# Logic of gauge

Alexander Afriat

May 18, 2013

**Abstract**

The logic of gauge theory is considered by tracing its development from general relativity to Yang-Mills theory, through Weyl's two gauge theories. A handful of elements—which for want of better terms can be called *geometrical justice*, *matter wave*, *second clock effect*, *twice too many energy levels*—are enough to produce Weyl's second theory; and from there, all that's needed to reach the Yang-Mills formalism is a *non-Abelian structure group* (say $\mathbb{SU}(N)$).

## 1 Introduction

Experience is less conspicuous here than the imaginative virtuosity of Hermann Weyl—who might have carried it even further, going beyond his second gauge theory (1929) to reach Yang-Mills theory, the generalisation from the Abelian structure group $\mathbb{U}(1)$ to $\mathbb{SU}(N)$ being most natural: why stop at varying just the phase of a spinor and not change its direction too? Besides, Weyl was already using a non-Abelian 'structure' group[1]—even larger than $\mathbb{SU}(2)$—for gravity, in other words for the covariant differentiation of spinors against a curved space-time background . . .

Since the transition (§2.1) from general relativity to Weyl's first gauge theory (1918, §2) has been amply discussed in Pais (1982), Vizgin (1984), Scholz (1994, 1995, 2001a, 2004, 2011b), Cao (1997), Hawkins (2000), Coleman & Korté (2001), Sigurdsson (2001), Ryckman (2003a,b, 2005, 2009) and Afriat (2009), I'll concentrate on the next step, that took Weyl to his second theory (§3)—mainly determined by the new undulatory ontology (§3.1) introduced by Louis de Broglie (1924), Dirac (1925), Schrödinger (1926) and others. Mainly but not wholly: Weyl had every reason to keep the electricity, gravity and gauge structure of his first theory; but now with three ingredients (matter, electricity, gravity), two different gauge relations were possible. Einstein's objection (1918, §2.3), the *second clock effect*, ruled out the old gauge relation ((3)-(5)) between electricity and gravity, leaving the new relation ((3)-(7)) between electricity and matter.

Weyl's objection that Dirac theory provided *twice too many energy levels* (§3.2) is only relevant to his own story, of how he reached his two-component theory of

---

[1] By this I mean not the gauge group $\mathscr{G}$ (vertical automorphisms) acting on the whole base manifold $M$ but the much smaller group acting (on the vector space $\mathbb{V}_x$) at a generic point $x \in M$. I have called this particular 'structure' group $\mathbb{W}(2, \mathbb{C})$; see §§3.3-3.5.

1929, and not to the ultimate derivation of Yang-Mills theory (§4)—which by no means favours Weyl's two-component theory over Dirac's four-component theory.

Three different gauge arguments are looked at in §3.5. The first (§3.5.1), though less of an explicit 'argument' than the other two, is enough to yield the compensation of (3) by (7). In the second (§3.5.2) Weyl extracts a *curved* connection from $\mathbb{U}(1)$. The third (§3.5.3) is the standard modern "gauge argument" or "gauge principle" which produces an exact connection without curvature, and shouldn't be blamed on Weyl.[2]

# 2 Weyl's first gauge theory

## 2.1 Geometrical justice

First, there was general relativity.[3] Levi-Civita (1917) saw that the connection determined by Einstein's covariant derivative transported the *direction* of a vector anholonomically, but not its *length*, which was left unchanged.[4] This was unfair, protested Weyl—length deserved the same treatment as direction.[5] To remedy he proposed a more general theory that propagated length just as anholonomically as direction. *Congruent* transport would also be governed by a connection, which Weyl defined as a bilinear mapping between neighbouring points: linear in the object propagated and in the direction of propagation.[6] A connection transporting the (squared) length $l$ from $a = \gamma(a)$ to its neighbour[7] $b = \gamma(b)$ along $\gamma : [a, b] \to M$ would therefore be a real-valued[8] one-form $A = A_\mu dx^\mu$ applied to the direction $\dot\gamma = \dot\gamma^\mu \partial_\mu \in T_a M$ and multiplied by the initial length $l_a$, yielding the small difference

$$\delta l = l_a - l_b = l_a \langle A, \dot\gamma \rangle$$

subtracted from $l_a$.[9] The final length $l_b$ is $l_a(1 - \langle A, \dot\gamma \rangle)$—unless $a$ and $b$ are too far apart for $\gamma$ to remain straight in between, in which case $l_b$ is

$$l_a \exp \int_\gamma A.$$

---

[2] See Afriat (2013).

[3] Einstein (1916)

[4] See Ryckman (2003b) p. 80, Ryckman (2009) p. 288.

[5] See Afriat (2009) for the details of this *geometrical justice*—which can also be understood in terms of group extensions (see Scholz (2004) pp. 183, 189, 191-2, Scholz (2011a) pp. 195, Scholz (2011b), third page of the paper): since a Levi-Civita connection subjects direction to $\mathbb{SO}^+(1, 3)$ but length to the (group containing only the) identity $\mathbb{1}$, it is only fair to extend the identity by the dilations, yielding $\mathbb{1} \times \mathbb{R}^+ = \mathbb{R}^+$—which (unlike $\mathbb{1}$) allows length anholonomies and therefore geometrical justice. The total group, for direction and length together, is the extension $\mathbb{SO}^+(1, 3) \times \mathbb{R}^+$ giving the relativistic similarities. Ryckman (2003a,b, 2005, 2009) provides an alternative account of Weyl's agenda.

[6] More on connections in §2.4.

[7] Which is so close to $a$ it practically belongs to the tangent space $T_a M$; see Weyl (1926) p. 28, Weyl (1931a) p. 52.

[8] Here the structure group is the multiplicative group $\mathbb{R}^+$ of dilations, generated by the Lie algebra $\langle \mathbb{R}, +, [\,\cdot\,,\cdot\,] \rangle$ or rather $\langle \mathbb{R}, + \rangle$; the Lie product $[\,\cdot\,,\cdot\,]$ vanishes since real numbers commute.

[9] *Cf.* Ryckman (2009) pp. 290-1.

To deal with the geometrical injustice that $A$ was introduced to remedy, the curvature[10]

(1) $$F = \frac{1}{2} F_{\mu\nu} dx^\mu \wedge dx^\nu = dA = \frac{1}{2}(\partial_\mu A_\nu - \partial_\nu A_\mu) dx^\mu \wedge dx^\nu$$

cannot vanish—unlike the three-form

(2) $$dF = d^2 A = \frac{1}{6}(\partial_\mu F_{\nu\sigma} + \partial_\nu F_{\sigma\mu} + \partial_\sigma F_{\mu\nu}) dx^\mu \wedge dx^\nu \wedge dx^\sigma,$$

which does. Seeing all this, Weyl couldn't help thinking[11] of the electromagnetic four-potential $A$, the Faraday two-form $F = dA$ and Maxwell's two homogeneous equations[12] $dF = 0$: he had unified gravity and electromagnetism, by mistake![13] And indeed Einstein would soon point out the mistake: the anholonomy on which Weyl had based his theory is not observed in nature, as we'll see in §2.3.

## 2.2 Gauge

Weyl sought to rectify general relativity using the curvature (1), which ensured geometrical justice.[14] Differentiation is destructive (or rather irreversible); what $d$ destroys is the freedom

(3) $$A \mapsto A' = A - d\lambda$$

invisible to $F = dA = dA'$, in the sense that the inverse image $[A] = d^{-1}F$ of $F$ under $d$ is the whole equivalence class given by the equivalence relation $A \sim (A + d\lambda)$. If $A$ only served to produce the curvature $F$, (3) would be vacuous; but $A$ appears elsewhere too, notably in the law of propagation[15]

(4) $$\nabla g = A \otimes g,$$

which is not indifferent to (3). To make (4) invariant, (3) therefore has to be balanced by

(5) $$g \mapsto g' = e^\lambda g.$$

Such compensation is typical[16] of a gauge theory: an invariant expression (here (4)) is sensitive to a first transformation, and to a second as well—but indifferent to the two together, if their variations are appropriately constrained, and balance one another.[17]

---

[10]Einstein's summation convention will sometimes be used.

[11]See Scholz (2001a) p. 75, Ryckman (2003a) p. 92, Ryckman (2005) p. 158.

[12]$\nabla \cdot \mathbf{B} = 0$ and $\nabla \times \mathbf{E} + \partial_t \mathbf{B} = 0$

[13]Ryckman (2003b) p. 61: "[...] Weyl did not start out with the objective of unifying gravitation and electromagnetism, but sought to remedy a perceived blemish in Riemannian 'infinitesimal' geometry. The resulting 'unification' was, as it were, serendipitous." See also p. 63, Ryckman (2003a) p. 86, Ryckman (2005) pp. 149-54, 158, Ryckman (2009) pp. 287-94.

[14]Vanishing *Streckenkrümmung* (length curvature) $F$ led to *holonomic* congruent transport which, alongside *anholonomic* parallel transport, offended Weyl's sense of geometrical justice.

[15]This generalises the length-preserving condition $\nabla g = 0$ satisfied by the Levi-Civita connection.

[16]Typical but mysterious, even for Weyl (1931a) p. 54: "insbesondere konnte ich nichts a priori Einleuchtendes vorbringen zugunsten der Koppelung des willkürlichen additiven Gliedes $\partial\lambda/\partial x_p$, das nach der Erfahrung in den Komponenten des elektromagnetischen Potentials steckt, mit dem von der klassischen Geometrie geforderten Eichfaktor $e^\lambda$."

[17]See Ryckman (2003a) p. 77.

So far we have two logical ingredients

1. GR: *general relativity*

2. GJ: *geometrical justice*

which together yield Weyl's theory of electricity and gravity; I'll write

$$\text{GR \& GJ} \rightarrow \text{W18}.$$

The next will be MW: *matter wave* and SC: (avoid) *second clock effect* ...

## 2.3 Einstein's objection

The tangent of a worldline's *image* $\bar{\gamma} \subset M$ only has a direction; the length $l = \|\dot{\gamma}\|^2 = g(\dot{\gamma}, \dot{\gamma})$ of the tangent $\dot{\gamma} = d\gamma/dt$ is given by the parameter rate $\partial\gamma/\partial t$. If the values of the parameter are identified with the readings of a clock describing $\gamma$, the length $l$ giving the proper ticking rate should remain constant—the hands of a good clock don't accelerate. But far from remaining constant, lengths in Weyl's theory aren't even integrable:

$$l_b(\gamma) = l_a \exp \int_\gamma A$$

depends on $\gamma$—whereas an exact connection $A = d\mu$ would of course give

$$l_b = l_a \exp \int_a^b d\mu = l_a \exp \Delta\mu$$

along any path, $\Delta\mu$ being the difference $\mu(b) - \mu(a)$ between the final and initial values of $\mu$. In addition to the *first* clock effect (Langevin's twins) already present in Einstein's theory, Weyl's theory therefore has a *second* clock effect expressed in the anholonomy of ticking rates.

Einstein objected that *nature provides integrable clocks*.[18] Two clocks trace out a loop $\bar{\gamma} = \partial\omega$ enclosing a region $\omega$ (without holes): starting from the same point $a$ they

---

[18]Letter to Weyl dated 15 April 1918: "So schön Ihre Gedanke ist, muss ich doch offen sagen, dass es nach meiner Ansicht ausgeschlossen ist, dass die Theorie die Natur entspricht. Das $ds$ selbst hat nämlich reale Bedeutung. Denken Sie sich zwei Uhren, die relativ zueinander ruhend neben einander gleich rasch gehen. Werden sie voneinander getrennt, in beliebiger Weise bewegt und dann wieder zusammen gebracht, so werden sie wieder gleich (rasch) gehen, d. h. ihr relativer Gang hängt nicht von der Vorgeschichte ab. Denke ich mir zwei Punkte $P_1$ & $P_2$ die durch eine Zeitartige Linie verbunden werden können. Die an $P_1$ & $P_2$ anliegenden zeitartigen Elemente $ds_1$ und $ds_2$ können dann durch mehrere zeitartigen Linien verbunden werden, auf denen sie liegen. Auf diesen laufende Uhren werden ein Verhältnis $ds_1 : ds_2$ liefern, welches von der Wahl der verbindenden Kurven unabhängig ist.—Lässt man den Zusammenhang des $ds$ mit Massstab- und Uhr-Messungen fallen, so verliert die Rel. Theorie überhaupt ihre empirische Basis." Another letter to Weyl, four days later: "wenn die Länge eines Einheitsmassstabes (bezw. die Gang-Geschwindigkeit einer Einheitsuhr) von der Vorgeschichte abhingen. Wäre dies in der Natur wirklich so, dann könnte es nicht chemische Elemente mit Spektrallinien von bestimmter Frequenz geben, sondern es müsste die relative Frequenz zweier (räumlich benachbarter) Atome der gleichen Art im Allgemeinen verschieden sein. Da dies nicht der Fall ist, scheint mir die Grundhypothese der Theorie leider nicht annehmbar, deren Tiefe und Kühnheit aber jeden Leser mit Bewunderung erfüllen muss."

describe worldlines $\gamma_1$, $\gamma_2$ that meet at $b$. They tick at the same rate if $A$ is exact, for then

$$\oint_{\partial\omega} d\mu = \iint_\omega d^2\mu$$

vanishes—in fact (without holes) it is enough for $A$ to be closed,

$$\oint_{\partial\omega} A = \iint_\omega dA$$

vanishes too provided $dA$ does. But if the loop encloses an electromagnetic field $F = dA$, one of the clocks will tick faster than the other once they're compared at $b$. In any case the theory didn't work: from the beginning it rested on an anholonomy not seen in nature.[19]

## 2.4 Connections

A few more words about connections, which provide a notion of constancy[20] or free fall or absence of force, even in the presence of influences that can complicate it. The fundamental operation is linear and infinitesimal, and can be extended by integration. It serves to displace an object, say an $N$-dimensional vector $\xi \in \mathbb{C}^N$ subject to $\mathbb{SU}(N)$,[21] from $a$ to a neighbouring point $b$. The operator $e^{i\mathbf{T}_{ab}} : \mathbb{C}_a^N \to \mathbb{C}_b^N$ will be unitary provided $\mathbf{T}_{ab}$ is Hermitian. One can also write $\xi_b = (\mathbb{1}_N - i\mathbf{T}_{ab})\xi_a$, where

$$i\mathbf{T}_{ab}\xi_a = \delta\xi_b = \xi_b - \xi_b' \in \mathbb{C}_b^N$$

is the difference between the constant or 'unperturbed' vector $\xi_b$ (given by the connection) and the value $\xi_b'$ of the 'perturbed' section $\xi'$ (defined independently of the connection). 'Free fall' corresponds to the identity of $\mathbb{SU}(N)$, to the origin of the Lie algebra $\mathfrak{su}(N)$—which makes the correction $\delta\xi_b$ vanish. But different fibers $\mathbb{C}_a^N$, $\mathbb{C}_b^N$ are best related with respect to bases $\varphi^m \subset \mathbb{C}^{N*}$, $\varphi_n \subset \mathbb{C}^N$: The operator $U = e^{i\mathbf{T}}$ becomes[22] $U_n^m = \langle \varphi^m, U\varphi_n \rangle$ and $\mathbf{T}_n^m$ equals $\langle \varphi^m, \mathbf{T}\varphi_n \rangle$, whereas the components of a vector are $\xi^m = \langle \varphi^m, \xi \rangle$; $m, n = 1, \ldots, N$. If the section $\xi'$ and basis are in free fall between $a$ and $b$, the difference

$$\delta\xi_b^m = \xi_b^m - \xi_b'^m = i\sum_{n=1}^N \mathbf{T}_n^m \xi_a^n$$

vanishes.[23] The basis can always, according to an appropriate 'equivalence principle,' be made to 'accelerate with' a section $\xi'$ at any point $x$ so that $\delta\xi_x^m$ vanishes.

---

[19]*Cf.* Ryckman (2009) p. 295.

[20]See Ryckman (2003b) p. 80, Ryckman (2005) pp. 151-2, Ryckman (2009) p. 290: "vector at $P'$ is "the same" as a given vector at $P$. Namely from the original vector at $P$, a new vector arises at $P'$ [ . . . ] affirmed to be "without change"."

[21]For one might as well go straight to the general non-Abelian (*i.e.* Yang-Mills) case.

[22]With all the indices: $(U_{ab})_n^m = \langle \varphi_b^m, U_{ab}\varphi_{an} \rangle$.

[23]*Cf.* Weyl (1918b) p. 392, Weyl (1988) pp. 113, 122, 126, Weyl (1921) p. 542, Weyl (1929b) p. 339, Scholz (1994), Ryckman (2003b) p. 79, Ryckman (2005) pp. 151-2, Ryckman (2009) p. 289.

If we now go from $a$ to $b$ then to its neighbour $c$ we have $\prod_{ac} e^{i\mathbf{T}} = e^{i\mathbf{T}_{bc}} e^{i\mathbf{T}_{ab}}$, and for arbitrary paths[24] $\prod_\gamma e^{i\mathcal{A}}$, where $\mathcal{A}$ is the Yang-Mills connection (see §4)

$$(6) \qquad \mathcal{A} = \sum_{\mu=0}^{3} dx^\mu \otimes \mathcal{A}_\mu = i \sum_{\mu=0}^{3} \sum_{k=1}^{N} \mathcal{A}_\mu^k dx^\mu \otimes \mathbf{T}_k,$$

which (unlike the infinitesimal generator $i\mathbf{T}$ on its own) takes account $\langle \mathcal{A}, \dot\gamma \rangle = i\mathbf{T}$ of the direction $\dot\gamma \in T_x M$ of propagation. The product $\prod$ is needed rather than the sum (or integral) because the exponential $e^{i(\mathbf{T}+\mathbf{T}'+\cdots)}$ only makes sense if the summands commute.

We have seen that the congruent transport of Weyl's first gauge theory is produced by a real-valued connection $A$, which, applied to a tangent vector $\dot\gamma \in T_x M$, yields a dilation generator $\langle A, \dot\gamma \rangle$ belonging to the Lie algebra $\mathbb{R}$ of the dilation group $\mathbb{R}^+$. In the second gauge theory[25] the structure group becomes $\mathbb{U}(1) = e^{i\mathbb{R}}$, whose Lie algebra is $i\mathbb{R} = \mathrm{Lie}\,\mathbb{U}(1)$.[26]

# 3   Weyl's second gauge theory

The setback of 1918 and Einstein's objection (his *preaching*![27]) had their benefits, they taught Weyl the experimental character of physics, an *empirical discipline founded directly on experience and not a geometrical fantasy deduced from æsthetic hunches*: "All these geometrical leaps-in-the-air [W18] were premature, we return [W29] to the solid ground of physical facts."[28] In 1929 he's understood, matured, and will have his "revenge";[29] he'll even claim that his new theory came straight out of experience,[30] directly derived from spectrographic data[31] . . .

For his new theory takes account of the electron's spin—which in fact got there through the Dirac equation; and in Dirac's argument (1928) spin does not come (straight)

---

[24]See Göckeler & Schücker (1987) p. 51, Healey (2007) p. 63, Gambini & Pullin (2011) pp. 68-9 for more details.

[25]Weyl (1929a,b,c). See Straumann (1987), O'Raifeartaigh (1997), O'Raifeartaigh & Straumann (2000), Brading (2002), Scholz (2004, 2005, 2006) for more recent accounts.

[26]This section owes much to discussions with Jean-Philippe Nicolas. I am also indebted to an anonymous referee, who spotted a mistake.

[27]Letter to Seelig—quoted in Seelig (1960) p. 274—in which Weyl quotes Einstein: "So – das heisst auf so spekulative Weise, ohne ein leitendes, anschauliches physikalisches Prinzip – macht man keine Physik!"

[28]Weyl (1931a) p. 56: "Alle diese geometrischen Luftsprünge waren verfrüht, wir kehren zurück auf den festen Boden der physikalischen Tatsachen." *Cf.* Scholz (2011a) pp. 190-1.

[29]"Rache"; Pauli (1979) p. 518: "Als Sie früher die Theorie mit $g'_{ik} = \lambda g_{ik}$ machten, war dies reine Mathematik und unphysikalisch. Einstein konnte mit Recht kritisieren und schimpfen. Nun ist die Stunde der Rache für Sie gekommen; jetzt hat Einstein den Bock des Fernparallelismus geschossen, der auch nur reine Mathematik ist und nichts mit Physik zu tun hat, und Sie können schimpfen!"

[30]Weyl (1931a) p. 57: "Das neue Prinzip ist aus der *Erfahrung* erwachsen und resümiert einen gewaltigen, aus der Spektroskopie entsprungenen Erfahrungsschatz." On Weyl's 'empirical turn' see Scholz (2004) pp. 165, 183, 191-3.

[31]Weyl (1931a) p. 57: "Dieses Transformationsgesetz der $\psi$ ist zuerst von PAULI aufgestellt worden und folgt mit unfehlbarer Sicherheit aus den spektroskopischen Tatsachen, genauer aus den Termdubletts der Alkalispektren und der Tatsache, daß die Dublettkomponenten nach Ausweis ihres Zeemaneffekts *halbganze* innere Quantenzahlen besitzen."

out of experience[32] but out of a mathematical, æsthetic, *a priori* principle, in much the same spirit as the geometrical justice that produced Weyl's first gauge theory.

## 3.1 The new undulatory ontology

But let us go back a few years. As mentioned in the Introduction, Louis de Broglie (1924), Dirac (1925), Schrödinger (1926) *et al.* had meanwhile produced an *undulatory* world. Weyl had no reason to get rid of electricity or gravitation; to those existing ingredients he therefore had to add a matter wave, to update his ontology. As long as there was only gravity and electricity, the gauge relation (3)-(5) could only hold between *them*; but now, with a third element, as many compensations were in principle possible, of which only two were plausible: the old relation between gravity and electricity, and a new one between electricity and the matter wave. With (3)-(5) the theory would have remained subject to Einstein's objection—which the presence of the electron's wavelength[33] $h/mc$ in the Dirac equation made even more convincing,[34] by providing an absolute standard of length allowing the distant comparisons Weyl wanted to prevent in 1918.[35] The other possibility was left: (3) with a quantum version of (5),[36] of which the simplest and most obvious[37] was

$$(7) \qquad\qquad \psi \mapsto \psi' = e^{i\lambda}\psi,$$

---

[32]On the logical priority of relativity (or spin) *cf.* Weyl (1931b) p. 193: "Da die Möglichkeit einer solchen relativitätsinvarianten Gleichung für ein skalares $\psi$ nicht vorhanden ist, erscheint *der spin als ein durch die Relativitätstheorie notwendig gefordertes Phänomen*."

[33]But here Planck's constant $h$ and the speed of light $c$—and even charge—are set equal to one.

[34]Weyl (1929c) p. 284: "By this new situation, which introduces an atomic radius into the field equations themselves—but not until this step—my principle of *gauge-invariance*, with which I had hoped to relate gravitation and electricity, is robbed of its support." Weyl (1931a) p. 55: "Die Atomistik gibt uns ja absolute Einheiten für alle Maßgrößen an die Hand. [...] So geht in die DIRACsche Feldgesetze des Elektrons die „Wellenlänge des Elektrons", die Zahl $h/mc$, als eine absolute Konstante ein. Damit fällt das Grundprinzip meiner Theorie, das Prinzip von der Relativität der Längenmessung, dem Atomismus zum Opfer und verliert seine Überzeugungskraft."

[35]See also Weyl (1929c) p. 290.

[36]Weyl (1929c) p. 284: "this principle has an equivalent in the quantum-theoretical field equations which is exactly like it in formal respects; the laws are invariant under the simultaneous replacement of $\psi$ by $e^{i\lambda}\psi$, $\varphi_\alpha$ by $\varphi_\alpha - \partial\lambda/\partial x_\alpha$ where $\lambda$ is an arbitrary real function of position and time."

[37]The conservation requirement $\|\psi'\|^2 = \|\psi\|^2$ being very natural. And transformation (7) isn't even observable (with respect to position at any rate); *cf.* Weyl (1928) p. 87.

where $\mathbb{U}(1)$ replaced the multiplicative group $\mathbb{R}$ of (5).[38] As $\psi$ was now part of a four-dimensional space-time theory, it could no longer obey the Schrödinger equation, which violates relativity by treating space and time very differently.[39] Weyl adopted what amounted to a Dirac equation,[40] but cut in half: deprived of mass and the associated crisscrossing of component pairs ...

We now have four logical ingredients:

1. GR: *general relativity*

2. GJ: *geometrical justice*

3. MW: *matter wave*

4. SC: *second clock effect*;

W29 ← W18 & MW & SC & ?

A final ingredient, EL: *twice too many energy levels*, will almost be enough to produce the second gauge theory.

## 3.2 Dirac-Weyl theory

We can take $H = p_1^2$ as the simplified Hamiltonian of a particle whose mass is one-half. Momentum $p$ in quantum mechanics is represented by differentiation, in the sense that[41]

(8) $$p \mapsto i\,d,$$

in components $p_\mu \mapsto i\partial_\mu$. Our quantum Hamiltonian will therefore be

$$-\partial_1^2 = -\left(\frac{\partial}{\partial x^1}\right)^2,$$

which means that Schrödinger's equation $i\partial_t \psi = \partial_1^2 \psi$ differentiates space twice as much as time. But by what should it be replaced? The d'Alembertian $\square = \partial_0^2 -$

---

[38] Weyl (1931a) p. 55: "In dem theoretischen Weltbild bedeutet die Verwandlung von $f_p$ in $-f_p$ eine objektive Änderung des metrischen Feldes; denn es ist etwas anderes, ob sich eine Strecke bei kongruenter Verpflanzung längs einer geschlossenen Bahn vergrößert oder verkleinert. Nach dem angenommenen Wirkungsgesetz aber ist die Entscheidung über das Vorzeichen der $f_p$ auf Grund der beobachteten Erscheinungen unmöglich. Hier enthält darum, in Widerstreit mit einem oben ausgesprochenen erkenntnistheoretischen Grundsatz, das theoretische Weltbild eine Verschiedenheit, welche sich auf keine Weise für die Wahrnehmung aufbrechen läßt." P. 57: "Die an der alten Theorie gerügte Unsicherheit des Vorzeichens $\pm f_p$ löst sich dadurch in das unbestimmte Vorzeichen der $\sqrt{-1}$ auf. Schon damals, als ich die alte Theorie aufstellte, hatte ich das Gefühl, daß der Eichfaktor die Form $e^{i\lambda}$ haben sollte; nur konnte ich dafür natürlich keine geometrische Deutung finden. Arbeiten von SCHRÖDINGER und F. LONDON stützten die Forderung durch die allmählich sich immer deutlicher abzeichnende Beziehung zur Quantentheorie." See also Weyl (1931b) p. 89. Scholz (2004) p. 193 associates the 'geometry to matter' transition from (3)-(5) to (3)-(7) with a transition from the *a priori* fantasies of 1918 to the sober empiricism of 1929.

[39] Weyl (1931b) pp. 187-8: "Es ist klar, daß man zu einer befriedigenden Theorie des Elektrons nur kommen wird, wenn es gelingt, das Grundgesetz seiner Bewegung in der von der Relativitätstheorie geforderten, gegenüber Lorentz-Transformationen invarianten Form zu fassen."

[40] See Scholz (2006) p. 470.

[41] See Weyl (1931b) p. 89.

$\partial_1^2 - \partial_2^2 - \partial_3^2$ and Klein-Gordon equation $(\Box - m^2)\psi = 0$ treat space about the same way as time, they have the right transformation properties; but $\Box$ is 'squared' and there are reasons to prefer a wave operator and especially a time derivative[42] that aren't. In seeking a square root $\sqrt{\Box}$ Dirac found $\partial\!\!\!/ = \gamma^\mu \partial_\mu$, where the $\gamma^\mu$'s have the algebraic properties needed to get rid of the cross terms that appear when squaring. He therefore proposed the *Dirac equation*[43]

$$(9) \qquad\qquad (m - i\partial\!\!\!/)\psi = 0$$

which not only treats the three spatial derivatives $\gamma^k \partial_k$ the same way as the time derivative $\gamma^0 \partial_0$, but differentiates with respect to time only once.[44] The $\gamma^\mu$'s, which do not commute, cannot be numbers; they admit for instance the canonical representations

$$(10) \qquad\qquad \gamma^0 \leftrightarrow \left( \begin{array}{cc} 0 & \sigma^0 \\ -\sigma^0 & 0 \end{array} \right) \qquad \gamma^k \leftrightarrow \left( \begin{array}{cc} 0 & \sigma^k \\ \sigma^k & 0 \end{array} \right),$$

where all four quaternions $\sigma^\mu : \mathbb{C}^2 \to \mathbb{C}^2$ are hermitian and unitary; $\sigma^0$ is the identity, and the three traceless operators $\sigma^k$ satisfy $2i\sigma^j = \varepsilon_{jkl}[\sigma^k, \sigma^l]$.

The wave $\psi$ on which the $\gamma^\mu$'s act therefore has four (complex) components—*embarras de richesses* which Weyl found most troubling : "doppelt zu viel Energieniveaus"! The anti-diagonality of the $\gamma^\mu$'s governs the embarrassing excess by swapping the two two-spinors making up $\psi$. As the embarrassment is due to the *sign* that distinguishes between the different interweavings[45] produced by the $\gamma^\mu$'s, Weyl deals with it by choosing the only mass—none at all—that doesn't distinguish between plus and minus.[46] Without mass and half the components, (9) becomes $\sigma^\mu \partial_\mu \zeta = 0$. The reduced wave $\zeta$ has two complex components $\zeta_1, \zeta_2$ but four ('lightlike'[47]) real ones: the squared length $x_0 = \|\zeta\|^2 = \bar\zeta \sigma^0 \zeta$ and the three Hermitian quadratic forms $x_k = \bar\zeta \sigma^k \zeta$.

[42]Weyl (1931b) p. 188: "Sie ist nicht im Einklang mit dem allgemeinen Schema der Quantenmechanik, welches verlangt, daß die zeitliche Ableitung nur in der ersten Ordnung auftritt." P. 193: "Legt man die de Brogliesche Wellengleichung für das skalare $\psi$ zugrunde, in welche die elektromagnetischen potentiale $[A_\mu]$ durch die Regel [(11)] eingeführt sind, so ergibt sich aber für die elektrische Dichte ein Ausdruck, der außer $\psi$ die zeitliche Ableitung $\partial\psi/\partial t$ enthält und nichts mit der Ortswahrscheinlichkeit zu tun hat; sein Integral ist überhaupt keine Einzelform. Dies ist nach *Dirac* das entscheidendste Argument dafür, daß die Differentialgleichungen des in einem elektromagnetischen Feld sich bewegenden Elektrons von 1. Ordnung in bezug auf die zeitliche Ableitung sein müssen."

[43]Dirac (1928)

[44]Weyl (1931b) p. 190: "Nach dem allgemeinen Schema der Quantenmechanik sollte, wie schon erwähnt, die Differentialgleichung für $\psi$ von 1. Ordnung hinsichtlich der zeitlichen Ableitung von $\psi$ sein. Gemäß dem Relativitätsprinzip kann sie aber dann auch nur die 1. Ableitungen nach den räumlichen Koordinaten enthalten."

[45]Symplectic for time but simply 'NOT' for space. The interweaving produced by a purely NOT $\gamma^0$ would be gratuitous; the symplecticity given by the sign difference is essential—with respect to the three $\gamma^k$'s with merely NOT anti-diagonality.

[46]Weyl (1929c) p. 292: "The [mass] term (5) of the Dirac theory is, however, more doubtful. It must be admitted that if we retain it we can obtain all details of the line spectrum of the hydrogen atom—of one electron moving in the electrostatic field of a nucleus—in accord with what is known from experiment. But we obtain twice too much; if we replace the electron by a particle of the same mass and positive charge $+e$ (which admittedly does not exist in nature) the Dirac theory gives, contrary to all reason and experience, the same energy terms as for a negative electron, except for a change in sign. Obviously an essential change is here necessary." P. 294: "Be bold enough to leave the term involving mass entirely out of the field equations."

[47]The condition $x_0^2 = x_1^2 + x_2^2 + x_3^2$ applies.

We now have five logical ingredients:

1. GR: *general relativity*

2. GJ: *geometrical justice*

3. MW: *matter wave*

4. SC: *second clock effect*

5. EL: *twice too many energy levels*;

W18 & MW & SC & EL → W29.[48] The foundations are in place, the rest will follow.

## 3.3 Tetrads and spinors

The need for tetrads follows from the presence of spinors (introduced, as we have seen, to satisfy relativistic transformation requirements).

As long as there were only tensors

$$B = B^{\mu\nu\cdots}_{\sigma\cdots}\partial_\mu \otimes \partial_\nu \otimes \cdots \otimes dx^\sigma \otimes \cdots,$$

their components $B^{\mu\nu\cdots}_{\sigma\cdots}$ and the oblique—but holonomic—frames $\partial_\mu$ and $dx^\sigma$ with respect to which they were represented were subject to appropriate representations of $\mathbb{GL}(4,\mathbb{R})$. But this is too comprehensive a group to ensure the preservation of scalar products required by (the spinors representing) the new arrival, the matter wave.[49] The two complex components of a normalised spinor $\phi$ are subject to a $2 \times 2$ complex representation of $\mathbb{SU}(2)$; the three real components $\bar{\phi}\sigma^k\phi$ to a $3 \times 3$ real representation of $\mathbb{SO}(3)$; the four real components $\bar{\zeta}\sigma^\mu\zeta$ of the arbitrarily long spinor $\zeta$ to a $4 \times 4$ real representation of $\mathbb{SO}^+(1,3)$. Weyl acts on the two complex components of $\zeta$ with a $2 \times 2$ complex representation of a group[50] one can call

$$\mathbb{W}(2,\mathbb{C}) = \{g \in \mathbb{GL}(2,\mathbb{C}) : |\det g| = 1\} = \mathbb{SL}(2,\mathbb{C}) \times \mathbb{U}(1),$$

which is the extension[51] of $\mathbb{SL}(2,\mathbb{C})$ by the phase transformations.

---

[48] W18 & MW & SC give something like Dirac-Maxwell theory in curved space-time.

[49] Weyl (1929c) p. 285: "The tensor calculus is not the proper mathematical instrument to use in translating the quantum-theoretic equations of the electron over into the *general theory of relativity*. Vectors and terms are so constituted that the law which defines the transformation of their components from one Cartesian set of axes to another can be extended to the most general linear transformation, to an affine set of axes. That is not the case for the quantity $\psi$, however; this kind of quantity belongs to a representation of the rotation group which cannot be extended to the affine group. Consequently we cannot introduce components of $\psi$ relative to an arbitrary coordinate system in general relativity as we can for the electromagnetic potential and field strengths. We must rather describe the metric at a point $P$ by local Cartesian axes $e(\alpha)$ instead of by the $g_{pq}$. The wave field has definite components $\psi_1^+, \psi_2^+; \psi_1^-, \psi_2^-$ [full Dirac theory] relative to such axes, and we know how they transform on transition to any other Cartesian axes in $P$. The laws shall naturally be invariant under arbitrary rotation of the axes in $P$, and the axes at different points can be rotated independently of each other; they are in no way bound together."

[50] Weyl (1929b) p. 333: "man beschränke sich auf solche lineare Transformationen $U$ von $\psi_1, \psi_2$, deren Determinante den absoluten Betrag 1 hat."

[51] See Scholz (2004) p. 189, Scholz (2011b), third page of the paper, and footnote 5 above. Here I am indebted to Julien Bernard, who spotted a mistake.

Spinors are therefore transformed in such a way as to preserve their nature; but tensors are subject to (diffeomorphisms, their components to) coordinate changes, whose tangent maps acting on $\partial_\mu$ or $dx^\mu$ belong to $\mathbb{GL}(4,\mathbb{R})$.[52] To adapt to the spinorial requirements which impose $\mathbb{SO}^+(1,3)$, Weyl introduced the (anholonomic) tetrads $\mathbf{e}^\mu$—his *Achsenkreuze*—whose orthonormality is preserved by the Lorentz group.[53] The *Achsenkreuze* and their orthonormality are more primitive than the metric

$$g = \eta_{\mu\nu}\mathbf{e}^\mu \otimes \mathbf{e}^\nu = \eta_{\mu\nu}(\Lambda^\mu_\sigma \mathbf{e}^\sigma) \otimes (\Lambda^\nu_\tau \mathbf{e}^\tau),$$

which derives from them by construction; the (second) equality holds for all transformations $\Lambda \in \mathbb{SO}^+(1,3)$.

## 3.4   The displacement of spinors

Spinors have directions, whose covariant differentiation and parallel propagation have to take account of the space-time curvature of the region crossed. We have seen that in 1929 Weyl abandons the congruent transport of his first gauge theory, to return to the old Levi-Civita transport of Einstein's theory, generated by the Lie algebra

$$\mathfrak{o}(1,3) = \mathrm{Lie}\,\mathbb{SO}^+(1,3).$$

The spin connection $\omega$ is a one-form which, applied to a direction $\dot\gamma \in T_a M$, yields a generator $\omega_{\dot\gamma} = \langle \omega, \dot\gamma \rangle : \mathbb{C}^2_a \to \mathbb{C}^2_b$ of transport belonging to the Lie algebra $\mathfrak{w}(2,\mathbb{C}) = \mathrm{Lie}\,\mathbb{W}(2,\mathbb{C})$. Subtracting the difference $\delta\psi_b = \omega_{\dot\gamma}\psi_a$ from the initial spinor $\psi_a$ we obtain the transported spinor $\psi_b = (\mathbb{1}_2 - \omega_{\dot\gamma})\psi_a$. But numbers are easier to compare than spinors at different points, so it is best to take components $\psi_x^m = \langle \varphi_x^m, \psi_x \rangle$:

$$\psi_b^m = \psi_a^m - \omega_{\mu n}^m \dot\gamma^\mu \psi_a^n = \psi_a^m - (\omega_{\dot\gamma})_n^m \psi_a^n,$$

$\varphi_x^m$ being a basis in the dual fiber $\mathbb{C}^{2*}_x$, $m = 1, 2$.

Even if the gauge group $\mathscr{G}$ given by (7) changes the direction of the whole wavefunction in Hilbert space, the structure group $\mathbb{U}(1)$ only changes the phase of the spinor, not its direction. But generalisation to a non-Abelian[54] structure group that also changes the *directions* of spinors is natural, and was accomplished by Yang & Mills (1954) as we'll see in §4.

---

[52]Here I am indebted to an anonymous referee, and to Jean-Philippe Nicolas.

[53]Weyl (1931b) p. 195: "Ferner bedarf man in der allgemeinen Relativitätstheorie an jeder Weltstelle $P$ eines aus vier Grundvektoren in $P$ bestehenden normalen Achsenkreuzes, um die Metrik in $P$ festzulegen und relativ dazu die Wellengröße $\psi$ durch ihre vier [full Dirac theory again] Komponenten $\psi_\varrho$ beschreiben zu können; die gleichberechtigten normalen Achsenkreuze in einem Punkte gehen durch die Lorentztransformationen auseinander hervor."

[54]Even if the structure group $\mathbb{U}(1)$ makes $\mathsf{W29}$ an Abelian gauge theory, we have seen that parallel transport against its curved background requires $\mathbb{W}(2,\mathbb{C})$ as well; as an electromagnetic gauge theory it is Abelian, but it is also a non-Abelian theory of spinors on curved space-time. One might say it is Abelian with respect to electromagnetism, non-Abelian with respect to gravity.

### 3.5 Three gauge arguments

#### 3.5.1 The inherited inexact connection

To reach, from Weyl's first gauge theory, the compensation of (7) by (3), with[55]

$$(11) \qquad\qquad d \mapsto D = d + iA$$

(or $\partial_\mu \mapsto D_\mu = \partial_\mu + iA_\mu$) we only need a handful of principles. In §3.1 we saw that Einstein's objection, supported by the absolute length $h/mc$, favours (7) over (5). The addition of the electromagnetic potential to momentum (and hence to the derivative) comes from analytical mechanics, where[56]

$$(12) \qquad\qquad p \mapsto p + A$$

(or $p_\mu \mapsto p_\mu + A_\mu$). Together (8) and (12) give (11), the one-form $A$ being the one that figures, by a natural identification, in (3). The compensation of (7) by (3) can be seen in the Lagrangian

$$\mathscr{L} = \bar\psi \sigma^\mu D_\mu \psi = \bar\psi' \sigma^\mu (D_\mu - i\partial_\mu \lambda)\psi'.$$

The inexact connection $A$ and its nonvanishing curvature (1) were there long before it even made sense to apply (7).

I'd say there's something of a 'gauge argument' here already. But Weyl has another gauge argument,[57] which extracts electromagnetism from the $\mathbb{U}(1)$ freedom left by the $h : \mathbb{W}(2,\mathbb{C}) \to \mathbb{SO}^+(1,3)$ homomorphism and expressed by

$$(13) \qquad\qquad h(e^{i\lambda}g) = h(g) \in \mathbb{SO}^+(1,3),$$

$g \in \mathbb{W}(2,\mathbb{C})$.[58]

#### 3.5.2 How Weyl extracts an inexact connection from $\mathbb{U}(1)$

$\mathbb{SO}^+(1,3) = G$ and $\mathbb{W}(2,\mathbb{C}) = G'$ are just 'structure' groups, acting at a generic space-time point. What about the corresponding gauge groups $\mathscr{G}$, $\mathscr{G}'$ acting on all of space-time $M$? In special relativity "there's just a single tetrad"; so there's just

---

[55] See Weyl (1929c) p. 283, Weyl (1931b) p. 89.

[56] See Weyl (1931b) p. 88.

[57] Weyl (1929b) p. 348, Weyl (1929c) p. 291, Afriat (2013)

[58] Weyl (1929c) p. 291: "It is my firm conviction that we must seek the origin of the electromagnetic field in another direction. We have already mentioned that it is impossible to connect the transformations of the $\psi$ in a unique manner with the rotations of the axis system; however we may attempt to accomplish this by means of invariants which can be used as constituents of an action quantity we always find that there remains an arbitrary "gauge factor" $e^{i\lambda}$. Hence the local axis-system does not determine the components of $\psi$ uniquely, but only within such a factor of absolute magnitude 1." Weyl (1931b) p. 195: "Aus der Natur, dem Transformationsgesetz der Größe $\psi$ ergibt sich, daß die vier Komponenten $\psi_\varrho$ relativ zum lokalen Achsenkreuz nur bis auf einen gemeinsamen Proportionalitätsfaktor $e^{i\lambda}$ durch den physikalischen Zustand bestimmt sind, dessen Exponent $\lambda$ willkürlich vom Orte in Raum und Zeit abhängt, und daß infolgedessen zur eindeutigen Festlegung des kovarianten Differentials von $\psi$ eine Linearform $\sum_\alpha f_\alpha dx_\alpha$ erforderlich ist, die so mit dem Eichfaktor in $\psi$ gekoppelt ist, wie es das Prinzip der Eichinvarianz verlangt."

one $\mathbb{SO}^+(1,3) = G = \mathcal{G}$, one $\mathbb{W}(2,\mathbb{C}) = G' = \mathcal{G}'$, and above all one $e^{i\lambda}$.[59] But with space-time curvature the tetrad varies,[60] and so does $\lambda$. This could mean the following:[61] Only a *flat* $\mathfrak{o}(1,3)$-valued connection $\mathfrak{A}$ allows the assignment of the *same* tetrad to distant points—only with flatness can there be *global* constancy or 'sameness.' With curvature it becomes meaningless to say that tetrads at distant points are the same. Where tetrads cannot remain constant, one has to suppose they *vary*. A flat real-valued phase connection $A$ alongside a curved $\mathfrak{A}$ can of course be countenanced, but it is in the spirit of Weyl's argument for both to be flat or both curved. So if the tetrad varies, $\lambda$ might as well too.[62]

The group homomorphism $h$ determines the Lie algebra homomorphism

$$\mathfrak{h} : \mathfrak{w}(2,\mathbb{C}) \to \mathfrak{o}(1,3),$$

where the Lie algebra $\mathfrak{w}(2,\mathbb{C}) = \mathrm{Lie}\,\mathbb{W}(2,\mathbb{C})$ is the direct sum $\mathfrak{sl}(2,\mathbb{C}) \oplus i\mathbb{R}\mathbb{1}_2$, and $i\mathbb{R} = \mathrm{Lie}\,\mathbb{U}(1)$. Doing away with the additive freedom $\lambda$ (or rather $i\lambda\mathbb{1}_2$) we're left with the isomorphism between $\mathfrak{w}(2,\mathbb{C})/i\mathbb{R}\mathbb{1}_2 = \mathfrak{sl}(2,\mathbb{C})$ and $\mathfrak{o}(1,3)$. Instead of the phase $e^{i\lambda} \in \mathbb{U}(1)$ we have $i\lambda\mathbb{1}_2 \in i\mathbb{R}\mathbb{1}_2$; instead of $\mathbb{U}(1)$ we have the Lie algebra $i\mathbb{R}\mathbb{1}_2$; and instead of (13),

$$\mathfrak{h}(\gamma \oplus i\lambda\mathbb{1}_2) = \mathfrak{h}(\gamma) \in \mathfrak{o}(1,3),$$

$\gamma \in \mathfrak{w}(2,\mathbb{C})$.[63]

The additive freedom $i\lambda\mathbb{1}_2$ is in the Lie algebra $\mathfrak{w}(2,\mathbb{C})$ where the spin connection has its values; and connections are there to generate parallel transport—*in a direction*.[64] A direction $\dot\gamma \in T_a M$ will therefore characterise the propagation of $\lambda$, whose

---

[59]Weyl (1929b) p. 348: "In der speziellen Relativitätstheorie muß man diesen Eichfaktor als eine Konstante ansehen, weil wir hier ein einziges, nicht an einen Punkt gebundes Achsenkreuz haben." Weyl (1929c) p. 291: "In the special theory of relativity, in which the axis system is not tied up to any particular point, this factor is a constant."

[60]The gauge groups become infinite-dimensional. Weyl (1929b) p. 348: "Anders in der allgemeinen Relativitätstheorie: jeder Punkt hat sein eigenes Achsenkreuz und darum auch seinen eigenen willkürlichen Eichfaktor; dadurch, daß man die starre Bindung der Achsenkreuze in verschiedenen Punkten aufhebt, wird der Eichfaktor notwendig zu einer willkürlichen Ortsfunktion." Weyl (1929c) p. 291: "But it is otherwise in the general theory of relativity when we remove the restriction binding the local axis-systems to each other; we cannot avoid allowing the gauge factor to depend arbitrarily on position."

[61]Here I am indebted to Johannes Huisman.

[62]*Cf.* Ryckman (2009) p. 295: "Weyl's argument for his correct conclusion is, in fact, flawed, resting on an unnecessary assumption about the representation of spinor matter fields within tetrad formulations of arbitrarily curved space-times."

[63]Weyl (1929b) p. 348: "Dann ist aber auch die infinitesimale lineare Transformation $dE$ der $\psi$, welche der infinitesimalen Drehung $d\gamma$ entspricht, nicht vollständig festgelegt, sondern $dE$ kann um ein beliebiges rein imaginäres Multiplum $i \cdot df$ der Einheitsmatrix vermehrt werden." Weyl (1929c) p. 291: "Then there remains in the infinitesimal linear transformation $dE$ of $\psi$, which corresponds to the given infinitesimal rotation of the axis-system, an arbitrary additive term $+id\varphi \cdot 1$."

[64]Weyl (1929b) p. 348: "Zur eindeutigen Festlegung des kovarianten Differentials $\delta\psi$ von $\psi$ hat man also außer der Metrik in der Umgebung des Punktes $P$ auch ein solches $df$ für jedes von $P$ ausgehende Linienelement $\overrightarrow{PP'} = (dx)$ nötig. Damit $\delta\psi$ nach wie vor linear von $dx$ abhängt, muß

$$df = f_P(dx)^p$$

eine Linearform in den Komponenten des Linienelements sein. Ersetzt man $\psi$ durch $e^{i\lambda}$, so muß man sogleich, wie aus der Formel für das kovariante Differential hervorgeht, $df$ ersetzen durch $df - d\lambda$." Weyl

13

infinitesimal variation $\delta\lambda$ has to be linear in $\lambda$ and in $\dot{\gamma}$. The object needed is a one-form; applied to the direction $\dot{\gamma}$ it yields the infinitesimal generator $\langle A, \dot{\gamma}\rangle \in \mathbb{R}$, which then multiplies $\lambda$ to produce the increment $\delta\lambda = \lambda\langle A, \dot{\gamma}\rangle$. So there's a connection for tetrads, another for spinors, *and a third one—A—for the residual* $\mathbb{U}(1)$ *freedom caught 'in between' tetrads and spinors.*

The whole point of allowing the propagation of $\lambda$ to depend on direction is to admit anholonomies. So the curvature (1) of $A$ will not necessarily vanish. In (1), $A$ and (2) Weyl again[65] saw[66] the electromagnetic field, its potential and Maxwell's two homogeneous equations.[67]

### 3.5.3 How the standard gauge argument extracts an exact connection from $\mathbb{U}(1)$

An alternative logic,[68] which is claimed to produce electromagnetism from the indifference of $\mathscr{L}$ to (7), is popular: The local phase transformation gives rise to a new Lagrangian

$$\mathscr{L}' = \bar{\psi}'\sigma^\mu\partial_\mu\psi' = \bar{\psi}e^{-i\lambda}\sigma^\mu(e^{i\lambda}\partial_\mu + e^{i\lambda}i\partial_\mu\lambda)\psi = \bar{\psi}\sigma^\mu(\partial_\mu + i\partial_\mu\lambda)\psi.$$

As the components $i\partial_\mu\lambda$ of $id\lambda$ account for the difference, invariance is restored once the same term is subtracted, thus producing the covariant derivative $D' = d - id\lambda$ and the invariant Lagrangian

$$\hat{\mathscr{L}} = \bar{\psi}'\sigma^\mu D'_\mu\psi'.$$

It is then argued that an interaction $F = d^2\lambda$ is thereby deduced,[69] whose potential $A$ is $d\lambda$. But since $d^2$ vanishes the interaction does too, as has often been pointed out.[70]

## 4 Yang-Mills theory

Here the structure group $\mathbb{SU}(N)$ replaces $\mathbb{U}(1)$. Weyl is no longer in the foreground, nor is his complaint that Dirac's theory had *twice too many energy levels*. The curved

---

(1929c) p. 291: "The complete determination of the covariant differential $\delta\psi$ of $\psi$ requires that such a $d\varphi$ be given. But it must depend linearly on the displacement $PP'$: $d\varphi = \varphi_p(dx)^p$, if $\delta\psi$ shall depend linearly on the displacement. On altering $\psi$ by multiplying it by the gauge factor $e^{i\lambda}$ we must at the same time replace $d\varphi$ by $d\varphi - d\lambda$ as is immediately seen from this formula of the covariant differential." Weyl's notation is confusing: whereas the one-form $d\lambda$ (which *is* a differential) is necessarily exact, $df$ and $d\varphi$ (my $A$) aren't.

[65] See §2.1.

[66] Weyl (1929b) p. 349, Weyl (1929c) pp. 291-2

[67] *Cf.* Ryckman (2009) p. 295: "Weyl derived the Maxwell equations from the requirement of local phase invariance, thus coupling charged matter to the electromagnetic field, and so originating the modern understanding of the principle of local gauge invariance ("*local symmetries dictate the form of the interaction*") that lies at the basis of contemporary geometrical unification programs in fundamental physics.

[68] See for instance Yang & Mills (1954) p. 192, Sakurai (1967) p. 16, Aitchison & Hey (1982) p. 176, Mandl & Shaw (1984) p. 263, Göckeler & Schücker (1987) p. 48, Ramond (1990) pp. 183-91, Ryder (1996) p. 93, O'Raifeartaigh (1997) p. 118.

[69] Ryder (1996) p. 95: "the electromagnetic field arises *naturally* by demanding invariance of the action [...] under *local* ($x$-dependent) rotations [...]."

[70] Auyang (1995) p. 58, Brown (1999) pp. 50-3, Teller (2000) pp. S468-9, Lyre (2001, 2004a,b), Healey (2001) p. 438, Martin (2002) p. S229, Martin (2003) p. 45

space-time from which Weyl's first theory arose can now, having done its bit,[71] be kept or dropped.

Let us go back to (7), which is indeed a natural choice to replace (5). But is it the *only* natural choice? The transformation on the Hilbert space $\mathscr{H}$ containing $\psi$ should of course be unitary, but there's a more general unitary transformation.

One thinks of the function $\lambda$ as 'real-valued': it assigns a real number to every $x \in M$. Since the wavefunction $\psi$ assigns not a complex number but a spinor $\psi \in \mathbb{C}_x^N$ to every $x$, the value $\lambda(x)$ is in fact the operator $\lambda(x) \cdot \mathbb{1}_N : \mathbb{C}_x^N \to \mathbb{C}_x^N$. But then why not take a *general* Hermitian operator $\Lambda(x) : \mathbb{C}_x^N \to \mathbb{C}_x^N$ rather than the very special Hermitian operator $\lambda(x) \cdot \mathbb{1}_N$? *Why stop halfway?* Legitimate question, which is enough to yield Yang-Mills theory. The structure group $\mathbb{SU}(N)$ being unitary, the operator $\mathscr{U} : \mathscr{H} \to \mathscr{H}$ (the 'direct integral' of all the $e^{i\Lambda(x)} \in \mathbb{SU}_x(N)$) representing the corresponding gauge group remains unitary.

In 1929 Weyl would have seen no *physical* reason to take the step from $\mathbb{U}(1)$ to $\mathbb{SU}(N)$. Though given to mathematico-physical speculation of uninhibited virtuosity, he didn't take the purely mathematical step either. The details of the physics that ultimately did produce the non-Abelian theory are in Yang & Mills (1954).

A glance at the formalism: Instead of the connection $A = A_\mu dx^\mu$ with values in the Lie algebra $i\mathbb{R}$ we have the Yang-Mills connection (6), with values in the Lie algebra $\mathfrak{su}(N)$ spanned by $\mathbf{T}_1, \ldots, \mathbf{T}_N$. Applied to a transport direction $\dot{\gamma} \in T_a M$, the connection $\mathcal{A}$ gives the infinitesimal generator

$$\langle \mathcal{A}, \dot{\gamma} \rangle = \sum_{\mu=0}^{3} \mathcal{A}_\mu \dot{\gamma}^\mu = i \sum_{\mu=0}^{3} \sum_{k=1}^{N} \mathcal{A}_\mu^k \langle dx^\mu, \dot{\gamma} \rangle \mathbf{T}_k = i\mathbf{T} : \mathbb{C}_a^N \to \mathbb{C}_b^N$$

which turns the initial spinor $\psi_a \in \mathbb{C}_a^N$ into the increment $\delta\psi_b = i\mathbf{T}\psi_a \in \mathbb{C}_b^N$. In components we have

$$\langle \varphi_b^m, \psi_b \rangle = \langle \varphi_a^m, \psi_a \rangle - i \sum_{n=1}^{N} \langle \varphi_b^m, \mathbf{T}_{ab}\varphi_{an} \rangle \langle \varphi_a^n, \psi_a \rangle,$$

in other words $\psi_b^m = \psi_a^m - i\mathbf{T}_n^m \psi_a^n$, where $\varphi_x^m$ is a basis in the dual fiber $\mathbb{C}_x^{N*}$, $m = 1, \ldots, N$. The curvature[72]

$$\mathcal{F} = d\mathcal{A} = i \sum_{\mu=0}^{3} \sum_{k=1}^{N} d\mathcal{A}_\mu^k \wedge dx^\mu \otimes \mathbf{T}_k = i\frac{1}{2} \sum_{\mu,\nu=0}^{3} \sum_{k=1}^{N} \mathcal{F}_{\mu\nu}^k dx^\mu \wedge dx^\nu \otimes \mathbf{T}_k$$

is a two-form with values in $\mathfrak{su}(N)$.

---

[71]The *geometrical justice* of §2.1 required a curved length connection $A$ to balance the curved directional connection. By adopting a flat space-time connection alongside a curved isospin connection Yang & Mills (1954) reversed the injustice of Einstein's theory—which has a curved directional connection and a flat length connection.

[72]For simplicity it is all concentrated in the coefficients $\mathcal{A}_\mu^k$, or rather $\mathcal{F}_{\mu\nu}^k$, and kept out of the basis fields $\mathbf{T}_k$—which (topology permitting) can be assumed integrable: here $(d\mathbf{T}_k)_\mu^\nu = \partial_\mu \mathbf{T}_k^\nu - \partial_\nu \mathbf{T}_k^\mu$ vanishes, $k = 1, \ldots, N$.

# 5 Logical summary

Summing up, Weyl's first gauge theory W18 was given by *geometrical justice* GJ applied to *general relativity* GR:

$$\text{GR \& GJ} \rightarrow \text{W18}.$$

To reach Weyl's second gauge theory W29, *matter wave* MW, *second clock effect* SC and *twice too many energy levels* EL were needed too:

$$\text{W18 \& MW \& SC \& EL} \rightarrow \text{W29}.$$

To obtain Yang-Mills theory YM from W29, *non-Abelian structure group* NA was enough:

$$\text{W29 \& NA} \rightarrow \text{YM},$$

or more precisely

$$\text{W18 \& MW \& SC \& NA} \rightarrow \text{YM}.$$

Weyl had the greatest creative freedom in 1918, when he applied *geometrical justice* to general relativity. The next moves were more constrained. In introducing a *matter wave* after the discoveries of Schrödinger *et al*. he had little choice; and it had to be relativistic, hence with spin—which led to the use of tetrads. Weyl's preference for (3)-(7) over (3)-(5) was dictated by Einstein's objection, the *second clock effect*. His reaction to the *twice too many energy levels* was less constrained, but also less right, less consequential, more idiosyncratic. The adoption of a *non-Abelian structure group* was mathematically so natural as to be almost inevitable; but ultimately the step was not taken for purely mathematical reasons, and as a physical move it seems more creative, less constrained.

# References

Afriat, A. (2009) "How Weyl stumbled across electricity while pursuing mathematical justice" *Studies in History and Philosophy of Modern Physics* **40**, 20-5

Afriat, A. (2013) "Weyl's gauge argument" *Foundations of Physics* **43**, 699-705

Aitchison, I. J. R. and A. J. G. Hey (1982) *Gauge theories in particle physics*, Hilger, Bristol

Auyang, S. Y. (1995) *How is quantum field theory possible?*, Oxford University Press

Brading, K. (2002) "Which symmetry? Noether, Weyl, and the conservation of electric charge" *Studies in History and Philosophy of Modern Physics* **33**, 3-22

Brading, K. and E. Castellani (editors) (2003) *Symmetries in physics: philosophical reflections*, Cambridge University Press

Broglie, L. de (1924) *Recherches sur la théorie des quanta*, Thèse, Paris

Brown, H. (1999) "Aspects of objectivity in quantum mechanics" pp. 45-70 in J. Butterfield and C. Pagonis (editors) *From physics to philosophy*, Cambridge University Press

Cao, T. (1997) *Conceptual developments of 20th century field theories*, Cambridge University Press

Coleman, R. and H. Korté (2001) "Hermann Weyl: mathematician, physicist, philosopher" pp. 161-388 in Scholz (2001b)

Dirac, P. A. M. (1925) "The fundamental equations of quantum mechanics" *Proceedings of the Royal society A* **109**, 642-53

Dirac, P. A. M. (1928) "The quantum theory of the electron" *Proceedings of the Royal society A* **117**, 610-24

Einstein, A. (1916) "Grundlage der allgemeinen Relativitätstheorie" *Annalen der Physik* **49**, 769-822

Gambini, R. and J. Pullin (2011) *A first course in loop quantum gravity*, Oxford University Press

Göckeler, M. and T. Schücker (1987) *Differential geometry, gauge theories, and gravity*, Cambridge University Press

Hawkins, T. (2000) *Emergence of the theory of Lie groups*, Springer, Berlin

Healey, R. (2001) "On the reality of gauge potentials" *Philosophy of Science* **68**, 432-55

Healey, R. (2007) *Gauging what's real: the conceptual foundations of contemporary gauge theories*, Oxford University Press

Levi-Civita, T. (1917) "Nozione di parallelismo in una varietà qualunque e conseguente specificazione geometrica della curvatura riemanniana" *Rendiconti del Circolo matematico di Palermo*" **42**, 173-205

Lyre, H. (2001) "The principles of gauging" *Philosophy of Science* **68**, S371-81

Lyre, H. (2004a) *Lokale Symmetrien und Wirklichkeit: eine Naturphilosophische Studie über Eichtheorien und Strukturenrealismus*, Mentis, Paderborn

Lyre, H. (2004b) "Holism and structuralism in $U(1)$ gauge theory" *Studies in History and Philosophy of Modern Physics* **35**, 643-70

17

Mandl, F. and G. Shaw (1984) *Quantum field theory*, Wiley, Chichester

Martin, C. (2002) "Gauge principles, gauge arguments and the logic of nature" *Philosophy of Science* **69**, S221-34

Martin, C. (2003) "On continuous symmetries and the foundations of modern physics" pp. 29-60 in Brading & Castellani (2003)

O'Raifeartaigh, L. (1997) *The dawning of gauge theory*, Princeton University Press

O'Raifeartaigh, L. and N. Straumann (2000) "Gauge theory: historical origins and some modern developments" *Reviews of Modern Physics* **72**, 1-23

Pais, A. (1982) *'Subtle is the Lord . . . ': the science and the life of Albert Einstein*, Oxford University Press

Pauli, W. (1921) *Relativitätstheorie*, Teubner, Leipzig

Pauli, W. (1979) *Wissenschaftlicher Briefwechsel, Band I: 1919-1929*, Springer, Berlin

Ramond, P. (1990) *Field theory: a modern primer*, Westview Press, Boulder

Ryckman, T. (2003a) "Surplus structure from the standpoint of transcendental idealism: the "world geometries" of Weyl and Eddington" *Perspectives on Science* **11**, 76-106

Ryckman, T. (2003b) "The philosophical roots of the gauge principle: Weyl and transcendental phenomenological idealism" pp. 61-88 in Brading & Castellani (2003)

Ryckman, T. (2005) *The reign of relativity: philosophy in physics 1915-1925*, Oxford University Press

Ryckman, T. (2009) "Hermann Weyl and "first philosophy": constituting gauge invariance" pp. 279-98 in M. Bitbol *et al.* (editors) *Constituting objectivity: transcendental perspectives on modern physics*, Springer Netherlands

Ryder, L. (1996) *Quantum field theory*, Cambridge University Press

Sakurai, J. J. (1967) *Advanced quantum mechanics*, Addison-Wesley, Reading

Scholz, E. (1994) "Hermann Weyl's contributions to geometry in the years 1918 to 1923" pp. 203-30 in J. Dauben *et al.* (editors) *The intersection of history and mathematics*, Birkhäuser, Basel

Scholz, E. (1995) "Hermann Weyl's "Purely Infinitesimal Geometry"" pp. 1592-1603 in S. D. Chatterji (editor) *Proceedings of the International congress of mathematicians, August 3-11, 1994 Zürich*, Birkhäuser, Basel

Scholz, E. (2001a) "Weyls Infinitesimalgeometrie, 1917-1925" pp. 48-104 in Scholz (2001b)

Scholz, E. (editor) (2001b) *Hermann Weyl's* Raum-Zeit-Materie *and a general introduction to his scientific work*, Birkhäuser, Basel

Scholz, E. (2004) "Hermann Weyl's analysis of the "problem of space" and the origin of gauge structures" *Science in Context* **17**, 165-97

Scholz, E. (2005) "Local spinor structures in V. Fock's and H. Weyl's work on the Dirac equation (1929)" pp. 284-301 in D. Flament *et al.* (editors) *Géométrie au vingtième siècle, 1930-2000*, Hermann, Paris

Scholz, E. (2006) "Introducing groups into quantum theory" *Historia mathematica* **33**, 440-90

Scholz, E. (2011a) "Mathematische Physik bei Hermann Weyl – zwischen „Hegelscher Physik" und „symbolischer Konstruktion der Wirklichkeit"" pp. 183-212 in K.-H. Schlote and M. Schneider (editors) *Mathematics meets physics: a contribution to their interaction in the 19th and the first half of the 20th century*, Harri Deutsch Verlag, Frankfurt

Scholz, E. (2011b) "H. Weyl's and E. Cartan's proposals for infinitesimal geometry in the early 1920s" *Boletim da Sociedada portuguesa de matemàtica*, **Numero especial A**, 225-45

Schrödinger, E. (1926) "Quantisierung als Eigenwertproblem (erste Mitteilung)" *Annalen der Physik* **79**, 361-76

Seelig, K. (1960) *Albert Einstein*, Europa Verlag, Zurich

Sigurdsson, S. (2001) "Journeys in spacetime" pp. 15-47 in Scholz (2001b)

Straumann, N. (1987) "Zum Ursprung der Eichtheorien bei Hermann Weyl" *Physikalische Blätter* **43**, 414-21

Teller, P. (2000) "The gauge argument" *Philosophy of Science* **67**, S466-81

Vizgin, V. (1984) *Unified field theories*, Birkhäuser, Basel

Weyl, H. (1918a) "Gravitation und Elektrizität" pp. 147-59 in *Das Relativitätsprinzip*, Teubner, Stuttgart, 1990

Weyl, H. (1918b) "Reine Infinitesimalgeometrie" *Mathematische Zeitschrift* **2**, 384-411

Weyl, H. (1921) "Feld und Materie" *Annalen der Physik* **65**, 541-63

Weyl, H. (1926) *Philosophie der Mathematik und Naturwissenschaft*, Oldenbourg, Munich

Weyl, H. (1928) *Gruppentheorie und Quantenmechanik*, Hirzel, Leipzig

Weyl, H. (1929a) "Gravitation and the electron" *Proceedings of the National academy of sciences, USA* **15**, 323-34

Weyl, H. (1929b) "Elektron und Gravitation" *Zeitschrift für Physik* **56**, 330-52

Weyl, H. (1929c) "Gravitation and the electron" *The Rice Institute Pamphlet* **16**, 280-95

Weyl, H. (1931a) "Geometrie und Physik" *Die Naturwissenschaften* **19**, 49-58

Weyl, H. (1931b) *Gruppentheorie und Quantenmechanik* (second edition), Hirzel, Leipzig

Weyl, H. (1988) *Raum Zeit Materie*, Springer, Berlin

Yang, C. N. and R. Mills (1954) "Conservation of isotopic spin and isotopic gauge invariance" *Physical Review* **96**, 191-5