

Einstein's misgivings about his 1905 formulation of special relativity

Harvey R Brown

Faculty of Philosophy, University of Oxford, 10 Merton Street, Oxford OX1 4JJ, UK

E-mail: harvey.brown@philosophy.ox.ac.uk

Received 29 June 2005, in final form 30 June 2005

Published 21 September 2005

Online at stacks.iop.org/EJP/26/S85

Abstract

When Einstein formulated his 1905 treatment of relativistic kinematics, the template in his mind was thermodynamics. This was because a more desirable 'constructive' account of the behaviour of moving rods and clocks, based on the detailed physics governing their microscopic constitution, was unavailable. The price to be paid was appreciated by Einstein and a handful of others since 1905.

The principle of relativity is a principle that narrows the possibilities; it is not a model, just as the second law of thermodynamics is not a model.
Albert Einstein¹

1. Special relativity as a 'principle theory'

In January 1908, roughly two and a half years after publishing his celebrated paper on special relativity (Einstein 1905), Albert Einstein wrote in a letter to Arnold Sommerfeld:

So, first to the question of whether I consider the relativistic treatment of, e.g., the mechanics of electrons as definitive. No, certainly not. It seems to me too that a physical theory can be satisfactory only when it builds up its structures from *elementary* foundations. The theory of relativity is not more conclusively and absolutely satisfactory than, for example, classical thermodynamics was before Boltzmann had interpreted entropy as probability. If the Michelson–Morley experiment had not put us in the worst predicament, no one would have perceived the relativity theory as a (half) salvation. (Einstein 1995)

Einstein is repeating here an analogy between special relativity (henceforth SR) and thermodynamics that he had mentioned in a published note addressed to Ehrenfest already in 1907, in which he compared SR with 'the second law of the theory of heat' (Einstein 1907). In both cases, Einstein was emphasizing the limitations of SR, not its strengths.

¹ This statement was made by Einstein in 1911 at a scientific meeting in Zurich; see Galison (2004), p 268. In 1911 Einstein was still using 'principle of relativity' to mean theory of relativity.

Why was SR only a half salvation? Let us dwell for a minute on the analogy with thermodynamics.

Think of an idealized single-piston heat engine undergoing a Carnot cycle, and consider the theoretical limits of its efficiency. Such limits can, in principle, be established by exploiting knowledge of the microstructure of the working substance of the engine, and in particular by using the principles of statistical mechanics that apply to the molecular structure of the gas in the piston. A much easier approach, however, is to fall back on the laws of classical thermodynamics to shed light on the performance of the engine—phenomenological laws which stipulate nothing about the deep structure of the working substance. According to this approach, the efficiency of the heat engine must depend in a certain way on the ratio of the temperatures of the two heat reservoirs simply because, whatever the gas in the piston is made up of, if it did not it would be possible for the engine to act as a perpetual motion machine of ‘the second kind’. And this possibility is simply ruled out by hypothesis in thermodynamics.

But why cannot such a perpetual motion machine exist? Although statistical mechanics in principle accounts for this (or so it is widely held, even if the details involved are controversial) thermodynamics cannot. The impossibility of perpetual motion machines of various kinds is the very starting point of thermodynamics. What this theory gains in practicality and in the evident empirical solidity of its premisses, it loses in providing physical insight.

Einstein considered thermodynamics as the archetypical example of what he would call in 1919 a ‘principle theory’ in physics, one which is based on well-verified, but unexplained observable regularities. On the other hand, statistical mechanics, or more specifically the kinetic theory of gases, was for Einstein the prime example of a ‘constructive theory’, one built on the ‘elementary foundations’ mentioned in his 1908 letter. These foundations involve hypotheses about unseen fundamental processes—normally involving the microstructure of bodies and its mechanical principles. Principle theories are typically employed when constructive theories are either unavailable, too difficult to build, or relatively unwieldy. For according to Einstein, ‘when we say we have succeeded in understanding a group of natural processes, we invariably mean that a constructive theory has been found which covers the processes in question’.² Yet, Einstein stressed that SR is a principle theory. Why then did he feel it necessary to sacrifice explanatory content in developing his theory of relativity?

2. Rods, clocks and the quantum

Recall the title of Einstein’s 1905 relativity paper: ‘On the electrodynamics of moving bodies’. One of the great challenges of late nineteenth century electrodynamics and optics was to predict the outcome of experiments involving electromagnetic phenomena being performed in a laboratory *moving with respect to the luminiferous ether*. After all, the Earth is in motion relative to the centre of mass of the solar system, and at least some of the time must be moving relative to the ether—the invisible seat of electromagnetic phenomena. But by the turn of the century, the ether had become in the minds of some experts a very shadowy entity indeed. Made of an obscure kind of ‘imponderable matter’, its main role was increasingly just that of providing the inertial frame of reference relative to which the fundamental electromagnetic field equations of Maxwell were postulated to hold. The question was now: what form do the field equations have in Earth-bound frames that are moving relative to this fundamental frame?

² The explicit distinction between principle theories and constructive theories was articulated in a popular article Einstein wrote in 1919 for the *London Times* (Einstein 1919). But as we see in this article, it was a theme that ran through his writings from at least 1907 up until his death.

Einstein is famous for claiming in 1905, on the basis of his relativity principle, that all laws of physics, including those of electrodynamics, take the same form in all inertial reference frames, so happily Maxwell's equations can be used just as well in the moving laboratory frame. But this conclusion, or something very close to it, had already been anticipated by several great ether theorists, including the Dutchman H A Lorentz, the Ulsterman Joseph Larmor and particularly the French polymath Henri Poincaré. This was largely because there had been from the middle of the nineteenth century all the way to 1905 a series of experiments involving optical and electromagnetic effects that failed to show any sign of the ether wind rushing through the laboratory: it was indeed as if the Earth was always at rest relative to the ether. (The most famous of these, and the most surprising, was the 1887 Michelson–Morley experiment.) Like the above-mentioned ether theorists, Einstein realized that the covariance of Maxwell's equations—the form invariance of the equations—is achieved when the relevant coordinate transformations take a very special form, but Einstein was unique in his understanding that these transformations, properly understood, encode new predictions as to the behaviour of rigid bodies and clocks in motion. That is why, in Einstein's mind, a new understanding of space and time themselves was in the offing.

Both the mathematical form of the transformations, and at least the non-classical distortion of moving rigid bodies were already known to Lorentz, Larmor and Poincaré—indeed a family of possible deformation effects was originally suggested independently by Lorentz and the Irish physicist G F FitzGerald to explain the Michelson–Morley result³. It was the connection between them, i.e., between the coordinate transformations and motion-induced deformation, that had not been fully appreciated before Einstein. In the first ('kinematical') part of his 1905 relativity paper, Einstein established the operational meaning of the so-called Lorentz coordinate transformations and showed that they lead not just to a special case of FitzGerald–Lorentz deformation (longitudinal contraction), but also to the 'slowing down' of clocks in motion—the phenomenon of time dilation. Now it is still not well known that Larmor and Lorentz had come tantalizingly close to predicting this phenomenon; they had independently seen just before the turn of the century how it must hold in certain very special cases. But as a general effect that does not depend on the constitution of a clock, its discovery was Einstein's own.

Einstein did something else that was new and important in the kinematical part of his paper. He derived the Lorentz transformations not from the symmetry properties of Maxwell's equations, but by using an argument inspired by thermodynamics. Why?

Several months before he wrote his paper on SR, Einstein had written a revolutionary paper claiming that electromagnetic radiation has a granular structure. The suggestion that radiation was made of quanta—or photons as they would later be dubbed—was the basis of Einstein's extraordinary treatment of the photoelectric effect in the same paper. This treatment would win its author the Nobel prize of 1921; acceptance of the photon by the physics community would take longer. But the immediate consequence of Einstein's commitment to the photon was to destabilize in his mind all the previous work on the electrodynamics of moving bodies.

All the work of the ether theorists was based on the assumption that Maxwellian electrodynamics is strictly true, and not just true on average. In the work of Lorentz, Larmor and Poincaré, the Lorentz transformations make their appearance as symmetry transformations (whether considered approximate or otherwise) of these equations. But Maxwell's equations are incompatible with the existence of the photon.

In his 1949 *Autobiographical Notes*, published when he was 67, Einstein was clear about the seismic implications of this conundrum.

³ For recent treatments of this episode, see Brown (2001, 2005).

Reflections of this type [on the dual wave–particle nature of radiation] made it clear to me as long ago as shortly after 1900, i.e., shortly after Planck’s trailblazing work, that neither mechanics nor electrodynamics could (except in limiting cases) claim exact validity. By and by I despaired of the possibility of discovering the true laws by means of constructive efforts based on known facts. (Einstein 1969), pp 51, 53

Already in the *Notes*, Einstein had pointed out that the general validity of Newtonian mechanics came to grief with the success of the electrodynamics of Faraday and Maxwell, which led to Hertz’s detection of electromagnetic waves—‘phenomena which by their very nature are detached from every ponderable matter’ *op cit*, p 25. Later, he summarized the nature of Planck’s 1900 derivation of his celebrated black-body radiation formula, in which quantization of absorption and emission of energy by the mechanical resonators is presupposed. Einstein noted that although this contradicted the received view, it was not immediately clear that electrodynamics—as opposed to mechanics—was violated. But now with the emergence of the light quantum, not even electrodynamics was sacrosanct.

All my attempts . . . to adapt the theoretical foundation of physics to this [new type of] knowledge failed completely. It was if the ground had been pulled out from under one, with no firm foundation to be seen anywhere, upon which one could have built. *Op cit*, p 45

Earlier in the *Notes*, Einstein had sung the praises of classical thermodynamics, ‘the only physical theory of universal content concerning which I am convinced that, within the framework of the applicability of its basic concepts, it will never be overthrown’. Now, he explains how the very structure of the theory was influential in the search for a way out of the turn-of-the-century crisis in physics.

The longer and more despairingly I tried, the more I came to the conviction that only the discovery of a universal formal principle could lead us to assured results. The example I saw before me was thermodynamics. The general principle was there given in the theorem⁴: the laws of nature are such that it is impossible to construct a *perpetuum mobile* (of the first and second kind). How, then, could such a universal principle be found? *Op cit*, p 53

3. Einstein’s misgivings

It is well known that Einstein based his derivation of the Lorentz transformations on a combination of the relativity principle (essentially the same as that defended by Newton) and his so-called light postulate. (The latter was the claim that relative to a certain inertial frame, the speed of light is independent of the speed of the source and isotropic—something every ether theorist took for granted when the frame in question is taken to be the fundamental ether rest frame⁵ and something which remarkably Einstein felt would survive the rise of quantum theory intact.) He showed that length contraction for rigid rods and time dilation for ideal clocks are consequences of these phenomenological assumptions, in the same way that, say, the existence of entropy and its non-decreasing behaviour over time for adiabatic systems are a consequence of the laws of thermodynamics.

⁴ The word ‘theorem’ for ‘Sätze’ in the translation by P A Schilpp is perhaps better rendered as ‘sentence’ or ‘statement’. I thank Thomas Müller for discussion of this point.

⁵ In 1921, Wolfgang Pauli would correctly describe Einstein’s light postulate as ‘true essence of the old aether point of view’; (Pauli 1981), p 5. It should also be noted that the derivation of the Lorentz transformations requires a third, admittedly innocuous, assumption: the isotropy of space.

Einstein would have preferred a constructive account of these relativistic effects, presumably based on the nature of the non-gravitational forces that hold the constituent parts of rods and clocks together. But as we have seen, for Einstein the elements of such an account were not to be had in 1905. The price to be paid for the resulting strategic retreat to a principle theory approach was not just loss of insight; Einstein became increasingly uneasy about the role played by rods and clocks in this approach. This unease is seen in a paper entitled 'Geometry and Experience' he published in 1921 (Einstein 1921), and in particular in his 1949 *Autobiographical Notes*:

One is struck [by the fact] that the theory [of special relativity] . . . introduces two kinds of physical things, i.e., (1) measuring rods and clocks, (2) all other things, e.g., the electromagnetic field, the material point, etc. This, in a certain sense, is inconsistent; strictly speaking measuring rods and clocks would have to be represented as solutions of the basic equations (objects consisting of moving atomic configurations), not, as it were, as theoretically self-sufficient entities. However, the procedure justifies itself because it was clear from the very beginning that the postulates of the theory are not strong enough to deduce from them sufficiently complete equations . . . in order to base upon such a foundation a theory of measuring rods and clocks. . . . But one must not legalize the mentioned sin so far as to imagine that intervals are physical entities of a special type, intrinsically different from other variables ('reducing physics to geometry', etc). (Einstein 1969 pp 59, 61)

These remarks are noteworthy for several reasons.

First, there is the issue of justifying the 'sin' of treating rods and clocks as primitive, or unstructured entities in SR. Einstein does not say in 1949, as he did in 1908 and 1921, that the 'elementary' foundations of a constructive theory of matter are still unavailable; rather he simply reminds us of the limits built into the very form of the 1905 theory. It is hardly any justification at all. Considerable progress in the relativistic quantum theory of matter *had* been made between 1905 and 1949. Was it Einstein's long-standing distrust of the quantum theory that held him back from recognizing this progress and its implications for his formulation of SR?

Second, consider the criticism Abraham Pais made of H A Lorentz in his acclaimed 1982 biography of Einstein: 'Lorentz never fully made the transition from the old dynamics to the new kinematics.' (Abraham 1982 p 167). As late as 1915 Lorentz thought that the relativistic contraction of bodies in motion can be explained if the known property of distortion of the electrostatic field surrounding a moving charge is supposed to obtain for all the other forces that are involved in the cohesion of matter. In other words, Lorentz viewed such kinematical effects as length contraction as having a dynamical origin, and it is this notion that Pais found reprehensible. Yet, when Einstein appeals to the nature of rods and clocks as 'moving atomic configurations', it seems that not even he ever fully made the transition from the old dynamics to the new kinematics. For to say that length contraction is intrinsically kinematical would be like saying that energy or entropy are intrinsically thermodynamical, not mechanical—something Einstein would never have accepted.

The limitations of Einstein's principle-theory approach to SR have been noted by a number of commentators since 1905, including Wolfgang Pauli and Arthur Eddington in the 1920s, W F G Swann in the 1940s, and Lajos Jánossy and John S Bell in the 1970s; see Pauli (1981), Eddington (1928 p 7), Swann (1941), Jánossy (1971), and Bell (1976, 1992). All of these authors called for a more constructive version of SR. It was perhaps Bell—whose name will be familiar to many readers through his famous inequality in quantum mechanics—who made the point in the clearest fashion.

If you are, for example, quite convinced of the second law of thermodynamics, of the increase of entropy, there are many things that you can get directly from the second law which are very difficult to get directly from a detailed study of the kinetic theory of gases, but you have no excuse for not looking at the kinetic theory of gases to see how the increase of entropy actually comes about. In the same way, although Einstein's theory of special relativity would lead you to expect the FitzGerald contraction, you are not excused from seeing how the detailed dynamics of the system also leads to the FitzGerald contraction. (Bell 1992)

What is remarkable is that Bell himself seemed to be unaware of Einstein's own distinction between principle and constructive theories, and his repeated references to the analogy between SR and thermodynamics. At any rate, Bell stressed that he had no 'reservation whatever about the power and precision of Einstein's approach'; his main point was that 'the longer road [a dynamical account of contraction and dilation] sometimes gives more familiarity with the country' (Bell 1976).⁶

References

- Bell J S 1976 How to teach special relativity *Prog. Sci. Cult.* **1** 67–80 (reprinted in (Bell 1987))
- Bell J S 1987 *Speakable and Unspeakable in Quantum Mechanics* (Cambridge: Cambridge University Press)
- Bell J S 1992 George Francis FitzGerald *Phys. World* **5** 31–5 (1989 lecture, abridged by D Weare)
- Brown H R 2001 The origins of length contraction: I. The FitzGerald–Lorentz deformation hypothesis *Am. J. Phys.* **69** 1044–54 (Preprints gr-qc/0104032, PITT-PHIL-SCI 218)
- Brown H R 2005 *Physical Relativity: Space-time Structure from a Dynamical Perspective* (Oxford: Oxford University Press)
- Brown H R and Pooley O 2001 The origins of the spacetime metric: Bell's Lorentzian pedagogy and its significance in general relativity *Physics Meets Philosophy at the Planck Scale* ed C Callender and N Huggett (Cambridge: Cambridge University Press) pp 256–72 (Preprint gr-qc/9908048)
- Callender C and Huggett N (ed) 2001 *Physics Meets Philosophy at the Planck Scale* (Cambridge: Cambridge University Press)
- Eddington A S 1928 *The Nature of the Physical World* (Cambridge: Cambridge University Press)
- Einstein A 1905 Zur Elektrodynamik bewegter Körper *Ann. Phys., Lpz.* **17** 891–921
- Einstein A 1907 Bemerkung zur Notiz des Herrn P. Ehrenfest: Translation deformierbarer Elektronen und der Flächensatz *Ann. Phys., Lpz.* **23** 206–8 (English translation in Einstein 1989 Doc 44, pp 236–7)
- Einstein A 1919 What is the theory of relativity? *The London Times* (reprinted in Einstein 1982 pp 227–32)
- Einstein A 1921 Geometrie und Erfahrung *Erweite Fassung des Festvortrages gehalten an der preussischen Akademie* (Berlin: Springer) (translated by S Bargmann as 'Geometry and Experience' in Einstein 1982 pp 232–46)
- Einstein A 1969 Autobiographical Notes *Albert Einstein: Philosopher-Scientist* vol 1, ed P A Schilpp (Peru, IL: Open Court) pp 1–94
- Einstein A 1982 *Ideas and Opinions* (New York: Crown)
- Einstein A 1989 *The Collected Papers of Albert Einstein, Vol 2: The Swiss Years: Writings, 1900–1909 (English Translation Supplement)* (Princeton, NJ: Princeton University Press) (translated by A Beck)
- Einstein A 1995 Letter to Arnold Sommerfeld, January 14, 1908. Document 73 *The Collected Papers of Albert Einstein, Vol 5, The Swiss Years: Correspondence, 1902–1914 (English Translation Supplement)* ed M J Klein, A J Kox and R Schulmann (Princeton, NJ: Princeton University Press) (translated by A Beck)
- Galison P 2004 *Einstein's Clocks, Poincaré's Maps: Empires of Time* (London: Hodder and Stoughton)
- Jánossy L 1971 *Theory of Relativity Based on Physical Reality* (Budapest: Akadémia Kiadó)
- Pais A 1982 'Subtle is the Lord. . .' *The Science and the Life of Albert Einstein* (New York: Oxford University Press)
- Pauli W 1981 *Theory of Relativity* (New York: Dover) originally published as
- Pauli W 1921 Relativitätstheorie *Encyklopädie der mathematischen Wissenschaften, mit Einschluss ihrer Anwendungen* vol 5. *Physik* ed A Sommerfeld (Leipzig: Tauber) pp 539–775
- Swann W F G 1941 Relativity, the FitzGerald–Lorentz contraction, and quantum theory *Rev. Mod. Phys.* **13** 197–202

⁶ Note that the longer road Bell has in mind does not involve commitment to the existence of any preferred inertial frame, or ether. For a discussion of Bell's 1976 treatment of SR by way of a 'Lorentzian pedagogy', see Brown and Pooley (2001). For a broader discussion of both the historical and conceptual issues raised in this paper, see Brown (2005).