-choice recognition test, how many other pairs might also be required for the test? Can an easy linkword (e.g., “fiddle”) be generated to help encode the pair? We contend that there are many processes used to determine a single number given in the JOL. When we look at the pattern of activity in the brain, it is difficult to correlate any activated area of the brain with one and only one of these potential sources of information for the JOL.

When one looks at the neuroscience literature, one finds that JOLs are correlated with different regions in the brain in different studies. For example, Kao et al. [18] found that JOLs were associated with activity in left ventromedial prefrontal cortex. However, Do Lam et al. [12] found that JOLs were associated with activation in medial PFC, orbital frontal, and anterior cingulate cortices. Thus, ventromedial cortex appears common to both studies, but additional areas were activated during JOLs in one study that were not in the other. Why does one judgment correlate with such different brain regions across studies? It likely has to do with the procedural differences between the two studies. Kao et al. had participants make JOLs on photographs of visual scenes (e.g., a mountain sunset) for eventual recognition whereas Do Lam et al. had participants make JOLs on photographs of faces for eventual cued recall of names. Thus, because of the different tasks, the JOLs were based on different sources of information and thus different areas of the brain were recruited. Because the processes underlying JOLs are sensitive to differences in tasks, the JOL is not a process-pure measure of metacognition.

Feeling-of-knowing judgments (FOKs) are generally defined as a feeling that an unrecalled item will eventually be recognized [32]. We also find that different studies find different regions of the brain are associated with FOK. For example, although Maril et al. [23] and Jing et al. [17] found that FOK was uniquely associated with activity in left inferior prefrontal cortex, Schnyer et al. [29] found that FOK was uniquely associated with ventromedial PFC. Moreover, Reggev et al. [28] found different areas of the prefrontal cortex were uniquely associated with FOK for episodic and semantic memory. The point is not these studies should all be the same, but that small changes in procedure elicit different processes that draw on different areas of the brain. Thus, the FOK task is not a process-pure measure of metacognition either.

This is not a pessimistic argument. We are not arguing against the use of neuroimaging or against the use of particular tasks to investigate metacognition or the association between brain regions and particular judgments. Our point is simple: metacognitive judgments draw on multiple underlying cognitive processes

that likely draw on different underlying brain processes. Small changes in procedure can therefore shift the relevance of different processes for the same judgment, leading to the pattern of data observed in JOLs and FOKs, in which small shifts in procedure yield different unique activity associated with a particular area. So to repeat, we must be cautious because metacognitive judgments are not necessarily process-pure.

2.7 Future Directions

We welcome and embrace metacognitive neuroscience. Already neuroscience research has contributed to our understanding of the underlying cognitive mechanisms involved in some metacognition paradigms (e.g., [24]). And we certainly agree that it is important to know which brain regions correlate with which cognitive processes. Our point here is a simple one—that any metacognitive judgment may map onto multiple cognitive processes and each of these processes may correlate with different neural networks. This leads to the conclusion that any attempt to map the neurocognition of metacognition will require us to start disentangling the processes that underlie a particular judgment. Thus, TOTs may be determined by cue familiarity, partial retrieval, and perhaps the fluency of the broken retrieval process. The TOT experience is an amalgamation of these different underlying processes. In looking at the brain, one must try to look for the regions and time course of activity associated with these different neural elements and networks (see [11]).

We close by returning to the issue of Tulving’s challenge to the notion of the doctrine of concordance. Tulving challenged the view that experience, behavior, and cognitive process were always perfectly correlated. Schwartz [31] modified this view to challenge the view that metacognitive experiences are based on the processes they are supposed to monitor. Indeed, much research now suggests that metacognitive judgments are largely a set of heuristics we use to infer what our cognitive processes are doing [20, 26]. We think it is clear from the arguments and data advanced here that it is not tenable to speak of metacognitive experience as being identical to the cognitive processes these experiences track. We think our data show that the TOT experience, for example, may be partially but not completely based on the processes that lead to retrieval failure and partially based on heuristics such as the amount of related information, with more retrieved related information leading to a greater likelihood of a TOT [35]. Thus, as neuroscience explores the nature of mental experience and metacognition, researchers must bear in mind the importance of distinguishing object-level processes from meta-level processes. We look forward to seeing continued work linking metacognition to brain function.

References

1. Bacon E, Schwartz BL, Paire-Ficout L, Izaute M (2007) Dissociation between the cognitive process and the phenomenological experience of the TOT: effect of the anxiolytic drug lorazepam on TOT states. *Cogn Conscious* 16:360–373
2. Benjamin AS (2005) Response speeding mediates the contribution of cue familiarity and target retrievability to metamnemonic judgments. *Psychon Bull Rev* 12:874–879
3. Benjamin AS, Bjork RA, Schwartz BL (1998) The mismeasure of memory: when retrieval fluency is misleading as a metamnemonic index. *J Exp Psychol Gen* 127:55–68
4. Botvinick M (2007) Conflict monitoring and decision making: reconciling two perspectives on anterior cingulate function. *Cogn Affect Behav Neurosci* 7:356–366
5. Brown AS (1991) A review of the tip-of-the-tongue experience. *Psychol Bull* 109:204–223. doi:[10.1037/0033-2909.109.2.204](https://doi.org/10.1037/0033-2909.109.2.204)
6. Brown AS (2012) *Tip of the tongue state*. Psychology Press, New York, Jan 2011
7. Buján A, Galdo-Álvarez S, Lindín M, Díaz F (2012) An event-related potentials study of face naming: evidence of phonological retrieval deficit in the tip-of-the-tongue state. *Psychophysiology* 49:980–990. doi:[10.1111/j.1469-8986.2012.01374.x](https://doi.org/10.1111/j.1469-8986.2012.01374.x)
8. Buján A, Lindín M, Díaz F (2009) Movement related cortical potentials in a face naming task: influence of the tip-of-the-tongue state. *Int J Psychophysiol* 72:235–245. doi:[10.1016/j.ijpsycho.2008.12.012](https://doi.org/10.1016/j.ijpsycho.2008.12.012)
9. Cleary AM, Brown AS, Sawyer BD, Nomi JS, Ajoku AC, Ryals AJ (2012) Familiarity from the configuration of objects in 3-dimensional space and its relation to déjà vu: a virtual reality investigation. *Conscious Cogn* 21:969–975
10. Díaz F, Lindín M, Galdo-Álvarez S, Facal D, Juncos-Rabadán O (2007) An event-related potentials study of face identification and naming: the tip-of-the tongue state. *Psychophysiology* 44(1):50–68. doi:[10.1111/j.1469-8986.2006.00483.x](https://doi.org/10.1111/j.1469-8986.2006.00483.x)
11. Díaz F, Lindín M, Galdo-Álvarez S, Buján A (in press) Neurofunctional correlates of the tip-of-the-tongue state. In: Schwartz BL, Brown AS (eds) *The tip-of-the-tongue and related phenomena*. Cambridge University Press, Cambridge
12. Do Lam ATA, Axmacher N, Fell J, Staresina BP, Gauggel S et al (2012) Monitoring the mind: the neurocognitive correlates of metamemory. *PLoS one* 7(1):e30009. doi:[10.1371/journal.pone.0030009](https://doi.org/10.1371/journal.pone.0030009)
13. Dunlosky J, Metcalfe J (2009) *Metacognition*. Sage Publications Inc, Thousand Oaks
14. Fechner GT (1860/1966) *Elements of psychophysics*. Holt, Rinehart, and Winston, New York
15. Galdo-Álvarez S, Lindín M, Díaz F (2009) Age-related prefrontal over-recruitment in semantic memory retrieval: evidence from successful face naming and the tip-of-the-tongue state. *Biol Psychol* 82(1):89–96. doi:[10.1016/j.biopsycho.2009.06.003](https://doi.org/10.1016/j.biopsycho.2009.06.003)
16. Galdo-Álvarez S, Lindín M, Díaz F (2009) The effect of age on event-related potentials (ERP) associated with face naming and with the tip-of-the-tongue (TOT) state. *Biol Psychol* 81:14–23. doi:[10.1016/j.biopsycho.2009.01.002](https://doi.org/10.1016/j.biopsycho.2009.01.002)
17. Jing L, Niki K, Xiaoping Y, Yue-jia L (2004) Knowing that you know and knowing that you don't know: a fMRI study on feeling-of-knowing (FOK). *Acta Psychol Sin* 36:426–433
18. Kao Y-C, Davis ES, Gabrieli JDE (2005) Neural correlates of actual and predicted memory formation. *Nat Neurosci* 8:1776–1783
19. Kikyo H, Ohki K, Sekihara K (2001) Temporal characterization of memory retrieval processes: an fMRI study of the 'tip of the tongue' phenomenon. *Eur J Neurosci* 14(5):887–892. doi:[10.1046/j.0953-816x.2001.01711.x](https://doi.org/10.1046/j.0953-816x.2001.01711.x)
20. Koriat A (1993) How do we know that we know? The accessibility account of the feeling of knowing. *Psychol Rev* 100:609–639
21. Lindín M, Díaz F (2010) Event-related potentials in face naming and tip-of-the-tongue state: further results. *Int J Psychophysiol* 77(1):53–58. doi:[10.1016/j.ijpsycho.2010.04.002](https://doi.org/10.1016/j.ijpsycho.2010.04.002)
22. Lindín M, Díaz F, Capilla A, Ortiz T, Maestú F (2010) On the characterization of the spatio-temporal profiles of brain activity associated with face naming and the tip-of-the-tongue

- state: a magnetoencephalographic (MEG) study. *Neuropsychologia* 48(6):1757–1766. doi:[10.1016/j.neuropsychologia.2010.02.025](https://doi.org/10.1016/j.neuropsychologia.2010.02.025)
23. Maril A, Simon JS, Mitchell JP, Schwartz BL, Schacter DL (2003) Feeling-of-knowing in episodic memory: an event-related fMRI study. *NeuroImage* 18:827–836
 24. Maril A, Simons JS, Weaver JJ, Schacter DL (2005) Graded recall success: an event-related fMRI comparison of tip of the tongue and feeling of knowing. *NeuroImage* 24:1130–1138
 25. Maril A, Wagner AD, Schacter DL (2001) On the tip of the tongue: an event-related fMRI study of semantic retrieval failure and cognitive conflict. *Neuron* 31:653–660
 26. Metcalfe J (1993) Novelty monitoring, metacognition, and control in a composite holographic associative recall model: interpretations for Korsakoff amnesia. *Psychol Rev* 100:3–22
 27. Metcalfe J, Schwartz BL, Joaquim SG (1993) The cue familiarity heuristic in metacognition. *J Exp Psychol Learn Mem Cogn* 19:851–861. doi:[10.1037/0278-7393.19.4.851](https://doi.org/10.1037/0278-7393.19.4.851)
 28. Reggev N, Zuckerman M, Maril A (2011) Are all judgments created equal? An fMRI study of semantic and episodic metamnemonic predictions. *Neuropsychologia* 49:1332–1343
 29. Schnyer DM, Nicholls L, Verfaellie M (2005) The role of VMPC in metamemorial judgments of content retrievability. *J Cogn Neurosci* 17:832–846
 30. Schwartz BL (1998) Illusory tip-of-the-tongue states. *Memory* 6:623–642
 31. Schwartz BL (1999) Sparkling at the end of the tongue: the etiology of tip-of-the-tongue phenomenology. *Psychon Bull Rev* 6:379–393
 32. Schwartz BL (2006) Tip-of-the-tongue states as metacognition. *Metacogn Learn* 1:149–158
 33. Schwartz BL, Bacon E (2008) Metacogn neurosci. In: Dunlosky J, Bjork R (eds) *Handbook of metamemory and memory*. Psychology Press, New York, pp 355–371
 34. Schwartz BL, Metcalfe J (1992) Cue familiarity but not target retrievability enhances feeling-of-knowing judgments. *J Exp Psychol Learn Mem Cogn* 18:1074–1083
 35. Schwartz BL, Metcalfe J (2011) Tip-of-the-tongue (TOT) states: retrieval, behavior, and experience. *Mem Cogn* 39(5):737–749. doi:[10.3758/s13421-010-0066-8](https://doi.org/10.3758/s13421-010-0066-8)
 36. Schwartz BL, Metcalfe J (in press) Tip-of-the-tongue (TOT) states: mechanisms and metacognitive control. In: Schwartz BL, Brown AS (eds) *The tip-of-the-tongue and related phenomena*. Cambridge University Press, Cambridge
 37. Shafto M, Stamatakis E, Tam P, Tyler L (2010) Word retrieval failures in old age: the relationship between structure and function. *J Cogn Neurosci* 22(7):1530–1540. doi:[10.1162/jocn.2009.21321](https://doi.org/10.1162/jocn.2009.21321)
 38. Shimamura AP (2008) A neurocognitive approach to metacognitive monitoring and control. In: Dunlosky J, Bjork RA (eds) *Handbook of memory and metamemory: essays in honor of Thomas O. Nelson*. Psychology Press, New York, pp 373–390
 39. Thomas AK, Bulevich JB, Dubois S (2011) The role of contextual information in episodic feeling of knowing. *J Exp Psychol Learn Mem Cogn* 38:96–108
 40. Tulving E (1989) Memory: performance, knowledge, and experience. *Eur J Cogn Psychol* 1:3–26

Chapter 3

Signal Detection Theory Analysis of Type 1 and Type 2 Data: Meta- d' , Response-Specific Meta- d' , and the Unequal Variance SDT Model

Brian Maniscalco and Hakwan Lau

Abstract Previously we have proposed a signal detection theory (SDT) methodology for measuring metacognitive sensitivity (Maniscalco and Lau, *Conscious Cogn* 21:422–430, 2012). Our SDT measure, meta- d' , provides a response-bias free measure of how well confidence ratings track task accuracy. Here we provide an overview of standard SDT and an extended formal treatment of meta- d' . However, whereas meta- d' characterizes an observer's sensitivity in tracking overall accuracy, it may sometimes be of interest to assess metacognition for a particular kind of behavioral response. For instance, in a perceptual detection task, we may wish to characterize metacognition separately for reports of stimulus presence and absence. Here we discuss the methodology for computing such a “response-specific” meta- d' and provide corresponding Matlab code. This approach potentially offers an alternative explanation for data that are typically taken to support the unequal variance SDT (UV-SDT) model. We demonstrate that simulated data generated from UV-SDT can be well fit by an equal variance SDT model positing different metacognitive ability for each kind of behavioral response, and likewise that data generated by the latter model can be captured by UV-SDT. This ambiguity entails that caution is needed in interpreting the processes underlying relative operating characteristic (ROC) curve properties. Type 1 ROC curves generated by combining type 1 and type 2 judgments, traditionally interpreted in

B. Maniscalco (✉)

National Institute of Neurological Disorders and Stroke, National Institutes of Health,
10 Center Drive, Building 10, Room B1D728, MSC 1065, Bethesda, MD 20892-1065, USA
e-mail: bmaniscalco@gmail.com

B. Maniscalco · H. Lau

Department of Psychology, Columbia University, 406 Schermerhorn Hall,
1190 Amsterdam Avenue MC 5501, New York, NY 10027, USA
e-mail: hakwan@gmail.com

H. Lau

Department of Psychology, UCLA, 1285 Franz Hall, Box 951563 Los Angeles,
CA 90095-1563, USA