



This is a chapter from the book

System Design, Modeling, and Simulation using Ptolemy II

This work is licensed under the Creative Commons Attribution-ShareAlike 3.0 Unported License. To view a copy of this license, visit:

<http://creativecommons.org/licenses/by-sa/3.0/>,

or send a letter to Creative Commons, 444 Castro Street, Suite 900, Mountain View, California, 94041, USA. Permissions beyond the scope of this license may be available at:

<http://ptolemy.org/books/Systems>.

First Edition, Version 1.0

Please cite this book as:

Claudius Ptolemaeus, Editor,
System Design, Modeling, and Simulation using Ptolemy II, Ptolemy.org, 2014.
<http://ptolemy.org/books/Systems>.

Heterogeneous Modeling

Edward A. Lee

Contents

<i>Sidebar: About the Term “Cyber-Physical Systems”</i>	4
1.1 Syntax, Semantics, and Pragmatics	5
1.2 Domains and Models of Computation	7
1.3 The Role of Models in Design	8
<i>Sidebar: Models vs. Realizations of Systems</i>	10
1.4 Actor Models	11
1.5 Model Hierarchy	12
1.6 Approaches to Heterogeneous Modeling	13
<i>Sidebar: About the Term “Actors”</i>	13
<i>Sidebar: Actors in UML, SysML, and MARTE</i>	14
<i>Sidebar: Plurality of Models</i>	16
<i>Sidebar: About Heterogeneous Models</i>	17
<i>Sidebar: Tools Supporting Heterogeneous Models</i>	18
1.7 Models of Time	19
1.7.1 Hierarchical Time	19
1.7.2 Superdense Time	20
1.7.3 Numeric Representation of Time	23
1.8 Overview of Domains and Directors	24
1.9 Case Study	29
1.10 Summary	37

Many of today's engineered systems combine heterogeneous and often complex subsystems. A modern car, for example, may combine a complex engine, electronic control units (ECUs), traction control systems, body electronics (for controlling windows and door locks), entertainment systems, climate control and ventilation, and a variety of safety subsystems (such as airbags). Each subsystem may be realized with a combination of software, electronics, and mechanical parts. Engineering such complex systems is quite challenging, in part because even the smallest subsystems span multiple engineering disciplines.

These complex systems also challenge the design tools that engineers use to specify, design, simulate, and analyze systems. It is no longer sufficient to sketch a mechanical structure and write down a few equations describing the interactions of the mechanical parts. Neither is it sufficient to rely entirely on software tools for 3D modeling of mechanical parts or tools for model-based design of software systems. The complex interplay across domains (mechanics, software, electronics, communication networks, chemistry, fluid dynamics, and human factors) reduce the usefulness of tools that address only a single domain.

The focus of this book is on **cyber-physical systems (CPS)** (Lee, 2008a, 2010a; Lee and Seshia, 2011), which combine computing and networking with physical dynamics. Cyber-physical systems require model combinations that integrate the continuous dynamics of physical processes (often described using differential equations) with models of software. Diverse models are most useful in applications where timed interactions between components are combined with conventional algorithmic computations.¹ They can also be used in traditional software systems that have concurrent² interactions between algorithmic components.

¹An **algorithm** is a finite description of a sequence of steps to be taken to solve a problem. Physical processes are rarely structured as a sequence of steps; rather, they are structured as continuous interactions between concurrent components.

²**Concurrency**, from the Latin verb *concurrere* meaning “run together,” is often taken in computer science to mean the arbitrary interleaving of two or more sequences of steps. However, this is a rather specialized interpretation of a basic concept. In this book, we take concurrency to mean simultaneous operation, with no implication of either interleaving nor sequences of steps. In particular, two continuous processes can operate concurrently without being directly representable as sequences of steps. Consider, for example, a resistive heating element immersed in a vat of water. Increasing the current through the heating element will cause the temperature of the water to rise. The electrical flow is one continuous process, as is the temperature of the water, and these processes are interacting. But neither process is reasonably representable as a sequence of steps, nor is the overall process an interleaving of such steps.

Sidebar: About the Term “Cyber-Physical Systems”

The term “cyber-physical systems” emerged around 2006, when it was coined by Helen Gill at the National Science Foundation in the United States. We are all familiar with the term “**cyberspace**,” attributed William Gibson, who used the term in the novel *Neuromancer* to refer to the medium of computer networks used for communication between humans. We may be tempted to associate the term cyberspace with CPS, but the roots of the term CPS are older and deeper. It would be more accurate to view the terms “cyberspace” and “cyber-physical systems” as stemming from the same root, “**cybernetics**,” rather than viewing one as being derived from the other.

The term “cybernetics” was coined by Norbert Wiener ([Wiener, 1948](#)), an American mathematician who had a huge impact on the development of control systems theory. During World War II, Wiener pioneered technology for the automatic aiming and firing of anti-aircraft guns. Although the mechanisms he used did not involve digital computers, the principles involved are similar to those used today in a huge variety of computer-based [feedback](#) control systems. Wiener derived the term from the Greek κυβερνητης (kybernetes), meaning helmsman, governor, pilot, or rudder. The metaphor is apt for control systems.

Wiener described his vision of cybernetics as the conjunction of control and communication. His notion of control was deeply rooted in closed-loop feedback, where the control logic is driven by measurements of physical processes, and in turn drives the physical processes. Even though Wiener did not use digital computers, the control logic is effectively a computation, and therefore cybernetics is the conjunction of physical processes, computation, and communication.

Wiener could not have anticipated the powerful effects of digital computation and networks. The fact that the term “cyber-physical systems” may be ambiguously interpreted as the conjunction of cyberspace with physical processes, therefore, helps to underscore the enormous impact that CPS will have. CPS leverages a phenomenal information technology that far outstrips even the wildest dreams of Wiener’s era.

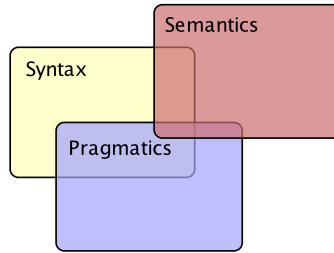


Figure 1.1: Today’s design tools involve complex combinations of syntax, semantics, and pragmatics.

1.1 Syntax, Semantics, and Pragmatics

At the time of this writing, we are in the midst of a dramatic transformation in engineering tools and techniques, which are evolving to enable us to adapt to increasing system complexity and heterogeneity. In the past, entire industries were built around providing design tools for a single engineering domain, such as digital circuits, software, 3D mechanical design, and heating and ventilation systems. Today we see a growing consolidation and combination of design tools; individual tools often expand into tool suites and provide capabilities outside of their traditional domain. This evolution often entails significant growing pains, where poorly integrated capabilities yield frustratingly unexpected behaviors. Tool integration often results in “frankeware,”³ brittle combinations of mostly incompatible tools that are extremely difficult to maintain and use effectively in combination.

Moreover, tools that have traditionally worked well within a relatively narrow domain are not as effective when used in broader domains. Today’s sophisticated design tools involve complex combinations of **syntax** (how a design is represented), **semantics** (what a design means and how it works), and **pragmatics** (Fuhrmann and von Hanxleden, 2008) (how an engineer visualizes, edits, and analyzes a design). When tools are used in domains for which they were not originally designed or in combination with other tools, awkwardness may arise from incompatible syntaxes, poorly understood semantics, and inconsistent human interfaces.

³The term “frankeware” is due to Christopher Brooks.

Incompatibilities in syntax may arise because the structure of designs is intrinsically different (software syntax has very little in common with 3D volumes, for example). But all too often, it arises because the tools were developed in different engineering communities using different techniques. Similarly, the pragmatics of tools, such as how design files are managed and how changes are tracked, are often starkly different by historical accident. Differences in semantics are often accidental as well, sometimes arising from simple misunderstandings. Semantics may not be intuitively obvious in a different domain. A block diagram, for example, may mean something completely different to a control engineer as to a software engineer.

This book examines key concepts in heterogeneous modeling using Ptolemy II, an open-source modeling and simulation tool.⁴ In contrast to most other design tools, Ptolemy II was developed from the outset to address heterogeneous systems. A key goal of the Ptolemy Project (an ongoing research effort at UC Berkeley) has been to minimize the accidental differences in syntax, semantics, and pragmatics between domains, and maximize the interoperability of designs expressed in different domains. As a consequence, Ptolemy II provides a useful laboratory for experimenting with design technologies for cyber-physical systems.

Ptolemy II integrates four distinct classes of syntaxes: block diagrams, bubble-and-arc diagrams, imperative programs, and arithmetic expressions. These syntaxes are complementary, and enable Ptolemy to address a variety of design domains. Block diagrams are used to express concurrent compositions of communicating components; bubble-and-arc diagrams are used to express sequencing of states or modes; imperative programs are used to express algorithms; and arithmetic expressions are used to express functional numeric computations.

Ptolemy II also integrates a number of semantic domains. For block diagrams, in particular, there are many distinct semantics possible. Connections between blocks represent interactions between components in a design, but what type of interaction? Is it an asynchronous message (like sending a letter)? Is it a rendezvous communication (like making a phone call)? Is it a clocked update of data (as in a synchronous digital circuit)? Does time play a role in the interaction? Is the interaction discrete or continuous? To enable heterogeneous modeling, Ptolemy II has been designed to support all of them, and is extensible to support more.

⁴Ptolemy II is available for download at <http://ptolemy.org>.

1.2 Domains and Models of Computation

A **semantic domain** in Ptolemy II, often just called a **domain**, defines the “laws of physics” for the interaction between components in a design. It provides the rules that govern concurrent execution of the components and the communication between components (such as those described above). A collection of such rules is called a **model of computation (MoC)**. We will use the terms “model of computation” and “domain” (nearly) interchangeably, though technically we think of a domain as being an *implementation* of a MoC. The MoC is an abstract model, whereas the domain is its concrete implementation in software.

The rules that constitute a model of computation fall into three categories. The first set of rules specifies what constitutes a component. In this book, a component is generally an **actor**, to be defined more precisely below. The second set of rules specifies the execution and concurrency mechanisms. Are actors invoked in order? Simultaneously? Nondeterministically? The third specifies the communication mechanisms. How do actors exchange data?

Each of the MoCs discussed in this book has many possible variants, many of which have been realized in other modeling tools. In this book, we focus only on MoCs that have been realized in Ptolemy II and that have well understood and documented semantics.⁵ For further context, we also provide brief descriptions and pointers to other useful MoCs that have not been realized in Ptolemy II, but have been realized in other tools.

To support the design of heterogeneous systems, Ptolemy II domains (and models of computation) interoperate with one another. This requires a level of agreement between semantic domains that is rare when tools are developed separately and then later integrated. The principles behind interoperation of domains in Ptolemy II are described in a number of papers (Eker et al., 2003; Lee et al., 2003; Goderis et al., 2009; Lee, 2010b; Tripakis et al., 2013). In this book, we focus on the practical aspects of domain interoperability, not on the theory.

Using a single, coherent software system lets us focus on domain interoperation rather than on less important incompatibilities that typically arise in tool integration. For example, the Ptolemy II type system (which defines the types of data that can be used with various computational components) is shared by all domains, by the **state machine** notation,

⁵In the electronic version of this book, most illustrations of models provide a hyperlink in the caption that enables you to browse the model online. If you are reading the book on a Java-capable machine, then you can edit and execute the models shown in most of the figures.

and by the [expression language](#). The domains are all capable of inferring and verifying appropriate data types; this functionality works seamlessly across heterogeneous models with multiple domains. Similarly, domains that include a notion of time in their semantics share a common representation of time and a (multiform) model of time. Finally, the same graphical editor spans domains, and the same XML schema is used to store design data. These agreements remove many of the practical obstacles to heterogeneous composition of models. They allow us to focus on the benefits of heterogeneous integration – most importantly, the ability to choose the domain that best matches the problem, even when the design is heterogeneous.

1.3 The Role of Models in Design

This book provides a framework for understanding and building models in Ptolemy II and, more broadly, for understanding key issues in modeling and simulating complex heterogeneous systems. This topic is broad enough that no single volume could possibly cover all of the techniques that could be useful to system designers. In this book, we focus on models that describe **dynamics**, or how a system or subsystem evolves in time. We do not cover techniques that focus primarily on the static structure of designs (such as UML class diagrams for software or 3D volumetric modeling). As a consequence, all of the models in this book are executable. We call the execution of a model a **simulation**.

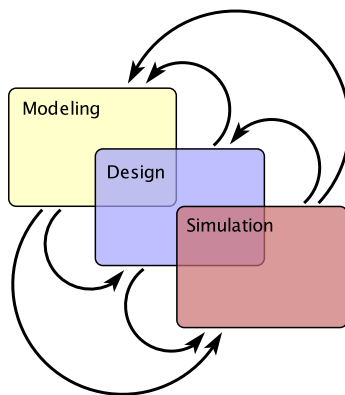


Figure 1.2: Iterative process of modeling, design, and simulation.

Figure 1.2 shows three major parts of the process of implementing systems: modeling, design, and simulation. **Modeling** is the process of gaining a deeper understanding of a system through imitation. Models imitate the system and reflect properties of the system. Models of dynamics specify *what* a system does; that is, how it reacts to stimulus from its environment, and how it evolves over time. **Design** is the structured creation of artifacts (such as software components) to implement specific functionality. It specifies *how* a system will accomplish the desired functionality. **Simulation** shows how models behave in a particular environment. Simulation is a simple form of design analysis; its goal is to lend insight into the properties of the design and to enable **testing** of the design. The models we discuss in this book can also be subjected to much more elaborate forms of analysis, including formal verification. In its most general form, **analysis** is the process of gaining a deeper understanding of a system through dissection, or partitioning into smaller, more readily analyzed pieces. It specifies *why* a system does what it does (or fails to do) what a model says it should do. We leave all analysis techniques except simulation to other texts.

As suggested in Figure 1.2, the three parts of the design process overlap, and the process iterates between them. Normally, the design process begins with modeling, where the goal is to understand the problem and to develop solution strategies.

Modeling plays a central role in modern design processes. The key principle of **model-based design** is to maximally leverage modeling to construct better designs. To be effective, models must be reasonably faithful, of course, but they also must be understandable and analyzable. To be understandable and analyzable, a model needs to have a clear meaning (a clear semantics).

Models are expressed in some **modeling language**. For example, a procedure may be expressed in Java or in C, so these programming languages are in fact modeling languages for procedures. A modeling language has a **strong semantics** if models expressed in the language have a clear and unambiguous meaning. Java, for example, has a stronger semantics than C, as illustrated by the following example.

Example 1.1: Suppose that the arguments to a procedure are of type *int*. In Java, this data type is well defined, but not in C. In C, *int* may represent a 16-bit integer or a 32-bit integer, for example. The behavior of the procedure may be quite different depending on which implementation is provided. Particularly, overflow occurs more easily with 16-bit integers than with 32-bit integers.

Sidebar: Models vs. Realizations of Systems

Models must be used with caution. The **Kopetz principle** (named after Prof. Dr. Hermann Kopetz of TU-Vienna, who taught us this principle), paraphrased, is: *Many properties that we assert about systems (determinism, timeliness, reliability) are in fact not properties of an implemented system, but rather properties of a model of the system.*

Golomb (1971) emphasizes understanding the distinction between a model and thing being modeled, famously stating “you will never strike oil by drilling through the map!” In no way, however, does this diminish the value of a map! Consider **determinism**. A model is **determinate** if it produces a uniquely defined output for each particular input. It is **nondeterminate** if there are multiple possible outputs for any particular input. Although this seems like a simple definition, there are many subtleties. What do we mean by a “particular input?” Does the time at which the input arrives matter? What do we mean by a “uniquely defined output?” Should we consider how the system behaves when its implementation hardware fails?

Any statement about the determinism of a physical “implemented” system is fundamentally a religious or philosophical assertion, not a scientific one. We may assert that no real physical system is determinate. How will it behave when it is crushed, for example? Or we may conversely assert that everything in the physical world is preordained, a concept that we find farfetched, difficult to refute, and not very useful.

For models, however, we can make definitive assertions about their determinism. For example, a procedure defined in a programming language may be determinate in that the returned value of the procedure depends only on the arguments. No actual realization of the procedure is actually determinate in an absolute sense (the hardware may fail and no returned value will be produced at all). The procedure is a **model** defined within a **formal framework** (the semantics of the language). It models the execution of a machine abstractly, omitting information. The time at which the inputs are provided makes no difference to the model, so time is not part of what we mean by a “particular input.” The inputs and outputs are just data, and the procedure defines the relationship between the inputs and outputs. This point about models is supported by **Box and Draper (1987)**, who state “Essentially, all models are wrong, but some are useful.” The usefulness of a model depends on the model **fidelity**, the degree to which a model accurately imitates the system being modeled. But models are *always* an approximation.

Many popular modeling languages based on block diagrams have quite **weak semantics**. It is common, for example, for modeling languages to adopt a block diagram notation without giving precise meaning to the lines drawn between blocks; they vaguely represent the fact that components interact. (For examples, see the sidebar on page 14.) Modeling languages with weak semantics are harder to analyze. Their value lies instead in their ability to informally communicate design concepts among humans.

1.4 Actor Models

Ptolemy II is based on a class of models called **actor-oriented models**, or more simply, **actor models**. **Actors** are components that execute concurrently and share data with each other by sending messages via ports.

Example 1.2: Consider, for example, the Ptolemy model shown in Figure 1.3. This model shows three actors, each of which has one port. Actor A sends messages to actors B and C via its port (the Relation diamond indicates that the output from A goes to both B and C).

The sum of all of the messages communicated via a port is referred to as a **signal**. The **Director** block in the example specifies the **domain** (and hence the **model of computation**). Most of this book is devoted to explaining the various domains that have been realized in Ptolemy II.

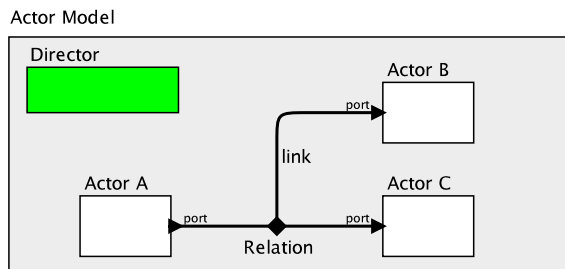


Figure 1.3: Visual rendition of a simple actor model.

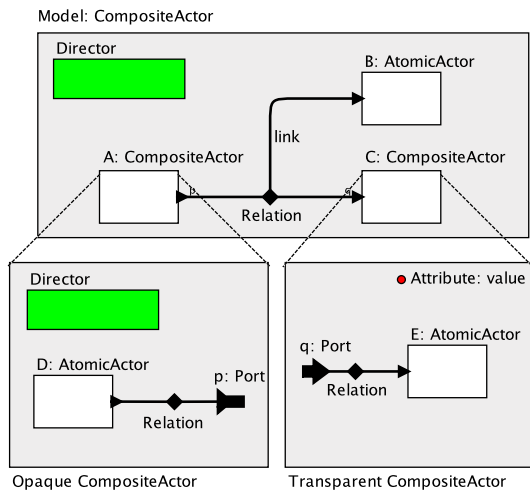


Figure 1.4: A hierarchical actor model consisting of a top-level composite actor and two submodels, each of which is also a composite actor.

1.5 Model Hierarchy

Models of complex systems are often complex. There is an art (the art of **model engineering**) to constructing good models of complex systems. A good model of a complex system provides relatively simple views of the system in order to facilitate easier understanding and analysis. A key approach to creating models with simplified views is to use modeling **hierarchy**, where what appears to be single component in one model is, internally, another model.

A hierarchical actor model is shown in Figure 1.4. It is an elaboration of Figure 1.3 where actors A and C are revealed to be themselves actor models. An **atomic actor** (where **atomic** comes from the ancient Greek **atomos**, meaning indivisible), is one that is not internally defined as an actor model. A **composite actor**, in contrast, is itself a composition of other actors. The ports p and q in the figure bridge the levels of hierarchy. A communication from D, for example, will eventually arrive at actor E after traversing ports and levels of the hierarchy.

1.6 Approaches to Heterogeneous Modeling

There are many approaches to heterogeneous modeling (Brooks et al., 2008). In **multi-view modeling**, distinct and separate models of the same system are constructed to model different aspects of a system. For example, one model may describe dynamic behavior, while another describes physical design and packaging. In **amorphous heterogeneity**, distinct modeling styles are combined in arbitrary ways within the same model without the benefit of structure. For example, some component interactions in a model may use rendezvous messaging (where both a sender and a receiver must be ready before a communication can occur), while others use asynchronous message passing (where the receiver receives the communication at some indeterminate time after the sender sends it). In **hierarchical multimodeling**, hierarchical compositions of distinct modeling styles are combined to take advantage of the unique capabilities and expressiveness of each style.

Sidebar: About the Term “Actors”

Our notion of actor-oriented modeling is related to the term “actor” as introduced in the 1970’s by Hewitt to describe the concept of autonomous reasoning agents (Hewitt, 1977). The term evolved through the work of Agha and others to describe a formalized model of concurrency (Agha et al., 1997). Agha’s actors each have an independent thread of control and communicate via asynchronous message passing. The term “actor” was also used in Dennis’s **dataflow** models (Dennis, 1974) of discrete atomic computations that react to the availability of inputs by producing outputs sent to other actors.

In this book, the term “actor” embraces a larger family of models of concurrency. They are often more constrained than general message passing and do not necessarily conform with a dataflow semantics. Our actors are still conceptually concurrent, but unlike Agha’s actors, they need not have their own thread of control. Unlike Dennis’ actors, they need not be triggered by input data. Moreover, although communication is still achieved through some form of message passing, it need not be asynchronous.

Actors are *components* in systems and can be compared to **objects**, software components in object-oriented design. In prevailing object-oriented languages (such as Java, C++, and C#), the interfaces to objects are primarily **methods**, which are procedures that modify or observe the state of objects. By contrast, the actor interfaces are primarily **ports**, which send and receive data. They do not imply the same sequential transfer of control that procedures do, and hence they are better suited to concurrent models.

Sidebar: Actors in UML, SysML, and MARTE

The **Object Management Group (OMG)** has standardized a number of notations that relate strongly to the block diagram syntax common in actor models. The actor models in this book relate to **composite structure diagrams** of UML 2 (the second version of the **unified modeling language**) (Bock, 2006; Booch et al., 1998), or more directly its derivative **SysML** (Object Management Group (OMG), 2008a). The **internal block diagram** notation of SysML, particularly with the use of flow ports, is closely related to actor models. In SysML, the actors are called “blocks.” (The term “actor” is used in UML for another purpose.)

SysML, however, emphasizes how model diagrams are rendered (their visual syntax), and leaves many details open about what the diagrams mean and how the models operate (their semantics). For example, although the SysML declares that “flow ports are intended to be used for asynchronous, broadcast, or send-and-forget interactions” (Object Management Group (OMG), 2008a), there is nothing like an **MoC** in SysML. Different SysML tools may give different behavior to flow ports and still be compliant with the standard. A single SysML model may represent multiple designs, and the behavior of the model may depend on the tools used to interpret the model. The emphasis of SysML is on standardizing the notation, not the meaning of the notation.

In contrast, the emphasis in Ptolemy II is on the semantics of models, rather than on how they are rendered (visually or otherwise). The visual notation is incidental, and in fact is not the only representation for a Ptolemy II model. Ptolemy II **directors** give models a very specific meaning. This concrete meaning ensures that a model means the same thing to different observers, and enables interoperation of heterogeneous models.

MARTE (modeling and analysis of real-time and embedded systems) puts more emphasis than SysML on the behavior of models (Object Management Group (OMG), 2008b). It avoids “constraining” the execution semantics, making the standard flexible, enabling representation of many prevalent real-time modeling techniques. In contrast, the emphasis in Ptolemy II is less on capturing existing design practices, and more on providing precise and well-defined models of system behavior. MARTE, interestingly, includes a **multiform time** model (André et al., 2007) not unlike that supported by Ptolemy II.

An example familiar to software engineers is Statecharts (Harel, 1987), which hierarchically combines *synchronous* concurrent composition with *finite state machines*. Another example of hierarchical modeling is **cosimulation**, where two distinct simulation tools are combined using a standardized interface such as the Simulink S function interface or the functional mockup interface (**FMI**) from the Modelica Association.⁶ It is also possible to support heterogeneous modeling by creating very flexible or underspecified modeling frameworks that can be adapted to cover the models of interest. The downside of this approach is *weak semantics*. The goal of Ptolemy II is to achieve *strong semantics*, yet embrace heterogeneity and provide mechanisms for heterogeneous models to interact concurrently.

As shown in Figure 1.4, one can partition a complex model into a hierarchical tree of nested submodels. At each level, the submodels can be joined together to form a network of interacting actors. Ptolemy II constrains each level of the hierarchy to be locally homogeneous, using a common *model of computation*. These homogeneous networks can then be hierarchically combined to form a larger heterogeneous model. The advantage to this approach is that each part of the system can be modeled using the model of computation that provides the best match for its processing requirements — yet each model of computation provides *strong semantics* to ensure that it is relatively easy to understand, analyze, and execute.

In Ptolemy II, a *director* defines the semantics of a model. In Figure 1.4, there are two directors. The one at the top level defines the interaction between actors A, B, and C. Since C does not internally contain a director, the same top-level director governs the interactions with actor E. Actor C is called a **transparent composite actor**; its contained model is visible to its director.

In contrast, actor A internally contains another director. That inside director governs the interaction of actors within the sub model (in this simple example, there is only one such actor, but there could be more). Actor A is called an **opaque composite actor**, and its contents are hidden from A's outside director. To distinguish the two directors, we call the outside director the **executive director**. To the executive director, actor A looks just like an atomic actor. But internally, A contains another model.

The directors at different levels of the hierarchy need not implement the same MoC. Opaque composite actors, therefore, are Ptolemy's way of realizing hierarchical multi-modeling and cosimulation.

⁶<http://www.functional-mockup-interface.org>

Sidebar: Plurality of Models

Occam's razor is a principle in science and engineering that encourages selection of those theories and hypotheses that require the fewest assumptions, postulates, or entities to explain a given phenomenon. The principle can be expressed as “entities must not be multiplied beyond necessity” (*entia non sunt multiplicanda praeter necessitatem*) or as “plurality should not be posited without necessity” (*pluralitas non est ponenda sine necessitate*) ([Encyclopedia Britannica, 2010](#)). The principle is attributed to 14th-century English logician, theologian and Franciscan friar William of Ockham.

Despite its compelling value, the principle has limitations. Immanuel Kant, for example, felt a need to moderate the effects of Occam's razor, stating “the variety of beings should not rashly be diminished.” (*entium varietates non temere esse minuendas*) ([Smith, 1929](#)). Einstein allegedly remarked, “everything should be made as simple as possible, but not simpler” ([Shapiro, 2006](#)).

When applied to design techniques, Occam's razor biases us towards using fewer and simpler design languages and notations. However, experience indicates that both redundancy and diversity can be beneficial. For example, there is benefit to using UML class diagrams even if the information they represent is already encoded in a C++ program. There is also value in UML use-case diagrams, which express concepts that are not encoded in the C++ program and are also not (directly) represented in the UML class diagram. The three representations serve different purposes, though they represent the same underlying process.

The fact that many different notations are used in UML and its derivatives runs counter to the principle in Occam's razor. Ironically, the [unified modeling language \(UML\)](#) originated in the 1990s to *reduce* the diversity of notations used to express object-oriented software architectures ([Booch et al., 1998](#)). So what is gained by this anti-razor?

Design of software systems is essentially a creative process; engineers create programs that did not previously exist. Occam's razor should be applied only cautiously to creative processes, because creativity often flourishes when there are multiple media with which to achieve the desired effect. UML facilitates the creative process by offering more abstract notations than C++ source code, and these notations encourage experimentation with design and communication of design ideas.

Sidebar: About Heterogeneous Models

Some authors use the term **multi-paradigm modeling** to describe approaches that mix models of computation (Mosterman and Vangheluwe, 2004). Ptolemy II focuses on techniques that combine actors with multi-paradigm modeling. An early systematic approach to such mixed models was realized in Ptolemy Classic (Buck et al., 1994), the predecessor to Ptolemy II (Eker et al., 2003). Influenced by the Ptolemy approach, SystemC is capable of realizing multiple MoCs (Patel and Shukla, 2004; Herrera and Villar, 2006). So are ModHel’X (Hardebolle and Boulanger, 2007) and ForSyDe (Jantsch, 2003; Sander and Jantsch, 2004).

Another approach supports mixing concurrency and communication mechanisms without the structural constraints of hierarchy (Goessler and Sangiovanni-Vincentelli, 2002; Basu et al., 2006). A number of other researchers have tackled the problem of heterogeneity in creative ways (Burch et al., 2001; Feredj et al., 2009).

It is also possible to use **tool integration**, where different modeling tools are combined either through interchange languages or through co-simulation (Liu et al., 1999; University of Pennsylvania MoBIES team, 2002; Gu et al., 2003; Karsai et al., 2005). This approach is challenging, however, and yields fragile tool chains. Many tools lack documentation on how and where they can be extended to enable cross-tool integration; implementing and maintaining integration requires considerable effort. Challenges include API incompatibilities, unstable or undocumented APIs, unclear semantics, syntactic incompatibilities, and unmaintainable code bases. Tool integration proves to be a painful way to accomplish heterogeneous design. A better approach is to focus on the semantics of interoperation, rather than on the software problems of tool integration. Good software architectures for interoperation will emerge only from a good understanding of the semantics of interoperation.

In Ptolemy, each model contains a director that specifies the MoC being used and provides either a code generator or an interpreter for the MoC (or both). An interesting alternative is given by “42” (Maraninchi and Bhouhadiba, 2007), which integrates a custom MoC with the model.

Sidebar: Tools Supporting Heterogeneous Models

Several widely used tools provide fixed combinations of a few MoCs. Commercial tools include Simulink/StateFlow (from The MathWorks), which combines continuous- and discrete-time actor models with finite-state machines, and LabVIEW (from National Instruments), which combines dataflow actor models with finite-state machines and a time-driven MoC. Statemate (Harel et al., 1990) and SCADE (Berry, 2003) combine finite-state machines with a synchronous/reactive formalism (Benveniste and Berry, 1991). Giotto (Henzinger et al., 2001) and TDL (Pree and Tempel, 2006) combine FSMs with a time-driven MoC. Several hybrid system modeling and simulation tools combine continuous-time dynamical systems with FSMs (Carloni et al., 2006).

The Y-chart approach supports heterogeneous modeling and is popular for hardware-software codesign (Kienhuis et al., 2001). This approach separates modeling of the hardware implementation from modeling of application behavior (a form of multi-view modeling), and provides mechanisms for bringing these disparate models together. These mechanisms allow developers to trade off hardware cost and complexity with software design. Metropolis is a particularly elegant tool for this purpose (Goessler and Sangiovanni-Vincentelli, 2002). It introduces a “quantity manager” that mediates interactions between the desired functionality and the resources required to implement that functionality.

Modelica (Fritzson, 2003; Modelica Association, 2009) also has actor-like semantics in the sense that components are concurrent and communicate via ports, but the ports are neither inputs nor outputs. Instead, the connections between ports declare equation constraints on variables. This approach has significant advantages, particularly for specifying physical models based on differential-algebraic equations (DAEs). However, the approach also appears to be harder to combine heterogeneously with other MoCs.

DESTTECS (design support and tooling for embedded control software) is a tool supported by a consortium from academia and industry that has a focus on fault-tolerant embedded systems (Fitzgerald et al., 2010). This tool integrates continuous-time models made in 20-sim (Broenink, 1997) and discrete-event models in VDM (Vienna Development Method) (Fitzgerald et al., 2008). DESTTECS synchronizes time and passes variables between the two tools.

1.7 Models of Time

Some models of computation have a notion of **time**. Specifically, this means that communication between actors and computation performed by actors occurs on a logical time line. Even more specifically, this means that there is a notion of two actions (communication or computation) being either ordered in time (one occurs before the other) or being simultaneous. A notion of time may also have a metric, meaning (loosely) that the time gap between two actions may be measured.

A key mechanism that Ptolemy II provides for interoperability of domains is a coherent notion of time. This mechanism has proven effective even for combining models of computation that have no notion of time (such as [dataflow](#) models and [finite state machines](#)), with models of computation that depend strongly on time (such as [discrete-event](#) models and [continuous-time](#) models). In this section, we outline key features of this mechanism.

1.7.1 Hierarchical Time

The model hierarchy discussed in Sections [1.5](#) and [1.6](#) is central to the management of time. Typically, only the top-level director advances time. Other directors in a model obtain the current model time from their enclosing director. If the top-level director does not implement a timed model of computation, then time does not advance. Hence, timed models always contain a top-level director that implements a timed model of computation.

Timed and untimed models of computation may be interleaved in the hierarchy. As we will discuss later, however, there are certain combinations that do not make sense, while other combinations are particularly useful, particularly the [modal models](#) discussed in Chapter [8](#).

Time can also advance non-uniformly in a model. In the [modal models](#) of Chapter [8](#), the advancement of time can be temporarily suspended in a submodel ([Lee and Tripakis, 2010](#)). More generally, as explained in Chapter [10](#), time may also progress at different rates at different levels of the hierarchy. This feature is particularly useful for modeling distributed systems where maintaining a perfectly coherent uniform time base is not physically possible. It is referred to as **multiform time**, and it enables highly realistic models that explicitly recognize that time can only be imperfectly measured.

1.7.2 Superdense Time

In addition to providing multiform time, Ptolemy II provides a model of time known as **superdense time** (Manna and Pnueli, 1993; Maler et al., 1992; Lee and Zheng, 2005; Cataldo et al., 2006). A superdense time value is a pair (t, n) , called a **time stamp**, where t is the **model time** and n is a **microstep** (also called an **index**). The model time represents the time at which some event occurs, and the microstep represents the sequencing of events that occur at the same model time. Two time stamps (t, n_1) and (t, n_2) can be interpreted as being **simultaneous** (in a weak sense) even if $n_1 \neq n_2$. A stronger notion of **simultaneity** would require the time stamps to be equal (both in model time and microstep). An example illustrates the value of superdense time.

Example 1.3: To understand the role of the microstep, consider Newton’s cradle, a toy with five steel balls suspended by strings, shown in Figure 1.5. If you lift the first ball and release it, it strikes the second ball, which does not move. Instead, the fifth ball reacts by rising.



Figure 1.5: Newton’s cradle. Image by Dominique Toussaint, made available under the terms of the [GNU Free Documentation License](#), Version 1.2 or later.

Consider the momentum p of the second ball as a function of time. The second ball does not move, so its momentum must be everywhere zero. But the momentum of the first ball is somehow transferred to the fifth ball, passing through the second ball. So the momentum cannot be always zero.

Let \mathbb{R} represent the real numbers. Let $p: \mathbb{R} \rightarrow \mathbb{R}$ be a function that represents the momentum of this second ball, and let τ be the time of the collision. Then

$$p(t) = \begin{cases} P & \text{if } t = \tau \\ 0 & \text{otherwise} \end{cases} \quad (1.1)$$

for some constant P and for all $t \in \mathbb{R}$. Before and after the instant of time τ , the momentum of the ball is zero, but at time τ , it is not zero. Momentum is proportional to velocity, so

$$p(t) = Mv(t),$$

where M is the mass of the ball. Hence, combining with (1.1),

$$v(t) = \begin{cases} P/M & \text{if } t = \tau \\ 0 & \text{otherwise.} \end{cases} \quad (1.2)$$

The position of a mass is the integral of its velocity,

$$x(t) = x(0) + \int_0^t v(\tau) d\tau,$$

where $x(0)$ is the initial position. The integral of the function given by (1.2) is zero at all t , so the ball does not move, despite having a non-zero momentum at an instant.

The above physical model mostly works to describes the physics, but has two flaws. First, it violates the basic physical principle of conservation of momentum. At the time of the collision, all three middle balls will simultaneously have non-zero momentum, so seemingly, aggregate momentum has magically increased. Second, the model cannot be directly converted into a discrete representation.

A **discrete** representation of a signal is a sequence of values that are ordered in time (for mathematical details, see the sidebar on page 334). Any such representation of the momentum in (1.1) or velocity in (1.2) is ambiguous. If the sequence does not include the value at the time of the collision, then the representation does not capture the fact that momentum is transferred through the ball. If the representation does include the value at the time of the collision, then the representation is

indistinguishable from a representation of a signal that has a non-zero momentum over some interval of time, and therefore models a ball that does move. In such a discrete representation, there is no semantic distinction between an instantaneous event and a rapidly varying continuous event.

Superdense time solves both problems. Specifically, the momentum of the second ball can be unambiguously represented by a sequence of samples where $p(\tau, 0) = 0$, $p(\tau, 1) = P$, and $p(\tau, 2) = 0$, where τ is the time of the collision. The third ball has non-zero momentum only at superdense time $(\tau, 2)$. At the time of the collision, each ball first has zero momentum, then non-zero, then zero again, all in an instant. The event of having non-zero momentum is weakly simultaneous for all three middle balls, but not strongly simultaneous. Momentum is conserved, and the model is unambiguously discrete.

One could argue that the physical system is not actually discrete. Even well-made steel balls will compress, so the collision is actually a continuous process, not a discrete event. This is true, but when building models, we do not want the modeling formalism to force us to construct models that are more detailed than is appropriate. Such a model of Newton's cradle would be far more sophisticated, and the resulting non-linear dynamics would be far more difficult to analyze. The [fidelity](#) of the model would improve, but at a steep price in understandability and analyzability.

The above example shows that physical processes that include instantaneous events are better modeled using functions of the form $p: \mathbb{R} \times \mathbb{N} \rightarrow \mathbb{R}$, where \mathbb{N} represents the natural numbers, rather than the more conventional $p: \mathbb{R} \rightarrow \mathbb{R}$. The latter is adequate for continuous processes, but not for discrete events. At any time $t \in \mathbb{R}$, the signal p has a sequence of values, ordered by their microsteps. This signal cannot be misinterpreted as a rapidly varying continuous signal.

We say that two time stamps (t_1, n_1) and (t_2, n_2) are **weakly simultaneous** if $t_1 = t_2$, and **strongly simultaneous** if, in addition, $n_1 = n_2$.

Thus we can represent causally-related, but weakly simultaneous events. A [signal](#) may have two *distinct* events at with time stamps (t, n_1) and (t, n_2) , where $n_1 \neq n_2$. A signal may therefore include weakly simultaneous, but distinct, events. Two distinct signals may contain strongly simultaneous events, but a single signal cannot contain two distinct strongly simultaneous events. This model of time unambiguously represents dis-

crete events, discontinuities in continuous-time signals, and sequences of zero-time events in discrete signals.

Superdense time is ordered lexicographically (like a dictionary), which means that $(t_1, n_1) < (t_2, n_2)$ if either $t_1 < t_2$, or $t_1 = t_2$ and $n_1 < n_2$. Thus, an event is considered to occur before another if its model time is less or, if the model times are the same, if its microstep is lower. Time stamps are a particular realization of **tags** in the tagged-signal model of Lee and Sangiovanni-Vincentelli (1998).

1.7.3 Numeric Representation of Time

Computers cannot perfectly represent real numbers, so a time stamp of form $(t, n) \in \mathbb{R} \times \mathbb{N}$ is not realizable. Many software systems approximate a time t using a double-precision floating point number. But such a representation has two serious disadvantages. First, the **precision** of a number (how close it is to the next smaller or large representable number) depends on its magnitude. Thus, as time increases in such systems, the precision with which time is represented decreases. Second, addition and subtraction can introduce quantization errors in such a representation, so it is not necessarily true that $(t_1 + t_2) + t_3 = t_1 + (t_2 + t_3)$. This significantly weakens the semantic notion of **simultaneity**, since whether two events are (weakly or strongly) simultaneous may depend on how their time stamps were computed.

Ptolemy II solves this problem by making the **time resolution** a single, global constant. Model time is given as $t = mr$, where m is an arbitrarily large integer, and the time resolution r is a double-precision floating point number. The multiple m is realized as a Java BigInteger (an arbitrarily large integer), so it will never overflow. The time resolution r , a *double*, is a parameter shared by all the directors in a model. A model, therefore, has the same time resolution throughout its hierarchy and throughout its execution, no matter how big time gets. Moreover, addition and subtraction of time values does not suffer quantization errors. By default, the time resolution is $r = 10^{-10}$, which may represent one tenth of a nanosecond. Then, for example, $m = 10^{11}$ represents 10 seconds.

In Ptolemy II, the microstep n in a time stamp (t, n) is represented as an *int*, a 32-bit integer. The microstep, therefore, is vulnerable to overflow. Such overflow may be prevented by avoiding models that have **chattering Zeno** behavior, as discussed in Chapter 7.

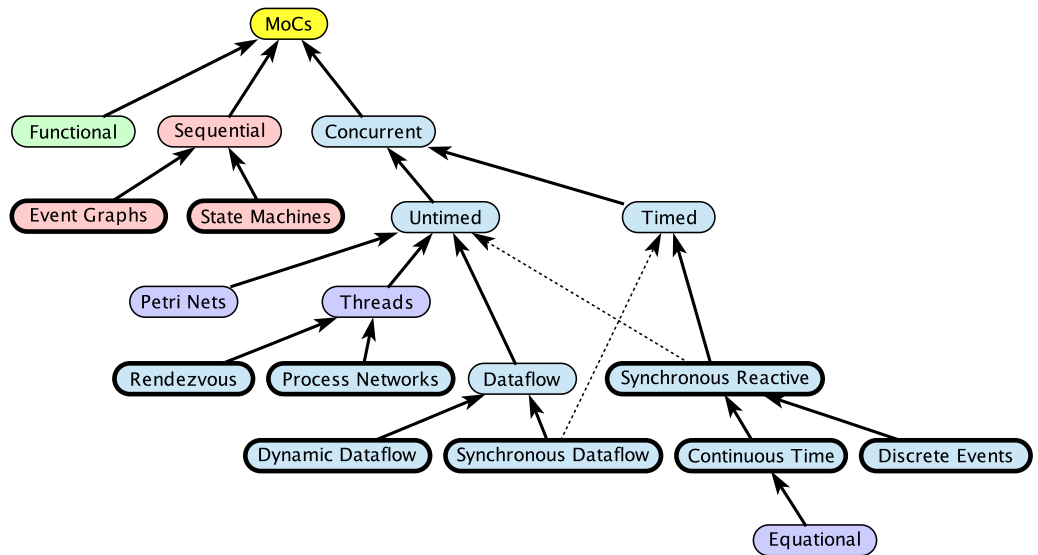


Figure 1.6: Summary of the relationship between models of computation. The ones with bold outlines are covered in detail in this Book.

1.8 Overview of Domains and Directors

In Ptolemy II an implementation of a model of computation is called a **domain**.⁷ In this section, we briefly describe domains that have been realized in Ptolemy II. This is not a complete list; the intent is to show the diversity of the models of computation under study. These and other domains are described in subsequent chapters in more detail. Figure 1.6 summarizes the relationships between these domains.

All of the domains discussed here ensure **determinism** unless the model explicitly specifies nondeterministic behavior. That is, nondeterminism, if desired, must be explicitly built into the models; it does not arise accidentally from **weak semantics** in the modeling framework. A domain is said to be **determinate** if the **signals** sent between actors —

⁷The term “domain” comes from a fanciful notion in astrophysics, that there are regions of the universe with different sets of laws of physics. A model of computation represents the “laws of physics” of the submodel it governs.

including data values carried by messages, their order, and their [time stamps](#) — do not depend on arbitrary scheduling decisions, despite the concurrency in the model. Ensuring determinism is far from trivial in concurrent MoCs, and providing reasonable nondeterminate mechanisms is also challenging. The goal is that, when a model includes nonde-terminate behavior, it should be explicitly specified by the builder of the model; it should not appear accidentally, nor should it surprise the user.

Dataflow. Ptolemy II includes several [dataflow](#) domains, described in Chapter 3. The execution of an actor in dataflow domains consists of a sequence of **firings**, where each firing occurs as a reaction to the availability of input data. A firing is a (typically small) computation that consumes the input data and produces output data.

The [synchronous dataflow](#) (SDF) domain ([Lee and Messerschmitt, 1987b](#)) is particularly simple, and is possibly the most used domain of all. When an actor is executed in SDF, it consumes a fixed amount of data from each input port, and produces a fixed amount of data to each output port. An advantage of the SDF domain is that (as described in Chapter 3) the potential for deadlock and boundedness can be statically checked, and schedules (including parallel schedules) can be statically computed. Communication in this domain is realized with **first-in, first-out (FIFO)** queues with fixed finite capacity, and the execution order of components is statically scheduled. SDF can be timed or untimed, though it is usually untimed, as suggested in Figure 1.6.

In contrast, the [dynamic dataflow](#) (DDF) domain is more flexible than SDF and computes schedules on the fly. In DDF, the capacity of the FIFO queues is not bounded. DDF is useful when communication patterns between actors are dependent on the data that is passed between actors.

Dataflow models are ideal for representing **streaming** systems, where sequences of data values flow in relatively regular patterns between components. Signal processing systems, such as audio and video systems, for example, are a particularly good match.

Process Networks. In the [process network](#) (PN) domain, described in Chapter 4, actors represent concurrent processes that communicate by (conceptually infinite capacity) FIFO queues ([Lee and Parks, 1995](#)). Writing to the queues always succeeds immediately, while reading from an empty queue blocks the reader process. The simple blocking-read, nonblocking-write strategy ensures the determinacy of the model ([Kahn and MacQueen, 1977](#)). Nevertheless, we have extended the model to support certain forms of nonde-terminism. Each actor executes in its own Java thread, so on multicore machines they

can execute in parallel. This domain is untimed. The PN domain realizes a generalization of dataflow where instead of discrete firings, actors represent continually executing processes (Lee and Matsikoudis, 2009).

PN is suitable for describing concurrent processes that communicate asynchronously by sending messages to one another. Messages are eventually delivered in the same order they are sent. Message delivery is presumed to be reliable, so the sender does not expect nor receive any confirmation. This domain has a “send and forget” flavor.

PN also provides a relatively easy way to get parallel execution of models. Each actor executes in its own thread, and most modern operating systems will automatically map threads onto available cores. Note that if the actors are relatively fine-grained, meaning that they perform little computation for each communication, then the overhead of multithreading and inter-thread communication may overwhelm the performance advantages of parallel execution. Thus, model builders should expect performance advantages only for coarse-grained models.

Rendezvous. The *Rendezvous* domain, also described in Chapter 4, is similar to PN in that actors represent concurrent processes. However, unlike PN’s “send and forget” semantics, in the *Rendezvous* domain, actors communicate by atomic instantaneous data exchanges. When one actor sends data to another, the sender will block until the receiver is ready to receive. Similarly, when one actor attempts to read input data, it will block until the sender of the data is ready to send the data. As a consequence, the process that first reaches a rendezvous point will stall until the other process reaches the same rendezvous point (Hoare, 1978). It is also possible in this domain to create multi-way rendezvous, where several processes must all reach the rendezvous point before any process can continue. Like PN, this domain is untimed, supports explicit nondeterminism, and can transparently leverage multicore machines.

The *Rendezvous* domain is particularly useful for modeling asynchronous resource contention problems, where a single resource is shared by multiple asynchronous processes.

Synchronous-Reactive. The *synchronous-reactive* (SR) domain, described in Chapter 5, is based on the semantics of synchronous languages (Benveniste and Berry, 1991; Halbwachs et al., 1991; Edwards and Lee, 2003a). The principle behind synchronous languages is simple, although the consequences are profound. Execution follows “ticks” of a global “clock.” At each tick, each variable (represented visually in Ptolemy II by the wires that connect the blocks) may or may not have a value. Its value (or absence of

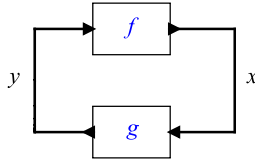


Figure 1.7: A simple feedback system.

value) is given by an actor whose output port is connected to the wire. The actor realizes a function that maps the values at its input ports to the values at its output ports (the function can vary from tick to tick). For example, in Figure 1.7, the variables x and y at a particular tick are related by

$$x = f(y), \text{ and } y = g(x).$$

The task of the domain’s director is to find, at each tick, values of x and y that solve these equations. This solution is called a **fixed point**. The SR domain is by default untimed, but it can optionally be timed, in which case there is a fixed time interval between ticks.⁸

The SR domain is similar to dataflow and PN in that actors send streams of data to one other. Unlike dataflow, however, the streams are synchronized; at a tick of the clock, every communication path either has a message, or the message is unambiguously absent. Dataflow models, by contrast, are more asynchronous; a message may be “absent” simply because it hasn’t arrived yet due to an accident of scheduling. To prevent nondeterminism, PN and dataflow have no semantic notion of an “absent” input. Inputs always have messages (or will have messages, in which case the actor is required to wait for the messages to arrive).

SR is well suited to situations with more complex control flow, where an actor may take different actions depending on whether a message is present or not. By synchronizing actions, the domain handles these scenarios without nondeterminism. SR is less concurrent than dataflow or PN, since each tick of the clock must be tightly orchestrated. As a consequence, it is harder to execute in parallel.

Finite-State Machine. The *finite state machine* (FSM) domain, described in Chapter 6, is the only domain discussed here that is not concurrent. The components in this domain

⁸If you need a variable time interval between ticks, you can accomplish this by placing an SR model within a DE model.

are not actors, but rather represent [states](#), and the relations represent not communication paths, but rather transitions between states. Transitions have [guards](#) that determine when state transitions occur.

An FSM can be used to define the behavior of an actor used in any of the other domains. The actor can have any number of input and output ports. When that actor executes, the FSM reads the inputs, evaluates the guards to determine which transition to take, and produces outputs as specified on the selected transition. FSMs can also have local variables whose values can be modified by transitions (providing a model of computation that is known as an [extended state machine](#)).

An FSM can also be used to create a rich class of hierarchical models known as [modal models](#), discussed in Chapter 8 ([Lee and Tripakis, 2010](#)). In a modal model, states of an FSM contain submodels that process inputs and produce outputs. Each state of the FSM represents a mode of execution, and the [mode refinement](#) defines the behavior in that mode. The mode refinement is a submodel with its own director that is active only when the FSM is in the corresponding state. When a submodel is not active, its local time does not advance, as explained above in Section 1.7.1.

Discrete Event. In the [discrete-event \(DE\)](#) domain, described in Chapter 7, actors communicate through events placed on a time line. Each event has a value and a [time stamp](#), and actors process events in chronological order. The output events produced by an actor are required to be no earlier in time than the input events that were consumed. In other words, actors in DE are causal.

The execution of this model uses a global event queue. When an actor generates an output event, the event is slotted into the queue according to its time stamp. During each iteration of a DE model, the events with the smallest time stamp are removed from the global event queue, and their destination actor is fired. The DE domain supports simultaneous events. At each time where at least one actor fires, the director computes a fixed point, similar to SR ([Lee and Zheng, 2007](#)). DE is closely related to the well-known DEVS (discrete event system specification) formalism ([Zeigler et al., 2000](#)), which is widely used for simulating large, complex systems. The semantics of the Ptolemy II variant of DE is given by [Lee \(1999\)](#).

DE is well suited for modeling the behavior of complex systems over time. It can model networks, digital hardware, financial systems, and human organizational systems, for example. Chapter 10 shows how DE can be extended to leverage [multiform time](#).

Continuous time. The [Continuous](#) time domain ([Lee and Zheng, 2005](#)), described in Chapter 9, models ordinary differential equations (ODEs), while also supporting discrete events. Special actors that represent *integrators* are connected in feedback loops in order to represent the ODEs. Each connection in this domain represents a continuous-time function, and the components denote the relations between these functions.

The Continuous model computes solutions to ODEs using numerical methods. As with SR and DE, at each instant, the director computes a fixed point for all signal values ([Lee and Zheng, 2007](#)). In each iteration, time is advanced by an amount determined by the ODE solver. To advance time, the director chooses a time stamp with the help of a solver and speculatively executes actors through this time step. If the time step is sufficiently small (key events such as level crossings, mode changes, or requested firing times are not skipped over, and the numerical integration is sufficiently accurate), then the director commits the time increment.

The Continuous director interoperates with all other timed Ptolemy II domains. Combining it with FSMs yields a particular form of modal model known as a [hybrid system](#) ([Lee and Zheng, 2005](#); [Lee, 2009](#)). Combinations with discrete-event and synchronous/reactive domains are also useful ([Lee and Zheng, 2007](#)).

Ptera. The [Ptera](#) domain, described in Chapter 11, realizes a variant of [event graphs](#). In Ptera, the components are not actors. Instead, the components are events, and the connections between components are triggering relations for events. A Ptera model represents how events in a system can trigger other events. Ptera is a timed model, and like FSM, it can be used to define the behavior of an actor to be used in another domain. In addition, events can be composites in that they have actions associated with them that are themselves defined by a submodel specified using another domain. Ptera is useful for specifying timed behaviors where input events may trigger chain reactions.

1.9 Case Study

[Cyber-physical systems](#) are intrinsically heterogeneous. CPS models, therefore, benefit from being able to combine models of computation. In this section, we walk through an example that uses several models of computation. The example is highly simplified, but with a little imagination, it is easy to see how the model can evolve to become an accurate and complete model of a large complex system. In particular, the large complex system we have in mind is an electric power system in a smart grid or on a vehicle (such as

an airplane or advanced ground vehicle). In such a system, there are multiple sources of electric power (windmills, solar panels, turbines, backup generators, etc.) that must be coordinated to provide power to a multiplicity of loads. The system includes controllers that regulate the generators to keep voltages and frequencies near constant, and supervisory controllers that connect and disconnect loads and generators to provide services and handle faults to protect equipment. Such a system also includes networks whose dynamics may affect the overall behavior of the system.

Here, we illustrate a highly simplified version of such a system to show how the various MoCs come into play.

Example 1.4: A simplified model of a gas-powered generator that may be connected to and disconnected from a load is shown in Figure 1.8. This is a **continuous-time** model, as indicated by the **Continuous** director, which is discussed in Chapter 9. The model has two inputs, a *drive* signal, and a *loadAdmittance*. The output is a *voltage* signal. In addition, the model has three parameters, a time constant T , an output impedance Z , and a drive limit L . The model gives the output voltage of a generator over time as the generator gets more or less gas (specified by the *drive* input), and as the load varies (as specified by the *loadAdmittance* input).

This model exhibits simplified linear and nonlinear dynamics. The nonlinear dynamics is realized by the **Limiter** actor (see the sidebar on page 57), which limits

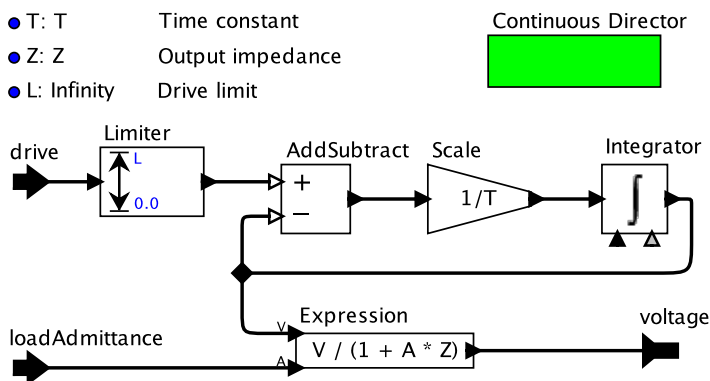


Figure 1.8: Simplified model of a gas-powered generator. [\[online\]](#)

the *drive* input. In particular, if the *drive* input becomes negative, it sets the drive to zero (you cannot extract gas from a generator). It also saturates the *drive* input at an upper bound given by the parameter *L*, which defaults to *Infinity*, meaning that there is no saturation (this generator can accept an arbitrarily large drive input).

The linear dynamics in this model is given by the small feedback loop, which includes an *AddSubtract* actor, a *Scale* actor, and an *Integrator*. If the output of the limiter is *D*, then this loop gives a value *V* that satisfies the following *ordinary differential equation*,

$$\frac{dV}{dt} = \frac{1}{T}(D - V) ,$$

where both *D* and *V* are functions of time (see Chapter 9 to understand how this model yields the above equation).

For our purposes here, understanding this equation is not important, since this part of such a model would typically be constructed by a mechanical engineer who is an expert in such models, but we can nevertheless make some intuitive observations. First, if $D = V$, then the derivative is zero, so the generator is stable and will produce an unchanging output. Second, when $D \neq V$, the feedback loop adjusts the value of *V* to make it closer to *D*. If $D > V$, then this equation makes the derivative of *V* positive, which means that *V* will increase. If $D < V$, then the derivative is negative, so *V* will decrease. In fact, the output *V* will converge to *D* exponentially with time constant *T*. A **time constant** is the amount of time that an exponential signal takes to reach $1 - 1/e \approx 63.2\%$ of its final (asymptotic) value.

The last part of the model is the part that models the effect of the load. This effect is modeled by the *Expression* actor (see Section 13.2.4), which uses Ohm's law to calculate the output voltage as a function of the value *V* (representing the generator's effort), the output impedance *Z*, and the load admittance *A*. An electrical engineer would recognize this calculation as the realization of a simple voltage divider.

For our purposes, it is sufficient to notice that if $A = 0$ (there is no load) or $Z = 0$ (the generator is an ideal voltage source with no output impedance), then the voltage output is equal to the effort *V*. A real generator, however, will have a non-zero output impedance. As the load admittance *A* increases from zero, the output voltage will drop.

The above model is about the simplest interesting model of **continuous dynamics**. To integrate this model with digital controllers, we could wrap the model in another one that defines the discrete interfaces, as shown next.

Example 1.5: The Generator model of Figure 1.8 is wrapped to provide a discrete interface in Figure 1.9. Here, the *drive* and *loadAdmittance* inputs go to instances of the *ZeroOrderHold* actor. These inputs, therefore, can be provided as *discrete events* rather than continuous-time signals. The *ZeroOrderHold* actor converts these discrete events into continuous-time signals by holding the value constant between arrivals of events (see Section 9.2).

The output voltage goes through a *PeriodicSampler* actor (see Section 9.2), which produces discrete events that are samples of the output voltage. The sample period is a parameter P of the model.

This model exposes the time constant T and output impedance Z , but hides the drive limit L . Of course, the model designer could make other choices about which parameters to expose.

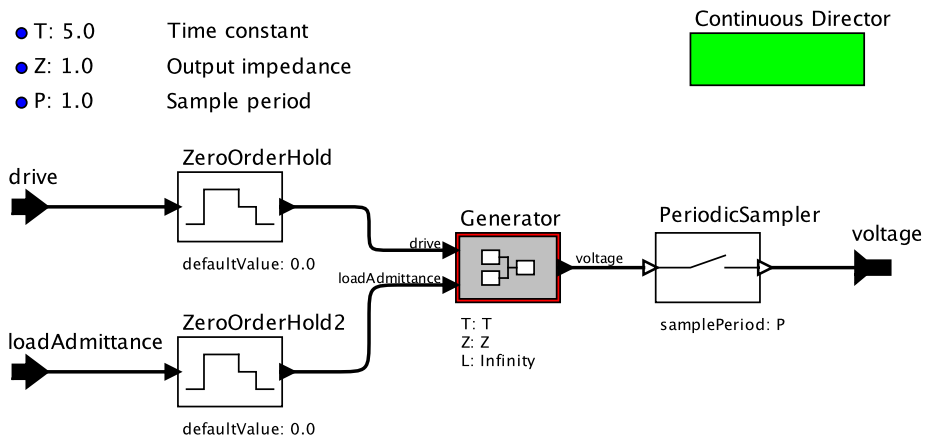


Figure 1.9: The Generator model of Figure 1.8 wrapped to provide a discrete interface. [\[online\]](#)

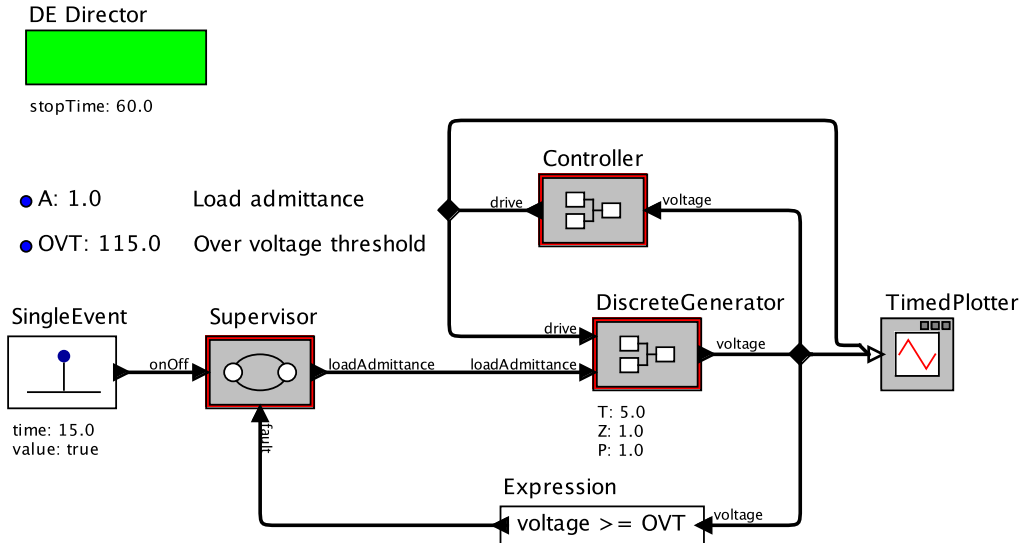


Figure 1.10: A discrete-event model with a generator, a controller, and an over-voltage protector. [\[online\]](#)

A continuous-time model may be embedded within a [discrete-event](#) model (see Chapter 7), as illustrated next.

Example 1.6: The DiscreteGenerator model of Figure 1.9 is embedded in a discrete-event model in Figure 1.10. This model has two parameters, the load admittance A and an over-voltage threshold OVT . The time constant T of the DiscreteGenerator is set to 5.0. This model includes two other components that we will explain below, a Supervisor, which provides the over-voltage protection, and a Controller, which regulates the *drive* input of the DiscreteGenerator based on measurements of the output voltage.

In addition, this model includes a simple test scenario, where a [SingleEvent](#) actor (see sidebar on page 241) requests that a load be connected at time 15.0, and a [TimedPlotter](#) actor (see Chapter 17), which displays the results of a run of the model, as shown in Figure 1.11.

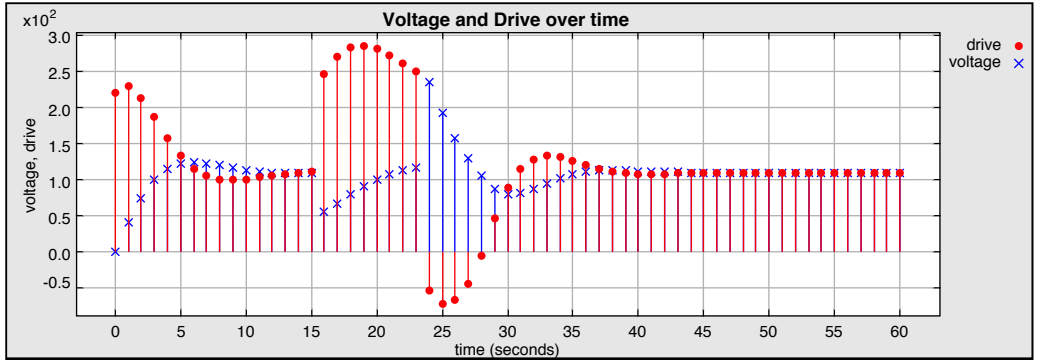
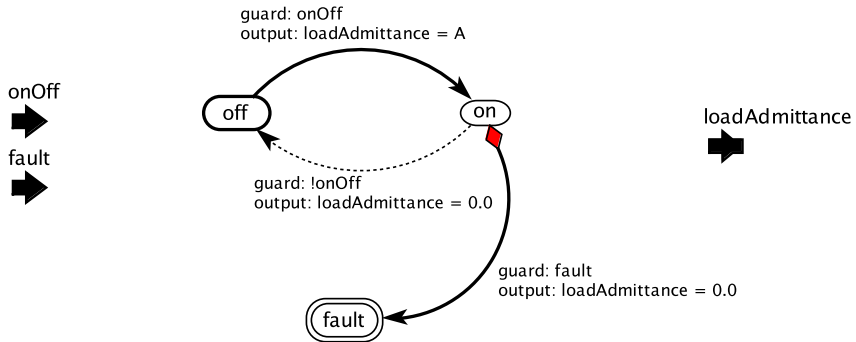


Figure 1.11: The plot produced by the model in Figure 1.10. [\[online\]](#)

In this test scenario, the load admittance is quite high (1.0) compared to the output impedance (also 1.0), so when the load is connected at time 15, the voltage abruptly drops to half its target value of 110 volts. The Controller compensates for this by substantially increasing the *drive*, but this causes the voltage to overshoot the target, and at time 24, to exceed the *OVT* threshold. The Supervisor reacts to this over-voltage condition by disconnecting the load, which causes the voltage to spike quite high, since the generator now has a substantial *drive* input. The Controller eventually brings the voltage back to the target level.

Notice further that when the load is disconnected, the Controller takes the *drive* signal negative. If this is a gas-powered generator, the Controller is trying to give the generator a negative flow of gas. Fortunately, our generator model includes a Limiter actor that prevents the model from actually providing that negative flow of gas.

The model in Figure 1.10 includes two very different kinds of controllers, a supervisory controller called Supervisor, and a low-level controller called simply Controller. These two controllers are specified using two additional MoCs, as explained next.

Figure 1.12: The Supervisor of Figure 1.10. [\[online\]](#)

Example 1.7: The Supervisor model of Figure 1.10 is a **finite state machine**, shown in Figure 1.12. The notation here is explained in Chapter 6, but we can easily grasp the general behavior.

This FSM has two inputs, *onOff* (a boolean that requests to connect or disconnect the load) and *fault* (a boolean that indicates that an over-voltage condition has occurred). It has one output, *loadAdmittance*, which will be the actually load admittance provided to the generator.

The initial state of the FSM is *off*. When an *onOff* input arrives that has value true, the FSM will transition from the *off* state to the *on* state and produce a *loadAdmittance* output with value given by *A*, a parameter of the model. This connects the load.

When the FSM is in state *on*, if a *fault* event arrives with value true, then it will transition to the **final state** *fault* and set the *loadAdmittance* to 0.0, disconnecting the load. If instead an *onOff* event arrives with value false, then it will transition to the state *off* and also disconnect the load. The difference between these two transitions is that once the FSM has entered the *fault* state, it cannot reconnect the load without a system reset (which will bring the FSM back to the initial state).

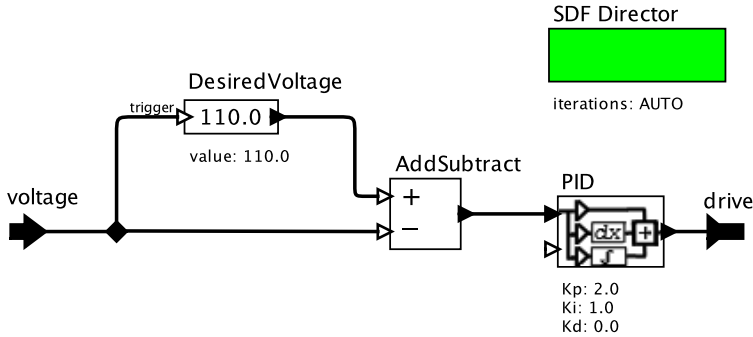


Figure 1.13: The Controller of Figure 1.10. [online]

Example 1.8: The Controller model of Figure 1.10 is the [dataflow](#) model shown in Figure 1.13. This model uses the [SDF](#) director (see Chapter 3), which is suitable for sampled-data signal processing. In this case, the controller compares the input *voltage* against a desired voltage (110 volts), and feeds the resulting error signal into a [PID](#) controller. A PID controller is a commonly used linear time-invariant system. A control engineer would know how to set the parameters of this controller, but in this case, we have simply chosen some parameters experimentally to yield an interesting test case.

Notice that the pieces of the model in Figure 1.10 are distinctly heterogeneous, touching on several disciplines within engineering and computer science. Typically, models of this type are the result of teams of engineers working together, and a framework that enables these teams to compose their models can become extremely valuable.

Many elaborations of this model are easy to envision. For example:

- The Generator could be defined as an [actor-oriented class](#), so that it can be instantiated multiple times, and yet developed and maintained in a single centralized definition (see Section 2.6).
- The Generator model could be elaborated to reflect more sophisticated linear and non-linear dynamics using the techniques discussed in Chapter 9.

- The Generator model could be elaborated to include frequency and phase effects, for example by using complex-valued impedances and admittances together with a phasor representation.
- Models with a variable size (e.g., n generators and m loads, where n and m are parameters) could be created using the [higher-order components](#) considered in Section 2.7.
- The effects of network timing, clock synchronization, and contention for shared resources could be modeled using the techniques in Chapter 10.
- Signal processing techniques such as machine learning and spectral analysis, (see Chapter 3), could be integrated into the control algorithms.
- A [units system](#) could be included to make the model precise about the units used to measure time, voltage, frequency, etc.
- An [ontology](#) could be included to make the model precise about which signals and parameters represent voltages, admittances, impedances, etc., or even to make distinctions between domain-specific concepts such as the internal voltage (effort) of a generator vs. the voltage exhibited at its output, which is affected by its output impedance and load.

1.10 Summary

Ptolemy II focuses on actor-oriented modeling of complex systems, providing a disciplined approach to heterogeneity and concurrency. The central notion in hierarchical model decomposition is that of a [domain](#), which implements a particular model of computation. Technically, a domain serves to separate the flow of control and data between components from the actual functionality of individual components. Besides facilitating hierarchical models, this separation can dramatically increase the reusability of components and models. The remainder of this book shows how to build Ptolemy II models and how to leverage the properties of each of the models of computation.