# Studies in Applied Philosophy, Epistemology and Rational Ethics

Volume 63

**Studies in Applied Philosophy, Epistemology and Rational Ethics (SAPERE)** publishes new developments and advances in all the fields of philosophy, epistemology, and ethics, bringing them together with a cluster of scientific disciplines and technological outcomes: ranging from computer science to life sciences, from economics, law, and education to engineering, logic, and mathematics, from medicine to physics, human sciences, and politics. The series aims at covering all the challenging philosophical and ethical themes of contemporary society, making them appropriately applicable to contemporary theoretical and practical problems, impasses, controversies, and conflicts. Our scientific and technological era has offered "new" topics to all areas of philosophy and ethics– for instance concerning scientific rationality, creativity, human and artificial intelligence, social and folk epistemology, ordinary reasoning, cognitive niches and cultural evolution, ecological crisis, ecologically situated rationality, consciousness, freedom and responsibility, human identity and uniqueness, cooperation, altruism, intersubjectivity and empathy, spirituality, violence. The impact of such topics has been mainly undermined by contemporary cultural settings, whereas they should increase the demand of interdisciplinary applied knowledge and fresh and original understanding. In turn, traditional philosophical and ethical themes have been profoundly affected and transformed as well: they should be further examined as embedded and applied within their scientific and technological environments so to update their received and often old-fashioned disciplinary treatment and appeal. Applying philosophy individuates therefore a new research commitment for the 21st century, focused on the main problems of recent methodological, logical, epistemological, and cognitive aspects of modeling activities employed both in intellectual and scientific discovery, and in technological innovation, including the computational tools intertwined with such practices, to understand them in a wide and integrated perspective. **Studies in Applied Philosophy, Epistemology and Rational Ethics** means to demonstrate the contemporary practical relevance of this novel philosophical approach and thus to provide a home for monographs, lecture notes, selected contributions from specialized conferences and workshops as well as selected Ph.D. theses. The series welcomes contributions from philosophers as well as from scientists, engineers, and intellectuals interested in showing how applying philosophy can increase knowledge about our current world. Initial proposals can be sent to the Editor-in-Chief, Prof. Lorenzo Magnani, lmagnani@unipv.it:

- A short synopsis of the work or the introduction chapter
- The proposed Table of Contents
- The CV of the lead author(s).

For more information, please contact the Editor-in-Chief at lmagnani@unipv.it.

Indexed by SCOPUS, zbMATH, SCImago, DBLP.

All books published in the series are submitted for consideration in Web of Science.

Vincent C. Müller

Editor

# Philosophy and Theory of Artificial Intelligence 2021

*Editor*
Vincent C. Müller ⓘD
Friedrich-Alexander University
of Erlangen-Nuremberg (FAU)
Erlangen, Germany

*Dedicated to the memory of our colleague
and friend
Ivica Crnkovic
(1955–2022)*

# Advisory Editors

# Preface

The papers in this volume result from the 4th conference on the 'Philosophy and Theory of Artificial Intelligence' (PT–AI) that we organised in Gothenburg on 27–28 September, 2021 (see http://www.pt-ai.org/). The local organisation was taken care of by Profs. Ivica Crnkovic and Gordana Dodig-Crnkovic, with support from Chalmers University in Gothenburg.

This conference had been planned for a much earlier date, but several things got in the way, most of all the COVID pandemic. Eventually, we decided that in some way or other we should run PT-AI in 2021, and for most people, this was the first conference where actual physical presence was possible again, though with significant safety measures. It was the first PT-AI conference that was run in a hybrid online/onsite fashion (earlier conferences had already allowed online live listening to keynote talks). As a result of all this, the conference was a bit smaller than usual, with 60 submissions, of which 25 were presented at the meeting. My thanks to the colleagues on the programme committee who worked hard on the double-blind reviewing and assured a very high academic level! The inspiring invited speakers were Virginia Dignum (Umeå, Sweden), Michael Levin (Tufts, USA), David Papineau (KCL, UK), and Shannon Vallor (Edinburgh, UK).

It was very good to see many new faces coming to the field, as well as some established philosophers from other areas showing an interest in the philosophy of AI—which is clearly moving into the mainstream now. The new people and the influences from many different directions are clearly enriching the field and I expect this to continue.

I have to end on a bitter note: In February 2022, our co-organiser Ivica died suddenly. It is extremely sad to see such a rich life cut short, and many people left behind with a huge gap in their lives—first of all, his wife and his children. At the same time, I am also grateful that I had the chance to know Ivica and learn from him, academically and as a human being.

Eindhoven, The Netherlands
April 2022

Vincent C. Müller
http://www.sophia.de/

# Contents