

# Characterizing Diagnoses and Systems<sup>1</sup>

Johan de Kleer  
Xerox Palo Alto Research Center  
3333 Coyote Hill Road, Palo Alto CA 94304 USA  
Alan K. Mackworth<sup>2</sup>  
University of British Columbia  
Vancouver, B.C. V6T 1W5, Canada  
Raymond Reiter<sup>2</sup>  
University of Toronto  
Toronto, Ontario M5S 1A4, Canada

## Abstract

Most approaches to model-based diagnosis describe a diagnosis for a system as a set of failing components that explains the symptoms. In order to characterize the typically very large number of diagnoses, usually only the minimal such sets of failing components are represented. This method of characterizing all diagnoses is inadequate in general, in part because not every superset of the faulty components of a diagnosis necessarily provides a diagnosis. In this paper we analyze the concept of diagnosis in depth exploiting the notions of implicate/implicant and prime implicate/implicant. We use these notions to consider two alternative approaches for addressing the inadequacy of the concept of minimal diagnosis. First, we propose a new concept, that of kernel diagnosis, which is free of this problem with minimal diagnosis. This concept is useful to both the consistency and abductive views of diagnosis. Second, we consider restricting the axioms used to describe the system to ensure that the concept of minimal diagnosis is adequate.

## 1 Introduction

The diagnostic task is to determine why a correctly designed system is not functioning as it was intended — the explanation for the faulty behavior being that the particular system under consideration is at variance in some way with its design. One of the main subtasks of diagnosis is to determine what could be wrong with a system given the observations that have been made.

Most approaches to model-based diagnosis [6] characterize all the diagnoses for a system as the minimal sets of failing components which explain the symptoms. Although this method of characterizing diagnoses is adequate for diagnostic approaches which model only the correct behavior of components, it does not generalize. For example, it does not necessarily extend to approaches which incorporate models of faulty behavior

[28] or which incorporate strategies for exonerating components [20]. In particular, not every superset of the faulty components of a diagnosis necessarily provides a diagnosis. In this paper we analyze the notion of diagnosis in depth and consider two approaches for addressing the inadequacy of minimal diagnoses. First, we consider an alternative notion, that of kernel diagnosis, which is free of this problem with minimal diagnosis. Second, we consider restricting the axioms used to describe the system to ensure that the concept of minimal diagnosis is adequate.

## 2 Problems with minimal diagnosis

Insofar as possible we follow Reiter's [23] framework.

**Definition 1** A system is a triple  $(SD, COMPS, OBS)$  where:

1.  $SD$ , the system description, is a set of first-order sentences.
2.  $COMPS$ , the system components, is a finite set of constants.
3.  $OBS$ , a set of observations, is a set of first-order sentences.

Although our framework does not require a distinction between  $SD$  and  $OBS$ , we do so because this is the convention in the diagnosis literature.

Most model-based diagnosis papers [2; 5; 8; 9; 14; 15; 20; 23; 27; 28] define a diagnosis to be a set of failing components with all other components presumed to be behaving normally. We represent a diagnosis as a conjunction which explicitly indicates whether each component is normal or abnormal. This representation of diagnosis captures the same intuitions as the previous definitions but generalizes more naturally.

The definition of diagnosis is built up from the notion of abnormal. We adopt Reiter's [23] convention that  $AB(c)$  is a literal which holds when component  $c \in COMPS$  is abnormal. (Some of the model-based diagnosis literature uses  $\neg OK(c)$  instead of  $AB(c)$  but this is just a trivial terminological shift and does not affect the results of this paper.) It is important to note that we neither define nor place any conditions whatsoever

<sup>1</sup>Originally appeared in *Artificial Intelligence* 56(1992).

<sup>2</sup>Fellow, Canadian Institute for Advanced Research.

on how  $AB$  is used. Researchers have used varying definitions of abnormality — each of which corresponds to a different policy for how  $AB$  appears in  $SD$ . Our results apply regardless of how  $AB$  is used. A few of the ways abnormality is used in current model-based research are:

- In GDE [8], if a component violates its behavioral model, then it must be abnormal. However, if it appears to be behaving normally, then it cannot logically distinguish whether it is abnormal or not. Instead, GDE uses probabilities to rank diagnoses.
- [21] extends GDE with a non-intermittency axiom which requires that a component's outputs are a function of its inputs even if it is abnormal. One of the consequences of this axiom is that if a component is behaving normally for all its inputs, then it cannot be abnormal.
- In [20] a component is abnormal only if it violates its behavioral model at the observation time of interest.
- [22] expands the preceding notion by requiring a component to be abnormal only if it violates its behavioral model at some known observation time.

Our general diagnosis framework encompasses all these notions of abnormality. Throughout this paper we use these differing policies in examples.

**Definition 2** Given two sets of components  $C_p$  and  $C_n$  define  $\mathcal{D}(C_p, C_n)$  to be the conjunction:

$$\left[ \bigwedge_{c \in C_p} AB(c) \right] \wedge \left[ \bigwedge_{c \in C_n} \neg AB(c) \right].$$

A diagnosis is a sentence describing one possible state of the system, where this state is an assignment of the status normal or abnormal to each system component.

**Definition 3** Let  $\Delta \subseteq COMPS$ . A diagnosis for  $(SD, COMPS, OBS)$  is  $\mathcal{D}(\Delta, COMPS - \Delta)$  such that:

$$SD \cup OBS \cup \{\mathcal{D}(\Delta, COMPS - \Delta)\}$$

is satisfiable.

The following important observation follows directly from the definition (similar to proposition 3.1 of [23]):

**Remark 1** A diagnosis exists for  $(SD, COMPS, OBS)$  iff  $SD \cup OBS$  is satisfiable.

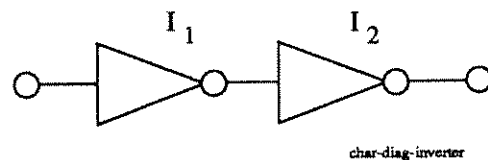
Unfortunately, there may be  $2^{|COMPS|}$  diagnoses. Therefore we seek a parsimonious characterization of the diagnoses of a system.

**Definition 4** A diagnosis  $\mathcal{D}(\Delta, COMPS - \Delta)$  is a minimal diagnosis iff for no proper subset  $\Delta'$  of  $\Delta$  is  $\mathcal{D}(\Delta', COMPS - \Delta')$  a diagnosis.

Thus a minimal diagnosis is determined by a minimal set of components which can be assumed to be faulty, while assuming the remaining components are functioning normally.

Note that these definitions subsume Reiter's [23]. Reiter's definition of the concept of diagnosis corresponds

Figure 1: Two inverters



to our notion of *minimal* diagnosis. Reiter provides no definition corresponding to our notion of a diagnosis. All the results of [23] therefore apply to our concept of a minimal diagnosis.

The following is an easy consequence of the above definitions:

**Remark 2** If  $\mathcal{D}(\Delta, COMPS - \Delta)$  is a diagnosis, then there is a minimal diagnosis  $\mathcal{D}(\Delta', COMPS - \Delta')$  such that  $\Delta' \subseteq \Delta$ .

Many approaches to model-based diagnosis have assumed that the converse holds:

**Hypothesis 1 (Minimal Diagnosis Hypothesis)** If  $\mathcal{D}(\Delta', COMPS - \Delta')$  is a minimal diagnosis and if  $\Delta' \subseteq \Delta \subseteq COMPS$ , then  $\mathcal{D}(\Delta, COMPS - \Delta)$  is a diagnosis.

As we see in section 7, the Minimal Diagnosis Hypothesis holds under the assumptions usually made. However, as we relax these assumptions, for example by allowing fault models or exoneration axioms, the Minimal Diagnosis Hypothesis fails to hold and we must explore alternative means for parsimoniously characterizing all diagnoses.

**Remark 3** The Minimal Diagnosis Hypothesis does not always hold: If  $\mathcal{D}(\Delta', COMPS - \Delta')$  is a minimal diagnosis and  $\Delta' \subseteq \Delta$ , then  $\mathcal{D}(\Delta, COMPS - \Delta)$  need not be a diagnosis.

Thus, not every superset of the faulty components of a minimal diagnosis need provide a diagnosis. To see why, consider the following two simple examples. The first example arises if we presume we know all the possible ways a component can fail such as in [28].

**Example 1** Consider the simple two inverter circuit of Fig. 1. If we are making observations at different times, then we must represent this in  $SD$  in some way. One scheme is to introduce observation time  $t$  as a parameter. Thus, a model for an inverter is:

$$INVERTER(x) \rightarrow \left[ \neg AB(x) \rightarrow [in(x, t) = 0 \equiv out(x, t) = 1] \right].$$

We assume that  $SD$  is extended with the appropriate axioms for binary arithmetic, etc. Suppose the input is 0 and the output is 1:  $in(I_1, T_0) = 0, out(I_2, T_0) = 1$ . There are three possible diagnoses:

$$\mathcal{D}(\{I_1\}, \{I_2\}) : AB(I_1) \wedge \neg AB(I_2)$$

$$\mathcal{D}(\{I_2\}, \{I_1\}) : AB(I_1) \wedge \neg AB(I_2)$$

$$\mathcal{D}(\{I_1, I_2\}, \{\}) : AB(I_1) \wedge AB(I_2)$$

These three diagnoses are characterized by the first two diagnoses, which are minimal. Suppose we know that the inverters we are using have only two failure modes: they short their output to their input or their output becomes stuck at 0. We model this as:

$$INVERTER(x) \wedge AB(x) \rightarrow [SA0(x) \vee SHORT(x)],$$

$$SA0(x) \rightarrow out(x, t) = 0,$$

$$SHORT(x) \rightarrow out(x, t) = in(x, t).$$

From these models we can infer that it is no longer possible that both  $I_1$  and  $I_2$  are faulted. Intuitively, if  $I_2$  is faulted and producing the observed 1, then it cannot be stuck at 0, and must have its input shorted to its output. But then  $I_1$  must be outputting a 1 and there is no faulty behavior of  $I_1$  which produces a 1 for an input of 0. Thus,  $AB(I_1) \wedge AB(I_2)$  is no longer a diagnosis, but the minimal diagnoses remain unchanged.

The only way to determine which of  $I_1$  or  $I_2$  is actually faulted is to make additional observations. For example, if we observed  $out(I_1, T_0)$ , we could distinguish whether  $I_1$  or  $I_2$  is faulted. Suppose  $I_1$  is faulted such that  $out(I_1, T_0) = 0$ . To identify the actual failure mode of  $I_1$  we have to observe  $out(I_1, T_1)$  or  $out(I_2, T_1)$  given  $in(I_1, T_1) = 1$ .

This example shows that the use of exhaustive fault models such as in [28] leads to difficulties with the usual definition of diagnosis. One way to avoid this difficulty is not to presume all the faulty behaviors are known as in [9]. However, if we do not know all the faulty behaviors, then nothing useful can ever be inferred from a component being abnormal which defeats the purpose of fault modes in the first place (this is addressed in [9] by introducing probabilities).

**Example 2** The usual definition of diagnosis encounters similar difficulties with the TRIAL framework of [20]. In this framework a component is considered faulty if it is actually manifesting a faulty behavior given the current set of inputs. If we are only concerned with one set of inputs, then every component is modeled with a biconditional. Thus, the inverters of Fig. 1 are instead described by:

$$INVERTER(x) \rightarrow [\neg AB(x) \equiv [in(x) = 0 \equiv out(x) = 1]].$$

Suppose the input and output are measured to be 0. There are only two diagnoses (the second of which is minimal):

$$AB(I_1) \wedge AB(I_2), \quad \neg AB(I_1) \wedge \neg AB(I_2).$$

It is not possible that one inverter is faulted and the other not. Each inverter exonerates the other. In terms of [20], each inverter is an alibi for the other. Thus, although  $\neg AB(I_1) \wedge \neg AB(I_2)$  is a minimal diagnosis,

neither  $\neg AB(I_1) \wedge AB(I_2)$  nor  $AB(I_1) \wedge \neg AB(I_2)$  are diagnoses. Again, we see that by including axioms which restrict faulty behavior in any way, the Minimal Diagnosis Hypothesis fails to hold.

In the remainder of this paper we explore two approaches to address this problem: (1) find an alternative means to characterize all diagnoses, and (2) restrict the form of  $SD \cup OBS$  such that the Minimal Diagnosis Hypothesis holds. We first require some preliminaries.

### 3 Minimal diagnoses

The minimal diagnoses are conveniently defined in terms of the familiar [18] notions of implicates and implicants (see [17; 24] for similar uses of these notions).

**Definition 5** An *AB-literal* is  $AB(c)$  or  $\neg AB(c)$  for some  $c \in COMPS$ .

**Definition 6** An *AB-clause* is a disjunction of *AB-literals* containing no complementary pair of *AB-literals*. A *positive AB-clause* is an *AB-clause* all of whose literals are positive.

Note that the empty clause is considered a positive *AB-clause*.

**Definition 7** A *conflict* of  $(SD, COMPS, OBS)$  is an *AB-clause* entailed by  $SD \cup OBS$ . A *positive conflict* is a conflict all of whose literals are positive.

If  $SD \cup OBS$  is propositional, then a conflict is any *AB-clause* which is an implicate of  $SD \cup OBS$ .

The conflicts provide an intermediate step in determining the diagnoses and are central to many diagnostic frameworks. The reason for this can be understood intuitively as follows. The diagnostic task is to determine malfunctions, and therefore the primary source of diagnostic information about a system are the discrepancies between expectations and observations. A conflict represents such a fragment of diagnostic information. For example, the conflict  $AB(A) \vee AB(B)$  might result from the discrepancy between observing  $x = 1$  while expecting it to be 2, if components  $A$  and  $B$  were normal. As a consequence, we infer that at least one of  $A$  or  $B$  is abnormal, i.e., the conflict  $AB(A) \vee AB(B)$ . Most researchers have focussed only on positive conflicts. (As most previous research has focused on the positive conflicts, they usually represented conflicts as sets of abnormal components.) However, as we see in Section 4, the non-positive conflicts are important when we model faults and do exoneration.

**Remark 4** A diagnosis exists for  $(SD, COMPS, OBS)$  iff the empty clause is not a conflict of  $(SD, COMPS, OBS)$ .

**Theorem 1** Suppose  $(SD, COMPS, OBS)$  is a system,  $\Pi$  is its set of conflicts, and  $\Delta \subseteq COMPS$ . Then  $\mathcal{D}(\Delta, COMPS - \Delta)$  is a diagnosis iff

$$\Pi \cup \{\mathcal{D}(\Delta, COMPS - \Delta)\}$$

is satisfiable.

*Proof.*  $\Rightarrow$  Consider a diagnosis  $D$ . Since  $SD \cup OBS \cup \{D\}$  is satisfiable, so is  $T \cup \{D\}$  for any set  $T$  of sentences entailed by  $SD \cup OBS$ . Since  $\Pi$  consists of clauses entailed by  $SD \cup OBS$ ,  $\Pi \cup \{D\}$  must be satisfiable.  
 $\Leftarrow$  Conversely, consider a  $\Delta \subseteq COMPS$  for which  $\Pi \cup \{D(\Delta, COMPS - \Delta)\}$  is satisfiable. Suppose  $SD \cup OBS \cup \{D(\Delta, COMPS - \Delta)\}$  is unsatisfiable. Therefore,

$$SD \cup OBS \models \neg D(\Delta, COMPS - \Delta).$$

But  $\neg D(\Delta, COMPS - \Delta)$  is an  $AB$ -clause so it must be in  $\Pi$ , contradicting the fact that  $\Pi \cup \{D(\Delta, COMPS - \Delta)\}$  is satisfiable.  $\square$

**Definition 8** A minimal conflict of  $(SD, COMPS, OBS)$  is a conflict no proper subclass of which is a conflict of  $(SD, COMPS, OBS)$ .

Thus, if  $SD \cup OBS$  is propositional, then a minimal conflict is any  $AB$ -clause which is a prime implicate of  $SD \cup OBS$ .

**Theorem 2** Suppose  $(SD, COMPS, OBS)$  is a system,  $\Pi$  is its set of minimal conflicts, and  $\Delta \subseteq COMPS$ . Then  $D(\Delta, COMPS - \Delta)$  is a diagnosis iff

$$\Pi \cup \{D(\Delta, COMPS - \Delta)\}$$

is satisfiable.

*Proof.*  $\Pi$  is logically equivalent to the set of conflicts of  $(SD, COMPS, OBS)$ . The result now follows from Theorem 1.  $\square$

**Remark 5** If all the minimal conflicts of  $(SD, COMPS, OBS)$  are non-empty and positive, then  $D(COMPS, \{\})$  is a diagnosis.

As the minimal conflicts determine the diagnoses, they play a central role in most diagnostic frameworks.

**Example 3** Consider the familiar circuit of Fig. 2. Suppose the component models are:

$$\begin{aligned} \text{ADDER}(x) &\rightarrow \\ &[\neg AB(x) \rightarrow \text{out}(x) = \text{in1}(x) + \text{in2}(x)] \\ \text{MULTIPLIER}(x) &\rightarrow \\ &[\neg AB(x) \rightarrow \text{out}(x) = \text{in1}(x) \times \text{in2}(x)] \end{aligned}$$

As before we assume that  $SD$  is extended with the appropriate axioms for arithmetic, etc. With the given inputs, there are two minimal conflicts:

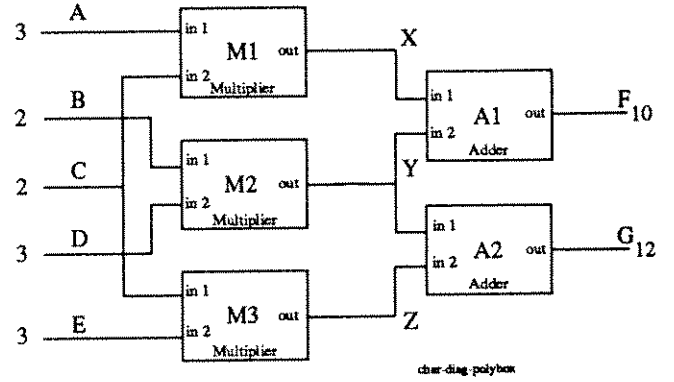
$$AB(A_1) \vee AB(M_1) \vee AB(M_2)$$

$$AB(A_1) \vee AB(M_1) \vee AB(M_3) \vee AB(A_2),$$

and four familiar minimal diagnoses:

$$\begin{aligned} D(\{A_1\}, \{A_2, M_1, M_2, M_3\}) : \\ &AB(A_1) \wedge \neg AB(A_2) \wedge \neg AB(M_1) \wedge \neg AB(M_2) \wedge \neg AB(M_3) \\ D(\{M_1\}, \{A_1, A_2, M_2, M_3\}) : \\ &AB(M_1) \wedge \neg AB(A_1) \wedge \neg AB(A_2) \wedge \neg AB(M_2) \wedge \neg AB(M_3) \\ D(\{M_2, M_3\}, \{A_1, A_2, M_1\}) : \\ &AB(M_2) \wedge AB(M_3) \wedge \neg AB(A_1) \wedge \neg AB(A_2) \wedge \neg AB(M_1) \end{aligned}$$

Figure 2:  $F = AC + BD, G = CE + BD$



$$D(\{A_2, M_2\}, \{A_1, M_1, M_3\}) : \\ AB(A_2) \wedge AB(M_2) \wedge \neg AB(A_1) \wedge \neg AB(M_1) \wedge \neg AB(M_3).$$

To prove the next two theorems we need the following lemma.

**Lemma 1** Suppose that  $\Pi$  is the set of minimal conflicts of  $(SD, COMPS, OBS)$ , and that  $\Delta$  is a minimal set such that,

$$\Pi \cup \left\{ \bigwedge_{c \in COMPS - \Delta} \neg AB(c) \right\}$$

is satisfiable. Then  $D(\Delta, COMPS - \Delta)$  is a minimal diagnosis.

*Proof.* By the minimality of  $\Delta$ , we have, for each  $c' \in \Delta$ , that

$$\Pi \cup \left\{ \bigwedge_{c \in COMPS - \Delta} \neg AB(c) \right\} \cup \{ \neg AB(c') \}$$

is unsatisfiable, i.e. for each  $c' \in \Delta$ ,

$$\Pi \cup \left\{ \bigwedge_{c \in COMPS - \Delta} \neg AB(c) \right\} \models AB(c')$$

so

$$\Pi \cup \left\{ \bigwedge_{c \in COMPS - \Delta} \neg AB(c) \right\} \models \bigwedge_{c' \in \Delta} AB(c').$$

Moreover, by hypothesis,

$$\Pi \cup \left\{ \bigwedge_{c \in COMPS - \Delta} \neg AB(c) \right\}$$

is satisfiable. Hence,  $\Pi \cup \{D(\Delta, COMPS - \Delta)\}$  is satisfiable, so by Theorem 2  $D(\Delta, COMPS - \Delta)$  is a diagnosis. It remains only to show that  $\Delta$  is a minimal set such that  $D(\Delta, COMPS - \Delta)$  is a diagnosis. But this is easy, for if  $D(\Delta', COMPS - \Delta')$  were a diagnosis for a strict subset  $\Delta'$  of  $\Delta$ , then

$$\Pi \cup \left\{ \bigwedge_{c \in COMPS - \Delta'} \neg AB(c) \right\}$$

would be satisfiable, contradicting the hypothesis of this lemma.  $\square$

**Definition 9** A conjunction  $C$  of literals covers a conjunction  $D$  of literals iff every literal of  $C$  occurs in  $D$ .

**Definition 10** Suppose  $\Sigma$  is a set of propositional formulas. A satisfiable conjunction of literals  $\pi$  (i.e., a conjunction containing no pair of complementary literals) is an implicant of  $\Sigma$  iff  $\pi$  entails each formula in  $\Sigma$ .  $\pi$  is a prime implicant of  $\Sigma$  iff the only implicant of  $\Sigma$  covering  $\pi$  is  $\pi$  itself.

**Theorem 3 (Characterization of minimal diagnoses)**  $\mathcal{D}(\Delta, COMPS - \Delta)$  is a minimal diagnosis of  $(SD, COMPS, OBS)$  iff  $\bigwedge_{c \in \Delta} AB(c)$  is a prime implicant of the set of positive minimal conflicts of  $(SD, COMPS, OBS)$ .

A proof of this theorem is given by Corollary 4.5 of [23]. The following is a direct proof in the terminology of this paper.

*Proof.*  $\Rightarrow$  Suppose  $\Pi^+$  is the set of positive minimal conflicts for  $(SD, COMPS, OBS)$ , and that  $\mathcal{D}(\Delta, COMPS - \Delta)$  is a diagnosis. By Theorem 2,  $\Pi \cup \{\mathcal{D}(\Delta, COMPS - \Delta)\}$  is satisfiable where  $\Pi$  is the set of minimal conflicts of  $(SD, COMPS, OBS)$ . Since  $\Pi^+ \subseteq \Pi$ ,  $\Pi^+ \cup \{\mathcal{D}(\Delta, COMPS - \Delta)\}$  is also satisfiable. Since  $\mathcal{D}(\Delta, COMPS - \Delta)$  contains every possible  $AB$ -literal or its negation, every clause of  $\Pi^+$  must contain a literal of  $\mathcal{D}(\Delta, COMPS - \Delta)$ . Therefore,

$$\bigwedge_{c \in \Delta} AB(c) \wedge \bigwedge_{c \in COMPS - \Delta} \neg AB(c) \models \Pi^+.$$

Since  $\Pi^+$  contains only positive literals, the negative literals are irrelevant:

$$\bigwedge_{c \in \Delta} AB(c) \models \Pi^+.$$

Since  $\mathcal{D}(\Delta, COMPS - \Delta)$  is a minimal diagnosis, no subset of  $\Delta$  has this property. Hence,  $\bigwedge_{c \in \Delta} AB(c)$  is not only an implicant but a prime implicant of  $\Pi^+$ .

$\Leftarrow$  Suppose  $\Pi$  and  $\Pi^+$  are the sets of minimal and positive minimal conflicts for  $(SD, COMPS, OBS)$ , and that  $\bigwedge_{c \in \Delta} AB(c)$  is a prime implicant of  $\Pi^+$ . We prove that  $\Delta$  is a minimal set such that

$$\Pi \cup \left\{ \bigwedge_{c \in COMPS - \Delta} \neg AB(c) \right\}$$

is satisfiable. The result will then follow from lemma 1. Suppose then that

$$\Pi \cup \left\{ \bigwedge_{c \in COMPS - \Delta} \neg AB(c) \right\}$$

is unsatisfiable, so that

$$\Pi \models \bigvee_{c \in COMPS - \Delta} AB(c)$$

which is a positive clause. Because  $\Pi$  consists of minimal conflicts, it follows that some clause of  $\Pi^+$  contains literals of  $\bigvee_{c \in COMPS - \Delta} AB(c)$ . But this cannot be since  $\bigwedge_{c \in \Delta} AB(c)$  is a prime implicant of  $\Pi^+$ . Hence

$$\Pi \cup \left\{ \bigwedge_{c \in COMPS - \Delta} \neg AB(c) \right\}$$

is satisfiable. We now prove  $\Delta$  is a minimal set with this property. Every conflict in  $\Pi^+$  has the form  $\bigvee_{c \in \Delta' \cup K} AB(c)$  for some  $\Delta' \subseteq \Delta$  and  $K \subseteq COMPS - \Delta$ . Moreover, for each  $\delta \in \Delta$ , some such conflict contains  $\mathcal{D}(\Delta, COMPS - \Delta)$  else  $\bigwedge_{c \in \Delta} AB(c)$  is not a prime implicant of  $\Pi^+$ . We prove that some conflict in  $\Pi^+$  containing  $\mathcal{D}(\Delta, COMPS - \Delta)$  must have the form  $\mathcal{D}(\Delta, COMPS - \Delta) \vee \bigvee_{k \in K} AB(k)$ . For if not, then every conflict in  $\Pi^+$  which contains  $\mathcal{D}(\Delta, COMPS - \Delta)$  must have the form  $\mathcal{D}(\Delta, COMPS - \Delta) \vee \mathcal{D}(\delta', COMPS - \delta') \vee \dots \vee \bigvee_{k \in K} AB(k)$ , where  $\delta' \in \Delta$  and  $\delta' \neq \delta$ . But then  $\bigwedge_{c \in \Delta - \{\delta\}} AB(c)$  is a smaller implicant than  $\bigwedge_{c \in \Delta} AB(c)$ , yielding a contradiction. Hence, for each  $\delta \in \Delta$  there is a conflict of the form  $\mathcal{D}(\Delta, COMPS - \Delta) \vee \bigvee_{k \in K} AB(k)$  where  $K \subseteq COMPS - \Delta$ . Hence, for each  $\delta \in \Delta$

$$\mathcal{D}(\Delta, COMPS - \Delta) \vee \bigvee_{c \in COMPS - \Delta} AB(c)$$

is a conflict so that,

$$\Pi \cup \left\{ \bigwedge_{c \in \{\delta\} \cup (COMPS - \Delta)} \neg AB(c) \right\}$$

is unsatisfiable. Since we have already proved that

$$\Pi \cup \left\{ \bigwedge_{c \in COMPS - \Delta} \neg AB(c) \right\}$$

is satisfiable,  $\Delta$  must be a minimal set with this property.  $\square$

This theorem underlies many model-based diagnostic algorithms. The first step, conflict recognition, finds positive minimal conflicts, and the second step, candidate generation, finds prime implicants. Clearly, if we were only interested in minimal diagnoses, then we would only be interested in identifying the positive minimal conflicts, but, in general, we must consider the non-positive minimal conflicts as well.

We now have the machinery to state precisely when the minimal diagnoses characterize all diagnoses.

**Theorem 4** The following are equivalent:

1. If  $\mathcal{D}(\Delta', COMPS - \Delta')$  is a minimal diagnosis for  $(SD, COMPS, OBS)$ , then  $\mathcal{D}(\Delta, COMPS - \Delta)$  is a diagnosis for  $(SD, COMPS, OBS)$  for every  $\Delta$  such that  $COMPS \supseteq \Delta \supseteq \Delta'$  (i.e., every superset of the faulty components of a minimal diagnosis provides a diagnosis).
2. All minimal conflicts of  $(SD, COMPS, OBS)$  are positive.

*Proof.*  $1 \Rightarrow 2$ . Suppose, for  $C_p, C_n \subseteq COMPS$ , that

$$\bigvee_{c \in C_p} AB(c) \vee \bigvee_{c \in C_n} \neg AB(c)$$

is a conflict of  $(SD, COMPS, OBS)$ . Then

$$\bigvee_{c \in C_p} AB(c) \vee \bigvee_{c \in COMPS - C_p} \neg AB(c)$$

is a conflict, so that the negation of this, which is  $\mathcal{D}(COMPS - C_p, C_p)$ , is not a diagnosis. We prove that  $\bigvee_{c \in C_p} AB(c)$  is a conflict, from which the result follows. Suppose not. Then  $SD \cup OBS \cup \{\neg AB(c) \mid c \in A\}$  is satisfiable. Let  $\Delta \supseteq C_p$  be a maximal subset of  $COMPS$  such that  $SD \cup OBS \cup \{\neg AB(c) \mid c \in \Delta\}$  is satisfiable. By lemma 1,  $\mathcal{D}(COMPS - \Delta, \Delta)$  is a minimal diagnosis. Since  $\Delta \supseteq C_p$ , then by property 1 of the theorem,  $\mathcal{D}(COMPS - C_p, C_p)$  is a diagnosis, contradicting our previously established result.

$2 \Rightarrow 1$ . Suppose  $\mathcal{D}(\Delta', COMPS - \Delta')$  is a minimal diagnosis and  $COMPS \supseteq \Delta \supseteq \Delta'$ . By Theorem 2, if  $\Pi$  is the set of minimal conflict of  $(SD, COMPS, OBS)$ , then  $\Pi \cup \{\mathcal{D}(\Delta', COMPS - \Delta')\}$  is satisfiable. Since for each  $c \in COMPS$  either  $AB(c)$  or  $\neg AB(c)$  occurs in  $\mathcal{D}(\Delta', COMPS - \Delta')$ , this means that every  $AB$ -clause of  $\Pi$  contains a literal of  $\mathcal{D}(\Delta', COMPS - \Delta')$  and this literal is positive since the  $AB$ -clauses are positive. Hence, because  $\Delta \supseteq \Delta'$ , each  $AB$ -clause of  $\Pi$  contains a positive literal of  $\mathcal{D}(\Delta, COMPS - \Delta)$ , so  $\Pi \cup \{\mathcal{D}(\Delta, COMPS - \Delta)\}$  is satisfiable, whence  $\mathcal{D}(\Delta, COMPS - \Delta)$  is a diagnosis.  $\square$

In Example 1,  $AB(I_1) \wedge \neg AB(I_2)$  was a diagnosis, but  $AB(I_1) \wedge AB(I_2)$ , which has more faulty components, was not. By Theorem 4 this must arise because one of the minimal conflicts is not positive. In this example, the negative clause,  $\neg AB(I_1) \vee \neg AB(I_2)$ , is a minimal conflict, which follows directly from the fault models of  $I_1$  and  $I_2$ .

## 4 Partial diagnoses

Suppose we have the following two diagnoses for a three component system:  $AB(c_1) \wedge AB(c_2) \wedge AB(c_3)$  and  $AB(c_1) \wedge AB(c_2) \wedge \neg AB(c_3)$ . We can interpret this as saying that  $c_1$  and  $c_2$  are faulty, and that  $c_3$  may or may not be faulty. Thus, the two diagnoses may be represented more compactly by  $AB(c_1) \wedge AB(c_2)$ . In fact, we can view this as a 'partial' diagnosis in which we are uncommitted to the status of  $c_3$ ; no matter what that status is, it leads to a diagnosis. This is the basis for Poole's observation [19] that a diagnosis need not commit to a status for each component whenever that status is a 'don't care'. Accordingly, we introduce the concept of a partial diagnosis. This concept also has the nice side effect of providing a convenient representation characterizing the set of all diagnoses.

**Definition 11** A partial diagnosis for  $(SD, COMPS, OBS)$  is a satisfiable conjunction  $P$  of  $AB$ -literals such that for every satisfiable conjunction  $\phi$

of  $AB$ -literals covered by  $P$ ,  $SD \cup OBS \cup \{\phi\}$  is satisfiable.

Notice that as every conjunction covers itself that all partial diagnoses are satisfiable.

The following is an easy consequence of this definition:

**Remark 6** If  $P$  is a partial diagnosis of  $(SD, COMPS, OBS)$  and  $C$  is the set of all components mentioned in  $P$ , then

$$P \wedge \bigwedge_{c \in COMPS - C} A(c)$$

is a diagnosis, where each  $A(c)$  is  $AB(c)$  or  $\neg AB(c)$ .

Thus, a partial diagnosis  $P$  represents the set of all diagnoses which contain  $P$  as a subconjunct. It is natural then to consider the minimal such  $P$ 's, which we call kernel diagnoses.

**Definition 12** A kernel diagnosis is a partial diagnosis with the property that the only partial diagnosis which covers it is itself.

The following easy result provides exactly the characterizing property we have been looking for:

**Theorem 5** (Characterization of diagnoses)  $\mathcal{D}(\Delta, COMPS - \Delta)$  is a diagnosis iff there is a kernel diagnosis which covers it.

Consider the example of Fig. 1. Without the introduction of fault models there were three diagnoses:  $AB(I_1) \wedge \neg AB(I_2)$ ,  $\neg AB(I_1) \wedge AB(I_2)$ ,  $AB(I_1) \wedge AB(I_2)$  which are characterized by the two kernel diagnoses:  $AB(I_1)$  and  $AB(I_2)$ . With the addition of the fault models, the kernel diagnoses become:  $AB(I_1) \wedge \neg AB(I_2)$  and  $\neg AB(I_1) \wedge AB(I_2)$ .

Partial and kernel diagnoses can be particularly easily characterized in terms of prime implicants and minimal conflicts. Recall that a conjunction of literals  $\pi$  containing no pair of complementary literals is an implicant of  $\Sigma$  iff  $\pi$  entails each formula in  $\Sigma$ .

**Theorem 6** The partial diagnoses of  $(SD, COMPS, OBS)$  are the implicants of the minimal conflicts of  $(SD, COMPS, OBS)$ .

*Proof.* Let  $\Pi$  be the set of all conflicts of  $(SD, COMPS, OBS)$ . Since  $\Pi$  is logically equivalent to the set of minimal conflicts of  $(SD, COMPS, OBS)$ , it is sufficient to prove that the partial diagnoses are the implicants of  $\Pi$ . As a further simplification, we appeal to the following analog of Theorem 2, whose proof is similar:

$K$  is a partial diagnosis iff  $\Pi \cup \Xi$  is satisfiable for every satisfiable conjunct  $\Xi$  of  $AB$ -literals covered by  $K$ .

$\Rightarrow$  Suppose  $K$  is a partial diagnosis. We prove  $K \models \pi$  for each  $\pi \in \Pi$ , whence  $K$  is an implicant of  $\Pi$ . Suppose not. Then  $K \not\models \pi$  for some  $\pi \in \Pi$ , which means that no literal of  $\pi$  occurs in  $K$ . Let  $\mathcal{L}$  be the set of those literals of  $\pi$  which are not complements of literals of  $K$ . Consider  $P = K \wedge \bigwedge_{l \in \mathcal{L}} \neg l$ .  $P \cup \pi$  is unsatisfiable. But  $K$  covers  $P$ , contradicting the fact that  $K$  is a partial diagnosis.

$\Leftarrow$  Suppose that  $K$  is an implicant of  $\Pi$ . We prove  $K$  is a partial diagnosis. Since  $K \models \Pi$  for each  $\pi \in \Pi$ ,  $C \models \pi$  for each satisfiable conjunct  $C$  of  $AB$ -literals covered by  $K$ . Hence  $\Pi \cup \{C\}$  is satisfiable for any such  $C$  so that  $K$  is a partial diagnosis.  $\square$

**Corollary 1** (*Characterization of kernel diagnoses*)  
The kernel diagnoses of  $(SD, COMPS, OBS)$  are the prime implicants of the minimal conflicts of  $SD \cup OBS$ .

*Proof.* Let  $\Pi$  be the set of minimal conflicts of  $(SD, COMPS, OBS)$ .

$\Rightarrow$  If  $K$  is a kernel diagnosis, then by Theorem 6 it is an implicant of  $\Pi$ . We prove it is prime. If not, then for some  $C$  distinct from  $K$  but covering  $K$ ,  $C \models \pi$  for each  $\pi \in \Pi$ . Hence, for every satisfiable conjunct  $D$  of  $AB$ -literals covered by  $C$ ,  $D \models \pi$ . Thus  $\Pi \cup \{D\}$  is satisfiable for each such  $D$ , which means that  $K$  is not a kernel diagnosis, contradiction.

$\Leftarrow$  Suppose  $K$  is a prime implicant of  $\Pi$ . Then by Theorem 6 it is a partial diagnosis. Suppose  $K$  is not a kernel diagnosis. Then there is a conjunct  $C$  covering  $K$  but distinct from  $K$  such that  $C$  is a partial diagnosis. By Theorem 6,  $C$  is an implicant of  $\Pi$ , contradicting the fact that  $K$  is a prime implicant of  $\Pi$ .  $\square$

As a consequence of this corollary and Theorem 3, if all minimal conflicts are positive, then there is a simple one-to-one correspondence between minimal diagnoses and kernel diagnoses.

Corollary 1 provides a direct way of computing the kernel diagnoses from the minimal conflicts (if  $SD$  is propositional then the minimal conflicts can be computed by a prime implicate algorithm, otherwise more sophisticated inferential machinery must be brought to bear). One way of doing this is to convert the CNF-form of the minimal conflicts to DNF and simplify as follows (we omit the proof):

1. 'Multiply' the minimal conflicts to give a disjunction of conjunctions.
2. Delete any conjunction containing a complementary pair of literals.
3. Delete any conjunction covered by some other conjunction.
4. The remaining conjunctions are the prime implicants of the original minimal conflicts, and hence the kernel diagnoses.

**Example 4a** Consider Example 3. There are two minimal conflicts:

$$AB(A_1) \vee AB(M_1) \vee AB(M_2)$$

$$AB(A_1) \vee AB(M_1) \vee AB(M_3) \vee AB(A_2),$$

and four kernel diagnoses:

$$AB(A_1)$$

$$AB(M_1)$$

$$AB(M_2) \wedge AB(M_3)$$

$$AB(M_2) \wedge AB(A_2).$$

As all minimal conflicts are positive, these diagnoses correspond one-to-one to the familiar minimal diagnoses.

**Example 4b** Suppose we used slightly different component models:

$$ADDER(x) \rightarrow$$

$$[\neg AB(x) \equiv [out(x) = in1(x) + in2(x)]]$$

$$MULTIPLIER(x) \rightarrow$$

$$[\neg AB(x) \equiv [out(x) = in1(x) \times in2(x)]].$$

In this case the minimal conflicts become:

$$AB(A_1) \vee AB(M_1) \vee AB(M_2)$$

$$AB(A_1) \vee AB(A_2) \vee AB(M_1) \vee AB(M_3)$$

$$AB(A_2) \vee \neg AB(M_2) \vee AB(M_3)$$

$$AB(A_2) \vee AB(M_2) \vee \neg AB(M_3)$$

$$\neg AB(A_2) \vee AB(M_3) \vee AB(M_2),$$

and the kernel diagnoses become:

$$\neg AB(A_2) \wedge AB(M_1) \wedge \neg AB(M_2) \wedge \neg AB(M_3)$$

$$AB(A_2) \wedge AB(M_1) \wedge AB(M_3)$$

$$AB(A_1) \wedge \neg AB(A_2) \wedge \neg AB(M_2) \wedge \neg AB(M_3)$$

$$AB(A_1) \wedge AB(A_2) \wedge AB(M_3)$$

$$AB(A_2) \wedge AB(M_2)$$

$$AB(M_2) \wedge AB(M_3).$$

Note that because the positive minimal conflicts are unchanged, the set of minimal diagnoses remains unchanged.

In this example there are only a few more kernel diagnoses than minimal diagnoses (6 vs. 4). However, one possible disadvantage of this approach is that there may sometimes be exponentially more kernel diagnoses than diagnoses.

It is interesting to note that the set of minimal conflicts may be redundant. In Example 4b, the first and third minimal conflicts entail the second:

$$AB(A_1) \vee AB(M_1) \vee AB(M_2)$$

$$AB(A_2) \vee \neg AB(M_2) \vee AB(M_3)$$

---


$$AB(A_1) \vee AB(A_2) \vee AB(M_1) \vee AB(M_3)$$

Therefore, the second minimal conflict is redundant. Such redundancy can only occur if there are non-positive minimal conflicts. Unfortunately, these observations do not seem to be of much practical use because there is no easy way to tell whether there are enough minimal conflicts without first finding them all.

**Definition 13** A set of kernel diagnoses is *irredundant* iff it is a smallest cardinality set with the property that every diagnosis is covered by at least one of its elements.

**Theorem 7** If all minimal conflicts are positive there is exactly one irredundant set of kernel diagnoses; namely the set of all kernel diagnoses.

A system can have multiple irredundant sets of kernel diagnoses.

**Example 5** Consider a circuit having three components  $A, B, C$  and the two minimal conflicts:

$$AB(A) \vee AB(B) \vee AB(C) \\ \neg AB(A) \vee \neg AB(B) \vee \neg AB(C)$$

These have six prime implicants (i.e., kernel diagnoses).

$$AB(A) \wedge \neg AB(B) \\ \neg AB(A) \wedge AB(C) \\ AB(B) \wedge \neg AB(C) \\ \neg AB(A) \wedge AB(B) \\ AB(A) \wedge \neg AB(C) \\ \neg AB(B) \wedge AB(C)$$

There are two irredundant sets of kernel diagnoses:

$$\{AB(A) \wedge \neg AB(B), \neg AB(A) \wedge AB(C), AB(B) \wedge \neg AB(C)\} \\ \{\neg AB(A) \wedge AB(B), AB(A) \wedge \neg AB(C), \neg AB(B) \wedge AB(C)\}.$$

Our analysis of kernel diagnoses corresponds to the classical analysis in switching theory of so-called two-level minimization of boolean functions (e.g., the Quine-McCluskey algorithm [16; 18]). The problem there is to synthesize a circuit realizing a given function as a disjunction of conjunctions of literals in such a way as to minimize the number of conjunctions and literals. Such circuits are characterized by irredundant sets of prime implicants of the given function. In the case of diagnosis, the given boolean function is specified by  $\Pi$ , the set of conflicts of  $SD \cup OBS$ . The kernel diagnoses are the prime implicants of  $\Pi$ , and the minimal sets of kernel diagnoses sufficient to cover every diagnosis are the irredundant sets of prime implicants of  $\Pi$ . It is well known from switching theory that the minimization problem is computationally intractable; there may be too many prime implicants, and even if there aren't, finding an irredundant subset of them is NP-hard. Therefore, designers of VLSI circuits have developed various approximation techniques [1]. Because of the correspondence with diagnosis, we can expect to profit from these techniques.

It can be useful to construct irredundant sets of partial diagnoses containing non-kernel diagnoses. For example, for probability calculations it is useful (as far as possible) to ensure that no two of the partial diagnoses have a common superset. The probability calculus of [8; 9; 10; 20] computes the probabilities of outcomes by combining the probabilities of partial diagnoses. For example, if some outcome holds in two diagnoses  $A$  and  $B$  then its probability is:

$$P(A \vee B) = P(A) + P(B) - P(A \wedge B)$$

If  $A$  and  $B$  have no common superset, then  $P(A \wedge B) = 0$ . This can result in an exponential speed up in the probability calculations.

## 5 Prime diagnoses

Raiman [20] proposes a notion of prime diagnosis to characterize diagnoses. In his TRIAL architecture, components are individually incriminated and exonerated. Therefore, he characterizes the diagnoses of a system in terms of the diagnoses involving its individual components. The following is a generalization of his definition.

**Definition 14** Given  $(SD, COMPS, OBS)$ , a prime diagnosis for  $c \in COMPS$  is a minimal diagnosis for  $(SD, COMPS, OBS \cup \{AB(c)\})$ .

Prime diagnoses characterize all diagnoses as follows.

**Theorem 8 (Raiman)** Suppose  $\mathcal{D}(\Delta, COMPS - \Delta)$  is a diagnosis. Then for each  $c_i \in \Delta$  there is a prime diagnosis  $\mathcal{D}(\Delta_i, COMPS - \Delta_i)$  for  $c_i$  such that  $\Delta = \bigcup_i \Delta_i$ .

Unfortunately, Example 1 shows that not every combination of prime diagnoses leads to a diagnosis. The prime diagnoses are:

$$P(I_1) = \{AB(I_1) \wedge \neg AB(I_2)\} \\ P(I_2) = \{AB(I_2) \wedge \neg AB(I_1)\}$$

However,  $AB(I_1) \wedge AB(I_2)$  is not a diagnosis. Thus, prime diagnoses are inadequate to characterize diagnoses.

Raiman [20] implicitly assumes all minimal conflicts contain at most one negative literal. In this case Raiman shows that the converse of Theorem 8 holds which makes prime diagnoses adequate for characterizing diagnoses. This useful property holds if  $SD \cup OBS$  is Horn, but we do not know of any more general practical condition on  $SD \cup OBS$  which ensures it.

## 6 Abductive diagnoses

An alternative to the consistency-based approach is to define diagnosis in terms of abduction [3; 4; 12; 19]. In order to do so we must differentiate those observations which are about inputs from those which are about outputs. The intuition is that we sometimes want the diagnoses not only to be consistent with the observations, but to also predict the outputs given the inputs. Using the logical framework we have laid out thus far, it is straight-forward to develop a characterization of abductive diagnoses.

In order to define the notion of abductive diagnosis we must distinguish between those sentences in  $OBS$  which are about inputs,  $I$ , from those which are about outputs,  $O$ . The terms "inputs", "outputs" and "diagnoses" are here being used generically. Abduction in general appeals to a built-in asymmetry based in part on a distinction between cause and effect. In performing abductive reasoning on causal systems, the observations to be explained are taken to be effects of causal factors; these causes are treated as though they are part of  $SD$ . So for circuits, outputs would be the results of measurements, while circuit inputs, which are the normal



causes of the outputs are treated as though they were in  $SD$ . In a medical setting, the "diagnoses" might be diseases (measles, malaria), while the "outputs" might be symptoms (fever, dizziness) and the "inputs" might be perturbations to the system, such as diet or lab tests. These observations about abduction are intended as a guide to formulating the contents of  $SD$ ,  $I$  and  $OBS$  in order to achieve intuitively satisfying results — but our framework and its conclusions apply whatever the contents of  $SD$ ,  $I$  and  $OBS$ .

**Definition 15** Let  $\Delta \subseteq COMPS$ ,  $OBS = I \cup O$ . An abductive diagnosis for  $(SD, COMPS, OBS)$  is  $\mathcal{D}(\Delta, COMPS - \Delta)$  such that:

$$SD \cup I \cup \{\mathcal{D}(\Delta, COMPS - \Delta)\}$$

is satisfiable, and

$$SD \cup I \cup \{\mathcal{D}(\Delta, COMPS - \Delta)\} \models O.$$

**Definition 16** A partial abductive diagnosis for  $(SD, COMPS, OBS)$  is a satisfiable conjunction  $P$  of AB-literals such that for every satisfiable conjunction  $\phi$  of AB-literals covered by  $P$ ,  $SD \cup I \cup \{\phi\}$  is satisfiable and  $SD \cup I \cup \{\phi\} \models O$ .

**Definition 17** A kernel abductive diagnosis is a partial abductive diagnosis with the property that the only partial abductive diagnosis which covers it is itself.

The following comes almost directly from the definitions:

**Remark 7** Every partial abductive diagnosis of  $(SD, COMPS, I \cup O)$  is also a partial diagnosis of  $(SD, COMPS, I \cup O)$ . The converse is not in general true.

**Remark 8** (Characterization of abductive diagnoses)  $\mathcal{D}(\Delta, COMPS - \Delta)$  is an abductive diagnosis iff there is a kernel abductive diagnosis which covers it.

**Definition 18** Suppose  $\Sigma$  is a set of first order sentences. A conjunction of ground literals  $\pi$  containing no pair of complementary literals is an implicant of  $\Sigma$  iff  $\pi$  entails each sentence in  $\Sigma$ . A satisfiable conjunction of ground literals  $\pi$  is a prime implicant of  $\Sigma$  iff the only implicant of  $\Sigma$  covering  $\pi$  is  $\pi$  itself.

Using the framework developed in this paper we can relate kernel abductive diagnoses to prime implicants, at least for finite axiomatizations:

**Theorem 9** Suppose  $SD$ ,  $I$  and  $O$  are finite sets (so that we can treat each of these as a sentence consisting of the conjunction of its elements). A conjunction  $K$  of AB-literals is a kernel abductive diagnosis of  $(SD, COMPS, I \cup O)$  iff  $K$  is a prime implicant of  $\Pi \wedge \{SD \wedge I \rightarrow O\}$ , where  $\Pi$  is the conjunction of the minimal conflicts of  $(SD, COMPS, I \cup O)$ .

*Proof.*  $\Leftarrow$  Consider any satisfiable conjunction  $\phi$  of AB-literals covered by  $K$ . Then  $\{\phi\} \models \Pi \wedge \{SD \wedge I \rightarrow O\}$ , in which case  $\{\phi\} \models \Pi$  and

$$\{\phi\} \models SD \wedge I \rightarrow O. \quad (1)$$

By Theorem 6,  $\phi$  is a partial diagnosis of  $(SD, COMPS, I \cup O)$ , and thus  $\{\phi\} \cup SD \cup I$  is satisfiable. Moreover, by (1),  $\{\phi\} \cup SD \cup I \models O$ . Hence, by definition,  $K$  is a partial abductive diagnosis. We must prove that  $K$  is a kernel abductive diagnosis. To that end, suppose  $K'$  is a partial abductive diagnosis of  $(SD, COMPS, I \cup O)$  which covers  $K$ . By Remark 7,  $K'$  is a partial diagnosis of  $(SD, COMPS, I \cup O)$ , whence by Theorem 6,  $\{K'\} \models \Pi$ . Moreover,  $\{K'\} \models SD \wedge I \rightarrow O$  by virtue of being a partial abductive diagnosis of  $(SD, COMPS, I \cup O)$ . Hence,  $K'$  is an implicant of  $\Pi \wedge \{SD \wedge I \rightarrow O\}$ . Since  $K$  is a prime implicant of  $\Pi \wedge \{SD \wedge I \rightarrow O\}$ ,  $K' = K$ . Thus  $K$  must be a kernel abductive diagnosis of  $(SD, COMPS, I \cup O)$ .

$\Rightarrow$  We first prove that  $K$  is an implicant of  $\Pi \wedge \{SD \wedge I \rightarrow O\}$ . Since  $K$  is a kernel (and hence partial) abductive diagnosis,  $K$  is a partial diagnosis of  $(SD, COMPS, I \cup O)$ , by Remark 7. By Theorem 6,  $\{K\} \models \Pi$ . Moreover,  $K \cup SD \cup I \models O$  since  $K$  is a partial abductive diagnosis. Hence,  $\{K\} \models SD \wedge I \rightarrow O$  and  $\{K\} \models \Pi$  so that  $K$  is an implicant of  $\Pi \wedge \{SD \wedge I \rightarrow O\}$ . Next we show that  $K$  is prime. Suppose  $K'$  is an implicant of  $\Pi \wedge \{SD \wedge I \rightarrow O\}$  which covers  $K$ , and let  $\phi$  be any satisfiable conjunction of AB-literals covered by  $K'$ . Then  $\phi$  is an implicant of  $\Pi \wedge \{SD \wedge I \rightarrow O\}$ . Since  $\phi$  is an implicant of  $\Pi$ , by Theorem 6,  $\{\phi\} \cup SD \cup I \cup O$  is satisfiable, whence so also is  $\{\phi\} \cup SD \cup I$ . Moreover,  $\{\phi\} \cup SD \cup I \models O$ . Therefore,  $K'$  is a partial abductive diagnosis. Since  $K'$  covers  $K$  and  $K$  is kernel,  $K' = K$  so that  $K$  is prime.  $\square$

Poole [19] has developed a very particular definition of "abductive diagnosis" which differs from that of definition 15. To prevent confusion we refer to his definition as P-abductive diagnoses.

**Definition 19** An P-abductive diagnosis of  $(SD, COMPS, I \cup O)$  is a conjunction  $P$  of AB-literals such that: (1)  $SD \cup I \cup P$  is satisfiable, (2)  $SD \cup I \cup P \models O$ , and (3) it is not covered by some other P-abductive diagnosis.

This definition is different than the three notions we have just seen. P-abductive diagnoses are not abductive diagnoses as they do not include an AB-literal for every component. Although partial diagnoses do not include an AB-literal for every component, they are not minimal. Although kernel diagnoses are minimal, Poole's definition does not require that every other conjunction of AB-literals covered by it is also an P-abductive diagnosis.

P-abductive diagnoses do not characterize the space of abductive diagnoses. Nevertheless, with the definitions we have developed it is possible to state precisely what P-abductive diagnoses are in terms of prime implicants.

**Theorem 10** A conjunction of AB-literals  $K$  is an P-abductive diagnosis of  $(SD, COMPS, I \cup O)$  iff  $K$  is a prime implicant of  $SD \cup I \rightarrow O$  and  $SD \cup I \cup \{K\}$  is satisfiable.

*Proof.*  $\Leftarrow$  Let  $K$  be a conjunction of AB-literals which is a prime implicant of  $SD \cup I \rightarrow O$  and  $SD \cup I \cup K$

is satisfiable. Consider any  $\phi$  covered by  $K$ . By the definition of cover,  $\{\phi\} \models SD \cup I \rightarrow O$ . If  $\{\phi\} \models SD \cup I \rightarrow O$ , then  $\{\phi\} \cup SD \cup I \models O$  by the deduction theorem.  $SD \cup I \cup \{\phi\}$  is satisfiable. As  $K$  is a prime implicant, it is not covered by any other conjunction of AB-literals meeting these two conditions. Thus,  $K$  meets the three conditions for P-abductive diagnosis.

$\Rightarrow$  Let  $\phi$  be a P-abductive diagnosis. By definition of P-abductive diagnosis we know that  $SD \cup I \cup \{\phi\}$  is satisfiable and that  $SD \cup I \cup \{\phi\} \models O$ . By the deduction theorem,  $\{\phi\} \models SD \cup I \rightarrow O$ . Hence,  $\phi$  is an implicant of  $SD \cup I \rightarrow O$ . As the only conjunction of AB-literals which covers  $\phi$  and meets these conditions is  $\phi$  itself,  $\phi$  is a prime implicant of  $SD \cup I \rightarrow O$ .  $\square$

## 7 Restricting the system description

Our overall objective is to find methods of characterizing all diagnoses. We saw that minimal diagnoses were inadequate for this task in general and we examined kernel and prime diagnoses as alternatives. Another approach is to restrict the form of the system so that the Minimal Diagnosis Hypothesis holds. We know from Theorem 4 that a necessary and sufficient condition ensuring that every superset of the faulty components of a minimal diagnosis provides a diagnosis is that all minimal conflicts be positive. Unfortunately, we are not aware of any simple necessary and sufficient condition on the syntactic form of a system which ensures that all minimal conflicts are positive. Clearly both  $OBS$  and  $SD$  need to be restricted because definition 1 allows non-positive  $AB$ -clauses to be part of  $OBS$  and  $SD$ . In this section we explore some commonly used practical restrictions on  $OBS$  and  $SD$  that suffice to ensure that the Minimal Diagnosis Hypothesis holds. In these definitions we assume that  $OBS$  and  $SD$  can be expressed as a set of first-order clauses.

**Definition 20** *The Ignorance of Abnormal Behavior (IAB) condition holds for a system  $(SD, COMPS, OBS)$  if in the clausal form of  $SD \cup OBS$  every occurrence of an  $AB$ -predicate is positive.*

For example, if all axioms of  $SD$  in which  $AB$  appears follow the schema:

$$\neg AB(x) \wedge A_1 \wedge \dots \wedge A_n \rightarrow C_1 \vee \dots \vee C_m,$$

which is equivalent to the clause,

$$AB(x) \vee \neg A_1 \vee \dots \vee \neg A_n \vee C_1 \vee \dots \vee C_m,$$

where the  $A_i$  and  $C_i$  are literals not mentioning  $AB$ , and if every  $AB$ -literal (if any) in  $OBS$  is positive, then IAB holds. The IAB condition is used in all of the model-based diagnosis frameworks which rely on knowing only the correct behavior of components (where the  $A_i$  specify the component type(s) and the  $C_i$  specify the various possible normal behavior modes for the component). For example,

$$\neg AB(x) \wedge TRANSISTOR(x) \rightarrow \\ ON(x) \vee OFF(x) \vee SATURATED(x).$$

**Theorem 11** *If  $(SD, COMPS, OBS)$  satisfies the IAB condition and  $\mathcal{D}(\Delta, COMPS - \Delta)$  is a diagnosis for  $(SD, COMPS, OBS)$ , then  $\mathcal{D}(\Delta', COMPS - \Delta')$  is a diagnosis for  $(SD, COMPS, OBS)$  for every  $\Delta' \supset \Delta$  where  $\Delta' \subseteq COMPS$ . In particular, the Minimal Diagnosis Hypothesis holds for  $(SD, COMPS, OBS)$ .*

*Proof.* If  $AB$  only appears positively in  $SD \cup OBS$ , then only positive minimal conflicts are possible. The result now follows from Theorem 4.  $\square$

The converse of this theorem is false. A less restrictive and more useful definition is:

**Definition 21** *The Limited Knowledge of Abnormal Behavior Condition (LKAB) holds for a system  $(SD, COMPS, OBS)$  if for every component  $c \in COMPS$  and any  $\mathcal{D}(C_p, C_n)$  where  $c \notin C_p$  and  $c \notin C_n$  and  $C_p, C_n \subseteq COMPS$  that if  $SD \cup OBS \cup \{AB(c)\}$  and  $SD \cup OBS \cup \{\mathcal{D}(C_p, C_n)\}$  are satisfiable, then  $SD \cup OBS \cup \{\mathcal{D}(C_p \cup \{c\}, C_n)\}$  is satisfiable.*

As shown later in Theorem 12, the LKAB condition provides a general characterization of a class of systems for which there is insufficient knowledge of abnormal behavior to rule out any diagnosis implicating a set of faulty components given a diagnosis implicating a subset of them.

**Remark 9** *If  $(SD, COMPS, OBS)$  satisfies the IAB condition, then it satisfies the LKAB condition.*

*Proof.* Consider each  $AB(c) \in COMPS$ . If  $AB$  occurs only positively in  $SD \cup OBS$ , then  $AB(c)$  cannot appear negatively in any minimal conflict. Thus,  $SD \cup OBS \cup \{AB(c)\}$  is always satisfiable. And, therefore, if  $SD \cup OBS \cup \{\mathcal{D}(C_p, C_n)\}$  is satisfiable where  $c \notin C_p$  and  $c \notin C_n$ , then  $SD \cup OBS \cup \{\mathcal{D}(C_p, C_n)\} \cup \{AB(c)\}$  is satisfiable.  $\square$

**Theorem 12** *If  $(SD, COMPS, OBS)$  satisfies the LKAB condition and  $\mathcal{D}(\Delta, COMPS - \Delta)$  is a diagnosis for  $(SD, COMPS, OBS)$ , then  $\mathcal{D}(\Delta', COMPS - \Delta')$  is a diagnosis for  $(SD, COMPS, OBS)$  for every  $\Delta' \supset \Delta$  where  $\Delta' \subseteq COMPS$  and for each  $c \in \Delta' \setminus \Delta$   $SD \cup OBS \cup \{AB(c)\}$  is satisfiable.*

*Proof.* Consider a diagnosis  $\mathcal{D}(\Delta, COMPS - \Delta)$  and each  $c \in COMPS - \Delta$  for which  $c \in \Delta' \setminus \Delta$   $SD \cup OBS \cup \{AB(c)\}$  is satisfiable. If  $\mathcal{D}(\Delta, COMPS - \Delta)$  is a diagnosis, then  $\{\mathcal{D}(\Delta, COMPS - \Delta)\} \cup SD \cup OBS$  is satisfiable by definition of diagnosis. Then, by LKAB  $\{\mathcal{D}(\Delta, COMPS - \Delta - \{c\})\} \cup AB(c) \cup SD \cup OBS$  is satisfiable and hence  $\{\mathcal{D}(\Delta \cup \{c\}, COMPS - \Delta - \{c\})\} \cup SD \cup OBS$  is also. By iterating this process we prove the theorem.  $\square$

Intuitively, this theorem shows that if a system obeys LKAB and no component can be proved correct, then the Minimal Diagnosis Hypothesis holds for that system.

Sherlock [9] exploits the LKAB condition. In Sherlock all axioms in  $SD$  mentioning  $AB$  have one of the following two forms:

$$\neg AB(x) \wedge A(x) \rightarrow G_1(x) \vee \dots \vee G_m(x)$$

$$AB(x) \wedge A(x) \rightarrow F_1(x) \vee \dots \vee F_m(x) \vee U(x)$$

where  $G_i(x)$  describes a possible normal behavior for component  $x$ ,  $F_i(x)$  describes a possible faulty behavior for a component  $x$ .  $U(x)$  specifies an unknown behavior so the only occurrences of the literal  $U(x)$  are in clauses of the form,  $\neg A(x) \vee \neg U(x) \vee \neg G_i(x)$  and  $\neg A(x) \vee \neg U(x) \vee \neg F_i(x)$ . Furthermore,  $G_i(x), F_j(x), U(x)$  only occur negatively in other clauses.

We show that by using resolution, a complete inference procedure, the LKAB conditions are met. Consider every  $AB(c) \in COMPS$ . We only need focus on those conclusions which follow from the axioms in which  $AB$  appears negatively. Notice that every axiom in which  $AB(c)$  appears negatively,  $U(c)$  appears positively. Consider the only two other types of clauses in which  $U(c)$  appears. The clause  $\neg A(c) \vee \neg U(c) \vee \neg F_i(c)$  contains the negations of two of the literals of the problematic  $AB$ -clause, therefore these two clauses do not resolve with each other. However, the problematic  $AB$ -clause can resolve with  $\neg A(c) \vee \neg U(c) \vee \neg G_i(c)$  to produce

$$AB(c) \wedge A(c) \rightarrow F_1(c) \vee \dots \vee F_m(c) \vee \neg G_i(c).$$

$G_i(c)$  only appears positively in clauses containing  $\neg AB(c)$ , therefore these clauses cannot resolve as well. As there are no other possible resolutions and the only sentences containing  $AB$  in OBS are atomic, the addition of  $AB(c)$  can never make some  $\mathcal{D}(Cp, Cn)$  unsatisfiable unless  $\neg AB(c) \in OBS$ . Thus LKAB holds.

For example, in Sherlock the axioms mentioning  $AB$  for an inverter are:

$$\neg AB(x) \wedge INVERTER(x) \rightarrow G(x),$$

$$AB(x) \wedge INVERTER(x) \rightarrow S1(x) \vee S0(x) \vee U(x).$$

And some of the other axioms for inverters are:

$$\neg INVERTER(x) \vee \neg G(x) \vee \neg S1(x)$$

$$\neg INVERTER(x) \vee \neg G(x) \vee \neg S0(x)$$

$$\neg INVERTER(x) \vee \neg G(x) \vee \neg U(x)$$

$$\neg INVERTER(x) \vee \neg S0(x) \vee \neg S1(x)$$

$$\neg INVERTER(x) \vee \neg S0(x) \vee \neg U(x)$$

$$\neg INVERTER(x) \vee \neg S1(x) \vee \neg U(x)$$

$$INVERTER(x) \wedge G(x) \rightarrow [IN(x) = 0 \equiv OUT(x) = 1]$$

$$INVERTER(x) \wedge S1(x) \rightarrow OUT(x) = 1$$

$$INVERTER(x) \wedge S0(x) \rightarrow OUT(x) = 0$$

From a purely logical point of view these clauses which mention  $U(x)$  convey no information, however, in the Sherlock framework every behavioral mode is assigned a probability and  $U(x)$  behavioral modes are typically assigned very small probability.

## 8 Summary

The notions of minimal and prime diagnosis are inadequate to characterize diagnoses generally. We argue that the notion of kernel diagnosis which designates some components as normal, others abnormal, and the remainder as being either, is a better way to characterize diagnoses. We avoid significant complexity if kernel diagnoses contain only positive literals (i.e., all minimal conflicts are positive). This can be achieved by limiting the description of the system to obey the IAB or LKAB condition which formalize the intuitions underlying many existing diagnosis systems.

Although there are many algorithms to compute prime implicate/implicants [13; 17; 26; 29], the task is NP-hard and experience has been that most diagnostic tasks have a large number of minimal conflicts and kernel diagnoses (or prime diagnoses, or minimal diagnoses). Therefore, the brute-force application of the techniques suggested by this paper is not practical. In practice, some focussing strategy must be brought to bear. One approach is to exploit hierarchical information as in [15]. Another approach is to focus the reasoning to identify the most relevant conflicts in order to find the most probable diagnoses [9; 11]. However, both of these approaches require additional information: the structural hierarchy and probabilistic information.

The central contribution of this paper is that it provides a clear formal framework for characterizing the space of diagnoses which also corrects some of the problems of [23]. It thus provides the specification for an ideal diagnostician and clarifies why systems such as GDE [8] work. This paper establishes the connection between diagnosis and the notions of prime implicate/implicant. The connection between prime implicates/implicants and the ATMS [7] has been presented elsewhere [24; 25]. Thus, we have constructed a logical bridge from a formal theory of diagnosis to the ATMS techniques that many diagnosis implementations use.

## 9 Acknowledgments

The contents of this paper benefitted from many discussions with Olivier Raiman. Daniel G. Bobrow, David Poole, Brian C. Williams, Vijay Saraswat and Jeffrey Siskind provided extensive insights on early drafts.

## References

- [1] Brayton, R.K., Hachtel, G.D., McMullen, C.T. and Sangiovanni-Vincentelli, A.L., *Logic minimization algorithms for VLSI Synthesis*, (Kluwer, 1984).
- [2] Brown, J.S., Burton, R. R. and de Kleer, J., Pedagogical, natural language and knowledge engineering techniques in SOPHIE I, II and III, in: D. Sleeman and J.S. Brown (Eds.), *Intelligent Tutoring Systems*, (Academic Press, New York, 1982) 227-282. An expansion of the relevant sections of this paper appears in these readings under the title, Model-based diagnosis in SOPHIE III.

- [3] Console, L., and Torasso, P., Integrating models of correct behavior into abductive diagnosis, in: *Proceedings ECAI-90*, Stockholm, Sweden (1990) 160-166.
- [4] Console, L., Theseider Dupre, D., and Torasso, P., On the relationship between abduction and deduction, *Journal of Logic and Computation* 1 (1991) 661-690.
- [5] Davis, R., Diagnostic Reasoning based on structure and behavior, *Artificial Intelligence* 24 (1984) 347-410. Appears in these readings.
- [6] Davis, R., and Hamscher, W., Model-based reasoning: Troubleshooting, in *Exploring artificial intelligence*, edited by H.E. Shrobe and the American Association for Artificial Intelligence, (Morgan Kaufmann, 1988), 297-346. Appears in these readings.
- [7] de Kleer, J., An assumption-based truth maintenance system, *Artificial Intelligence* 28 (1986) 127-162. Also in *Readings in NonMonotonic Reasoning*, edited by Matthew L. Ginsberg, (Morgan Kaufmann, 1987), 280-297. March 23, 1992
- [8] de Kleer, J. and Williams, B.C., Diagnosing multiple faults, *Artificial Intelligence* 32 (1987) 97-130. Also in *Readings in NonMonotonic Reasoning*, edited by Matthew L. Ginsberg, (Morgan Kaufmann, 1987), 372-388. Appears in these readings.
- [9] de Kleer, J. and Williams, B.C., Diagnosis with behavioral modes, in: *Proceedings IJCAI-89*, Detroit, MI (1989) 1324-1330. Appears in these readings.
- [10] de Kleer, J., Focusing on probable diagnoses, in: *Proceedings AAAI-91*, Anaheim, CA (1991) 842-848. Appears in these readings.
- [11] Dressler, O., and Farquhar, A., Focusing ATMS-based problem solvers, Siemens Report INF-2-ARM 13, 1989.
- [12] El Ayeb, B., Marquis, P., and Rusinowitch, M., A new diagnosis approach by deduction and abduction, in: *Proceedings: International workshop on expert systems in engineering*, Vienna (1990), Lecture notes in artificial intelligence 462, Springer Verlag, 32-46. Long version available as: Report 91-LE-LAB-002, Department of Mathematics, Université de Sherbrooke, Canada and as a Technical Report 91-R-146, CRIN-CNRS/INRIA-Lorraine, 1991.
- [13] Forbus, K. and J. de Kleer, *Building Problem Solvers*, (MIT Press, 1992).
- [14] Genesereth, M.R., The use of design descriptions in automated diagnosis, *Artificial Intelligence* 24 (1984) 411-436. Appears in these readings.
- [15] Hamscher, W.C., Modeling Digital Circuits for Troubleshooting, *Artificial Intelligence* 51 (1991) 223-271. Appears in these readings.
- [16] Hill, F.J. and Peterson, G.R., *Introduction to Switching Theory and Logical Design* (John Wiley and Sons, New York, 1974).
- [17] Kean, A. and Tsiknis, G., An incremental method for generating prime implicants/implicates, *Journal of Symbolic Computation* 9 (1990) 185-206.
- [18] Kohavi, Z., *Switching and Finite Automata Theory* (McGraw-Hill, 1978).
- [19] Poole, D., Representing knowledge for logic-based diagnosis, *Proc. Int. Conf. on Fifth Generation Computer Systems* (1988) 1282-1290.
- [20] Raiman, O., Diagnosis as a trial: The alibi principle, IBM Scientific Center, 1989. Appears in these readings.
- [21] Raiman O., de Kleer, J., Saraswat, V., Shirley M., Characterizing Non-Intermittent Faults, in: *Proceedings AAAI-91*, Anaheim, CA (1991) 849-854. Appears in these readings.
- [22] Raiman, O., A circumscribed diagnosis engine, in: *Proceedings: International workshop on expert systems in engineering*, Vienna (1990), Lecture notes in artificial intelligence 462, Springer Verlag, 90-101.
- [23] Reiter, R., A theory of diagnosis from first principles, *Artificial Intelligence* 32 (1987) 57-95. Also in *Readings in NonMonotonic Reasoning*, edited by Matthew L. Ginsberg, (Morgan Kaufmann, 1987), 352-371. Appears in these readings.
- [24] Reiter, R. and de Kleer, J., Foundations of Assumption-Based Truth Maintenance Systems: Preliminary Report, *Proceedings of the National Conference on Artificial Intelligence*, Seattle, WA (July, 1987), 183-188.
- [25] Selman, B. and H.J. Levesque, Abductive and Default Reasoning: A Computational Core, in: *Proceedings AAAI-90* Boston, MA (1990) 343-348.
- [26] Slagle, J.R., C.L. Chang, and Lee, R.C.T., A new algorithm for generating prime implicants, *IEEE Transactions on Computers* C-19 (1970) 304-310.
- [27] Struss, P., Extensions to ATMS-based Diagnosis, in: J.S. Gero (ed.), *Artificial Intelligence in Engineering: Diagnosis and Learning*, Southampton, 1988.
- [28] Struss, P., and Dressler, O., "Physical negation" — Integrating fault models into the general diagnostic engine, in: *Proceedings IJCAI-89* Detroit, MI (1989) 1318-1323. Appears in these readings.
- [29] Tison, P., Generalized consensus theory and application to the minimization of boolean functions, *IEEE transactions on electronic computers* 4 (August 1967) 446-456.