

The Architecture Tradeoff Analysis Method

Rick Kazman, Mark Klein, Mario Barbacci,
Tom Longstaff, Howard Lipson, Jeromy Carriere

Software Engineering Institute
Carnegie Mellon University
Pittsburgh, PA 15213

kazman@sei.cmu.edu

Abstract

This paper presents the Architecture Tradeoff Analysis Method (ATAM), a structured technique for understanding the tradeoffs inherent in the architectures of software intensive systems. This method was developed to provide a principled way to evaluate a software architecture's fitness with respect to multiple competing quality attributes: modifiability, security, performance, availability, and so forth. These attributes interact—improving one often comes at the price of worsening one or more of the others—as is shown in the paper, and the method helps us to reason about architectural decisions that affect quality attribute interactions. The ATAM is a spiral model of *design*: one of postulating candidate architectures followed by analysis and risk mitigation, leading to refined architectures.

Keywords

Software Architecture, Architecture Analysis, Quality Attributes

1. Architecture Tradeoff Analysis

Quality attributes of large software systems are principally determined by the system's software architecture. That is, in large systems, the achievement of qualities such as performance, availability, and modifiability depends more on the overall software architecture than on code-level practices such as language choice, detailed design, algorithms, data structures, testing, and so forth. This is not to say that choice of algorithms or data structures is unimportant, but rather that such choices are less crucial to a system's success than its overall software structure, its architecture. Thus, it is in our interest to try and determine, before it is built, whether a system is destined to satisfy its desired qualities.

Although methods for analyzing specific quality attributes exist (e.g. [4], [5], [8]), these analyses have typically been performed in isolation. In reality, however, the attributes of a system *interact*. Performance impacts modifiability. Availability impacts safety. Security affects performance. Everything affects cost. And so forth. While experienced designers know that these tradeoffs exist, there is no principled method for characterizing them and, in particular, for characterizing the interactions among attributes.

For this reason, software architectures are often designed “in the dark”. Tradeoffs are made—they must be made if the system is to be built—but they are made in an *ad hoc* fash-

ion. Imagine a sound engineer being given a 28 band graphic equalizer, where each of the equalizer's controls has effects that interact with some subset of the other controls. But the engineer is not given a spectrum analyzer, and is asked to set up a sound stage for optimal fidelity. Clearly such a task is untenable. The only difference between this analogy and software architecture is that software systems have far more than 28 independent but interacting variables to be “tuned”.

There are techniques that designers have used to try to mitigate the risks in choosing an architecture to meet a broad palette of quality attributes. The recent activity in cataloguing design patterns and architectural styles is an example of this. A designer will choose one pattern because it is “good for portability” and another because it is “easily modifiable”. But the analysis of patterns doesn't go any deeper than that. A user of these patterns does not know how portable, or modifiable, or robust an architecture is until it has been built.

To address these problems this paper introduces the Architecture Tradeoff Analysis Method (ATAM). ATAM is a method for evaluating architecture-level designs that considers multiple quality attributes such as modifiability, performance, reliability and security in gaining insight as to whether the fully fleshed out incarnation of the architecture will meet its requirements. The method identifies trade-off points between these attributes, facilitates communication between stakeholders (such as user, developer, customer, maintainer) from the perspective of each attribute, clarifies and refines requirements, and provides a framework for an ongoing, concurrent process of system design and analysis.

The ATAM has grown out of work at the Software Engineering Institute on architectural analysis of individual quality attributes: SAAM (Software Architecture Analysis Method) [4] for modifiability, performance analysis [5], availability analysis, and security analysis [6]. SAAM has already been successfully used to analyze architectures from a wide variety of domains: software tools, air traffic control, financial management, telephony, multimedia, embedded vehicle control, and so on.

The ATAM, as with SAAM, has both social and technical aspects. The technical aspects deal with the kinds of data to be collected and how it is analyzed. The social aspects deal with the interactions among the system's stakeholders and area-specific experts, allowing them to communicate using a common framework, to make the implicit assumptions in their analyses explicit, and to provide an objective basis for

negotiating the inevitable architecture tradeoffs. This paper will demonstrate the use of the method, and its benefits in clarifying design issues along multiple attribute dimensions, particularly the tradeoffs in design.

2. Why Use Architecture Tradeoff Analysis?

All design, in any discipline, involves tradeoffs; this is well accepted. What is less well understood is the means for making informed, and possibly even optimal tradeoffs. Design decisions are often made for non-technical reasons: strategic business concerns, meeting the constraints of cost and schedule, using available personnel, and so forth.

Having a structured method helps ensure that the right questions will be asked *early*, during the requirements and design stages when discovered problems can be solved cheaply. It guides users of the method—the stakeholders—to look for conflicts in the requirements and for resolutions to these conflicts in the software architecture.

In realizing the method, we assume that attribute-specific analyses are *interdependent*, and that each quality attribute has connections with other attributes, through specific architectural *elements*. An architectural element is a component, a property of the component, or a property of the relationship between components that affects some quality attribute. For example, the priority of a process is an architectural element that could affect performance. The ATAM helps to identify these dependencies among attributes; what we call *tradeoff points*. This is the principal difference between the ATAM and other software analysis techniques—that it explicitly considers the *connections* between multiple attributes, and permits principled reasoning about the tradeoffs that inevitably result from such connections. Other analysis frameworks, if they consider connections at all, do so only in an informal fashion, or at a high level of abstraction (e.g. [7], [8]). As we will show, tradeoff points arise from architectural elements that are affected by multiple attributes.

3. The ATAM

The ATAM is a spiral model of *design* [3], depicted in Figure 1. The ATAM is like the standard spiral model in that each iteration takes one to a more complete understanding of the system, reduces risk, and perturbs the design. It is unlike the standard spiral in that no implementation need be involved: each iteration is motivated by the results of the analysis and results in new, more elaborated, more informed designs.

Analyzing an architecture involves manipulating, controlling, and measuring several sets of architectural elements, environment factors and architectural constraints. The primary task of an architect is to lay out an architecture that will lead to system behavior which is as close as possible to the requirements within the cost constraints. For example, performance requirements are stated in terms of latency and/or throughput. However, these attributes depend on the architectural elements pertaining to resource allocation: the policy for allocating processes to processors, scheduling concurrent processes on a single processor or managing access to shared data store. The architect must understand the impact of such architectural elements on the ability of the system to meet its requirements and manipulate those elements appropriately.

This task is typically approached with a dearth of tools however. The best architects use their hunches, their experience with other systems, and prototyping to guide them. Occasionally an explicit modeling step is also included as a design activity, or an explicit formal analysis of a single quality attribute is performed.

3.1 The Steps of the Method

The method is divided into four main areas of activities. These are: scenario and requirements gathering, architectural views and scenario realization, model building and analysis, and tradeoffs. The method works, in broad brush, as follows: once a system's initial set of scenarios and requirements have been elicited and an initial architecture (or small set of

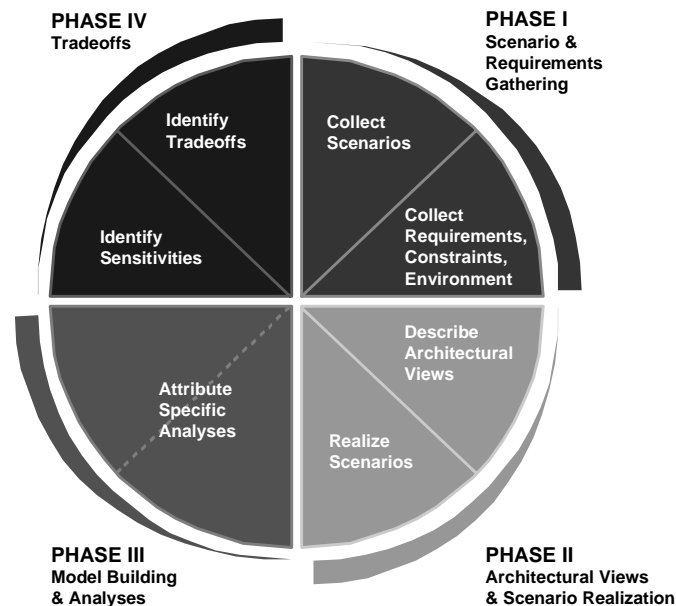


Figure 1: Steps of the Architecture Tradeoff Analysis Method

architectures) is proposed, subject to environment and other considerations, each quality attribute will be evaluated in turn, and *in isolation*, with respect to any proposed design. After these evaluations comes a critique step. During this step tradeoff points are found: elements that affect multiple attributes. After the critique we can either: refine the models and re-evaluate; refine the architectures, change the models to reflect these refinements and re-evaluate; or change some requirements. We now look at each these steps in more detail.

Step 1 — Collect Scenarios

The first step in the method is to elicit system usage scenarios from a representative group of stakeholders. This serves the same purposes as it does in SAAM: to operationalize both functional and quality requirements, to facilitate communication between stakeholders, and to develop a common vision of the important activities the system should support.

Step 2— Collect Requirements/Constraints/Environment

The second step in the method is to identify the *attribute-based* requirements, constraints, and environment of the system. A requirement can have a specific value or can be described via scenarios of hypothetical situations. The environment must be characterized for subsequent analyses (e.g. performance or security) and constraints on the design space, as they evolve, are recorded as these too affect attribute analyses. This step places a strong emphasis on revisiting the scenarios from the previous step to ensure that they account for important quality attributes.

Step 3 — Describe Architectural Views

The requirements, scenarios, and engineering design principles together generate candidate architectures and constrain the space of design possibilities. In addition, design almost never starts from a clean slate: legacy systems, interoperability, and the successes/failures of previous projects all constrain the space of architectures.

Moreover, the candidate architectures are described in terms of the architectural elements that are relevant to each of the important quality attributes. For example, voting schemes are an important element for reliability; concurrency decomposition and process prioritization are important for performance; firewalls and intruder models are important for security, and encapsulation is important for modifiability.

Throughout our presentation of the method, we assume that *multiple, competing* architectures are being compared. However, designers typically consider themselves to be working on only a *single* architecture at a time. Why are these views not aligned? From our perspective, an architecture is a collection of functionality assigned to a set of structural elements, with constraints on the coordination model—the control flow and data flow among those elements. Almost any change will mutate one of these aspects, thus resulting in a new architecture. While this point might seem like a splitting of hairs, these are important hairs to split in this context for the following reason. The ATAM requires building and maintaining attribute models (both quantitative models and qualitative) that reflect and help to reason about the architecture. To change any aspect of an architecture—functionality,

structural elements, coordination model—will change one or more of the models. Once a change has been proposed, the new and old architectures are “competing”, and must be compared. Hence the need for new models that mirror those changes. Using the ATAM, then, is a continual process of choosing among competing architectures, even when these look “pretty much the same” to a casual observer.

Step 4 — Attribute-Specific Analyses

Once a system’s initial set of requirements and scenarios has been elicited and an initial architecture (or small set of architectures) is proposed, each quality attribute must be analyzed *in isolation*, with respect to each architecture. These analyses can be conducted in any order; no individual critique of attributes against requirements or interaction between attributes is done at this point. Allowing separate (concurrent) analysis is an important separation of concerns that allows individual attribute experts to bring their expertise to bare on the system.

The result of the analyses leads to statements about system behavior with respect to values of particular attributes: “requests are responded to in 60 ms. average”, “the mean time to failure is 2.3 days”, “the system is resistant to known attack scripts:”, “the hardware will cost \$80,000 per platform”, “the software will require 4 people per year to maintain”, and so forth.

Step 5 — Identify Sensitivities

Here, the sensitivity of individual attribute analyses to particular architectural elements is determined. That is, one or more attributes of the architecture are varied, the models are then varied to capture these design changes, and the results are evaluated. Any modeled values that are significantly affected by a change to the architecture are considered to be *sensitivity points*.

Step 6 — Identify Tradeoffs

The next step of the method is to critique the models build in step 4 and to find the architectural tradeoff points. Although it is standard practice to critique designs, significant additional leverage can be gained by focussing this critique on the *interaction* of attribute-specific analyses, particularly the location of tradeoff points. Here is how this is done.

Once the architectural sensitivity points have been determined, finding *tradeoff* points is simply the identification of architectural elements to which multiple attributes are sensitive. For example, the performance of a client-server architecture might be highly sensitive to the number of servers (performance increases, within some range, by increasing the number of servers). The availability of that architecture might also vary directly with the number of servers. However, the security of the system might vary inversely with the number of servers (because the system contains more potential points of attack). The number of servers, then, is a tradeoff point with respect to this architecture. It is an element, potentially one of many, where architectural tradeoffs will be made, consciously or unconsciously.

3.2 Iterations of the ATAM

When we have completed the above steps, we are then in a position to compare the results of the analyses to the require-

ments. When the analyses show that the system’s predicted behavior comes adequately close to its requirements, the designers can proceed to a more detailed level of design or to implementation. In practice, however, it is useful to continue to track the architecture with analytic models; to support development, deployment, and beyond to maintenance. Design never ceases in a system’s life cycle, and neither should analysis.

In the event that the analysis reveals a problem, we now develop an action plan for changing the architecture, the models, or the requirements. The action plan will draw on the attribute-specific analyses and identification of tradeoff points. This then leads to another iteration of the method.

It should be made clear that we do not expect these steps to be followed linearly. They can and do interact with each other in complex ways: an analysis can lead to the reconsideration of a requirement; the building of a model can point out places where the architecture has not been adequately though out or documented. This is why we depict the steps as wedges in a circle: at the center of the circle every step touches (and exchanges information with) every other step.

4. An Example Analysis

To exemplify the ATAM, we have chosen an example that has already been extensively analyzed in the research literature, that of a remote temperature sensor (discussed in [8] and elsewhere). We have chosen this example precisely because it has already been heavily scrutinized. The existence of other analyses focuses attention on the differences of the ATAM. We will analyze this system with respect to its availability, security, and performance attributes.

4.1 System Description

The RTS (remote temperature sensor) system exists to measure the temperatures of a set of furnaces, and to report those temperatures to an operator at some other location. In the original example the operator was located at a “host” computer. The RTS sends *periodic temperature updates* to the host computer, and the host computer sends *control requests* to the RTS, to change the frequency at which periodic updates are sent. These requests and updates are done on a furnace by furnace basis. That is, each furnace can be reporting its temperature at a different frequency. The RTS is presumably part of a larger process control system. The control part of the system is not discussed in this example, however.

We are interested in analyzing the RTS for the qualities of performance, security, and availability. To illustrate these analyses we have made the model problem richer and more complex than its original manifestation. In addition to the original set of functional requirements, we have embedded the RTS into a system architecture based on the client-server idiom. The remote temperature sensor functionality is encapsulated in a server, that serves some number of clients. To remain consistent with the original problem, our analysis will assume that there are 16 clients, one per furnace.

The RTS server hardware includes an analog to digital converter (ADC), that can read and convert a temperature for one furnace at a time. Requests for temperature readings are

queued and fed, one at a time, to the ADC. The ADC measures the temperature of each furnace at the frequency specified by its most recently received control request.

4.2 Architectural Options

Since the ATAM was created to illustrate architectural tradeoffs, we need some architectures to analyze. We will consider three options: a simple *Client-Server* architecture, a more complex version of this architecture, called *Client-Server-Server*, where the server has been replicated, and finally an option called *Client-Intelligent Cache-Server*. Each of these architectures will use the identical *server* architecture—all that changes is the ways in which the rest of the system interacts with the server (or servers).

The server’s architecture contains three kinds of components: furnace tasks (independently scheduled units of execution), that schedule themselves to run with some period; a shared communication facility task, that accepts messages from the furnace tasks and sends them to a specified client; and the ADC task, which accepts requests from the furnace tasks, interfaces with the physical furnaces to determine their temperatures, and passes the result back to the requesting furnace task.

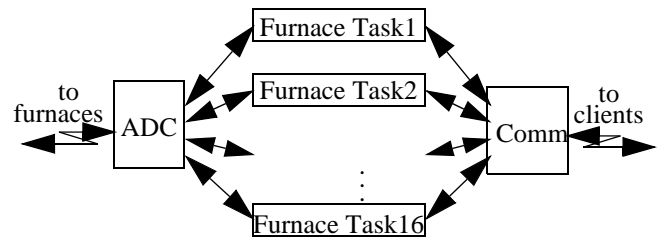


Figure 2: The Architecture of a Furnace Server

Now that the server architecture has been described, we will present the overall system architectures. In each of the systems a set of 16 clients interacts with one or more servers, communicating via a network.

4.3 Architectural Option 1 (Client-Server)

Option 1 is the baseline; a simple and inexpensive client-server architecture, with a single server serving all 16 clients, as shown in Figure 3.

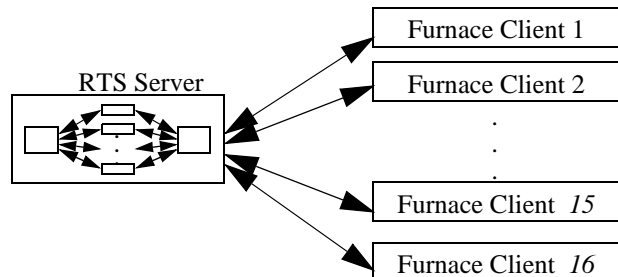


Figure 3: Option 1’s Architecture

4.4 Architectural Option 2 (Client-Server-Server)

Option 2 differs from option 1 in that it adds a second server to the system architecture. These servers interact with clients

as a “primary” server (indicated by the solid lines between servers and clients) or as a “backup” server (indicated by the dashed lines). As shown in Figure 4, each server has its own set of independent furnace tasks, ADC and Comm, but communicates with the same furnaces and with the same set of clients, although under normal operation each server only serves 8 of the 16 clients.

Every client knows the location of both servers and if they detect that the server is down (because it has failed to respond for a prescribed period of time), they will automatically switch to their specified backup.

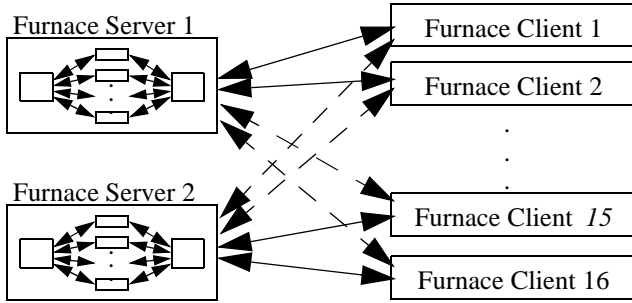


Figure 4: Option 2's Architecture

4.5 Architectural Option 3 (Client-Intelligent Cache-Server)

Option 3 differs from option 1 in only one way: each client has a “wrapper” that intercedes between it and the server. This wrapper is an “intelligent cache”, shown as IC in Figure 5. The cache works as follows: it intercepts periodic temperature updates from the server to the client, builds a history of these updates, and then passes each update to the client. In the event of a service interruption, the cache synthesizes updates for the client. It is an *intelligent* cache because the updates it provides take advantage of historical temperature trends to extrapolate plausible values into the future. This intelligence may be nothing more than linear extrapolation or it might be a sophisticated model that analyzes changes in temperature trends, or takes advantage of domain-specific knowledge on how furnaces heat up and cool down.

As long as the furnaces exhibit regular behavior in terms of temperature trends, then the cache's extrapolated updates will be accurate. Obviously, the cache's synthesized updates will become less meaningful over time.

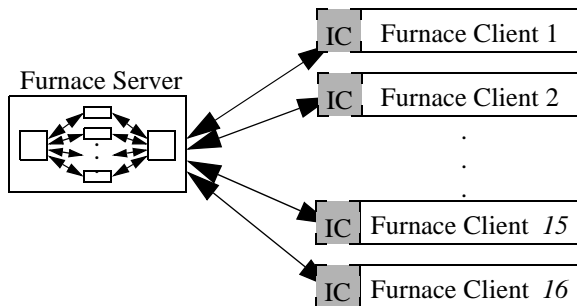


Figure 5: Option 3's Architecture (with Cache)

These then are our three initial architectural alternatives. To understand and compare them, we will analyze them using the ATAM. This method will aid us in understanding not only the relative strengths and weaknesses of each architecture, but will also provide a structured framework for eliciting and clarifying requirements. This is because each analysis technique incorporates (often implicit) assumptions. The use of several analysis techniques together helps to uncover these assumptions and make them explicit.

5. Performance Analyses

In the analyses that follow, we will not show the details of doing performance, availability, security, or any other kind of analysis in detail. We do this for two reasons. First, these details can be found in [2] and the references cited therein. Second, this paper is not intended to propose or exemplify any particular analysis technique. Indeed, any technique that meets the information requirements of the ATAM would do just as well. Our interest is in how the techniques interact, and how this interaction minimizes risk in a rational, documented design process.

In doing a performance analysis, we will consider requirements that typically are derived from scenarios generated through interviews with the stakeholders. In this case the performance requirements are:

PR1: Client must receive a temperature reading *within F seconds* of sending a control request.

PR2: Given that Client X has requested a periodic update every $T(i)$ seconds, it must receive a temperature *on the average every $T(i)$ seconds*.

PR3: The interval between consecutive periodic updates must be *not more than $2T(i)$ seconds*.

In addition to these requirements, we will assume that the behavior patterns and execution environment as follows:

- Relatively infrequent control requests
- Requests are not dropped
- No message priorities
- Server latency = de-queuing time ($Cdq = 10$ ms) + furnace task computation ($Cfnc = 160$ ms)
- Network latency between client and server ($Cnet = 1200$ ms)

Because attributes “trade off” against each other, each assumption is subject to inspection, validation, and questioning as a result of the ATAM.

5.1 Performance Analysis of Option 1

The performance characteristics of architectural option 1 are summarized in Table 1.¹

WCCL	ACPL	Jitter
41,120 ms	5,100 ms	20,400 ms

Table 1: Performance Analysis for Option 1

A worst case control latency of 41.12 seconds sounds like a

bad thing. However, is it? To answer this question we must understand the requirement better. How often will the worst case occur? Is it *ever* tolerable to have the worst case occur? For a safety-critical application, the answer might be “no”. For an interactive Web-based application, the answer might be “yes”, because the price of ensuring a smaller worst case is prohibitive. Doing an analysis of a single quality attribute forces one to consider such requirements issues.

The worst case periodic latency is 37.12 seconds. However, the worst case scenario is unlikely: it assumes that all furnaces are queried at the maximum rate ($T(1) = 10$), that all periodic updates are issued simultaneously, and that the update being measured (the worst case update) is the last one in the queue. More importantly, in this application the cost of a missed update is not great—another one will arrive in the next $T(i)$ seconds. Given these facts, we calculate the *average case* latency, to see if the system can meet its deadlines under more normal conditions, and accept the fact that an occasional periodic update might be missed.

Finally, we turn to PR3, the “Jitter” requirement. Jitter is the amount of time variation that can occur between consecutive periodic temperature updates. The requirement is that the jitter be not more than $2T(i)$, which is a minimum of 20 seconds for $T(i) = 10$. The interval between consecutive readings will be not more than $2T(i)$ if the difference between best case and worst case latency is not more than $2T(i)$, for this is an expression of jitter. So, the worst case jitter = $BCPL - WCPL = 21,760 - 1,360 = 20,400$ ms. This is greater than the minimum $2T(1)$ of 20 seconds, and so option 1 cannot meet PR3.

However, in evaluating architectural option 1’s response to PR3, we must ask “What is the cost of a missed update?”. Is it ever acceptable to violate this requirement? In some safety-critical applications the answer would be “no”. In most applications, the answer would be “yes”, providing that this occurrence was infrequent. The results of this evaluation force one to reconsider the importance of meeting PR3.

5.2 Performance Analysis of Option 2

The performance characteristics of architectural option 2 are summarized in Table 2.

WCCL	ACPL	Jitter
20,560 ms	2,550 ms	9,520 ms

Table 2: Performance Analysis for Option 2

One point should be noted here, and will be returned to later in this discussion: if one of the servers fails, option 2 has the performance and availability characteristics of option 1.

5.3 Performance Analysis of Option 3

The performance characteristics of architectural option 3 are

summarized in Table 3.

WCCL	ACPL	Jitter
41,120 ms	5,200 ms	$\leq 20,400$ ms

Table 3: Performance Analysis for Option 3

For this analysis, we have added a new factor: servicing the intelligent cache (adding a new update and recalculating the extrapolation model) takes 100 ms. In this case, the worst case jitter is exactly the same as for option 1, 20,400 ms. *However*, the intelligent cache exists to protect the client against some amount of lost data. As a consequence, it can bound the worst case jitter. When some pre-set time period elapses, the intelligent cache can pass a synthesized update to the client. When the actual update arrives, the cache updates its state accordingly. Thus, if *we trust the intelligent cache*, we can bound the worst case jitter to any desired value. The smaller the bounding value the more likely a given update will be synthesized by the intelligent cache rather than coming directly from the server.

5.4 Critique of the Analysis

This simple performance analysis gives insight into the characteristics of each solution early in the design process, as befits an *architectural level* analysis. As more details are required, the analyses can be refined, using techniques such as RMA [5], SPE [8], simulation, or prototyping. More importantly, a high-level analysis guides our future investigations, highlighting potential performance “hot-spots”, and allowing us to determine areas of architectural sensitivity to performance, which lead us to the location of tradeoff points.

The ATAM thus promotes analysis at multiple resolutions as a means of minimizing risk at acceptable levels of cost. Areas of high risk are analyzed more deeply (perhaps simulated or prototyped) than the rest of the architecture. And each level of analysis helps determine where to analyze more deeply in the next iteration.

6. Availability Analyses

We will initially only consider a single availability requirement for the RTS system:

AR1: System must not be unavailable for more than 60 minutes per year.

The availability analysis considers a range of component failure rates, from 0 to 24 per year. We only present the results for the case of 24 failures per year. We also consider two classes of repairs, depending on the type of failure:

- major failures, such as a burned-out power supply, that require a visit by a hardware technician to repair, taking 1/2 a day; and
- minor failures, such as software bugs, that can be “repaired” by rebooting the system, taking 10 minutes.

To understand the availability of each of the architectural options, we built and solved a Markov model. In this analysis, we only considered server availability.

1. WCCL = *worst-case control latency*, ACPL = *average-case periodic latency*, and BCPL = *best-case periodic latency*.

6.1 Availability Analysis of Option 1

Solving the Markov model for option 1 gives the results shown in Table 4: 279 hours of down time per year for the burned-out power supply and almost 4 hours down per year for the faulty operating system.

Repair Time	Failures/yr	Availability	Hrs down/yr
12 hours	24.	0.96817	278.833
10 minutes	24.	0.99954	3.9982

Table 4: Availability of Option 1

6.2 Availability Analysis of Option 2

We would expect option 2 to have better availability than option 1, since each server acts as a backup for the other, and we expect the probability of both servers being unavailable to be small. Solving the Markov model for this architecture, we get the results shown in Table 5.

Repair Time	Failures/yr	Availability	Hrs down/yr
12 hours	24.	0.99798	17.7327
10 minutes	24.	~1.0	0.0036496

Table 5: Availability of Option 2

Table 5 shows that option 2 now suffers almost 18 hours of down time per year in the burned-out power supply case. This indicates that architectural option 2 might still suffer outages if it encounters frequent hardware problems. On the other hand, option 2 shows near-perfect availability in the operating system reboot scenario. The availability is shown as perfect 1.0 (the calculations were performed to 5 digits of accuracy). In the worst case of 24 annual failures option 2 exhibits only 13 *seconds* of down time per year.

6.3 Availability Analysis of Option 3

Considering architectural option 3, we expect that it will have better availability characteristics than option 1, but worse than option 2. This is because the intelligent cache, while providing some resistance to server failures, is not expected to be as trustworthy as an independent server. Solving the Markov model, we get the results shown in Table 6 for a cache that is assumed to be trustworthy for 5 minutes.

Repair Time	Failures/yr	Reliability	Hrs down/yr
12 hours	24.	0.96839	276.91
10 minutes	24.	0.9997	2.66545

Table 6: Availability of Option 3

The results in Table 6 show that the 5 minute intelligent cache does little to improve option 3 over option 1 in the scenario with the burnt-out power supply. Option 3 still suffers over 277 hours of down time per year. However, the results for the reboot scenario look more encouraging. The cache reduces down time to 2.7 hours per year. Thus, it appears that the intelligent cache, if its extrapolation was improved, might provide high availability at low cost (since this option uses a single server, compared with the replicated servers

used in option 2). We return to this issue shortly.

7. Critique of the Options

Now that we have seen two different attribute analyses, one part of the method can be commented on: the level of granularity at which a system is analyzed. The ATAM advocates analysis at multiple levels of resolution as a means of minimizing risk at acceptable investments of time and effort. Areas that are deemed to be of high risk are analyzed and evaluated more deeply than the rest of the architecture. And each level of analysis helps to determine “hot spots” to focus on in the next iteration. We will illustrate this point next.

The three architectures can be partially characterized and understood by the measures that we have just derived. From this analysis, we can conclude the following:

- Option 1 has poor performance and availability. It is also the least expensive option (in terms of hardware costs; the detailed cost analyses can be found in [2]).
- Option 2 has excellent availability, but at the cost of extra hardware. It also has excellent performance (when both servers are functioning), and the characteristics of option 1 when a single server is down.
- Option 3 has slightly better availability than option 1, better performance than option 1 (in that the worst case jitter can be bounded), slightly greater cost than Option 1, and lower cost than Option 2.

The conclusions that our analyses lead us to also cause us to ask some further questions.

7.1 Further Investigation of Option 2

For example, we need to consider the nature of option 2 with a server failure. Given that option 2 is identical to option 1 when one server fails, and we have already concluded that option 1 has poor performance and availability, it is important to know how much time option 2 will be in that reduced state of service. When we calculate the availability of *both* servers, using our worst-case assumption of 24 failures per year, we expect to suffer over 22 days of reduced service.

7.2 Action Plan

Given this understanding of options 2 and 3, we see that none of these completely meet their requirements. While option 2 meets its availability target (for failures that involve rebooting the server), it leaves the system in a state where its performance targets can not be met for more than 22 days per year. Perhaps a combination of options 2 and 3—dual servers and intelligent cache on clients—will be a better alternative. This option will provide the superior availability and performance of option 2, but during the times when one server has failed, we mitigate the jitter problems of the single remaining server by using the intelligent cache.

We could not have made these design decisions without knowledge gained from the analysis. Performing a multi-attribute analysis allows one to understand the strengths and weaknesses of a system, and of the parts of a system, within a framework that supports making design decisions.

8. Sensitivity Analyses

Given that the performance and availability of option 2 were so much better than option 1, we would suspect that these attributes are sensitive to the number of servers. Sensitivity analysis confirms this: performance increases linearly as the number of servers increases (up to the point where there is 1 server per client) and availability increases by roughly an order of magnitude with each additional server [2].

Given that option 3 has some desirable characteristics in terms of cost and jitter, we might ask if we can improve the intelligent cache sufficiently to make this option acceptable from an availability perspective. To answer this, we plot option 3’s availability against the length of time during which the intelligent cache’s data is trusted. This plot is shown in Figure 6.

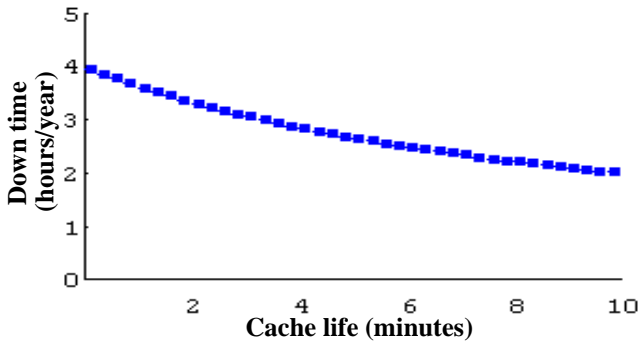


Figure 6: Down time vs. Intelligent Cache Life

As we can see, an improved intelligent cache does improve availability. However, the rate of improvement in availability as a function of cache life is so small that no reasonable, *achievable* amount of cache improvement will result in the kind of availability demonstrated for option 2. In effect, the intelligent cache is an architectural barrier with respect to availability, because it can not be made to achieve the levels of utility required of it. To put it another way, the availability of option 3 is not *sensitive* to cache life. To increase the availability substantially, other paths must be investigated.

9. Security Analyses

Although we could have been conducting security analyses with the performance and availability analyses from the start, the ATAM does not require that all attributes be analyzed in parallel. The ATAM allows the designer to focus on those attributes that are considered to be primary, and then introduce others later on. This can lead to cost benefits in applying the method, since what may be costly analyses for some secondary attributes need not be applied to architectures that were unsuitable for the primary attributes. Though all analyses need not occur “up-front and simultaneously”, the analyses for the secondary attributes can still occur well before implementation begins.

We will now analyze our three options in terms of their security. In particular, we will examine the connections between the furnace servers and clients, since this could be the subject of an attack. The object at risk is the temperature sent

from the server to the client, since this information is used by the client to adjust the furnace settings. If the temperature is tampered with it could be a significant safety concern. Thus we have the security requirement:

SRI: The temperature readings must not be corrupted before they arrive to the client.

Our initial security investigation of the architectural options must, once again, make some environmental assumptions. These assumptions are dependent on the operational environment of the delivered system and include factors such as operator training and patch management. These dependencies are out of scope for the analysis at this level of detail, but must be considered later in the design process.

So, to calculate the probability of a successful attack within an acceptable window of opportunity for an intruder, we define initial values that are *reasonable* for the functions provided in the RTS architectures. These are:

Attack Components		Value
Attack Exposure Window		60 minutes
Attack Rate		0.05 systems/min
Server failure rate		24 failures/year
Prob of server failure within 1 hour		0.0027
Prob of successful	TCP Intercept	0.5
	Spoof IP address	0.9
	Kill Connection	0.75
	Kill Server	0.25

Table 7: Environmental Security Assumptions

In addition, we will posit two attack scenarios: one where the intruder uses a “man in the middle” (MIM) attack, and one where the intruder uses a “spoof server” attack.

For the MIM attack, the attacker uses a TCP intercept tool to modify the values of the temperatures during transmission. Since there are no specific security countermeasures to this attack, the only barrier is the 60 minute window of opportunity and the 0.5 probability of success for the TCP intercept tool. Thus the rate of successful attack is 0.025 systems/minute, or about 1.5 successful attacks expected in the window of opportunity.

For the spoof-server attack, there are three possible ways to succeed. The intruder could wait for a server to fail, then spoof that server’s address and take over the client connections. This presumes that the intruder can determine when a server has failed and can take advantage of this before the clients time out. Another successful method would be to cause the server to fail (the “kill server” attack), then take over the connections. A third is to disrupt the connections between the client and server, then establish new connections as a spoofed server (the “kill connection” attack). For this analysis, it is presumed that the intruder is equally likely to attempt any of these methods in a given attack and the results are summarized in Table 8. Of course, these numbers appear precise, but must be treated as estimates given the

subjective nature of the environmental assumptions.

Attack Type	Expected Intrusions in 60 Mins
Kill Connection	2.04
Kill Server	0.66
Server Failure	0.0072

Table 8: Anticipated Spoof Attack Success Rates

It should be noted that if the system must deal with switching servers and reconnecting clients when a server goes in and out of service, it will be easier for an intruder to spoof a server and convince a client to switch to the bogus server. We will return to this point in the sensitivity analysis.

The results of this analysis show that in each case, it is expected that a penetration will take place within 60 minutes. For the MIM scenario, the expected number of successful attacks is 1.5, indicating that an intruder would have more than enough time to complete the attack before detection. For the spoof attack, the number of successful attacks ranges from 0.0072 to just over 2, again showing that a penetration using this technique is also likely.

9.1 Refined Architectural Options

To address the apparent inadequacy of the three options, we need to cycle around the steps of the ATAM, proposing new architectures. The modified versions of the options include the addition of encryption/decryption and the use of the intelligent cache as an intrusion detection mechanism, as shown in Figure 7.

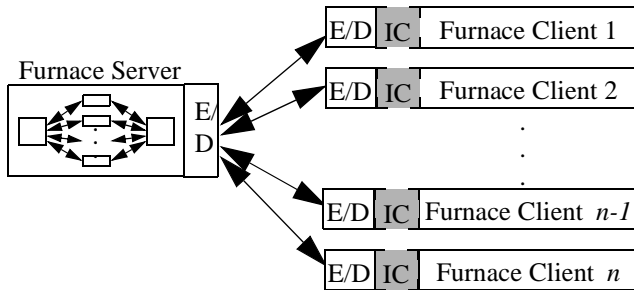


Figure 7: Security Modifications

Encryption/decryption needs little explanation; it is the most common security “bolt on” solution. The other security enhancement is not a topological change, but rather a change in the function of the intelligent cache. In this design, the cache uses its predictive values to determine if the temperatures supplied by the network are reasonable. A temperature that is significantly outside a reasonable change range is deemed to be generated by an intruder and thus the cache aids the operator in detecting an intrusion. Adding encryp-

tion adds some new environmental assumptions. These are:

Attack Components	Value	
Prob of successful	Decrypt	0.0005
	Replay	0.05
	Key Distribution	0.09

Table 9: Additional Security Assumptions

Based on these assumptions, we can calculate the expected number of intrusions. Not surprisingly, the addition of encryption has reduced these substantially—by at least an order of magnitude—in each option:

Attack Type	Expected Intrusions in 60 Mins
Kill Connection	0.18225
Kill Server	0.03375
Server Failure	0.0006

Table 10: Spoof Attack Success Rates with Encryption

Our analysis of the intelligent cache changes only one environmental assumption: the “Attack Exposure Window” goes down to 5 minutes, since we assume that an operator can detect and respond to an intrusion in that time. Using this form of intrusion detection reduces the number of expected intrusions by 1-2 orders of magnitude, giving a result comparable to encryption, but at substantially lower performance and software/hardware costs:

Attack Type	Expected Intrusions in 60 Mins
Kill Connection	0.16875
Kill Server	0.05625
Server Failure	0.005

Table 11: Spoof Attack Success Rates with Intrusion Detection

At this point, new performance and availability analyses will need to be run to account for the additional functionality and hardware required by the intelligent cache or encryption modifications, thus instigating another trip around the spiral.

10. Sensitivities and Tradeoffs

Following the ATAM, we are now in a position to do further sensitivity analysis. In particular, we noted earlier that both availability and performance highly positively correlated to the number of servers. A sensitivity analysis respect to security shows just the opposite: security is negatively related to the number of servers. This is for two reasons:

- going from one to multiple servers requires additional logic within the clients, so that they are able to switch between servers. This provides an entry point for spoofing attacks that does not exist when a client is “hard-wired” to a single server;
- the probability of a server failure within 1 hour increases linearly with the number of servers, thus increasing the

opportunities for server spoofing.

At this point we have discovered an *architectural* tradeoff point, in the number of servers. Performance and availability are positively correlated, while security and presumably cost are negatively correlated, with the number of servers. We cannot maximize cost, performance, availability, and security simultaneously. Using this information, we can make informed tradeoff decisions regarding the level of the various attributes that we can achieve at an acceptable cost, and we can do so within an analytic framework.

11. The Implications of the ATAM

For every assumption that we make in a system's design, we trade cost for knowledge. For example, if a periodic update is supposed to arrive every 10 seconds, do we want it to arrive exactly every 10 seconds, on average every 10 seconds, some time within each 10 second window? To give another example, consider the requirement detailing the worst case latency of control packets. As discussed earlier, is this worst case *ever* acceptable? If so, how frequently can we tolerate it? The process of analyzing architectural attributes forces us to try to answer these questions. Either we understand our requirements precisely or we pay for ignorance by over-engineering or under-engineering the system. If we over-engineering, we pay for our ignorance by making the system needlessly expensive. If we under-engineer, we face system failures, losing customers or perhaps even lives.

Can we believe the numbers that we generated in our analyses? No. However we can believe the *trends*—we have seen differences among designs in terms of orders of magnitude—and these differences, along with sensitivity analysis, tell us where to investigate further, where to get better environmental information, where to prototype, which will get us numbers that we can believe. Every analysis step that we take precipitates new questions. While this seems like a daunting, never-ending prospect, it is manageable because these questions are posed and answered within an analytic attribute framework, and because in architectural analysis we are more interested in finding large effects than in precise estimates.

In addition to concretizing requirements, the ATAM has one other benefit: it helps to uncover implicit requirements. This occurs because attribute analyses are, as we have seen, *interdependent*—they depend, at least partially, on a common set of elements, such as the number of servers. However, in the past, they have been modeled as though they were *independent*. This is clearly not the case.

Each analyzed attribute has implications for other attributes. For example, although the availability analysis was only focussed on servers availability, in a complete analysis we would look at potential failures of all components, including the furnaces, the clients, and the communication lines, and we would look at the various failure types. One such failure is dropping a message. If we assume that the communication channel is not reliable, then we might want to plan for re-sending messages. To do this involves additional computation (to detect and re-send lost messages), storage (to store the messages until they have been successfully transmitted),

and time (for a time-out interval and for message re-transmission). Thus one of the major implications of this availability concern is that the *performance* models of the options under consideration need to be modified.

To recap, we discover attribute interactions in two ways: using sensitivity analysis to find tradeoff points, and by examining the assumptions that we make for analysis *A* while performing analysis *B*. The “no dropped packets” assumption is one example of such an interaction. This assumption, if false, may have implications for safety, security, and availability. A solution to dropping packets will have implications for performance.

In the ATAM attribute experts independently create and analyze their models, then they exchange information (clarifying or creating new requirements). On the basis of this information they refine their models. The *interaction* of attribute-specific analyses, and the identification of tradeoffs, has a greater effect on system understanding and stakeholder communication than any of those analyses could do on their own.

The complexity inherent in most real-world software design implies that an architecture tradeoff analysis will rarely be a straightforward activity that allows you to proceed linearly to a perfect solution. Each step of the method answers some design questions, and brings some issues into sharper focus. However, each step often raises new questions and reveals new interactions between attributes which may require further analysis, sometimes at different levels of abstraction. Such obstacles are an intrinsic part of a detailed methodical exploration of the design space and cannot be avoided. Managing the conflicts and interactions that are revealed by the ATAM places heavy demands on the analysis skills of the individual attribute experts. Success largely depends upon the ability of those experts to transcend barriers of differing terminology and methodology to understand the implications of inter-attribute dependencies, and to jointly devise candidate architectural solutions for further analysis. As burdensome as this may appear to be, it is far better to intensively manage these attribute interactions early in the design process than to wait until some unfortunate consequences of the interactions are revealed in a deployed system.

12. Conclusions

The ATAM was motivated by a desire to make rational choices among competing architectures, based upon well-documented analyses of system attributes at the architectural level, concentrating on the identification of tradeoff points. The ATAM also serves as a vehicle for the early clarification of requirements. As a result of performing an architecture tradeoff analysis, we have an enhanced understanding of, and confidence in, a system's ability to meet its requirements. We also have a documented rationale for the architectural choices made, consisting of both the scenarios used to motivate the attribute-specific analyses and the results of those analyses.

Consider the RTS case study: we began with vague requirements and enumerated three architectural options. The analytical framework helped determine the useful characteristics

of each of the architectural options and highlighted the costs and benefits of the architectural features. More importantly, the ATAM helped determine the locations of architectural tradeoff points, which helped us understand the limits of each option. This helped us develop informed action plans for modifying the architecture, leading to new evaluations and new iterations of the method.

13. References

- [1] M. Barbacci, M. Klein, C. Weinstock, "Principles for Evaluating the Quality Attributes of a Software Architecture", CMU/SEI -96-TR-36, 1996.
- [2] M. Barbacci, J. Carriere, R. Kazman, M. Klein, H. Lipson, T. Longstaff, C. Weinstock, "Steps Toward an Architecture Trade-off Analysis Method: Quality Attribute Models and Analysis", CMU/SEI -97-TR-29, 1997.
- [3] B. Boehm, "A Spiral Model of Software Development and Enhancement", *ACM Software Eng. Notes*, 11(4), 22-42, 1986.
- [4] R. Kazman, G. Abowd, L. Bass, P. Clements, "Scenario-Based Analysis of Software Architecture", *IEEE Software*, Nov. 1996, 47-55.
- [5] M. Klein, T. Ralya, B. Pollak, R. Obenza, M. Gonzales Harbour, *A Practitioner's Handbook for Real-Time Analysis*, Kluwer Academic, 1993.
- [6] H. Lipson, T. Longstaff (eds.), *Proceedings of the 1997 Information Survivability Workshop*, IEEE CS Press, 1997.
- [7] J. McCall, "Quality Factors", in (J. Marciniak, ed.), *Encyclopedia of Software Engineering*, Vol. 2, Wiley: New York, 1994, 958-969.
- [8] C. Smith, L. Williams, "Software Performance Engineering: A Case Study Including Performance Comparison with Design Alternatives", *IEEE Transactions on Software Engineering*, 19(7), 720-741.