

## A reconstruction of Sigmund Freud's early theory of the unconscious

Balzer, Wolfgang; Marcou, Ph.

Veröffentlichungsversion / Published Version  
Sammelwerksbeitrag / collection article

### Empfohlene Zitierung / Suggested Citation:

Balzer, W., & Marcou, P. (1989). A reconstruction of Sigmund Freud's early theory of the unconscious. In H. Westmeyer (Ed.), *Psychological theories from a structuralist point of view* (pp. 13-31). Berlin u.a.: Springer. <https://nbn-resolving.org/urn:nbn:de:0168-ssoar-39851>

### Nutzungsbedingungen:

Dieser Text wird unter einer CC BY-NC-ND Lizenz (Namensnennung-Nicht-kommerziell-Keine Bearbeitung) zur Verfügung gestellt. Nähere Auskünfte zu den CC-Lizenzen finden Sie hier:  
<https://creativecommons.org/licenses/by-nc-nd/4.0/deed.de>

### Terms of use:

This document is made available under a CC BY-NC-ND Licence (Attribution-Non Commercial-NoDerivatives). For more information see:  
<https://creativecommons.org/licenses/by-nc-nd/4.0>

## A Reconstruction of Sigmund Freud's Early Theory of the Unconscious

W.Balzer and P.Marcou

University of Munich

### Abstract

A version of Freud's early theory of the unconscious — which is a completely qualitative theory — is reconstructed in precise terms. The important structuralist components: models, potential models, partial potential models, constraints, links, measuring models, intended applications are identified, and an explanation for neurosis is given. The theory's empirical claim and confirmation are considered. The particular position of therapeutically relevant theories is discussed and clarified. The reconstruction shows that Freud's early theory lends itself to a completely precise logical treatment and has all the features of empirical theories in the social sciences in general.

In this paper a version of Freud's early theory of the unconscious is presented in the format of an empirical theory which was developed in structuralist philosophy of science.<sup>1</sup> The present account emends our earlier treatment in Balzer (1982) which was already taken up in Stegmüller (1986). We concentrate on Freud's first writings from the years 1892-99 in which his theoretical ideas are expressed in few sentences dispersed in those texts. These sentences together with his descriptions of various cases (many of which occur in later works) form the "data" on which our reconstruction is based.<sup>2</sup> We want to show that Freud's ideas can be stated so as to form a consistent and (in a certain sense) complete picture which - though qualitative - allows for a precise description. By this, we hope to contribute constructively to the never ending debate about his writings. If it should turn out that our reconstruction misses some point this will come out by precisely locating that point; by so doing the reconstruction may be emended, and the controversy had some definite result.<sup>3</sup>

From the theory's exemplifying all important features of an empirical theory, and from the assumption that our reconstruction is adequate it follows that the theory satisfies the essential requirements empirical theories outside the realm of natural science are generally expected to satisfy. On the level of psychology this is a point in favour of "bold" theorizing (as contrasted to the dominant, operational attitude). On the meta-level we provide a further successful application of the structuralist view of theories.

Moreover, we take side on the issue of how to look at medical (or therapeutically relevant) theories in general: sickness (e.g. neurosis) is represented by specializing a respective model

---

<sup>1</sup>Compare Balzer, Moulines, & Sneed (1987) for a recent, comprehensive account of the structuralist meta-theory.

<sup>2</sup>We use Freud (1967) as main text. The works mentioned are found in volume 1, pp. 2-459. Further references always refer to this volume, if not stated otherwise.

<sup>3</sup>In fact, two such emendations of our original account have been induced by criticism of W.Stegmüller and M.Perrez (which we gratefully acknowledge), and are built in in the present version.

of health. In formalizing this picture a decisive difference emerges between medical theories on the one hand and theories from the natural sciences as well as from the social sciences on the other hand.

In order to avoid misunderstandings it may be useful also to state what we are not claiming. We are not claiming that our reconstruction is the only possible one. As the word indicates, any re"construction" involves some constructive element which may well go beyond the original or beyond what the original author wanted to express. We are not so much interested in what a particular author meant by writing his papers, the interesting point for a reconstruction is whether his theories have some content independent of his intentions. Other reconstructions of the same matter may look differently, and by comparison it may be possible to eliminate one in favour of another one (if they are equivalent, or if the former reduces to the latter, for instance).

This leads to a second point. We do not have independent, "absolute" criteria for a reconstruction being adequate. The business of reconstructing scientific theories is very recent, not much older than 50 years. No criteria ready for use have been put forward in this time. This situation is by no means exceptional, it also obtains in every other domain where scientific theories are applied to real systems. In physics, to mention the most admired case, 400 years of experimental application did not lead to rough and ready criteria of when the application of a theory to a system, or in an experiment, is adequate. The criteria for adequacy of a reconstruction, as well as those for successful application of a theory still are mainly those of performance in comparison with alternative attempts.

Third, we do not want to discuss in detail, and to put forward a kind of decision on, the issue of whether psychology is scientific or not.<sup>4</sup> We feel that the standards for "scientific" activities as put forward by logical empiricism or by Popper are too narrowly drawn from the theoretical branches of natural science. On these standards too many scientific activities (among them many firmly established at universities) fall off the board. We do not object to drawing distinctions but we fear great damage from eliminating psychology (together with most social sciences) from the realm of science by insisting on criteria mainly drawn from classical physics.

Fourth, we do not claim that our reconstruction covers the overwhelming richness present in Freud's writings. There are two aspects to this richness. First, there are many details, data, reported in the descriptions of real cases. These details are parts of the descriptions of intended applications of the theory, and do not occur in the following because for reasons of space we cannot state the intended applications in detail. Second, Freud states, or partially indicates a large number of hypotheses; far more than will be treated here. With respect to these the axiomatic method proves its force. It reduces a large number (infinitely many, to be sure) of hypotheses to a small number of axioms from which all the former are derivable. Not all the hypotheses indicated in Freud's writings will follow from ours. We believe, however, that a substantial part of them can be formulated as specializations of the basic models introduced here. Thus our reconstruction only provides the frame, the basis on which a full net of specializations (that is, in structuralist terminology, a theory-net<sup>5</sup>) can be worked out comprising the full range of Freudian ideas. Still some hypotheses will not be treatable in

---

<sup>4</sup>Compare Perrez (1979) for a review and critique of discussions before 1979, and later on Gruenbaum (1984).

<sup>5</sup>See Balzer, Moulines, & Sneed (1987), Chap. IV.

this way : they "get lost" in the process of reconstruction. We do not have to say new things about other Freudian topics, like his account of dreams or his later theories of the ego, id, and superego.

## Freud's Basic Picture

Freud often uses physical metaphors which may serve as a guide through the abstract passages. For his early theories the "steam pot picture" seems particularly relevant: a closed pot filled with water is constantly heated up, and so from time to time the cover will be lifted and steam will escape. This picture by way of analogy models his view of the unconscious. Man in some respects is like the steam pot. It is constantly "heated up" by an ongoing flow of "drivings" ("Affekte"<sup>6</sup>). The steam is set analogous to "ideas" ("Vorstellungen"<sup>7</sup>), pressure to "suppressing" ("verdrängen"<sup>8</sup>), and escaping to "abreaction" ("Abreagieren", "assoziatives Verarbeiten"<sup>9</sup>). The term "abreaction" ("Abreagieren") is used in connection with drivings as well as with ideas<sup>10</sup>, we restrict it to drivings here. Thus the drivings cause the production of ideas which usually give rise to abreactions, or actions. But what if the cover is locked, i.e. if those actions are not feasible? According to Freud the drivings then are blocked (his term is "eingeklemmt"<sup>11</sup>), pressure rises, and the ideas get suppressed, they become what he later calls unconscious. Here are some important passages.<sup>12</sup> "She [psychotherapy] takes back the effect of an idea originally not abreacted by allowing the blocked affect of it to run off through speech, and causes its associative correction by pulling it into ordinary consciousness..." (p. 97), "The pictures of "screening" an incompatible idea, of the origin of hysterical symptoms by conversion of psychical to bodily excitement, of the formation of a separate psychical group by means of the act of the will which leads to the screening, all this was seizably put before my eyes in this moment." (p. 222), "The screening reaches its aim of pressing the incompatible idea out of consciousness, if in the respective person, healthy up to then, sexual scenes are present as unconscious remembrances, and if the idea to be suppressed can be brought into logical or associative connection to such an infantile experience." (p. 447-8), "If in a person with disposition there is no ability for conversion while nonetheless for the sake of screening an unbearable idea this idea is cut off from its affect then this affect has to remain in the psychic field. The idea thus weakened remains in consciousness aside of all associations, its freed affect however joins other ideas, which by themselves are not incompatible but, by means of this 'false connection' become compulsive ideas." (p. 65-6). In later periods, the unconscious acquires its central position: "We learned from psychoanalysis that the essence of the process of suppression does not consist in the destruction of an idea representing the driving. Rather it consists in keeping the idea away from becoming conscious. We then say the idea is in a state of 'unconsciousness', and we can offer good proofs that, though unconscious, it can have effects; even effects that eventually reach 'consciousness'." <sup>13</sup>

---

<sup>6</sup>See, e.g. pp. 65, 66, 85, 97.

<sup>7</sup>E.g. pp. 63, 66, 75, 90, 97, 174.

<sup>8</sup>E.g. pp. 174, 181, 234, 386, 387, 388.

<sup>9</sup>E.g. pp. 87, 89, 90, 94, 97, 224.

<sup>10</sup>On pp. 87, 89, 94 it is used for drivings, on pp. 90, 97 for ideas.

<sup>11</sup>For instance on p. 97.

<sup>12</sup>Translations are ours.

<sup>13</sup>Freud (1963), p. 7. Translation ours.

In order to give a precise account of this, terminology needs to be fixed. We will use the term "event" in a broad sense including ideas, as well as all kinds of experiences a person may have, and actions. In particular, we will use the term such as to include possible events. It will be convenient to regard events as tokens rather than as types. The models to be considered will refer to a person's (possible) concrete experiences, ideas, and actions over some fixed period of time. These experiences, ideas and actions are represented by a set  $E$  of events. Besides event-tokens we also use event-types which we treat as classes of event-tokens similar in certain respects. The similarity-relation is left unspecified, and even implicit; it has to be chosen in each application appropriately. Event-types are denoted by  $k, k'$ , and their collection by  $K$ , so that  $K$  is a subset of the power set of  $E$ . For reasons of simplicity we assume that event-types are disjoint. With a little extra effort at other places this assumption may be relaxed.

Some events are of special importance, for instance those which are conscious or those which are particularly horrible ("negative"). We introduce consciousness as a binary relation between instants and events:  $C(t, e)$  means that the person at time  $t$  is conscious of event  $e$ .

Drivings are treated as ontologically different from events, they are entities of a new kind occurring in the psyche of humans. The set of drivings relevant for a person we denote by  $D$ . Drivings also are tokens rather than types. It is the particular driving, say, to eat at 3 a.m. under given circumstances rather than the type of driving we call "hunger" that occurs in the models.

Next, we use the notion of a driving being abreacted ("abreagiert"). We prefer to use the more neutral term "realization" instead. Realization is treated as a four-place relation between instants, drivings, events and event-types.  $REAL(t, d, e, k)$  means that, at  $t$ , driving  $d$  is realized (abreacted) by event  $e$  of type  $k$ .

Since each model will refer only to one period in a single person's life there will be no need to represent that person in the model. All the concepts used in describing a model are always understood as referring just to one person's case.

A further primitive we use is that of a horrible event, horrible for the person under consideration. We will speak more neutrally of "negative" events, and represent them by a binary relation between instants and events:  $N(t, e)$  means that, at  $t$ , the person experiences the negative event  $e$ . Roughly, negative events are those which cause drivings to be blocked from getting realized.

Our most important primitive is that of the unconscious which we introduce as a binary relation  $U$  between instants and events.  $U(t, e)$  means that, at  $t$ , event ("idea")  $e$  is unconscious or suppressed. We need no special term for blocked drivings; these will be defined by means of the above terms.

Next, we use a function  $f$  to assign to each driving  $d$  the type of events by means of which  $d$  ordinarily, or naturally, is realized. We assume that each driving has exactly one such natural type of realizations. This assumption is more technical in nature, for different event-types may be joined to form one new, "bigger" type.

In order to round up the picture we use an ordering relation  $\leq$  among instants, and a binary relation  $AS$  of association among events:  $AS(e, e')$  means that events  $e$  and  $e'$  are

associated for the person considered. Association plays an important role in psychoanalysis. However, from its role in Freud's theory as well as from the development of psychology later on we conclude that Freud's theory is not intended to give a special meaning to that term. We treat association as not explicitly depending on time. This indicates that the "mechanism" of association essentially works "across" time, depending only on the internal structure of the events involved.

It may be noted that our use of events to cover the actions, experiences, as well as ideas has no analogue in Freud's writing; it is an essential ingredient of our reconstruction. We chose events for their great unifying power. It also may be noted that consciousness and the corresponding actions might be omitted without essential loss. We include it, mainly in order to have a "complete" set of terms for the psychic structure: drivings-consciousness-unconsciousness, which motivates the label of unconsciousness so central in Freud.

Altogether we thus arrive at the following list of primitives:  $T$ , a set of instants;  $E$ , a set of events;  $D$ , a set of drivings;  $K$ , a set of event-types;  $\leq$ , an ordering-relation among instants;  $AS$ , an association-relation among events;  $C$ , consciousness;  $U$ , unconsciousness;  $N$ , negative events;  $REAL$ , relation of realization;  $f$ , assignment of event-types to drivings.

## Potential Models and Models

By stating the precise "grammar" of these primitives, as well as some further, trivial requirements, we obtain the so called potential models of the theory.<sup>14</sup> These are possible systems in which all the primitives have some interpretation but which need not satisfy the theory's central axioms. Among such possible systems there may be real systems as well as purely abstract ones defined in terms of numbers or abstract sets. The class  $M_p$  of all potential models may be regarded as the set of "possible worlds" for Freud's theory. Any "world" (=potential model) which contains interpretations for all primitives is possible.

- D1  $x$  is a *potential model* of Freud's theory ( $x \in M_p$ ) iff there exist  $T, E, D, K, \leq, AS, C, N, U, REAL$  and  $f$  such that  $x = \langle T, E, D, K, \leq, AS, C, N, U, REAL, f \rangle$  and
- 1)  $T, E, D$  are finite, non-empty sets, and pairwise disjoint;
  - 2)  $K$  is a partition of  $E$  (i.e.  $K \subseteq Po(E)$ , all  $k \in K$  are non-empty and pairwise disjoint, and  $\bigcup \{k : k \in K\} = E$ );
  - 3)  $\leq$  is a weak order on  $T$  (i.e.  $\leq \subseteq T \times T$  is reflexive, transitive, and connected);
  - 4)  $AS \subseteq E \times E$ ;
  - 5)  $C \subseteq T \times E, N \subseteq T \times E$ , and  $U \subseteq T \times E$ ;
  - 6)  $REAL \subseteq T \times D \times E \times K$ ;
  - 7)  $f : D \Rightarrow K$  is injective.

The set  $T$  of instants must not be confused with instants of physical time (which often are represented by real numbers). The elements of  $T$  simply serve as indices for ordering events and drivings as they occur one after another. Note that by D1-1,  $T$  is a finite set. It is best to imagine members of  $T$  as those (physical) instants at which something important

<sup>14</sup>See Balzer, Moulines, & Sneed (1987), Chap.I.

happens in the person's history: she experiences an event important from the point of view of psychology, or some sufficiently strong driving is present. The choice of  $T$  in an application is not much restricted theoretically but a wrong choice may lead to an uninteresting model. If  $T$  contains too many instants, the model will be overloaded with redundant data, if  $T$  is "too small" nothing interesting will happen in the model. It may be noted that  $T$  needs not cover the whole life of the person described.

In order to get used to this kind of representation, consider the following trivial, contrived example of a person  $B$  during an ordinary working day, say November 4, 1977. The instants will make up a list of the following kind: 7 a.m., 7.15, 7.35, 7.50, 8, 9, 11, 11.30, ... Of course, these instants are ordered in the natural way. The set  $E$  of events will comprise those events which in some sense are important, and to each event-token we may assign a corresponding event-type as given by the verbal description (which never can be as fine as a token).  $E$  may contain, among others, the event of being waked up (at 7 a.m.), of taking a shower (at 7.15), having breakfast (at 7.35), seeing an accident while sitting in the bus (7.50), entering office (at 8), getting an order (at 9) etc. The set of drivings contains those occurring during that day, like a little sexual excitement under the shower, a small feeling of hunger before breakfast, nervousity and anxiety while meeting the boss at 9, hunger at 11 etc. Some of the events experienced will get associated with other events in  $B$ 's history: the event  $e_1$  of listening to some advertisement at breakfast will get associated with other, past, experiences  $(e_1)_2, \dots, (e_n)_2$  of the same silly sentences, the event  $e_3$  of seeing a broken limb at the accident at 7.50 will get associated with the past event  $e_4$  of  $B$  himself breaking a leg. The association-relation  $AS$  will consist of a set of such pairs of associated events:

$$AS = \{ \langle e_1, (e_1)_2 \rangle, \langle e_1, (e_2)_2 \rangle, \dots, \langle e_1, (e_n)_2 \rangle, \langle e_3, e_4 \rangle, \dots \}.$$

Consciousness is represented by those pairs  $\langle t, e \rangle$  for which  $e$  is some event of importance which gets into the person's mind, and  $t$  is the instant at which it occurs.  $\langle 7.50, e_3 \rangle$  and  $\langle 9, e_5 \rangle$  are examples, if  $e_5$  is the event of  $B$ 's getting a special order from his boss. So  $C = \{ \langle 7.50, e_3 \rangle, \langle 9, e_5 \rangle, \dots \}$ . Negative events in such an ordinary case will not occur, so  $N$  will be empty. The realization-relation will consist of events which abreact present drivings. One example here might be  $\langle 11, d, e_6, k \rangle$  where  $d$  is "hunger" and  $e_6$  the event of having lunch. Thus  $REAL = \{ \langle 11, d, e_6, k \rangle, \dots \}$ . On the assumption that  $B$  is a psychically healthy person, there will be no important unconscious, suppressed events, i.e.  $U = \emptyset$ . If we put together all these sets in a 11-tuple we obtain a potential model which describes  $B$  during the particular day from Freud's psychological point of view.

Such a description is called a model of the theory if it satisfies the following axioms.

- A1 If, at  $t$ , event  $e$  of type  $k$  realizes driving  $d$  then, at  $t$ ,  $e$  is conscious, and of type  $k$ .
- A2 Any negative event at  $t$  is conscious at  $t$ .
- A3 Any event is associated with itself.

These are rather trivial, "analytic" statements concerning the use of the primitives. A1 rules out the case of  $e$  realizing  $d$  without  $e$  being conscious. Similarly, A2 rules out negative events from not being conscious at the time when they occur. A3 is simply a convention. In addition, there are four more substantial axioms. We say that  $e$  is a *natural realization* of driving  $d$  if  $e \in f(d)$ .

- A4** An event  $e$  cannot be unconscious and a natural realization at the same time.
- A5** If some driving  $d$  is realized by two events  $e$  and  $e'$  (possibly at different times) then these events are associated.
- A6** For any two events which are associated and which both are natural realizations of the same driving: if one of the two is negative at  $t$  then the other will be unconscious at any later time.
- A7** Any driving becomes realized at some time.

A4 imposes a restriction on the realization-relation. Only those events may occur as realizations of some driving which are not unconscious. In other words, unconscious, suppressed events cannot serve for natural realization. A5 says that associations of events are caused by their being linked to a common driving. In this form A5 certainly contains a strong element of idealization. We note, however, that realization of  $d$  by  $e$  implies that  $e$  is conscious (by A1), so  $e$  has to be rather strong and important. A6 represents the "mechanism" of suppression: if  $e$  at  $t$  is negative ("horrible") then any natural realization of the driving realized by  $e$  which is associated with  $e$  will get unconscious after  $t$ . A7 is the central axiom concerning the steam pot picture: any driving will lead to a corresponding realization.

By collecting these axioms and eliminating some ambiguities still present in the verbal formulation we obtain

- D2**  $x$  is a *model* of Freud's theory ( $x \in M$ ) iff there exist  $T, E, D, K, \leq, AS, C, N, U, REAL, f$  such that  $x = \langle T, E, D, K, \leq, AS, C, N, U, REAL, f \rangle$ ,  $x \in M_p$  and
- 1) for all  $t \in T$ ,  $d \in D$ ,  $e \in E$  and  $k \in K$  : if  $REAL(t, d, e, k)$  then  $C(t, e)$  and  $e \in k$ ;
  - 2) for all  $t \in T$  and  $e \in E$ : if  $N(t, e)$  then  $C(t, e)$ ;
  - 3) for all  $e \in E$  :  $AS(e, e)$ ;
  - 4) for all  $t \in T$ ,  $e \in E$  and  $k \in K$  such that  $e \in k$ : not [ $U(t, e)$  and (there exists  $d \in D$  such that  $REAL(t, d, e, k)$  and  $k = f(d)$ )];
  - 5) for all  $t, t' \in T$ ,  $e, e' \in E$ ,  $d \in D$  and  $k, k' \in K$  : if  $REAL(t, d, e, k)$  and  $REAL(t', d, e', k')$  then  $AS(e, e')$ ;
  - 6) for all  $t, t' \in T$ ,  $e, e' \in E$  : if  $N(t, e)$ ,  $t < t'$ ,  $AS(e, e')$ , and if there exists  $d \in D$  such that  $e \in f(d)$  and  $e' \in f(d)$  then  $U(t', e')$ ;
  - 7) for all  $d \in D$  there exist  $t \in T$ ,  $e \in E$  and  $k \in K$  such that  $REAL(t, d, e, k)$ .

The expression " $t < t'$ " in A6 is of course defined in terms of  $\leq$  in the usual way:  $t < t'$  iff ( $t \leq t'$  and not( $t = t'$ )). Obviously,  $M \subseteq M_p$ .



## Neurosis

The formal models just introduced basically describe "normal" healthy persons but they also can be applied to psychically sick persons. One may wonder how this is possible because "sick" and "healthy" exclude each other, and a theory allowing for both would seem inconsistent. The first thing to note here is that the axioms are, in fact, consistent.

**Theorem 1**  $M$  is not empty.

**Proof :** It is easy to check that  $x = \langle \{t\}, \{e, e'\}, \{d\}, \{\{e\}, \{e'\}\}, \{\langle t, t \rangle\}, \{\langle e, e \rangle, \langle e', e' \rangle\}, \{\langle t, e \rangle, \langle t, e' \rangle\}, \emptyset, \emptyset, \{\langle t, d, e, \{e\} \rangle\}, \{\langle d, \{e\} \rangle\} \rangle$  is a model. #

Secondly, it has to be noted that one single model cannot describe a sick and a healthy person at the same time. The distinction between sickness and health cannot be drawn within a single model; it has to be drawn in the class of models. We may regard  $M$  as being split up into two subsets, one subset, *HEALTH*, containing exactly the models describing healthy persons, and a second one, *SICK*, describing sick persons. Still the problem remains to characterize *SICK*, for how can a person be psychically sick (be described by an  $x \in SICK$ ) and at the same time satisfy all the axioms for healthy, "normal" persons ( $x \in M$ )?

In order to see this let us reflect on how a person gets sick (neurotic) according to the scheme expressed in the axioms of D2. At some instant  $t$  the person realizes some driving  $d$  by means of event  $e$  of type  $k$  :  $REAL(t, d, e, k)$ . This causes sickness if  $e$  turns out as a horrible experience,  $N(t, e)$ , due to the reactions of other persons which our person could not anticipate. By D2-6, any event  $e'$  which would be a natural realization of  $d$  and which is associated with  $e$  at any later time  $t'$  will be unconscious:  $U(t', e')$ . In the light of D2-7 this is a source of inner stress for the person because, according to D2-7, the driving  $d$  "presses for" realization, but realization, by D2-4, is excluded by  $e$  being unconscious. More precisely, D2-7 yields some  $e'$  and  $t'$  such that  $REAL(t', d, e', k')$ . If  $k'$  is the natural type of realizations for  $d$  then  $k' = k$  which, by D2-4, implies: not  $U(t', e')$ . But if  $t < t'$  this contradicts to  $U(t', e')$  above, so  $t' \leq t$ , and  $d$  can be realized by a natural event from  $f(d)$  only at  $t$ , or later by events of a different type (neurotic symptoms). In other words, the negative event is conscious at the time  $t$  of its first being experienced, at which time it also is a realization of some corresponding driving. But this driving cannot get realized by natural realizations at any later time, it is blocked in this sense, and the person gets sick (neurotic) because of this suppression. The symptoms of neurosis on Freud's account occur because the driving gets dissociated from its natural realizations and, in order to find some abreaction at all, gets realized by those symptoms (this feature of the theory is not made explicit here). The cause of sickness thus is given by a pair  $\langle d, e \rangle$  of a driving  $d$  and an event  $e$  such that  $REAL(t, d, e, k)$  and  $N(t, e)$ . This derivation is central to Freud's theory, it shows how neurosis is caused, and therefore explains neurosis and psychical sickness. Because of its importance let us state the derivation as a formal proof.

**Theorem 2** Let  $x = \langle T, E, D, K, \leq, AS, C, N, U, REAL, f \rangle \in M$ , and let  $t \in T$ ,  $d \in D, e \in E$  and  $k \in K$  be such that  
 (1)  $N(t, e), REAL(t, d, e, k)$  and  $f(d) = k$ .  
 Then

- (2) there is no  $t' \in T$  and no  $e' \in E$  such that  $t < t'$  and  $REAL(t', d, e', k)$ .

**Proof:** Suppose there exist  $t' \in T$  and  $e'$  such that  $t < t'$  and

$REAL(t', d, e', k)$ . By D2-4,

(3) not  $U(t', e')$ . By D2-5, (1), and the assumption,

(4)  $AS(e, e')$ . By (1), D2-1, and the assumption,

(5)  $e \in k$  and  $e' \in k$ . From this and (1) we obtain  $e, e' \in f(d)$  which by (1), (4) and D2-6 yields

(6)  $U(t', e')$ . From (5), (1) and the assumption we get

(7)  $e' \in k$  and there exists  $d \in D$  such that  $REAL(t', d, e', k)$  and  $k = f(d)$ . But (7) and D2-4 yield: not  $U(t', e')$ , in contradiction to (6).#

The conclusion (2) of theorem 2 may be rephrased by saying that after  $t$ , driving  $d$  is blocked and the natural realizations for  $d$  are suppressed. Theorem 2 therefore may be restated as follows. If driving  $d$  gets realized at  $t$  by a negative event then any natural realization for  $d$  will be suppressed after  $t$ .

These considerations show how psychical sickness may occur in a model even if all axioms for a healthy person are satisfied. Sickness is caused by negative events. If there are no such events the person normally (not provably) will be healthy. If there are negative events they are realized at the moment of first being experienced (by which the axioms for normal development are satisfied), and afterwards the corresponding driving is blocked (the person gets sick). To put it differently: the axioms describing the "normal" case still provide room for sickness, they represent necessary, but not sufficient conditions for health.

We define a person to be psychically sick or neurotic precisely if some driving is blocked after a certain instant in the sense just introduced:  $d$  is *blocked after*  $t$  iff there is no  $t' \in T$ , no  $e \in E$  and no  $k \in K$  such that  $t < t'$ ,  $REAL(t', d, e, k)$  and  $k = f(d)$ , and if there is some  $t^* \in T$  such that  $t < t^*$ . The final clause here excludes the definition to be trivialized by taking  $t$  as the "last" point of time occurring in  $T$ . The models describing persons with blocked drivings form a subset *SICK* of  $M$ .

**D3**  $x$  is a *model of a psychically sick person* ( $x \in SICK$ ) iff

- 1)  $x = \langle T, E, D, K, \leq, AS, C, N, U, REAL, f \rangle \in M$ ;
- 2) there exist  $t \in T$  and  $d \in D$  such that  $d$  is blocked after  $t$ .

Clearly,  $SICK \subseteq M$ . This characterization of neurosis is theoretical, and does not refer to the real symptoms. In applications of the theory, one will start of course from the symptoms, infer "backwards" that the person is sick in the sense of D3, and then start to find the crucial items of the model:  $D, U, REAL$ , and  $f$ . Note that suppression takes place "immediately" after the negative event occurs (A6). This does not preclude, however, that neurotic symptoms may occur only much later, for we do not establish here the connections between the unconscious and neurotic symptoms in explicit terms. The unconscious remains in a rather theoretical status without explicit links to neurotic symptoms. Note also that negative events on this account are strongly theory-laden. The only indicator for the occurrence of such events is neurosis. If a person is not neurotic even the most dramatic events will not be counted as negative.

## Constraints, Links, Measuring Models

According to structuralist meta-theory an empirical theory in general consists of a *formal core*  $K$  and a set  $I$  of *intended applications*. The core itself is made up of five items: classes  $M_p$  and  $M$  of potential models and models, constraints  $C$  and links  $L$ , some approximation apparatus  $U$ , and a distinction between theoretical and non-theoretical terms.

Constraints express assumptions of identity or stability "across" different models. In the present case such an assumption may be stated for association: if two events are associated in one potential model (i.e. for the person described by that model) then they also are associated in any other potential model describing any other person. Two features of this requirement should be noted. First, it is "across" different potential models. This is due to our representing different persons by different models. Second, the requirement applies only in cases where two persons "have" or experience identical events. This comes out clearly once we formulate the constraint in precise terms. We write  $E^x, AS^x, E^y$  etc. in order to denote the entities  $E, AS$  etc. occurring in system  $x = \langle T^x, E^x, \dots, AS^x, \dots \rangle$  and system  $y = \langle T^y, \dots, AS^y, \dots \rangle$ .

D4  $X$  satisfies the *constraint* of Freud's theory ( $X \in C$ )

iff

- 1)  $X \subseteq M_p$  and  $X$  is not empty;
- 2) for all  $x, y \in X$  and all  $e, e'$  :  
if  $e \in E^x \cap E^y$  and  $e' \in E^x \cap E^y$   
then  
( $AS^x(e, e')$  iff  $AS^y(e, e')$ ).

Any  $X$  satisfying the constraint may be called an admissible combination (of potential models). It is admissible insofar as it represents a set of persons in which association is stable in the sense of D4. Of course, this constraint is very idealized, and will be satisfied only to some degree, but it nevertheless is crucial for applying the theory in psychoanalysis. For it allows to infer certain associations from other persons' associations or from earlier ones of the person under treatment. A second constraint which we do not formulate explicitly, requires stability of the function  $f$  which assigns types of natural events to drivings. The constraint says that "identical" drivings in different models (persons) get identical  $f$ -values (types of events).

The role of links in a theory is to "import" data, information and meaning from "outside", that is, from other, "surrounding" theories or from everyday experience. A first link for Freud's theory concerns the ordering of instants. This ordering is constituted mainly in physical terms. We may regard each potential model as being linked to a model of some theory of space-time such that the instants in the former correspond to instants in the latter, and the physical ordering in the latter determines the ordering of the corresponding instants in the psychological model. Note that this does not enforce an ontological identity of corresponding instants, so the existence of this link does not conflict with what was said about the interpretation of  $T$  above. Formally, the link may be represented by a class of pairs of potential models  $\langle x, y \rangle$  ( $x \in M_p$  and  $y \in M_p(SP)$ , the class of potential models of a theory of space-time) in which corresponding instants are linked.

D5 If  $SP$  is a theory of space-time with class  $M_p(SP)$  of

potential models such that the elements of  $M_p(SP)$  contain an ordering relation  $\leq'$  then

$\langle x, y \rangle$  satisfies the link to  $SP$  iff  $x \in M_p$ ,  
 $y \in M_p(SP)$ ,  $x = \langle T, \dots, \leq, \dots \rangle$ ,  $y = \langle T', \dots, \leq', \dots \rangle$  and there  
 is a one-one mapping  $g : T \Rightarrow T'$  such that for all  $t, t' \in T$  :  
 $t \leq t'$  iff  $g(t) \leq' g(t')$ .

A second link can be established to theories of association. It imports information about the association-relation from association theory. A third link may be used which connects Freud's theory with theories dealing with consciousness. From these different links a *global* link  $L$  may be defined which consists of all systems  $x$  that are linked by each of the "individual" links to some system from another theory.

We skip the discussion about which terms of the theory are theoretical. Also, we cannot go into the details of describing a full-fledged approximation apparatus to be used in applying the above models. As the theory is completely qualitative the approximation apparatus has to be built up by exploiting the notion of qualitative similarity. Definitions like that proposed by Tversky<sup>15</sup> may be used to form neighbourhoods by means of which the idealized claims associated with the previous idealized items have to be blurred in order to yield realistic, non-trivial empirical claims.<sup>16</sup> We note that, in general, the need for approximation comes from idealizing features inherent in quantitative theories. If quantities (and therefore numbers) are involved, the theoretical pictures become much finer than reality, and therefore the theoretical picture can be claimed to represent reality ("to be true") at most approximatively. But Freud's theory is completely qualitative. Still, it contains some idealizing features. Instead of blurring these, it also seems possible to weaken some axioms. In order to obtain a more liberal version of D2-5, for instance, we could use an approximation of the following kind. The if-clause "if  $REAL(t, d, e, k)$  and  $REAL(t', d, e', k')$ " might be strengthened to refer to repeated such experiences "if  $REAL(t_i, d, e, k)$  and  $REAL(t_i, d, e', k')$  at different times  $t_i, i = 1, \dots, n$ ". Approximation then could be "tacked on" at the number  $n$  of repetitions.

In addition to the items introduced so far it is of interest to investigate the means of measurement formally possible within the present theory. Such investigation, even if performed purely abstractly, may be helpful in evaluating the empirical performance or adequacy of the theory. Any model in which one of the primitives  $t$  is uniquely determined by the others we call a measuring model for  $t$ . This notion is helpful in two respects. First, if a term has no measuring model at all in the theory, this indicates that the term obtains its meaning from "outside", either from other theories or from everyday experience. Second, measuring models are an important tool to derive predictions or other singular statements to be used in confirming the theory.<sup>17</sup> We write  $x\langle t/t' \rangle$  to denote the result of substituting  $t'$  in  $x$  for  $t$ .

**D6** If  $x = \langle T, E, D, K, \leq^x, AS^x, C^x, \dots, f^x \rangle \in M_p$   
 and  $t \in \{\leq^x, AS^x, C^x, \dots, f^x\}$  we say that  $x$  is a  
*measuring model for  $t$*  iff  $x \in M$  and for all  $t'$  of the same

<sup>15</sup>Tversky (1977).

<sup>16</sup>A detailed account of such an approximative apparatus which works well for quantitative theories is found in Balzer, Moulines, & Sneed (1987), Chap.VII.

<sup>17</sup>Compare Balzer (1985) for a fuller account of measuring models.

type as  $t$  :

(\*) if  $x(t/t') \in M$  then  $t = t'$ .

(\*) expresses that  $t$  in  $x$  is uniquely determined (by  $M$ ): any possible other  $t'$  substituted for  $t$  in  $x$  by means of the axioms characterizing  $M$  is forced to be identical with  $t$ . So there is just one possibility for  $t$  in  $x$ .

We may ask what kinds of measuring models are possible in Freud's theory. A first answer - which some may find a bit astonishing - is that the theory provides measuring models for all its terms. Let us exemplify this for the unconscious,  $U$ .

**Theorem 3** There exist measuring models for  $U$  in Freud's theory.

**Proof:** We define a potential model as follows.  $T = \{t_0, t_1\}$ ,  $E = \{e\}$ ,  $D = \{d\}$ ,  $K = \{\{e\}\}$ , and for the relations exactly the following atomic statements are true:  $t_0 \leq t_1$ ,  $AS(e, e)$ ,  $C(t_0, e)$ ,  $N(t_0, e)$ ,  $REAL(t_0, d, e, \{e\})$ ,  $f(d) = \{e\}$ . We leave  $U$  undetermined. From these definitions it follows by D2-6 that  $U(t_1, e)$ , and by D2-4 that not( $U(t_0, e)$ ). That is,  $U = \{\{t_1, e\}\}$ , and therefore is uniquely determined by the other components. It is easy to check that all axioms of D1 and D2 are satisfied. So  $x = \langle T, \dots, U, \dots, f \rangle$  is a measuring model for  $U$ .#

The minimal example constructed in the proof exemplifies a general strategy for determining  $U$ . Find all relevant negative events, and use knowledge about association to determine all positive instances of unconsciousness by means of D2-6. Second, exclude as many events as possible from being unconscious by means of D2-4, using knowledge about the realization-relation. If this is to be carried out in a real case, the hypothetical steps will consist in making sure that the drivings occurring in D2-4 and 6 have the respective events as natural realizations. These hypotheses will refer to further theoretical knowledge about which drivings cause which natural realizations. It seems well possible that such knowledge be expressed by specializations of the present theory. The point of theorem 3 is only to show that if sufficient knowledge is available about the components different from  $U$  then  $U$  may be determined with the help of the axioms, and logics. Also, the knowledge about the other components used in the proof of theorem 3, and outlined in the previous discussion, seems to be of the sort realistically needed in applications.

Similar theorems can be proved for the other relations. This result is not too surprising. It shows that Freud's theory provides measuring models for all of its primitive relations, even for its time-ordering  $\leq$ . (With respect to  $\leq$  this once again stresses that  $\leq$  is not just the order of physical time.) This is to be expected from a theory which for the first time structures a new domain of phenomena, and which is not much concerned with operationalizing its primitives.

A final point to be discussed in connection with the formal concepts of the theory are the various forms of specializations by means of which a whole theory-net may be established on the basis of the present models. Each specialization corresponds to some special hypothesis which can be formulated in our vocabulary or in some extension of it. Such special hypotheses usually have a range of applications much more narrow than the basic axioms. We cannot go into the details of describing the specializations, we only can briefly indicate the forms they may take. The most important form of specialization concerns drivings. These may be analyzed and differentiated into various sorts. Together with this, different forms of natural realizations may be considered, and the special laws arising express which drivings have which

natural realizations. In formal terms such specializations define special forms of the function  $f$ . A second way of specializing is to introduce "weaker" forms of the unconscious, the most famous being sublimation. In these cases, negative events are not suppressed, getting unconscious, but also they are not abreacted by natural types of events: they give rise to sublimation. A third form which is not clearly stated in Freud but may be important for psychology consists of weakening the notion of a negative event. Instead of one strong event we may consider a long sequence of "weakly" negative events (like little frustrations of the same type each day) which in the long run also may lead to some form of suppression (like depression). Finally, the realization-relation may be specialized by introducing a time-scale for drivings and connecting the occurrences of drivings to those of their realizations. In the present formulation, these connections are left very weak and general.

## Intended Applications

A theory is empirical only if there exist real systems to which it is applied. These are called intended applications of the theory. In Freud's theory all the cases he reports himself will be intended applications. In addition, other cases, sufficiently similar to those mentioned by Freud usually also will be treated as intended applications. It is characteristic of empirical theories that the set  $I$  of all intended applications can not be described very sharply. The picture of how  $I$  is determined in first approximation may be drawn as follows. First, a very small subset  $I_0$  of paradigm intended applications will be described ostensively (the cases mentioned by Freud himself), and second, intended applications in general are characterized as those real systems which are sufficiently similar to members of  $I_0$ . Often, the theory itself is used to decide which systems are "sufficiently similar". In Freud's case, as in medical theories in general, intended applications always are cases of sick persons. Only the subsumption of a case of sickness under the theory will count as success, the explanation of persons being healthy is not intended. This situation creates a difficulty when we want to formulate an empirical claim. Roughly, a theory's empirical claim is that its intended applications can be subsumed under the theoretical picture as represented by the models, that is, all intended applications "are" models. But now we have two classes that might serve as the relevant theoretical models:  $M$  and  $SICK$ . If we choose  $M$  as representing the theoretical image the claim would be that all intended cases fall under the general picture described in D1 and D2. This claim clearly is inadequate, it might be true even if all the persons considered were healthy. This claim has no implications for neurosis. We obtain an adequate claim only if we use  $SICK$  as "theoretical picture". The resulting claim is that all intended cases are cases of psychically sick persons.

But now the problem occurs of what to do with  $M$ . If the empirical claim is formulated with  $SICK$ , could we not simply omit  $M$  altogether? Nothing will get lost - so it seems - because  $SICK \subseteq M$ .

This move seems tempting but there is a decisive objection. If we eliminate  $M$  we lose the "standard of health" which is an essential part of the theory's identity. For if we look at the models of psychically sick persons, how do we justify that they are models of "sick" persons? Of course, this is justified by pointing to the mechanism described by the axioms for "healthy" persons (D2) and the additional axiom for sickness (D3-2). But how can we

separate these two groups of axioms from each other? Clearly, we have to refer to the content of the respective axioms. In general, therefore, it seems necessary to include some extra standard - axioms covering healthy cases - in order to draw the distinction between healthy and sick cases.

This distinction constitutes a clear formal difference between medical theories as contrasted to theories in the natural as well as in the social sciences. In medical or therapeutical theories a further component not present in other types of theories is necessary in order to identify the theory: a component characterizing health forming the standard or the background against which we may sensibly talk about sickness. We will therefore treat both classes *M* and *SICK* as basic constituents of the theory. By summarizing all the items described by now we obtain the formal core  $K(FREUD)$  of Freud's theory.

$$K(FREUD) = \langle M_p, M, SICK, C, L \rangle$$

The core comprises all those features of the theory that can be formulated in a precise way, it represents the "theoretical picture".

The empirical claim which may be formulated with the help of this picture in first approximation is that all intended applications "are" models of psychically sick persons and in addition satisfy the constraints and links. But if *I* is a set of real systems, what is the meaning of "are" in such a statement? How can a real system "be" a set theoretic structure? The answer is: it cannot. In order to make "real systems" on the one hand and "models" of *SICK* on the other hand compatible there is only one way. We have to make further assumptions on the structure of the real systems (but of course assumptions which do not already imply that the systems are elements of *SICK*, for this would turn the empirical claim into a tautology).

Let's consider two "average" cases reported by Freud.<sup>18</sup> Katharina, a peasant girl in the alpes, suffers from depressions and anxiety. Her history turns out as follows. The girl lived with her uncle, and had to work in the house. Repeatedly, she was threatened by her uncle, but the character of the threat (sexual or punishing or otherwise) did not become clear to her. Twice she found him unexpectedly during intercourse with the maidservant. After some time she got depressed, and she suffers from anxiety on her long way to the village which she often has to go alone.

A second report is about a woman of 38 suffering from attacks of agoraphobia and mortal fear. Her history reveals a first such attack 21 years earlier, at her age of 17. Before that event she was out for preparations concerning a ball to which she was invited. Some days before, her girl-friend had died, and the event coincided with her critical days, the only ones she had in that year. When she passes by her friend's house she gets a first attack of giddiness, anxiety and swoon. She believes she will die. In the following she had similar attacks several times.

We may format both cases according to our terminology, thus specifying the instants and events important in both cases. Little information is given about association; we may assume that all important events were conscious when occurring for the first time. The negative events are clearly pointed out in both cases. By collecting these data we obtain two

<sup>18</sup>The peasant girl is found on pp. 184-95, the agoraphobia on p. 170f.

fragments of potential models, in which information, especially about  $D, U, REAL, f$ , but also about  $AS$ , is missing. Such fragments we call partial potential models.

**D7**  $x$  is a *partial potential model* of Freud's theory ( $x \in M_p$ )

iff there exist  $T, E, D, K, \leq, AS, C, N, U, REAL, f$  and

$T', E', D', K', \leq', AS', C', N', U', REAL', f'$ , such that

1)  $\langle T', \dots, REAL', f' \rangle \in M_p$ ;

2)  $x = \langle T, E, D, \dots, REAL, f \rangle$ ;

3)  $T \subseteq T', E \subseteq E', D \subseteq D', K \subseteq K', \leq \subseteq \leq', AS \subseteq AS', C \subseteq C',$   
 $N \subseteq N', U \subseteq U', REAL \subseteq REAL',$  and  $f = f'/K$   
 $(= f', \text{ restricted to } K).$

In order to close the gap between "real systems" and "models" we now assume that all intended applications are partial potential models:  $I \subseteq M_{pp}$ . This assumption has two parts. First, it amounts to assuming that the real systems considered can be described in our vocabulary; all the primitives should be interpretable in such a system. Second, it means that other features of the system, which cannot be expressed in terms of the chosen primitives, are neglected as irrelevant. The assumption does *not* imply any theoretical connection, or structure, in the system. Roughly, we may regard intended applications as descriptions of all those data that can be obtained from the underlying real system.

By adding a set  $I(FREUD)$  of intended applications, we complete our description. Freud's theory, which we call  $FREUD$ , then has the following form

$$FREUD = \langle K(FREUD), I(FREUD) \rangle.$$

Note that  $M_{pp}$  is explicitly defined in terms of  $M_p$  and therefore needs not to be mentioned in  $K(FREUD)$ .

## Empirical Claim and Confirmation

Our assumption of representing intended applications as partial potential models leads to a relatively simple formulation of the theory's empirical claim. If the "given" systems are fragments of potential models the appropriate relation to proper models is that of extension, and the empirical claim will be that each intended application can be extended to a model of a psychically sick person (such that constraints and links are satisfied). The formal definition of extension is like in D7:  $y$  is an extension of  $x$  in case all components of  $x$  are subsets of the corresponding components of  $y$ .

**D8** If  $x = \langle u_1, \dots, u_{11} \rangle$  and  $y = \langle v_1, \dots, v_{11} \rangle$  are partial potential

models of Freud's theory then  $y$  is an *extension* of  $x$  ( $x \sqsubset y$ )

iff, for all  $i \leq 11 : u_i \subseteq v_i$ . For  $X \subseteq M_{pp}$  and  $Y \subseteq M_p$  we write

" $Y \in e(X)$ " as a shorthand for "for all  $x \in X$  there is  $y \in Y$

such that  $x \sqsubset y$ ", and we then say that  $Y$  is an *extension* of  $X$ .

Note that  $M_p \subseteq M_{pp}$ , so an extension of  $x \in I(FREUD)$  can be a full potential model. We now may state the empirical claim of  $FREUD$  pertaining to sick persons in precise terms.

**D9** The *empirical claim* of  $FREUD$  is this:

There exists a set  $X$  such that



- 1)  $X \in e(I(FREUD))$ ,
- 2)  $X \subseteq SICK$ ,
- 3)  $X \in C$ ,
- 4)  $X \subseteq L$ .

That is, the intended applications in  $I(FREUD)$  can be extended such that the resulting set of extensions  $X$  ( $X \in e(I(FREUD))$ ) is a set of models of psychically sick persons ( $X \subseteq SICK$ ), satisfies the constraints ( $X \in C$ ), and satisfies the links ( $X \subseteq L$ ). This apparently simple statement in fact expresses very complex and holistic relations.

In the light of the two examples given in the previous section an intended application  $x$  typically contains data about times, events and negative events. Sometimes there may be some information about  $K, C$ , and  $AS$ , and, more rarely, about  $U, REAL$ , and  $f$ . If we assume such data to be collected for many different cases (persons) we get (part of) the set of intended applications. The empirical claim then is that further items like  $U, D$ , etc. can be "filled in" in order to obtain full models, and in a way which satisfies constraints and links. In other words, the data may be extended to full structures which satisfy all the axioms stated in D1, D2 and D4. In addition, the different association-relations and  $f$ -functions are constrained as described earlier, and all extensions have to be linked to appropriate structures (from space-time theories, and theories about association and consciousness).

The question now may be raised whether such a claim is an empirically non-trivial statement, a statement which can be tested, confirmed, and which might turn out false. This crucially depends on the particular systems making up the set  $I(FREUD)$ . If members of  $I(FREUD)$  contain "large" parts of full potential models it may be difficult to extend them to proper models of  $SICK$ . In principle there even might occur full potential models in  $I(FREUD)$  for which the question of extendability to models reduces to the question of whether the axioms of D2 are satisfied. Since we do not have at hand a complete list of all intended applications and the corresponding data which have been investigated by now, we have to resort to hypothetical reasoning: "What about the empirical claim if  $I(FREUD)$  had this and that form?" Obviously, there exist "potential falsifiers" (i.e. sets  $X \subseteq M_{pp}$  for which  $SICK$  is not an extension): take an appropriate set of full potential models which do not satisfy the axioms of D2. As long as we do not worry about where the data in  $I(FREUD)$  come from, such cases exist. However, reflection on real examples shows that, normally, only a small fragment of all the data making up a full potential model will be known. In such cases extension is a substantial procedure. If we go still one step further and assume that certain primitives, like  $U$ , will never be represented in  $I(FREUD)$ , and that extension consists in adding such missing relations, we obtain more interesting forms of possible falsifiers. In theorem 4 below we write

$$r(y) = \langle T, E, \emptyset, K, \leq, AS, C, N, \emptyset, \emptyset, \emptyset \rangle$$

if

$$y = \langle T, E, D, K, \leq, AS, C, N, U, REAL, f \rangle.$$

**Theorem 4** There exists  $X \subseteq M_{pp}$  such that

- 1) for all  $x \in X$  there is  $y \in M_p$  such that  $x = r(y)$ ;
- 2) there is no  $Y \subseteq SICK$  such that  $\{r(y) : y \in Y\} = X$ .

**Proof :** Let  $X = \{x\}$  where  $x = \langle T, E, \emptyset, K, \leq, AS, C, N, \emptyset, \emptyset, \emptyset \rangle \in M_{pp}$  such that

- (1)  $T = \{t_0, t_1\}$ ,  $E = \{e_0, e_1\}$ ,  $K = \{E\}$ , and  $t_0 < t_1$ ,
- (2) not  $(C(t_0, e_0))$ , not  $(C(t_0, e_1))$  and not  $(C(t_1, e_0))$ .  
 Suppose there exists some  $y \in SICK$  such that  $r(y) = x$ . It follows that  $y$  has the form  
 $y = \langle T, E, D, K, \leq, AS, C, N, U, REAL, f \rangle$   
 and, by D3, satisfies
- (3) there exist  $t \in T$ ,  $d \in D$  such that
  - (3.1) there is  $t^* \in T$  such that  $t < t^*$ ,
  - (3.2) for all  $t', e, k$  : not  $(t < t' \text{ and } REAL(t', d, e, k) \text{ and } f(d) = k)$ .  
 Let  $t$  and  $d$  be as in (3). From (3) and (1) it follows that
- (4)  $d \in D$  and  $t = t_0$ . As  $y \in SICK$  we obtain from (4) and D1-7 that there is some  $k$  with  $f(d) = k$ , and by (1):  $f(d) = E$ . From this, (4) and (3.2) we obtain: for all  $t', e$ , if  $t_0 < t'$  then not  $REAL(t', d, e, E)$  and from this by (1):
- (5) not:  $REAL(t_1, d, e, E)$  for all  $e \in E$ . On the other hand, (4) and D2-7 yield
- (6) there exist  $t_2, e_2, k_2$  such that  $REAL(t_2, d, e_2, k_2)$ , and this, by D2-1 yields  $C(t_2, e_2)$ . From this by (1) and (2) we obtain  $t_2 = t_1$  and  $e_2 = e_1$ , so, by (6) and (1):
- (7)  $REAL(t_1, d, e_1, E)$ . But by (5): not  $REAL(t_1, d, e_1, E)$  in contradiction to (7).#

Theorem 4 shows that there may be data which cannot be extended to proper models of *SICK* by adding suitable  $D, U, REAL$  and  $f$ . If such data occur in  $I(FREUD)$  then the empirical claim will become false. Note that the data assumed to be present in the partial potential model  $x$  in the proof of theorem 4 are very "thin". They amount to knowing that just two instants and two events are present, that there is only one type of events and that both events are not conscious at the first point of time, and one of them also is not conscious at the second. Though the example itself is contrived, the *kind* of data used certainly is not. In the light of these considerations it seems difficult to maintain that *FREUD* is not empirical.

Still, the objection may be held up that it cannot be seen how the empirical claim could be confirmed. For, it may be said, the crucial point is not whether there are potential falsifiers but rather how the data making up  $I(FREUD)$  are collected. More precisely, if all the data in  $I(FREUD)$  can be obtained only by assuming the theory to hold then the empirical claim of course will be true, and its truth will spring from the way the data are collected. Roughly, if we collect only those items that fit with the theory and call them data, we are sure that those data can be extended to models (i.e. satisfy the empirical claim). But then it seems impossible to test and to confirm the theory.

To this objection there are two replies. First, we do not think that the data are collected in a way that presupposes the theory in the sense of fully implying it. To be sure, the general picture drawn by the theory may be used but this only amounts to focussing on the "right" things, on those things to which the theory refers to. This does not mean that any determination of a datum is understood as being explicitly dependent on all the complex axioms of the theory being true in the situation of determination. We believe that a distinction

has to be made between the act of determining a datum, and "presupposing" the theory in that act. "Presupposing the theory" here is of the same kind as "presupposing a concept" (say "table") in order to be able to "see" tables. In this trivial sense of course the theory is presupposed in any act of determination, and this holds for *all* empirical theories, in physics as well as, say, in economics. The interesting question is whether the theory is fully implied by the assumptions used in order to determine some datum. In some cases this is so but in others it is not.<sup>19</sup> In Freud's theory it seems natural to use the theory's axioms in order to determine the drivings, the unconscious, the realization-relation and the function  $f$  of the models (compare theorem 3 above). If this is done those parts should not occur in the descriptions of the intended applications, and in the light of the examples considered it is unlikely that they actually do occur.

Second, it has to be pointed out that the assumption underlying the previous discussion, namely that a theory can be tested or confirmed only on the basis of independent data, does not express the dominant view of confirmation to-day. The mainstream view concedes that also coherence may be used in order to confirm and justify a theory. A precise view of confirmation along these lines exists.<sup>20</sup> According to this view, confirmation consists in the coincidence of values (of the "same" function for the "same" argument) determined or computed along different lines, irrespective of whether the theory under consideration was used in the course of these computations. This view has some credit in the area of comprehensive theories in the natural sciences.<sup>21</sup> So, even if one is not convinced by our first reply, there still remains the possibility of changing one's view about confirmation.

By way of summary let us state that there is no proof that *FREUD*'s empirical claim is correct and well confirmed. Such proof is impossible for any empirical theory. But there are hypothetical falsifiers (even of a rather realistic kind, see theorem 4), and the difficulties in collecting independent data for *FREUD* seem not very different from such difficulties in general. We therefore tend to conclude that the empirical claim of *FREUD* is not very different from such claims in general.<sup>22</sup> Of course, there are other differences between *FREUD* and theories from the natural sciences in connection with the reliability of measurement and repetition which we did not address. But these *FREUD* shares with all theories in the social sciences.

## References

- Balzer, W. (1982). *Empirische Theorien: Modelle, Strukturen, Beispiele*. Braunschweig-Wiesbaden : Vieweg.
- Balzer, W. (1985). *Theorie und Messung*. Berlin etc.: Springer.
- Balzer, W. (1986). Theoretical Terms. A New Perspective. *The Journal of Philosophy*, 83, 71-90.
- Balzer, W., Moulines, C.U., & Sneed, J.D. (1987). *An Architectonic for Science*. Dordrecht: Reidel.

<sup>19</sup>An attempt to clarify the distinction discussed here is made in Balzer (1986).

<sup>20</sup>See Glymour (1980).

<sup>21</sup>Glymour even claims that the bootstrap method was applied by Freud, see Glymour (1980), pp. 263.

<sup>22</sup>A popular argument against *FREUD* consists in pointing out the damage caused by psychotherapists in many cases. This argument, however, rests on grave misunderstanding. It does not show that *FREUD* is badly performing empirically. It just shows that it is badly used or even misused, in the sense in which physics may be said to be misused for the construction of atomic missiles.

- Freud, S. (1967). *Gesammelte Werke*. Volume 1, 5th ed. Frankfurt/Main: S.Fischer.
- Freud, S. (1963). *Das Unbewußte*. Frankfurt/Main: S.Fischer.
- Gruenbaum, A. (1984). *The Foundations of Psychoanalysis*. Berkeley: Univ. of California Press.
- Glymour, C. (1980). *Theory and Evidence*. Princeton: University Press.
- Perrez, M. (1979). *Ist die Psychoanalyse eine Wissenschaft?*. 2nd ed. Bern: Huber.
- Stegmüller, W. (1986). *Theorie und Erfahrung. Dritter Teilband*. Berlin: Springer.
- Tversky, A. (1977). Features of Similarity. *Psychological Review*, 84, 327-352.