

Fault-Tolerant Computer System Design

ECE 695/CS 590

Topic 7: Modeling

Saurabh Bagchi

ECE/CS

Purdue University

Lecture by Paul Wood (pwood@purdue.edu)

ECE 695/CS 590

1



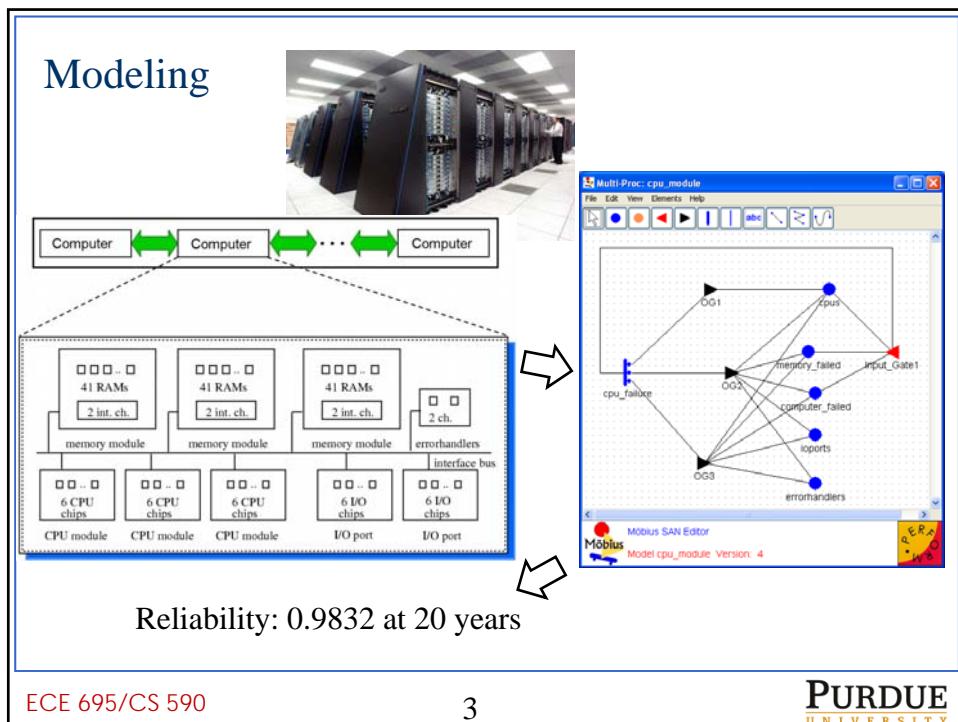
Logistics

- Use Möbius from Dr. Bill Sanders' group at UIUC
- There will be a programming assignment on this
 - ~~Nov. 8 assigned, due Nov. 21~~
- Request account from <https://www.mobius.illinois.edu/>
 - Mention Prof. Bagchi and that you are registered students for this class
 - Distribution available for Windows, OS X, and Ubuntu Linux
- Takes 2-3 days for approval
- We have acquired an academic Möbius license
- Video lectures exist on this topic

ECE 695/CS 590

2





ECE 695/CS 590

3

PURDUE
UNIVERSITY

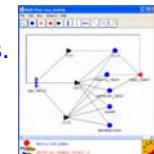
Stochastic Activity Networks

Stochastic activity networks, or SANs, are a convenient, graphical, high-level language for describing system behavior. SANs are useful in capturing the stochastic (or random) behavior of a system.

Examples:

- The amount of time a program takes to execute can be computed precisely if all factors are known, but this is nearly impossible and sometimes useless. At a more abstract level, we can approximate the running time by a random variable.
- Fault arrivals almost always must be modeled by a random process.

We begin by describing a subset of SANs: stochastic Petri nets.



ECE 695/CS 590

4

PURDUE
UNIVERSITY

Stochastic Petri Net Overview

One of the simplest high-level modeling formalisms is called *stochastic Petri nets*. A stochastic Petri net is composed of the following components:

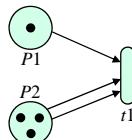
- Places: which contain tokens, and are like variables
 **NumCPU**
- tokens: which are the “value” or “state” of a place
 **NumCPU=3**
- transitions: which change the number of tokens in places
 **Timed** **Untimed**
every 30 seconds
- input arcs: which connect places to transitions
 **CPU consumed every 30 seconds**
- output arcs: which connect transitions to places
 **CPU is placed in busy place**

Firing Rules for SPNs

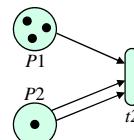
A stochastic Petri net (SPN) executes according to the following rules:

- A transition is said to be *enabled* if for each place connected by input arcs, the number of tokens in the place is \geq the number of input arcs connecting the place and the transition.

Example:



Transition t_1 is enabled.

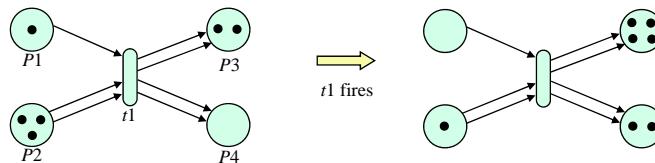


Transition t_2 is not enabled.

Firing Rules, cont.

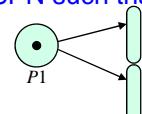
- A transition may *fire* if it is enabled. (More about this later.)
- If a transition fires, for each input arc, a token is removed from the corresponding place, and for each output arc, a token is added to the corresponding place.

Example:



Specification of Stochastic Behavior of an SPN

- A stochastic Petri net is made from a Petri net by
 - Assigning an exponentially distributed time to all transitions.
 - Time represents the “delay” between enabling and firing of a transition.
 - Transitions “execute” in parallel with independent delay distributions.
- Since the minimum of multiple independent exponentials is itself exponential, time between transition firings is exponential.
- If a transition t becomes enabled, and before t fires, some other transition fires and changes the state of the SPN such that t is no longer enabled, then t *aborts*, that is, t will not fire.
- Since the exponential distribution is memoryless, one can say that transitions that remain enabled continue or restart, as is convenient, without changing the behavior of the network.



SPN Example: Readers/Writers Problem

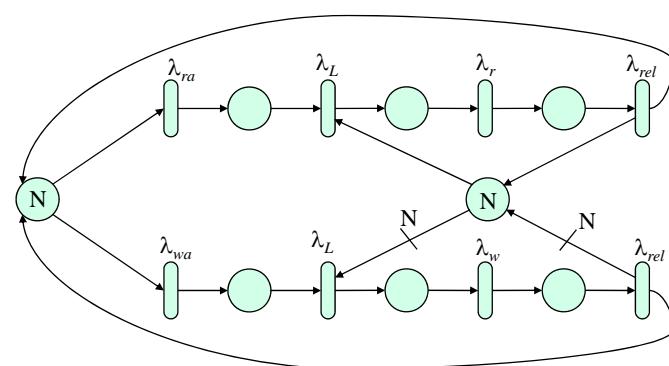
There are at most N requests in the system at a time. Read requests arrive at rate λ_{ra} and write requests at rate λ_{wa} . Any number of readers may read from a file at a time, but only one writer may write at a time. A reader and writer may not access the file at the same time.

Locks are obtained with rate λ_L (for both read and write locks); reads and writes are performed at rates λ_r and λ_w respectively. Locks are released at rate λ_{rel} .

Note:

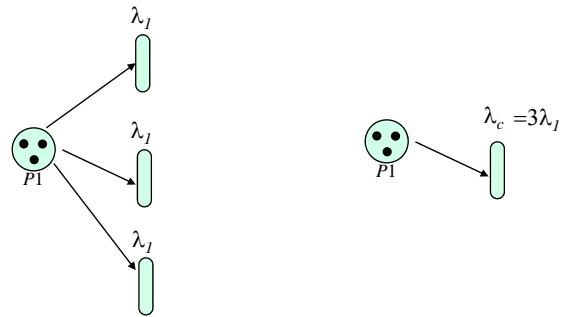


SPN Representation of Reader/Writers Problem



SPN Example: Expected Timings

- Since the minimum of multiple independent exponentials is itself exponential, time between transition firings is exponential.



Stochastic Activity Networks

The need for more expressive modeling languages has led to several extensions to stochastic Petri nets. One extension that we will examine is called *stochastic activity networks*. Because there are a number of subtle distinctions relative to SPNs, stochastic activity networks use different words to describe ideas similar to those of SPNs.

Stochastic activity networks have the following properties:

- A general way to specify that an activity (transition) is enabled
- A general way to specify a completion (firing) rule
- A way to represent zero-timed events
- A way to represent probabilistic choices upon activity completion
- State-dependent parameter values
- General delay distributions on activities

SAN Symbols

Stochastic activity networks (hereafter SANs) have four new symbols in addition to those of SPNs:

- Input gate:  used to define complex enabling predicates and completion functions
- Output gate:  used to define complex completion functions
- Cases:  (small circles on activities) used to specify probabilistic choices
- Instantaneous activities:  used to specify zero-timed events

SAN Enabling Rules

An input gate has two components:

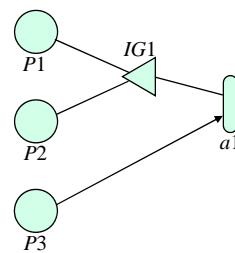
- enabling_function(state) → boolean; also called the *enabling predicate*
- input_function(state) → state; rule for changing the state of the model

An activity is *enabled* if for every connected input gate, the enabling predicate is true, and for each input arc, the number of tokens in the connected place \geq number of arcs.

We use the notation $MARK(P)$ to denote the number of tokens in place P .

Example SAN Enabling Rule

Example:



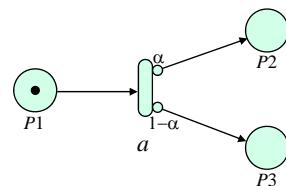
IG1 Predicate:

```
if( (MARK(P1)>0 && MARK(P2)==0) ||  
     (MARK(P1)==0 && MARK(P2)>0) )  
    return 1;  
else return 0;
```

Activity a_1 is enabled if $/G1$ predicate is true (1) and $MARK(P3) > 0$.
(Note that "1" is used to denote true.)

Cases

Cases represent a probabilistic choice of an action to take when an activity completes.



When activity a completes, a token is removed from place P_1 , and with probability α , a token is put into place P_2 , and with probability $1 - \alpha$, a token is put into place P_3 .

Note: cases are numbered, starting with 1, from top to bottom.

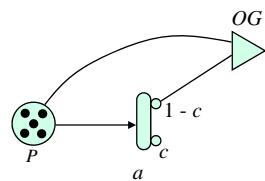
Output Gates

When an activity completes, an output gate allows for a more general change in the state of the system. This output gate function is usually expressed using pseudo-C code.

Example

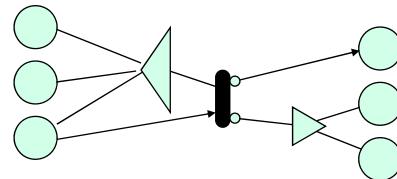
OG Function

MARK(P) = 0;



Instantaneous Activities

Another important feature of SANs is the instantaneous activity. An *instantaneous activity* is like a normal activity except that it completes in zero time after it becomes enabled. Instantaneous activities can be used with input gates, output gates, and cases.

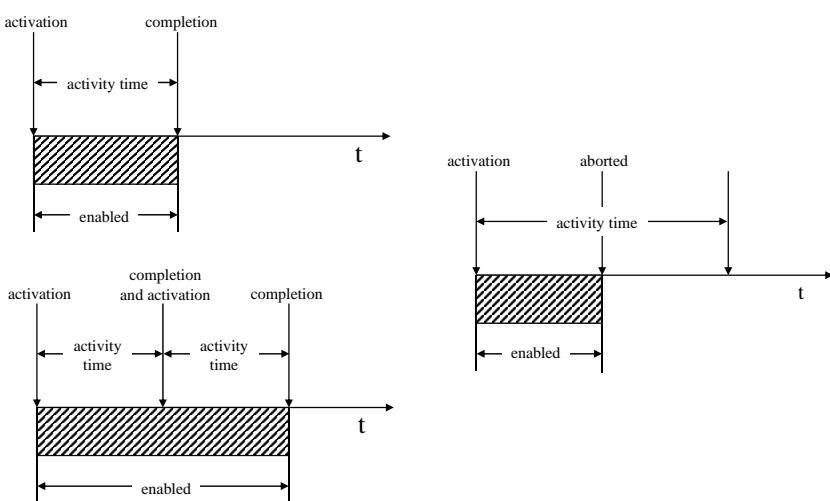


Instantaneous activities are useful when modeling events that have an effect on the state of the system, but happen in negligible time, with respect to other activities in the system, and the performance/dependability measures of interest.

SAN Terms

1. *activation* - time at which an activity begins
2. *completion* - time at which activity completes
3. *abort* - time, after activation but before completion, when activity is no longer enabled
4. *active* - the time after an activity has been activated but before it completes or aborts.

Illustration of SAN Terms



Completion Rules

When an activity *completes*, the following events take place (in the order listed), possibly changing the marking of the network:

1. If the activity has cases, a case is (probabilistically) chosen.
2. The functions of all the connected input gates are executed (in an unspecified order).
3. Tokens are removed from places connected by input arcs.
4. The functions of all the output gates connected to the chosen case are executed (in an unspecified order).
5. Tokens are added to places connected by output arcs connected to the chosen case.

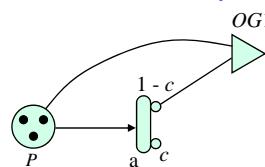
Ordering is important, since effect of actions can be marking-dependent.

Marking Dependent Behavior

Virtually every parameter may be any function of the state of the model. Examples of these are

- rates of exponential activities
- parameters of other activity distributions
- case probabilities

An example of this usefulness is a model of three redundant computers where the coverage (probability that a single computer crashing crashes the whole system) increases after a failure.



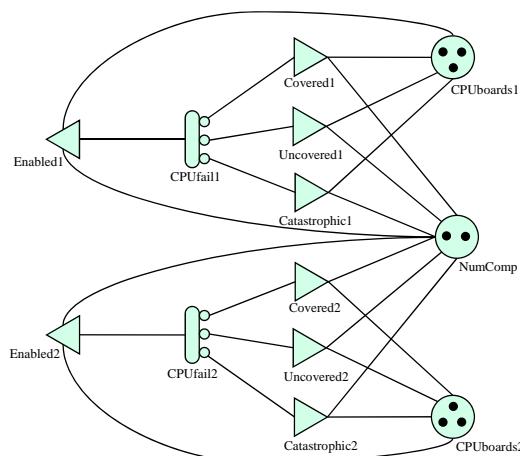
	a
case 1	$0.1 + 0.02 * \text{MARK}(P)$
case 2	$0.9 - 0.02 * \text{MARK}(P)$

Fault-Tolerant Computer Failure Model Example

A fault-tolerant computer system is made up of two redundant computers. Each computer is composed of three redundant CPU boards. A computer is operational if at least 1 CPU board is operational, and the system is operational if at least 1 computer is operational.

CPU boards fail at a rate of $1/10^6$ hours, and there is a 0.5% chance that a board failure will cause a computer failure, and a 0.8% chance that a board will fail in a way that causes a catastrophic system failure.

SAN computer for Computer Failure Model



Activity Case Probabilities and Input Gate Definition

<i>Activity</i>	<i>Case</i>	<i>Probability</i>
<i>CPUfail1</i>	1	0.987
	2	0.005
	3	0.008

<i>Gate</i>	<i>Definition</i>
<i>Enabled1</i>	<u>Predicate</u> $MARK(CPUboards1 > 0) \&\& MARK(NumComp) > 0$
	<u>Function</u> $MARK(CPUboards1)--;$

Output Gate Definitions

<i>Gate</i>	<i>Definition</i>
<i>Covered1</i>	<u>Function</u> $if(MARK(CPUboards1) == 0)$ $MARK(NumComp)--;$
<i>Uncovered1</i>	<u>Function</u> $MARK(CPUboards1) = 0;$ $MARK(NumComp)--;$
<i>Catastrophic1</i>	<u>Function</u> $MARK(CPUboards1) = 0;$ $MARK(NumComp) = 0;$

Reward Variables

- Reward variables are a way of measuring performance- or dependability-related characteristics about a model.
- Examples:
 - Expected time until service
 - System availability
 - Number of misrouted packets in an interval of time
 - Processor utilization
 - Length of downtime
 - Operational cost
 - Module or system reliability
- Types:
 - A model may be in a certain state or states for some period of time, for example, “CPU idle” states. This is called a *rate reward*.
 - An activity may complete. This is called an *impulse reward*.

Reward Variables for Example

Reliability	
	<u>Rate rewards</u> <i>Subnet = computer</i> <u>Predicate:</u> <i>MARK(NumComp) > 0</i> <u>Function:</u> <i>I</i>
	<u>Impulse reward</u> <i>none</i>
NumBoardFailures	
	<u>Rate reward</u> <i>none</i>
	<u>Impulse reward</u> <i>Subnet = computer</i> <i>activity = CPUfail1, value = 1</i> <i>activity = CPUfail2, value = 1</i>

Reward Variables for Example

Performability	
	<u>Rate rewards</u> <i>Subnet = computer</i> <u>Predicate:</u> <i>I</i> <u>Function:</u> <i>MARK(NumComp)</i>
	<u>Impulse reward</u> <i>none</i>
NumBoards	
	<u>Rate reward</u> <i>Subnet = computer</i> <u>Predicate:</u> <i>I</i> <u>Function:</u> <i>MARK(CPUBboards1) + MARK(CPUboards2)</i>
	<u>Impulse reward</u> <i>none</i>

Mobius Software

- Graphically builds and simulates SAN's
- Includes abstraction:
 - Models
 - These can be either atomic models or other composed models
 - Cycles are not allowed
 - Join nodes
 - Join nodes define a set of shared variables that are joined across the child models and nodes
 - Shared variables must all have the same type and initial value
 - Rep nodes
 - Rep nodes define a set of identical replicas and a set of shared variables across those replicas

Mobius Examples (v 2.5)

- Running an Existing Example:
 - Projects -> Unarchive (Multi-Proc)
 - Projects -> Open (Multi-Proc) , re-save if necessary
 - Documentation:
https://www.mobius.illinois.edu/wiki/index.php/Fault-Tolerant_Multiprocessor_Model
- Note the rewards for tracking performance
- Note the study variable sweeps for the experiments
- Open either the simulation or transformer
- Collect the output (ASCII) and map it to a plot (Excel)

Duplicating Example in Presentation

- Projects -> New
- Double click Atomic
 - SAN Model, give it a name “combined_pc_model”
 - Replicate the system diagram
 - (Use Variable->Mark() instead of MARK(Variable))
 - Setup the reward
 - Create an experiment/study
 - Use a simulated solver

Markov Chains

- **Markov process:** Probability distribution for future state only depends on the present state, not on how the process arrived at the present state
- **Markov chain:** State space is discrete
- **Discrete time Markov chain:** Time is discrete
- $X_0, X_1, \dots, X_n, \dots$: observed state at discrete times $t_0, t_1, \dots, t_n, \dots$
- $X_n = j \Rightarrow$ system state at time step n is j .
- $P(X_n = i_n | X_0 = i_0, X_1 = i_1, \dots, X_{n-1} = i_{n-1})$
 $= P(X_n = i_n | X_{n-1} = i_{n-1})$ (Markov Property)
- $p_{jk}(m,n) \equiv P(X_n = k | X_m = j), 0 \leq m \leq n$ (conditional pmf)
- $p_j(n) \equiv P(X_n = j)$ (unconditional pmf)

Homogeneous Markov Chain

- Homogeneous Markov chain: $p_{jk}(m,n)$ only depends on $n-m$
- n-step transition probability: $p_{jk}(n) = P(X_{m+n} = k | X_m = j)$
 - 1-step transition probability: $p_{jk} = p_{jk}(1) = P(X_n = k | X_{n-1} = j)$
- Initial probability row vector: $\mathbf{p}(0) = [p_0(0), p_1(0), \dots, p_k(0), \dots]$
- Transition probability matrix:

$$P = [p_{ij}] = \begin{bmatrix} p_{00} & p_{01} & p_{02} & \dots & \dots \\ p_{10} & p_{11} & p_{12} & \dots & \dots \\ \vdots & \vdots & \vdots & \ddots & \ddots \end{bmatrix}$$

Markov Chain Example

1. System is in state 0 if it is operational, state 1 if it is undergoing repair. Homogeneous DTMC process.

$$P = \begin{bmatrix} 1-a & a \\ b & 1-b \end{bmatrix}, 0 \leq a, b \leq 1$$

2. A sequence of successive software runs, each run has two possible outcomes: success or failure. X_n denotes outcome of n^{th} run (0 for success, 1 for failure). Occurrence of failure in a run depends on the outcome of previous run.

$$\begin{aligned} P(X_{n+1}=1 | X_n=0) &= 1-p \\ P(X_{n+1}=1 | X_n=1) &= q \end{aligned}$$

ECE 695/CS 590

35



Computation of n-step Transition Probabilities

- For a DTMC, find $p_{ij}(n) = P(X_{m+n} = j | X_m = i)$
- Events: State reaches k (from i) & reaches j (from k) are independent due to the Markov property (i.e. no history)
- Invoking the theorem of total probability:
$$p_{ij}(m+n) = \sum_k p_{ik}(m) p_{kj}(n)$$
- Let $P(n)$: n -step prob. transition matrix (i,j) th entry is $p_{ij}(n)$.
Making $m=1$, $n=n-1$ in the above equation:

$$P(n) = P \cdot P(n-1) = P^n$$

$$p_j(n) = P(X_n = j) = \sum_i P(X_0 = i) P(X_n = j | X_0 = i) = \sum_i p_i(0) p_{ij}(n)$$

ECE 695/CS 590

36



Computation of n-step Transition Probabilities

- The pmf of the random variable X_n is

$$p_j(n) = P(X_n = j) = \sum_i P(X_0 = i)P(X_n = j|X_0 = i) = \sum_i p_i(0)p_{ij}(n)$$

- j , in general can assume countable values, 0,1,2,
Defining,

$$\mathbf{p}(n) = [p_0(n), p_1(n), \dots, p_j(n), \dots]$$

- $p_j(n)$ for $j=0,1,2,\dots,k,\dots$ can be written in the vector form as,
 $[p_0 \ p_1 \ p_2 \ \dots \ p_j \ \dots](n) =$

$$\underbrace{[p_0 \ p_1 \ p_2 \ \dots \ p_j \ \dots](0)}_{\mathbf{p}(0)} \begin{bmatrix} p_{00}(n) & p_{01}(n) & \dots & p_{0j}(n) & \dots \\ p_{10}(n) & p_{11}(n) & \dots & p_{1j}(n) & \dots \\ p_{20}(n) & p_{21}(n) & \dots & p_{2j}(n) & \dots \\ \vdots & \vdots & & \vdots & \vdots \\ p_{i0}(n) & p_{i1}(n) & \dots & p_{ij}(n) & \dots \\ \vdots & \vdots & & \vdots & \vdots \end{bmatrix}$$

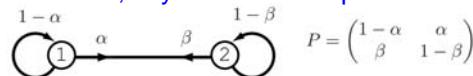
ECE 695/CS 590

37

PURDUE
UNIVERSITY

Two state Markov chain example

- For a two state DTMC, say the transition prob. matrix is



The eigenvalues are 1 and $1 - \alpha - \beta$. So we can write

$$P = U \begin{pmatrix} 1 & 0 \\ 0 & 1 - \alpha - \beta \end{pmatrix} U^{-1} \implies P^n = U \begin{pmatrix} 1 & 0 \\ 0 & (1 - \alpha - \beta)^n \end{pmatrix} U^{-1}$$

So $p_{11}^{(n)} = A + B(1 - \alpha - \beta)^n$ for some A and B .

But $p_{11}^{(0)} = 1 = A + B$ and $p_{11}^{(1)} = 1 - \alpha = A + B(1 - \alpha - \beta)$. So $(A, B) = (\beta, \alpha)/(\alpha + \beta)$, i.e.

$$P_1(X_n = 1) = p_{11}^{(n)} = \frac{\beta}{\alpha + \beta} + \frac{\alpha}{\alpha + \beta}(1 - \alpha - \beta)^n.$$

Note that this tends exponentially fast to a limit of $\beta/(\alpha + \beta)$.

Other components of P^n can be computed similarly, and we have

$$P^n = \begin{pmatrix} \frac{\beta}{\alpha+\beta} + \frac{\alpha}{\alpha+\beta}(1 - \alpha - \beta)^n & \frac{\alpha}{\alpha+\beta} - \frac{\alpha}{\alpha+\beta}(1 - \alpha - \beta)^n \\ \frac{\beta}{\alpha+\beta} - \frac{\beta}{\alpha+\beta}(1 - \alpha - \beta)^n & \frac{\alpha}{\alpha+\beta} + \frac{\beta}{\alpha+\beta}(1 - \alpha - \beta)^n \end{pmatrix}$$

Source: <http://www.statslab.cam.ac.uk/~rrw1/markov/M.pdf>

ECE 695/CS 590

38

PURDUE
UNIVERSITY

Two state Markov chain example

- Then, the n-step transition probability matrix $P(n) = P^n$ is
$$P^n = \begin{pmatrix} \frac{\beta}{\alpha+\beta} + \frac{\alpha}{\alpha+\beta}(1-\alpha-\beta)^n & \frac{\alpha}{\alpha+\beta} - \frac{\alpha}{\alpha+\beta}(1-\alpha-\beta)^n \\ \frac{\beta}{\alpha+\beta} - \frac{\beta}{\alpha+\beta}(1-\alpha-\beta)^n & \frac{\alpha}{\alpha+\beta} + \frac{\beta}{\alpha+\beta}(1-\alpha-\beta)^n \end{pmatrix}$$
- Example: Consider a communication network consisting of a sequence of binary communication channels, with $\alpha=1/4$, $\beta=1/2$. Here X_n denotes the bit leaving the n th stage and X_0 represents the bit entering the system.
 - If a bit enters as 1, what is the probability of being correctly transmitted over (a) two stages (b) three stages?
 - Assuming initial probabilities $P(X_0=0) = 1/3$ and $P(X_0=1) = 2/3$, what is $p(n)$?