
A Subsumption Architecture for Theorem Proving? [and Discussion]

Alan Bundy, D. Dennett, M. Sharples, M. Brady and D. Partridge

Phil. Trans. R. Soc. Lond. A 1994 **349**, 71-85

doi: 10.1098/rsta.1994.0114

Email alerting service

Receive free email alerts when new articles cite this article - sign up in the box at the top right-hand corner of the article or click [here](#)

To subscribe to *Phil. Trans. R. Soc. Lond. A* go to:
<http://rsta.royalsocietypublishing.org/subscriptions>

A subsumption architecture for theorem proving?

BY ALAN BUNDY

*Department of Artificial Intelligence, University of Edinburgh,
80 South Bridge, Edinburgh, EH1 1HN, U.K.*

Brooks has criticized traditional approaches to artificial intelligence as too inefficient. In particular, he has singled out techniques involving search as inadequate to achieve the fast reaction times required by robots and other AI products that need to work in the real world. Instead he proposes the *subsumption architecture* as an overall organizing principle. This consists of layers of behavioural modules, each of which is capable of carrying out a complete (usually simple) task. He has employed this architecture to build a series of simple mobile robots, but he claims that it is appropriate for all AI products. Brooks's proposal is usually seen as an example of *nouvelle* AI, in contrast to good old-fashioned AI (GOFAI).

Automatic theorem proving is the archetypal example of GOFAI. The resolution theorem proving technique once served as the engine of AI. Of all areas of AI it seems the most difficult to implement using Brooks's ideas. It would thus serve as a keen test of Brooks's proposal to explore to what extent the task of theorem proving can be achieved by a subsumption architecture.

Tactics are programs for guiding a theorem prover. They were introduced as an efficient alternative to search-based techniques. In this paper I compare recent work on tactic-based theorem proving with Brooks's proposals and show that, surprisingly, there is a similarity between them. It thus seems that the distinction between *nouvelle* AI and GOFAI is not so great as is sometimes claimed. However, this exercise also identifies some criticisms of Brooks's proposal.

1. Intelligence without reason?

Rodney Brooks has championed a new approach to robotics, which he calls behaviour-based (see, for example, Brooks 1991). John Hallam and Chris Malcolm describe this approach in more detail elsewhere in this volume, so we will only summarize it here. Its key features are the following.

Situatedness. The robot must operate in and on the real world.

Embodiment. The robot must have sensors and actuators.

Subsumption architecture. The robot must be organized as a network of modules. Each of these modules is behavioural, i.e. it is capable of performing a whole robot behaviour, e.g. walking, grasping. This is in contrast to a functional module, where each module performs one part of a sequence of tasks, e.g. perceiving, planning, acting. Behavioural modules are organized into layers, but operate independently and in parallel. Overall behaviour emerges from the combination of

the behaviours of the modules. The only interaction allowed is that higher-level modules may inhibit the action of lower-level ones. There is no hierarchical organization, no message passing between modules and no central model of the world.

Some critics have seen this as a promising approach to dealing with low-level robot control, but as inadequate in general. They advocate a hybrid approach in which behaviour-based and traditional robotics are combined (see, for example, Malcolm *et al.* 1989). Brooks repudiates this. For instance, in Brooks (1991, § 6) he says

I think that the new approach can be extended to cover the whole story, both with regards to building intelligent systems and to understanding human intelligence.

In this paper we attempt to test this strong hypothesis. We have deliberately adopted the toughest test we could devise: to see whether Brooks's approach works on a task for which it seems least suited, namely mathematical theorem proving. We will see that some aspects of behavioural-based robotics are inherently unsuited, or irrelevant, to the theorem proving task. However, perhaps surprisingly, the essential idea of the subsumption architecture does seem to echo some current approaches to theorem proving, namely the use of *tactics*.

2. Automatic theorem proving

Work on automatic theorem proving goes back to the earliest days of artificial intelligence. A theorem prover for propositional logic was one of the first AI programs (Newell *et al.* 1957). The roots of the field are in mathematical logic. Most early work was based on a theorem of the logician Herbrand, which implicitly defined a procedure for exhaustively searching for a proof. A sequence of systems implementing and refining this procedure (by Gilmore, Davis & Putnam, Prawitz, etc.) led to the resolution method (Robinson 1965) which then became the basis for most future work. Many of these early papers are reprinted in Siekmann & Wrightson (1983).

Both resolution theorem proving and many of the alternative mechanisms can be seen as manipulating an initial conjecture by the application of rules of inference. Each rule breaks the conjecture into one or more subgoals. Rules can then be applied to each of these subgoals to break them into further subgoals. The process ends when the conjecture is reduced to a collection of trivial subgoals, e.g. axioms of the mathematical theory.

Several different rules may apply to each goal. To be sure of finding a proof it is necessary to try each of the alternative rules at some stage. This situation is illustrated in figure 1. Thus automatic theorem proving is a *search* problem. This is where the problems start. For non-trivial theorems the number of rule combinations which must be explored becomes astronomically large. Storing all the possibilities can easily exhaust the memory capacity of even the largest computers. Exploring all the possibilities can take more time than anybody is willing to wait. This phenomenon is the *combinatorial explosion*.

Most work in automatic theorem proving has been directed to controlling the combinatorial explosion.

(a) Resolution-based systems have been refined so that fewer rules apply to each goal or so that each rule application makes more progress.

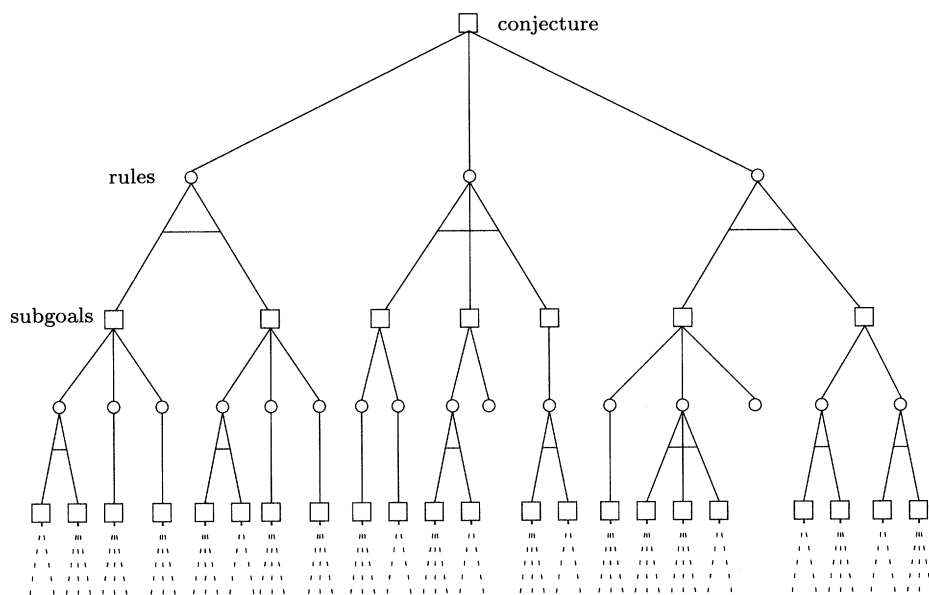


Figure 1. Theorem proving as search. Search is represented as exploration of a tree with the root at the top and the leaves at the bottom. The root is labelled with the original conjecture. Each rule applying to this conjecture is represented as a circle connected to the conjecture node. Each of these rules gives rise to a (possibly empty) set of subgoals. Each subgoal is represented as a square connected to the rule. Further rules apply to each subgoal, creating further subgoals, etc.

(b) Alternatives to resolution, requiring less search, have been invented, e.g. term rewriting (explained below).

(c) Heuristics have been devised to select the most promising rule applications first or to decide that some rule applications should not be tried at all.

(d) More efficient means of storing and applying rules have been devised so that the effect of the combinatorial explosion is reduced.

(e) Interactive systems have been developed in which a human user helps to direct the search.

Proving theorems in most areas of mathematics is a *semi-decidable* problem. This means that: if the conjecture is a theorem then a *complete* theorem prover will eventually find a proof (though this could take a *very* long time); if the conjecture is not a theorem then a complete theorem prover might never terminate in its fruitless search for a proof. It is possible to prove that this is the best that can be done in general. A few simple areas of mathematics are *decidable*. This means that a *decision procedure* can be built for them. Given any conjecture this procedure will eventually terminate and say whether or not the conjecture is a theorem. Unfortunately, many decision procedures are very inefficient, i.e. taking a time that is much worse than exponential in the size of the conjecture. Some work in automatic theorem proving has gone into improving the efficiency of these decision procedures.

Techniques have been discovered which are incomplete (i.e. may not find a proof although one exists), but are efficient and which work in many situations of practical interest. An example is term rewriting. A mathematical theory is expressed as a set of equations, e.g. $x + 0 = x$. These equations are oriented to

form *rewrite rules*, i.e. they are given a direction from left to right or right to left, e.g. $x + 0 \rightarrow x$, where the = sign is turned into the arrow to indicate the direction. If the left-hand side of a rule matches a subexpression of the conjecture then it is replaced by the right-hand side of the rule. To prove a conjecture it is exhaustively rewritten with these rules. It is often possible to show that this process will terminate and that when it does the final result will be trivially either true or false.

Progress has been made since the 1950s. Some automatic provers can generate millions of goals in a feasible timescale and can find the proofs of open problems by exhaustive search (Wos 1993). However, these open problems are typically in new and obscure areas of mathematics, for which human mathematicians have not yet developed strong intuitions. Interactive systems have been used to prove some very hard theorems (Boyer & Yu 1992). However, these are typically in application areas, e.g. program verification, in which the proofs are straightforward, but long and complicated. Despite these advances, the totally automatic proof of hard mathematical theorems still seems a long way off.

3. Criticisms of uniform proof procedures

In the early days of artificial intelligence, automatic theorem proving promised to be a central engine of intelligent systems. Many AI problems (in robot planning, natural language understanding, visual perception, etc.) could be formulated as theorem proving in a theory of common sense knowledge. Theorem proving was a unifying factor in AI.

High expectations ended in disappointment. Combinatorial explosion proved to be an obstacle, even in common sense reasoning. Theorem provers were far too slow. Many critics identified the uniform nature of resolution-like systems as the source of the problem. They proposed the use of inference mechanisms more closely tuned to their applications, i.e. containing features which would limit the amount of search to that strictly required. In expert systems, for instance, production rules were a popular choice. In natural language, semantic nets, frames and *isa* hierarchies were used for representing the meaning of utterances. In vision, constraint propagation was used for object recognition. In practice, these specialized inference mechanisms were often essentially similar to the uniform proof procedures they replaced, but more limited in the kinds of inference they could perform.

Brooks's proposals can be seen as a continuation of this criticism. He characterizes search-based mechanisms as a dead-end in the progress of AI. For instance, in Brooks (1991, §3.2) he says:

In the early days of the formal discipline of Artificial Intelligence, search was adopted as a basic technology. It was easy to program on digital computers. It lead [*sic*] to reasoning systems which are not easy to shoe-horn into situated agents.

So, in the 1970s, AI turned away from theorem proving mechanisms because of the problems of search. (In recent years, though, there has been a reawakening interest in automatic theorem proving within AI as a result of the increasing use of logic for knowledge representation.) But proving theorems remains as a task performed by humans and requiring intelligence. Can AI emulate it, despite the problems caused by the combinatorial explosion? In particular, does Brooks's

proposal of behaviour-based robots suggest an alternative mechanism to uniform proof procedures? Put another way, can we prove theorems without search?

4. Tactic-based theorem proving

Divorced from AI, automatic theorem proving turned to computer programming as its main application area. Computer programming problems can be turned into theorem proving problems in the following way.

1. A computer program can be viewed as a mathematical theory. In logic programming a computer program consists of a set of logical formulæ which are the axioms of a logical theory. In functional programming a computer program is a set of equations which are the axioms of an algebra.

2. A simple automatic theorem prover can be used to run such a computer program. A resolution-based theorem prover can run a logic program. A rewrite-rule-based theorem prover can run a functional program.

3. Automatic theorem provers can also be used to reason about computer programs. A program can be proved to meet a specification of its intended behaviour. A program can be synthesized which meets a specification. One program can be transformed into another more efficient one, meeting the same specification. A program can be proved to terminate.

The specification of a program is usually represented as a logical formula. This formula defines a relation between the input and output of a program, but without describing how the output is to be constructed from the input. For instance, a sorting program might be specified by saying that the output is an ordered permutation of the input. Automatic theorem proving can then be used to give machine assistance to the development of programs, e.g. by helping to synthesize a program to meet the specification.

Researchers in *automated program development* have developed interactive theorem provers to tame the combinatorial explosion. They use automatic theorem proving for the straightforward but long and tedious phases of the proof. Human users direct the proof between these phases. An automated phase might be the application of a decision procedure or the simplification of an expression. A human directed phase might be a split into cases or the use of a key lemma.

A major development came with the introduction of *tactics* to direct the automated phases of the proof (Gordon *et al.* 1979). A tactic is a computer program whose effect is to apply the rules of a mathematical theory to a conjecture. The simplest tactics merely apply one rule. Tactics can be composed using *tacticals*. A tactical might apply two tactics in sequence, apply one tactic repeatedly or use a test to decide which of two tactics should be applied. Compound tactics can be built, for instance, to implement simplification and decision procedures.

5. Examples of tactics

In my group at Edinburgh we have built a number of tactics for different areas of mathematics. Our methodology is to study a number of similar proofs, try to abstract from these to identify the common structure, and then implement this as a tactic. We also try to identify the conditions under which these tactics should be applied. Our theorem provers use these tactics and their preconditions

to search, not at the level of low level proof rules, but at the higher level of tactics. At this higher level the proof steps are bigger and the number of choices is fewer. This helps to defeat the combinatorial explosion and make a practical theorem prover.

Some of our tactics implement the traditional simplification and decision procedures, but others go beyond this to implement larger and more difficult proof steps. For instance, we have built tactics which make case splits, apply induction, generalize conjectures and introduce key lemmas. Some of our tactics can direct a whole proof. This section gives a flavour of what is possible by describing some of these tactics.

(a) *Equation solving: isolation*

Our PRESS system solved equations taken from ‘high school’ examination papers (Sterling *et al.* 1989). Typical equations solved are

$$\begin{aligned}\log_2 x + \log_x 2 &= 5, \\ 3 \tan 3x - \tan x + 2 &= 0 \\ 4^{2x+1} 5^{x-2} &= 6^{1-x}.\end{aligned}$$

In each case the solver is asked to find the values of unknown quantity, x , which will make the equation true.

For PRESS we developed a collection of tactics which together are highly successful in solving equations. For instance, on examples, such as those above, drawn from the GCE A Level (a pre-University examination) papers, PRESS typically solves 80–90%.

One of these tactics is *isolation*. Isolation applies whenever the equation contains only a single occurrence of the unknown. It works by stripping off all functions surrounding this single occurrence until it is isolated on one side of the equation. Consider, for instance, the following simple equation

$$2x^2 - 3 = 5.$$

Since it contains only one occurrence of x , isolation applies and solves it as follows

$$\begin{aligned}2x^2 - 3 &= 5, \\ 2x^2 &= 5 + 3, \\ x^2 &= (5 + 3)/2, \\ x &= \pm\sqrt{[(5 + 3)/2]}, \\ x &= \pm 2.\end{aligned}$$

Each isolation step removes the outermost function surrounding x and applies its inverse to the other side of the equation. First -3 is taken from the left-hand side and replaced by $+3$ on the right; then 2 is replaced by $/2$ and, finally, square is replaced by square root. Arithmetic is then used to simplify the result.

Isolation can be implemented by selective term rewriting. The rewrite rules all have the same essential form

$$f(U) = V \rightarrow U = f^{-1}(V)$$

where f^{-1} is the inverse of f . Note that U must be instantiated to an expression which contains the single occurrence of the unknown x . The isolation rules which

are used in our example are

$$\begin{aligned}U - W = V &\rightarrow U = V + W, \\WU = V &\rightarrow U = V/W, \quad \text{where } W \neq 0, \\U^2 = V &\rightarrow U = \pm\sqrt{W}.\end{aligned}$$

Isolation will not work on an equation containing more than one occurrence of the unknown. It isolates only one occurrence on the left-hand side of the equation. The right-hand side of the equation will then contain the other occurrences, which means it fails the criteria for being a solution.

To solve equations containing more than one occurrence of x requires additional tactics, e.g. *collection* and *attraction*. Collection works by reducing the number of occurrences of the unknown, x , until isolation applies. Attraction works by bringing occurrences of the unknown closer together until collection applies. Thus attraction prepares equations for collection which prepares them for isolation. In PRESS, six major tactics and a few lesser ones interact in this way to produce solutions to a wide variety of equations. All these tactics are implemented as the selective application of rewrite rules.

(b) *Inductive proofs: rippling*

Our oyster-clam system proves inductive theorems of the kind required to reason about recursive computer programs (Bundy *et al.* 1991). Typical theorems proved are

$$\begin{aligned}x + (y + z) &= (x + y) + z, \\ \text{len}(\text{app}(x, y)) &= \text{len}(x) + \text{len}(y), \\ \text{rev}(\text{rev}(x)) &= x,\end{aligned}$$

where $\text{app}(x, y)$ appends the lists x and y , $\text{len}(x)$ is the length of the list x and $\text{rev}(x)$ is the list x in the reverse order. We have developed a collection of tactics which are highly successful in proving such theorems.

Many of the functions which occur in these theorems are defined recursively, i.e. in terms of themselves. For instance, $+$ is defined on the natural numbers (the non-negative integers: $0, 1, 2, 3, \dots$) as follows:

$$\begin{aligned}0 + y &= y, \\ s(x) + y &= s(x + y),\end{aligned}$$

where $s(x) = x + 1$, e.g. $s(s(s(0)))$ represents the number 3. s is the *successor function*. It is a convenient trick for representing all the natural numbers with just one constant, 0, and one function, s .

To prove theorems about recursive functions it is necessary to use mathematical induction. There are many different rules of induction (in fact, infinitely many), but the most common is

$$\frac{P(0), \quad P(n) \rightarrow P(s(n))}{P(n)}$$

i.e. to prove P for any natural number n , first prove it for 0, and second, assuming it for n , prove it for $s(n)$. The first subgoal is the *base case* and the second subgoal is the *step case*. In the step case $P(n)$ is called the *induction hypothesis* and $P(s(n))$ is called the *induction conclusion*.

$$\begin{aligned}
 s(x) + (y + z) &= (s(x) + y) + z \\
 s(x + (y + z)) &= (s(x + y)) + z \\
 s(x + (y + z)) &= s((x + y) + z) \\
 x + (y + z) &= (x + y) + z
 \end{aligned}$$

Figure 2. Rippling in the associativity of $+$. The first equation is the induction conclusion annotated by wave-fronts. Subsequent equations show the effect of rippling these wave-fronts outwards using the rewrite rules (5.1) (three times) and finally the cancellation rule for s (5.2). The wave-fronts ripple outwards until they finally disappear altogether. At this point the induction hypothesis can be used to complete the proof. In general, the wave-fronts do not disappear but are moved to the outside of the induction conclusion leaving an expression inside that matches the induction hypothesis. When rippling is complete the induction hypothesis is used to prove the induction conclusion. The tactic which does this is called fertilization.

Rippling is a tactic we have developed to guide the step case of inductive proofs (Bundy *et al.* 1993). It tries to rewrite the induction conclusion so that the induction hypothesis can be applied to it. In particular, it moves the innermost $s(\dots)$ parts of the induction conclusion outwards to form a copy of the induction hypothesis within the induction conclusion, i.e. it rewrites $P(s(n))$ to $Q[P(n)]$. The expression $s(\dots)$ is an example of what we call a *wave-front*. The word ‘rippling’ comes from the following analogy.

Imagine you are in Scotland standing beside a loch. The surrounding mountains are reflected in the loch. You throw something in the loch. The waves it makes disturb the reflection. The wave-fronts ripple outwards leaving the reflection intact again. The mountains are the induction hypothesis, the reflection is the induction conclusion and the wave-fronts are the expressions introduced into the induction conclusion by the induction rule.

We illustrate rippling in figure 2 with the example of the associativity of $+$. The wave-fronts are shaded to emphasize their movement.

Rippling is implemented as selective rewriting using rewrite rules of the form

$$f(w_1(U)) \Rightarrow w_2(f(U))$$

where the wave-fronts in the rule must match wave-fronts in the induction conclusion. These are called *wave-rules*. Examples of wave-rules are

$$s(U) + V \Rightarrow s(U + V) \quad (5.1)$$

$$s(U) = s(V) \Rightarrow U = V \quad (5.2)$$

Note that in wave-rule (5.2) the right-hand side wave-front, w_2 , is empty.

When wave-rules are applied, not only must the mathematical expressions match, but so too must the wave-fronts. This cuts down the amount of search quite drastically (Bundy *et al.* 1993, §2.5), i.e. there is rarely more than one wave-rule applicable to a sub-expression. Similar remarks hold for other tactics described in this paper. As a result, our tactic-based approach avoids the combinatorial explosion which plagues other approaches to automatic theorem proving.

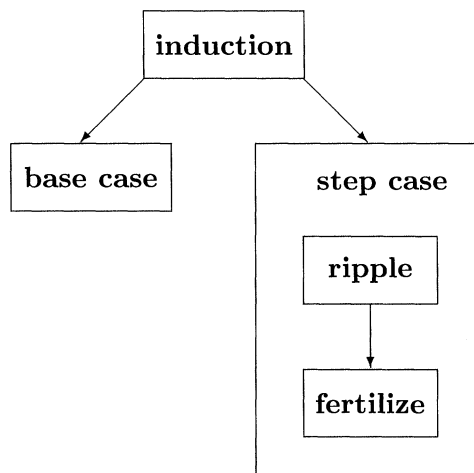


Figure 3. The induction strategy. The induction strategy consists of the application of a rule of induction followed by one or more base and step cases. Only one of each is shown here. The step case consists of rippling until the induction hypothesis appears within the induction conclusion and then fertilizing to use the induction hypothesis in the proof of the induction conclusion.

(c) *An induction strategy*

By combining tactics together it is possible to make a tactic that can direct the complete application of an induction rule, i.e. including the proof of the base and step cases. We have implemented such a tactic in oyster-clam, which we call the *induction strategy*. It is summarized in figure 3. For some simple theorems the induction strategy can direct the whole proof, e.g. the associativity of $+$. Some slightly more complex theorems can be proved by nested applications of the induction strategy.

The induction strategy shows that the word ‘tactic’ was, perhaps, ill-chosen. Theorem proving tactics are not restricted to directing only parts of proofs, as their name might suggest. The induction strategy alone can direct a complete proof. Nor are they restricted to directing only the simplest parts of proofs, e.g. decision procedures and simplification. For instance, rippling is capable of making, so called, ‘eureka’ steps, i.e. those proof steps that seem to call for human intervention.

(d) *Proof critics*

However, in tactic-based systems there are nearly always theorems which cannot be proved solely by the use of available tactics. For instance, the ripple tactic will fail if none of the available wave-rules is applicable to the current subgoal. Then the induction strategy will also fail.

To cope with these failures we have developed a system of *proof critics*. With each tactic is associated a description of what it is trying to achieve and a collection of critics. Each critic tries to capture one of the ways in which the tactic might fail to reach its goal and to suggest a patch to the partial proof. The critic has a precondition, which is a variant of the precondition of its associated tactic, e.g. it might be missing one key part. This precondition is used to determine which critic to use to patch a failed proof.

The critics associated with rippling can suggest any of the following patches, according to the precise way in which rippling fails.

1. Conjecturing a wave-rule which would apply to the current subgoal, proving it and then using it.

2. Attempting a more complex form of induction which would produce wave-fronts at this point in the proof which would enable a wave-rule to apply that currently does not.

3. Generalizing the conjecture in such a way as to allow a wave-rule that currently does not apply to do so.

4. Splitting the proof into cases so that a conditional wave-rule, which was previously inapplicable owing to the failure to prove its condition, will apply.

All of these patches represent proof steps which have previously been regarded as hard to automate. For instance, suggesting appropriate lemmas and generalizing conjectures are classic examples of ‘eureka’ steps.

(e) *Proof plans*

Our oyster–clam system works in two phases. The clam system first tries to find a plan of the proof. This is represented as a custom-built tactic constructed from the general-purpose tactics known to the system and modified by the proof critics. The oyster system then runs this custom-built tactic to produce a formal proof. The preconditions of the tactics are used during the proof planning phase together with predictions of the effects of the tactics. Fairly standard AI plan formation techniques are used to construct the customized tactic from the preconditions and effects of the general-purpose tactics and proof critics.

We find this two phase process useful. For instance, it enables critics to work on a partial proof and to suggest radical rearrangements of it, e.g. going back to the beginning and proving a more general theorem but with a similar proof. However, it is not an essential ingredient of a tactic-based system. Tactics are most commonly used in interactive systems, in which the user decides when and which tactic to call. It is also possible to implement a totally automated one-phase process in which search is conducted at the level of tactics using their preconditions to determine which tactics are applicable. Christian Horn has built a rival version of oyster which uses tactics in this way (Horn 1992). This one phase approach works quite well for organizing tactics but it does not lend itself so well to the use of critics for patching failed proof attempts.

6. Comparison of tactics and subsumption architectures

We now turn to the main question addressed in this paper; to what extent are tactic-based theorem provers similar to the subsumption architecture advocated by Brooks. To facilitate this comparison I have translated my characterization of Brooks’s subsumption architecture from § 1 into tactic-based terminology.

The theorem prover must be organized as a network of tactics. Each of these tactics is behavioural, i.e. it is capable of performing a whole theorem proving behaviour, e.g. decision making, simplification, induction. This is in contrast to a functional module, where each module performs one part of a sequence of tasks, e.g. formula retrieval, matching, search control. Behavioural tactics are organized into layers, but operate independently and in parallel. Overall behaviour emerges from the combination of the behaviours of the tactics.

The only interaction allowed is that higher-level tactics may inhibit the action of lower-level ones. There is no hierarchical organization, no message passing between tactics and no central model of the world.

We start by considering in what ways tactic-based theorem proving fits this description and then discuss in what ways it does not fit.

(a) *Similarities between tactics and subsumption architectures*

A tactic-based system is organized as a network of behavioural modules. Each tactic is a module. A tactic is 'behavioural' in the sense that it performs part of the theorem proving task unaided. For instance, a decision procedure tactic will prove whole theorems if they happen to fall into its decidable class. More usually it will 'finish off' the cases of a larger proof. Other tactics perform the earlier parts of a proof, e.g. setting up an induction. That these tactics only perform part of the overall theorem proving task does not make them less behavioural than Brooks's modules. Some of Brooks's modules only perform part of a task. For instance, a Brooks's walking module might perform the early part of a task by getting the robot to the right place, e.g. so that it can grasp an object.

Although we have organized them hierarchically, tactics could be organized in a non-hierarchical way. A theorem prover could consist of a pool of tactics together with their preconditions. The conjecture to be proved could be stored in some central pool accessible to all tactics. Any tactic whose preconditions are satisfied would fire, manipulate the conjecture, and put any resulting subgoals in the central pool. The other tactics would try to solve these subgoals. Each tactic would be built on the assumption that the other tactics exist and will do the right thing with the subgoals it produced, but it would not need to communicate with them directly. A special module would need to inspect the goal central pool to determine when the theorem was proved and terminate the process. Horn's alternative version of oyster already works in much this way, except for the parallelism, and this omission could be readily rectified.

It is worth noting that this kind of non-hierarchical organization is quite common in GOFAL, e.g. production rules, blackboard architectures. However, it usually introduces lots of search and causes a combinatorial explosion. The preconditions of our tactics preclude too much competition and we rarely have much search.

The central store of conjecture and subgoals can be regarded as analogous to the real world in which the robot moves. Manipulations of these formulæ correspond to actions on the real world. This central store is not the forbidden (by Brooks) 'model' of the real world. These are not simulations of real formulæ, they *are* real formulæ.

The fact that we did not choose this kind of organization in our proof plans work is irrelevant to the argument in this paper. It could be done this way, Horn has shown.

(b) *Differences between tactics and subsumption architectures*

(i) *Real world interaction*

The most obvious difference between Brooks's behaviour-based robotics proposal and any automatic theorem prover is the lack of any real world in theorem proving. The objects of pure mathematics are symbols and expressions. These are inherently abstract. Thus there is nothing for the prover to be situated or

embodied in. This difference seems to be unavoidable, i.e. this part of Brooks's proposal seems unachievable.

However, this disconnectedness with the real world could be interpreted as an implicit criticism of automated mathematical reasoning as being too restricted to proving theorems. Real mathematicians are not concerned only with proving theorems in particular theories. They also formalize informally-stated problems, design mathematical theories, decide which conjectures are interesting to prove, etc. Much of this mathematical activity is informed by real world problems. The classic example is the interaction between physics and mathematics, which led to the invention of the calculus, Fourier series, modern analysis and much, much more. A Brooks-style behaviour-based automatic mathematician might, for instance, use a mathematics module to help solve engineering problems prior to constructing mechanical devices.

This said, it is still true that much of mathematics can and is carried out in a disembodied manner, even by human mathematicians. Such disembodied mathematics remains as a challenge to the Brooks proposal.

(ii) *Hierarchical organization*

In our proof planning approach to theorem proving we chose a very hierarchical organization. This takes two different forms.

1. Large tactics are built from smaller ones, e.g. our induction strategy is built from the base and step case tactics; the step case is built from rippling and fertilization.

2. Our clam program is a centralized planner. It constructs a customized tactic for the current conjecture by combining general-purpose tactics modified by proof critics.

Let us call these type 1 and type 2 hierarchy, respectively. I argued in §6(a) that type 2 hierarchy is an inessential part of tactic-based theorem proving. A flat organization of tactics is possible and has been implemented (see, for example, Horn 1992). We have chosen not to do this because it does not facilitate the use of proof critics to patch failed proof attempts. Critics must investigate the failure of a tactic, which requires message passing between tactic and critic. Critics have proven a powerful technique in avoiding the inherent limitations of the tactic-based approach and for suggesting key 'eureka' steps in proofs.

An alternative, more Brooksonian, implementation of critics is possible.† Each critic would stand alone, monitoring the state of the proof. However, such an implementation would be much clumsier. Currently each critic is cued by the failure of a particular method. Without message passing between a failing method and its critics, it would be quite complicated for a critic to work out that its time had come. Thus, type 2 hierarchy can be avoided, but at a heavy price; at best we would have to settle for an inferior implementation of critics.

Type 1 hierarchy is harder to avoid. Building compound tactics from smaller ones is an essential part of the philosophy of tactic-based theorem proving. Fortunately, this kind of hierarchy seems quite compatible with Brooks's subsumption architecture. He imposes no constraint on how each behavioural module is built. Hierarchical module *construction* seems to be merely good programming prac-

† I thank Geraint Wiggins for suggesting this.

tice. Brooks only rules out a hierarchical *organization* of the resulting module network, any message passing between them or any central model of the world. Indeed Brooks's proposal for 'layering' seems to come close to type 1 hierarchy[‡] (Brooks 1986, §II.B). The main departure from layering is that we have not yet discovered a need to have 'higher level' tactics inhibit the behaviour of 'lower level' ones; unless you count the termination of the theorem proving process by the module that recognizes that the proof is complete.

Deviation from Brooks's proposal might arise if we want a tactic both to be part of some larger tactic and also to be a member of the tactic network in its own right. From our own experience, examples of this seem to be quite rare, but they do occasionally arise, e.g. rippling is both part of our induction strategy and useful stand-alone.

7. Conclusion

In this paper we have tried to test Brooks's hypothesis that his behaviour-based robotics proposal can be extended to cover the 'whole story' of 'building intelligent systems'. We have subjected it to the toughest test we could devise, to see if it could cover an aspect of intelligent behaviour for which it seemed very unpromising: automatic theorem proving. This test has been surprisingly successful. The tactic-based approach to theorem proving is remarkably similar to Brooks's subsumption architecture. Tactics are very similar to behaviour-based modules. A non-hierarchical organization of these tactics is possible, although it remains to be seen whether this is the most effective organization. There are reasons for suspecting that a hybrid tactic-based and central planning model may be more successful. For instance, this allows the efficient use of proof critics. Central reasoning is also required for analysing and modifying the existing collection of tactics, i.e. for learning.

Brooks's requirement that the theorem prover be situated and embodied in the real world seems inherently impossible to achieve. It is not clear what constitutes the real world in the case of mathematics. It is necessary to have some internal representation of the conjecture to be proved and the subgoals generated during a proof attempt. These must be either stored centrally or passed between tactics to enable them to interact. Unfortunately, either a central model of the world or message passing between tactics violate Brooks's proposal. One solution to this is to view the central goal pool as the real world, rather than as a model of it. This solution is clearly controversial.

Situatedness and embodiment might be achieved by building a robot that used mathematics to solve real world problems. However, there is still a challenge to emulate the behaviour of pure mathematicians who appear to do a lot of theorem proving without direct reference to the real world.

It is significant that it is the non-hierarchical organization part of Brooks's proposal that has been most severely tested by our experiment. This has always seemed the most unsatisfactory aspect of his work. His proposal that modules be combined by 'layering' is vague and in practice seems to result in rather *ad hoc* programming rather than a modular and principled combination. Similar

[‡] I thank Fausto Giunchiglia for pointing this out.

remarks apply to the restriction on message passing. For instance, in one case, Brooks avoids the letter of this law by having one module record a message in the external world and another module read it from there.

Thus these parts of the proposal frequently impose unwanted restrictions which the implementer is forced to hack around. It seems timely to examine them and decide whether they are serving a useful purpose or whether Brooks's proposal would be better without them.

I thank Margaret Boden, Rodney Brooks, Fausto Giunchiglia, Andrew Ireland, Sean Matthews and Geraint Wiggins for feedback on an earlier draft of this paper. Ian Green and Gordon Reid helped me fix some tricky L^AT_EX problems. The research reported here was supported by SERC grant GR/H/23610.

References

- Boyer, R. S. & Yu, Y. 1992 Automated correctness proofs of machine code programs for a commercial microprocessor. In *Proc. 11th Conf. on Automated Deduction* (ed. D. Kapur), pp. 416–430, Saratoga Springs, NY, USA. Published as Springer Lecture Notes in Artificial Intelligence, No 607.
- Brooks, R. A. 1986 A robust layered control system for a mobile robot. *J. Robotics Automation*, RA-2(1).
- Brooks, R. A. 1991 Intelligence without reason. In *Proc 12th IJCAI* (ed. J. Mylopoulos & R. Reiter), pp. 569–595.
- Bundy, A., van Harmelen, F., Hesketh, J. & Smaill, A. 1991 Experiments with proof plans for induction. *J. Automated Reasoning* **7**, 303–324.
- Bundy, A., Stevens, A., van Harmelen, F., Ireland, A. & Smaill, A. 1993 Rippling: A heuristic for guiding inductive proofs. *Artificial Intelligence* **62**, 185–253.
- Gordon, M. J., Milner, A. J. & Wadsworth, C. P. 1979 Edinburgh LCF – A mechanised logic of computation, volume 78 of *Lecture Notes in Computer Science*. Springer Verlag.
- Horn, Ch. 1992 Oyster-2: Bringing type theory into practice. *Information Processing* **1**, 49–52.
- Malcolm, C., Smithers, T. & Hallam, J. 1989 An emerging paradigm in robot architecture. In *Intelligent autonomous systems 2* (ed. T. Kanade, F. C. O. Groen & L. O. Hertzberger). Amsterdam.
- Newell, A., Shaw, J. C. & Simon, H. A. 1957 Empirical explorations with the logic theory machine. In *Proc. West. Joint Comp. Conf.*, pp. 218–239. (Reproduced 1963 in *Computers and thought* (ed. Feigenbaum & Feldman), pp.109–133. New York: McGraw-Hill.)
- Robinson, J. A. 1965 A machine oriented logic based on the resolution principle. *J. Assoc. Comput. Mach.* **12**, 33–41.
- Siekman, J. & Wrightson, G. 1983 *Automation of reasoning 1 & 2*. Springer-Verlag.
- Sterling, L., Bundy, A., Byrd, L., O'Keefe, R. & Silver, B. 1989 Solving symbolic equations with PRESS. *J. Symbolic Computation* **7**, 71–84.
- Wos, L. 1993 Automated reasoning answers open questions. *Notices AMS* **5**, 15–26.

Discussion

D. DENNETT (*Tufts University, U.S.A.*). If we compare some classical architectures with behaviour-based robotics, we find similarities; for example, blackboard architectures or SOAR.

A. BUNDY. I agree – and production rules.

M. SHARPLES (*University of Sussex, U.K.*). One difference between Professor *Phil. Trans. R. Soc. Lond. A* (1994)

Bundy's theorem proving system and Brooks's approach is the role of mediating representations between the different modules. Brooks rejects these.

A. BUNDY. Mediating representations are inherent in the theorem proving task. But one can think about my architecture so that this difference disappears. Brooks appeals to the world as the medium for communication between independent behaviour-producing modules. 'The world' in my system is the set of partial proofs used by the independent sub-systems in the theorem prover.

M. BRADY (*Oxford University, U.K.*). In Brooks's architecture, but not yours, each layer may suppress the activity of the layer below. Secondly, don't the complex difficulties faced in theorem proving or robotics mean that, in practice, the problems can't be solved by either extremes (GOFAI or subsumption techniques)? Don't we need hybrid approaches?

A. BUNDY. Agreed on both points.

D. PARTRIDGE (*University of Exeter, U.K.*). Theorem proving is unlike the problems typical of behaviour-based robotics. It's an intellectual task, not one of general intelligence.

A. BUNDY. Agreed. However, Brooks does claim that subsumption architectures can solve high-level cognitive tasks. I've shown that a subsumption architecture can do theorem proving. That's what's interesting.