
APPROXIMATE REASONING USING ANYTIME ALGORITHMS

Shlomo Zilberstein and Stuart Russell*

*Department of Computer Science
University of Massachusetts, Amherst, MA 01003*

** Department of EECS, Computer Science Division
University of California, Berkeley, CA 94720*

ABSTRACT

The complexity of reasoning in intelligent systems makes it undesirable, and sometimes infeasible, to find the optimal action in every situation since the deliberation process itself degrades the performance of the system. The problem is then to construct intelligent systems that react to a situation after performing the “right” amount of thinking. It is by now widely accepted that a successful system must trade off decision quality against the computational requirements of decision-making. Anytime algorithms, introduced by Dean, Horvitz and others in the late 1980’s, were designed to offer such a trade-off. We have extended their work to the construction of complex systems that are composed of anytime algorithms. This paper describes the compilation and monitoring mechanisms that are required to build intelligent systems that can efficiently control their deliberation time. We present theoretical results showing that the compilation and monitoring problems are tractable in a wide range of cases, and provide two applications to illustrate the ideas.

1 INTRODUCTION

A fundamental problem in computer science and artificial intelligence is the construction of systems that can operate robustly in a variety of real-time environments. A real-time environment can be characterized by a *time-dependent* utility function. In almost all cases, the deliberation required to select optimal actions will degrade the system’s overall utility. It is by now well-understood that a successful system must trade off decision quality for deliberation cost [1, 12, 14, 19, 21, 22].

The problem of deliberation cost has been widely discussed in artificial intelligence, economics, engineering and philosophy. In artificial intelligence in particular, researchers have proposed a number of meta-level architectures to control the cost of base-level reasoning [4, 7, 9, 12, 18]. One promising approach is to use *anytime* [3] or *flexible* [10] algorithms, which allow the execution time to be specified, either as a parameter or by an interrupt, and exhibit a time/quality tradeoff defined by a *performance profile*. They provide a simple means by which a system can control its deliberation without significant overhead.

Soon after the introduction of anytime algorithms, it became apparent that their composition presents a vital, non-trivial problem [3]. We show that modular composition of anytime algorithms can be achieved, hence the advantages of anytime algorithms can be extended to the design of complex real-time systems with many components. In standard algorithms, the fixed quality of the output allows for composition to be implemented by a simple call-return mechanism. When algorithms have resource allocation as a degree of freedom, and can be interrupted at any time, the situation becomes more complex. Consider the following simple example: a real-time medical expert system containing a diagnosis component which passes its results to a treatment-planning component. The following issues arise:

1. How can the individual components be designed as anytime algorithms?
2. How can their performance be described as a function of time and the nature of the inputs?
3. How does the output quality of the treatment component depend on the accuracy of the diagnosis it receives?
4. What sort of programming language constructs are needed to specify how the system is built from its components?
5. For any given amount of time, how should that time be allocated to each of the components?
6. What if the condition of the patient suddenly requires intervention while the diagnosis component is still running and no treatment has been considered?
7. How should the execution of the composite system be managed so as to optimize overall utility, particularly when the total execution time is not known in advance?

In several publications, particularly [23], we address these issues in some depth. Here, we sketch our general approach, focusing in particular on the *compilation* problem. Given a system composed of anytime algorithms, compilation determines off-line the optimal allocation of time to the components for any given total allocation. The crucial meta-level knowledge for solving this problem is kept in the *anytime library* in the form of *conditional performance profiles*. These profiles characterize the performance of each elementary anytime algorithm as a function of run-time and input quality.

In Section 2, we define the basic properties of anytime algorithms. We show how to construct anytime algorithms and how to characterize the trade-off that they offer between quality of results and computation time. Section 3 explains the benefits and difficulties involved in the composition of anytime algorithms. Sections 4 and 5 describe the two main components of our solution to the composition problem, namely off-line compilation and run-time monitoring. Section 6 describes briefly two applications of this approach. Finally, in Section 7, we summarize the benefits of our approach and discuss some directions for further work in this field.

2 ANYTIME ALGORITHMS

The term “anytime algorithm” was coined by Dean in the late 1980’s in the context of his work on time-dependent planning. Anytime algorithms are algorithms whose quality of results improves gradually as computation time increases, hence they offer a tradeoff between resource consumption and output quality.

Various metrics can be used to measure the quality of a result produced by an anytime algorithm. From a pragmatic point of view, it may seem useful to define a *single* type of quality measure to be applied to all anytime algorithms. Such a unifying approach may simplify the meta-level control. However, in practice, different types of anytime algorithms tend to approach the exact result in completely different ways. The following metrics have been proved useful in anytime algorithm construction:

1. **Certainty** – this metric reflects the degree of certainty that the result is correct. The degree of certainty can be expressed using probabilities, certainty factors, or any other approach.

2. **Accuracy** – this metric reflects a measure of the difference between the approximate result and the exact answer. Many anytime algorithms can provide a *guarantee* a bound on the error, where the bound is reduced over time.
3. **Specificity** – this metric reflects the level of detail of the result. In this case, the anytime algorithm always produces *correct* results, but the level of *detail* is increased over time.

Many existing programming techniques produce useful anytime algorithms. Examples include iterative deepening search, iterative improvement algorithms in numerical computation, variable precision logic, and randomized techniques such as Monte Carlo algorithms or fingerprinting algorithms. For a survey of such programming techniques and examples of algorithms see [23].

The notion of interrupted computation is almost as old as computation itself. Traditionally, interruption was used primarily for two purposes: aborting the execution of an algorithm whose results are no longer necessary, or suspending the execution of an algorithm for a short time because a computation of higher priority must be performed. Anytime algorithms offer a third type of interruption: interruption of the execution of an algorithm whose results are considered “good enough” by their consumer. Similar ideas motivated the development of the CONCORD system for imprecise computation [15].

2.1 Conditional performance profiles

To allow for efficient meta-level control of anytime algorithms, we characterize their behavior by *conditional performance profiles* (CPP) [24]. A conditional performance profile captures the dependency of output quality on time allocation as well as on input quality. In [23], the reader can find a detailed discussion of various types of conditional performance profiles and their representation. To simplify the discussion of compilation, we will refer only to the *expected* CPP that maps computation time and input quality to the expected output quality.

Definition 1 *The conditional performance profile (CPP), of an algorithm A is a function $CPP_A : Q_{in} \times \mathcal{R}^+ \rightarrow Q_{out}$ that maps input quality and computation time to the expected quality of the results.*

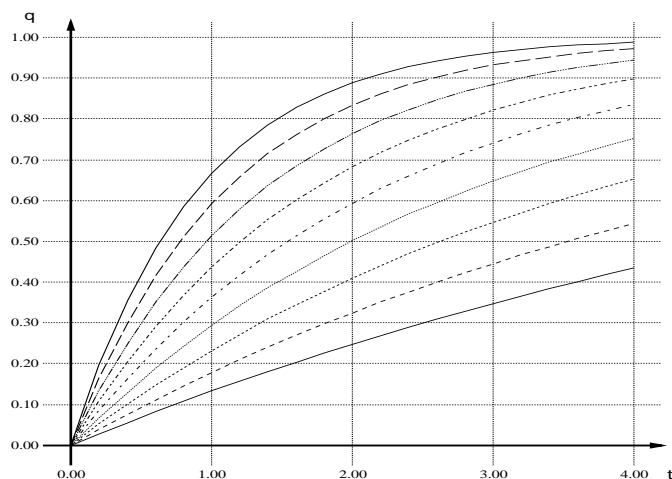


Figure 1 Graphical representation of a CPP

Figure 1 shows a typical CPP. Each curve represents the expected output quality as a function of time for a *given* input quality.

2.2 Interruptible and contract algorithms

In [20] we make an important distinction between two types of anytime algorithms, namely interruptible and contract algorithms. An interruptible algorithm can be interrupted at any time to produce results whose quality is described by its performance profile. A contract algorithm offers a similar trade-off between computation time and quality of results, but it must know the total allocation of time in advance. If interrupted at any point before the termination of the contract time, it may yield no useful results. Interruptible algorithms are in many cases more appropriate for the application, but they are also more complicated to construct. In [20] we show that a simple, general construction can produce an interruptible version for any given contract algorithm. Furthermore, the interruptible algorithm requires at most $4t$ seconds to produce results of the quality achieved by the contract algorithm in t seconds, for all t . This theorem allows us to concentrate on the construction of contract algorithms for complex decision-making tasks and then convert them into interruptible algorithms using a standard transformation.

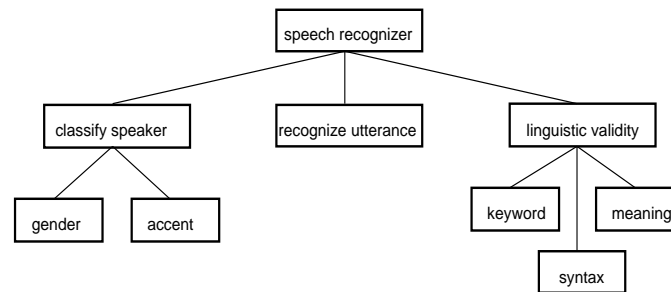


Figure 2 A composite module for speech recognition

3 COMPOSING ANYTIME ALGORITHMS

Modularity is widely recognized as an important issue in system design and implementation. However, the use of anytime algorithms as the components of a modular system presents a special type of scheduling problem. The question is how much time to allocate to each component in order to maximize the output quality of the complete system. We refer to this problem as the anytime algorithm *composition problem*.

Consider for example a speech recognition system whose structure is shown in Figure 2. Each box represents an elementary anytime algorithm whose conditional performance profile is given. The system is composed of three main components. First, the speaker is classified in terms of gender and accent. Then a recognition algorithm suggests several possible matching utterances. And finally, the linguistic validity of each possible utterance is determined and the best interpretation is selected. The composition problem is the problem of calculating how much time to allocate to each elementary component of the composite system, so as to maximize the quality of the utterance recognition.

Solving the composition problem is important for several reasons. First, it introduces a new kind of modularity into real-time system development by allowing for separation between the development of the performance components and the optimization of their performance. In traditional design of real-time systems, the performance components must meet certain time constraints that are not always known at design time. The result is a hand-tuning process that may, or may not, culminate with a working system. Anytime computation offers an alternative to this approach. By developing performance components that are responsive to a wide range of time allocations, one

avoids the commitment to a particular performance level that might fail the system.

The second reason why the composition problem is important relates to the difficulty of programming with anytime algorithms. To make a composite system optimal (or even executable), one must control the activation and interruption of the components. In solving the composition problem, our goal is to minimize the responsibility of the programmer regarding this optimization problem. Our solution is described in the following two sections.

4 COMPILATION

Given a system composed of anytime algorithms, the compilation process is designed to: (a) determine the optimal performance profile of the complete system; and (b) insert into the composite module the necessary code to achieve that performance. The precise definition and solution of the problem depend on the following factors:

1. **Composite program structure** – what type of programming operators are used to compose anytime algorithms?
2. **Type of performance profiles** – what kind of performance profiles are used to characterize elementary anytime algorithms?
3. **Type of anytime algorithms** – what type of elementary anytime algorithms are used as input? what type of anytime algorithm should the resulting system be?
4. **Type of monitoring** – what type of run-time monitoring is used to activate and interrupt the execution of the elementary components?
5. **Quality of intermediate results** – what access does the monitoring component have to intermediate results? is the actual quality of an intermediate result known to the monitor?

Depending on these factors, different types of compilation and monitoring strategies are needed. To simplify the discussion in this paper, we will consider only the problem of producing contract algorithms when the conditional performance profiles of the components are given. We will assume that no

active monitoring is allowed once the system is activated. A broader, in-depth analysis of compilation and monitoring can be found in [23].

Let \mathcal{F} be a set of anytime functions. Assume that all function parameters are passed by value and that functions have no side-effects (as in pure functional programming). Let \mathcal{I} be a set of input variables. Then, the notion of a composite expression is defined as follows:

Definition 2 *A composite expression over \mathcal{F} with input \mathcal{I} is:*

1. *An expression $f(i_1, \dots, i_n)$ where $f \in \mathcal{F}$ is a function of n arguments and $i_1, \dots, i_n \in \mathcal{I}$.*
2. *An expression $f(g_1, \dots, g_n)$ where $f \in \mathcal{F}$ is a function of n arguments and each g_i is a composite expression or an input variable.*

For example, the expression $A(B(x), C(D(y)))$ is a composite expression over $\{A, B, C, D\}$ with input $\{x, y\}$. Suppose that each function in \mathcal{F} has a conditional performance profile associated with it that specifies the quality of its output as a function of time allocation to that function and the qualities of its inputs. Given a composite expression of size n , the main part of the compilation process is to determine a mapping:

$$\mathcal{T} : t \rightarrow (t_1, \dots, t_n) \quad (1.1)$$

This mapping determines for each total allocation, t , the allocation to the components that maximizes the output quality.

4.1 A compilation example

Let us look first at a simple example of compilation involving only two anytime algorithms. Suppose that one algorithm takes the input and produces an intermediate result. This result is then used as input to another anytime algorithm which, in turn, produces the final result. Many systems can be implemented by a composition of a sequence of two or more algorithms. For example, an automated repair system can be composed of two algorithms: diagnosis and treatment. This can be represented in general by the following expression:

$$Output \leftarrow \mathcal{A}_2(\mathcal{A}_1(Input))$$

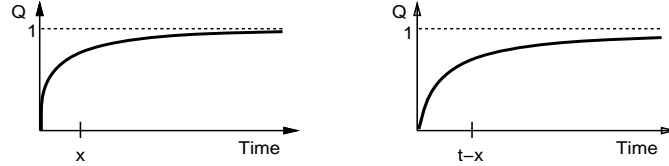


Figure 3 The performance profiles of \mathcal{A}_1 and \mathcal{A}_2

Figure 3 shows the performance profiles of \mathcal{A}_1 and \mathcal{A}_2 . These performance profiles are defined by:

$$Q_1(t) = 1 - e^{-\lambda_1 t} \quad Q_2(t) = 1 - e^{-\lambda_2 t}$$

Assume that the output quality is the sum of the qualities of \mathcal{A}_1 and \mathcal{A}_2 and is expressed by:

$$Q(x) = 1 - e^{-\lambda_1 x} + 1 - e^{-\lambda_2(t-x)} \quad (1.2)$$

The maximal quality is achieved when $\frac{\partial Q}{\partial x} = 0$. In other words:

$$\lambda_1 e^{-\lambda_1 x} - \lambda_2 e^{-\lambda_2(t-x)} = 0 \quad (1.3)$$

The solution of this equation yields the following optimal time allocation mapping:

$$\mathcal{T} : t \rightarrow \left(\frac{\ln \lambda_1 - \ln \lambda_2 + \lambda_2 t}{\lambda_1 + \lambda_2}, \frac{\ln \lambda_2 - \ln \lambda_1 + \lambda_1 t}{\lambda_1 + \lambda_2} \right) \quad (1.4)$$

To complete the compilation process, the compiler needs to insert code in the original expression for proper activation of \mathcal{A}_1 and \mathcal{A}_2 as contract algorithms with the appropriate time allocation. This is done by replacing the simple function call by an anytime function call [23]. The implementation of an anytime function call depends on the particular programming environment and will not be discussed in this paper.

4.2 The complexity of compilation

The compilation problem is defined as an optimization problem, that is, a problem of finding a schedule of a set of components that yields maximal output quality. In order to analyze its complexity, it is more convenient to refer to the decision problem variant of the compilation problem. Given a composite

expression e , the conditional performance profiles of its components, and a total allocation B , the decision problem is whether there exists a schedule of the components that yields output quality greater than or equal to K . To begin, consider the general problem of global compilation of composite expressions, or GCCE. In [23], we prove the following result:

Theorem 3 *The GCCE problem is NP-complete in the strong sense.*

The proof is based on a reduction from the PARTIALLY ORDERED KNAP-SACK problem which is known to be NP-complete in the strong sense. The meaning of this result is that the application of the compilation technique may be limited to small programs. To address the complexity problem of global compilation, we have developed an efficient local compilation technique.

4.3 Local compilation

Local compilation is the process of finding the best performance profile of a module based on the performance profiles of its *immediate* components. If those components are not elementary anytime algorithms, then their performance profiles are determined using local compilation. Local compilation replaces the global optimization problem with a set of simpler, local optimization problems and reduce the complexity of the whole problem. Unfortunately, local compilation cannot be applied to every composite expression. If the expression has repeated subexpressions, then computation time should be allocated only once to evaluate all identical copies. Local compilation cannot handle such cases. However, the following three assumptions make local compilation both efficient and optimal [23]:

1. **The tree-structured assumption** – the input composite expression has no repeated subexpressions, thus its DAG (directed acyclic graph) representation is a tree.
2. **The input-monotonicity assumption** – the output quality of each module increases when the quality of the input improves.
3. **The bounded-degree assumption** – the number of inputs to each module is bounded by a constant, b .

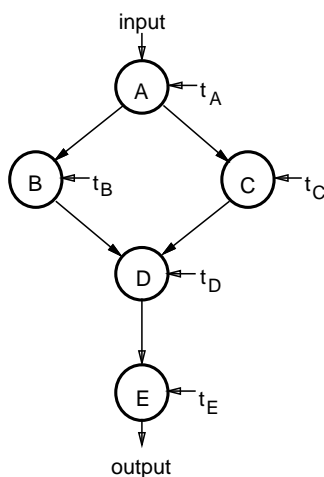


Figure 4 DAG representation of F

Under these assumptions, local compilation is both efficient and yields optimal results [23]. The first assumption is needed so that local compilation can be applied. The second assumption is needed to guarantee the optimality of the resulting performance profile. And the third assumption is needed to guarantee the efficiency of local compilation. Using an efficient tabular representation of performance profiles, we could perform local compilation in constant time and reduce the overall complexity of compilation to be linear in the size of the program.

4.4 Repeated subexpressions

While the input-monotonicity and the bounded-degree assumptions are quite reasonable (and also desirable from a methodological point of view), the tree-structured assumption is somewhat restrictive. We want to be able to handle the case of repeated subexpressions. To understand the problem, consider the following expression:

$$F = E(D(B(A(x)), C(A(x))))$$

Figure 4 shows the DAG representation of F. Recall that the purpose of compilation is to compute a time allocation mapping that specifies for each input quality and total allocation of time the best apportionment of time to

the components so as to maximize the expected quality of the output. But local compilation is only possible when one can repeatedly break a program into sub-programs whose execution intervals are disjoint, so that allocating a certain amount of time to one sub-program does not affect in any way the evaluation and quality of the other sub-programs. This property does not hold for DAGs. In the example shown in Figure 4, B and C are the ancestors of D , but their time allocations cannot be considered independently since they both use the same sub-expression, $A(x)$.

To address this problem we have developed a number of *approximate* compilation techniques that work efficiently on DAGs, but do not guarantee optimality of the schedule [23]. We have also analyzed the compilation of additional programming constructs, such as conditional statements and loops, and derived compilation techniques for those constructs.

5 RUN-TIME MONITORING

Monitoring plays a central role in anytime computation as it complements anytime algorithms with a mechanism that determines their run-time. We have examined the monitoring problem in two types of domains [23]. One type is characterized by the predictability of utility change over time. High predictability of utility allows an efficient use of contract algorithms modified by various strategies for contract adjustment. The second type of domains is characterized by rapid change and a high level of uncertainty. In such domains, active monitoring, that schedules interruptible algorithms based on the value of computation criterion, becomes essential.

Two primary sources of uncertainty affect the operation of real-time intelligent systems. The first source is internal to the system. It is caused by the unpredictable behavior of the system itself. The second source is external. It is caused by unpredictable changes in the environment. These two sources of uncertainty are characterized by two separate knowledge sources. Uncertainty regarding the performance of the system is characterized by the performance profile of the system (in particular, we use *performance distribution profiles* to represent the probability distribution of quality of results). Uncertainty regarding the future state of the environment is characterized by the model of the environment. Obviously, the type of active monitoring may vary as a function of the source of uncertainty and the degree of uncertainty. To demonstrate the operation of active monitoring we consider in this paper only domains that

are characterized by non-deterministic rapid change. Medical diagnosis in an intensive care unit, trading in the stock exchange market, and vehicle control on a highway are examples of such domains. Since accurate projection into the future is very limited in such domains, they require interruptible decision making.

Consider a system whose main decision component is an interruptible anytime algorithm, \mathcal{A} . The conditional probabilistic performance profile of the algorithm is $Q_{\mathcal{A}}(q, t)$ where q is the input quality and t is the time allocation. $Q_{\mathcal{A}}(q, t)$ is a probability distribution and $Q_{\mathcal{A}}(q, t)[q_i]$ denotes the probability of output quality q_i .

Let S be the current state of the domain. Let S_t be the state of the domain at time t . And, let q_t represent the quality of the result of the interruptible anytime algorithm at time t . $U_{\mathcal{A}}(S, t, q)$ represents the utility of a result of quality q in state S at time t . The purpose of the monitor is to maximize the expected utility by interrupting the main decision procedure at the “right” time. Due to the high level of uncertainty in rapidly changing domains, the monitor must constantly assess the value of continued computation by calculating the net expected gain from continued computation given the current best results and the current state of the domain. This is done in the following way:

Due to the uncertainty concerning the quality of the result of the algorithm, the expected utility of the result in a given future state S_t at some future time t is represented by:

$$U'_{\mathcal{A}}(S_t, t) = \sum_i Q_{\mathcal{A}}(q, t)[q_i] U_{\mathcal{A}}(S_t, t, q_i) \quad (1.5)$$

The probability distribution of future output quality is provided by the performance profile of the algorithm. Due to the uncertainty concerning the future state of the domain, the expected utility of the results at some future time t is represented by:

$$U''_{\mathcal{A}}(t) = \sum_S p(S_t = S) U'_{\mathcal{A}}(S, t) \quad (1.6)$$

The probability distribution of the future state of the domain is provided by the model of the environment.

Finally, the condition for continuing the computation at time t for an additional Δt time units is therefore $VOC > 0$ where:

$$VOC = U''_{\mathcal{A}}(t + \Delta t) - U''_{\mathcal{A}}(t) \quad (1.7)$$

Monitoring of interruptible systems can be simplified when it is possible to separate the value of the results from the time used to generate them. In such cases, one can express the comprehensive utility function, $U_{\mathcal{A}}(S, t, q)$, as the difference between two functions:

$$U_{\mathcal{A}}(S_t, t, q) = V_{\mathcal{A}}(S, q) - Cost([t_c, t]) \quad (1.8)$$

where $V_{\mathcal{A}}(S, q)$ is the intrinsic utility function, S is the current state, t_c is the current time, and $Cost([t_c, t])$ is the cost of the time interval $[t_c, t]$. Under this separability assumption, the intrinsic value of allocating a certain amount of time t to the interruptible system (resulting in domain state S) is:

$$V'_{\mathcal{A}}(S, t) = \sum_i Q_{\mathcal{A}}(q, t)[q_i]V_{\mathcal{A}}(S, q_i) \quad (1.9)$$

Hence, the intrinsic value of allocating a certain time t in the current state is:

$$V''_{\mathcal{A}}(t) = \sum_S p(S_t = S)V'_{\mathcal{A}}(S, t) \quad (1.10)$$

And the condition for continuing the computation at time t for an additional Δt time units is again $VOC > 0$ where:

$$VOC = V''_{\mathcal{A}}(t + \Delta t) - V''_{\mathcal{A}}(t) - Cost([t, t + \Delta t]) \quad (1.11)$$

We have shown that monitoring interruptible algorithms using the value of computation criterion is optimal when the intrinsic value function is monotonically increasing and concave down and the time cost function is monotonically increasing and concave up [23]. This assumption is identical to the assumption of Dean and Wellman (See [5], Chapter 8, page 364) that performance profiles have the property of *diminishing returns*.

6 APPLICATIONS

The advantages of compilation and monitoring of anytime algorithms have been demonstrated through a number of applications. In this section we briefly describe two such applications.

6.1 Mobile robot navigation

One of the fundamental problems facing any autonomous mobile robot is the capability to plan its own motion using noisy sensory data. A simu-

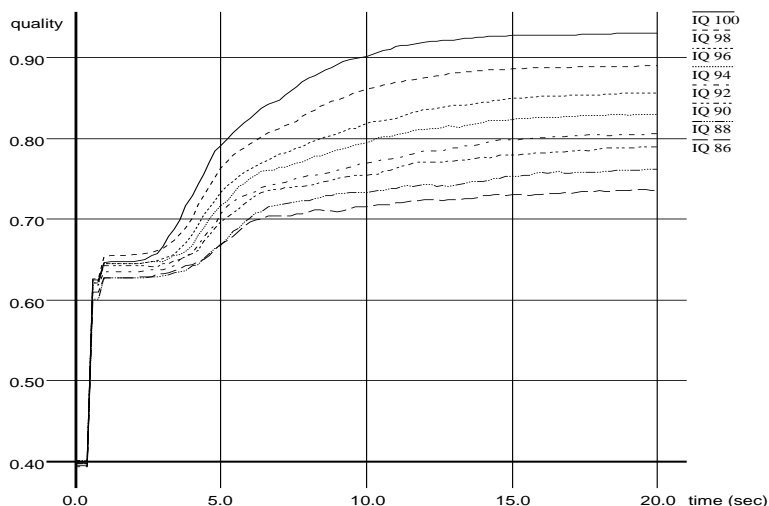


Figure 5 The CPP of the anytime planner

lated robot navigation system has been developed by composing two anytime modules [25]. The first module, a vision algorithm, creates a local domain description whose quality reflects the probability of correctly identifying each basic position as being free space or an obstacle. The second module, a hierarchical planning algorithm, creates a path between the current position and the goal position. The quality of a plan reflects the ratio between the shortest path and the path that the robot generates when guided by the plan.

Anytime hierarchical planning is based on performing coarse-to-fine search that allows the algorithm to find quickly a low quality plan and then repeatedly refine it by replanning a segment of the plan in more detail. Hierarchical planning is complemented by an execution architecture that allows for the execution of abstract plans – regardless of their arbitrary level of detail. This is made possible by using plans as advice that directs the base level execution mechanism but does not fully specify a particular behavior. In practice, uncertainty makes it impossible to use plans except as a guidance mechanism.

The conditional performance profile of the hierarchical planner is shown in Figure 5. Each curve shows the expected plan quality as a function of run-time for a particular quality of the vision module. Finally, an active monitoring scheme was developed to use the compiled performance profile of this system

and the time-dependent utility function of the robot in order to allocate time to vision and planning so as to maximize overall utility.

One interesting observation of this experiment was that the anytime abstract planning algorithm produced high quality results (approx. 10% longer than the optimal path) with time allocation that was less than 30% of the total run-time of a standard search algorithm. This shows that the flexibility of anytime algorithms does not necessarily require a compromise in overall performance.

6.2 Model-based diagnosis

Model-based diagnostic methods identify defective components in a system by a series of tests and probes. Advice on informative probes and tests is given using diagnostic hypotheses that are based on observations and a model of the system. The goal of model-based diagnosis is to locate the defective components using a small number of probes and tests.

The General Diagnostic Engine [6] (GDE) is a basic method for model-based diagnostic reasoning. In GDE, observations and a model of a system are used in order to derive *conflicts* (A conflict is a set of components of which at least one has to be defective). These conflicts are transformed to *diagnoses* (A diagnosis is a set of defective components that might explain the deviating behavior of the system). The process of observing, conflict generation, transformation to diagnoses, and probe advice is repeated until the defective components are identified. GDE has a high computational complexity – $O(2^n)$, where n is the number of components. As a result, its applicability is limited to small-scale applications. To overcome this difficulty, Bakker and Bourseau have developed a model-based diagnostic method, called Pragmatic Diagnostic Engine (PDE), whose computational complexity is $O(n^2)$. PDE is similar to GDE, except for omitting the stage of generating all diagnoses before determining the best measurement-point. Probe advice is given on the basis of the most relevant conflicts, called *obvious* and *semi-obvious* conflicts (An obvious (semi-obvious) conflict is a conflict that is computed using no more than one (two) observed outputs).

In order to construct a real-time diagnostic system, Pos [17] has applied the model of compilation of anytime algorithms to the PDE architecture. PDE can be analyzed as a composition of two anytime modules. In the first module, a subset of all conflicts is determined. Pos implements this module by a contract form of breadth-first search. The second module consists of a repeated loop that

determines which measurement should be taken next, takes that measurement and assimilates the new information into the current set of conflicts. Finally, the resulting diagnoses are reported.

Two versions of the diagnostic system have been implemented: one by constructing a contract algorithm and the other by making the contract system interruptible using our reduction technique. The actual slow down factor of the interruptible system was approximately 2, much better than the worst case theoretical ratio of 4.

7 CONCLUSION

We presented a model for meta-level control of approximate reasoning that is based on compilation and monitoring of anytime algorithms. The technique has several important advantages that include: (1) simplifying the design and implementation of complex intelligent systems by separating the design of the performance components from the optimization of performance; (2) mechanizing the composition process and the monitoring process; and (3) constructing machine independent real-time systems that can automatically adjust resource allocation to yield optimal performance.

The study of anytime computation is a promising and growing field in artificial intelligence and real-time systems. Some of the primary research directions in this field include: (1) Extending the scope of compilation by studying additional programming structures and producing a large library of anytime algorithms; (2) Extending the scope of anytime computation to include the two other aspects of intelligent agents, namely sensing and action; and (3) Developing additional, larger applications that demonstrate the benefits of this approach. The ultimate goal of this research is to construct robust real-time systems in which approximate deliberation, perception and action are governed by a collection of anytime algorithms.

REFERENCES

- [1] M. Boddy and T. L. Dean, Solving time-dependent planning problems, In *Proceedings of the Eleventh International Joint Conference on Artificial Intelligence*, Detroit, Michigan (1989) 979–984.

- [2] M. Boddy, Anytime problem solving using dynamic programming, In *Proceedings of the Ninth National Conference on Artificial Intelligence*, Anaheim, California (1991) 738–743.
- [3] T. L. Dean and M. Boddy, An analysis of time-dependent planning, In *Proceedings of the Seventh National Conference on Artificial Intelligence*, Minneapolis, Minnesota (1988) 49–54.
- [4] T. L. Dean, Intractability and time-dependent planning, In *Proceedings of the 1986 Workshop on Reasoning about Actions and Plans*, M. P. Georgeff and A. L. Lansky, eds., Los Altos, California (Morgan Kaufmann, 1987).
- [5] T. L. Dean and M. P. Wellman. *Planning and Control*. San Mateo, California (Morgan Kaufmann, 1991).
- [6] J. de Kleer and B. C. Williams. Diagnosing multiple faults. *Artificial Intelligence* **32** (1987) 97–130.
- [7] J. Doyle, Rationality and its roles in reasoning, In *Proceedings of the Eighth National Conference on Artificial Intelligence*, Boston, Massachusetts (1990) 1093–1100.
- [8] A. Garvey and V. Lesser, Design-to-time real-time scheduling, In *IEEE Transactions on Systems, Man and Cybernetics*, **23**(6) (1993).
- [9] M. R. Genesereth, An overview of metalevel architectures, In *Proceedings of the Third National Conference on Artificial Intelligence*, Washington, D.C. (1983) 119–123.
- [10] E. J. Horvitz, Reasoning about beliefs and actions under computational resource constraints, In *Proceedings of the 1987 Workshop on Uncertainty in Artificial Intelligence*, Seattle, Washington (1987).
- [11] E. J. Horvitz, H. J. Suermondt and G. F. Cooper, Bounded conditioning: Flexible inference for decision under scarce resources, In *Proceedings of the 1989 Workshop on Uncertainty in Artificial Intelligence*, Windsor, Ontario (1989) 182–193.
- [12] E. J. Horvitz and J. S. Breese, *Ideal partition of resources for meta-reasoning*, Technical Report KSL-90-26, Stanford Knowledge Systems Laboratory, Stanford, California (1990).
- [13] R. E. Korf, Depth-first iterative-deepening: An optimal admissible tree search, *Artificial Intelligence* **27** (1985) 97–109.

- [14] V. Lesser, J. Pavlin and E. Durfee, Approximate processing in real-time problem-solving, *AI Magazine* **9**(1) (1988) 49–61.
- [15] K. J. Lin, S. Natarajan, J. W. S. Liu and T. Krauskopf, Concord: A system of imprecise computations, In *Proceedings of COMPSAC '87*, Tokyo, Japan (1987) 75–81.
- [16] R. S. Michalski and P. H. Winston, Variable precision logic, *Artificial Intelligence* **29**(2) (1986) 121–146.
- [17] A. Pos, *Time-Constrained Model-Based Diagnosis*, Master Thesis, Department of Computer Science, University of Twente, The Netherlands (1993).
- [18] S. J. Russell and E. H. Wefald, Principles of metareasoning, In *Proceedings of the First International Conference on Principles of Knowledge Representation and Reasoning*, R.J. Brachman *et al.*, eds., San Mateo, California (Morgan Kaufmann, 1989).
- [19] S. J. Russell and E. H. Wefald, *Do the Right Thing: Studies in limited rationality*, Cambridge, Massachusetts (MIT Press, 1991).
- [20] S. J. Russell and S. Zilberstein, Composing real-time systems, In *Proceedings of the Twelfth International Joint Conference on Artificial Intelligence*, Sydney, Australia (1991) 212–217.
- [21] H. A. Simon, *Models of bounded rationality, Volume 2*, Cambridge, Massachusetts (MIT Press, 1982).
- [22] S. V. Vrbsky and J. W. S. Liu, Producing monotonically improving approximate answers to database queries, In *Proceedings of the IEEE Workshop on Imprecise and Approximate Computation*, Phoenix, Arizona (1992) 72–76.
- [23] S. Zilberstein, *Operational Rationality through Compilation of Anytime Algorithms*, Ph.D. dissertation, Computer Science Division, University of California, Berkeley, California (1993).
- [24] S. Zilberstein and S. J. Russell, Efficient resource-bounded reasoning in AT-RALPH, In *Proceedings of the First International Conference on AI Planning Systems*, College Park, Maryland (1992) 260–266.
- [25] S. Zilberstein and S. J. Russell, Anytime sensing, planning and action: A practical model for robot control, In *Proceedings of the Thirteenth International Joint Conference on Artificial Intelligence*, Chambéry, France (1993) 1402–1407.