

# Scientific Theories as Intervening Representations<sup>1</sup>

Andoni IBARRA, Thomas MORMANN

BIBLID [0495-4548 (2006) 21: 55; pp. 21-38]

**ABSTRACT:** In this paper some classical representational ideas of Hertz and Duhem are used to show how the dichotomy between representation and intervention can be overcome. More precisely, scientific theories are reconstructed as complex networks of intervening representations (or representational interventions). The formal apparatus developed is applied to elucidate various theoretical and practical aspects of the *in vivo/in vitro* problem of biochemistry. Moreover, adjoint situations (Galois connections) are used to explain the relation between empirical facts and theoretical laws in a new way.

**Key words:** Representation, adjoint situations, *in vitro/in vivo* problem, Hertz, Duhem.

## 1. Introduction

The concept of representation has not yet been secured on the agenda of philosophy of science. Some philosophers flatly deny that it could be of any use in epistemology or philosophy of science. Instead, they claim, the concept of representation leads us into a hopeless maze of pseudo-questions without answers. This is the case of Rorty and his antirepresentationalist followers. According to them, epistemology based on the notions of negotiation and interpretation should replace epistemological accounts based on ‘representation’. In this paper, we will not address this kind of radical anti-representationalism. But suffice to say, it is based on a rather primitive conception of representation identifying representation with some kind of copying or mirroring.

In this paper we want to elaborate some classical representational ideas of Hertz and Duhem in order to show that a diagrammatical or combinatorial account of representations can be useful for elucidating the role of representations in describing the practice of representational reasoning in science.

The outline of this paper is as follows: in section 2, we outline some ideas of Hertz and Duhem concerning the structure of scientific reasoning that can be used to understand how representations in science work. More precisely, following Hertz the idea of a commutative diagram of interconnected representations is introduced, and Duhem’s account of empirical theories will lead us to the idea that the theoretical and the empirical are correlated in a so called adjoint situation. In section 3, the rudiments of a combinatorial theory of representations are introduced, and are put to use in section 4 for the representational elucidation of the *in vitro/in vivo* problem in biochemistry. In section 5, it is shown that Duhem’s account of an empirical theory as a cor-

---

<sup>1</sup> We would like to thank two anonymous referees for their detailed and penetrating criticisms that helped us to correct some major blunders in the first version of this paper. Further we’d like to express our sincere gratitude to the guest editor José Antonio Díez, who pointed out some conceptual obscurities and infelicitous formulations in the original text. Of course, we are responsible for all remaining errors.



relation of symbolical and empirical facts leads to the conception of an empirical theory as a Galois connection (or, more generally, an adjoint situation) in the sense of mathematical category theory. We close with some general remarks on the role of representational concepts in philosophy of science.

## 2. Classical Ideas of Representations

Let us start with some basic ideas on scientific representations put forward by the classical philosopher-scientists Hertz and Duhem. These ideas naturally lead towards an interesting conception of scientific theories as representations.

As our first classical intuition pump for the development of a comprehensive account of representation, we take Hertz's well known 'symbolical account' put forward in his *The Principles of Mechanics presented in a New Form* (Hertz 1894) where he described the general procedure of scientific representations as follows:

We form for ourselves images or symbols of external objects; and the form which we give them is such that the necessary consequents of the images in thought are always the images of the necessary consequents in nature of the things pictured. In order that this requirement may be satisfied, there must be certain conformity between nature and our thought. Experience teaches us that the requirement can be satisfied, and hence that such a conformity does in fact exist. [...]

The images, which we may form of things, are not determined without ambiguity by the requirement that the consequents of the images must be the images of the consequents. Various images of the same objects are possible, and the images may differ in various respects. [...]

Of two images of equal distinctness the more appropriate is the one which contains, in addition to the essential characteristics, the smaller number of superfluous or empty relations, —the simpler of the two. Empty relations cannot be altogether avoided: they enter into the images because they are simply images ... produced by our mind and necessarily affected by the characteristics of its mode of portrayal. (Hertz 1894, pp. 1f.)

We propose to translate Hertz's informal description of the representational activity of science in a diagrammatical language as follows: let the set of "external objects" be denoted by  $E$ , and denote the set of "images" by  $S$ . The following diagram may be used to capture the essential structure of Hertz's account:

$$\begin{array}{ccc}
 E & \xrightarrow{t} & S \\
 f \downarrow & & \downarrow g \\
 E & \xrightarrow{t} & S
 \end{array}$$

The details are as follows: the horizontal arrow  $t$  corresponds to Hertz's formation of mental images. More precisely, if  $e \in E$  is an external object,  $t(e) \in S$  is the image corresponding to it. In other words,  $t(e)$  may be considered as the theoretical counterpart of  $e$ . The left vertical arrow  $f$  in Hertz's diagram is to be conceived as a process or an experiment that 'necessarily' brings about the external fact that  $e$  is changed to another

external fact  $f(e) \in E$ . In Hertz's terms,  $f(e)$  is the 'necessary consequent' of  $e$ . Analogously, the vertical arrow  $g$  on the right may be interpreted as a mathematical calculation or a logical argument that leads from a 'symbol'  $s \in S$  to another symbol  $g(s)$ . It is to be interpreted as the result or the conclusion of the symbolical transaction  $g$ . In Hertz's terms,  $g(s)$  is the 'necessary consequent' of  $s$ . These ingredients of Hertz's diagram are of course not independent of each other; rather, as is informally stated in his *Principles*, they form a commutative diagram in 'that necessary consequents of the images in thought are always the images of the necessary consequents in nature', which in our diagrammatical language just amounts to the commutativity of the diagram:

(2.1) *Commutativity of Hertz's Diagram.* Assume  $t, f$ , and  $g$  as characterized above. They are assumed to satisfy the following concatenation law:

$$g \bullet t = t \bullet f$$

This equation is to be interpreted as follows: If we start with an empirical fact  $e$  in the left upper corner of the Hertz diagram, translate it to its theoretical counterpart  $t(e)$ , and use  $t(e)$  as the input for a calculation or a logical argument that leads to  $g \bullet t(e)$ , then this outcome is the same as if we had submitted the empirical fact  $e$  to an experimental transformation  $f$  arriving at  $f(e)$ , and translated this experimental fact  $f(e)$  by  $t$  finally yielding  $t \bullet f(e) = g \bullet t(e)$ . In other words, the two paths in Hertz's diagram are strictly equivalent in that they may be considered as paths that lead to one and the same destination. As an elementary example consider  $e$  to be some chemical substance that is submitted to a certain chemical experiment  $f$  which, say, oxidizes  $e$  thereby yielding as outcome another chemical substance  $f(e)$ . For this transaction a chemical theory has to provide a chemical formula  $t(e)$  for  $e$ , and a theoretical transformation  $g(t(e))$  of  $t(e)$  such that  $t(f(e)) = g(t(e))$ . As is emphasized by Hertz, given  $E$  there may be different 'symbolic completions'  $S, S'$ . The choice between them is a pragmatic matter of simplicity and local usefulness. It may be that for different purposes different 'images' may be appropriate (cf. Hertz 1894, p. 3).

Second, let us come to Duhem's contribution to a modern representational account of scientific theorizing, which is found in his classic *The Aim and Structure of Physical Theory* (Duhem 1906). At various occasions in his *opus magnum* he asserts that scientific theories are to be conceived as representations. More precisely, he considers a physical theory 'as an economical representation' that

establishes an order and a classification among [the experimental laws]. It brings some laws together, closely arranged in the same group; it separates some others by placing them into two groups very far apart. Theory gives, so to speak, the table of contents and the chapter headings under which the science to be studied will be methodologically divided. (Duhem 1906, pp. 23f).

Later he goes on to explain this 'representation' as a correspondence between 'practical facts' and 'theoretical' or 'symbolical facts'. It is certainly not too far fetched to consider Duhem's account as presented up to now as just another version of Hertz's structural approach. But there is one feature in Duhem's representational ap-

proach that is novel and not present in Hertz. In describing a physical theory as a correspondence between practical and symbolical facts he insists that

a symbolic formula ... can be translated into concrete facts in an infinity of different ways, *because all these disparate facts admit the same theoretical interpretation.* (*Ibid.*, p. 150)

And, in an analogous vein:

The same practical fact may correspond to an infinity of logically incompatible theoretical facts; the same group of concrete facts may be made to correspond in general not with a single symbolic judgment but with an infinity of judgments different from one another and logically in contradiction with one another. (*Ibid.*, p. 152)

Duhem's account is rather informal, and he is not very clear about what is to be understood by 'theoretical fact'. In particular, one should not interpret him as conceiving a 'theoretical fact' as a fact 'belonging' to a specific theory. Rather, the most appropriate interpretation of Duhemian theoretical facts is to take a theoretical fact as one that asserts a physical state of affairs in precise mathematical terms, as is explained by Duhem. A typical example of a theoretical fact (or statement) is the following: 'An increased pressure of 100 atmospheres causes the electromotive force of a given gas battery to increase by 0.0844 volts.' (*Ibid.*, p. 152) Other 'logically incompatible' theoretical statements would be obtained by replacing '0.0844' by '0.0845' or '0.0846'. Hence, Duhem's account of an empirical theory can be formulated in relational terms as follows:

(2.2) *Duhem's Relational Account of Empirical Theories.* Denote the class of symbolic facts by  $S$  and the class of practical or empirical facts by  $E$ . Then a theory  $T$  is to be conceived as a relation

$$T \subseteq E \times S.$$

If  $(e, s) \in T$  then this is to be interpreted as the empirical fact that  $e$  is related to  $s$ , or, to put it the other way round, that the symbolic fact  $s$  is related to the empirical fact  $e$ .

It is important to note that Duhem insisted that this relation is multi-valued: to a single  $e$  there may correspond many symbolic facts  $s$ , and, vice versa, to a single  $s$ , there may correspond many empirical facts  $e$ . This double ambiguity of the relation between empirical and symbolical facts is characteristic of Duhem's account and has no counterpart in Hertz's approach. As we shall show in the next section, this feature may be combined with the representational insights of Hertz to yield a complex representational account of empirical theories.

### 3. Representational Combinatorics

Following Hertz and Duhem in conceiving the practice of science as engaged in producing and manipulating representations of various kinds, the impression that comes to mind is that scientific representations do not live in isolation, rather they may be combined and concatenated in various ways (Ibarra, Mormann 2000). Hence, investigating these combinatorial aspects of representations is a central task of a general theory of representation (Ibarra, Mormann 1997 a, b).

Regardless of what kind of representations we consider, they are not unconnected with each other, rather, they form a representational network. One and the same entity  $A$  may be represented by several different entities  $B, C, D$  etc. such that we have representations  $A \xrightarrow{r} B, A \xrightarrow{s} C, A \xrightarrow{t} D$ , etc. On the other hand, it may happen that one and the same entity  $E$  appears as the representative of several different entities  $A, B, C$  etc. That is to say we have representations  $A \longrightarrow E, B \longrightarrow E, C \longrightarrow E$ . Furthermore, it can be the case that representations such as  $A \xrightarrow{r} B$  and  $B \xrightarrow{s} C$  are concatenated yielding an indirect or combined representation  $A \xrightarrow{s \cdot r} C$ .

As the result of these considerations, we can see that any theory of representations should comprise a combinatorial part, which describes the various possibilities of combinations and iterations of representations. In the following we shall assume that this combination or concatenation of representations is associative, i.e. representations  $f, g$  and  $h$ , which 'match', satisfy the following law of associativity:

$$(3.1) \quad f \cdot (g \cdot h) = (f \cdot g) \cdot h.$$

The combination or iteration of representations is of utmost importance for the practice of science. For instance, in the standard representational theory of measurement the numerical measurement of an empirical domain  $D$  is conceptualized as a representation  $r: D \longrightarrow \mathfrak{R}$  of  $D$  into the real numbers  $\mathfrak{R}$ . This is a rather idealized description. Actually, by a closer inspection the representation  $D \xrightarrow{r} \mathfrak{R}$  should be regarded as a more or less extended chain of representations

$$(3.2) \quad D \longrightarrow E \longrightarrow F \longrightarrow \dots \longrightarrow \mathfrak{R}.$$

In most cases, numerical or, more generally, mathematical representations of empirical data cannot be 'read off' directly; usually they have to be considered as constructs which have been built by a more or less complicated constructional processes. The long way from data to theory shows that the standard dichotomy is, at best, a very idealized picture. Dealing with an example from general relativity theory, Laymon gives a detailed account of the 'long contrafactual path from data to theory' (cf. Laymon 1982). Other examples of complex 'long distance' representations are discussed in detail in Latour (1999): Latour tells us in detail the long story from raw findings to theoretically digestible data in the case of 'botanical pedology' (*ibid.*, chapter 2). Notwithstanding important differences, all these accounts rely—implicitly or explicitly—on what may be called a combinatorics of representations.

The combination of representations is not restricted, however, to linear combinations. As will be shown by the *in vivo/in vitro* example of biochemistry, the point of the combinatorial account of representations only comes to the fore if we do not restrict

our attention to linear chains of representations but, instead, also take into account non-linear net-like configurations of representations.

The importance of representational nets or diagrams is evidenced by the fact that in the last forty years or so mathematics (and parts of other sciences as well) has been successfully reformulated in terms of representational networks. Here we refer, of course, to the mathematical theory of categories founded by Eilenberg and Mac Lane in the forties and presented for the general scientific public in books such as Mac Lane's *Mathematics – Form and Function* (1986) or Lawvere and Schanuel's *Conceptual Mathematics – A First Introduction to Categories* (1996). In category theory, representations appear under the names morphisms, functors, and natural transformations. In the last decades it has been shown that not only the bulk of mathematics can be reconstructed in these terms, but also that this representational reconstruction has led to new and fruitful lines of mathematical research. We take this fact, together with the representational ideas of Hertz and Duhem as evidence that combinations of various kinds of representations play an indispensable role for a representational theory of scientific knowledge. This claim is substantiated in the next section in which we propose to study in some detail various combinations of representations that arise from the so-called *in vitro/in vivo* problem in biochemistry.

#### 4. *A Representational Account of the In Vivo/In Vitro Problem*

In this section we are going to apply the formal apparatus sketched so far to a specific problem of a scientific discipline that up to now has not received too much attention from philosophy of science, to wit, the so called '*in vitro/in vivo*' problem of biochemistry (cf. Strand, Fjelland, and Flatmark 1996 and Strand 1999). For the information on biochemical matters we heavily rely on these papers. Our purpose is to show that the rudiments of a theory of meaningful representations set out in the previous sections may be used to elucidate the problems of the representational practice biochemistry have to cope with. We chose the approach of Strand *et al.* as our starting point since it seemed to us particularly well suited for our purposes: on the one hand, it is sufficiently complex to require the employment of some non-trivial representational tools; on the other hand, it is conceptually not too complex as to be inaccessible for non-experts in biochemistry.

First, let us recall the basic ingredients of the *in vitro/in vivo* problem as it presents itself in biochemistry. The first point to note is that although biochemistry may be defined as 'the field of science concerned with the chemical substances and processes that occur in plants, animals, and microorganisms' it would be misleading to assume that 'biochemists study processes that occur in living organisms' (cf. Strand 1999, p. 273). The reason is that normally

it is impossible to perform a chemical analysis of an intact organism. A biochemical analysis is typically preceded by an isolation procedure, in which the organism of interest is disrupted and a specific component of it is isolated. To put another way, almost all biochemical evidence is obtained *in vitro* under artificial experimental conditions. ... [Nevertheless] biochemists are concerned with the chemistry of the living organism, *in vivo*. (Strand 1999, p. 273)

Hence one may even assert:

It would be wrong to say biochemists *observe* or *describe* or *study* processes that occur in living organisms, because they very rarely do so. Normally, it is impossible to perform a chemical analysis of an intact organism. (*Ibid.*, p. 273).

Almost all biochemical evidence is obtained *in vitro*, under artificial experimental conditions.

An *in vivo* system is a biologically interesting but experimentally inaccessible system, and the corresponding *in vitro* system as a related accessible, but biologically less interesting system. Although analogous situations also occur in other sciences, the difference between *in vitro* and *in vivo* systems is particularly striking in biochemistry. Now we may define the *in vitro/in vivo problem* (or the IVIV-problem henceforth) as the problem of justifying knowledge claims about *in vivo* systems on the basis of evidence obtained in ‘corresponding’ *in vitro* systems. Or, on a more descriptive level, the IVIV-problem may be said to be the problem of describing as clearly as possible the various methods used by biochemists to extract the information on *in vivo* systems they are seeking from the evidence they have obtained from *in vitro* systems.

One has to note that the IVIV distinction is a relative distinction. That is to say, in one context a system may play the role of the *in vitro* part, in another context the same system may be considered as the *in vivo* component. Of course, one may say that in many sciences one finds analogous distinctions to that of the IVIV distinction in biochemistry. Nevertheless, the case of biochemistry is special since the IVIV is thus central for this discipline, as is convincingly pointed out by Strand (1999, pp. 274f). His discussion may be summed up in the contention that the concept of artifact is central in any biochemical discussion. A biochemical artifact is a chemical reaction that occurs between biomolecules *in vitro*, but not *in vivo*. Now, the problem of artifacts is a central problem of meaningful representations everywhere. Given a representation  $A \xrightarrow{r} B$  the problem arises to interpret elements and relations defined on  $B$  in terms of  $A$ , for instance, if  $r(a) = r(a')$  one may ask if this identity on the representing domain  $B$  may be pulled back to  $A$ , i.e., one asks if it is possible to infer  $a = a'$ . Consider first the special case that  $r$  is a function. Then, of course, this inference is not valid in general. It is only admissible to infer from  $r(a) \neq r(a')$  that  $a \neq a'$ . Now let us consider the general case that the representation  $r$  is any relation between  $A$  and  $B$ . Denote the power set of  $B$  by  $PB$ . Then, committing an innocent abuse of language,  $r$  may be conceived as a function  $A \xrightarrow{r} PB$  defined by  $r(a) := \{b; (a, b) \in r\}$ . In the same vein as above one can infer from  $r(a) \neq r(b)$  that  $a \neq b$ . In other words, the inequality on  $B$  (or  $PB$ ) can be pulled back to an inequality on  $A$ . Of course, the problem of artifacts is not restricted to this kind of artifacts. Other, more complicated relations on  $B$  such as  $R(f(a_1), \dots, R(f(a_n)))$  may be considered and tested for their  $A$ -meaningfulness.

We take the fact that the problem of artifacts can be naturally couched in representational terms as evidence that the IVIV problem should be treated in terms of a theory of meaningful representations. Thus we propose to conceive the relation between

an *in vivo* system  $S$  and a corresponding *in vitro* system  $S^*$  as a representational relation  $S \xrightarrow{d} S^*$  contending that the *in vivo* system  $S$  is represented by the *in vitro* system  $S^*$ .

First, it should be noted that this representation is a material long distance representation par excellence: Usually the representing system  $S^*$  is obtained from  $S$  by a variety of massive, often destructive interventions of various kinds (cf. Strand *et al.* 1996, Strand 1999). The representing system  $S^*$  is far from being similar to  $S$ , and it is neither natural nor necessary to represent  $S$  by  $S^*$ . There may be many other ways of representing  $S$  by other  $S^*$ ,  $S^{**}$ , ... depending on the representational interests and capacities of those who are engaged in the construction of these *intervening* representations. Thus, as the first outcome of considering the IVIV problem in biochemistry we contend that the dichotomy between representing and intervening put forward by some philosophers such as Hacking is pointless in the case of biochemistry, and, regarding biochemistry as a paradigmatic case for science in general, for other sciences as well (cf. Hacking 1983).

As lucidly explained by Strand *et al.*, there is much more in the IVIV problem than the statement that it gives rise to an intervening representation  $S \longrightarrow S^*$ . To deal with these more fine-grained aspects of the IVIV problem, let us introduce the following terminological conventions: properties, objects, relations, procedures etc. belonging to the realm of *in vivo* systems are denoted by  $E, F, a, b, R, p, \dots$ , while the corresponding properties, objects, etc. belonging to the *in vitro* realm are denoted by  $E^*, F^*, a^*, b^*, \dots$ . Our first purpose is to show that IVIV problems give rise in a natural way to a plurality of Hertz's diagrams. Given systems  $S$  and  $S^*$ , and important task of the biochemist's work is to study how these systems behave under certain perturbations  $p$  and  $p^*$ . Here, a perturbation  $p$  of  $S$  may be considered as a map:  $S \xrightarrow{p} S$ . More precisely,  $p(s)$  is to be understood that for  $s \in S$  the state  $p(s) \in S$  is the state that resulted from  $s$  when submitted to the perturbation  $p$ . Analogously for *in vitro* states  $S^*$  and *in vitro* perturbations  $p^*$ :  $S^* \longrightarrow S^*$ . Then the systems and perturbations  $S, p, S^*, p^*$  may be said to be optimally correlated if the following Hertz diagram commutes:

$$\begin{array}{ccc}
 S & \xrightarrow{d} & S^* \\
 p \downarrow & & \downarrow p^* \\
 S & \xrightarrow{d} & S^*
 \end{array}
 \quad p^* \bullet d = d \bullet p.$$

By definition, an artifact is an *in vitro* perturbation  $d(s) \neq p^*(d(s))$  such that  $s = p(s)$ . If the Hertz diagram commutes, artifacts can be shown not to exist: Assume  $d(s) \neq p^*(d(s))$  and  $s = p(s)$ . From Hertz we get  $p^*(d(s)) = d(p(s))$ . Hence we get the following proposition:

*Proposition 1.* If Hertz commutes, then there are no artifacts.



In a similar vein, one obtains that the non-existence of artifacts implies that the Hertz diagram commutes for states  $s$  that are invariant under the perturbation  $p$ , i.e., states for which  $s = p(s)$ :

*Proposition 2.* If  $s$  is invariant under  $p$  AND there are no artifacts, then HERTZ commutes for  $s$ .

*Proof.* Assume  $s = p(s)$ . Then  $d(s) = d(p(s))$ . Assume that HERTZ does not commute for  $s$ . That is to say  $p^*(d(s)) \neq d(p(s))$ . Then  $p^*(d(s)) \neq d(s)$ . Since there are no artifacts one infers  $s \neq p(s)$ . This is a contradiction. ■

In sum, the diagrammatically natural requirement that Hertz diagrams commute is a bit stronger than the claim that no artifacts exist. The existence of artifacts is, however, not the only problem that may arise when studying the relation between *in vivo* and *in vitro* systems. It may well happen that the combination of *in vitro* perturbation  $p^*: S^* \longrightarrow S^*$  and the intervening representation  $d: S \longrightarrow S^*$  are jointly too invasive and too coarse, such that a salient *in vivo* perturbation  $p$  fails to be detected by them. This is the case if it happens that  $s \neq p(s)$  but  $d(s) = p^*(d(s))$ . This may be called an artificial null effect. Artificial null effects and the commuting of the Hertz diagram are related as follows:

*Proposition 3.* If the Hertz diagram commutes and the representation  $d: S \longrightarrow S^*$  is mono, i.e.,  $d(a) = d(b)$  implies  $a = b$ , then no artificial null effects occur. ■

In this implication, the second clause of the antecedent is clearly necessary. This may be more conspicuously expressed by contraposition:

*Proposition 4.* If artificial null effects occur, then either the Hertz diagram does not commute or the IVIV representation  $d: S \longrightarrow S^*$  is not mono. ■

One may ask whether the converse holds: If no artificial null effects occur, does the Hertz diagram commute and is  $d$  mono? As is easily checked by examples, this is not the case. In other words, the conjunctive assumption that the Hertz diagram is commutative and the IVIV representation  $d$  is mono is strictly stronger than the non-existence of artificial null effects.

As has been pointed by Strand *et al.*, the IVIV problem is not completely described by a Hertz diagram connecting an *in vivo* systems  $S$  and an *in vitro* systems  $S^*$ . Usually these systems are accompanied by what may be called their model systems  $M$  and  $M^*$  respectively. That is to say, for the *in vivo* system  $S$  there is a theoretical (or maybe sometimes a computer) model  $M$ , and for the *in vitro* system  $S^*$  there is a theoretical (computer model) model  $M^*$ . Then it is natural to assume that  $M$  is an appropriate representation of  $S$ , and  $M^*$  is an appropriate representation of  $S^*$ . These may be ex-

plicated by the assumption that the representations  $S \xrightarrow{t} M$  and  $S^* \xrightarrow{t} M^*$  have Hertz diagrams of the following kind:

$$(4.1) \quad \begin{array}{ccc} S & \xrightarrow{t} & M \\ \downarrow & & \downarrow \\ S & \xrightarrow{t} & M \end{array} \quad \begin{array}{ccc} S^* & \xrightarrow{t} & M^* \\ \downarrow & & \downarrow \\ S^* & \xrightarrow{t} & M^* \end{array}$$

For each of these diagrams one may study the various ways in which artifacts may influence the reliability of surrogative reasoning dealing with  $M$ ,  $S^*$ , and  $M^*$  and finally bound to obtain information about the *in vivo* system  $S$ .

For dealing in a reasonable way with problems of this kind it is not sufficient, however, to assume that Hertz diagrams for  $(S, S^*)$ ,  $(S, M)$ , and  $(S^*, M^*)$  exist. One has to assume the existence of a further ‘purely theoretical’ Hertz diagram for  $(M, M^*)$  such that the following ‘3-dimensional’ or ‘cubical’ diagram commutes:

$$(4.2) \quad \begin{array}{ccccc} & & S^* & \xrightarrow{t^*} & M^* \\ & \nearrow & \vdots & \nearrow & \downarrow \\ S & \xrightarrow{t} & M & \nearrow & \\ \downarrow & & \downarrow & & \downarrow \\ S & \xrightarrow{t} & M & \nearrow & \\ & \searrow & S^* & \xrightarrow{t^*} & M^* \\ & \downarrow & \vdots & \downarrow & \\ & S & \xrightarrow{t} & M & \nearrow \end{array}$$

Of course, it can hardly be expected that in reality the cube (4.2) is fully commutative. Rather, there will exist various sources of non-commutativity, which show that the various kinds of systems and models only match approximately. Nevertheless, the cube presentation (4.2) may be useful as an idealized model to spot where precisely commutativity and thereby the validity of surrogative reasoning via models and systems of various kinds may fail.

Let us consider a particularly simple theoretical model of (*in vivo* or *in vitro*) systems, which, at first, may not appear as models at all. Assume that for a given system  $S$  the possible states  $s$  of  $S$  may have certain properties. This assumption may be cast in a representational framework by stipulating that there is a map  $F: S \longrightarrow C$ ,  $C$  being a structure whose elements are to be interpreted as properties belonging to a certain property type. In other words,  $F(s)$ ,  $s \in S$ , is to be conceived as the assertion that the state  $s$  has the property  $F(s)$ . Given a perturbation  $p: S \longrightarrow S$  one may ask, if  $F$  is in-

variant with respect to  $p$ , i.e., if  $F(p(s)) = F(s)$ , or not. Analogously, for *in vitro* properties  $F^*: S^* \longrightarrow C^*$  a corresponding *in vitro* perturbation  $p^*: S^* \longrightarrow S^*$  is defined that may or may not be invariant under the *in vitro* property  $F^*: S^* \longrightarrow C^*$ . The properties  $F$  and  $F^*$  are correlated by  $d$  and  $d_c$  iff there is a commutative diagram of the following kind:

$$\begin{array}{ccc}
 \begin{array}{c} \curvearrowright \\ S \end{array} & \xrightarrow{F} & C \\
 \downarrow d & & \downarrow d_c \\
 S^* & \xrightarrow{F^*} & C^* \\
 \begin{array}{c} \downarrow \\ \curvearrowright \\ p^* \end{array} & & 
 \end{array}$$

This diagram describes the (ideal) relation between *in vivo* properties  $F$  and *in vitro* properties  $F^*$ . From now on, let us assume that  $F$  and  $F^*$  are such that there exist  $d$  and  $d_c$  so that the diagram commutes. This means that  $F$  and  $F^*$  are reasonably correlated with each other. This is to ensure that assertions dealing with  $F^*$  may be possibly translated into assertions dealing with  $F$ , that is to say that  $F$  and  $F^*$  can be correlated by surrogate reasoning. On this non-trivial assumption about  $F$  and  $F^*$  is based the entire *in vivo/in vitro* argumentation.

Usually, the domain of values  $C$  of a property  $F$  is not just a set, but has some structure. For instance, often  $C$  is assumed to be endowed with an order relation  $\leq$ . Then we may define an order relation on  $S$  by pulling back the order defined on  $C$  by the following definition:

$$s \leq s' := F(s) \leq F(s').$$

In an analogous way the state space  $S^*$  of an *in vitro* system may be endowed with an order via a map  $S^* \xrightarrow{F^*} C^*$  of  $S^*$  into an ordered property space  $C^*$ . The *in vivo* property  $F$  is stable under the *in vivo* perturbation  $p: S \longrightarrow S$  iff  $s \leq s'$  implies  $p(s) \leq p(s')$ , i.e., iff  $F(s) \leq F(s') \Rightarrow F(p(s)) \leq F(p(s'))$ . Analogously for *in vitro* perturbation  $p^*$  and an *in vitro* feature  $F^*$ . Of course, the two properties  $F$  and  $F^*$  and the perturbations  $p$  and  $p^*$  should not be unrelated to each other. Rather, the *in vivo* property  $F$  and the *in vitro* property  $F^*$  should obey the following relation:

$$F(p(s)) \leq F(p(s')) \Rightarrow F^*(p^*(d(s))) \leq F^*(p^*(d(s'))).$$

In this case, from the accessible *in vitro* relation  $\text{NOT}(F^*(p^*(d(s))) \leq F^*(p^*(d(s'))))$  one can infer  $\text{NOT}(F(p(s)) \leq F(p(s')))$ .

Discussing the structural features of the IVIV problem shows that in the biochemical practice the concepts of representation and intervention are intimately related. More precisely, they are correlated by a variety of commutative diagrams that combine *in vivo* systems, *in vitro* systems, *in vivo* models, and *in vitro* models by a com-

plex net of representational and intervening links. The basic building block of this net, which intertwines theoretical representations and practical interventions are various kinds of Hertz diagrams. Thus, taking the IVIV problem in biochemistry as paradigmatical for empirical theories in general we contend that representations and interventions should be treated together, since both may be characterized as moves in the complex network of an empirical theory.

The IVIV problem of biochemistry is particularly interesting for a representational philosophy of science as it shows the necessity of considering iterations and combinations of various kinds of interventions and representations. The language of representational diagrams is particularly apt for dealing with the various kinds of connections. We think that the opposition between the representative and the performative perspective in philosophy of science is an artifact of a misinformed philosophy of science. One does not have to choose between them. Indeed, in some sense, every representation has an interventional aspect, at least indirectly, and every intervention leads to a representation.

### 5. *Adjoint Situations*

In this section we are going to show that Duhem's relational account of theories that conceives a theory  $T$  as a relation  $T \subseteq S \times E$  between symbolic and empirical facts may be elucidated by using so called Galois connections or, more generally, adjoint situations in the sense of category theory. This part of the paper is the most speculative one, and some readers may object that we introduce a heavy formal apparatus without real justification. Thus the following preliminary remark may be in order. Our point is this: conceiving an empirical theory as a certain relation between empirical and theoretical facts seems to us quite a natural and intuitive approach. Otherwise Duhem, who certainly was not interested in formal technicalities, would not have endorsed it. Now, as soon as a theory is given as a relation  $T \subseteq S \times E$ , the whole apparatus of Galois relations is available. One may even say that a Galois relation between  $PS$  and  $PE$  is nothing but a relation. Since Galois connections have turned out to be a useful tool in the study of binary relations in mathematics, computer science and elsewhere. Hence, one may suspect that they could do some useful work in formal philosophy of science as well. This conjecture is further supported by the fact that Galois connections are just a very special case of adjoint situations that may be characterized as *the* fundamental concept of category theory. Hence, there is some hope that these conceptual tools have some applications in philosophy of science as well.

By conceiving a theory as a relation  $T \subseteq E \times S$  of empirical and symbolical facts in the sense of Duhem's *The Aim and Structure of Physical Theory*, it is not claimed, of course, that any relation  $X \subseteq S \times E$  counts as a genuine theory. There are countless relations between the two classes of facts that make no sense at all. Further restrictions will have to be imposed on  $T$  in order that  $T$  can be acknowledged as a genuine theory. As will be shown later, for this task the representational ideas of Hertz turn out to be useful.

For the moment we only want to emphasize Duhem's main point, to wit, that for any given empirical fact  $f \in E$  there may be many symbolic facts  $s \in S$  such that  $f$  and  $s$  are theoretically correlated, i.e., that  $(f, s) \in T$ , and that vice versa for any  $s \in S$  there may be many empirical facts  $f \in E$  such that  $(f, s) \in T$  (cf. Duhem 1906, pp. 152ff). Formally, this means that  $T \subseteq E \times S$  is a relation and not a function.

This multivalued correlation between empirical and symbolical facts renders it plausible that a single fact, be it symbolical or empirical, hardly makes sense as such. That is to say, a single  $s \in S$  or  $f \in E$  is an object that in real science hardly occurs. Rather, what shows up in the practice of real science are *clusters* or *complexes* of empirical and theoretical facts. Thus, we propose to consider appropriate *sets*  $A \subseteq S$  and  $B \subseteq E$  as the real building blocks of scientific theories; single empirical facts  $f \in E$  or symbolic facts  $s \in S$  are auxiliary concepts introduced for methodological reasons. Replacing elements by subsets in this way is a natural generalization in so far as the 'elementary' facts of type  $s$  and  $f$  may be considered as special cases of facts of type  $A$  and  $B$  by identifying  $s$  and  $f$  with their singletons  $\{s\}$  and  $\{f\}$ . This technical move from elementary facts to subsets of elementary facts resembles the approach Duhem's Austrian colleague Ernst Mach proposed long time ago: according to Mach, it was the task of science to describe the functional relations of appropriate complexes or clusters of elements in the most economical way possible. In any case, the move from elements to subsets facilitates to get started the formal apparatus we are going to apply in order to elucidate Duhem's relational account of scientific theories. After these preparatory remarks we are now ready to set up the formal apparatus we need in order to cast Duhem's relational account of empirical theories in the framework of Galois connections. First, let us deal with the necessary technicalities.

Denote by  $PS$  and  $PE$  the power sets of  $S$  and  $E$ , respectively. For the moment, let us assume that  $PS$  and  $PE$  are endowed with their natural (set-theoretical) order structures  $(PS, \subseteq)$  and  $(PE, \subseteq)$ . A theory  $T \subseteq E \times S$  gives rise to order-preserving maps between  $PS$  and  $PE$  by the following recipe:

(5.1) *Proposition.* Let  $T \subseteq E \times S$  be a theory. Define maps  $PE \xrightarrow{t} PS$  and  $PS \xrightarrow{e} PE$  by:

- (a) For  $Y \in PE$  define  $e(Y)$  by  $e(Y) := \{s; \exists y (y \in Y \text{ AND } (y, s) \in T)\}$
- (b) For  $X \in PS$  define  $t(X)$  by  $t(X) := \{y; (e(\{y\}) \subseteq X)\}$ .

Then the maps  $e$  and  $t$  are order preserving.

*Proof.* Check the definitions of  $e$  and  $t$ .

Obviously,  $e$  and  $t$  are not unrelated to each other. Indeed, it can be shown that  $t$  is completely determined by  $e$ , and vice versa. Actually, much more is true, as is shown by the following proposition:

(5.2) *Proposition.* Let  $e$  and  $t$  be defined as above. Then for all  $X \subseteq S$  and  $Y \subseteq E$  the following holds:

$$X \subseteq t(Y) \text{ IFF } e(X) \subseteq Y.$$

In technical jargon, the ordered pair  $(t, e)$  is called a *Galois connection* between the order structures  $PS$  and  $PE$  (cf. Gierz *et al.* 2003). More precisely,  $t$  is called the upper (or right) adjoint, and  $e$  is called the lower (or left) adjoint. One should note that a Galois connection  $(t, e)$  is *not* a symmetric notion, i.e., if  $(t, e)$  is a Galois connection, usually  $(e, t)$  fails to be a Galois connection. The difference between upper and lower adjoint is reflected in the notational convention that  $t$  as the upper adjoint is on the right or ‘upper’ side of  $\leq$ , while  $e$  as the lower adjoint is on the ‘lower’ side of the order relation  $\leq$ . This asymmetry is essential in the following to set up an asymmetric relation between the domain of empirical facts  $E$  and the domain of symbolical facts  $S$ .

*Proof (5.2).* The proof naturally splits into two parts: (i) assume  $X \subseteq t(Y)$  and  $z \in e(X)$ . Then one has to show  $z \in Y$ . By definition of  $e(X)$  there is an  $s \in X$  with  $(s, z) \in T$ . That is to say  $z \in e(s)$ . By presupposition  $s \in t(Y)$ . This means  $e(s) \subseteq Y$ , and therefore  $z \in Y$ ; (ii) Assume  $e(X) \subseteq Y$  and  $s \in X$ . One has to show  $s \in t(Y)$ . But  $e(s) \subseteq e(X) \subseteq Y$ , and this just means  $s \in t(Y)$ .

(5.3) *Corollary.* The map  $PS \xrightarrow{e \bullet t} PS$  is a kernel operator, i.e.,  $e \bullet t(X) \subseteq X$ , for all  $X \subseteq S$ , and the map  $PE \xrightarrow{t \bullet e} PE$  is a closure operator, i.e.  $Y \subseteq t \bullet e(Y)$  for all  $Y \subseteq E$ .

After having presented these rudiments of the theory of Galois connections, let us start now with the task of elucidating the intuitive meaning of this gadget. This amounts to an interpretation of the components  $t$  and  $e$ , which form the Galois connection  $(t, e)$ , and an explanation of their most important properties in informal terms of philosophy of science.

For this task it is expedient to start with the map  $e: PS \longrightarrow PE$ . Recall that subsets  $X \subseteq S$  and subsets  $Y \subseteq E$  are to be interpreted as symbolic (theoretical) and empirical facts. By definition  $e(X)$  is the collection of all ‘atomic’ empirical facts  $z$  that are empirically correlated to at least one ‘atomic’ symbolic fact  $s \in X$ . This may be interpreted as that the empirical fact  $e(X)$  provides an empirical realization of  $X$  in a broad sense, i.e., it may be that the empirical facts  $z$  realizing the symbolic facts of  $X$  may have theoretical correlates  $s$  that do not belong to  $X$  but at least  $X$  is covered by the empirical facts of  $e(X)$  in the sense that  $t(e(X)) \supseteq X$ .

Analogously, the map  $t$  may be interpreted as a recipe to translate an empirical fact  $Y \subseteq E$  into a related theoretical fact  $t(Y)$  such that each theoretical fact  $s$  belongs to  $t(Y)$ , i.e.,  $s \in t(Y)$  if and only if all empirical correlates  $z$  of  $s$  belong to  $Y$ . In other

words,  $t(Y)$  is the most comprehensive theoretical fact for which  $Y$  provides a complete empirical realization.

We hasten to add that this relational account of empirical theories as a relation  $T \subseteq E \times S$  is seriously incomplete. Its essential flaw is that it does not allow us to distinguish between approximately true theories and false theories, i.e., theories that are completely off the mark. If a theory  $T$  is just a relation  $T \subseteq E \times S$  relating symbolic and empirical facts, there is no room for asking if  $T$  is (approximately) correct or not. This is clearly not sufficient to model the way of how theories relate theoretical facts to often recalcitrant empirical facts. To overcome this shortcoming, it is expedient to rely once more on the insights encapsulated in Hertz's diagram. In other words, we propose to combine the insights of Hertz and Duhem to obtain a better model of scientific theorizing that comprises the advantages of both the Hertzian and the Duhemian accounts.

This is done as follows: Let us start over again from the domains  $PS$  and  $PE$  of theoretical facts and symbolic facts, respectively, endowed with maps  $e: PS \xrightarrow{e} PE$  and  $PE \xrightarrow{t} PS$  as before. That is to say,  $e$  and  $t$  are to be interpreted as Duhemian maps correlating symbolic facts and empirical facts as explained above. The new ingredient we are going to introduce in order to distinguish between (approximately) true theories and those that are plainly false is provided by the replacement of the trivial set theoretical order relation  $\subseteq_S$  on  $S$  and  $\subseteq_E$  on  $E$  by appropriate non-trivial order relations  $\leq_S$  and  $\leq_E$  on  $PS$  and  $PE$ , respectively, which reflect some theoretical or empirical intervention and processes as explained in our discussion of the Hertz diagram in section 2. More precisely this is explained in the following definition:

(5.4) *Definition.* (a) Assume  $Y, Y^* \in PE$ . Assume that there is an empirical process  $P$  or intervention such that the empirical fact  $Y$  is the initial state  $P(i)$  of  $P$ , and  $Y^*$  is the final state  $P(f)$  of  $P$ . It is further assumed that processes or interventions  $P, P', P''$  can be concatenated associatively. Define  $Y \leq Y^* :=$  there is a process  $P$  with initial state  $Y$  and final state  $Y^*$ .

(b) Assume  $X, X^* \in PS$ . Assume that there is a symbolic process  $P$  or intervention such that the symbolic fact  $X$  is the initial state  $P(i)$  of  $P$ , and  $X^*$  is the final state  $P(f)$  of  $P$ . It is further assumed that processes or interventions  $P, P', P''$  can be concatenated associatively. Define  $X \leq X^* :=$  there is a process  $P$  with initial state  $X$  and final state  $X^*$ .

The class of processes or interventions defined for symbolic and empirical facts render  $PS$  and  $PE$  order structures, to be denoted by  $(PS, \leq_S)$  and  $(PE, \leq_E)$ , respectively. From now on,  $PS$  and  $PE$  are assumed to be endowed with these interventional orders which differ from the set-theoretical orders  $\subseteq_S$  and  $\subseteq_E$ . In Hertz's terms, then,  $X \leq X'$  is to read as ' $X$  is a necessary consequent of  $X'$ ', and analogously  $Y \leq Y'$  is to

be read as ‘ $Y$  is a necessary consequent of  $X$ ’. Then the following Duhem-Hertz requirement makes sense:

(5.5) *Definition.* Let  $T \subseteq E \times S$  be a relation of symbolical and empirical facts. Assume  $PE$  and  $PS$  endowed with interventional orders  $\leq_E$  and  $\leq_S$  respectively,  $X \in PS$ ,  $Y \in PE$ . Let  $PS \xrightarrow{e} PE$  and  $PE \xrightarrow{t} PS$  defined by  $T$ . Then the theory  $T$  is said to satisfy the Duhem-Hertz condition iff for all  $X \in PS$ ,  $Y \in PE$  the following equivalence holds:

$$(5.6) \quad e(X) \leq_E Y \text{ IFF } X \leq_S t(Y).$$

In other words, the pair  $(t, e)$  is a Galois connection between  $(PS, \leq_S)$  and  $(PE, \leq_E)$ . More precisely,  $t$  is the upper (or right) adjoint, and  $e$  is the lower (or left) adjoint of this Galois connection.

Before we explain in some detail why theories satisfying (5.6) should be considered as (approximately) true let us note that instead of the set theoretical structures  $PS$  and  $PE$  it may be more expedient, more intuitive, and even less clumsy to replace  $PS$  and  $PE$  by ordered domains  $(U, \leq_U)$  and  $(V, \leq_V)$ . Then the Duhem-Hertz condition (5.6) simply requires that there are order-preserving maps  $U \xrightarrow{e} V$  and  $V \xrightarrow{t} U$  such that  $(t, e)$  defines a Galois connection between  $U$  and  $V$  in the sense of (5.2). This may be even further generalized by the assumption that  $U$  and  $V$  are categories in an adjoint situation (cf. Goldblatt 1978). That is to say, conceiving an empirical theory as an adjoint situation  $(F, G)$  between a category of symbolic facts  $U$  and a category  $V$  of empirical facts combines in a neat and natural way the classical insights of Hertz and Duhem.

As a summary of this section let us reformulate once more the basic thesis in somewhat different terms, assuming that an empirical theory is given as a Galois connection  $(t, e)$  between an ordered domain  $(U, \leq)$  of symbolic facts and an ordered domain  $(V, \leq)$  of empirical facts, i.e. the maps  $U \xrightarrow{e} V$  and  $V \xrightarrow{t} U$  satisfy the Galois equivalence

$$(5.7) \quad e(x) \leq a \quad \text{IFF} \quad x \leq t(a), \quad x \in U, a \in V.$$

Then we may conceive  $x$  as a *theoretical law* that may be considered as the blueprint for the building of a nomological machine or experimental apparatus  $e(x)$  that produces the empirical fact  $a$  as its outcome. Then the Galois connection states:



The nomological machine  $e(x)$  brings about the empirical fact  $a$

IFF

The theoretical law  $x$  implies an idealized version  $t(a)$  of  $a$ .

This yields another interpretation of the formal apparatus of Galois connection that renders plausible the claim why theories which satisfy the Galois connection should be considered as (approximately) true theories: such theories are approximately true since they ensure a relation between the empirical and the theoretical that captures the idea that an approximately true theory should approximately correspond to the facts.

### 6. Concluding Remarks

The leitmotif of this paper was the thesis that scientific theories are to be considered as *representations*, and, more generally, that the practice of science may be conceptualized as a representational practice. This idea is not new, and many have put forward it in many different ways. Philosopher-scientists such as Hertz and Duhem provide distinguished examples. Tapping some of their essential insights we hope to have rendered plausible the following theses: (i) representation is a complex concept in need of a theory, (ii) representations do not live in isolation. Rather, they may be *iterated* and *combined* in various ways, and (iii) representations do not ‘speak for themselves’. Rather, representations are in need of interpretation. A large part of scientific practice consists in interpreting and reinterpreting representations.

### REFERENCES

- Duhem, P. (1906). *The Aim and Structure of Physical Theory*. Princeton: Princeton University Press, 1954.
- Gierz, G.; Hofmann, K.H.; Keimel, K.; Lawson, J.D.; Mislove, M.; Scott, D.S. (2003). *Continuous Lattices and Domains*. Cambridge: Cambridge University Press.
- Goldblatt, R. (1978). *Topoi – The Categorical Analysis of Logic*. Amsterdam and New York: North-Holland.
- Hacking, I. (1983). *Representing and Intervening*. Cambridge: Cambridge University Press.
- Hertz, H. (1894). *The Principles of Mechanics Presented in a New Form*. New York: Dover, 1956.
- Ibarra, A.; Mormann, T. (1997a). *Representaciones en la ciencia. De la invariancia estructural a la significatividad pragmática*. Barcelona: Ediciones del Bronce.
- (1997b). “Theories as Representations”, in A. Ibarra, T. Mormann (eds.), *Representations of Rationality*. Poznan Studies in the Philosophy of Science and the Humanities 61, Amsterdam and Atlanta, 65-92.
- (2000). “Una teoría combinatoria de la representación científica”, *Crítica*, vol. XXXII, N.º. 95, 3-46.
- Latour, B. (1999). *Pandora’s Hope. Essays on the Reality of Science Studies*, Cambridge/Massachusetts: Harvard University Press.

- Lawvere, F. W.; Schanuel, S.H. (1996). *Conceptual Mathematics. A First Introduction to Categories*. Cambridge: Cambridge University Press.
- Laymon, R. (1982). "Scientific Realism and the Hierarchical Counterfactual Path from Data to Theory", in P. Asquith and T. Nickles (eds.), *PSA 1982*, vol. 1, Philosophy of Science Association, East Lansing, 107-121.
- Mac Lane, S. (1986). *Mathematics. Form and Structure*. New York: Springer.
- Melton, A.; Schmidt, D.A.; Strecker, G.E. (1985). *Galois Connections and Computer Science Applications*. Springer Lecture Notes in Computer Science 240. Berlin: Springer, 299-312.
- Strand, R.; Fjelland, R.; Flatmark, T. (1996). "In Vivo Interpretations of In Vitro Effect Studies with a Detailed Analysis of the Method of *in vitro* Transcription in Isolated Cell Nuclei", *Acta Biotheoretica* 44, 1-21.
- Strand, R. (1999). "Towards a Useful Philosophy of Biochemistry: Sketches and Examples", *Foundations of Chemistry* 1, 271-294.

**Andoni IBARRA** and **Thomas MORMANN** have worked together in a number of projects dealing with the role of representations in science, philosophy of science, and general epistemology. They published a book on this topic (*Representaciones en la Ciencia*, Barcelona, 1997) and served as editors of a number of anthologies on representational themes.

Andoni Ibarra moreover is presently coordinator of the Sanchez-Mazas Chair at the University of the Basque Country. Thomas Mormann holds a PhD in Mathematics and obtained his Habilitation in philosophy from the University of Munich (Germany).

**ADDRESS:** Dpt. of Logic and Philosophy of Science, University of the Basque Country, Av. Tolosa 70, 20018 Donostia-San Sebastián, Spain. E-mails: andoni.ibarra@ehu.es; ylxmomot@sf.ehu.es.