Modeling Concurrency in Parallel Debugging

Wenwey Hseush Gail E. Kaiser Columbia University Department of Computer Science New York, NY 10027 (212) 854-8123 hseush@cs.columbia.edu

Abstract

We propose a debugging language, Data Path Expressions (DPEs), for modeling the behavior of parallel programs. The debugging paradigm is for the programmer to describe the expected program behavior and for the debugger to compare the actual program behavior during execution to detect program errors. We classify DPEs into five subclasses according to syntactic criteria, and characterize their semantics in terms of a hierarchy of extended Petri Net models. The characterization demonstrates the power of DPEs for modeling parallelism. We present predecessor automata as a mechanism for implementing the third subclass of DPEs, which expresses bounded parallelism. Predecessor automata extend finite state automata to provide efficient event recognizers for parallel debugging. We briefly describe the application of DPEs to race conditions, deadlock and starvation.

keywords: debugging, formal models, Petri nets, path expressions, synchronization

1. Introduction

We are concerned with debugging parallel programs. One approach to locating the causes of program misbehavior is for the programmer to provide a high-level description of the expected behavior and for the debugger to compare the expected and actual behavior during execution. Expected behavior is specified abstractly in terms of control flow, data flow and/or synchronization events. In this approach, defining an appropriate notation for modeling program behaviors is a crucial prerequisite to developing a debugger. The conventional debugging approach, exemplified by dbx [Linton 81], also models program behavior but at a lower-level, in terms of source code entities such as subroutine names and line numbers; the programmer is responsible for comparing expected with actual behavior during execution. We refer to the first approach as problem-oriented, and the second as program-oriented. Both are necessary in practical debugging, just as both specification-based testing and program-based testing are required for practical testing [Howden 87].

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© 1990 ACM 089791-350-7/90/0003/0011 \$1.50

We have developed a style of problem-oriented debugging for parallel programs called *data path debugging*. Program behaviors are described in a formal notation called *Data Path Expressions* (DPEs), an extension of Bruegge and Hibbard's generalized path expressions for debugging sequential programs [Bruegge 85, Bruegge 83]. This work is in turn an application to debugging of Campbell and Habermann's classical work on path expressions for describing process behavior in operating systems [Campbell 74]. Other researchers also advocate a problem-oriented approach to parallel debugging (e.g., Bates [Bates 88a], Miller and Choi [Miller 88]). The primary advantage of our DPEs is that they model *true concurrency*.

Subclass	Semantic Model	
Sequential DPEs	Finite State Automata	
Multiple DPEs K-Safe Nets (subset)		
Safe DPEs	K-Safe Nets	
General DPEs Petri Nets		
Extended DPEs	tended DPEs Extended Petri Nets (subset)	

Figure 1-1: DPE Hierarchy

In our previous paper [Hseush 88], we informally described preliminary work on DPEs and discussed how they could be used in debugging parallel programs. The goal of this paper is to formally define several subclasses of DPEs in terms of their syntax and semantics. We define five subclasses according to syntactic criteria, and characterize the semantics of each subclass using a hierarchy of extended Petri net models [Peterson 81] (see Fig. 1-1). Extended Petri nets are equivalent to Turing machines [Hack 75, Thomas 76]. The first subclass expresses only sequential behavior. The second subclass expresses limited concurrency, in which process splitting (e.g., fork or para-do) is not permitted following a program branch (e.g, if-then-else), while program branching is permitted following a process split. The third expresses general bounded parallelism. The fourth permits unbounded parallelism, but without the ability to join an unknown number of threads. The fifth subclass describes general concurrency.

We propose *predecessor automata* as an implementation vehicle for the third subclass, safe DPEs, which subsumes the first and second subclasses. Predecessor automata extend finite state automata to represent predecessor events, and thus can recognize or generate partial ordering graphs [Lamport 78] as well as strings. The concurrent composition [Milner 80] of two predecessor automata preserves causal independence (*i.e.*, true concurrency), while the concurrent composition of two finite state automata loses this information. The expected program behavior described by a programmer as DPEs is translated into a predecessor automaton for efficiently recognizing concurrent events during execution. The DPE debugger will be useful for parallel applications where race conditions, deadlocks and starvation are concerned, and several small examples are given in this paper. We are in the process of implementing safe DPEs as a debugging language for both a concurrent extension of C and for the Meld concurrent object-oriented programming language [Kaiser 89].

Section 2 introduces DPEs and explains other background material necessary to understand the remainder of the paper. In section 3, we show the power of the five subclasses of DPEs in terms of extended Petri Net models. Section 4 presents the predecessor automata model for efficient implementation of safe DPEs in a debugging system. Section 5 discusses the practicality of DPEs for parallel debugging.

2. Background

A DPE consists of up to three components: one or more events, zero or more relations among events and zero or more actions. Events and the relations among them specify the behavior of program execution, while actions are performed by the debugger on program or debugger variables (or input/output) when the particular behavior is recognized during execution. A set of operators like sequencing (;), exclusive selection (+) and repetition (*) express the basic relationships among events.

2.1. Events and Actions

There are four kinds of events: control, data compound and conditional. Control events represent control activities, such as procedure entry and exit. Since there is a simple mapping from program execution to source code, control events can be specified using the corresponding identifiers in the program's source code, notably procedure names. Function.enter is the entering to Function, and Function.exit is the exiting from Function.

Data events occur when the specified program states become true. Data events are denoted as "[condition]", where the condition is an expression in the target programming language, augmented with the ability to express the history of program states and activities associated with data such as read and write. For example, "[X = 0]" is the event that variable X becomes equal to zero, " [X = X' + 1] " is the event that X is incremented (X' refers to the previous value of X). Data events are not associated with any particular control thread when defined, even though they are eventually caused by specific control events during execution. The programmer need only specify the effects on program entities without the knowledge of which control activities cause them; the debugger detects when the effects occur and reports which control activities cause them. It is difficult to efficiently recognize data events without either hardware support or significant modifications to the compiler and/or run-time support of the programming language, but we do not address this here.

Compound events are data path expressions. A compound event is defined by associating an identifier with a DPEs, permitting new DPEs to be defined in terms of the identifier in the style of context-free grammars, except that recursive and empty event definitions are not allowed. The format is "event-id = dpe", where dpe is a data path expression as defined in the next section.

Conditional events are control, data or compound events with predicates attached. The format is "event [condition]", where condition is a predicate. Conditional events are recognized when the event occurs in a context where the condition is already satisfied. For example, "READ [lock = 1]" is the condition where the READ procedure is called while lock is equal to one.

Actions may be attached to events. The format is "event { statements }", where statements is treated as a single action. The action is evaluated when the event is recognized. For example,

"[X'] { counter = counter + 1; }" means that every time X is updated, counter is incremented by one. Statements may involve program variables and/or debugger variables or functions, such as input/output and break.

2.2. Safe Concurrency

The term safe concurrency refers to the case of bounded parallelism and unsafe concurrency to unbounded parallelism. Language constructs designed for expressing concurrency (e.g., fork-join) often permit unsafe concurrency. Examples of safe and unsafe programs are shown in Figs. 2-1 and 2-2. The semantics of safe concurrency is characterized as a subclass of Petri nets, k-safe nets [Peterson 81], where the maximum number of tokens in a place is bounded by k. A k-safe net assures bounded parallelism. Every program with safe concurrency can be represented by a k-safe net, and every k-safe net is equivalent to a program with safe concurrency. The corresponding k-safe and unsafe nets for the safe and unsafe programs are also shown in Figs. 2-1 and 2-2.



Figure 2-1: A safe program



Figure 2-2: An unsafe program

3. DPE Hierarchy

DPEs are classified into five subclasses by the operators employed and some other syntactic restrictions. Each subclass is also defined in terms of a semantic model and the corresponding programming domain. The first subclass is sequential DPEs, which expresses only sequential behavior. The second is multiple DPEs, which subsumes the first and expresses limited safe concurrency, in which process splitting (e.g, fork or para-do) is not permitted following a program branch (e.g, if-then-else), while program branching is permitted following a process split. The third subclass is safe DPEs, which expresses safe concurrency. The fourth subclass is general DPEs, which expresses limited unsafe concurrency, where unbounded parallel threads never join. The fifth subclass is extended DPEs, which subsumes the fourth one and expresses unsafe concurrency.

DPEs employ five operators: sequencing (;), selection (+), repetition(*), concurrency (&) and concurrent closure (@). Examples are shown in Table 3-1. Table 3-2 summarizes the DPE hierarchy, including equivalence proofs and related work.

Expression	Description	
A;B	A causally precedes B.	
A+B	Either A or B occurs, but not both.	
A*	ε + A + A;A + A;A;A +	
A&B	A and B occur causally independently.	
A@	ε + A + A&A + A&A&A +	

Table 3-1: Operators

The first subclass is well known as regular expressions or path expressions. The path expression

"open ; (write + read) * ; close" states that a file has to be opened, before an arbitrary sequence of reads and writes is performed, and then closed.

Multiple DPE	Partial Order
(a;b)&(c;d) a;b and c:d occur	a 🌑 — 🕨 b
causally independently.	c 🚱 — — d
(a;s;b)&(c;s;d) a;s;b and c:s:d	a s b
synchronize at s.	c d
(a;s;b)&(c;s^;d) a;s;b and c:s^:d	a ⊕→● ^s →● b
occur causally independently.	$c \bullet \bullet \bullet \bullet^{s^{\wedge}} \bullet \bullet d$

Table 3-3: Multiple DPEs

Subclass	Sequential DPEs	Multiple DPEs	Safe DPEs	General DPEs	Extended DPEs
Expresses	sequential behavior	limited safe concurrency	safe concurrency	limited unsafe concurrency	unsafe concurrency
Syntax	dpe ₁ : EVENT (dpe ₁) dpe ₁ + dpe ₁ dpe ₁ ; dpe ₁ dpe ₁ *	dpe ₂ : dpe ₁ dpe ₂ & dpe ₁	dpe ₃ : EVENT (dpe ₃) dpe ₃ + dpe ₃ dpe ₃ ; dpe ₃ dpe ₃ * dpe ₃ & dpe ₃	dpe ₄ : dpe ₃ dpe ₄ + dpe ₄ dpe ₄ & dpe ₄ dpe ₄ @	$dpe_{5}: EVENT$ $ (dpe_{5}) $ $ dpe_{5} + dpe_{5} $ $ dpe_{5}; dpe_{5} $ $ dpe_{5} *$ $ dpe_{5} & dpe_{5} $ $ dpe_{5} & @$
Semantic Model	finite state automata	k-safe nets subset	k-safe nets	Petri nets [Garg 88]	extended Petri nets subset
Proof	Hopcroft & Ullman [Hopcroft 79]	Lauer & Campbell [Lauer 75]	Hseush & Kaiser [Hseush 89]	Garg [Garg 88]	Hseush & Kaiser [Hseush 89]
Example	open;(read+write)*; close	(a;s;b) & (c;s;d)	(enq;deq)+(enq&deq)	(fork;parent)* & (fork;child)@	(enq;update)@;deq; display
Limitations	no concurrency	no process splitting following program branching	no unbounded parallelism	no joining for unbounded parallelism	open question
Related Work	Generalized path expressions [Bruegge 85]	COSY [Lauer 81]	EBBA [Bates 83]	Concurrent regular expressions [Garg 88]	

 Table 3-2:
 DPE Hierarchy

The second subclass, multiple DPEs, expresses only global-level concurrency, where no nested concurrency (&) is allowed. Three examples are shown in Table 3-3. When the same event name appears in multiple subparts of the DPE, it is treated as a synchronization event and renaming is necessary to avoid this synchronization convention. We use ($^{\text{A}}$) to distinguish two distinct events with the same name (see the third example). In concurrent programming, a synchronization event usually involves two events in different threads, as explained in subsection 3.1.

The third subclass, safe DPEs, allows multi-level concurrency. One example that can be expressed by safe DPEs but not multiple DPEs is "enq ; deq + (enq & deq)", which states that if the queue size is equal to zero, then enqueuing must precede dequeuing; otherwise, enqueuing and dequeuing can operate concurrently. Safe DPEs are equivalent to k-safe nets.

The fourth subclass, general DPEs, expresses unbounded parallelism by using the concurrent closure operator (@), but disallows an event causally succeeding an unknown (unbounded) number of concurrent events. The general DPE "(fork;parent)* & (fork;child)@" models the unsafe program and the corresponding unsafe net mentioned in Fig. 2-2. Some programming examples are: (1) an unbounded number of messages may arrive at an object and each message activates a control thread for executing the same function without waiting for the prior activations to finish (e.g, process servers); and (2) an unknown number of signals arise and each signal invokes an unmasked signal handling routine. General DPEs are also known as concurrent regular expressions. Concurrent regular expressions have been proved equivalent to Petri nets [Garg 88]. The limitations on general DPEs are the same as those on Petri nets: no zero testing [Keller 72]. Zero testing is the ability to test for zero tokens in an unbounded place of a Petri net. For example, "(A;B)@;C" is not expressible in general DPEs or Petri nets. It states that an unknown (unbounded) number of threads A; B are created, and when all B events in the concurrent threads are complete, then C occurs; note that C can occur while the number of non-processed B's is tested equal to zero, because of the concurrent closure operator. This expression cannot be described by a Petri net, but is expressible by an extended Petri net [Peterson 81].

The fifth subclass of DPEs, extended DPEs, allows an event causally succeeding an unknown (unbounded) number of concurrent events, as modeled by extended Petri nets. For example, "(enq ; update)@ ; deq ; display" represents the case where an unbounded number of signals arise and each signal invokes a signal handling routine without disabling further signal invocations. The handling routine puts one character into a global queue and updates some information (the enqueue operation is atomic). After the control eventually returns to the main program, further signal invocation is disabled and all characters are dequeued and displayed. Extended DPEs express extended Petri nets, but whether extended DPEs are equivalent to extended Petri nets is an open question.

3.1. Synchronization Events

The behavior of language-specific synchronization primitives can be described using DPEs. Systems programmers or debugger users instruct the debugger to recognize the event patterns that constitute synchronizations among threads. For example, the pattern of sending a message X followed by receiving a message X constitutes a synchronization between the sender and the receiver. The description is

send(M); receive(M) { sync_event(\$1, \$2); }

which instructs the debugger that send is causally related to receive by message M. Then a synchronization event from the send event (\$1) to the receive event (\$2) can be established by the debugger based on the information from the sender and the receiver, once both are recognized during execution. Otherwise, the debugger would have no knowledge that send and receive matched as a synchronization event. Another example,

V(X).exit; P(X).exit { sync_event(\$1, \$2); }

, where P and V are the basic semaphore operations, instructs the debugger that V.exit is causally related to P.exit by the shared datum X and constitutes a synchronization event. Say the set of events is $P_1.enter$, $P_1.exit$, $P_2.enter$, $V_1.enter$, $V_1.exit$, $P_2.exit$, $V_2.enter$, $V_2.enter$, $V_2.exit$; the debugger uses the synchronization directive to establish the synchronization event ($V_1.exit$, $P_2.exit$).

This approach requires the same knowledge as in other approaches, but it provides the flexibility that users can easily invent and debug new synchronization primitives. In contrast, other debugging systems (e.g., [Goldszmidt 89]) retrieve such information through either source-to-source program transformation or augmenting the compiler with particular knowledge about synchronization primitives as related to parallel debugging.

4. Predecessor Automata

The problem-oriented debugging paradigm assumes that the programmer provides a description of expected program behavior, and the debugger compares this description to actual behavior at run-time to detect discrepancies. In our case, the debugger must be able to recognize sets of concurrent events matching DPEs. The debugger itself consists of support added (in hardware or software) to each executing thread or processor that submits messages representing primitive events (*i.e.*, control or data events) to a centralized DPE recognition process. The sequence of events it receives are treated as a string of tokens and are compared with the user-specified DPEs. Since the recognition of primitive events has been addressed in traditional debugging, where breakpointing is typically adopted, we are here concerned with the central recognition process, which is to recognize the program behavior specified in DPEs.

One key issue is the tradeoff between the efficiency of recognition and the memory space needed to represent the DPEs in a suitable internal form. In the case where minimizing memory space is most important, Petri nets are probably the best choice. Petri nets can represent sequential and concurrent behavior in a compact form, but they are relatively inefficient for recognizing events at runtime. In contrast, finite state automata (FSAs) are efficient recognizers for sequential behavior, but they cannot represent concurrent events that are causally independent. FSAs express interleaving semantics, but not true concurrency. The concurrent composition of two FSAs involves combining two FSAs into one such that all possible states and all possible interleavings of two sets of transitions are preserved [Milner 80]. This process loses the information regarding which events occur causally independently and there is no way to reverse the process to recover the original two FSAs. With or without concurrent composition, FSAs cannot distinguish two causally independent events interleaved with each other from two sequential events.

We present an implementation model, *predecessor automata* (PAs), that has the clean and efficient structure of FSAs, but also

the capability of representing true concurrency as in safe Petri nets. PAs can recognize behavior with safe concurrency and possibly detect the situation of unsafe concurrency, and thus implement our third subclass, safe DPEs.

4.1. Definition of Predecessor Automata

- A predecessor automaton is a 5-tuple $(Q, \Sigma, \delta, q_0, F)$.
 - Q is a finite set of states.
 - Σ is a finite set of *events*.
 - δ is the transition function mapping $Q \times \Sigma \times P$ to Q, where P is the predecessor set, $P \subseteq \Sigma^*$.
 - q_0 is the initial state, $q_0 \in Q$.
 - F is the set of final states, $F \subseteq Q$.

The definition of a PA is the same as an FSA except for the transition function, which not only carries the information about the expected events, but also the information about their predecessors.

The predecessor p ($\in P$) of an event e is a list of events (w_1, w_2, \ldots, w_n) , where (1) *n* is a non-negative integer, (2) event w_i causally precedes e and (3) w_i and w_i occur causally independently, $1 \le i, j \le n, i \ne j$. If $p = \varepsilon$, e is an original event (ε is also represented as '.'). The occurrence of event e implies that all its predecessor events $w_i \in p, 1 \le i \le n$, have occurred. The input to a PA is not a string of events, but a string of event-predecessor pairs, $(e_0 \ p_0)$, $(e_1 \ p_1)$, ..., $(e_n \ p_n)$, where $e_i \in \Sigma$, $p_i \in P$, $0 \le i \le n$ and for $w \in p_i$, there exists e_j , $j \ne i$, $0 \le j \le n$, such that $e_i = w$. That is, every event mentioned as a predecessor event must occur. The string of event-predecessor pairs can be considered as a partial ordering graph, which represents a program execution. The vertics in the partial ordering graph are events and the directed edges indicate the relations of causal dependency (i.e., sequencing (;)). Remember that no relationships like selection (+) and repetition (*) can be represented in the partial ordering graph. In the case that j < i, it means the receiving order preserves the occurrence (partial) order. This is discussed in more detail later in this section.



Figure 4-1: An example of PA

A PA moves from one state q to another state r on an input (e p), according to the transition function $\delta(q, (e p)) = r$. That is, a move is made by examining the incoming event and its predecessors. For example, a predecessor automaton is shown in Fig.4-1, where the transition function is represented by a set of labeled edges. Transition (a ϵ) (or (a .)) indicates the state transition of the occurrence of an original event a (without predecessors), and (d (c b)) (or (d c b)) indicates the occurrence of event d with two predecessor events c and b. One possible input that will be accpeted by the automaton is (a .) (b .) (c a) (d c b). The state of the automaton will change from 1 to 2 with (a .), from 2 to 5 with (b .), from 5 to 6 with (c a), and from 6 to the final state with (d c b). The input represents an execution of a concurrent program, and its partial ordering graph is also shown in Fig.4-1.

In the case where generators are concerned, a partial ordering graph can be constructed by giving a *PA* and a sequence of transitions, $\delta(q_0, (e_0 p_0)) = q_1$, $\delta(q_1, (e_1 p_1)) = q_2$, ..., $\delta(q_{n-1}, (e_{n-1} p_{n-1})) = q_n$, where q_0 is the initial state. For each $(e_i p_i)$, $0 \le i \le n$, (1) create a vertex labeled with e_i and (2) for each event $w \in p_i$, create a directed edge from w to e_i .

Two problems arise when constructing partial ordering graphs: (1) ambiguity and (2) instability. A PA is *ambiguous* if and only if there exists a sequence of transitions, $\delta(q_0, (e_0 p_0)) = q_1, \delta(q_1, (e_1 p_1)) = q_2, ..., \delta(q_{n-1}, (e_{n-1} p_{n-1})) = q_n$, where q_0 is the initial state, such that more than one partial ordering graph can be constructed. Fig. 4-2 illustrates an ambiguous situation.



Figure 4-2: An ambiguous situation



Figure 4-3: An unstable situation

The first step to eliminate the ambiguous situation is to rename some events e_i in the ambiguous PA, such that there does not exist an event e_j , $i \neq j$ and $e_i = e_j$. The second step is to modify the graph construction procedure to eliminate the ambiguous situations due to cyclic paths of transitions in PAs. When a directed edge $w \rightarrow e$ is constructed, the vertex labeled w is the one that was added to the graph most recently and labeled with w. These two steps eliminate all possible ambiguous situations.

Given a PA and a sequence of inputs, $(e_0 p_0)$, $(e_1 p_1)$, ..., $(e_n p_n)$, a situation is *unstable* at input $(e_i p_i)$ if and only if there exists a predecessor event $w \in p_i$, such that $e_i \neq w$ for all $j, 0 \leq j \leq i-1$. Informally, a situation is unstable if an event mentioned as a predecessor has not arrived so far or the event (vertex) is missing

in the constructed graph. One example is shown in Fig. 4-3.

4.2. Event Recognition

Recognition involves two components: (1) a target system that reflects the actual program behavior and provides the information about primitive events and their predecessors, and (2) a recognizer that represents the expected program behavior in some internal form and collects and processes the information. The target system is a concurrent system, with messages representing primitive events and their predecessors generated from different processors and sent to the centralized recognizer.

The recognizer is a sequential machine that receives messages representing primitive events from different threads one message at a time, compares them with the expected ones, and eventually reports the results. The message receiving order is assumed independent from the order of the event occurrences, since the sending order may be different from the receiving order. The recognizer has two parts, a stabilizer and a PA. The stabilizer has two functions: (1) filtering incoming event/predecessors messages such that only the "interesting" events (*i.e.*, the primitive events mentioned by users in DPEs and the synchronization events) and their similarly "interesting" predecessor events go into the automaton, and (2) regulating the incoming event/predecessors messages such that the ordering of event messages that go into the automaton preserves the partial ordering of event occurrences in the target system. For example, if the input messages to the stabilizer are $(e_1 \ \epsilon) \ (e_2 \ e_3)$ $(e_3 \ e_1) \ (e_4 \ e_1) \ (e_5 \ e_4)$, where all events are "interesting" events except e_4 . The output of the stabilizer is $(e_1 \ \epsilon) \ (e_3 \ e_1) \ (e_2 \ e_3) \ (e_5 \ e_1)$. The output messages of the stabilizer are the input messages of the PA, which will compare the input messages (the actual behavior) with the DPEs (the expected behavior) provided by the users as represented by a PA. A general structure for such a debugging system is shown in Fig. 4-4.



Figure 4-4: Event recognizer

The PA is in the initial state before receiving any messages. Every time a message describing an event and its predecessors arrives from the stabilizer, the PA compares the received information with the transitions directed from the current state. If both the event and its predecessors match one of the transitions, the automaton moves to the next state according to the matched transition. An example is illustrated in Figs. 4-5 and 4-6. One important assumption in our event recognition framework is that the target system (eventually) has full knowledge about every event that occurs and its predecessor events, where these events appear in some DPE used to construct the PA and/or reflect synchronization events.



Figure 4-5: An event recognizer for a;(b&(c;d));e

Targeted system	stabilizer	РА
events	events	state transitions
(a .)	(a .)	0>1
(b a)	(b a)	1>2
(d c)		
(c a)	(c a)	2>4
	(d c)	4>6
(e b d)	(e b d)	6>7

Figure 4-6: PA Description

4.3. Constructing Predecessor Automata From Safe DPEs

Given a safe DPE, a predecessor automaton can be constructed. There are two steps, involving transformations of subexpressions and translation using an attribute grammar [Knuth 68]. The first step is to transform each expression into a new expression where there are no subexpressions \mathbb{R}^* , such that $\varepsilon \in \mathbb{R}$. For example, $(e^*)^*$ can be transformed into e^* . This guarantees that the constructed automaton has no transition cycles $\delta(q_1, (e_1 p_1)) = q_2$, $\delta(q_2, (e_2 p_2)) = q_3, ..., \delta(q_n, (e_n p_n)) = q_1$, such that $e_i = \varepsilon$, for all *i*, $0 \le i \le n$. The transformation is based on an extension to Foster's conversion theorem [Foster 86]. For any DPE **R**, there is a DPE **N**(**R**) such that (1) **N**(**R**) does not contain the empty string, and (2) $\mathbf{R}^* = (\mathbf{N}(\mathbf{R}))^*$. If $\varepsilon \subseteq \mathbf{R}$, $\mathbf{N}(\mathbf{R}) = \mathbf{R}$. Otherwise, there are four cases.

- 1. If $\mathbf{R} = \mathbf{P}^*$, $\mathbf{N}(\mathbf{R}) = \mathbf{N}(\mathbf{P})$
- 2. If R = P+Q, N(R) = N(P) + N(Q)
- 3. If $\mathbf{R} = \mathbf{P}; \mathbf{Q}, \mathbf{N}(\mathbf{R}) = \mathbf{N}(\mathbf{P}) + \mathbf{N}(\mathbf{Q})$
- 4. If R = P & Q, N(R) = (N(Q) & N(P)) + N(Q) + N(P)

The second step applies an attribute grammar that specifies how to construct a PA. A DPE is first parsed into an abstract syntax tree, where three attributes are attached to each node of the tree, **AUTO**, **PRED** and **LAST**. The **AUTO** attribute of a node *n* will contain an automaton that represents the subtree (subexpression) rooted at node *n*. A subtree can be considered as a subexpression or a PA. The **PRED** attribute of *n* represents its predecessors, the events that might precede any event occurring in the subtree rooted with *n*. The **LAST** attribute of *n* refers to the events without successors in the subtree. The values of **PRED** and **LAST** have the form $(e_{0,0} \land \dots \land e_{0,m}) \lor \dots \lor (e_{n,0} \land \dots \land e_{n,m})$, where the events related with (\land) occur concurrently and the events related with (\checkmark) occur exclusively. The semantic rules associated with the grammar are shown in Fig. 4-7.

This is not a syntax-directed translation system like YACC [Johnson 78]. Instead, the semantic rules describe the relations between a node in the abstract syntax tree and its parent node, and

between the node and its children nodes. A semantic rule is evaluated only when its dependent attribute(s) is changed [Reps 84]. instead of at the time of parsing. For example, the first semantic rule,

"dpe.AUTO = new_PA (EVENT, dpe.PRED)", which is associated with a leaf node, is evaluated when its **PRED** attribute is changed.

```
dpe : EVENT
```

```
dpe.AUTO = new PA(EVENT, dpe.PRED);
 dpe.LAST = last events(dpe.AUTO);
 '(' dpe<sup>1</sup> ')'
 dpe^{1}.PRED = dpe.PRED;
 dpe.AUTO = dpe^{1}.AUTO;
 dpe.LAST = dpe^{1}.LAST;
dpe<sup>1</sup> ';' dpe<sup>2</sup>
 dpe^{1}.PRED = dpe.PRED;
 dpe^2.PRED = dpe^1.LAST;
 dpe.AUTO = concat(dpe<sup>1</sup>.AUTO,dpe<sup>2</sup>.AUTO);
 dpe.LAST = dpe^2.LAST;
|dpe^1' + dpe^2
 dpe^1.PRED = dpe.PRED;
 dpe^2.PRED = dpe.PRED;
 dpe.AUTO = union(dpe<sup>1</sup>.AUTO,dpe<sup>2</sup>.AUTO);
 dpe.LAST = dpe<sup>1</sup>.LAST \vee dpe<sup>2</sup>.LAST;
 dpe<sup>1</sup> '*'
 dpe<sup>1</sup>.PRED = dpe.PRED;
 dpe.AUTO = repeat(dpe<sup>1</sup>.AUTO);
 dpe.LAST = dpe<sup>1</sup>.LAST \vee \varepsilon;
dpe^1 '&' dpe^2
 dpe^{1}.PRED = dpe.PRED;
 dpe^2.PRED = dpe.PRED;
 dpe.AUTO = compose(dpe^1.AUTO,dpe^2.AUTO);
 dpe.LAST = last events(dpe.AUTO);
}
```

Figure 4-7: Attribute Grammar for DPEs

The function last events, with a PA as an input parameter, obtains the last events that might occur in the PA. The return value has the same form as LAST and PRED. The function new PA creates a new automaton with two input parameters, an event e and its predecessors p. The new automaton has one start state p, one final state q and one transition $\delta(p (e \text{ PRED}(e))) = q$. The attribute grammar evaluation is started by setting the PRED attribute of the root to ε ; every node will eventually be visited a few times, as changes are propagated around the tree. The root is the first node visited, since its **PRED** is changed. For each node evisited, if e is a leaf, AUTO is assigned a new PA and LAST is set to e. Since the values of PRED and AUTO are changed, its parent node will be visited again according to the semantic rules associated with the parent. If the node is not a leaf, it propagates the value of PRED down to its child nodes, and when the node is eventually visited again, it constructs a new PA from its children's PAs according to the operators and properly sets the value of its LAST attribute. The functions concat is to concatenate two PAs, union is the union of two PAs, and

repeat is the Kleene closure of a PA. These fuctions are the same as those for FSAs. The function compose concurrently composes two PAs into one, as explained in the next section. When the evaluation is complete, the AUTO attribute of the root contains the PA for the given DPE.

4.4. Concurrent Composition

The concurrent composition of two PAs creates a new PA that preserves all possible states and all possible transitions as if the two original automata operate concurrently. As explained above, the concurrent composition of two finite state automata will lose the concurrency information, while the concurrent composition of two PAs will not. An example is shown in Fig. 4-8.



Figure 4-8: Concurrent composition of two FSAs



Figure 4-9: Concurrent composition of two PAs

Composition of two PAs can be divided into two cases, those that do and do not involve synchronization. Synchronization occurs when two automata have common events (or reflect components of synchronization events). Assume the first automaton has nstates, s_0 , s_1 , ..., s_{n-1} , s_0 is the initial state, and the second automaton has m states, z_0 , z_1 , ..., z_{m-1} , z_0 is the initial state. In the case that two PAs have no synchronization, the composed automaton will have $n \times m$ states, $q_{1,1}, q_{1,2}, \dots, q_{n,m-1}, q_{n,m}$. The state $q_{i,j}$ is the combined state of the state s_i in the first automaton and the state z_i in the second automaton. The transitions from $q_{i,j}$ to $q_{k,i}$ in the composed automaton are the transitions from s_i to $\tilde{s_k}$ in the first automaton, and the transitions from $q_{i,j}$ to $q_{i,k}$ in the composed automaton are the transitions from z_i to z_k in the second automaton. There exist no transition between $q_{i,j}$ and $q_{h,k}$, $i \neq h$ and $j \neq k$. In the case where two PAs do have synchronization, there must exist a transition $(e_i p_i)$ in the first automaton and transition $(e_i p_i)$ in the second automaton, such that $e_i = e_i$. Assume s_k and s_j are the pre-state and the post-state for transition $(e_i p_i)$ in the first automaton, and z_g and z_h for $(e_j p_j)$ in the second. The composed PA will have a combined state $q_{k,g}$ of s_k and z_g , a combined state $q_{1,h}$ of s_1 and z_h , and a transition

4.5. Related Work

EBBA [Bates 88b] employs *shuffle automata* [Bates 87] as a formal model for event recognition in distributed systems. Shuffle automata recognize concurrent events based on the interleaving semantics. That is, shuffle automata cannot distinguish two causally independent events interleaving with each other from two causally dependent events.

Shuffle automata are an FSA-like formalism that consist of a set of states and a finite state control that effects transitions from an initial state to some final state. An important difference between the shuffle automaton and an FSA is that in order to make transitions in the shuffle automaton, the finite state control examines sets of input symbols, rather than individual symbols. At run-time, the recognizer will accumulate the incoming events in a set. Whenever a subset of the accumulated event set becomes sufficient to make a transition, the finite control then goes from the current state to another state.

5. Debugging Concurrent Programs

Most concurrency-related bugs involve problems with synchronization among multiple threads, which may share information in a number of different ways, including shared memory, message passing, files and devices, and human interaction. In this section, we demonstrate that DPEs are useful for aiding detection and correction of three typical kinds of synchronization errors: race conditions, deadlocks and starvations.

A race condition happens when two or more concurrent threads interact with some common resources without properly constraining the ordering of interactions, resulting in a computation that is nondeterministic and incorrect. To eliminate the race conditions, appropriate synchronization must be added to the program so that the crucial interactions are properly ordered. Two types of synchronization mechanisms are frequently adopted: (1) wait-resume and (2) rollback-retry. Wait-resume constrains the ordering of interactions by blocking threads from competing for resources, but may lead to a deadlock situation when two or more threads wait for each other indefinitely due to lack of knowledge of the global situation. In the rollback-retry type of synchronization, a thread constrains the ordering of interactions by expecting other threads to complete their crucial interactions while temporarily releasing its resources. This may lead to a starvation situation where one or more threads repeats the rollback-retry cycle indefinitely. In the dining philosophers example, there is a deadlock when every philosopher has a fork in his right hand and is waiting for the fork on his left-hand side; there is starvation when a philosopher repeatedly picks up the forks on his right-hand side and then puts down the fork because the fork on his left-hand side is always unavailable.

It is difficult to debug programs with race conditions, deadlocks or starvations, where bugs may be embedded in (1) the synchronization primitives and/or (2) the program units that apply the synchronization primitives. It is also difficult for programmers to detect, by observing the external program behavior, whether the error is caused by buggy synchronization primitives or buggy program units. We assume in this paper that synchronization primitives are always correct, and are thus concerned only with (2). One concern in debugging is reproducibility, since it is desirable for the identical program behavior to be replayed by re-execution or simulation over and over again until the bugs are located. We assume this is possible, but do not address the mechanism here.

5.1. Debugging Race Conditions

There are two necessary conditions for race conditions: (1) concurrent threads share common resources, and (2) the particular events within these threads that compete for the common resources are causally independent. Therefore, debugging a program with race conditions can be treated as a process of establishing relations of causal dependence and detecting whether the critical events that access the common resources occur causally independently.

program producer_consumer; var s: semaphore := 1;deposited: semaphore := 0; procedure producer; var next: integer; begin while true do begin next = calculate(); P(s); -----> (1) enqueue(next); -----> (2) V(s); V(deposited); end; end; procedure consumer; var next: integer; begin while true do begin P(deposited); P(s);next = dequeue(); V(s);print(next); end; end; begin para-do producer(); consumer(); para-end end

Figure 5-1: Producer-Consumer Program

For example, Fig. 5-1 shows a producer-consumer program, where the producer thread puts numbers in a queue, and the consumer thread gets and prints the numbers from the queue when the queue is not empty. A semaphore s and its operations P(s) and V(s) are used for synchronization. Assume the P(s) at point (1) is missing from the program. During execution, the queue data structure may become inconsistent. In order to debug the program, the first step is to define, using DPEs, the synchronization events in the program (see section 3.1).

Then, in the case where a race condition between producer and consumer is suspected, the second step is to describe, in DPEs,

the expected misbehavior that enqueue and dequeue occur concurrently. The expression

"enqueue() & dequeue() { print(s); break; }" instructs the debugger to print the value of semaphore s and stop the execution when enqueue and dequeue occur concurrently. The third step is to replay the program execution. The program execution will stop at (2) and the value of s is printed out. The debugger will detect the true concurrency of enqueue and dequeue, no matter how the event messages interleave with each other. Some interleavings might accidently produce correct results and others produce the wrong results; in both cases, the debugger will detect the race condition.

5.2. Debugging Deadlocks

There are four necessary conditions for deadlock [Coffman 71]: (1) Threads claim exclusive control of the resources they require (mutual exclusion condition), (2) Threads hold resources already allocated to them while waiting for additional resources (wait for condition), (3) Resources cannot be removed from the threads holding them until completion (no preemption condition), and (4) A circular chain of threads exists in which each holds one or more resources that are requested by the next thread in the chain (circular wait condition). Debugging a program with deadlock requires the same description of synchronization events as in debugging a program with race conditions, but has a more complicated expected program behavior.

One example is that lock and unlock are used to allocate resources before reference to the data. The first three conditions are determined by the synchronization primitives, and the fourth condition can be established by constructing a wait-for graph during debugging. The synchronization events can be described as "unlock(X).exit; lock(X).exit { sync_event(\$1, \$2); }". The expected program behavior can then be described as "lock(X).exit {hold(\$1.pid, X)}; unlock(X).exit {unhold(\$2.pid, X)}" and "lock(X).enter; wait() {wait_for(\$1.pid,X); check_deadlock()}; resume(); lock(X).exit {release(\$4.pid,X)}", where (1) the hold() function informs the debugger that the thread of the event (\$1.pid) holds the resource X, (2) unhold () tells the debugger that the associated thread (\$3.pid) does not hold the resource X any more, (3) wait for () means that the associated thread (\$1.pid) waits for resource X, (4) release() that the associated thread no longer waits for the resource X, and (5) check deadlock asks the debugger to check whether a deadlock exists according to the information provided by the first four functions.

5.3. Debugging Starvations

Starvation is a special type of race condition where a set of causally independent events might repeat indefinitely. In the example of dining philosophers, every philosopher might repeat picking up the fork on his left-hand side and putting it down. One possibility for detecting this is to store the program state every time a philosopher picks up his right fork and compare it with the previous states. If there exists an identical previous state and between them no progress has been made, there may (or may not) be an error. Detecting starvation is probably more amenable to program verification than debugging, but DPEs can check the correctness of verification assertions during execution.

6. Conclusions

We have defined a formal notation, DPEs, for modeling concurrent behavior in the context of debugging parallel programs. There are five subclasses of DPEs, four equivalent in power to a member of a hierarchy of Petri net models and the fifth a subset of extended Petri nets. We have developed an efficient implementation vehicle for the third subclass of DPEs, which models safe concurrency. We have briefly described the application of DPEs to practical concurrent debugging problems, from a viewpoint of problem-oriented behavior. DPEs must be combined with conventional debugging mechanisms to observe program-oriented behavior, for example, to support singlestepping among statements and modification of the program state at a breakpoint.

Acknowledgments

Hseush is supported in part by the Center for Telecommunications Research. Part of this research was conducted while Hseush was a summer employee of the IBM T.J. Watson Research Center. Kaiser is supported by National Science Foundation grants CCR-8858029 and CCR-8802741, by grants from IBM, AT&T, Siemens, Sun and Xerox, by the Center for Advanced Technology and by the Center for Telecommunications Some of the ideas presented here originated in Research. discussions with Timothy Balraj, who has been working on Petri net models [Balraj 86] and path expressions in the context of silicon compilers. We would also like to thank Janice Stone, Bowen Alpern, Felix Wu and Dannie Durand for their helpful discussion, and Colin Harrison for his support of DPE debugging. Krish Ponamgi is working with us on an implementation of DPEs for C on the IBM 8CE multiprocessor running the Mach operating system [Rashid 87]. Krish is a co-op MS student at IBM under the supervision of Colin Harrison. Yi-wun Lu and Taka Ishizuka have previously worked with us on the Meld Debugger (MD) implementation of DPEs on Sun 3 workstations.

References

[Balraj 86]	 T.S. Balraj and M.J. Foster. Miss Manners: A Specialized Silicon Compiler for Synchronizers. In Proceedings of the Fourth MIT Conference, pages 3-20. The MIT Press, April, 1986.
[Bates 83]	Peter Bates and Jack C. Wileden. An Approach to High-Level Debugging of Distributed System. In ACM SIGSoft/SIGPlan Software Engineering Symposium on High-Level Debugging, pages 107-111. Pacific Grove, CA, March, 1983. Special issue of Software Engineering Notes, 8(4), August 1983.
[Bates 87]	Peter C. Bates. Shuffle Automata: A Formal Model for Behavior Recognition in Distributed Systems. Technical Report COINS 87-27, University of Massachusetts at Amherst, January, 1987.
[Bates 88a]	 Peter Bates. Distributed Debugging Tools for Heterogeneous Distributed Systems. In 8th International Conference on Distributed Computing Systems, pages 308-315. Computer Society Press, San Jose CA, June, 1988.
[Bates 88b]	 Peter Bates. Distributed Debugging Tools for Heterogeneous Distributed Systems. In ACM SIGPLAN/SIGOps Workshop on Parallel and Distributed Debugging, pages 11-22. Madison WI, May, 1988. Special issue of SIGPLAN Notices, 24(1), January 1989.

[Bruegge 83]	 Bernd Bruegge and Peter Hibbard. Generalized Path Expressions: A High-Level Debugging Mechanism. The Journal of Systems and Software 2(3):265-276, 1983.
[Bruegge 85]	Bernd Bruegge. Adaptability and Portability of Symbolic Debuggers. PhD thesis, Carnegie Mellon University, 1985. CMU-CS-85-174.
[Campbell 74]	 R.H. Campbell and A.N. Habermann. The Specification of Process Synchronization by Path Expressions. In G. Goos and J. Hartmanis (editors), Lecture Notes in Computer Science. Volume 16: Operating Systems, pages 89-102. Springer-Verlag, Berlin, 1974.
[Coffman 71]	E. G. Coffman, Jr., M. Elphick; and A. Shoshani. Systems Deadlocks. <i>Computing Surveys</i> 3(2):71-76, June, 1971.
[Foster 86]	 M.J. Foster. Avoiding Latch Formation in Regular Language Recognizers. In Proceedings of the Allerton Conference on Communication, Control, and Computing, pages 740-748. University of Illinois, Urbana-Champaign IL, October, 1986. This paper also appears in the IEEE VLSI Technical Bulletin, 1(2), September 1986.
[Garg 88]	Vijay Kumar Garg. Specification and Analysis of Distributed Systems With a Large Number of Processes.
[Goldszmidt 89]	 German S. Goldszmidt, Shmuel Katz and Shaula Yemini. High Level Language Debugging for Concurrent Programs. Technical Report RC14341, IBM Research Division, T.J.Watson Research Center, Yorktown Heights, N.Y. 10598, January, 1989.
[Hack 75]	 M. Hack. Decidability Questions for Petri Nets. PhD thesis, Department of Electrical Engineering, Massachusetts Institute of Technology, 1975. Technical Report 161.
[Hopcroft 79]	John E. Hopcroft and Jeffrey D. Ullman. Introduction to Automata Theory, Languages and Computation. Addison-Wesley Publishing Company, 1979, pages 28-35.
[Howden 87]	William E. Howden. Software Engineering and Technology: Functional Program Testing & Analysis. McGraw-Hill Book Co., New York, 1987.
[Hseush 88]	 Wenwey Hseush and Gail E. Kaiser. Data Path Debugging: Data-Oriented Debugging for a Concurrent Programming Language. In ACM SIGPLAN/SIGOps Workshop on Parallel and Distributed Debugging, pages 236-246. Madison WI, May, 1988. Special issue of SIGPLAN Notices, 24(1), January 1989.
[Hseush 89]	 Wenwey Hseush and Gail E. KAISER. Modeling Concurrency in Parallel Debugging. Technical Report CUCS460, Computer Science Department, Columbia University, NY, NY, October, 1989.
[Johnson 78]	S.C. Johnson and M.E. Lesk. Language Development Tools. The Bell System Technical Journal 57(6):2155-2175, July-August, 1978.
[Kaiser 89]	 Gail E. Kaiser, Steven S. Popovich, Wenwey Hseush and Shyhtsun Felix Wu. Melding Multiple Granularities of Parallelism. In Stephen Cook (editor), 3rd European Conference on Object- Oriented Programming, pages 147-166. Cambridge University Press, Nottingham, UK, July, 1989.
[Keller 72]	 R. Keller. Vector Replacement Systems: A Formalism for Modeling Asynchronous Systems. Technical Report 117, Computer Science Laboratory, Princeton University, December, 1972.

[K1	nuth 68]	Donald E. Knuth. Semantics of Context-Free Languages. Mathematical Systems Theory 2(2):127-145, June, 1968.
[La	import 78]	Leslie Lamport. Time, Clocks and the Ordering of Events in a Distributed System. CACM 21(7):558-564, July, 1978.
[La	wer 75]	 P. E. Lauer and R. H. Campbell. Formal Semantics of a Class of High-Level Primitives for Coordinating Concurrent Processes. Acta Informatica 5(4):297-332, 1975.
[La	wer 81}	 P. E. Lauer and M. W. Shields. Formal behavioural specification of concurrent systems without globality assumptions. In J. Diaz and I.Ramos (editor), Lecture Notes in Computer Science. Number 107: Proceedings of Internation Colloquium on Formalization of Programming Concepts, pages 115-151. Springer-Verlag, Berlin, 1981.
[Li	nton 81]	M. Linton. A Debugger for the Berkeley Pascal System. Master's thesis, University of California at Berkeley, June, 1981.
[M	iller 88]	 Barton P. Miller and Jong-Deok Choi. Breakpoints and Halting in Distributed Programs. In 8th International Conference on Distributed Computing Systems, pages 316-323. Computer Society Press, San Jose CA, June, 1988.
(M	ilner 80]	 Robin Milner. A Calculus of Communicating Systems. In G. Goos and J. Hartmanis (editors), Lecture Notes in Computer Science (92). Springer-Verlag, Berlin, 1980.
[Pe	eterson 81]	James L. Peterson. Petri Net Theory and The Modeling of Systems. Prentice-Hall, Inc., Englewood Cliffs, NJ 07632, 1981.
[Ra	ashid 87]	 Richard Rashid, Avadis Tevanian, Michael Young, David Golub, Robert Baron, David Black, William Bolosky and Jonathan Chew. Machine-Independent Virtual Memory Management for Paged Uniprocessor and Multiprocessor Architectures. In 2nd International Conference on Architectural Support for Programming Languages and Operating Systems, pages 31-39. Palo Alto CA, October, 1987. Special issue of SIGPIan Notices, 22(10), October 1987.
[Re	eps 84]	Thomas Reps. Generating Language-Based Environments. The MIT Press, Cambridge MA, 1984.
[Tז	10mas 76]	 P. Thomas. The Petri Net: A Modeling Tool for the Coordination of Asynchronous processes. Master's thesis, University of Tennessee, 1976.