

# Approximations / Hierarchical Modelling of Stochastic Networks

## Hierarchy

Micro models ex. Queueing networks

Macro Fluid nets

Meso Diffusion nets

## Approximations (strong: Hungarian)

F-approx : 1<sup>st</sup> order, via FSLLN

D-approx : 2<sup>nd</sup> order, FCLT

Framework DCP (Oblique Reflection)

# SCALING and CENTERING of DATA

Here: Visual + intuitive convergence.

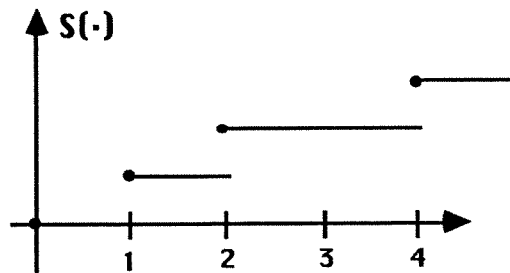
(Rigorous: Weak convergence, or convergence of stochastic processes *in distribution*).

Bernoulli process:

$$\begin{aligned}\Delta_n &= 2 \quad \text{wp } 1/2 \quad n = 1, 2, \dots, \text{ iid} \\ &= 0 \quad \text{wp } 1/2 .\end{aligned}$$

Define

$$\begin{aligned}S_n &= \Delta_1 + \Delta_2 + \dots + \Delta_n, \quad n \geq 1, \\ S(t) &= \sum_{k=1}^{\lfloor t \rfloor} \Delta_k \quad (= S_{\lfloor t \rfloor}), \quad t \geq 0:\end{aligned}$$



- Fit  $\{S(t), 0 \leq t \leq n\}$ ,  $n$  large, into a  $1 \times 1$  frame (a computer screen).

What to expect?

## FLUID and DIFFUSION Approximations (Donsker's Theorem)

$$\Delta_n = \begin{cases} 2 & 1/2 \\ wp & \\ 0 & 1/2 \end{cases} ; S_n = \sum_{k=1}^n \Delta_k , S(t) = \sum_{k=1}^{\lfloor t \rfloor} \Delta_k , t \geq 0.$$

Fluid approximation:  $S(t) \approx t, t \geq 0$  ; supported by

$$\text{FSLLN } \frac{1}{n} S(nt) \xrightarrow{wp1} t, t \geq 0.$$

Diffusion refinement:  $S(t) \approx t + B(t), t \geq 0$ , supported by

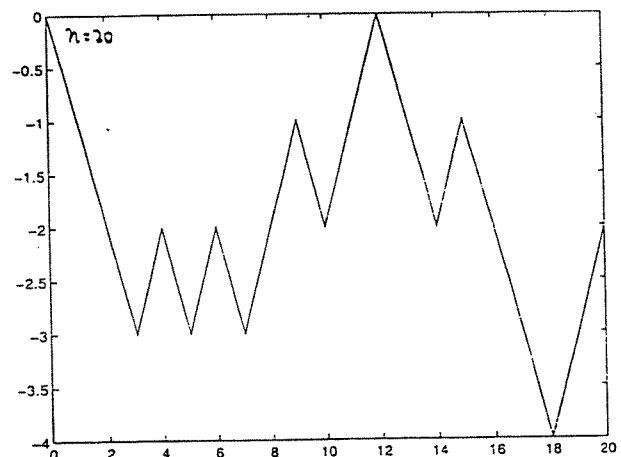
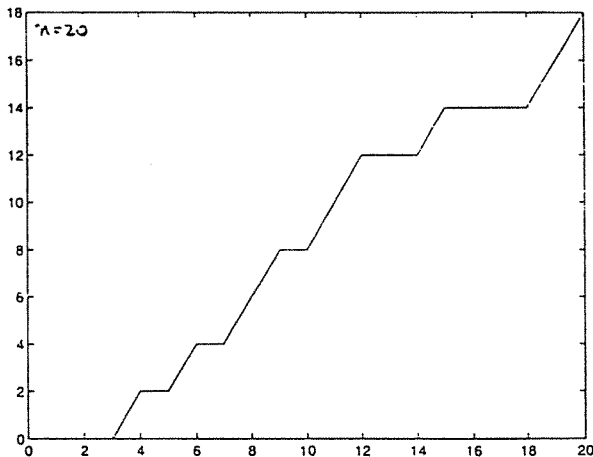
$$\text{FCLT } \sqrt{n} \left[ \frac{1}{n} S(nt) - t \right] \xrightarrow{d} B(t), t \geq 0,$$

where  $B = \{B(t), t \geq 0\}$  is Levy with normal (Gaussian) marginals, indeed SBM. □

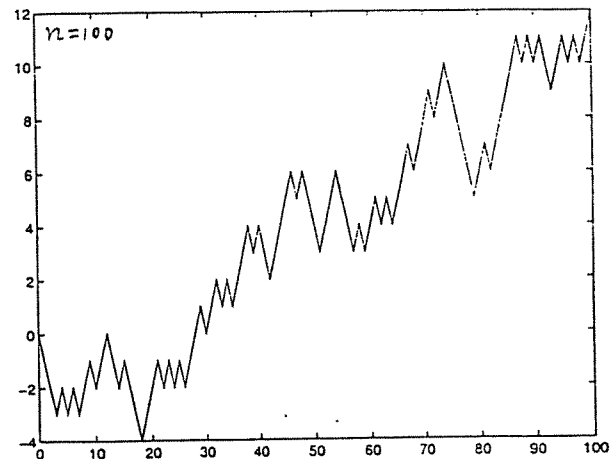
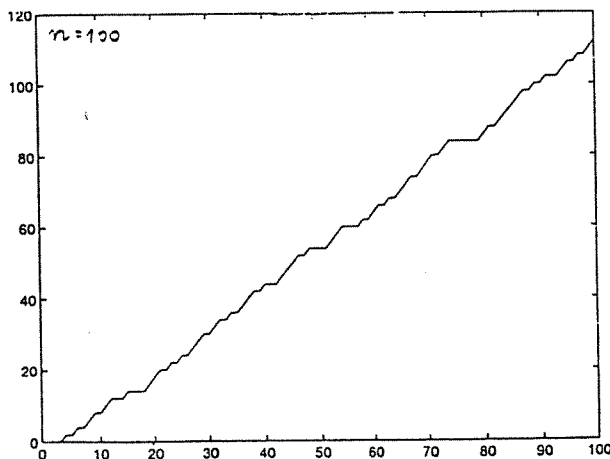
21

Invariance Principle:  $\{\Delta_k\}$  mean  $\mu$ , std.  $\sigma$ , then  $S(t) \approx \mu t + \sigma B(t), t \geq 0$ .

20



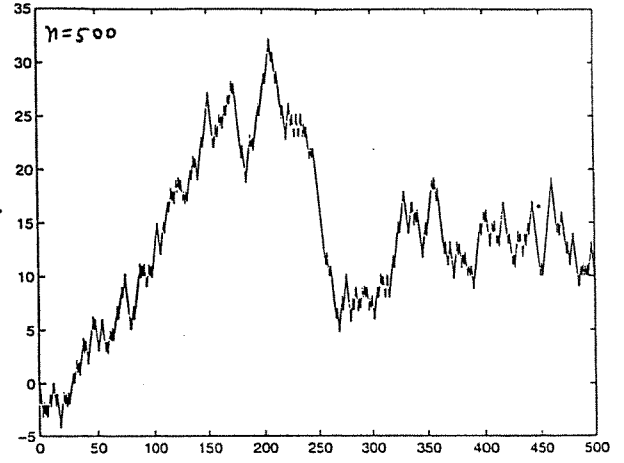
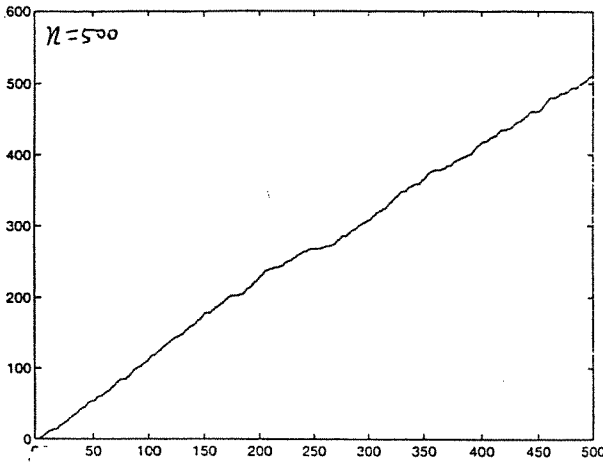
100



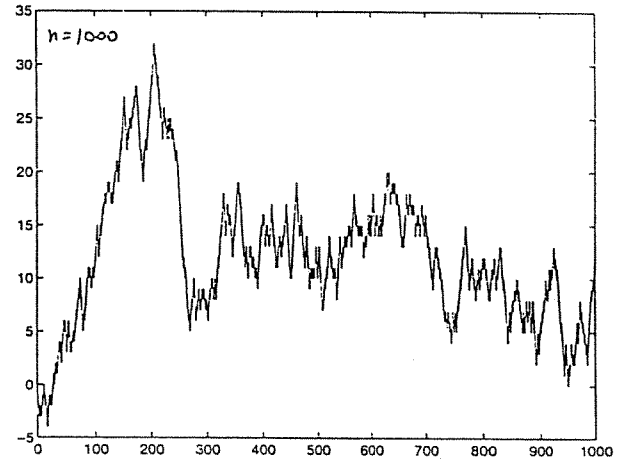
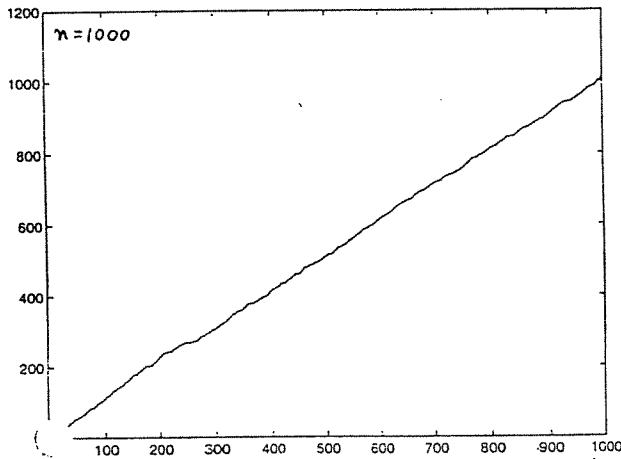
*Fluid approximation: Deterministic  
Functional  
Strong Law of Large Numbers (FSLLN)*

*Diffusion deviation: Stochastic  
Functional  
Central Limit Theorem (FCLT)*

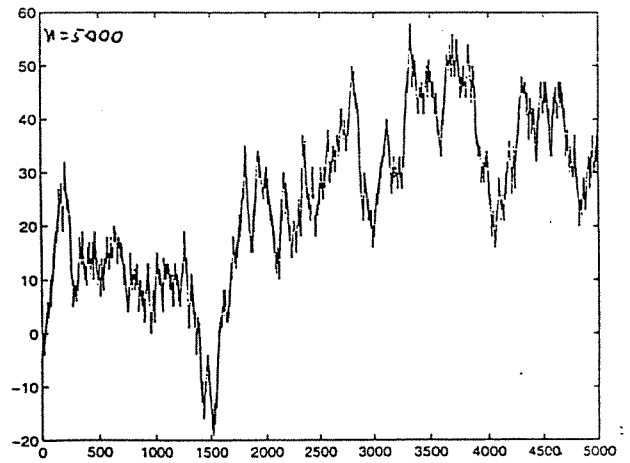
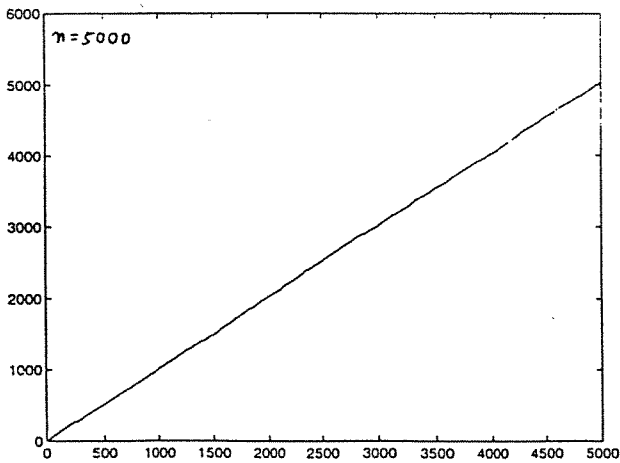
500



1000



5000



Interpolating  $S(0), S(1), \dots, S(n)$

Interpolating deviations from trend

## SLLN (Strong Law of Large Numbers)

$$\frac{1}{n}S(n) \xrightarrow{wp1} E\Delta_1 = 1, n \uparrow \infty.$$

1. **Rescaling time:**  $\{S(nt), 0 \leq t \leq 1\}$ ;  
at  $t = 1$ :  $ES(n) \approx n$ ;

2. **Aggregating space:**  $\{\frac{1}{n}S(nt), t \geq 0\}$   
 $\xrightarrow{wp1} \{t, t \geq 0\}, n \uparrow \infty.$

(But recall the Cauchy distribution)

## CLT (Central Limit Theorem)

$$\sqrt{n}\left[\frac{1}{n}S(n) - 1\right] \xrightarrow{d} N(0, 1), \quad n \uparrow \infty.$$

3. **Centering**  $\frac{1}{n}S(nt) - t$ : deviations from trend

4. **Amplifying**  $\sqrt{n}\left[\frac{1}{n}S(nt) - t\right] \xrightarrow{d} N(0, t), \quad n \uparrow \infty.$

Obtained:

$$\begin{aligned} S(nt) &\approx n\left[t + \frac{1}{\sqrt{n}}N(0, t)\right] = nt + \sqrt{n}N(0, t) \\ &\stackrel{d}{=} nt + N(0, nt) \end{aligned}$$

Expect  $S(t) \stackrel{d}{\approx} t + B(t), \quad t \geq 0,$

where  $B = \{B(t), t \geq 0\}$  is a **Levy process**

with normal Gaussian marginals, namely a

**Brownian Motion**  $BM(0, 1) = SBM$  (= Standard

Brownian Motion).

## Invariance Principle (Donsker's Theorem)

Let  $\Delta_1, \Delta_2, \dots$  be *iid* with mean =  $\mu$  and variance =  $\sigma^2$   
(*general* distribution).

For  $S = \{S(t), t \geq 0\}$  as before, we have

$$S(t) \stackrel{d}{\approx} \mu \cdot t + \sigma B(t), \quad t \geq 0,$$

where  $B = \{B(t), t \geq 0\}$  is an SBM.

Equivalently,  $S \stackrel{d}{\approx} BM(\mu, \sigma^2)$ .

**Terminology:**  $\mathcal{S}(t)$

$\bar{\mathcal{S}}(t) = \mu t$       fluid approximation (deterministic);

$\hat{\mathcal{S}}(t) = \sigma B(t)$       diffusion refinement (stochastic).

- **FSLLN** (Functional SLLN)

$$\bar{S}^n(t) = \frac{1}{n}S(nt) \longrightarrow \bar{S}(t), \quad t \geq 0.$$

Convergence wp1 (strong) to a deterministic fluid limit.

- **FCLT** (Donsker's Theorem)

$$\hat{S}^n(t) = \sqrt{n} \left[ \frac{1}{n}S(nt) - \mu t \right] \Longrightarrow \hat{S}(t), \quad t \geq 0.$$

Convergence in distribution (weak) to a stochastic diffusion limit.

- **Strong Approximation** (Kolmos, Major, Tusnady)

Assume moments. Then  $\exists$  prob. space supporting

$B = SBM$  and  $\tilde{S}$  such that

1.  $\tilde{S} \stackrel{d}{=} S$

2.  $\sup_{0 \leq t \leq \tau} |\tilde{S}(t) - (\mu t + \sigma B(t))| = o(\sqrt{\tau}), \quad \tau \uparrow \infty$

Note: This implies both FSLLN and FCLT

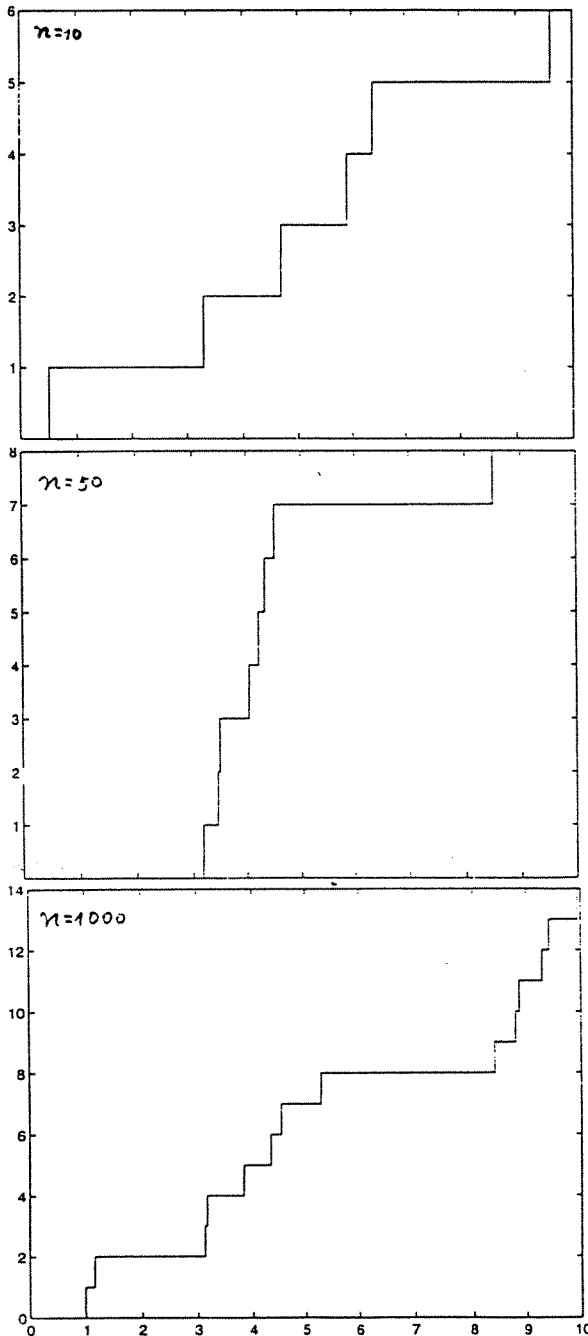
(but requires "too many" moments).



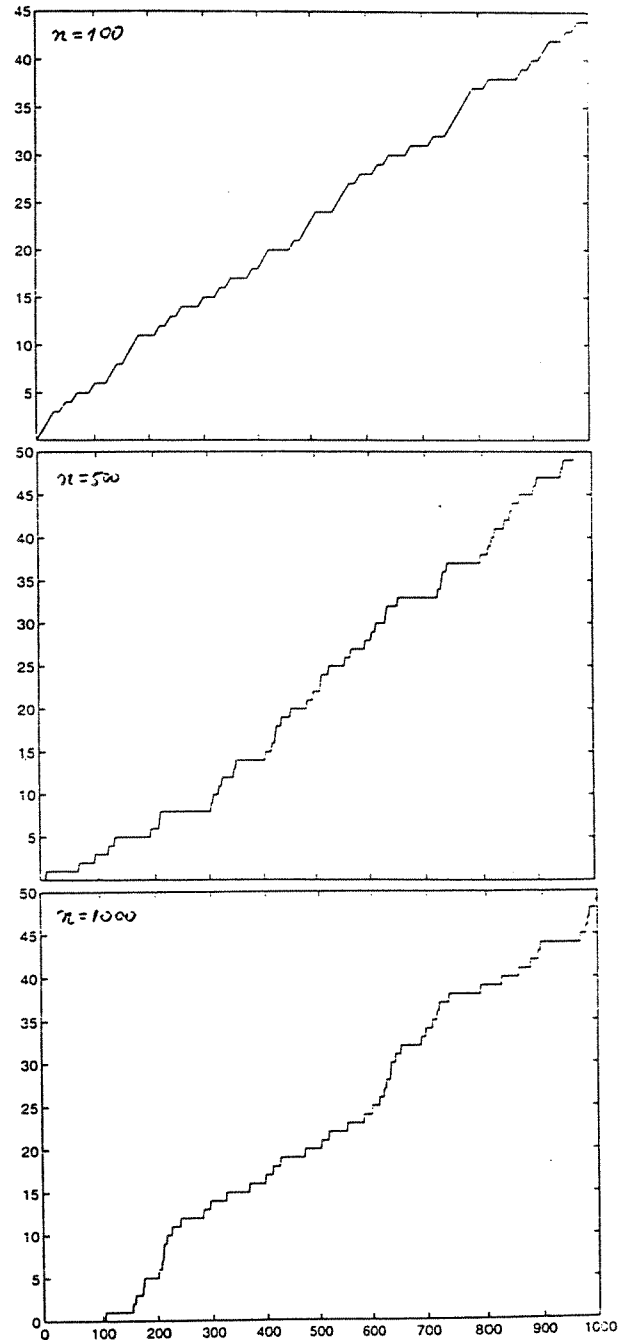
## Bernoulli $\Rightarrow$ Poisson

$$\forall n, iid \quad \Delta_k^n = \begin{cases} 1 & \text{wp } p^n = \frac{\lambda}{n} \\ 0 & \text{wp } q^n = 1 - p^n \end{cases}, \quad S^n(t) = \sum_{k=1}^{\lfloor t \rfloor} \Delta_k^n, \quad t \geq 0.$$

Rare Events:  $S^n(nt) \sim \text{Binomial}(\lfloor nt \rfloor, \frac{\lambda}{n}) \Rightarrow \text{Poisson}(\lambda t), t \geq 0.$



$S^n(nt), 0 \leq t \leq 10; \lambda = 1.$



$S^n(t), 0 \leq t \leq n; \lambda = 50.$

# Stochastic-Process Limits

## An Introduction to Stochastic-Process Limits And their Application to Queues

Ward Whitt

AT&T Labs - Research  
The Shannon Laboratory  
Florham Park, New Jersey

Draft  
June 13, 2001

Copyright ©info

Marginals

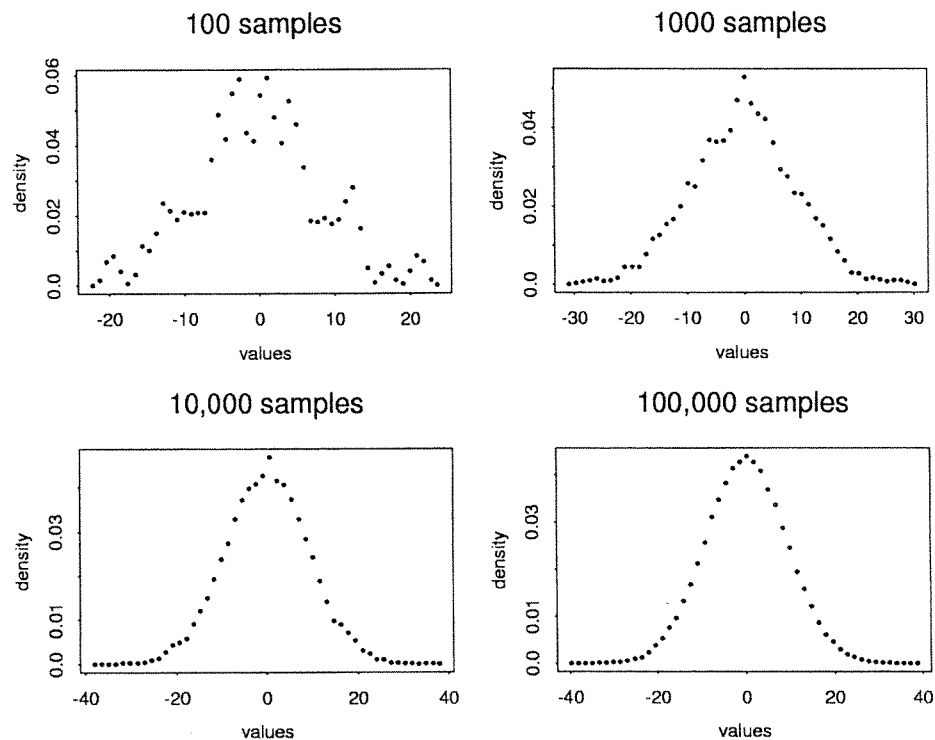


Figure 1.7: Estimates of the probability density of the final position of the random walk, obtained from  $10^j$  independent samples of the centered partial sum  $S_{1000} - 500$  for  $j = 2, \dots, 5$ , for the case in which the steps  $U_k$  are uniformly distributed in the interval  $[0, 1]$ , based on the nonparametric density estimator *density* from *S*.

density estimates converge to a normal pdf as  $n \rightarrow \infty$ .

It is not our purpose to delve deeply into statistical issues, but it is worth remarking that we obtain new interesting plots, like the random walk plots, when we do. Our brief examination of the distribution of the final position of the random walk suggests looking for a more precise statistical test to determine whether or not the final position of the random walk is indeed approximately normally distributed. To evaluate whether some data can be regarded as an independent sample from any specified probability distribution, it is natural to carefully investigate how the empirical distribution of a sample from that probability distribution tends to differ from the underlying probability distribution itself.

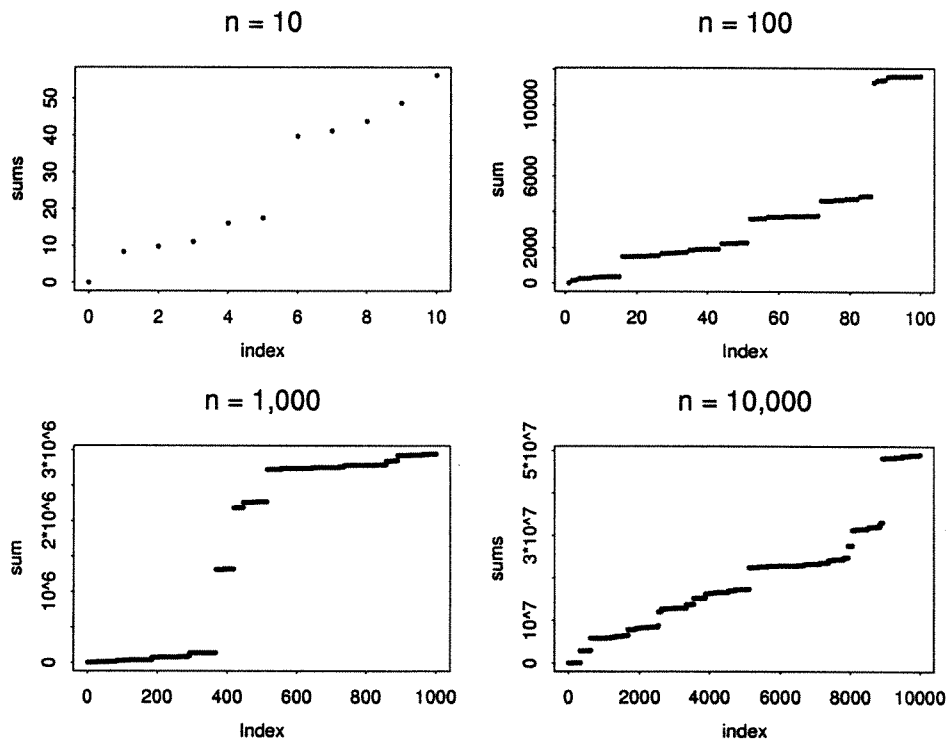


Figure 1.19: Possible realizations of the first  $10^j$  steps of the uncentered random walk  $\{S_k : k \geq 0\}$  with steps distributed as  $U_k^{-1/p}$  in case (iii) of (3.5) for  $p = 1/2$  and  $j = 1, \dots, 4$ .

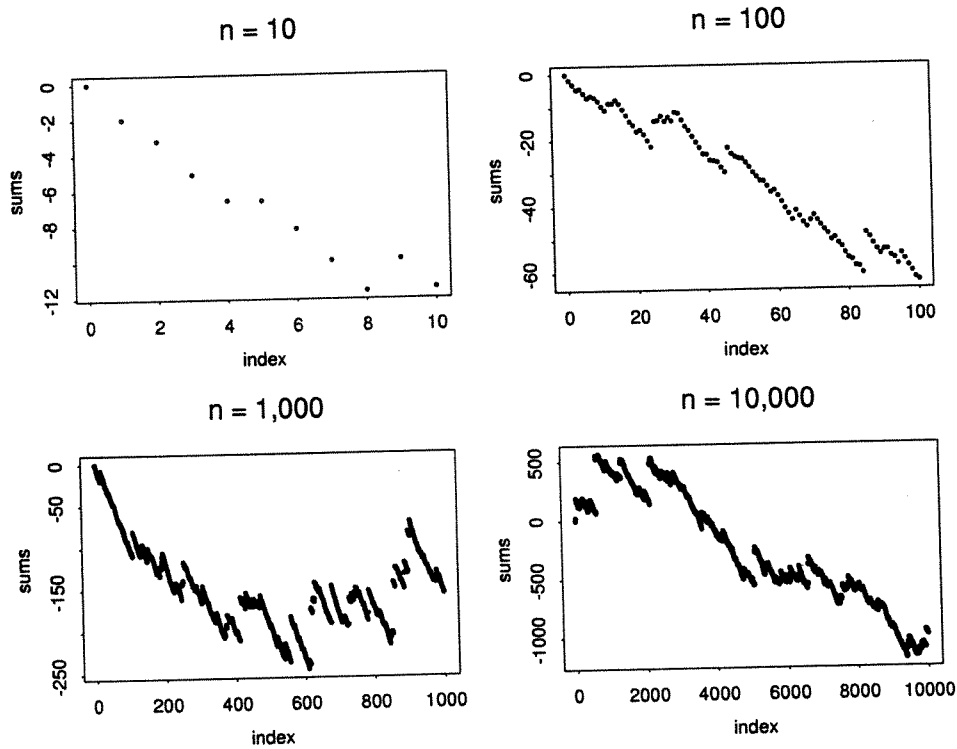


Figure 1.20: Possible realizations of the first  $10^j$  steps of the centered random walk  $\{S_k - 3k : k \geq 0\}$  associated with the Pareto steps  $U_k^{-1/p}$  for  $p = 3/2$ , having mean 3 and infinite variance, for the cases  $j = 1, \dots, 4$ .

1.20. As before, the centering causes the plotter to automatically blow up the picture. However, now the slight departures from linearity for large  $n$  in Figure 1.18 are magnified. Now, just as in Figure 1.19, we see jumps in the plot!

Once again, probability theory offers an explanation. Just as the SLLN ceases to apply when the IID summands have infinite mean, so does the (classical) CLT cease to apply when the IID summands have finite mean but infinite variance. Such a case occurs with the Pareto( $p$ ) summands in case (iii) in (3.5) when  $1 < p \leq 2$ . Thus, consistent with what we see in Figure 1.18, the SLLN holds, but the CLT does not, for the Pareto( $p$ ) random variable  $U^{-1/p}$  in case (iii) when  $p = 3/2$ .

We have arrived at another critical point, where an important intellectual step is needed. We need to recognize that, *even though the sample paths are*

## Dynamic Randomness: The Poisson Process

Hall, Chapter 3: The *Arrival* Process

( with the framework of Levy Processes )

Review:

- Introduction to Services and Queues (Service Nets = Queueing Nets).

Why queues? Scarce resources, coordination gaps/design constraints.  
There is a "future for queues" in service systems.

↳ Poisson  
↳ Brownian  
↳ Gamma

- Measurements; Empirical Models I; Scenario Analysis.

Little's Law  $L = \lambda W$  &  $H = \lambda G$ ;  
Capacity analysis (via first order data);  
There is a "future for models" in service systems.

- The Processing Network Paradigm (ReEngineering = BPR)  
via Dynamic Stochastic Project Networks (DSP-nets).

- Can we do it? Bottleneck analysis  $\leftrightarrow$  the fluid view;
- How long will it take? Typically stochastic networks;
- Can we do better? Parametric/sensitivity/what-if.

- A Deterministic Service Station; Empirical Models II

Skorohod's model (cumulatives a necessity):  $z = x + y$   
 $z \geq 0, y \uparrow 0$   
 $zdy = 0$

Modelling scope: Lindley's equation  
Workload

Towards Stochastic Service Station:

arrivals' epochs;  
service durations.

Today: Model for Completely Random Arrivals.

# Approximating Stochastic Networks

$$Z = f(X) \quad X \text{ Point Process, Levy, ...}$$

The mapping  $f$  continuous / Lipschitz

$$\Rightarrow Z \approx f(X^\circ), \quad X^\circ \approx X \quad \text{simpler}$$

Fluid approx.  $X^\circ \in BV$  Bounded Variation

Diffusion approx.  $X^\circ \in BM$  Brownian Motion

Both

Long-run  $Z \approx f(X^n), n \uparrow \infty$  horizon

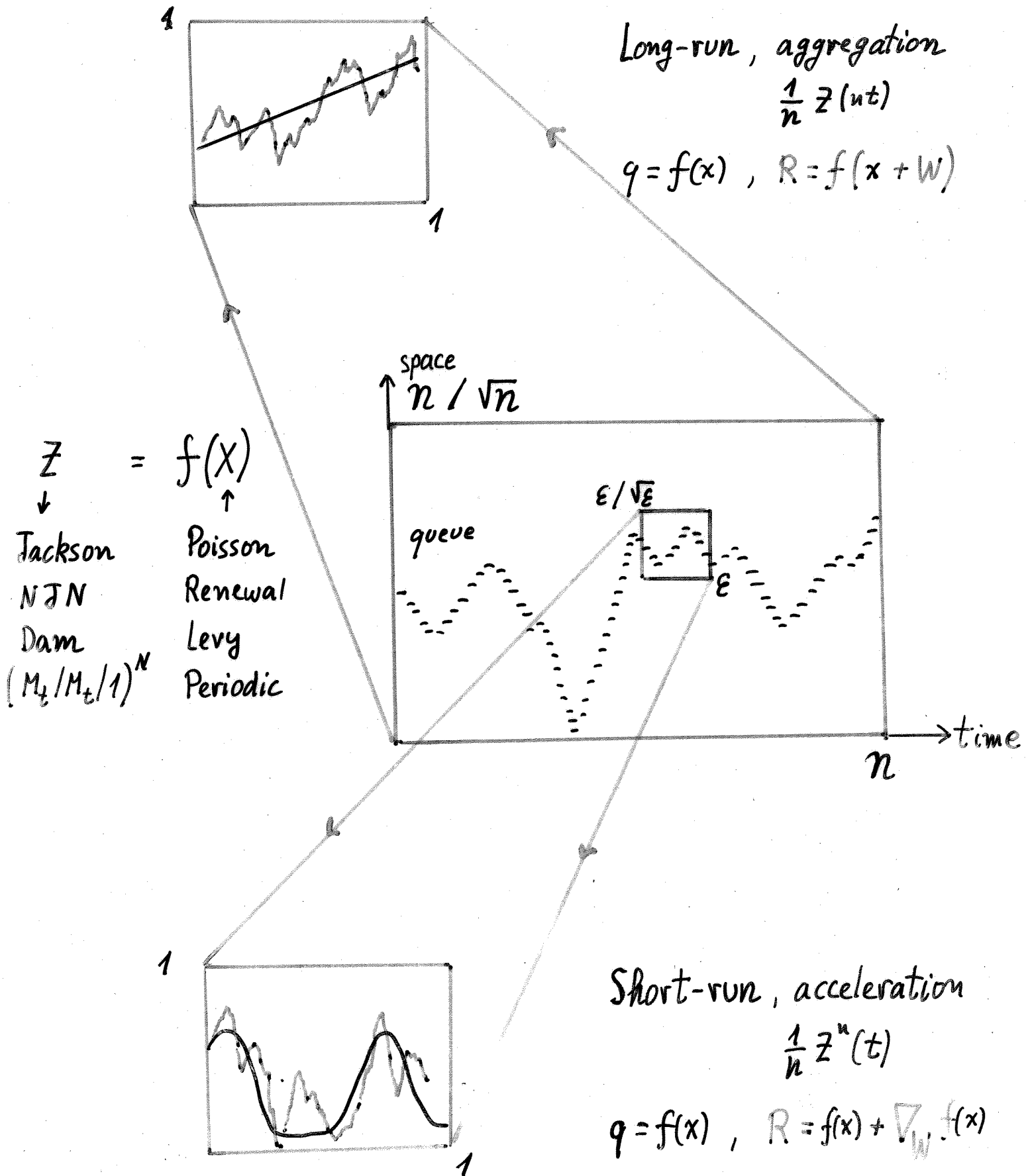
or strong approx., weak converg.

Short-run  $Z \approx f(X^\epsilon), \epsilon \downarrow 0$  neighbourhood

acceleration, perturbation

$\Rightarrow$  Five-level Hierarchical Framework

# Hierarchical Framework





# Hierarchical Modelling: a Framework

Ex.  $Q$ -length, open periodic  $Q$ -net  $Q = f(X)$

Strategic compress long-range evolution of network  
into short-range horizon of observer

$$\bar{Q} = f(\bar{X})$$

$f$  macro  
bottlenecks, stability

$$\hat{Q} = f(\hat{X})$$

$d$  meso  
sojourn times

Tactical

$$Q = f(X)$$

$q$ -net micro

$$\hat{Q} = f(\hat{X})$$

$d$ -net meso  
envelopes

$$\bar{Q} = f(\bar{X})$$

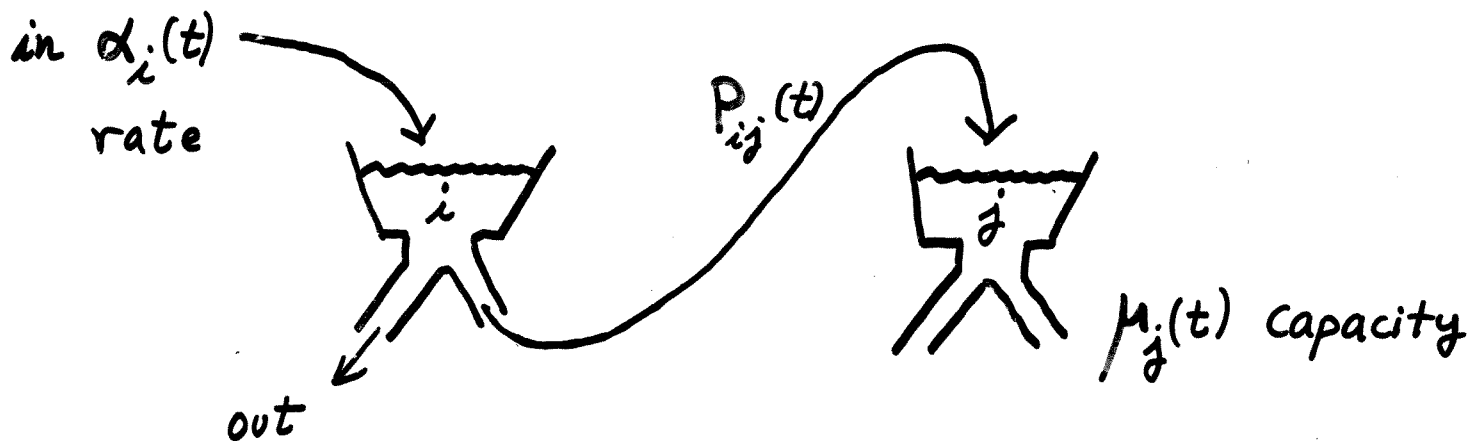
$f$ -net macro  
rush-hour

Operational amplify short-range horizon of observer  
to long-range evolution of a virtual net

Rescaling  $\rightarrow$  limit theorems (state-space collapse)

Multiple scales

# A Fluid Network : Animating $f$



Outflow rates  $\delta_j(t) \leq \mu_j(t)$

Efficient "  $< \mu_j(t) \Rightarrow z_j(t) = 0$

Inflow  $\lambda_j = \alpha_j + \sum_i \delta_i P_{ij}$

Content  $z = z(0) + \int \lambda - \int \sigma = f(X)$

The Mapping  $\nearrow$  data ( $\alpha, p, P$ )

Reflection

Regulator

DCP( $X$ )

$$f \begin{cases} z(t) = X(t) + \int_0^t d\gamma(s) [I - P(s)] , t \geq 0 \\ z \geq 0 , \gamma \uparrow 0 , z \cdot d\gamma = 0 \end{cases}$$

$\nwarrow$  eff.

where  $X = z(0) + \int (\lambda + pP) - \int \mu$  netflow

$\gamma = \int \mu - \int \sigma$  cum. lost capacity

$$f: X \rightarrow Z$$

Outflow

$$\delta_j(t)$$

Inflow

$$\lambda_j(t) = \alpha_j(t) + \sum_i \delta_i(t) P_{ij}(t)$$

Eg. constant rates  $\alpha(t) \equiv \alpha, \mu, P$  : Linear  $X$

$\Rightarrow$  Equilibrium  $\lambda_j = \alpha_j + \sum_i (\lambda_i \wedge \mu_i) P_{ij}$  Traffic eq.

Traffic intensity  $\rho_j = \frac{\lambda_j}{\mu_j} \begin{cases} > 1 & \text{Bottleneck} \\ = 1 & \text{critical} \\ < 1 & \text{"stable"} \end{cases}$

General (eg. periodic)

~~$$\frac{\lambda_j(t)}{\mu_j(t)}$$~~

$$\rho_j(t) = \sup_{0 \leq s \leq t} \frac{\int_s^t \lambda_j(u) du}{\int_s^t \mu_j(u) du}$$

$\rho_j(t) > 1 \Leftrightarrow Z_j(t) > 0$  super-critical

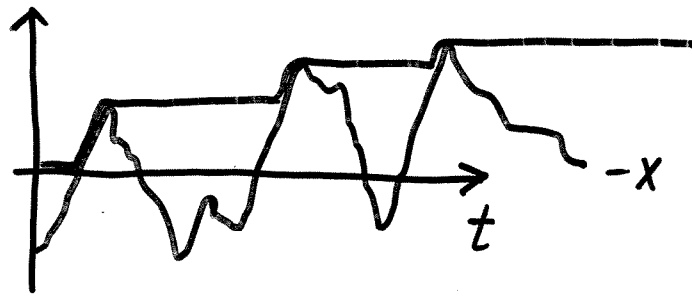
$= 1 \Leftrightarrow Z_j(t) = 0, \delta_j(t) = \mu_j(t)$  critical

$< 1 \Leftrightarrow Z_j(t) = 0, \delta_j(t) < \mu_j(t)$  sub-critical

Phases Time-dependent 19

# Geometric Interpretation: Oblique Reflection

Single buffer  $Z = X + Y \gg 0, Y \uparrow 0, Z dy = 0$

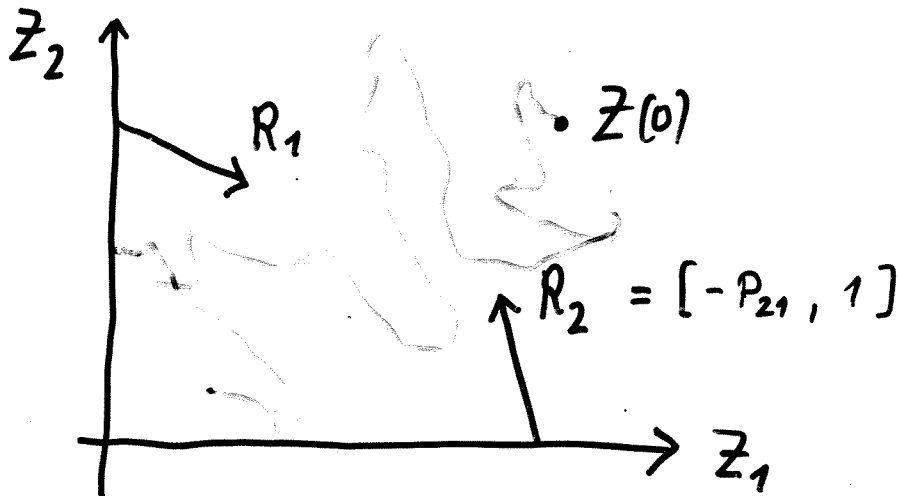


$$Y = \overline{(-X)^+}$$

$$= \underline{-X} \text{ when } X(0) = 0$$

(Skorohod)

Two buffers  $Z = X + Y [I - P] = X + Y_1 R_1 + Y_2 R_2$



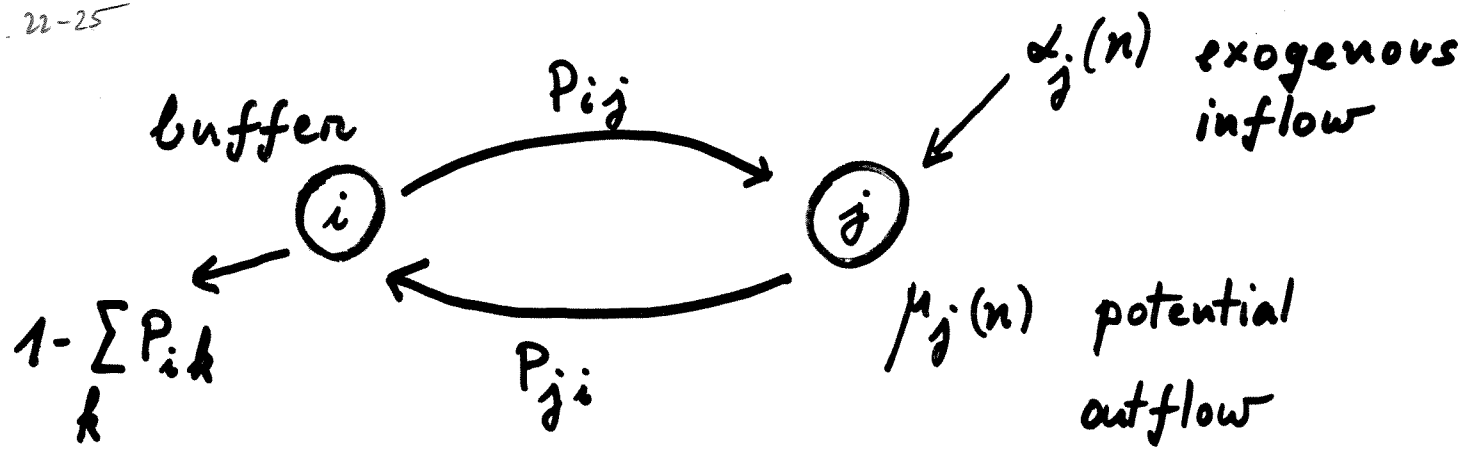
(Harrison + Reiman)

General  $f: X \rightarrow Z$  Lipschitz

The detailed version

# A Flow Network in Discrete Time

pg. 22-25



At end of periods  $n = 1, 2, \dots$

Content 
$$z_j(n) = z_j(n-1) + \alpha_j(n) - \delta_j(n) + \sum_i \delta_i(n) P_{ij}$$

Actual outflow  $\delta_j(n) \leq \mu_j(n)$  : efficient?

Data  $z(0) = (z_j(0)), \alpha = \{\alpha_j(n)\}, \mu, P = [P_{ij}]$

Balance eq.  $z(n) = z(n-1) + \alpha(n) + \delta(n)[P - I]$

Feasible  $z \geq 0, 0 \leq \delta \leq \mu$  ; efficient?

Open spec  $\text{rad}(P) < 1$  ; Closed P stoch. +  $\alpha \geq 0$

# Efficient Operation

Content 
$$z(n) = z(n-1) + \alpha(n) + \delta(n)[P-I]$$

lost potential 
$$\eta_j(n) = \mu_j(n) - \delta_j(n) \geq 0$$

$$z(n) = z(n-1) + \underbrace{\{\alpha(n) + \mu(n)[P-I]\}}_{\dot{z}(n) \text{ data}} + \underbrace{\eta(n)[I-P]}_{\text{control}}$$

$$\begin{cases} z(n) = z(n-1) + \dot{z}(n) + \eta(n)[I-P] \\ z \geq 0, \eta \geq 0 \end{cases} \text{ feasible}$$

Fact:  $\exists$  least feasible  $\eta^*$

potential lost only when buffer empty

$$\eta_j^*(n) > 0 \Rightarrow z_j^*(n) = 0$$

Complementarity 
$$z_j^*(n) \cdot \eta_j^*(n) = 0 \quad \text{least } \eta$$

# Geometric Interpretation : Discrete Time

$$z(n) = z(n-1) + f(n) + \underbrace{\eta(n) [I - P]}_{\eta_1(n)R_1 + \eta_2(n)R_2 + \dots}, n=1,2,\dots$$

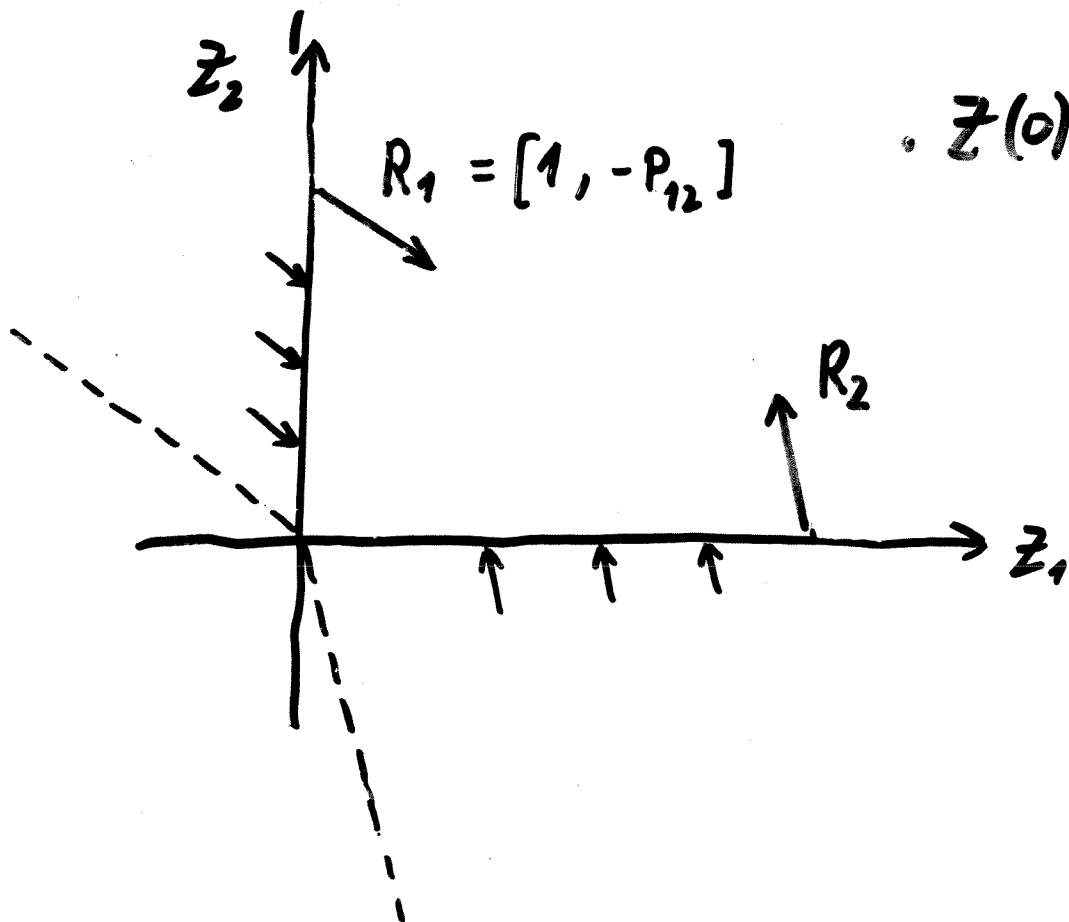
Sequentially

$z(0), z(1), z(2), \dots$

$\eta(1) \quad \eta(2) \quad \dots$  least (no \*)

$$\eta_1(n)R_1 + \eta_2(n)R_2 + \dots$$

↑  
2<sup>nd</sup> row



Indeed  $z_j(n) \cdot \eta_j(n) = 0$  LCP

# Continuous Time

$$\begin{cases} z(n) - z(n-1) = f(n) + \eta(n)[I-P] & , n=1,2,\dots \\ z \gg 0, \eta \gg 0 & \text{feasible} \\ z_j \eta_j = 0 & \text{least} \end{cases}$$

$$\text{cum. } z(t) = z(0) + \underbrace{\sum_{n \leq t} f(n)}_{X(t) \text{ data}} + \underbrace{\sum_{n \leq t} \eta(n)[I-P]}_{Y(t) \text{ control}}$$

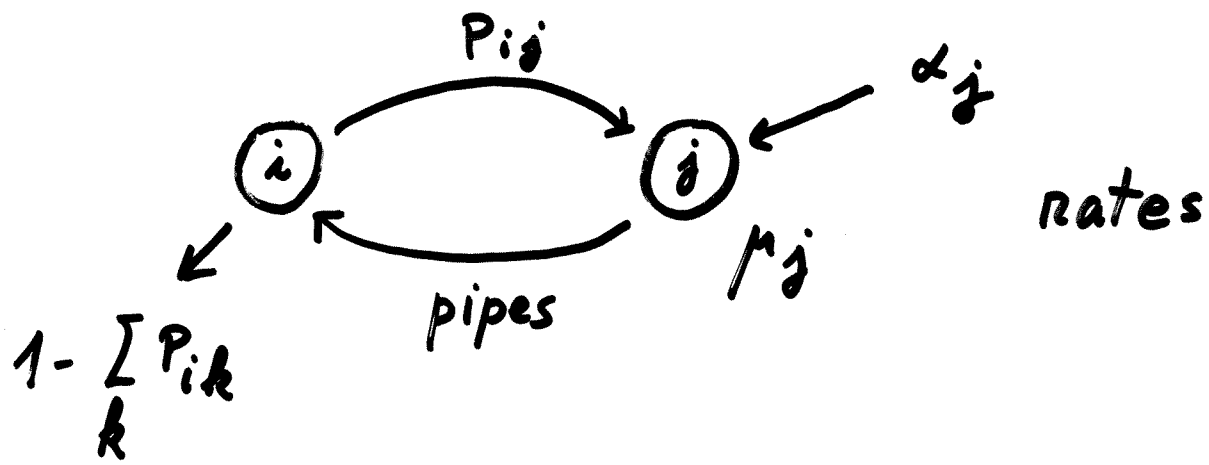
DCP Given  $X = \{X(t), t \gg 0\} \in \text{RCLL}$   
 find  $y, z \ni$

$$\begin{cases} z(t) = X(t) + Y(t)[I-P] & , t \gg 0 \\ z \gg 0, Y \uparrow 0 & \text{feasible} \\ z_j dy_j = 0 & \text{least } (dy_j(t) > 0 \Rightarrow z_j(t) = 0) \end{cases}$$

DCP animated by a flow network  $\Leftrightarrow X \in BV$  24



# Linear X: Fluid Network



Data	$Z_j(0)$	initial content	$Z(0)$
	$\alpha_j$	exogenous inflow rates	$\alpha$
	$\mu_j$	potential outflow rates	$\mu$
	$P_{ij}$	fractions	$P$

Most efficient operation via DCP with data

$$X(t) = Z(0) + \theta t, \quad \theta = \alpha + \mu[P - I]$$

Equilibrium inflow rates : solve traffic eq.

$$d_j = \alpha_j + \sum_i (d_i \mu_i) P_{ij} \begin{cases} < \mu_j & \text{non-bottlenecks} \\ = \mu_j & \text{balanced} \\ > \mu_j & \text{strict bottlenecks} \end{cases}$$

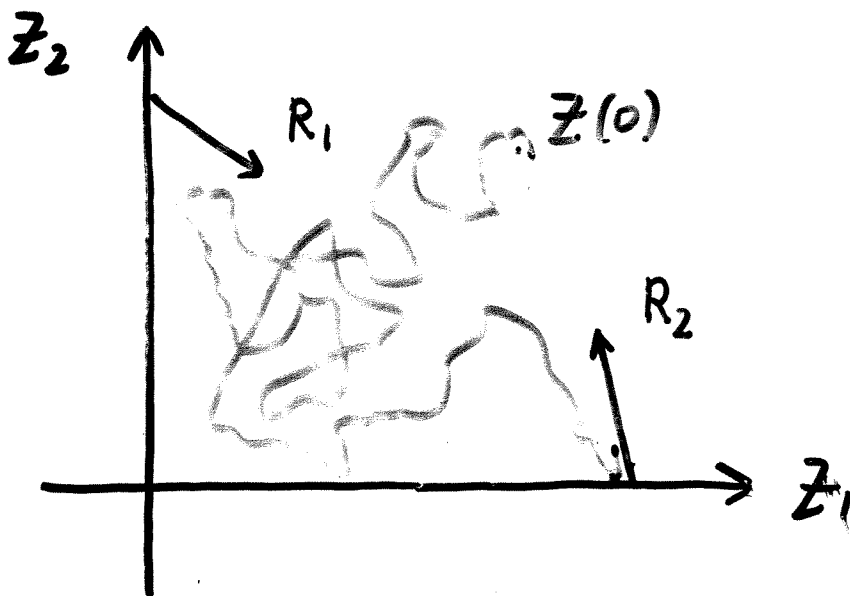
$$d = \alpha + (d \mu) P$$

# Geometric Interpretation: Continuous Time

$$Z(t) = X(t) + \underbrace{Y(t)}_{Y_1(t)R_1 + Y_2(t)R_2 + \dots} [I - P]$$

$$Y_1(t)R_1 + Y_2(t)R_2 + \dots$$

$X$  continuous  $\Rightarrow Y, Z$  as well



$dY_1(t) > 0 \Rightarrow Z_1(t) = 0$ , reflect in direction  $R_1$

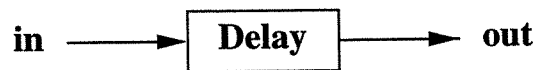
$f: X \rightarrow Z$  oblique reflection

DCP:  $X \in RCLL$

Lipshitz

**A Deterministic Model of a Service Station (Fluid View)****Primitives**

$Z(0)$	initial content
$\alpha(t)$	input rate (inflow rate)
$\mu(t)$	potential service rate

**Model: (Think cumulants)**

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

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

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

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

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

$$\Rightarrow \text{Outflow: } D = M - Y \quad (\mathbf{A} \text{ arrivals; } \mathbf{D} \text{ departures})$$

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

$$\text{Model} \quad Z = X + Y$$

$$\text{Feasible} \quad Z \geq 0, Y \uparrow 0 \quad (Y(0) = 0);$$

$$\text{Efficient} \quad Y \text{ least} \quad (\text{hence, } Y \text{ unique});$$

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

$$\underline{X}(t) = \inf_{0 \leq u \leq t} X(u), \text{ which is called the lower envelope of } X; \text{ similarly, } \bar{X} \text{ upper envelope.})$$

# THE DYNAMIC COMPLEMENTARITY PROBLEM

Running Head: *DCP*

Avi Mandelbaum †

Graduate School of Business  
Stanford University  
Stanford, Cal. 94305, USA

and

Industrial Engineering and Management  
Technion – Israel Institute of Technology  
Haifa, Israel

**Abstract** A dynamic version of the Linear Complementarity Problem (*LCP*) is proposed. We call it the Dynamic Complementarity Problem (*DCP*). There are versions of *DCP* both in discrete and continuous time. (The latter, which constitutes a natural limit of the former, was actually conceived first; it emerged as a framework for fluid and diffusion approximations of stochastic queueing networks.) Existence and uniqueness results for discrete-time *DCP*'s are precisely analogous to their classical *LCP* counterparts. In contrast, continuous-time *DCP*'s exhibit  $\mathcal{P}$ -matrices for which uniqueness fails to hold. Such examples will be constructed by introducing and solving a certain system of differential inequalities. This system implicitly characterizes the matrices for which uniqueness prevails, but their explicit characterization remains a challenging open problem.

**Keywords and phrases:** Linear and non-linear complementarity;  $\mathcal{Q}$ -matrices,  $\mathcal{P}$ -matrices,  $\mathcal{M}$ -matrices; Reflected/Regulated stochastic processes; Convex processes; Brownian networks, Diffusion and fluid approximations; Elasto-plastic models.

First Version: June 1987

Revised: October 1989

( Accepted to MOR )

---

*AMS 1980 subject classifications (1985 Revision):* Primary 90C33. Secondary 90B22, 60F15, 60F17, 60K20, 90C48.

*OR/MS Index subject classifications:* Primary Programming: Complementarity, Infinite dimensional; Queues: Diffusion models, Networks; Secondary Mathematics: Matrices; Networks: Stochastic.

† Research supported in part by the Stanford Business School Trust Fund, the Israeli Bat-Sheva Foundation and the Fund for Promotion of Research at the Technion.

# Derivative of Reflection

Massey, Pats  
Ramanan, Whitt

Fluid  $\frac{1}{n} X^n \xrightarrow{\text{a.a.}} \bar{X}$

$$X^n \sim n \bar{X}$$

Diffusion  $\sqrt{n} \left( \frac{1}{n} X^n - \bar{X} \right) \xrightarrow{d} \hat{X}$

$$X^n \approx n \bar{X} + \sqrt{n} \hat{X}$$

Consider  $Z^n = f(X^n)$

Fluid  $\frac{1}{n} Z^n = \frac{1}{n} f(X^n) \underset{\substack{\uparrow \\ f \text{ homog.}}}{=} f\left(\frac{1}{n} X^n\right) \underset{\substack{\uparrow \\ f \text{ cont}}}{\rightarrow} f(\bar{X}) = \bar{Z}$

$$Z^n \sim n f(\bar{X}) = n \bar{Z}$$

Diffusion  $\sqrt{n} \left( \frac{1}{n} Z^n - \bar{Z} \right) = \sqrt{n} \left[ f\left(\frac{1}{n} X^n\right) - f(\bar{X}) \right]$

$$\approx \sqrt{n} \left[ f\left(\bar{X} + \frac{1}{\sqrt{n}} \hat{X}\right) - f(\bar{X}) \right]$$

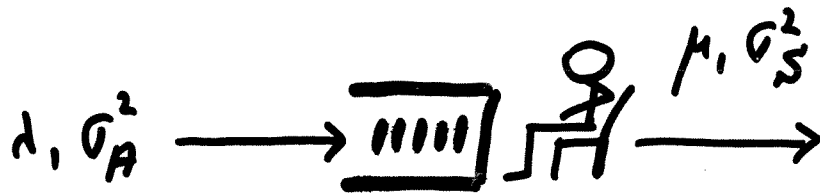
$$= \frac{1}{\varepsilon} \left[ f(\bar{X} + \varepsilon \hat{X}) - f(\bar{X}) \right] \Rightarrow \nabla_{\hat{X}} f(\bar{X}) = \hat{Z}$$

$$Z^n \approx n f(\bar{X}) + \sqrt{n} \nabla_{\hat{X}} f(\bar{X}) = n \bar{Z} + \sqrt{n} \hat{Z}$$

Subtlety:  $\Rightarrow \nabla_{\hat{X}} f(\bar{X})$  in  $M_q$

Common:  $\bar{X} = 0, f(0) = 0 \Rightarrow \bar{Z} = 0, \hat{Z} = f(\hat{X})$  Critical

# The GI/GI/1 Queue



Arrival process  $A = \{A(t), t \geq 0\}$

Renewal process: iid interarrival times

$$\text{mean} = \frac{1}{\lambda}, \text{var} = \sigma_A^2 \quad (C_A^2 = \lambda^2 \sigma_A^2)$$

Service times  $S_1, S_2, \dots$

iid services:  $\frac{1}{\mu}, \sigma_S^2 \quad (C_S = \mu \sigma_S)$

$$L(t) = \sum_{i=1}^{A(t)} S_i \quad \text{offered load}$$

$$X(t) = L(t) - t \quad \text{netflow}$$

$$V = X - \underline{X} \quad \text{unfinished work}$$

(virtual waiting time)

$V$  obtained from  $X$  via reflection: RBV

(RBM)

# G/G/1 : Strong Approximations

$$S(t) = \sum_1^t S_i \approx \frac{1}{\mu} t + \sigma_S B_S(t) \quad \text{Donsker}$$

$$A(t) \approx \lambda t + \lambda^{3/2} \sigma_A B_A(t) \quad \text{Inverse D.}$$

$$L(t) = S[A(t)] \quad \text{Time-change}$$

$$\approx \frac{1}{\mu} [\lambda t + \lambda^{3/2} \sigma_A B_A(t)] + \sigma_S B_S(\lambda t) \quad \text{B-fluct.}$$

$$= \frac{\lambda}{\mu} t + \frac{\lambda^{1/2}}{\mu} \left[ C_A B_A(t) + C_S \frac{1}{\sqrt{\lambda}} B_S(\lambda t) \right]$$

$$\stackrel{d}{=} \frac{\lambda}{\mu} t + \frac{\lambda^{1/2}}{\mu} (C_A^2 + C_S^2) B(t) \quad \text{self-similarity}$$

$$X(t) = L(t) - t \stackrel{d}{\approx} -(1-\rho)t + \sigma B(t)$$

$$\rho = \frac{\lambda}{\mu} < 1 \quad \text{traffic intensity} ; \quad \sigma^2 = \frac{1}{\mu} \rho (C_A^2 + C_S^2)$$

$$\text{RBM} \quad V \stackrel{d}{\approx} \text{RBM} \left( -(1-\rho) \lambda, \sigma \right) ; \quad V(\infty) \stackrel{d}{=} \exp \left( \text{mean} = \frac{\sigma^2}{2(1-\rho)} \right)$$

$$\text{G/G/1} \quad \frac{V(\infty)}{1/\mu} \stackrel{d}{\approx} \exp \left( \text{mean} = \frac{\rho}{1-\rho} \frac{C_A^2 + C_S^2}{2} \right)$$

# Teaching Note (soon to be installed) on the G/G/1 Queue

## The Congestion Index (Understanding Khinchine-Pollatschek)

$$M/G/1 \Rightarrow E(W_q) = E(S) \frac{\rho}{1-\rho} \frac{1+C^2(S)}{2}$$

G/G/1:

$\frac{E(W_q)}{E(S)} \approx \frac{\rho}{1-\rho} \cdot \frac{C^2(A) + C^2(S)}{2}$	Allen-Cunneen (Hall (5.70))
-----------------------------------------------------------------------------------	--------------------------------

↑  
pure  
(unitless)

↑  
"utilization"  
availability

↑  
stochastic variability

**Fact:** Right-hand side upper bound, becoming asymptotically exact as  $\rho \uparrow 1$  (Heavy traffic).

**Kingman's Exponential Law of Congestion**  
(Invariance principle)

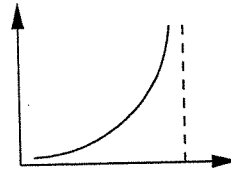
$\frac{W_q}{E(S)} \approx \begin{cases} \exp\left(\text{mean} = \frac{1}{1-\rho} \cdot \frac{C^2(A) + C^2(S)}{2}\right) & , \text{wp } \rho, \\ 0 & , \text{wp } 1-\rho, \end{cases}$
---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

asymptotically exact as  $\rho \uparrow 1$ . (Later, at tails.)

*Understanding the formula:*

Assume  $C^2(A) = C^2(S) = 1$  (as in M/M/1):  $\frac{E(W_q)}{E(S)} = \frac{\rho}{1-\rho}$

Substitute  $\rho = 0.5, 0.9, 0.95, 0.99$ .



**Other MOP's:**

$$E(W) = E(S) + E(W_q)$$

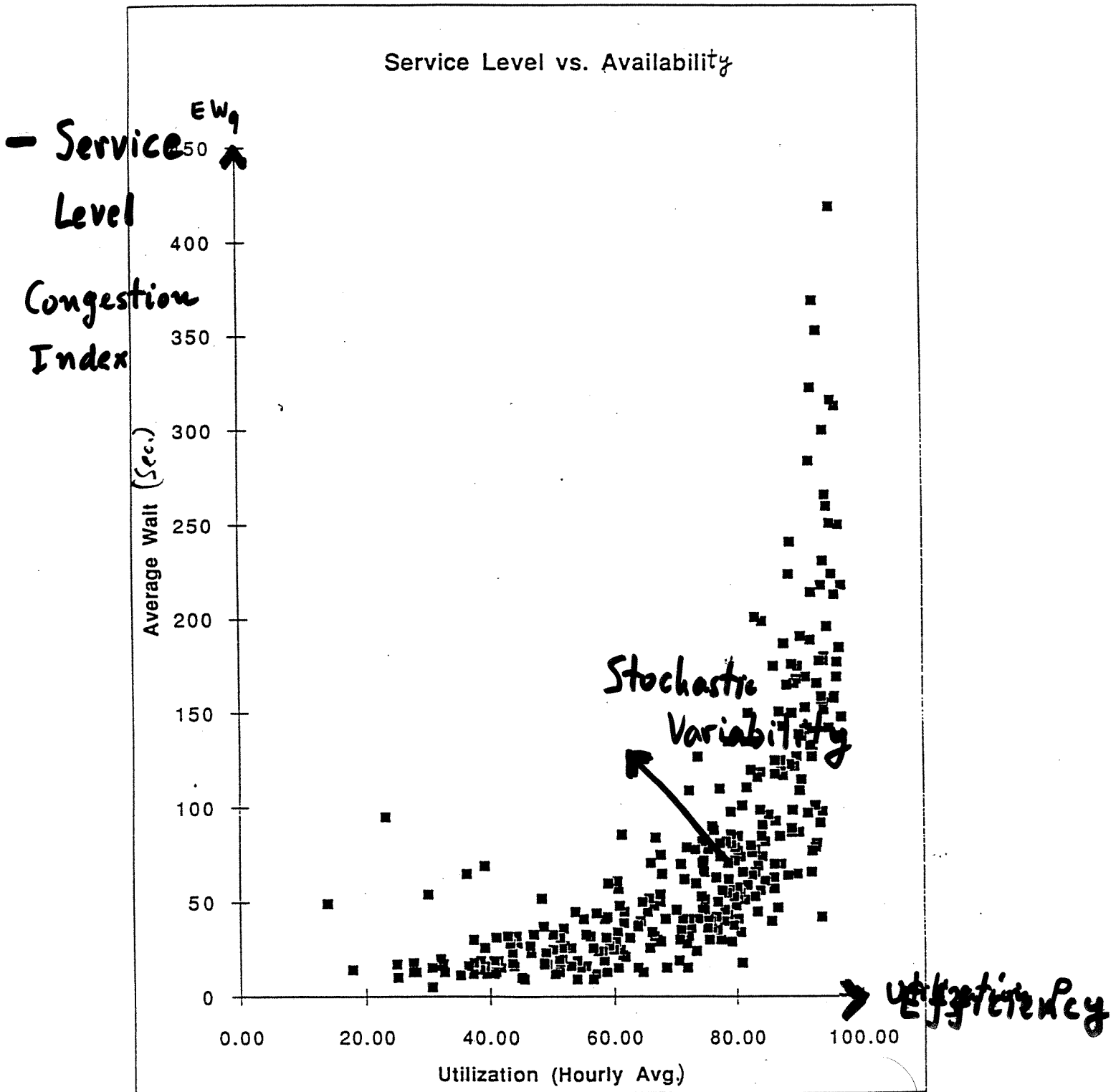
$$E(L_q) = \lambda E(W_q)$$

$$E(L) = \lambda E(W) = E(L_q) + \rho$$



# K-T 2nd LAW OF LONGESTION

## Congestion Curves

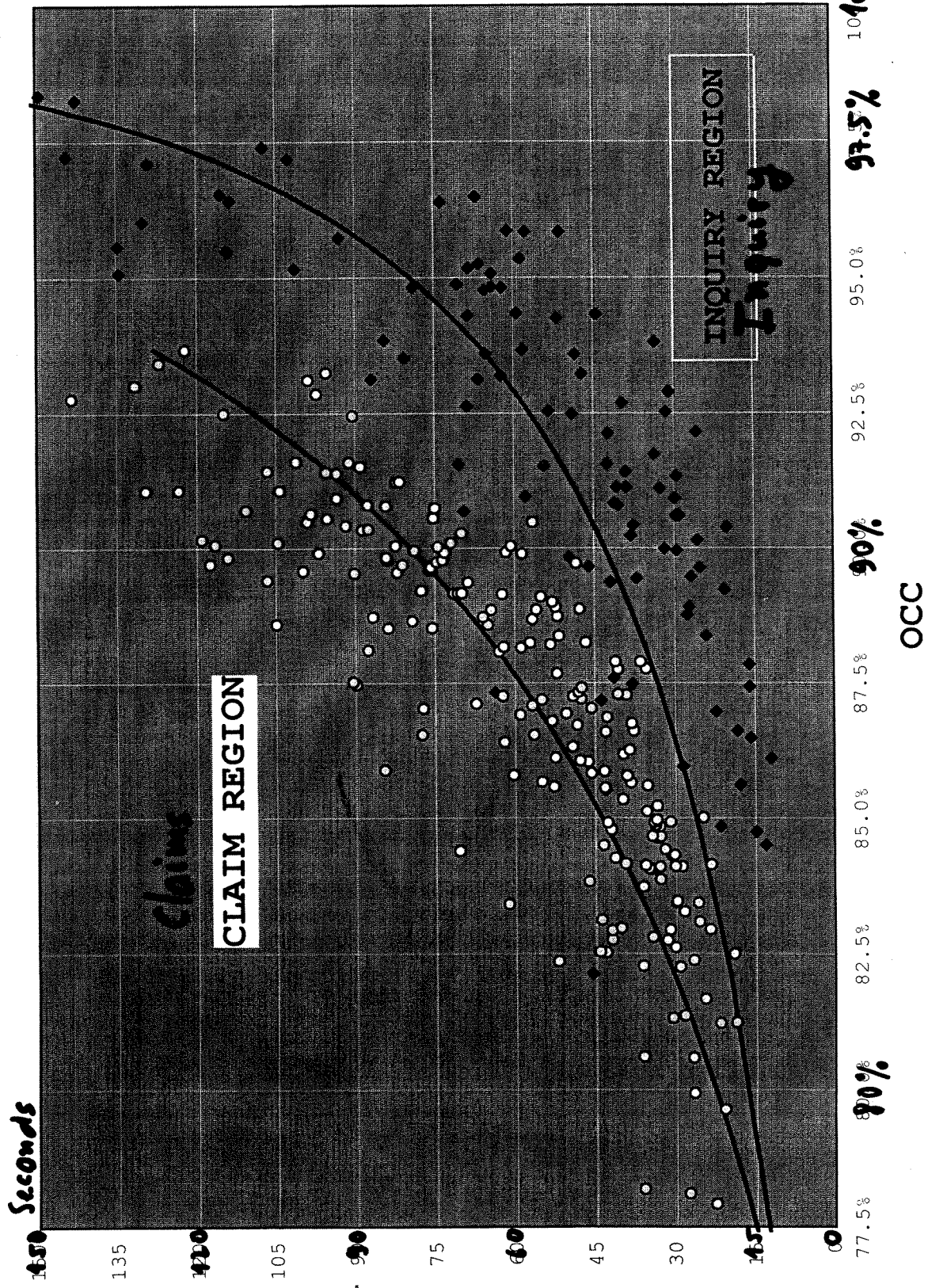


**A** Congestion Index:  $\frac{E(W_q)}{E(S)} \approx \frac{1}{M} \cdot \frac{\rho}{1-\rho} \cdot \frac{C_a^2 + C_s^2}{2}$  ( $M = \# \text{ Servers}$ )

$$= \frac{1}{M} \cdot \frac{\rho}{1-\rho} C^2$$

# Health Insurance Call Center

Daily 06  
SBR



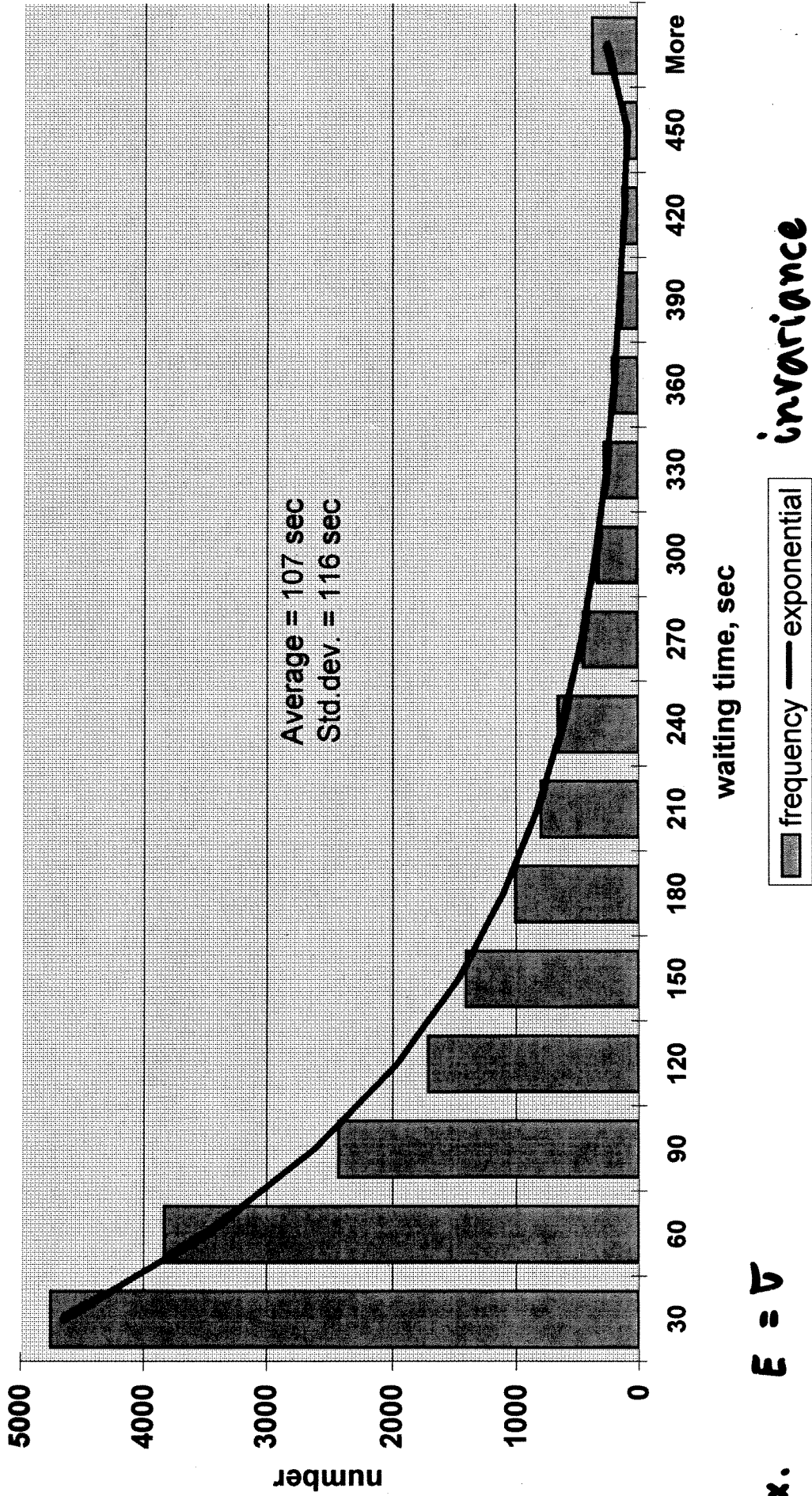
Average  
Speed  
of  
Answer

Staffing?  
Specialization?

Occupancy

# Kingman's 3rd Law of Congestion

November. Waiting times of AGENT customers.

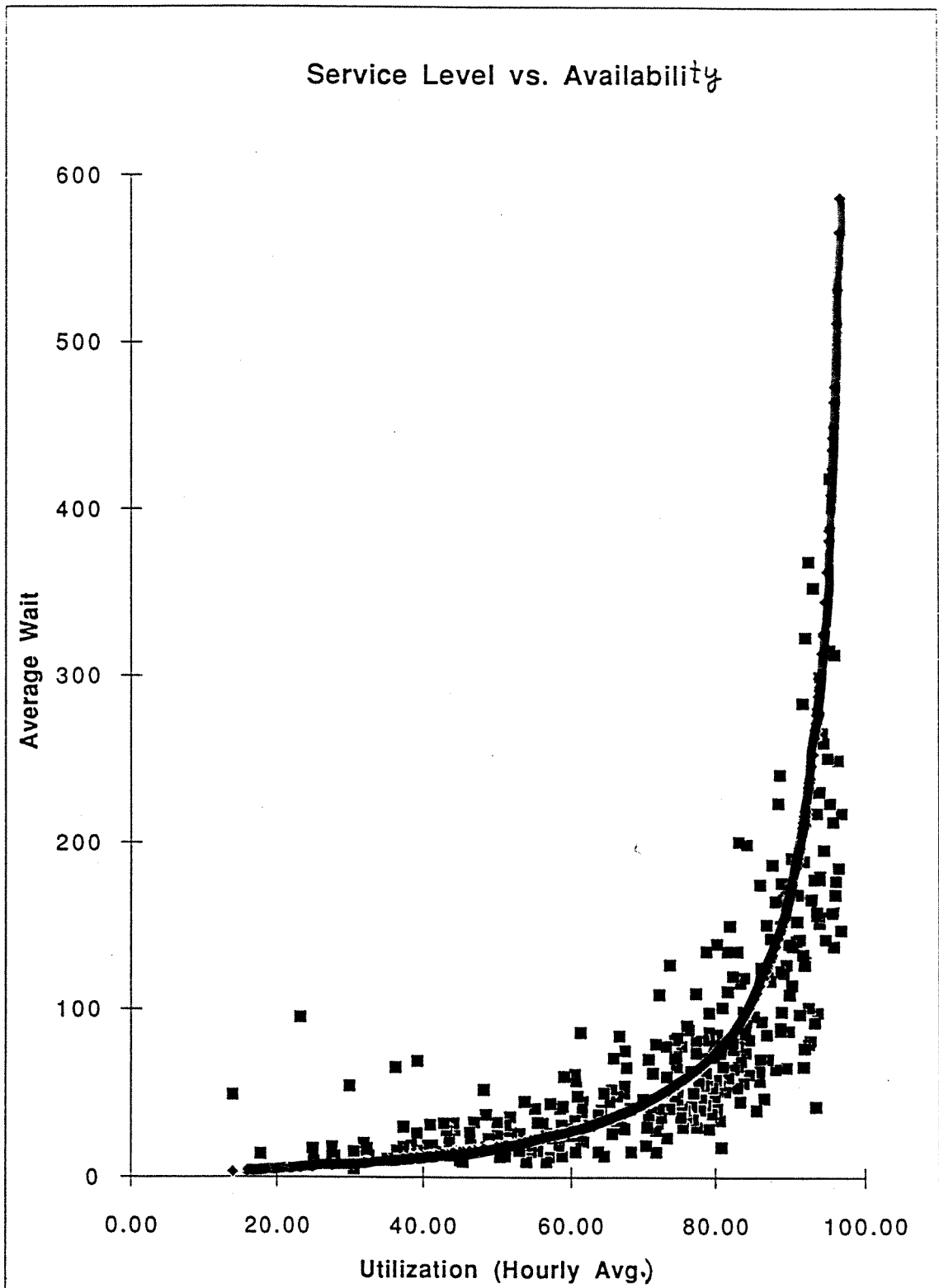


Ex.  $E = \sigma$

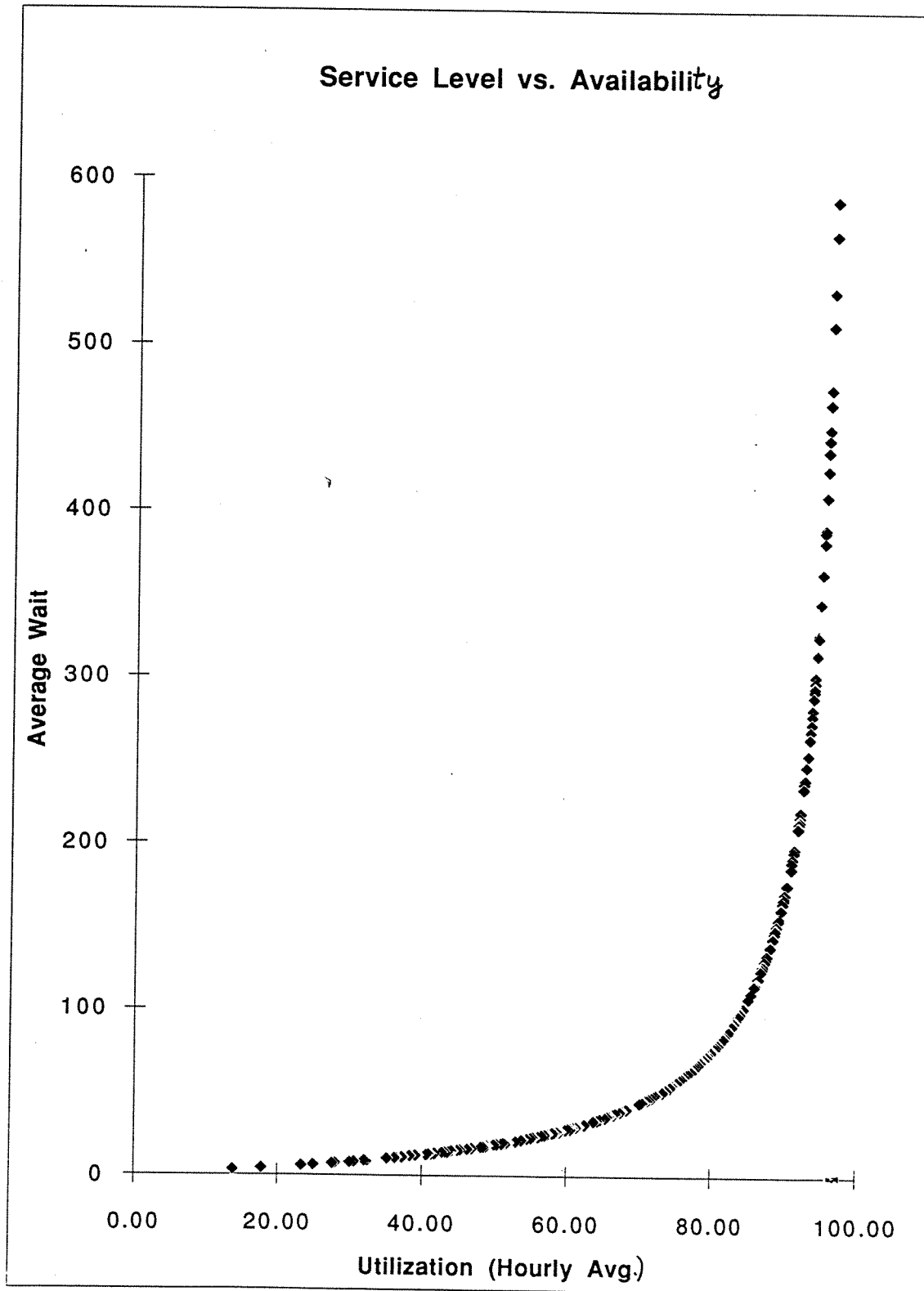
$$P(\text{wait} > E) = e^{-1} \approx 37\%$$

$$P(\text{wait} > 3E) = 5\%$$

Fit  $W_q \approx M \cdot \frac{\rho}{1-\rho} \cdot \boxed{?}$



# Queueing Science



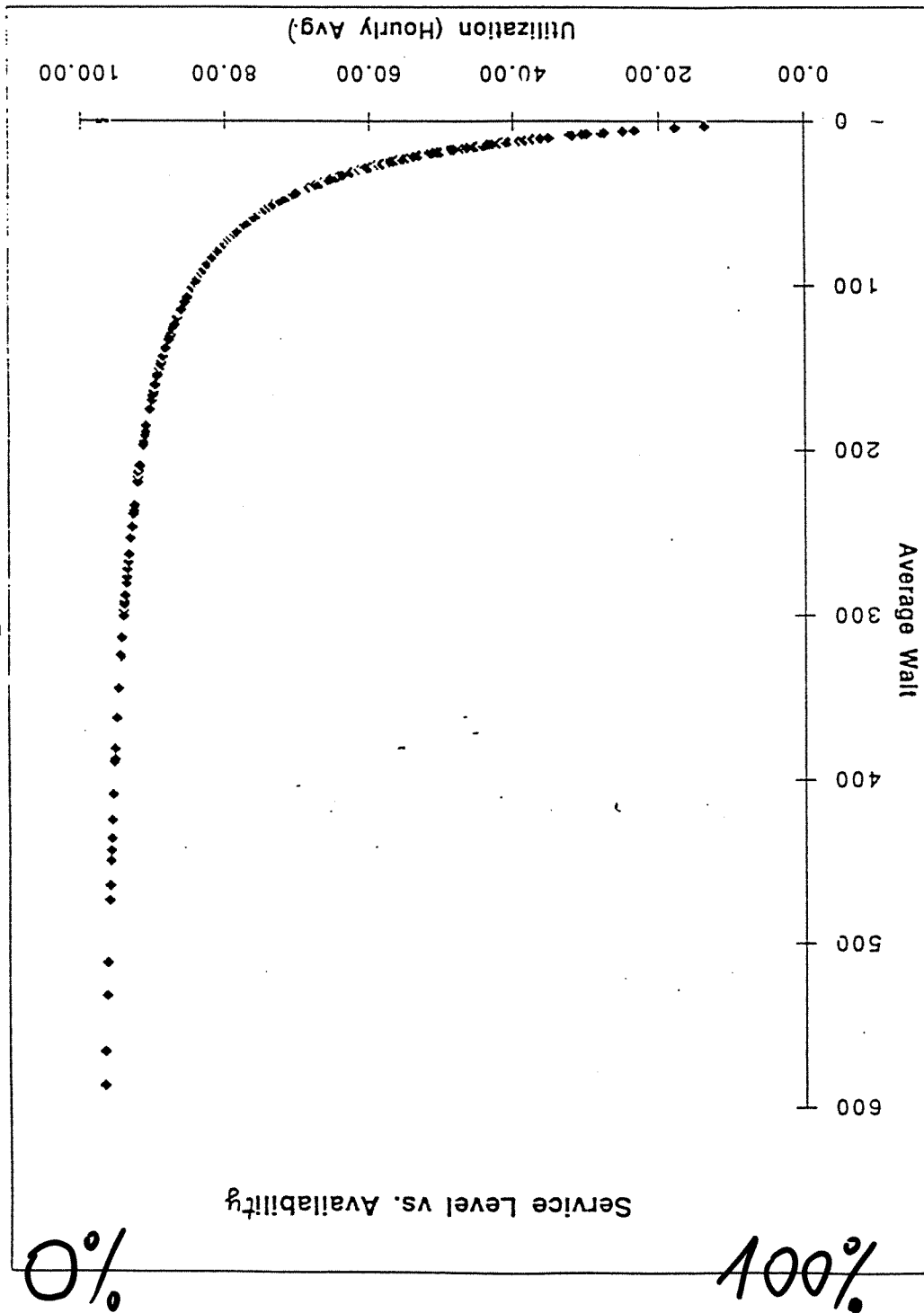
$$W_q = M \times \frac{P}{1-P} \times 0.093$$

Can use

# Service Level, as a function of Availability

$$W_q = M \times \frac{\rho}{1-\rho} \times 0.093$$

Can use



Congestion = f (Availability) ↓

Service Level ↑ with availability

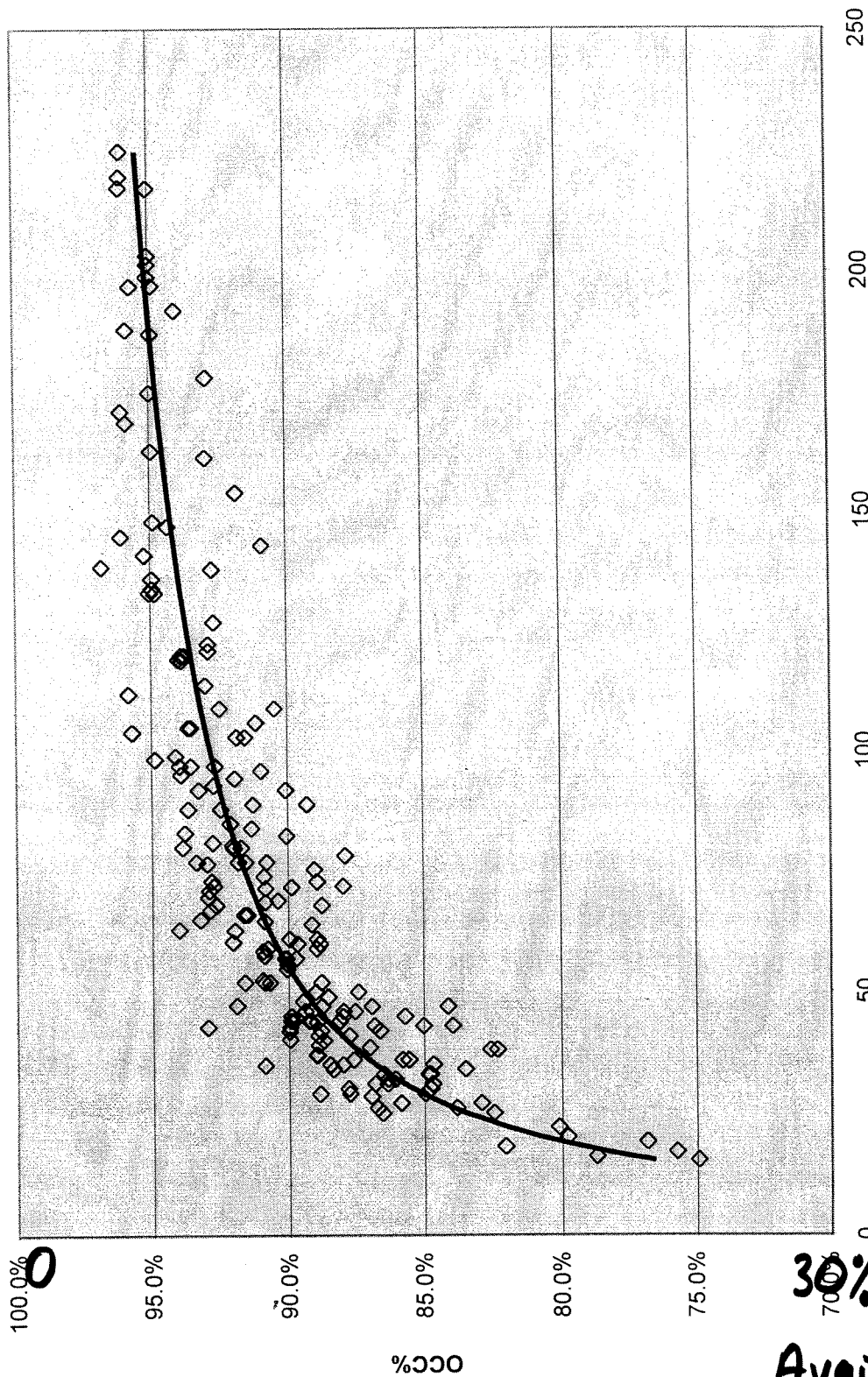
◇ OCC%  
— Model

ASA

BRISTOL

ASA

ASA



OCC

Availability  
(Idleness)

דוגמה: השפעת גודל המערכת

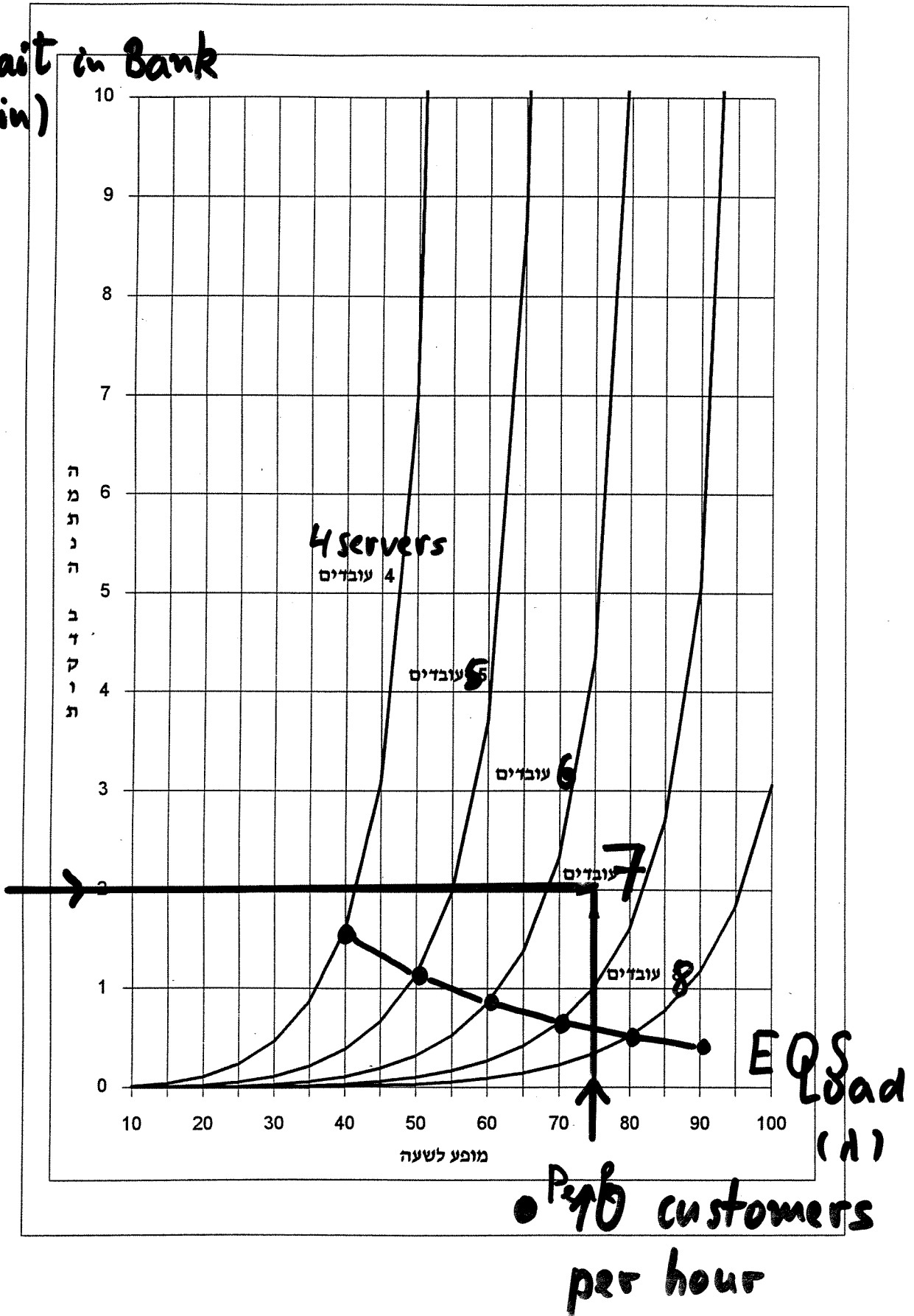
@ given load/efficiency, service level improves w/ size

נספח 6 - אשנב כל מחלקת שירותים בנקאיים - קביעת רמת איוש לפי אמד שירות זמן

המתנה כללי

Avg. Wait in Bank (min)

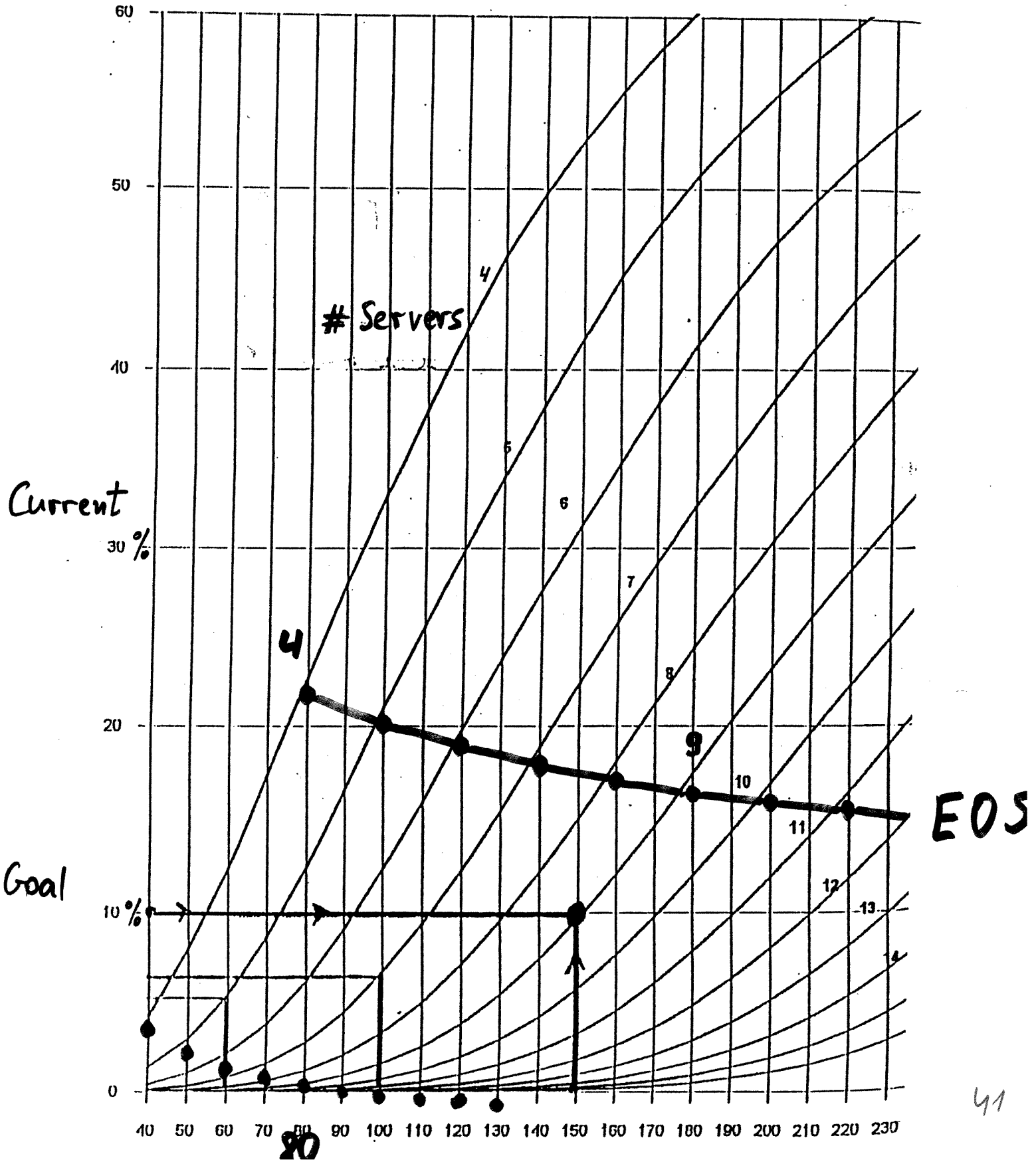
Service level: Goal



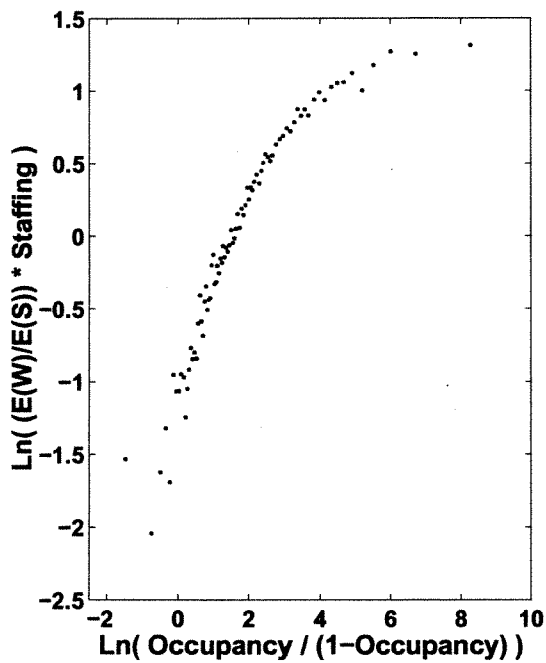
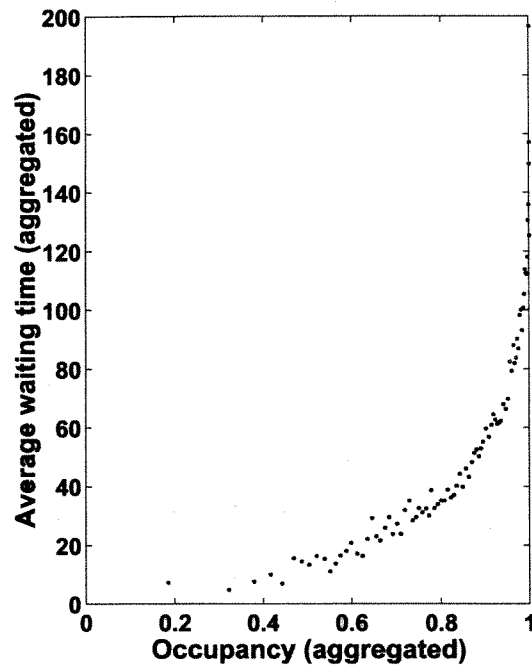
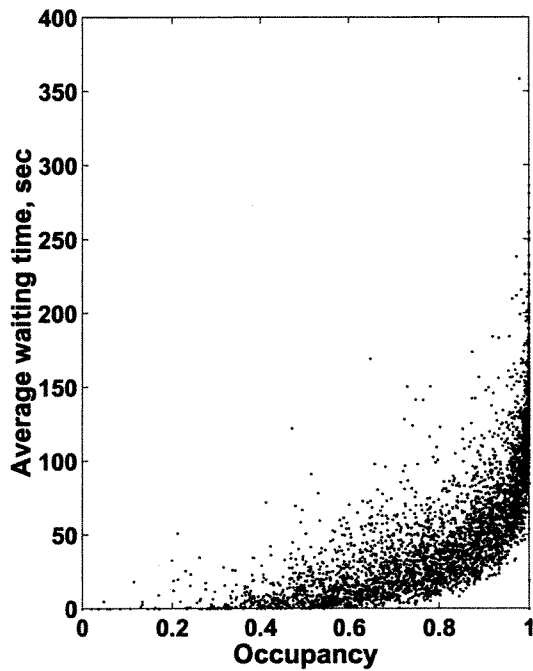


# Abandonment Rate (% Arrivals = Enter)

@ given <sup>MMQ</sup> <sub>203</sub> efficiency): 20 cust's / server / hr  
 service level (abandons) improves with size

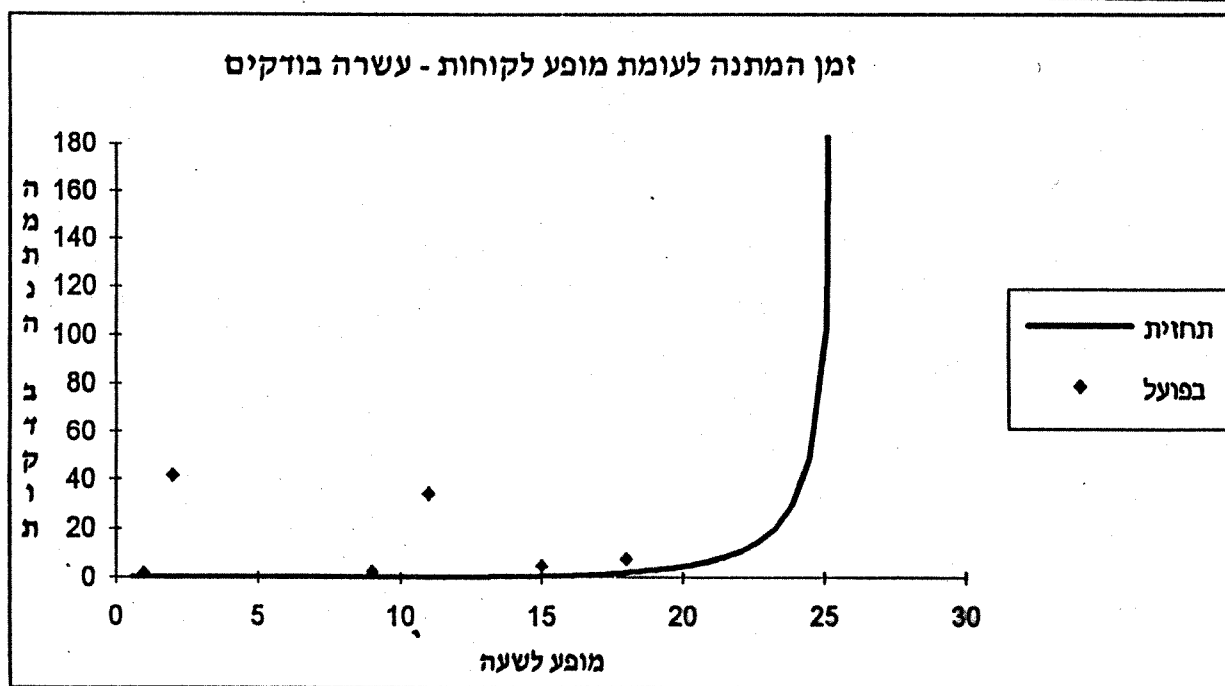
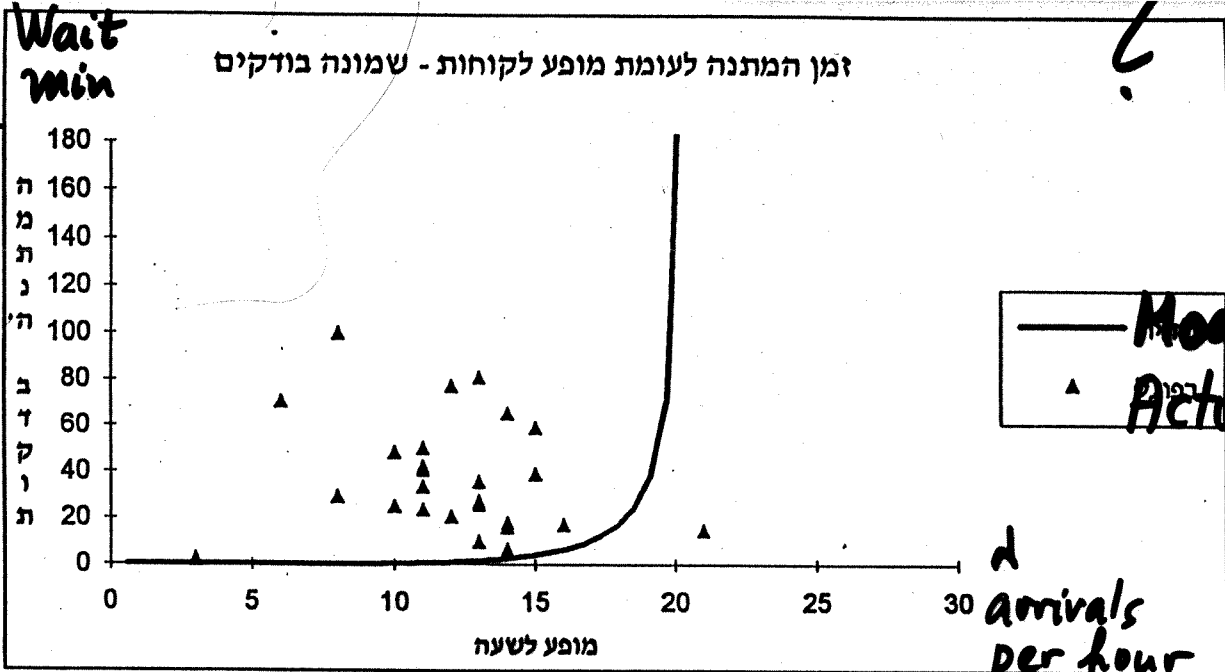


# Average Waiting Time vs. Agents' Occupancy (Khintchine-Pollaczek)



W

Wait  
min



# The $M_t/M_t/1$ Queue

Arrival Time-varying Poisson

" " arrival rate  $\lambda_t, t \geq 0$ .

Services Think exponential( $\mu$ )

( Time-varying service rate  $\mu_t, t \geq 0$ .)

Stable? when periodic,  $\exists$  analysis (Harrison & Lemoine, ...)

Must resort to approx. : short-run

Representation  $Z = \{Z_t, t \geq 0\}$  number in system

$Z = f(X)$   $f$  reflection

$$X_t = N_+ \left( \int_0^t \lambda_u du \right) - N_- \left( \int_0^t \mu_u du \right)$$

Std Poisson indep.

# Hierarchical Modelling: Short-Run (Massey)

$$M_t / M_t / 1 \quad X_t = N_+ \left( \int_0^t \alpha_u du \right) - N_- \left( \int_0^t \mu_u du \right)$$

eg. periodic
std. Poisson

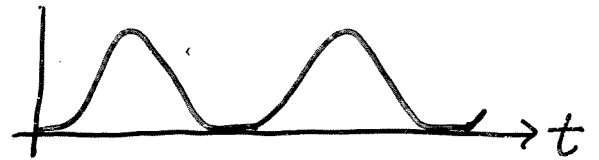
Micro  $Z = f(X)$

Acceleration  $Z^\epsilon = f(X^\epsilon), \quad \alpha_t^\epsilon \leftrightarrow \frac{1}{\epsilon} \alpha_t, \quad \mu_t^\epsilon \leftrightarrow \frac{1}{\epsilon} \mu_t, \quad \epsilon \downarrow$

Strong Approx  $X^\epsilon \approx \underbrace{\frac{1}{\epsilon} [\int \alpha - \int \mu]}_{\frac{1}{\epsilon} x} + \underbrace{W_+ \left( \frac{1}{\epsilon} \int \alpha \right) - W_- \left( \frac{1}{\epsilon} \int \mu \right)}_{\frac{1}{\sqrt{\epsilon}} W \text{ in dist.}}$

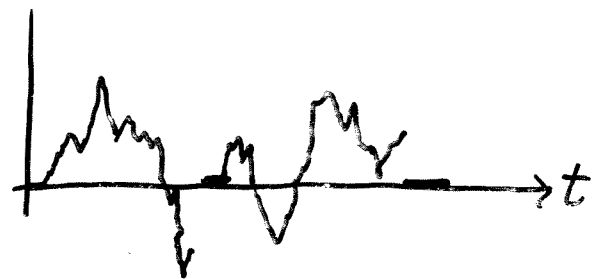
FLLN  $\epsilon Z^\epsilon \approx f(x + \sqrt{\epsilon} W) \xrightarrow{\text{a.s.}} f(x) = q \text{ fluid}$

Macro  $Z \sim \frac{1}{\epsilon} q$



FCLT  $\frac{1}{\sqrt{\epsilon}} [\epsilon Z^\epsilon - q] \xrightarrow{M_1} \nabla_W f(x) = V \text{ diffusion}$

Meso  $Z \sim \frac{1}{\epsilon} q + \frac{1}{\sqrt{\epsilon}} V$



Phase Transitions: over  $\rightarrow$  under  $\rightarrow$  critical 45

Fluid determines time-varying phases

$$q(t) > 0 \quad \text{super-critical} \quad \Leftrightarrow \rho(t) > 1$$

$$q(t) = 0 \quad dy > 0 \quad \text{sub-critical} \quad \Leftrightarrow \rho(t) < 1$$

$$q(t) = 0 \quad dy = 0 \quad \text{critical} \quad \Leftrightarrow \rho(t) = 1$$

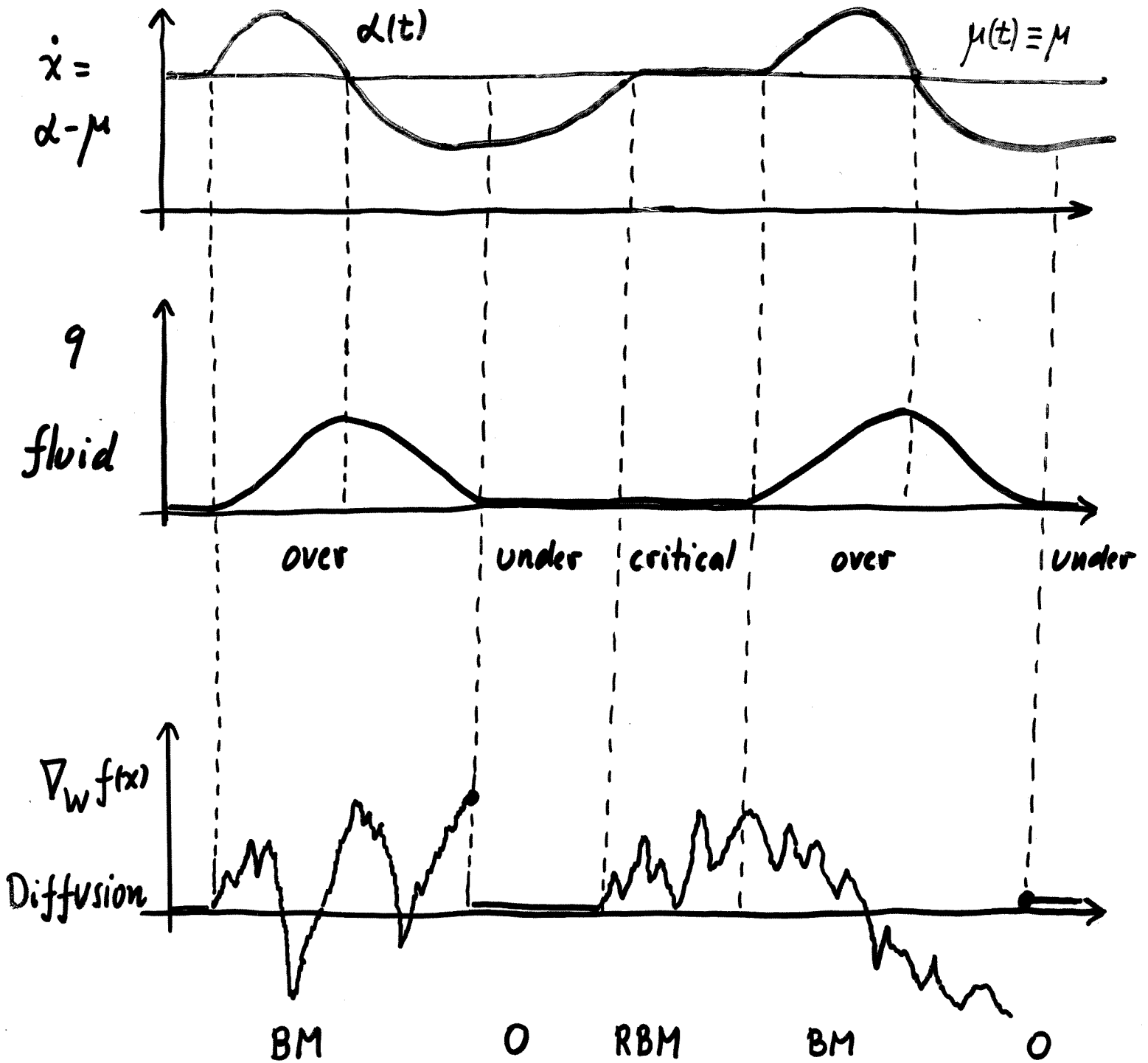
$$\rho(t) = \sup_{0 \leq s \leq t} \frac{\int_s^t du \, du}{\int_s^t \mu \, du}$$

$$\left( \rho(t) \neq \frac{\dot{L}}{\dot{r}} \right)$$

Diffusion approx. changes by phases

# Phase Transitions

$$p(t): < 1, = 1, > 1$$



$M_1$  queue of size  $\frac{1}{\sqrt{\epsilon}} = \sqrt{n}$  depletes during  $\sqrt{n}$

Dynamic acceleration: slow down  $\pm \sqrt{n}$  around jumps. but accelerated by  $n$

# Phase Transitions

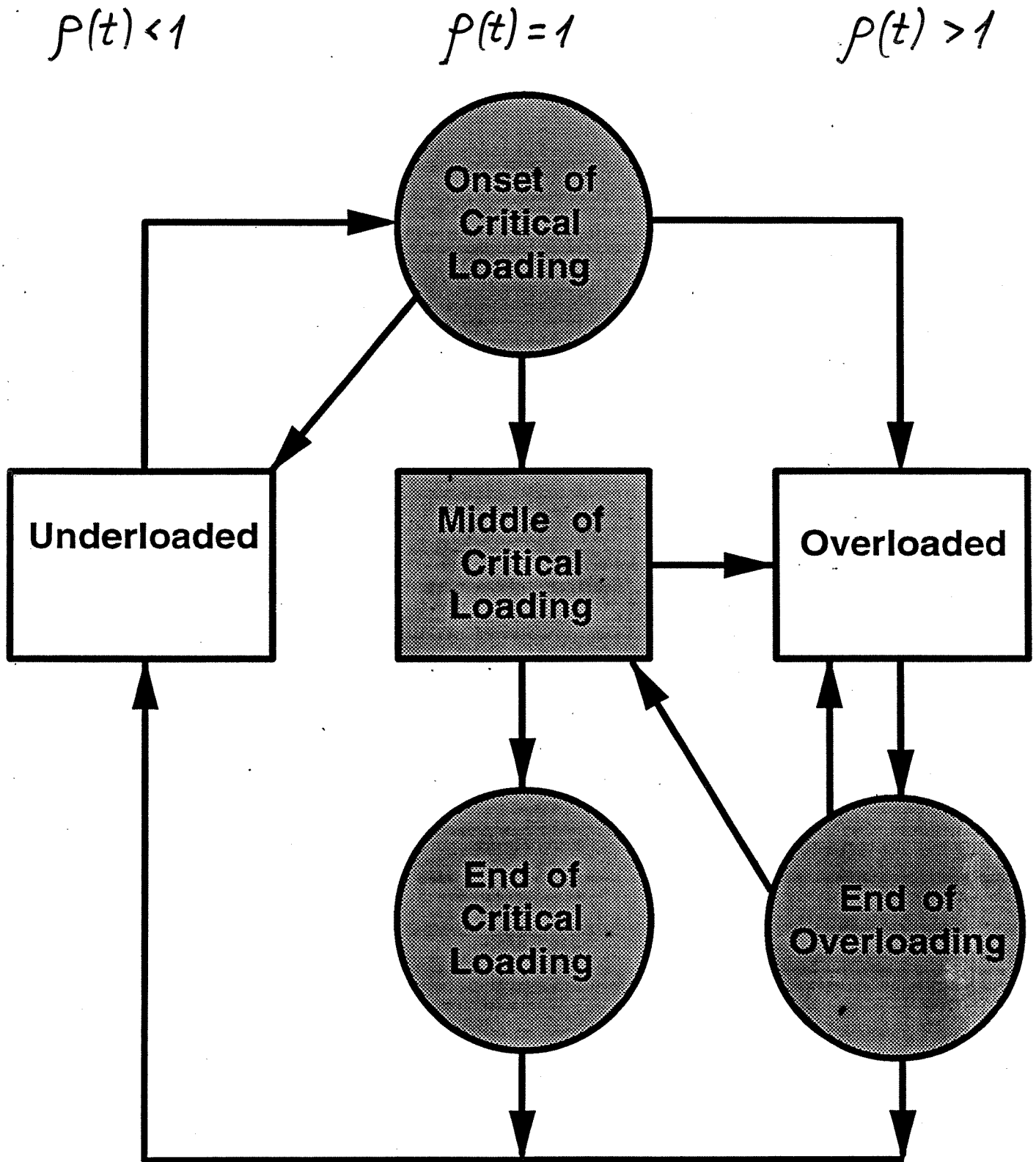


Figure 3.1: Phase transition diagram for the asymptotic regions

Refined regions for asymptotic analysis



# Approximations - Examples

Middle + End of Critical Loading

$$Q \stackrel{d}{=} \sqrt{\epsilon} V + o(\sqrt{\epsilon}), \quad V = \nabla_W f(x)$$

Middle : continuous, reflected BM

End :  $V$  left-continuous, not right-cont. w.p.  $> 0$

Overloading

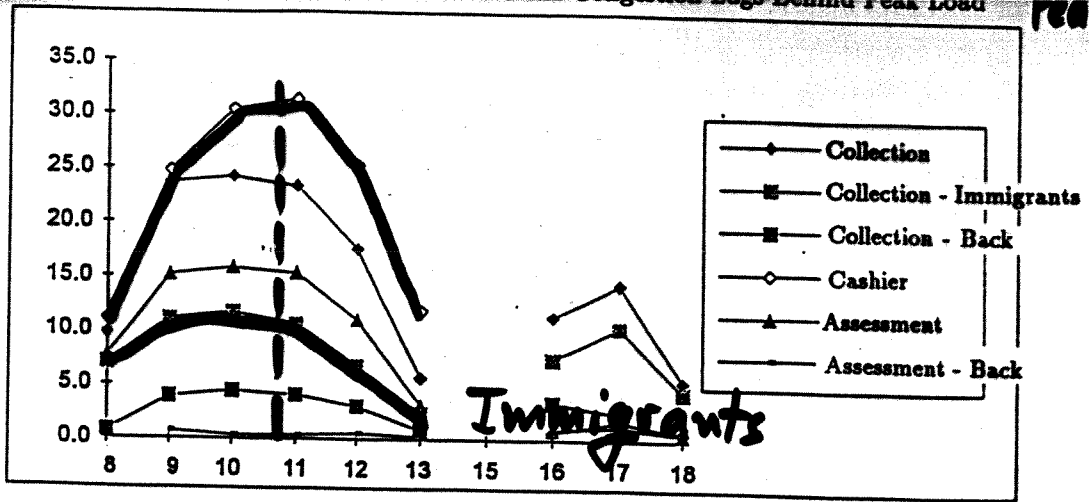
$$Q(t) \stackrel{d}{=} \int_s^t (\alpha - \mu) ds + \sqrt{\epsilon} [V(s) + \Delta W(s, t)] + o(\sqrt{\epsilon})$$

$\uparrow$   
onset

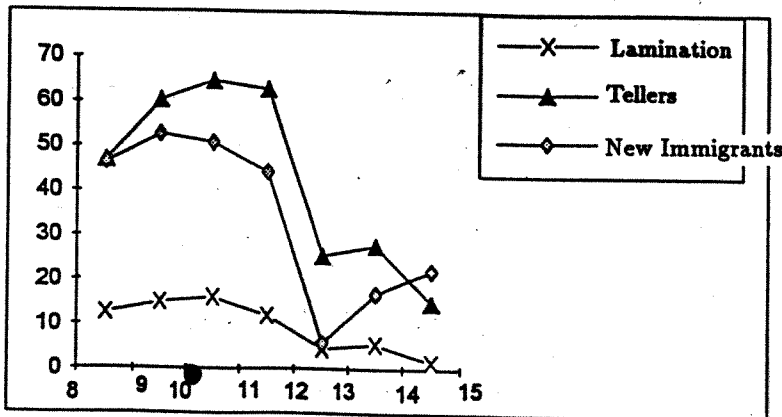
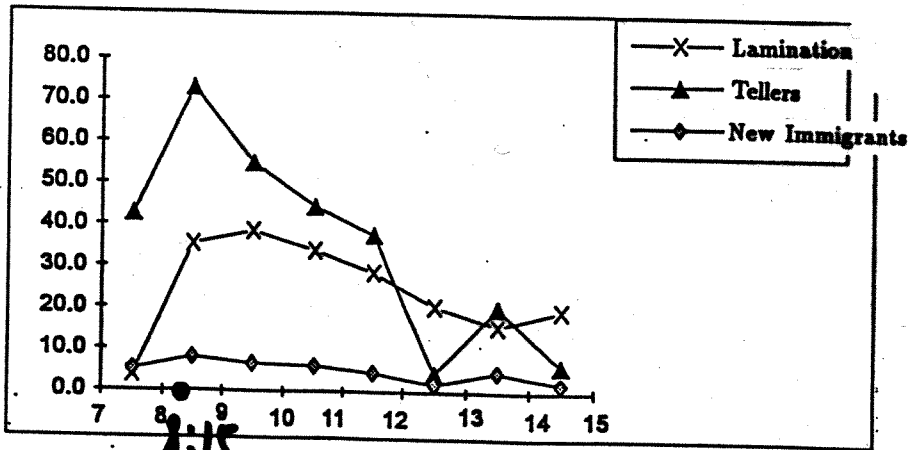
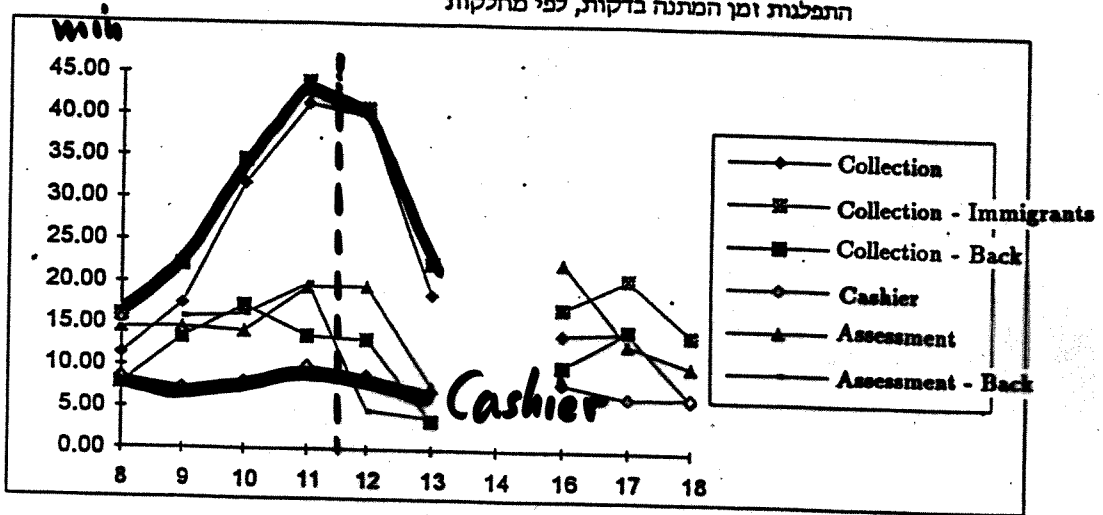
Eg. Peak load < Peak congestion (queue)

$$t = \operatorname{argmax} q(\cdot) \text{ over } [s, t]$$
$$\Rightarrow \exists u = \operatorname{argmax} d(\cdot) \text{ over } [s, t]$$

Arrival Rates:  
6 Dept's



Waiting Times



10:15

# Short-Run: State-Dependent

$$Z : \alpha(z), \mu(z)$$

$$Z : \alpha^\varepsilon(z) = \frac{1}{\varepsilon} \alpha(\varepsilon z), \mu^\varepsilon(z) = \dots, \varepsilon \downarrow 0 \text{ Acceleration}$$

$$\text{FSLLN} \quad \varepsilon Z^\varepsilon \xrightarrow{\text{a.s.}} q = f(x) \quad \text{Fluid}$$

$$\text{D.E.} \quad \dot{x} = [\alpha(q) + \mu(q)P] - \mu(q)$$

$$\text{FCLT} \quad \frac{1}{\sqrt{\varepsilon}} [\varepsilon Z^\varepsilon - q] \xrightarrow{M_1} V = \nabla_W f(x) \quad \text{Diffusion}$$

$$\text{S.D.E.} \quad dW = a(q) V dt + b(q) dB$$

Phases  $\rho(t) = \text{traffic intensity, based on } q$   
 $< 1, = 1, > 1$

1 Station Repairman (Iglehart)

Multiple Servers (Halfin, Whitt)

Finite Population (Kogan, Krichagina, Lipster)

Reneging (Coffman, Pukhalski, Reiman, Wright)

Networks  $P(\text{state})$  - Routing

## REVIEW: MARKOV JUMP-PROCESS (MJP)

**MJP**  $X = \{X_t, t \geq 0\}$  on  $\mathcal{S} = \{i, j, \dots\}$  countable.

Markov property:  $P_r\{X_t = j | X_r, r < s; X_s = i\} = P_{ij}(s, t), \forall s < t, \forall i, j \in \mathcal{S}$ .

Time homogeneity:  $P_r\{X_{s+t} = j | X_s = i\} = P_{ij}(t), \forall s, t, i, j$ , transition probabilities.

Characterization:  $\pi^0 =$  initial distribution and  $P(t) = [P_{ij}(t)], t \geq 0$ , stochastic.

Finite-dimensional distributions:

$$P_r\{X_0 = i_0, X_{t_1} = i_1, \dots, X_{t_n} = i_n\} = \pi^0(i_0)P_{i_0, i_1}(t_1) \dots P_{i_{n-1}, i_n}(t_n - t_{n-1}).$$

$P(t)$  : stochastic ;  $P(s+t) = P(s)P(t), \forall s, t$  (Chapman Kolmogorov);

$$\exists P(0) = I ; \exists \dot{P}(0) = Q = [q_{ij}], \text{ infinitesimal generator } \left( \sum_{j \in \mathcal{S}} q_{ij} = 0 \right).$$

Micro to Macro :  $\dot{P}(t) = P(t)Q (=QP(t))$  and  $P(0) = I$   
Forward (Backward) equations.

$$\text{Solution : } P(t) = \exp[tQ] = \sum_{n=0}^{\infty} \frac{t^n}{n!} Q^n, t \geq 0.$$

Animation:  $i \xrightarrow{q_{ij}} j; \forall i, j \in \mathcal{S} \exists$  exponential clock at rate  $q_{ij}$ , call it  $(i, j)$ .

Given  $i$ , consider clocks  $(i, j), j \in \mathcal{S}$ ; move to the "winner" when rings.

Thus: stay at  $i \sim \exp(q_i = \sum_{j \neq i} q_{ij})$  and switch to  $j$  with probability  $P_{ij} = q_{ij}/q_i$   
( $q_{ij} = q_i P_{ij}, i \neq j; q_{ii} = -q_i$ ).

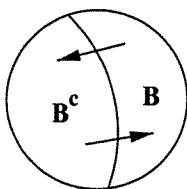
Transient analysis vs. long-run/limit stability/steady-state  
 $\exists \lim_{t \uparrow \infty} P_{ij}(t) = \pi_j, \forall i; \pi = \pi P(t), \forall t.$

$$\text{Calculation via steady-state equations: } \dot{P}(\infty) = P(\infty)Q \Rightarrow \left\{ \begin{array}{l} 0 = \pi Q \\ \sum_i \pi_i = 1, \pi_i \geq 0 \end{array} \right\}$$

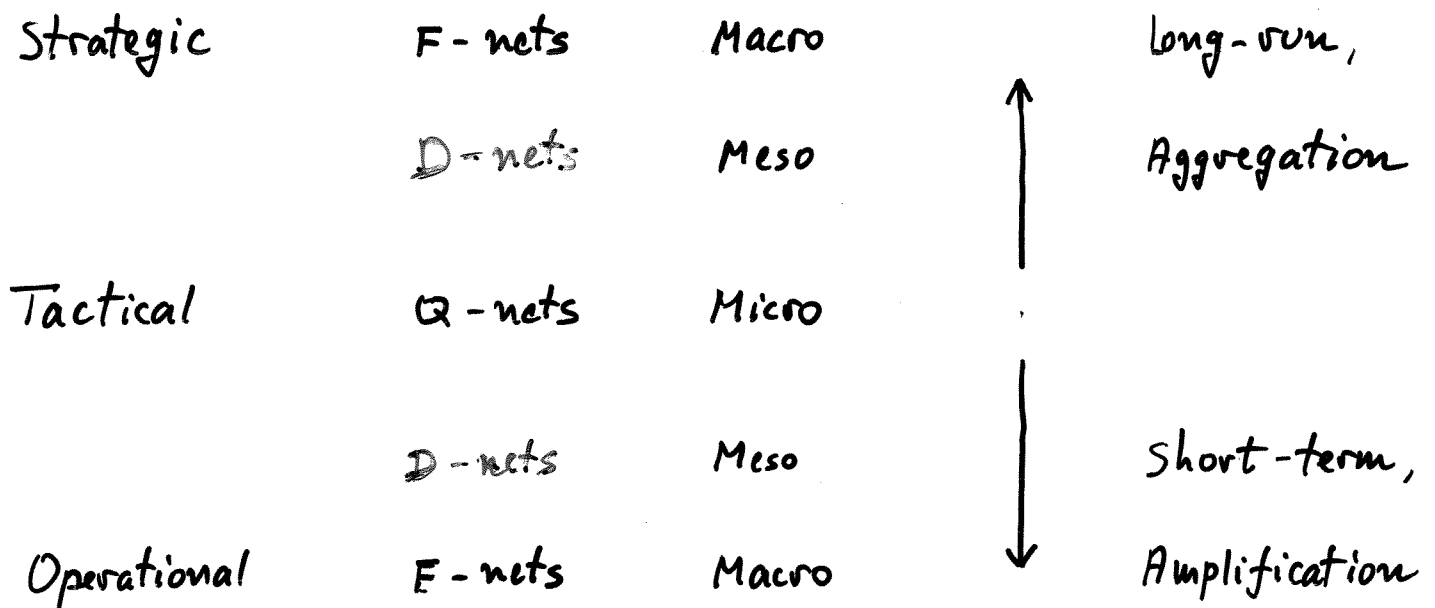
or balance equations:  $\sum_{i \neq j} \pi_i q_{ij} = -\pi_j q_{jj} = \sum_{i \neq j} \pi_j q_{ji}, \forall j.$

Transition rates:  $\pi_i q_{ij} =$  long-run average number of switches from  $i$  to  $j$ .

$$\text{Cuts: } \sum_{i \in B} \sum_{j \in B^c} \pi_i q_{ij} = \sum_{i \in B^c} \sum_{j \in B} \pi_i q_{ij}, \forall B \subset \mathcal{S}.$$



# Hierarchical Modelling: a Framework



Anthony, R.N. "Planning and Control Systems:  
A Framework for Analysis", 1965

"I kept recalling the point that President James B. Conant demonstrated so vividly in *On Understanding Science*, namely

The development of a framework or a conceptual scheme

often has led to progress, even though the framework

turns out to be wrong ... the framework will have served

a useful purpose if it prepares the way for a better one."