

# Abstract Theorem Proving \*

Fausto Giunchiglia  
Mechanised Reasoning Group  
IRST  
Povo, I 38100 Trento  
Italy  
fausto@irst.uucp

Toby Walsh  
Department of Artificial Intelligence  
University of Edinburgh  
80 South Bridge, Edinburgh  
Scotland  
Toby\_Walsh@uk.ac.edinburgh

## Abstract

*Informally*, abstraction can be described as the process of mapping a representation of a problem into a new representation. The aim of the paper is to propose a *theory* of abstraction. The generality of the framework is tested by formalizing and analyzing some work done in the past; its efficacy by giving a procedure which solves the "*false proof*" problem by avoiding the use of inconsistent abstract spaces.

## 1 Introduction

Abstraction has been suggested as a very powerful technique for constraining search in automated reasoning. *Informally*, abstraction can be described as the process of mapping a representation of a problem (also called the "*ground?*" representation) into a new representation (also called the "*abstract?*" representation) which *preserves certain desirable properties* and *is simpler to handle*. The "*desirable properties*" amount to requiring that the abstract solution be of help in solving the problem in the original search space. The notion of "*simplicity*" depends on the application, it may mean decidability or lower complexity.

As far as we know, no comprehensive theory of abstraction has been given. The only work in this direction [Plaisted, 1981] is concerned with one form of abstraction and is limited to the area of resolution theorem proving. This has caused the lack of a satisfactory characterization and general understanding of abstraction.

The aim of the work (partially) described in this paper is to provide a theory of abstraction and use it to: (1) understand what "to abstract" means and how it can be formalized; (2) classify the various forms of abstraction; (3) investigate their formal properties and the operations which can be defined on them; (4) analyze and classify past work; (5) define ways of building "useful

\*This work was begun when the first author was working at the Department of Artificial Intelligence at Edinburgh University supported by SERC grant GR/E/4459.8. The second author is supported by a SERC studentship. The research described in this paper owes a lot to the openness and sharing of ideas which exists in the Mathematical Reasoning group. The authors thank Alan Bundy, Enrico Giunchiglia, Alex Simpson and Richard Weyhrauch for the many discussions on the topic. Alan Bundy is also thanked for reading early versions of the paper.

abstractions" and (6) study how the proof in the abstract space can be used to aid the proof in the ground space.

In this paper, for lack of space, only some issues are discussed and proofs are only outlined and the simplest not given (for a more complete treatment see [Giunchiglia and Walsh, Forthcoming 89]). Only topics (1), (2), (3) and (4) are (partially) dealt with. The following of the paper can be structured in three main parts.

In the first part (section 2), the formal framework is presented and some of the underlying motivations are discussed (topics (1) and (2)). Abstraction is first defined as a mapping between formal systems and then classes of abstraction are identified depending on how certain properties (ie. provability, inconsistency) are preserved by the mapping.

In the second part (section 3), some examples/ case studies of previous work in abstraction are presented [Sacerdoti, 1974, Plaisted, 1981, Hobbs, 1985] (topic (4)). The goal here is to motivate the formal framework by showing how it can be used to capture and formalise most of the relevant previous work in various areas of AI<sup>1</sup>. This allows us to get an unified view of work which, on the surface, seems very different. For instance it is proven that the theory of granularity presented by Hobbs in [Hobbs, 1985] and described in example 3 can be formalised as a particular case of the weak and ordinary abstractions defined by Plaisted in [Plaisted, 1981] and described in example 2.

In the third and last part (section 4), it is then shown how the framework can be actually used to understand the properties of abstraction mappings and to find solutions to existing problems (topic (3)). In particular the "*false proof*" problem is investigated. Intuitively stated, the false proof problem is as follows. In order to build an abstract space "simpler" than the ground one, the trick is to forget some "irrelevant" details<sup>2</sup>. This, Plaisted noticed [Plaisted, 1981], may cause problems. In particular the abstract space may be inconsistent even if the ground space is not. It is proven that this problem cannot be avoided as it is always possible to find a set of

<sup>1</sup>This case study analysis is performed in much more depth in [Giunchiglia and Walsh, Forthcoming 89]; in [Giunchiglia and Walsh, 1989] it is shown how the framework can be effectively used to build global strategies for the unfolding of definitions.

<sup>2</sup>Where "irrelevant" should be read as "irrelevant according to same theorem proving strategy".

consistent axioms whose abstraction is inconsistent. A solution is then proposed based on the ordering of abstractions into a partial order.

## 2 The formal framework

**Definition 1 (Formal system)** : A formal system  $\Sigma$  is a triple  $\langle \Lambda, \Delta, \Omega \rangle$ , where  $\Lambda$  is the Language,  $\Omega$  is the set of axioms and  $\Delta$  is the Deductive Machinery of  $\Sigma$ .

The language  $\Lambda$  is composed of an alphabet, the set of (well formed) terms and the set of well formed formulae (wffs from now on).  $\Omega$  is a subset of the wffs of  $\Lambda$ . The deductive machinery is a set of rules of inference for deriving theorems from axioms.

**Definition 2 (Abstraction)** If  $\Sigma_1 = \langle \Lambda_{\Sigma_1}, \Omega_{\Sigma_1}, \Delta_{\Sigma_1} \rangle$  and  $\Sigma_2 = \langle \Lambda_{\Sigma_2}, \Omega_{\Sigma_2}, \Delta_{\Sigma_2} \rangle$  are two formal systems, an abstraction mapping  $f$ , written also  $f : \Sigma_1 \mapsto \Sigma_2$ , is a triple of total functions  $\langle f_{\Lambda}, f_{\Omega}, f_{\Delta} \rangle$  such that:

$$\begin{aligned} f_{\Lambda} : \Lambda_{\Sigma_1} &\mapsto \Lambda_{\Sigma_2} \\ f_{\Omega} : \Omega_{\Sigma_1} &\mapsto \Omega_{\Sigma_2} \\ f_{\Delta} : \Delta_{\Sigma_1} &\mapsto \Delta_{\Sigma_2} \end{aligned}$$

If not explicitly stated to the contrary we assume that  $f_{\Lambda}$  and  $f_{\Omega}$  agree on the axioms; that is, for any wff  $\omega$ , if  $\omega \in \Omega$ , then  $f_{\Lambda}(\omega) = f_{\Omega}(\omega)$ <sup>3</sup>. When no confusion arises we drop the subscripts. Given a deduction tree  $\Pi_{\Sigma_1}$  of  $\vdash_{\Sigma_1} \varphi_{\Sigma_1}$  in  $\Sigma_1$ , we indicate by  $f(\Pi_{\Sigma_1})$  a deduction tree  $\Pi_{\Sigma_2}$  of  $\vdash_{\Sigma_2} f(\varphi_{\Sigma_1})$ .

**Definition 3 (T\*-abstractions)** : An abstraction  $f : \Sigma_1 \mapsto \Sigma_2$  is said to be a<sup>4</sup>

**TC-Abstraction** iff, for any wff  $\varphi_{\Sigma_1}$ ,  $\vdash_{\Sigma_1} \varphi_{\Sigma_1}$  iff  $\vdash_{\Sigma_2} f(\varphi_{\Sigma_1})$ ;

**TD-Abstraction** iff, for any wff  $\varphi_{\Sigma_1}$ , if  $\vdash_{\Sigma_2} f(\varphi_{\Sigma_1})$  then  $\vdash_{\Sigma_1} \varphi_{\Sigma_1}$ ;

**TI-Abstraction** iff, for any wff  $\varphi_{\Sigma_1}$ , if  $\vdash_{\Sigma_1} \varphi_{\Sigma_1}$  then  $\vdash_{\Sigma_2} f(\varphi_{\Sigma_1})$ .

It can be easily proved that any TC-abstraction is both a TD-abstraction and a TI-abstraction.

We write "T\*-abstraction" to mean any of the above abstractions,  $TH(\Sigma)$  to mean the set of wffs provable in  $\Sigma$  and  $NTH(\Sigma)$  to mean the set of wffs which, if added to the axioms of  $\Sigma$ , make it inconsistent.

TC-abstractions map all the elements of  $TH(\Sigma_1)$  into elements of  $TH(\Sigma_2)$  and these are all and only the elements of  $TH(\Sigma_2)$ . Herbrand's theorem [Goldfarb, 1971] can be formalized as a TC-abstraction. TC-abstractions are used, for instance, in decision theory, under the name of *reduction methods*, to prove the decidability of and to build deciders for the validity problem for certain subclasses of the first order calculus [Dreben and Goldfarb, 1979]. The trick is to find a class whose decidability is known and prove that there is a proof of a wff iff there is a proof of the "translated" wff in the new class.

<sup>3</sup>To be precise, since we distinguish between wffs occurring as axioms and as anything else, we should consider occurrences of wffs and not wffs. Since, in this paper, for any wff  $\omega$ , if  $\omega \in \Omega$ , then  $f_{\Lambda}(\omega) = f_{\Omega}(\omega)$ , to make things simpler, we consider  $f_{\Lambda}$  and  $f_{\Delta}$  to range over wffs.

<sup>4</sup>"T" stands for theorem, "C" for constant, "D" for decreasing and "I" for increasing.

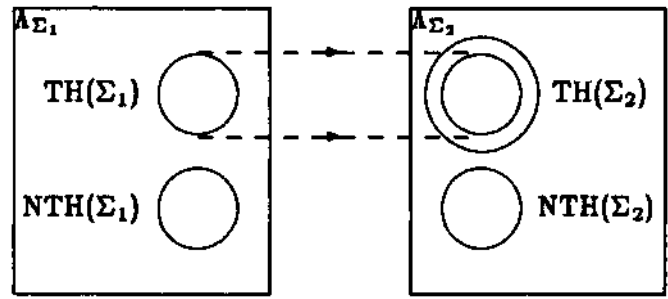


Figure 1: TI-abstraction

In TD-abstractions a subset of the elements of  $TH(\Sigma_1)$  is mapped into  $TH(\Sigma_2)$  and these are all the elements of  $TH(\Sigma_2)$ . TD-abstractions are used, for instance, to implement derived inference rules [Giunchiglia and Giunchiglia, 1988] and, as alternatives to TI-abstractions, to overcome some of their problems [Tenenber, 1987] (see later).

In TI-abstractions all the elements of  $TH(\Sigma_1)$  are mapped into a subset of  $TH(\Sigma_2)$ . TI-abstractions have been mostly used in "abstract theorem proving" (see figure 1). The main idea underlying the use of these abstractions is to prove the abstracted theorem in  $\Sigma_2$  (which, supposedly, should be simpler than in  $\Sigma_1$ ) and then to use the structure of the proof in  $\Sigma_2$  to shape the proof in  $\Sigma_1$ . The fact that there is a proof in  $\Sigma_2$  does not guarantee that there is a proof in  $\Sigma_1$ .

T\*-abstractions are classified on how provability is preserved between the ground space and the abstract space; they are thus useful when the deductive machinery is defined to generate theorems. On the other hand there are formal systems (i.e. resolution) whose deductive machinery determines inconsistency. In these cases, abstractions must be classified on how inconsistent formal systems are mapped. This requires the definition of new classes of abstractions, called NT\*-abstractions. Thus, for instance, NTI-abstractions are defined as follows<sup>5</sup>:

**Definition 4** ; An abstraction  $f : \Sigma_1 \mapsto \Sigma_2$  is an NTI-Abstraction iff, for any wff  $\varphi_{\Sigma_1}$ , if adding  $\varphi_{\Sigma_1}$  to the axioms of  $\Sigma_1$  yields an inconsistent formal system, then adding  $f(\varphi_{\Sigma_1})$  to the axioms of  $\Sigma_2$  yields an inconsistent formal system.

Various properties, equivalences, and relationships among T\*- and NT\*-abstractions can be proved [Giunchiglia and Walsh, Forthcoming 89]. NTI-abstractions\* behaviour can be represented as in figure 2.

In this paper we mention only one result which allows us to prove how and to what extent T\*-abstractions (and in particular TI-abstractions) can be used in resolution based theorem provers. Note that, if a formal system  $E$  has negation, then, for any wff  $a$ ,  $a \in TH(L)$  iff  $\neg a \in NTH(E)$ . Thus trivially:

**Corollary 1** : If  $\Sigma_1$  and  $\Sigma_2$  are two formal systems with negation and if  $f : \Sigma_1 \mapsto \Sigma_2$  is a TI-abstraction

<sup>5</sup>NTC-abstractions and NTD-abstractions are defined analogously to TC-abstractions and TD-abstractions respectively, but preserving inconsistency instead of theoremhood (see definitions 3, 4).

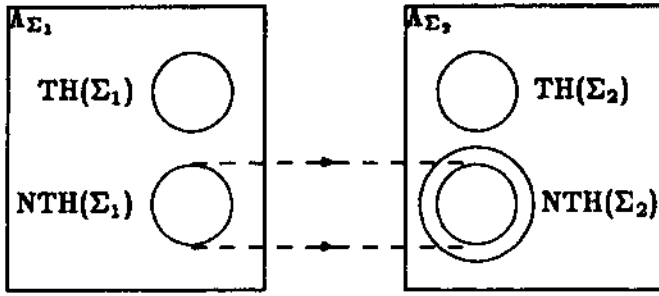


Figure 2: NTI-abstractions (Falseful abstractions)

then, for any  $\alpha$ , if  $\alpha \in TH(\Sigma_1)$  then  $\neg\alpha \in NTH(\Sigma_1)$  and  $\neg f(\alpha) \in NTH(\Sigma_2)$ .

TI-abstractions can be used in resolution theorem proving as long as, in the abstract space, the wff added to the axioms is the negation of the abstraction of the wff whose negation is added to the axioms of the ground space. This makes TI-abstractions very useful since often refutation systems, take as input a goal formula  $\alpha$  (usually automatically) negate it, add the result to the axioms and try to prove that the resulting theory is inconsistent. In the abstract space, instead of  $f(\neg\alpha)$ ,  $\neg f(\alpha)$  is added as an axiom.

On the other hand there are TI-abstractions which can be used in resolution-based systems independently of whether the negation in front of the goal is abstracted.

**Definition 5 (Negation preserving)** : An abstraction  $f : \Sigma_1 \mapsto \Sigma_2$  is negation preserving iff  $f(\neg\alpha) = \neg f(\alpha)$ .

**Theorem 1** : If  $\Sigma_1 = \langle \Lambda_1, \Omega_1, \Delta_1 \rangle$  and  $\Sigma_2 = \langle \Lambda_2, \Omega_2, \Delta_2 \rangle$  are two formal systems with negation, then a negation preserving abstraction  $f : \Sigma_1 \mapsto \Sigma_2$  is a TI-abstraction iff  $f' : \Sigma'_1 \mapsto \Sigma'_2$  is a NTI-abstraction, where  $\Sigma'_1 = \langle \Lambda_1, \Omega_1, \Delta'_1 \rangle$ ,  $\Sigma'_2 = \langle \Lambda_2, \Omega_2, \Delta'_2 \rangle$  and  $\Delta'_1$  and  $\Delta'_2$  are such that  $TH(\Sigma'_1) = TH(\Sigma_1)$  and  $TH(\Sigma'_2) = TH(\Sigma_2)$ .

Examples are  $f = f'$  with  $\Sigma_1, \Sigma_2, \Sigma'_1$  and  $\Sigma'_2$  being natural deduction, and  $f \neq f'$  with  $\Sigma_1, \Sigma_2$  being natural deduction and  $\Sigma'_1, \Sigma'_2$  being resolution. As far as we know, all the abstractions proposed to work in resolution systems are negation preserving. However, there are useful abstractions which are not negation preserving (for instance when negation is not part of the language of  $\Sigma_1$  or  $\Sigma_2$  [Newell et al., 1963], or only partially preserved by the mapping).

### 3 Some examples of abstraction

The purpose of these examples, together with providing a rational reconstruction of the work described, is to convince the reader that the framework is very powerful and allows us to present an unified view of the work done in different areas and with different goals. For lack of space, only three examples are reported, more "historical" examples are reported in [Giunchiglia and Walsh, Forthcoming 89].

**EXAMPLE 1 [Planning]: Abstrips** [Sacerdoti, 1974] was one of the first noticeable applications of abstraction. In Abstrips the preconditions of operators were abstracted according to their *criticality*. To formalize

Strips-like planning we shall adopt a situation calculus in a natural deduction formal system. Let us consider the abstraction  $f_{AB}$  where  $\Sigma_1$  and  $\Sigma_2$  are situation calculi with a first order language,  $\Omega$  consists of frame, operator and theoretic axioms and  $\Delta$  consists of natural deduction rules of inference. Operators are wffs of the form " $\forall s. (\bigwedge_{1 \leq i \leq n} p_i(s) \rightarrow q(f(s)))$ " where  $p_i$  is a precondition,  $s$  is a state of the world,  $f$  is some action, and  $q$  describes the new state of the world. Goals are wffs of the form " $\exists s \tau(s)$ ".  $f_{AB}$  is applied to wffs and axioms as follows:

$f_{AB}(\alpha) = \alpha$  if  $\alpha$  is an atomic formula.

$f_{AB}(\neg\alpha) = \neg f_{AB}(\alpha)$ ;

$f_{AB}(\alpha \circ \beta) = f_{AB}(\alpha) \circ f_{AB}(\beta)$ , where " $\circ$ " is " $\wedge$ " or " $\vee$ ";

$f_{AB}(\|x \alpha) = \|x f_{AB}(\alpha)$ , where " $\|$ " is " $\exists$ " or " $\forall$ ";

$f_{AB}(\alpha \rightarrow \beta) = f_{AB}(\alpha) \rightarrow f_{AB}(\beta)$ , provided " $\alpha \rightarrow \beta$ " is not an operator;

$f_{AB}(\bigwedge_{1 \leq i \leq n} p_i(s) \rightarrow \tau) = \bigwedge_{i \in \text{crit}(\kappa)} p_i(s) \rightarrow f_{AB}(\tau)$ , provided that " $\bigwedge_{1 \leq i \leq n} p_i(s) \rightarrow \tau$ " is an operator, where  $i \in \text{crit}(\kappa)$  if the criticality of  $p_i$  is greater than  $\kappa$ .

**Theorem 2** :  $f_{AB}$  is TI, namely, if  $\vdash_{\Sigma_1} \varphi_{\Sigma_1}$ , then  $\vdash_{\Sigma_2} f_{AB}(\varphi_{\Sigma_1})$ .

**Proof**[Outline]: By proving that given a deduction tree  $\Pi_{\Sigma_1}$  of  $\vdash_{\Sigma_1} \varphi_{\Sigma_1}$ , we can build a deduction tree  $\Pi_{\Sigma_2} = f_{AB}(\Pi_{\Sigma_1})$  of  $\vdash_{\Sigma_2} f_{AB}(\varphi_{\Sigma_1})$ . The proof proceeds by induction on the weight  $N^6$  of  $\Pi_{\Sigma_1}$ . For proofs of weight 1,  $f_{AB}$  is applied to the single wff; this generates a valid proof in  $\Pi_{\Sigma_2}$ . Assume that we have a deduction tree up to weight  $N$ . Any rule application that is not modus ponens involving an operator translates unmodified, in the sense that, for instance, an " $\forall I$ " on  $\varphi$  in  $\Pi_{\Sigma_1}$  becomes an " $\forall I$ " on  $f_{AB}(\varphi)$  in  $\Pi_{\Sigma_2}$ . For an operator application, the following transformation is performed:

$$\frac{\frac{\frac{\Pi_1}{\bigwedge_{1 \leq i \leq n} p_i} \quad \bigwedge_{1 \leq i \leq n} p_i \rightarrow q}{q}}{f_{AB}\left(\frac{\Pi_1}{\bigwedge_{1 \leq i \leq n} p_i}\right)} \quad \frac{\bigwedge_{i \in \text{crit}(\kappa)} p_i \quad \bigwedge_{i \in \text{crit}(\kappa)} p_i \rightarrow f_{AB}(q)}{f_{AB}(q)}}{f_{AB}(q)}$$

The abstract proof is valid since  $f_{AB}(\dots)$  is a valid deduction tree from the induction hypothesis, and the hypothesis of the (abstract) operator axiom is obtained from " $\bigwedge_{1 \leq i \leq n} p_i$ " by a (possibly empty) sequence of applications of "and-elimination".  $\square$

Note that the abstract proof is longer than the one in the ground space. The purpose of abstracting is not to find these longer proofs; we hope that there are also going to be shorter proofs. These shorter proofs are those that don't try to satisfy  $p_j$  for  $j \notin \text{crit}(\kappa)$ . There is no guarantee that there will be a shorter proof than the one exhibited; we will always be able to devise an obtuse theory in which to prove the  $p_i$  for  $i \in \text{crit}(\kappa)$  we have to prove all the other  $p_j$  for  $j \notin \text{crit}(\kappa)$ .  $\spadesuit$

<sup>6</sup>The *weight* of a deduction tree is the number of its formula occurrences.

**EXAMPLE 2** [Resolution theorem proving, logic programming]: The work by Plaisted is closest in spirit to ours. Plaisted proposes two classes of abstraction, *ordinary abstractions* and *weak abstractions* [Plaisted, 1981], which map a set of clauses onto a set of clauses and preserve inconsistency. His work is less general than ours as: *he restricts his attention to resolution systems and his classes of abstraction fail to capture all NTI-abstraction mappings that preserve inconsistency between resolution systems*. In other words, Plaisted's abstractions are NTI, but not all NTI-abstractions are weak or ordinary<sup>7</sup>. Moreover we claim that our definitions of abstraction are "more natural" in the sense that better reflect and capture the functionalities they are given for.

Ordinary abstractions are described as taking both  $\Sigma_1$  and  $\Sigma_2$  to be first order calculi with  $\Lambda_\Sigma$  allowing clausal form,  $\Delta_\Sigma$  being resolution and  $\Omega_\Sigma$  being arbitrary. Any ordinary abstraction mapping  $f$  maps a clause in  $\Lambda_{\Sigma_1}$  onto a set of clauses in  $\Lambda_{\Sigma_2}$  subject to the following conditions:

- a)  $f(\perp) = \{\perp\}$ ;
- b) if  $\alpha_3$  is a resolvent of  $\alpha_1$  and  $\alpha_2$  in  $\Sigma_1$ , and  $\beta_3 \in f(\alpha_3)$  then there exist  $\beta_2 \in f(\alpha_2)$  and  $\beta_1 \in f(\alpha_1)$  such that a resolvent of  $\beta_1$  and  $\beta_2$  subsumes  $\beta_3$  in  $\Sigma_2$ ;
- c) if  $\alpha_1$  subsumes  $\alpha_2$  in  $\Sigma_1$ , then for every  $\beta_2 \in f(\alpha_2)$  there exists  $\beta_1 \in f(\alpha_1)$  such that  $\beta_1$  subsumes  $\beta_2$  in  $\Sigma_2$ .

Weak abstractions are identically defined to ordinary abstractions except condition b) is weakened to the property that if  $\alpha_3$  is a resolvent of  $\alpha_1$  and  $\alpha_2$  in  $\Sigma_1$ , and  $\beta_3 \in f(\alpha_3)$  then there exist  $\beta_2 \in f(\alpha_2)$  and  $\beta_1 \in f(\alpha_1)$  such that either  $\beta_1$  subsumes  $\beta_3$ , or  $\beta_2$  subsumes  $\beta_3$ , or a resolvent of  $\beta_1$  and  $\beta_2$  subsumes  $\beta_3$  in  $\Sigma_2$ .

**Theorem 3** : *Weak and ordinary abstractions are NTI.*

**Proof:** The proof is a corollary to Theorem 2.5 on page 55 of [Plaisted, 1981].  $\square$

**Theorem 4** : *There exist NTI-abstractions between resolution systems that are not weak or ordinary abstractions.*

**Proof**[Outline]: We can find NTI-abstractions that fail every one of the three conditions in the definition of weak and ordinary abstractions. Condition a) is failed by the NTI-abstraction  $f$  such that, for any wff  $\varphi$  in  $\Sigma_1$ ,  $f(\varphi) = \{\varphi \vee \perp\}$ . The problem with condition b) is that we may also need to resolve with an axiom of the theory. Consider, for instance, the abstraction defined by  $f(p \vee q) = \{p \vee r\}$  and  $f(\varphi) = \{\varphi\}$  otherwise. If  $\Sigma_1$  contains the axioms,  $\neg q$ , and  $\neg r$  then  $f$  is NTI. In particular,  $p \vee q$  resolves with  $\neg p$  in  $\Sigma_1$  to give  $q$ . However, no clause in the abstraction of  $p \vee q$ , or  $\neg p$  (or their resolvent) subsumes the clause  $q$  found in the abstraction of  $q$ . For condition c), consider the abstraction defined by  $f(p \vee q) = \{r, p \vee q\}$  and  $f(\varphi) = \varphi$  otherwise. Now  $f$  is NTI. However,  $f$  fails condition c) of the definition of weak and ordinary abstractions as  $p$  subsumes  $p \vee q$  but no clause in the abstraction of  $p$  subsumes  $r$  which is in  $f(p \vee q)$ .  $\square$

<sup>7</sup>All Plaisted's examples of abstraction are negation preserving and thus also TI.

The definition of weak and ordinary abstractions could be extended to overcome the first counter-example by replacing condition a) with the more general requirement that  $\exists \varphi \in f(\perp). \vdash_{\Sigma_1} \neg \varphi$ . However, this still leaves useful NTI-abstractions that fail conditions b) and c). For example, if  $p_0 \leftrightarrow p_i$  for many  $i$  we might abstract many clauses of the form  $p_i \vee q$  onto the one clause  $\{p \vee q\}$ . One could argue that ordinary and weak abstractions have the advantage, over NTI-abstractions, that they always map into simpler theories, in the sense that there is always an abstract proof that is no longer than the shortest proof of the unabstracted theorem [Plaisted, 1981]. This does not seem a good point since, first of all, we intuitively expect NTI-abstractions (that are not NTC) to have this or similar properties and, second and more importantly, there are NTI-abstractions, which are not weak or ordinary, which build simpler theories (the last example is one possible case).  $\spadesuit$

**EXAMPLE 3** [Common sense reasoning]: In [Hobbs, 1985], Hobbs presents a theory of granularity in which a complex theory is abstracted onto a simpler, more "coarse-grained" theory with a smaller domain; for example, we could map the real world of continuous time and positions onto a (micro)world of discrete time and positions. Hobbs' granularity theory can be formalized as a mapping (let us call it " $f_{gran}$ ") that can be proved to be TI. Let us suppose that both  $\Sigma_1$  and  $\Sigma_2$  are calculi with a first order language, an arbitrary set of axioms and any complete deductive machinery for first order logic.  $f_{gran}$  maps different objects in  $\Sigma_1$  into (not necessarily different) objects in  $\Sigma_2$  according to an *indistinguishability relation*, defined by the (second-order) axiom:

$$\forall x, y. x \sim y \leftrightarrow \forall p \in R. p(x) \leftrightarrow p(y)$$

where  $R$  is the subset of the predicates of the theory determined to be *relevant* to the situation at hand<sup>8</sup>. Thus  $f_{gran}$  keeps the same logical structure of wffs and translates any constant into its equivalence class, namely  $f_{gran}(p(a)) = p(\kappa(a))$  where  $a$  is any constant symbol and  $\kappa(a)$  is the equivalence class of the constant  $a$  wrt the indistinguishability relation; that is  $\kappa(x) = \{y : x \sim y\}$ .

**Theorem 5** :  *$f_{gran}$  is TI/NTI.*

**Proof** By mapping a proof tree  $\Pi_{\Sigma_1}$  in  $\Sigma_1$  of  $\varphi_{\Sigma_1}$  (possibly of  $\perp$ ) onto a proof tree  $\Pi_{\Sigma_2}$  in  $\Sigma_2$  of  $f(\varphi_{\Sigma_1})$ . The proof proceeds by induction on the weight of  $\Pi_{\Sigma_1}$ . We merely apply  $f_A$  to every wff in  $\Pi_{\Sigma_1}$ .  $\square$

(Like any TI abstraction, see next section)  $f_{gran}$  can map a consistent theory into an inconsistent theory. For example, if the constants  $a$  and  $b$  are "indistinguishable", then a consistent theory with equality and the axiom  $\neg(a = b)$  maps into an inconsistent theory with the axiom  $\neg(\{a, b\} = \{a, b\})$ . However, the following result holds:

**Theorem 6** :  *$f_{gran}$  preserves consistency if indistinguishability is defined over all predicates.*

<sup>8</sup>As in [Hobbs, 1985], we define indistinguishability for unary predicates; it can, however, be easily generalized to  $n$ -ary predicates.

**Proof[Outline]:** By contradiction. Assume that a consistent theory,  $\Sigma_1$  maps onto an inconsistent theory,  $\Sigma_2$ . That is, we can find a proof tree,  $\Pi_{\Sigma_2}$ , of  $\perp$ . We show how you can construct a proof tree,  $\Pi_{\Sigma_1}$  of  $\perp$ , contradicting the assumption that  $\Sigma_1$  is consistent. For every equivalence class,  $\kappa(a)$  we pick one member of that class,  $b$ ; to every wff,  $\varphi$  in  $\Pi_{\Sigma_2}$ , we apply the substitution  $\{\kappa(a)/b\}$ . This will generate a proof tree,  $\Pi_{\Sigma_1}$  whose assumptions will either be axioms of  $\Sigma_1$ , or will be derivable from them using the indistinguishability relation and substitution of equivalences. If indistinguishability is not defined over all predicates, this last fact will not necessarily be true.  $\square$

Note that  $f_{gran}$  is a special case of the example of weak/ordinary abstractions (given in [Plaisted, 1981]) where function symbols are renamed in a systematic (but not necessarily 1-to-1) way.  $\spadesuit$

#### 4 The false proof problem

A major problem with the use of TI-abstractions <sup>9</sup> is that, even if  $\Sigma_1$  is consistent,  $\Sigma_2$  may be inconsistent. An example has already been given for  $f_{gran}$ . With  $f_{AB}$  it is sufficient to consider abstracting the operators " $\alpha_1 \wedge \alpha_2 \rightarrow \alpha_3$ " and " $\alpha_1 \wedge \alpha_4 \rightarrow \alpha_3$ " onto " $\alpha_1 \rightarrow \alpha_3$ " and " $\alpha_1 \rightarrow \alpha_3$ " when  $\alpha_1$  is a theorem but  $\alpha_2$  and  $\alpha_4$  are not both theorems. This problem was noticed by Plaisted [Plaisted, 1981] who called it the "false proof problem". A symptom of the problem comes from the following theorem:

**Lemma 1 :** *If  $f : \Sigma_1 \mapsto \Sigma_2$  is an abstraction and  $\Sigma_2$  is inconsistent then  $f$  is a TI-abstraction.*

Lemma 1 holds independently of the consistency of  $\Sigma_1$ . Example can be easily found with  $\Sigma_1$  consistent or with  $\Sigma_1$  inconsistent. Once  $f$  has been proved to be TI it may happen that  $\Sigma_2$  is inconsistent. This is a major blow to the use of TI-abstractions to guide the proof in the ground space. When  $\Sigma_2$  is inconsistent the structure of the proof in  $\Sigma_2$  could still be used to shape the proof in  $\Sigma_1$ . However, any wff in  $\Sigma_2$  is a theorem and thus  $\Sigma_2$  does not filter out any of the wffs which are not theorems in  $\Sigma_1$ . In a way  $\Sigma_2$  gives too little information. To make matters worse, in general it is not possible to decide in a finite amount of time whether a formal system is consistent.

When working with a fixed formal system (i.e. set theory + first order logic) a solution is to build abstractions which are proved a priori to have consistent  $\Sigma_2$ . In many cases, however, (i.e. planning, logic programming, knowledge based systems), while the set of inference rules of  $\Sigma_1$  is fixed, its axioms may vary and depend on the application. Tenenberg [Tenenberg, 1987, Tenenberg, 1988] proposed some solutions to the problem for predicate abstractions <sup>10</sup> in a resolution-based system. However, the abstractions he proposes have many drawbacks: in one case the abstraction is TI but  $\Omega_{\Sigma_2}$  is not recursive (even if recursively enumerable) and it may take an infinite amount of time to generate it; the

<sup>9</sup>Everything stated in this section holds dually for NTI-abstractions.

<sup>10</sup>Predicate abstractions are abstractions where distinct predicate symbols in  $\Sigma_1$  are mapped onto (possibly not distinct) predicate symbols in  $\Sigma_2$ .

other two types of abstraction are TD or similar to TD [Giunchiglia and Walsh, Forthcoming 89] <sup>11</sup>. This means that completeness is lost since there is at least one theorem in  $\Sigma_1$  whose abstraction is not a theorem in  $\Sigma_2$ . We consider completeness the one property you do not want to lose.

The ideal solution would be to generalize the concept of abstraction mapping to be parameterized on the axioms of  $\Sigma_1$  and then to find sufficient conditions which guarantee that a TI-abstraction maps  $\Sigma_1$  into a consistent  $\Sigma_2$ , independently of the axioms of  $\Sigma_1$  (as long as  $\Sigma_1$  is consistent). This seems a reasonable request since there are abstractions which, fix  $\Lambda_{\Sigma_1}, \Lambda_{\Sigma_2}, \Delta_{\Sigma_1}, \Delta_{\Sigma_2}, f_A, f_{\Omega}$  and  $f_{\Delta}$  are TI for any choice of the theoretic axioms (this is, for instance, true for the abstractions of the three examples) <sup>12</sup>. This request turns out to be unsatisfiable. In fact, even if there are abstractions which are TI- for any choice of the theoretic axioms, for any such TI-abstraction there always exists a set of consistent axioms whose abstraction is inconsistent. The argument goes as follows.

Let  $\Lambda_{\Sigma_1}$  and  $\Lambda_{\Sigma_2}$  be two languages,  $\Delta_{\Sigma_1}$  and  $\Delta_{\Sigma_2}$  two deductive machineries. Then, if  $f_A : \Lambda_{\Sigma_1} \mapsto \Lambda_{\Sigma_2}$ ,  $g : \Lambda_{\Sigma_1} \mapsto \Lambda_{\Sigma_2}$  and  $f_{\Delta} : \Delta_{\Sigma_1} \mapsto \Delta_{\Sigma_2}$  are three total functions,  $F = \langle f_A, g, f_{\Delta} \rangle$  is an abstraction from  $\Sigma_1 = \langle \Lambda_{\Sigma_1}, \Lambda_{\Sigma_1}, \Delta_{\Sigma_1} \rangle$  to  $\Sigma_2 = \langle \Lambda_{\Sigma_2}, \Lambda_{\Sigma_2}, \Delta_{\Sigma_2} \rangle$ . Then for any  $\Omega_{\Sigma_1} \subset \Lambda_{\Sigma_1}$ , if by " $g \uparrow \Omega_{\Sigma_1}$ " we indicate  $g$  restricted to apply to  $\Omega_{\Sigma_1}$ ,  $F^{\Omega_{\Sigma_1}} = \langle f_A, g \uparrow \Omega_{\Sigma_1}, f_{\Delta} \rangle$  is an abstraction from  $\Sigma_1^{\Omega_{\Sigma_1}} = \langle \Lambda_{\Sigma_1}, \Omega_{\Sigma_1}, \Delta_{\Sigma_1} \rangle$  to  $\Sigma_2^{\Omega_{\Sigma_1}} = \langle \Lambda_{\Sigma_2}, \Omega_{\Sigma_2}, \Delta_{\Sigma_2} \rangle$ , with  $\Omega_{\Sigma_2} = g(\Omega_{\Sigma_1})$ .

**Theorem 7 :** *Let  $\Lambda_{\Sigma_1}$  and  $\Lambda_{\Sigma_2}$  be two languages,  $\Delta_{\Sigma_1}$  and  $\Delta_{\Sigma_2}$  two deductive machineries,  $f_A : \Lambda_{\Sigma_1} \mapsto \Lambda_{\Sigma_2}$ ,  $g : \Lambda_{\Sigma_1} \mapsto \Lambda_{\Sigma_2}$  and  $f_{\Delta} : \Delta_{\Sigma_1} \mapsto \Delta_{\Sigma_2}$  three total functions. Then there exists  $\Omega_{\Sigma_1} \subset \Lambda_{\Sigma_1}$  such that, if the abstraction  $F^{\Omega_{\Sigma_1}} = \langle f_A, g \uparrow \Omega_{\Sigma_1}, f_{\Delta} \rangle$  is TI and NTI but not NTC, then  $\Sigma_1^{\Omega_{\Sigma_1}}$  is consistent and  $\Sigma_2^{\Omega_{\Sigma_1}}$  is inconsistent.*

**Proof[Outline]:** By constructing  $\Omega_{\Sigma_1}$ . Because  $F^{\Omega_{\Sigma_1}}$  is NTI- but not NTC-, there exists a wff  $\varphi$  such that adding  $F^{\Omega_{\Sigma_1}}(\varphi)$  as an axiom to  $\Sigma_2$  makes an inconsistent formal system, but that adding  $\varphi$  as an axiom to  $\Sigma_1$  doesn't.  $\square$

Theorem 7 can actually be proved in more powerful forms; however the hypotheses hold for most TI-abstractions. For instance negation preserving abstractions that are TI are also NTI and vice versa (theorem 1). Theorem 7 proves that we cannot find a TI-abstractions which maps a consistent  $\Sigma_1$  into a consistent  $\Sigma_2$  independently of the axioms of  $\Sigma_1$ . One possible solution is to define abstractions which guarantee to build consistent abstract spaces for certain (syntactic characterisations of) classes of theories. However such a solution would not solve the problem in its generality and could not be used in a general purpose, theory independent "abstract theorem prover". A different approach is to be able to vary (in an automated way) the TI-abstraction until we can

<sup>11</sup>Note that it can be proved that, more generally, for any TD-abstraction, if  $\Sigma_1$  is consistent, so is  $\Sigma_2$ .

<sup>12</sup>Of course theory independent TI-abstractions are in general less efficient than the ones geared towards one single theory as they do not exploit the structure of theoretic axioms.

decidably show that  $\Sigma_2$  is consistent. This can be done as follows.

TI-abstractions applied to the same  $\Sigma_1$  can be classified into a weak partial order, indicated by " $\sqsubseteq$ ".

**Definition 6 ( $\sqsubseteq$ )** : If  $f_i : \Sigma_1 \mapsto \Sigma_2^i$  and  $f_j : \Sigma_1 \mapsto \Sigma_2^j$  are two TI-abstractions then  $f_i \sqsubseteq f_j$  iff for all wffs  $\varphi_{\Sigma_1}$ , if  $\vdash_{\Sigma_1} f_i(\varphi_{\Sigma_1})$  then  $\vdash_{\Sigma_2^j} f_j(\varphi_{\Sigma_1})$ . We say that  $f_i$  is weaker than  $f_j$  or, dually, that  $f_j$  is stronger than  $f_i$ .

If  $f_i \sqsubseteq f_j$ , then  $f_j$  is stronger than  $f_i$  in the sense that there are fewer wffs which are theorems in  $\Sigma_2^i$  and not in  $\Sigma_1$  than wffs which are theorems in  $\Sigma_2^j$  and not in  $\Sigma_1$ . " $\sqsubseteq$ " is in general a weak partial order (respecting transitivity, antisymmetry and reflexivity) but not a total order. If, however, we have a set of totally ordered abstractions then the following result holds:

**Theorem 8** : If  $f_1 : \Sigma_1 \mapsto \Sigma_2^1, \dots, f_n : \Sigma_1 \mapsto \Sigma_2^n$  are TI-abstractions and  $f_1 \sqsubseteq \dots \sqsubseteq f_n$  ( $f_1, \dots, f_n$  are totally ordered), then if  $\Sigma_2^n$  is consistent so is  $\Sigma_2^i$  for any  $1 \leq i \leq n$ .

Theorem 8 suggests the following process:

- build sets of abstractions,  $F_i = \{f_1^i, \dots, f_{n_i}^i\}$  where  $f_1^i \sqsubseteq \dots \sqsubseteq f_{n_i}^i$ , and  $f_{n_i}^i(\Sigma_1)$  is decidable (eg. it is propositional).
- find a set,  $F_j$  in which the codomain of the strongest abstraction  $f_{n_j}^j(\Sigma_1)$  is consistent. Note that, since  $f_{n_j}^j(\Sigma_1)$  is decidable, its consistency or inconsistency can be proved in a finite amount of time.
- starting with the strongest abstraction (that is with  $l = n_j$ ), until  $l > 1$  use the proof that the abstracted wff is a theorem in  $f_l^j(\Sigma_1)$  to help construct a proof in  $f_{l-1}^j(\Sigma_1)$ . If, in any of the  $f_l^j(\Sigma_1)$ , the abstracted wff is not a theorem, then the wff cannot be a theorem in  $\Sigma_1$  (since  $f_l^j$  is a TI-abstraction).
- when  $l = 1$ , use the proof in  $\Sigma_2^1$  to help construct the proof in  $\Sigma_1$ . If success, then the goal is a theorem.

Of course there is no guarantee that all the steps in unabstracting back to  $\Sigma_1$  will go through or terminate. The overall performance depends on how the various abstractions in the total order are built and on how the process of unabstracting is performed. For instance, computing the consistency of  $f_{n_j}^j(\Sigma_1)$  can be optimized by building a very simple, "minimal"  $f_{n_j}^j(\Sigma_1)$ . Further time can also be saved, when  $f_{n_j}^j(\Sigma_1)$  is proved inconsistent by introducing (in an automated way) small variations in  $f_{n_j}^j$ , that are tuned to the source of the inconsistency.

## 5 Conclusions

In this paper we have proposed a theory of abstraction which extends the notions of abstraction previously used. We have focused on abstract theorem proving and have suggested that a certain class of provability preserving abstractions<sup>13</sup>, TI-abstractions (which are not TC) are

the correct abstractions to use. TC-abstractions are in general too strong, and the goal of having "simpler" abstract proofs does not seem achievable except in very special and limited forms (for instance, if  $f : \Sigma_1 \mapsto \Sigma_2$  is a TC-abstraction then if  $\Sigma_1$  is undecidable then  $\Sigma_2$  cannot be decidable). The dual class of provability preserving abstractions, TD-abstractions (which are not TC) are of less use as they lose completeness; that is, there is at least one theorem whose abstraction is not a theorem.

Unfortunately, TI-abstractions are subject to the false proof problem; they can map a consistent theory into an inconsistent abstract theory. This problem has been tackled in the last section, where a new general purpose solution has been proposed.

## References

- [Dreben and Goldfarb, 1979] B. Dreben and W.D. Goldfarb. *The Decision problem - Solvable classes of quantificational formulas*. Addison-Wesley Publishing Company Inc., 1979.
- [Giunchiglia and Giunchiglia, 1988] F. Giunchiglia and E. Giunchiglia. Building complex derived inference rules: a decider for the class of prenex universal-existential formulas. In *Proc. 1th ECAI*, 1988. Extended version available as DAI Research Paper 359, Dept. of Artificial Intelligence, Edinburgh.
- [Giunchiglia and Walsh, 1989] F. Giunchiglia and T. Walsh. Theorem Proving with Definitions. In *Proc. AISB 89*, 1989.
- [Giunchiglia and Walsh, Forthcoming 89] F. Giunchiglia and T. Walsh. A Theory of Abstraction. Research paper, Dept. of Artificial Intelligence, University of Edinburgh, Forthcoming 89.
- [Goldfarb, 1971] W.D. Goldfarb. *Jacques Herbrand Logical writings*. D. Reidel Publishing Company, Dordrecht, Holland, 1971. A translation of the 'Ecrit logiques', edited by J.V. Heijnoort.
- [Hobbs, 1985] J.R. Hobbs. Granularity. In *Proc. 9th IJCAI conference*, pages 432-435. International Joint Conference on Artificial Intelligence, 1985.
- [Newell et al., 1963] A. Newell, J.C. Shaw, and H.A. Simon. Empirical explorations of the logic theory machine. In Fiegenbaum and Feldman, editors, *Computers & Thought*, pages 134-152. McGraw-Hill, 1963.
- [Plaisted, 1981] D.A. Plaisted. Theorem proving with abstraction. *Artificial Intelligence*, 16:47-108, 1981.
- [Sacerdoti, 1974] E.D. Sacerdoti. Planning in a hierarchy of abstraction spaces. *Artificial Intelligence*, 5:115-135, 1974.
- [Tenenbergs, 1987] J.D. Tenenbergs. Preserving Consistency across Abstraction Mappings. In *Proc. 10th IJCAI conference*, pages 1011-1014. International Joint Conference on Artificial Intelligence, 1987.
- [Tenenbergs, 1988] J.D. Tenenbergs. *Abstraction in Planning*. PhD thesis, Computer Science Department, University of Rochester, 1988. Also TR 250.

<sup>13</sup>Everything said in this section holds dually for NTI-abstractions.