# CHAPTER 5: SIMULATION MODELS

Consider the problem of simulating a stochastic process $\{X_t; t\in T\}$ with the goal of evaluating one (or several) performance functors

$$J = \mathbb{E}(\phi(X_t; t\in T)),$$

where $\phi$ is a functional of the whole trajectory $X_t, t\in T$, that takes values in $\mathbb{R}$.

$(\Omega, \mathcal{F})$ and $\{\mathcal{F}_t; t\in T\} = \mathbb{F}$ a filtration, $\mathcal{F}_t \subset \mathcal{F}$.

Remark: Because each trajectory $(X_t(\omega); t\in T)$ is uniquely defined for each $\omega$, it follows that $\phi(X_t(\omega); t\in T)$ is a well defined random variable on $(\Omega, \mathcal{F})$.

Example: Let $X_n$ be the amount of apples kept in store for sale in the cafeteria. The manager wishes to evaluate the policies of ordering and reducing prices for quick sales. Assume that at the end of a period, any remaining apples must be discarded. Demand $\{d_n(u)\}$ is assumed to be iid,

given a price $u$ per apple (the mean depends on price).
Ordering cost are of the form $K+p y_n$, where $y_n$ is the number of apples ordered at the end of period $n$.

Here $X_{n+1} = \cancel{y_n}\ y_n$, and the cost is:

$$C_n = K + X_n p - \max(X_n, d_n(u)) \cdot u$$

A simulation can be done by generating $\{(X_n, d_n(u))\}$ and then the running costs can be calculated using

$$\frac{1}{N}\sum_{n=1}^{N} C_n(X_n, u).$$

- Continuous model for simulation (tick-based)
- Discrete event model (event-based)
- Standard clock model
- Reduced models (Petri-nets, transformations)

Tick-based

(a) Continuous Simulation Model (Tick-simulation)

Origin of model in deterministic context: numerical reproduction of trajectories of a dynamic physical system (computer animation for billard or pool games). Here the model is described via ODE's (or PDE's more generally):

$$\frac{dx_t}{dt} = v(x,t) \qquad x_t\in\mathbb{R}^d; \quad T=(0,T]$$

$x_0\in\mathbb{R}^d$ is known as the initial position.

$v: \mathbb{R}^d\times T \to \mathbb{R}^d$ is called a vector-field or "drift"

and in mechanical systems it represents the velocity $v$ at each point in space.

Step-by-step or tick-wise animation uses a discretization of time into "ticks" or small units of length $h>0$:

$$x(i+h) = x(i) + v(x(i), ih)$$
$$i=1,\dots \lfloor T/h \rfloor, \quad x(0)=x_0.$$

Thm: If $v$ is continuous and bounded, then the piecewise linear interpolation $x_t^h$ of the sequence $\{x(i)\}$ converges in the sup-norm to the solution of the ODE for every initial condition $x_0\in\mathbb{R}^d$.

The proof of this result follows from Ascoli–Arzelà Theorem. ②

Example: Consider the Black-Scholes model of a geometric Brownian motion for the stock price:

$$S_t = S_0 e^{\mu t + \sigma B_t} \quad ; \quad S_0 \in \mathbb{R}$$

where $B(\cdot)$ denotes the standard Brownian motion.

[ Def: A stochastic process $\{B(t) ; t > 0\}$ on $(\Omega, \mathcal{F}, \mathbb{P})$ is called a standard Brownian motion if:

(i) $B(0) = 0$ a.s.

(ii) $\{B(t+s) - B(t)\}$ is independent of $\mathcal{F}_t \ \forall t, s > 0$

(iii) for all $t, s > 0$ $\{B(t+s) - B(t)\} \sim N(0, s)$ (normal dist.)

]

Using a discretization, the simulation by ticks considers the discrete-time sampling:

$$S(i+1) = S(i) \exp\{\mu h + \sigma h Z_i\} \ ; \ i = 1, \ldots \lfloor T/h \rfloor$$

where $\{Z_i\}$ are iid $N(0,1)$.

A financial option or derivative on the asset is of the form:

$$\max(0, \Phi(\{S_t ; t \le T\}))$$

and can be approximated using the discretization.

A European option: $\Phi(\{S_t ; t \le T\}) = S_T - K$

An American option:

A Bermudan option: $\dfrac{1}{T} \int_0^T S_t\, dt$ // Barrier: $(S_T - K)\mathbb{1}(S_T > B ; t \le T)$

$\dfrac{1}{T} \sum^{\lfloor T/h \rfloor} S_{nh}$ // etc.

$N = \text{INT}(T/h)$, $S[0] = S_0$, $c = \exp(\mu h)$

FOR $i = 1, \ldots N$

Generate $Z \sim N(0,1)$

$S[i] = S[i-1] * c * \exp(\sigma h \bar{Z})$

To make it more efficient, work with log-prices to avoid exponentiation calculation:

$X[i] = X[i-1] + \sigma * h * \bar{Z}$ in loop.

then $S[i] = C * \exp(X[i])$ can be calculated after

↯ How do we know that this will work?

Def: Let $(X(t) ; t > 0)$ be a stochastic process on $(\Omega, \mathcal{F}, \mathbb{P})$ with $X(0) = x_0$ a given random variable. A tick-based (continuous) simulation model for the process is a family of stochastic processes in discrete-time, indexed by $h > 0$: $\{X^h(n) ; n \in \mathbb{N}\}$ with $X^h(0) = x_0 \ \forall h$, that satisfies:
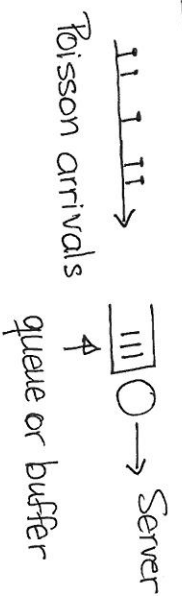
$$\forall t \in (0, T] \quad \lim_{\substack{n \to \infty \\ nh = t}} X^h(n) \Rightarrow X(t).$$

The symbol "⇒" means convergence in distribution.

Remark: most performance functions are Lipschitz continuous and convergence of $\{X^h(\cdot)\}$ is sufficient to ensure convergence of the corresponding discrete performance function, but in general, this requires verification.

# Example: Queueing model FCFS

Poisson arrivals → [ ||| ] ○ → Server
                    queue or buffer

$N(t)$ : number of arrivals up to time $t$  ($\sim$ Poisson($\lambda t$))

$\{\xi_i\}$ iid $\sim G$ are the consecutive service times.

Example of performance or objective functions are : mean queue length, probability that the queue size is larger than a given value, probability that the waiting times of clients is larger than a certain value, mean idle time of server, etc.

How do we build a continuous simulation?

Thm: Let $\{A_n^h\} \sim$ iid Bernoulli random variables with parameter $p = \lambda h$, for $h > 0$ small enough so that $\lambda h < 1$. Let $X_n^h$ be the number of arrivals, $X_n^h \sim \text{Bin}(n, \lambda h)$ then for any $t \in \mathbb{R}$,

$$X_n^h \underset{h \to 0}{\Longrightarrow} \text{Poisson}(\lambda t).$$
$$n = t/h$$

[proof in references, wikipedia, etc...].

For fixed $h$ (we will drop the subscript notation) we have $\int n = 1, \ldots \lfloor T/h \rfloor$
$A_n \sim \text{Ber}(\lambda h) \in \{0, 1\}$
$D_n$: number of service completions in $(nh, (n+1)h]$.
$Q(n)$: queue size at start of period.

Consecutive queue lengths satisfy:

$$Q(n+1) = Q(n) + (A_n - D_n).$$

Is this enough information? How can we know what $D_n$ is?

At arrival time of $i$-th client, that is

$$\check{i} = \min \left( n : \sum_{k=1}^{n} A_k = i \right)$$

we can generate the random variable $\xi_i$. But this will complicate the code with all concurrent arrivals in queue having to store values. Another possibility: Augmentation of state space to include memory in service :

Let $R_n$ = residual service time at period $n$.

Then:

$$R_{n-1} = \begin{cases} R_n - 1 & \text{if } R_n > 1 \\ \lfloor \xi_{n+1}/h \rfloor & \text{otherwise, if } R_n = 0 \ \& \ Q_n \geq 1 \end{cases}$$

because if $R_n = 0$, there is a service completion during this period and it is immediately followed by a new customer entering service.

$$Q(n+1) = Q(n) + A_n - \mathbb{1}(R_n = 0 \ \& \ Q(n) \geq 1)$$
$$R(n+1) = (R(n) - 1)\mathbb{1}(R(n) \geq 1)$$
$$+ \lfloor \xi_{n+1}/h \rfloor \mathbb{1}(R(n) = 0 \ \& \ Q(n) \geq 1)$$

The above system of equations can be simulated iteratively. Because $\{A_n, \xi_n\} \sim$ iid, then it follows that $\{Q(n), R(n)\}$ is Markovian.

## Exercises!

Let $X_n^h = \{(Q(n), R(n))\}$ for given $h > 0$.

• Show that $\{X_n^h\}$ is a Markov chain in two dimensions.
How many classes? Recurrence?

• Show that as $h \to 0$, for each $t \in \mathbb{R}$ the random sequence $\{Q(\lfloor t/h \rfloor)\}$ converges in distribution to the original queue process.

Ask students about initializing.

% service first

If $R(n) = 0$    % completion of service
if $Q(n) > 0$
{
  $Q(n+1) = Q(n) - 1$
  Generate $S \sim G$
  $R(n+1) = \lfloor S^*/h \rfloor$
}

% arrivals:

Generate $A \sim Ber(\lambda h)$

$Q(n+1) = Q(n) + 1$

If $Q(n) = 1$    % first customer
{
  Generate $S \sim G$
  $R(n+1) = \lfloor S/h \rfloor$
}

% clock-tick advance clock:
$R(n+1) = R(n-1) - 1h$

Important Questions: Why are limits valid? How small should $h$ be? How can we estimate approximation errors? How small should $h$ be? All of these questions are studied in the field of Simulation.

---

④

## A PARENTHESIS FOR MORE PROBABILITY

Def: Let $\{X_n\}_{n \in \mathbb{N}}$ be a sequence of random variables on a common probability space $(\Omega, \mathcal{F}, \mathbb{P})$, and $X$ a random variable. We say that:

• $X_n$ converges almost surely to $X$ if
$$\mathbb{P}(\omega : \lim_{n \to \infty} X_n(\omega) = X(\omega)) = 1. \quad X_n \to X \text{ a.s.}$$

• $X_n$ converges weakly or in distribution to $X$ if
$$\lim_{n \to \infty} \mathbb{P}(X_n \leq x) = \mathbb{P}(X \leq x) \text{ for every } x \in \mathbb{R}$$
that is a point of continuity of $F(x) = \mathbb{P}(X \leq x)$.

Result: Convergence in distribution: $X_n \Rightarrow X$ is equivalent to the condition that for every continuous and bounded function $g : \mathbb{R} \to \mathbb{R}$,
$$\lim_{n \to \infty} \mathbb{E}[g(X_n)] = \mathbb{E}[g(X)].$$

There are a number of important theorems that can be used to establish convergence: Dominated Convergence Theorem, Monotone Convergence Theorem, etc.

Motivation for event-based simulation (notes P. 44-45)

service $S_n$
$S_n$: service time
$\xi_n$: service time
time between arrivals

⟹ many loops "do nothing"
idea: jump ahead in one go

⟹ small $h$ for accuracy

# (0) Discrete Event Model

Physical state $S$ (usually but not always countable)

Possible event set (finite number of distinct events) $E$, $|E| = d$

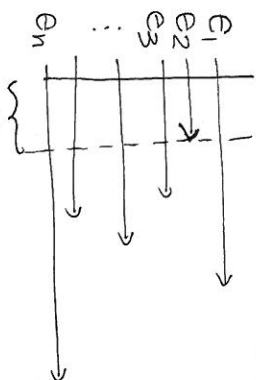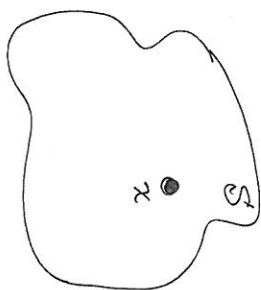For each $x \in S$, $T(x) \subset E$ is the set of possible events @ $x$.

Main dynamics: by "jumps"

$$\mathbb{P}(X\text{"new"} = j \mid X\text{"old"} = x, \text{event} = e) = P_S(j; x, e)$$

Clock dynamics: also known and Markovian:

$$\mathbb{P}(\text{time for new event's } t \mid X\text{"new"} = \mathbf{x}, \text{event} = e) = \overline{F_e}(t, x)$$
given distribution (known)



$e_1$
$e_2$
$e_3$
...
$e_n$

next event in $e_2$

$z$ = elapsed time

In $z$ units of time the state will jump from $x$ to another state, according to the probability

$$\mathbb{P}(X\text{"new"} \in \cdot \mid x, e_2)$$



→ GENERATE NEW TIME FOR
NEXT EVENT OF THIS TYPE

advance clock and update

---

**Def:** A discrete event process $Z_t$ on $(\Omega, \mathcal{F}, \mathbb{P})$ has a physical component $X_t$ and a clock component $Y_t$.

Let:

$e(x,y) = \operatorname{argmin}(y_i, i \in T(x))$ next event

$\tau(x,y) = \min(y_i, i \in T(x))$ time to next event

$z = (x,y) \subset \mathbb{R}^n \times \mathbb{R}^d$

$$\text{Prob}(X_{t+z} \in A \mid Z_{t} = (x,y)) = P_S(A; x, e(x,y))$$
for any $A \in \mathcal{B}(S)$

If $i \neq e(x,y)$ then

$y_{i, t+\tau} = y_{i,t} - \tau$ to update clocks

$$\mathbb{P}(Y_{e(z),t+\tau} \leq t \mid Z_t = (x,y) = z) = \overline{F}_{e(z)}(t, x)$$

where $P_S(\cdot; x, e)$ are well defined probability distributions for any $x \in S$ and $e \in E$, and for each event $e \in E$, $\overline{F}_e(\cdot, x)$ is a well defined distribution.

**Remark:** also called "stochastic timed automata"

**Result:** It is left as an exercise for students to show that the embedded discrete-time process

$$Z_n = \{Z_{\tau_n}\}, \text{ where } \tau_i = \text{time of } i\text{-th event, is a Markov Chain on } S \times \mathbb{R}^d.$$

Also called "Generalized Semi-Markov Process"

$\{Y_1, Y_2\}$ list of residual arrival (1) and departure (2) times

(initialize $Y_1$ = Generate new arrival)

$\{T_i\} \sim$ iid dist. $F$ (not necessarily a Poisson arrival)

$\{\xi_i\} \sim$ iid dist $G$

```
main loop:
  if Q=0 => e=1 else
  e = argmin(Y1, Y2) ∈ {1,2} ;
  τ = min(Y1, Y2)
case(e=1) % arrival
    Q = Q+1
    T ~ F ,  Y1 = T   % new inter-arrival

    if Q=1 => Y2 = ξ ~ G
case(e=2) % service
    Q --
    If Q > 0 => Y2 = ξ ~ G  % new service
```

What is the simulation model?

Variables: $Q_n, Y_{1,n}, Y_{2,n} = Z_n$

Notice that $\{Z_n\}$ is a Markov chain and consider the natural filtration $\sigma(Z_1, \ldots Z_n) = \mathcal{F}_n$.

Ex: Consider the physical process $\{Q_t;\ t \geq 0\}$ and let $\{\tau_1, \tau_2, \ldots\}$ be the consecutive event or jump times. Let $\tilde{\mathcal{F}}_t = \sigma(Q_s;\ s \leq t)$. Explain the difference between $\tilde{\mathcal{F}}_{\tau_n}$ and $\mathcal{F}_n$.

Simulation package SSJ highly recommended

## (C) Standard Clock

Example in telecom router

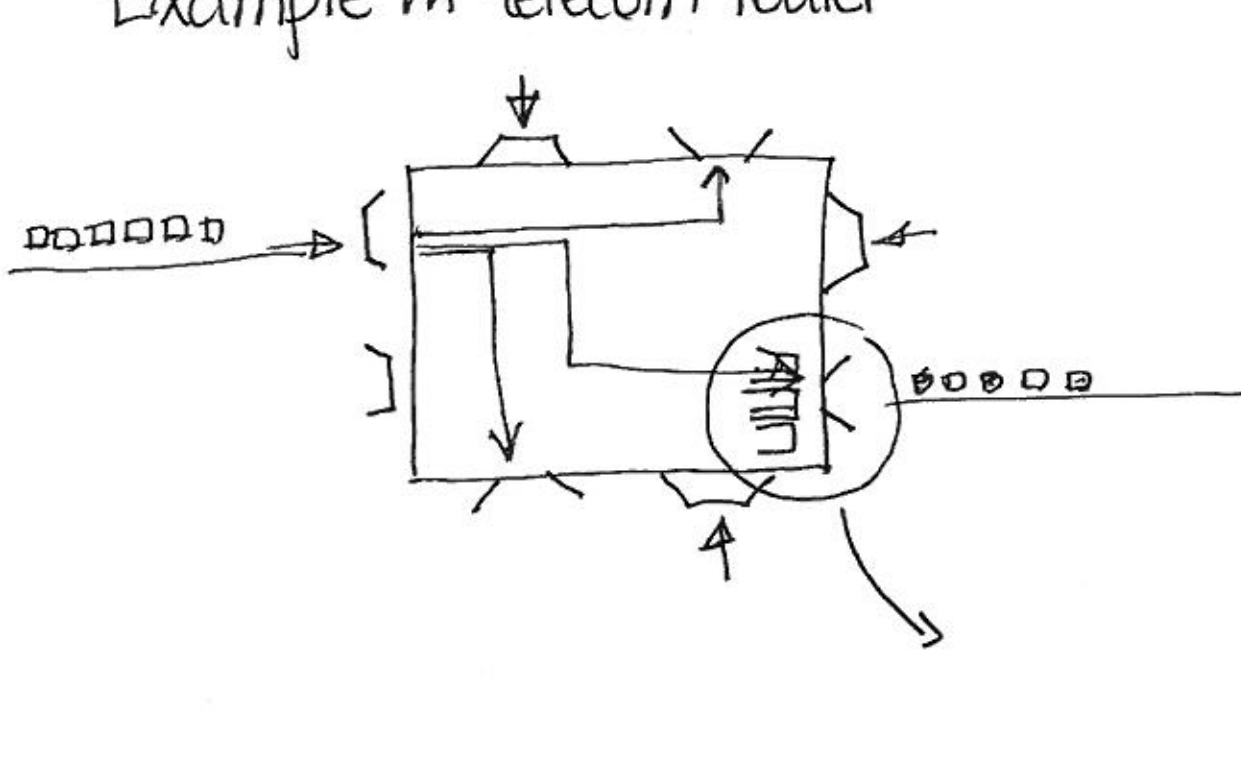

Exit queues are classified according to priority customers (voice, video, real-time, etc). $C$ classes

How many queues in network?

Each "switch" has $CN$ queues ($N$ is the number of exit ports). If there are $M$ identical switches in the network then there are $CNM$ number of queues, $MN$ #servers

How many events in $E$?

@each server : 1 residual service
$C$ residual arrival times

Needs $MN(C+1)$ clocks in list

List search can be slow, even if heap or other methods are used to accelerate.

Assume that residual times are exponentially distributed & independent.

**Result:** Let $X \perp\!\!\!\perp Y$ be exponential rv's, on a common space $(\Omega, \mathcal{F}, \mathbb{P})$ with intensities $\lambda_1, \lambda_2$ rep. Then

$$e = \min(X, Y) \stackrel{d}{=} \exp(\lambda_1, \lambda_2)$$

**Proof:**
$$\mathbb{P}(e \le t) = 1 - \mathbb{P}(e > t) = 1 - \mathbb{P}(X > t, Y > t)$$
$$= 1 - \mathbb{P}(X > t)\mathbb{P}(Y > t) =$$
$$= 1 - e^{-\lambda_1 t} e^{-\lambda_2 t} = 1 - e^{(\lambda_1 + \lambda_2)t}$$

$$\Rightarrow e \stackrel{d}{=} \exp(\lambda_1 + \lambda_2).$$

**Proposition:** Let $(Y_1, \ldots Y_d)$ d independent rv's with exponential distribution of intensities $(\lambda_1, \ldots \lambda_d)$ respectively. Then

$$\tau = \min(Y_1, \ldots, Y_d) \stackrel{d}{=} \exp(\Lambda)$$

$$\Lambda = \sum_{i=1}^{d} \lambda_i.$$

**Proposition:** Let $Y \stackrel{d}{=} \exp(\Lambda)$ ~~event arrival~~ be the next event time, $\tau = \min(Y_1, \ldots Y_d)$. ~~Then that~~ and $e = \operatorname{argmin}(Y_1, \ldots Y_d)$ then $\mathbb{P}(e = i \mid Y) = \dfrac{\lambda_i}{\Lambda}$.

The proof is left as an exercise (start with d=2).

---

such that:

**Def:** Given $E, (\lambda_e, e \in E)$ a standard clock model is a stochastic process $\{X_n\} \, \{Z_n\}$

~~$\tau_n \stackrel{d}{=} \eta \text{ mp } (\Lambda_n)$~~ $Z_n \stackrel{d}{=} \exp(\Lambda_n)$ $\Lambda_n = \sum_{e \in T(X_n)} \lambda_e$ → time to next event

$$\mathbb{P}(e_{n+1} = e) = \frac{\lambda_e}{\Lambda_n}, \ \forall e \in T(X_n) \text{ next event type}$$

$$\mathbb{P}(X_{n+1} \in A \mid X_n) = P_S(A; x, e)$$

There is no search through a list $\Rightarrow$ speed up.

**Example:** Let $O(m)$ the set of exit ports at each switch $m$, then we generate $\tau$ using

$$\sum_{m=1}^{M} \left( \Lambda = \sum_{c=1}^{C} \lambda_{c, m t} \sum_{i \in O(m)} M_i \ \mathbb{1}(Q_i > 0) \right)$$

arrivals — state dependency
$X_n = Q_n.$

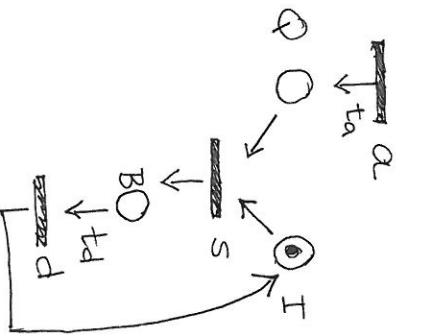**Remark:** generalization to other than exponential

**(4)** Reduced Models

- Process-oriented simulation schemes (Cassandras) and Petri-net representation,
- state aggregation, retrospective simulation models (paper on intelligent subways and bus fleet as examples)

Process-oriented models can also be very useful. Instead of describing the dynamics in terms of the time evolution of the state of the system, these models represent the various processes affecting the system, in a similar way that a Gauntt or PERT chart orders the various stages in project management ("tasks").

Petri nets are used to represent the relationship of precedence, recurrence and interactions as follows.

· Nodes indicate completed stages

· Arcs represent precedence and causal relations.

Example of the FCFS queue



▬ transitions

○ states, called "places"

→ flow relations, called "arcs"

• the "tokens"/"marks" that flow within the Petri net.

a: arrival of clients, always active "timed" $t_a$ are the {Ti} iid inter-arrival times   S: service initiatio

Q: queue size

I: indicator of idle server, busy or not (server)

d: departure, timed
$t_d$ : {$S_n$} ~ G

(a)(s)(d) are called "transitions" (See wikipedia)
Initial marking: one token @ I.
To activate it is necessary that there be tokens in the nodes that point to the transition.

Go through logic: "activate" arrivals.
Call {$A_n$} the consecutive firing times for transition (a) or "epochs"

[note that they correspond to customer arrival times). Call {$D_n$} the consecutive firing times for transition (d); these are customer departure times. Finally, call {$W_n$} the waiting times in the Q place (the queue) for the n-th token. comment]

If $D_n \leq A_{n+1}$ the net is at same as initial marking so that $W_{n+1} = 0$ (as soon as (a) fires, the tokens enable transition s). Otherwise, the time that the n+1-st token has to wait to enable transition s is exactly

$$W_{n+1} = A_{n+1} - D_n.$$ That is:

$$W_{n+1} = \max(0, D_n - A_n).$$

On the other hand, following the customer's process, clearly
$D_n = A_n + W_n + S_n$, where {$S_n$} are the consecutive service times required to fire transition s. This yields:

$$W_{n+1} = \max(0, (A_{n+1} - A_n) - (W_n + S_n))$$

this expression is also known as Lindley equation.

If $X(n) = W_n + S_n$ denotes total time that customer

n spends in the system, then $\{X(n)\}$ is a Markov ⑨
chain on the space $\mathbb{R}^+$.

$X[1] \sim G$    %o service distribution

for $n = 1$ to $N$ do

$\quad T = $ Generate Inter Arrival $(\lambda)$ ;

$\quad S = $ Generate Service $\sim G$ ;

$\quad$ If $(X[n] < T)$    $X[n+1] = S$ ;

$\quad$ else

$\quad\quad\quad\quad X[n+1] = S + (X[n] - T)$ ;

If we are interested in evaluating statistics about the waiting
times then this simulation model is much simpler than ~~this~~
tick- or event - based simulation models.

Example : airport car park.

- - - - - - - - - - - -

| End of Class |

why do we simulate?

$\quad \circ$ Research question

$\quad \circ$ Performance functions

Experimental Design

$\quad \circ$ Methodology and proposed scenarios

Output Analysis

Efficiency of a Simulation