

To the Editors:

We submit our manuscript, "A Way Around the Exploration-Exploitation Dilemma", for consideration at *Science*.

Since Schultz et al's (*J Neurophys*, 1996; > 4000 citations) seminal report that dopamine neurons behave quantitatively like a reinforcement learning signal, reinforcement learning has been central to the study of learning as well as important to our conception of numerous neurological and psychiatric disorders, including autism and schizophrenia. Reinforcement learning is also undergoing a renaissance in artificial intelligence research, achieving superhuman performance on Atari video games (Minh et al, *Nature* 2015; >6000 citations), chess (Silver et al, *Science* 2018; >400 citations), and Go (Silver et al, *Nature* 2016; >5000 citations). However, despite this progress all of reinforcement learning shares a common dilemma.

The exploration-exploitation dilemma is summarized by a simple question: "Should I exploit an available reward, or explore to try out a new uncertain action?" Unfortunately, it's been proven that this dilemma, when stated as a mathematical problem, is intractable and so can't be solved directly. This fundamentally limits our ability to predict optimal naturalistic behavior during foraging and exploration, and to optimally drive learning in artificial agents.

To overcome this field-wide limitation, we took a fresh look at the dilemma. Our goal was simple: when one mathematical problem can't be solved, it's often good to find another related problem that can be and use that to make progress on both.

We show, for the first time, that nearly any dilemma problem can also be viewed as competition, between exploiting known rewards or exploring to learn a *world model*--a simplified concept of memory borrowed from computer science. We prove this competition can be perfectly solved using the simplest of all value learning algorithms: a deterministic greedy policy. To ensure this solution is as broad as possible we also derived a new universal theory of information value which complements--but is independent of--Shannon's Information Theory.

To a popular audience, we have to emphasize the algorithm we discovered translates both colloquially and exactly to, "*What do you value more right now, getting a reward or learning something new? Pick the biggest*". This means our analysis can have broad application both to scientific research--in fields as diverse as psychology, neuroscience, biology, data science, artificial intelligence, game theory, economics, demography, and political science--and directly to everyday life.

With our *Report* submission we include the text of the primary manuscript with 3 figures as well as a full mathematical appendix, supplementary discussion, and methods section. The main text contains about 2600 words and with figures and figure legends will be approximately 3 pages in the printed publication. This work is not under consideration for publication elsewhere. Code and data are available at <https://github.com/CoAxLab/infomercial>.

Erik J Peterson, Ph.D.
Research Fellow

Tim Verstynen, Ph.D.
Associate Professor

Carnegie Mellon University, Pittsburgh

Carnegie Mellon University, Pittsburgh