# A <u>Deterministic</u> Model of a Service Station (Fluid View)

**Primitives**

$Z(0)$     initial content

$\alpha(t)$     input rate

$\mu(t)$     *potential* service rate

$$\text{in} \longrightarrow \boxed{\textbf{Delay}} \longrightarrow \text{out}$$

**Model**: (Think cumulants)

$$\text{Inflow:} \quad A(t) = \int_0^t \alpha(u)du, \quad t \geq 0\,;$$

$$\text{Potential Outflow:} \quad M(t) = \int_0^t \mu(u)du, \quad t \geq 0\,.$$

- We could start with primitives $A, M$, in which case they need not be continuous; for example, they could be counting processes.

Netflow:     $X(t) = Z(0) + A(t) - M(t), \quad t \geq 0.$

Introduce     $Y(t) =$ cumulative potential lost during $[0, t]$.

$\Rightarrow$ Outflow:     $D = M - Y$     (**A** arrivals; **D** departures)

$\Rightarrow$ Balance:

$$
\begin{aligned}
Z(t) &= Z(0) + A(t) - D(t) \\
&= Z(0) + A(t) - [M(t) - Y(t)] \\
&= X(t) + Y(t), \quad t \geq 0\,.
\end{aligned}
$$

**Model**     $Z = X + Y$

Feasible     $Z \geq 0, \ Y \uparrow 0$     $(Y(0) = 0)$;

Efficient     $Y$ least     (hence, $Y$ unique);

Existence:     $Y = \overline{(-X)^+}$     $(Y = -\underline{X}$, when $Z(0) = 0)$;

$\underline{X}(t) = \inf_{0 \leq u \leq t} X(u)$, which is called the **lower envelope** of $X$.
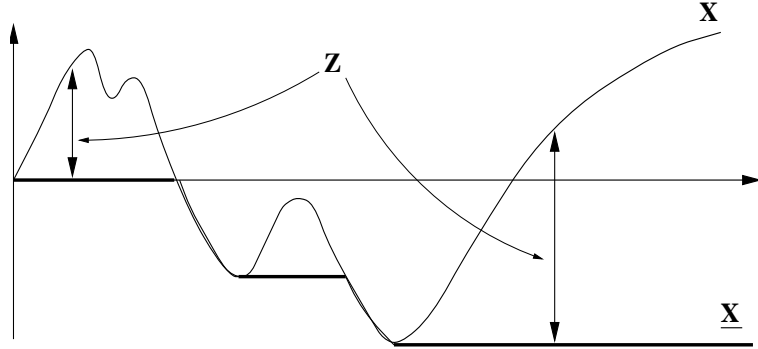
**"Proof"**

Least $Y \uparrow 0$
s.t. $Y \geq -X$



When $Z(0) = 0$:

$Z = X - \underline{X}$,
$\underline{X} = $ lower envelope.



Equivalent characterization via complementarity: (LCP/DCP)

$$Y \text{ least} \iff ZdY = 0, \text{ i.e. } Y \text{ increases at } t$$
$$\text{only when } Z(t) = 0.$$

In words: potential lost due to idleness.

**Claim** (Skorohod)    Given $X \in$ RCLL (**R**ight **C**ontinuous **L**eft **L**imit),

there exists a unique $(Y, Z)$ such that

$$
\begin{aligned}
Z &= X + Y, \\
Z &\geq 0, \quad Y \uparrow 0, \\
ZdY &= 0.
\end{aligned}
$$

**Proof**    Existence by checking    $Y = \overline{(-X)^+} \quad (= -\underline{X} \wedge 0)$.

Uniqueness by Lyapunov-function argument:

(Note: if minimality is established, then uniqueness is automatic.)

If $(Y_i, Z_i)$, $i = 1, 2$, are two solutions, then consider

$$\eta = \frac{1}{2}(Y_1 - Y_2)^2.$$

2

Assume, for simplicity, continuous $Y_i$'s, in which case differentiate:

$$d\eta = (Y_1 - Y_2)(dY_1 - dY_2) = (Z_1 - Z_2)(dY_1 - dY_2)$$
$$= -Z_1 dY_2 - Z_2 dY_1 \leq 0 .$$

Deduce that $\eta$ decreases, but also

$$\eta(0) = 0 \quad \Rightarrow \quad \eta \equiv 0$$
$$\Rightarrow \quad Y_1 \equiv Y_2.$$

**Outflow** $\qquad D(t) = M(t) - Y(t) = \int_0^t \delta(u)du, \quad$ where $\delta(u) =$ outflow rate,

$$\Rightarrow \qquad Y(t) = \int_0^t [\mu(u) - \delta(u)]du .$$

In terms of rates: $dY \geq 0$ implies $\delta \leq \mu$.

Now, either

$\delta = \mu$ or

$\delta < \mu \iff dY > 0,$

$\qquad \Rightarrow Z = 0$ (since $ZdY = 0$),

$\qquad \Rightarrow d(X + Y) = 0$ (consider a neighbourhood and differentiate),

$\qquad \Rightarrow (\alpha - \mu) + (\mu - \delta) = \alpha - \delta = 0.$

Thus (Hall, pg. 190, Def. 6.6),

$$\delta(t) = \begin{cases} \mu(t) & \text{when } Z(t) > 0, \\ \alpha(t) & \text{when } Z(t) = 0 . \end{cases}$$

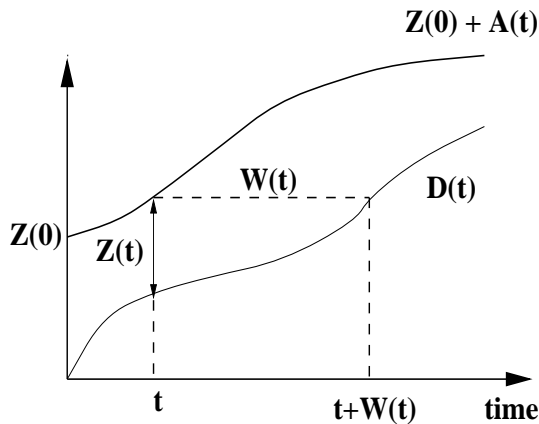**Note** that the above is *not* a direct definition of $\delta$, since it uses $Z$, which is defined in terms of $\delta$.

How to calculate **Delay**?

Define

$$W(t) \;=\; \text{work-load at time } t$$
$$(= \text{time to process all that is present at time } t)$$
$$=\; \text{under FCFS, virtual waiting time.}$$

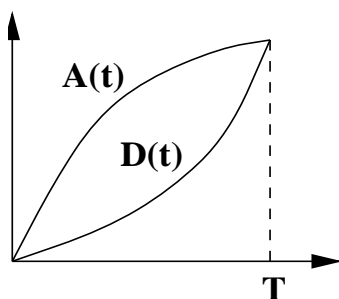Defining relation for $W$:

$$D(t + W(t)) = Z(0) + A(t)$$



Hence, $Z(t + W(t)) = Z(0) + A(t + W(t)) - A(t)$.

**MOP**'s over a finite horizon $T$:

*Averages*   **Inflow:**     $\bar{\alpha} = \frac{1}{T} \int_0^T \alpha(t)dt;$

           **Outflow:**     $\bar{\delta} = \frac{1}{T} \int_0^T \delta(t)dt;$

           **Throughput:**     $\lambda,$ defined when $\bar{\alpha} = \bar{\delta}$ as their common value.



eg.   $\lambda = \frac{1}{T}A(T) = \frac{1}{T}D(T).$

**Queue length** (Inventory):     $\bar{Z} = \frac{1}{T} \int_0^T Z(t)dt = \frac{1}{T} \times \text{Area}.$

**Delay:**     $\bar{W} = \frac{1}{A(T)} \int_0^T W(t)dA(t) \quad \left( = \frac{\int_0^T W(t)\alpha(t)dt}{\int_0^T \alpha(t)dt} \right).$

                      $\uparrow$

            Rieman-Stiltjes

4

*Intuition:*

- Discrete arrivals $\Rightarrow \bar{W} = \frac{1}{A(T)} \sum_{n=1}^{A(T)} W_n$    (as in Hall, Chap. 2);

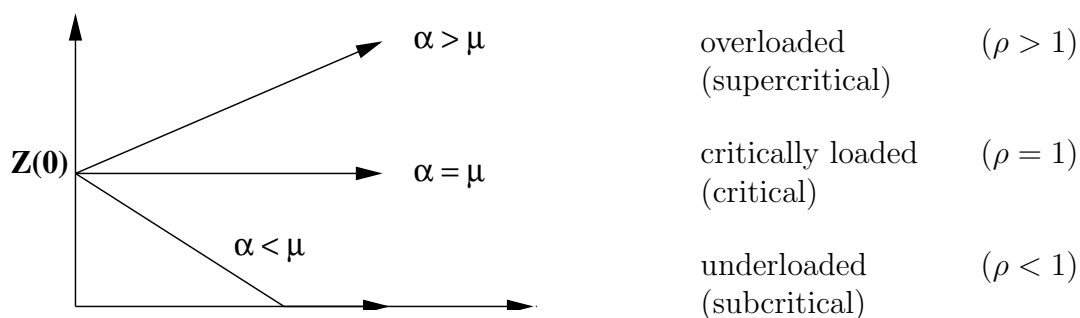- Absolutely continuous: $\alpha(t)dt$ arrivals during $(t, t+dt)$, each suffering a delay of $W(t)$.

**Little's Conservation Law:**    $\bar{Z} = \lambda \cdot \bar{W}$.

**Cumulative lost potential** $Y(T)$.

**Efficiency**    $\varepsilon(T) = 1 - \frac{Y(T)}{M(T)} =$

$$\overset{\text{actual} \searrow}{\underset{\text{potential} \nearrow}{= \frac{D(T)}{M(T)}}} \left( = \frac{\int_0^T \delta(t)dt}{\int_0^T \mu(t)dt}, \text{ when applicable} \right).$$

**Example**    *constant rates*    $\alpha(t) \equiv \alpha$ ,    $\mu(t) \equiv \mu$.
           (linear model)



| | |
|---|---|
| overloaded (supercritical) | $(\rho > 1)$ |
| critically loaded (critical) | $(\rho = 1)$ |
| underloaded (subcritical) | $(\rho < 1)$ |

Definition: $\rho = \alpha/\mu$ **traffic (flow) intensity**.

Natural *extension*: piecewise constant rates, as in National Cranberry (HBS case).

**Example**    *periodic rates* e.g.



(If $\alpha$ has a period $T_\alpha = 8$, $\mu$ has a period $T_\mu = 3$, take period $T = T_\alpha \cdot T_\mu = 24$.)

Long-run: $\quad \bar{\alpha} = \frac{1}{T}\int_0^T \alpha(t)dt; \quad \bar{\mu} = \frac{1}{T}\int_0^T \mu(t)dt;$
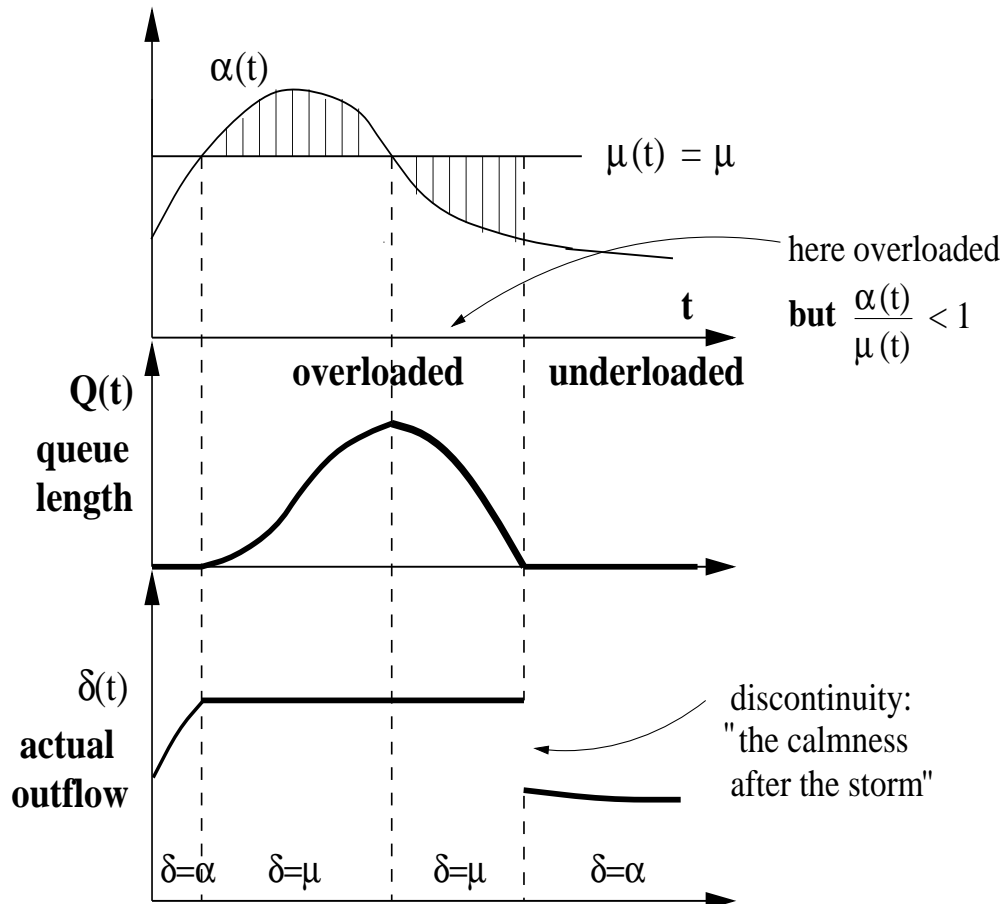
$\qquad\qquad \rho = \bar{\alpha}/\bar{\mu}$ (Heyman-Whitt).

Short-run: $\quad$ Phase-transitions (different from Hall, pg. 189–190, that has stagnant $\to$ growth $\to$ decline $\to$ stagnant).

**Short-Run Phase Transitions**

$\qquad$ Overloaded at $t$ $\quad:\quad Z(t) > 0$;

$\qquad$ Underloaded $\qquad:\quad Z(t) = 0 \quad$ and $\quad \delta(t) < \mu(t) \quad$ (excess capacity, $dY(t) > 0$);

$\qquad$ Critically loaded $\quad:\quad Z(t) = 0 \quad$ and $\quad \delta(t) = \mu(t) \quad$ (balanced capacity, $dY(t) = 0$).



The analogue of $\rho$, traffic intensity, is here (assume $Z(0) = 0$):

$$\rho(t) = \sup_{0 \le s \le t} \frac{\int_s^t \alpha(u)du}{\int_s^t \mu(u)du} \qquad \begin{cases} > 1 & \text{overloaded} \\ = 1 & \text{critically loaded} \\ < 1 & \text{underloaded} \end{cases}$$

6

For finer approximations, we must acknowledge more phases, as depicted in the following figure.



Phase transition diagram for the asymptotic regions.
(Massey & Mandelbaum.)

**References:**

– Hall, R.W., "*Queueing Methods for Service and Manufacturing*", Prentice Hall, 1991.

– Harrison, J.M., "*Brownian Motion and Stochastic Flow Systems*", Wiley, 1985.

– Mandelbaum, A. and Massey, William, A., "Strong approximations for time-dependent queues", *Math. of Operations Research,* 20, 33-64, 1995.

## Mathematical Framework

*Reflection* Mapping $\qquad X \to X - \underline{X} \wedge 0$
(Regulator)

$$(X \to X - \underline{X}, \quad \text{when } X(0) = 0).$$

*Fundamental*:

- Flow analysis (Fluid Models);

- Economics;

- Stochastic Processes;

  - Skorohod (needed cumulant $Y$!);
  - Queueing Models (later);

- Approximations.

*Idea* of Approximations: $\quad Z = f(X), \quad f$ continuous (Lipshitz).

$\qquad$ Hence, $\quad X \approx \tilde{X}$ implies $Z \approx \tilde{Z} = f(\tilde{X})$

$\qquad\quad X \;\approx\; \bar{X} \qquad\quad \text{fluid} \qquad\quad \Rightarrow \; \bar{Z} = f(\bar{X}) \qquad\quad \text{fluid approximations.}$

$\qquad\quad X \;\approx\; \bar{X} + \hat{X} \; \text{diffusion} \quad \Rightarrow \; \hat{Z} = f(\bar{X} + \hat{X}) \;\; \text{diffusion refinements.}$

*Reference*: Harrison, Chapter 2 (which covers also finite buffers, and two-node networks).