

The Dreams of Theory

James P. Crutchfield

SFI WORKING PAPER: 2010-12-033

SFI Working Papers contain accounts of scientific work of the author(s) and do not necessarily represent the views of the Santa Fe Institute. We accept papers intended for publication in peer-reviewed journals or proceedings volumes, but not papers that have already appeared in print. Except for papers by our external faculty, papers must be based on work done at SFI, inspired by an invited visit to or collaboration at SFI, or funded by an SFI grant.

©NOTICE: This working paper is included by permission of the contributing author(s) as a means to ensure timely distribution of the scholarly and technical work on a non-commercial basis. Copyright and all rights therein are maintained by the author(s). It is understood that all persons copying this information will adhere to the terms and constraints invoked by each author's copyright. These works may be reposted only with the explicit permission of the copyright holder.

www.santafe.edu



SANTA FE INSTITUTE

OPINION

The Dreams of Theory

Computing power and sophisticated data acquisition mask the fact that, in many sciences and engineering, the balance between theory and experiment is getting increasingly out of whack, says **Jim Crutchfield**.

The twenty-one year old Werner Heisenberg, already a rising star in quantum theory, was mortified when his doctoral exam committee awarded a nearly failing grade. He had passed, but excused himself early from the celebration party that evening put on by his advisor Arnold Sommerfeld. Boarding the midnight train, he abruptly left Munich to take up a previously arranged job in Göttingen, humiliated by his mediocre exam and concerned that his boss, Max Born, would no longer have him.

Responding to criticism that Heisenberg's brilliance in theory was eclipsing a well-rounded appreciation of physics, Sommerfeld had required him to take Wilhelm Wien's course in experimental physics and focus his dissertation on turbulence in fluid flows. Both Wien and Sommerfeld were on his oral exam committee: Sommerfeld gave Heisenberg the highest grade, Wien failed him. In the early days of the 20th century, physics was experimental physics. Heisenberg was, plainly to Wien, no experimentalist [1]. The rest, they say, is history. Heisenberg's towering contributions to modern physics are unassailable and, after a time, led to many key experimental discoveries. Heisenberg stands as a pre-eminent example of the card-carrying theorist—a profession new to 20th century science.

Today, looking more broadly across all of science, theorists still have reason to worry. Advances in computing power combined with increasingly sophisticated data acquisition technologies have led some to claim that theory is obsolete [2,3].

Computing technology suffuses our lives and has unalterably changed scientific practice. And, there is no hint of a let up. The power available to the natural sciences and engineering have grown substantially, via unanticipated innovations. Cobbling together outlet-store PCs in large clusters has brought supercomputing-level performance to even small research groups and colleges. At the high-end, focusing on environmental concerns rather than following the US's focus on compliance with

the nuclear Test Ban Treaty, Japan's Earth Simulator was the fastest machine on the planet in 2002. Designed especially to accelerate climate studies and more generally high-resolution modeling and simulation in geophysics, it was eclipsed just two years later by IBM's Blue Gene/L machine, doubling its performance. At that time, jumping by another factor of 5 in power to petascale computing was forecast to occur in 2010. We reached the petascale in 2008. We search for extraterrestrial life and fold proteins at home. We regularly attend block-buster, feature-length, all-digitally computed films; a commercial success the computer graphics community of the 1980s thought would be impossible [4]. We expect split-second responses to searches of billions of documents. We could very well have more computing than we know how to productively use. That said, we now know there are problems of daunting complexity that must be tackled and understood. Computing power is most likely not the bottleneck to the required innovation; indeed, its an essential driver.

A parallel acceleration has occurred in extracting data from natural and engineered systems. The development of tetrodes in neurobiology that simultaneously record dozens of neural spike trains, sensor networks that monitor fault segments for immanent earthquakes, digital sky surveys that produce three-dimensional maps of billions of stars, atomic-scale microscopes that show the structure of nanotubes, and scanners that record the dynamics of the brain are just a few notable successes in the advance of measurement technology. In this setting, the familiar story of gene sequencing is barely worth highlighting these days. Now, its only one among many in the new era of the Data Deluge. By any measure, empiricism has come to dominate science.

This concerns me. Data, whether produced by massive simulations or automated experimentation, is not understanding. Wrestling knowledge from data is theory's role. No surprise, the technological wizardry of the Data

Deluge has its seductions. My concern is that we, happy victims of technological success, are on the verge of permanently confusing the trees for the forest.

Theory's most important value lies in something we were taught in school. Scientists, obviously enough, have constructed models for centuries. While guess and insight are key, model building is also not a random process. We all know, for example, that in building a model parsimony is helpful. This is captured by Occam's Razor: Out of the range of all models consistent with observations, pick the smallest. Often interpreted as a convenience, we now know that parsimonious models are more than this, they capture a key aspect of understanding: Compact models reveal the mechanisms that produce a system's behavior and how those mechanisms interact [5,6]. One important methodological conclusion is that theory is not mere data compression [7]. There is meaning in our good models. They give insight, which is substrate for innovation.

We also should not forget the pragmatic aspects of theory's healthy critical role in science—aspects that complement this new view of theory building. First, predictive theories allow for a level of conceptual hygiene. Intuitions, hunches, and hypotheses that can be articulated in a mathematical theory can be tested for consistency. Conceptual hygiene is also key when designing and implementing effective and interpretable simulations. Second, theory provides necessary calibrations for complicated simulations. It's often essential in computational science to know at least those behavioral regimes that can be solved analytically, so that the code can be benchmarked against those cases. Finally, theory is frugal. It simply does not cost as much as either experiment or supercomputing. And, perhaps more importantly, in a discipline with a healthy balance between theory and experiment, theory helps reduce costs by precluding unnecessary experiment.

For all these reasons, both practically and methodologically, theory should have a primary role in most all of the sciences. Despite the recent vocal attacks and the more-worrying mission creep away from it, theory may still retake its proper place. For example, the centuries-old observation of allometric scaling in biological organisms—that metabolic rate is the three-quarters power of an organism's mass—was only recently put on a firm theoretical foundation. The underlying mechanism identified was the self-similar organization of nutrient transport networks. This showed in a rather transparent way why the previous area-to-volume explanation, proposed by Francis Galton in the nineteenth century, for allometric scaling was wrong [8]. And, the mechanistic insight suggested that even the organization of cities scales, too [9].

Another example comes from Heisenberg's *bête noire*—fluid turbulence. The basic equations of motion of fluid flow have been known for well over a century. Using these, a significant fraction of all supercomputer time is currently spent on flow simulations. And modern measurement technique facilitates collecting vast amounts of data on the temporal evolution of experimental fluid flows. It was only with the advent of nonlinear physics and mathematics, however, that an understanding of emergent flow structures, and the mechanisms that produce them, has now come tantalizingly within reach [10-15].

So, theory can play a positive role, but aren't massive data sets and computing power replacing theory? Can data and computing alone lead to understanding? Advocates interpret the evidence as suggesting that they can.

Today, it's a stunning feat that language translation engines have reached their levels of usability. Genuine semantic context is preserved in automated translation, when this was not previously the case just a few years ago. This performance is achieved not by new theoretical insights into human linguistic structures, rather they succeed by applying machine-learning techniques to mine massive corpora of translated written texts. In effect, the algorithms organize the original corpora into large, efficiently searchable tables. The "translation" returned is the closest match in the table to the input text [3]. This Linguist-Free approach is very telling on several fronts.

First, of course, is that historical language usage—language in the wild—appears more germane to how humans verbally communicate than any extant linguistic theory. Second, and this introduces the overarching philosophical concern here, these translation engines reduce to practice an attack against intelligent machines. The philosopher John Searle imagined a Chinese Room into which one slid paper with written Chinese text and out of which slipped a piece of paper with the translated English. With sufficient trials one would conclude that someone in the Room understood Chinese. Searle then noted that inside was a mechanism which used a huge reference dictionary for the translation. Clearly, the mechanism does not understand Chinese and so the intelligence is only imputed by the user. Moreover, Searle argued, not only does the device not understand Chinese, neither does it's builder [16]. Such is the state of affairs with current automated translation.

The very real possibility now exists that the Data Deluge will drive scientific explanation to become a Chinese Room. More precisely, my concern is that we scientific users will come to accept this brand of operationalism as understanding and this, in turn, will cost us our creativity.

Likewise with computing power, we can predict the weather more than a week ahead and forecast the population dynamics of spreading viral pandemics. Having powerful computers at hand, though, colors scientific practice. With modern software engineering tools we can now build (and manage) very large, “realistic” codes to simulate complicated natural phenomena. But is 30,000 lines of LISP code a theory of how the mind makes analogies? Having written that code does it mean we understand the brain’s process of analogy-making? Not hardly. In fact, even 10 lines of computer code, such as that for the now-famous chaotic Logistic Map, can produce complex behavior that half a century of hard mathematical work has yet to completely crack. Second and ironically, the very advances in computing power we champion now translate into large, exquisitely detailed models as complicated as natural experiment. Hundreds of components and parameters obscure the essential mechanisms responsible for a system’s organization and behavior and, practically, make it well-nigh impossible to systematically explore the range of possible system behaviors. Finally, the vast amounts of data automatically generated can be as rich as any empirical data set. And so, a computational scientist is often left with a data analysis task as daunting as any from experiment.

Let’s take stock. We have two trends, each driven by inexorably improving technology. On the one hand, we have *Data Literalism* born of the Data Deluge: All we need is data. “Data describes nature”, full stop. This rides tandem with the skepticism experimentalists hold of theory: “In any case, theory leaves out essential details of Nature and is, perforce, incomplete.” Its a short step to the production of large data sets becoming a goal unto itself. On the other hand, we have *Computationalism* born of high-performance computing: A computer code is a theory of the phenomenon it simulates; the programmer understands the mechanisms that produce the natural system’s behavior; and to be believable a simulation model must include all of a phenomenon’s detailed components.

In short, Data Literalism conflates science with data gathering and Computationalism conflates it with detailed simulation. Where is understanding in all of this? With Data Literalism pressing in from the side of experiment and Computationalism attempting an eclipse from the technology side, we are poised to squeeze out the path to understanding through theory. This is the challenge.

A century of gathering data on metabolic rates across hundreds of species did not lead to the new laws of biological scaling. The solution relied on the theory of fractional dimensions, which in turn has its origins in the mathematical theory of infinity developed by Georg Cantor. Thousands of experiments on turbulent fluids did not lead to our new views of emergent flow patterns.

Rather these arose from the theory of qualitative dynamics invented by the mathematician Henri Poincaré to explain the chaotic dynamics of the solar system [17].

There may be a saving grace, though. Recent progress in understanding complex systems suggests a new role for theory, one that relies essentially on the concepts of computation and information—at levels deeper than the tools (simulation and data acquisition) they engender. The hopeful view is that we are now on the verge of a new era, what in 1989 I called “artificial science”: the automated extraction of theories from data [4]. The goal is to analyze how we, intelligent agents that we are, discover novel patterns in nature—patterns that we’ve never seen before.

Contrast pattern recognition. What happens when the airline reservation system responds to your prompted verbal responses for a destination? The machine has a built-in vocabulary and asks intentionally focused questions to reduce the range of your responses and so increase the accuracy of its estimate of your intentions. In modern speech recognition systems there is an internal vocabulary of patterns—patterns that have been hand-designed by a speech engineer. The spectral data of your utterance is matched to the closest template in the pre-engineered vocabulary.

Many natural systems, however, even those for which we’ve established the underlying principles, produce organization spontaneously at spatial and temporal scales not directly determined by the microscopic balance of forces or equations of motion. Of late, we refer to this process as emergence [18,19]. For natural systems with emergent properties it simply begs the question to appeal to pattern recognition, since we don’t know ahead of time at which spatial and temporal scales patterns appear, let alone what they will be.

More to the point, what are “patterns” in the first place? How does nature form them? When is a measurement value due to some patternedness in the data or to some random component? How do we discover novel patterns? And, once identified, how do we build models of the mechanisms that produced them? These are the questions we need to answer to realize artificial science. Constructive answers will address what I consider the most fundamental and abiding problem of twenty-first century science: Pattern discovery. Discovering the unknown is a conundrum that is distinct from pattern recognition. Nonetheless, we as scientists do this all the time. To make progress on pattern discovery, a first step is to understand what patterns are and how to quantify them. (Is system A more “patterned” than system B?) The answer to this comes in the realization that nature’s patterns capture how nature stores and processes information [5].

With the conceptual challenge of emergence and a deeper understanding of the role of theory, there is now a new and very real possibility for a novel synthesis of advances in experimental technique, high-performance computing, and theory: automatically building theories from data [20,21]. That is, to the extent we understand pattern, we can use machines to find emergent organization in our vast data sets. And from there, the machines can build our theories, most likely with guidance from a new generation of theorists. This suggests a new target for scientific theory: A theory of theory building. Note that success in this will not put theorists out of work. Rather, it will allow them to work at a higher level, to be more productive, and to tackle systems which are that much more complex. The pace of progress will accelerate, yet again.

And this takes us back to May 1923 and the twenty-one year old budding quantum theorist. If you must have a doctoral near-disaster, you can do no better than on the topic Sommerfeld assigned Heisenberg: fluid turbulence. In his later years, Heisenberg would opine that turbulence was one of the most fundamental and difficult problems of contemporary physics [22]. Heisenberg was stymied by the richness and diversity of what was, after all, an inanimate system. As we look forward, perhaps we're in a similar humbled circumstance as we strive to understand the biological and then the social worlds. How are we to understand their emergent structures? How will we come to understand the underlying mechanisms well enough for social goals, such as human sustainability? Heisenberg's doctorate challenge to understand fluid turbulence might end up providing a lesson for twenty-first century science: a balanced interplay of experiment, computing, and theory will be required. ■

James P. Crutchfield is Director of the Complexity Sciences Center, Professor of Physics at the University of California at Davis, and External Faculty of the Santa Fe Institute; <http://cse.ucdavis.edu/~chaos/>.

References

- David Cassidy, "Uncertainty: The Life and Science of Werner Heisenberg", W. H. Freeman, New York (1992). See also "Werner Heisenberg—Student Years: 1920-1927", <http://www.aip.org/history/heisenberg/p06.htm>.
- Chris Anderson, "The End of Theory: The Data Deluge Makes the Scientific Method Obsolete", WIRED Magazine 16.07 (23 June 2008). See http://www.wired.com/science/discoveries/magazine/16-07/pb_theory.
- Alon Halevy, Peter Norvig, and Fernando Pereira, "The Unreasonable Effectiveness of Data", IEEE Computer (March-April 2009) 8-12.
- Alvy Ray Smith, personal communication (1986).
- James P. Crutchfield and Karl Young, "Inferring Statistical Complexity", Physical Review Letters **63** (1989) 105-108.
- Suzanne Still and James P. Crutchfield, "Structure or Noise?", arxiv.org:0708.0654 [physics.gen-ph].
- James P. Crutchfield, "Semantics and Thermodynamics" in Nonlinear Modeling and Forecasting, M. Casdagli and S. Eubank, editors, Addison-Wesley (Reading, Massachusetts, 1992) 317-359.
- Geoffery B. West, James H. Brown, and Brian J. Enquist, "A General Model for the Origin of Allometric Scaling Laws in Biology", Science **276**: 5309 (1997) 122-126.
- Luis M. A. Bettencourt, Jose Lobo, Dirk Helbing, Christian Kuhnert, and Geoffrey B. West, "Growth, Innovation, Scaling, and the Pace of Life in Cities", Proceedings of the National Academy of Sciences 104:17 (2007) 7301.
- Edward N. Lorenz, "Deterministic Nonperiodic Flow", J. Atmos. Sci. **20** (1963) 130.
- James P. Crutchfield and Kunihiko Kaneko, "Are Attractors Relevant to Turbulence?", Physical Review Letters **60** 1988: 2715-2718.
- Gregory Falkovich and Katepalli R. Sreenivasan, "Lessons from Hydrodynamic Turbulence", Physics Today (April 2006) 43-49.
- Tomas Bohr, Mogens H. Jensen, Giovanni Paladin and Angelo Vulpiani. **Dynamical Systems Approach to Turbulence**, Cambridge Nonlinear Science Series (No. 8) Cambridge University Press (Cambridge, UK, 1998).
- Paul Manneville, **Dissipative Structure and Weak Turbulence**, Academic Press (New York 1990).
- Bjorn Hof, Alberto de Lozar, Dirk Jan Kuik, and Jerry Westerweel, "Repeller or Attractor? Selecting the Dynamical Model for the Onset of Turbulence in Pipe Flow", Physical Review Letters **101** (2008) 214501.
- John Searle, "Consciousness and Language", Cambridge University Press, Cambridge, UK (2002).
- Peter Galison, **Einstein's Clocks, Poincare's Maps: Empires of Time**, W.W. Norton & Co. (New York, 2004).
- James P. Crutchfield, "Is Anything Ever New? Considering Emergence", in **Complexity: Metaphors, Models, and Reality**, G. Cowan, D. Pines, and D. Melzner, editors, Santa Fe Institute Studies in the Sciences of Complexity **XIX** (1994) 479-497.
- Robert Laughlin, **A Different Universe: Reinventing Physics from the Bottom Down**, Basic Books (2006).
- James P. Crutchfield and Bruce McNamara, "Equations of Motion from a Data Series", Complex Systems **1** (1987) 417-452.
- David Waltz and Bruce G. Buchanan, "Automating Science", Science **324** (3 April 2009) 43-44; Michael Schmidt and Hod Lipson, "Distilling Free-Form Natural Laws from Experimental Data", Science **324** (2009) 81-85.
- Werner Heisenberg, "Nonlinear Problems In Physics", Physics Today **20** (1967) 23-33.

This article originally was solicited in Spring 2009 as an opinion piece for Nature Magazine. It is based on presentations "Theory Theory—How We Come to Understand Our World", Santa Fe Institute Spring 2002 Trustee/Business Network Symposium, Santa Fe, New Mexico, 3 May 2002, and "Complex Systems Theory?" Santa Fe Institute Faculty Retreat, Bishop's Lodge, Santa Fe, New Mexico, 24-25 October 2003.