

Service Engineering (Science, Management)

Avi Mandelbaum
Technion IE&M

Course Contents

- Introduction to “Services” and “Service-Engineering”
- The Two Prerequisites: Measurements, Models (Operational)
- Empirical (Data-Based) Models
- Fluid (Deterministic) Models
- Stochastic Framework: Dynamic-Stochastic PERT/CPM
- The Building Blocks of a Basic Service Station:
 - Arrivals; Forecasting
 - Service Durations; Workload
 - (Im)Patience; Abandonment
 - Returns (During, After; Positive, Negative)
- Stochastic Models of a Service Station
 - Markovian Queues: Erlang $B/C/A, \dots, R$, Jackson
 - Non-Parametric Queues: $G/G/n, \dots$
- Operational Regimes and Staffing: ED, QD, QED
- Heterogeneous Customers and Servers (CRM, SBR)

Background Material

Downloadable from the **References** menu in
<http://ie.technion.ac.il/serveng/References>

Gans (U.S.A.), Koole (Europe), and M. (Israel):
“Telephone Call Centers: Tutorial, Review and Research Prospects.”
MSOM, 2003.

Brown, Gans, M., Sakov, Shen, Zeltyn, Zhao:
“**Statistical** Analysis of a Telephone Call Center: A Queueing-
Science Perspective.” JASA, 2005.

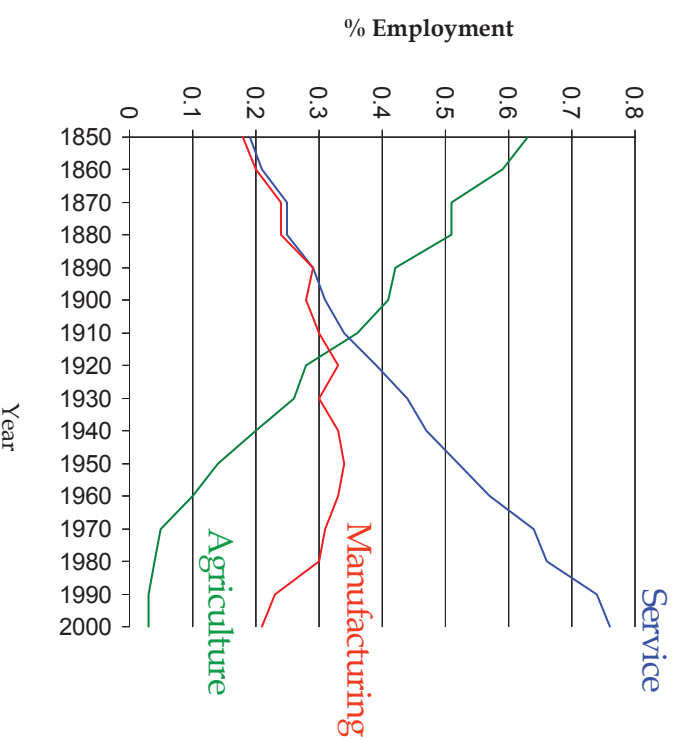
Trofimov, Feigin, M., Ishay, Nadjharov:
“**DataMOCCA**: Models for Call/Contact Center Analysis. (Model
Description and Introduction to User Interface.)” Technion Report,
2004-2006.

Technion’s “**Service-Engineering**” course lectures: Measure-
ments, Arrivals, Service Times, (Im)Patience, Fluid Models, QED
Q’s.

M. “Call Centers: Research **Bibliography** with Abstracts.”
Version 7, December 2006.

Introduction to “Services”

U.S. Employment by Sector, 1850 - 2000+



We focus on:

- Function: **Operations** (vs./plus IT, HRM, Marketing)
- Dimension: Accessibility, **Capacity** (vs. RM, SCM,...)
- Modelling Framework: **Queueing** Theory (plus Science)
- Applications: **Call/Contact Centers** (Healthcare,...)

3

Scope of the Service Industry

- Wholesale and retail trade;
- Government services;
- Healthcare;
- Restaurants and food;
- Financial services;
- Transportation;
- Communication;
- Education;
- Hospitality business;
- Leisure services.

Our Application Focus: **telephone call centers**, which play an important role in most of these sectors.

4

Services: Subjective Trends

“Everything is Service”

Rather than buying a **product**, why not **buy only the service it provides**? For example, **car leasing**; or, why setup and run a **help-desk** for technical support, with its costly fast-to-obsolete hardware, growing-sophisticated software, high-skilled peopleware and ever-expanding infoware, rather than let **outsourcing** do it all for you?

“Data; Technology and Human Interaction

Far too little reliance on **data**, the **language of nature**, in formulating models for the **systems and processes of the deepest importance to human beings**, namely those in which **we are actors**. Systems with fixed rules, such as physical systems, are relatively simple, whereas systems involving human beings expressing their microgoals ... can exhibit incredible complexity; there is yet the hope to devise tractable models through **remarkable collective effects** ...

(Robert Herman: “Reflection on Vehicular **Traffic Science**”.)

Fusion of Disciplines: POM/IE, Marketing, IT, HRM

The highest challenge facing banks with respect to efficient and effective innovation lies in the “**New Age Industrial Engineer**” that must combine technological knowledge with process design in order to create the delivery system of the future.

(Frei, Harker and Hunter: “**Innovation in Retail Banking**”).

Service-Engineering

Goal (Subjective):

Develop scientifically-based design principles (**rules-of-thumb**) and tools (**software**) that support the balance of service **quality**, process **efficiency** and business **profitability**, from the (often conflicting) views of customers, servers and managers.

Contrast with the traditional and prevalent

- **Service** Management (U.S. Business Schools)
- Industrial **Engineering** (European/Japanese Engineering Schools)

Additional **Sources** (all with websites):

- Fraunhofer **IAO** (Service Engineering, 1995): ... application of engineering science know-how to the service sector ... models, methods and tools for systematic development and design of service products and service systems ...
- **NSF SEE** (Service Enterprise Engineering, 2002): ... Customer Call/Contact Centers ... staff scheduling, dynamic pricing, facilities design, and quality assurance ...
- **IBM SSME** (Services Science, Management and Engineering, 2005): ... new discipline brings together computer science, operations research, industrial engineering, business strategy, management sciences, social and cognitive sciences, and legal sciences ...

Staffing: How Many Servers?

Fundamental problem in service operations: Healthcare, . . . , or **Call Centers**, as a representative example:

- People: $\approx 70\%$ operating costs; $\geq 3\%$ U.S. workforce.
- Business-Frontiers but also **Sweat-Shops** of the 21st Century.

Reality

- **Complex** and becoming more so
- Staffing is Erlang-based (1913!)

\Rightarrow Solutions urgently needed

- Technology can accommodate smart protocols
- Theory lags significantly behind needs

\Rightarrow Ad-hoc methods prevalent: heuristics- or simulation-based.

Research Progress based on

- **Simple Robust Models**, for theoretical insight into complex realities. Their analysis requires and generates:
- Data-Based **Science**: Model, Experiment, Validate, Refine.
- **Management** Principles, Tools: **Service Engineering**.

The First Prerequisite: Data & Measurements

Robert Herman (“Father” of Transportation Science): Far too little reliance on **Data**, the **language of nature**, in formulating models for the systems of the deepest importance to human beings, namely those in which we are actors.

Empirical “Axiom”: The Data One Needs is **Never** There
For One To Use (Always Problems with Historical Data).

Averages do NOT tell the whole story

Individual-Transaction Level Data: Time-Stamps of Events

- **Face-to-Face**: T, C, S, I, O, F (QJE, RFID)
- **Telephone**: ACD, CTI/CRM, Surveys
- **Internet**: Log-files
- **Transportation**: measuring devices on highways/intersections

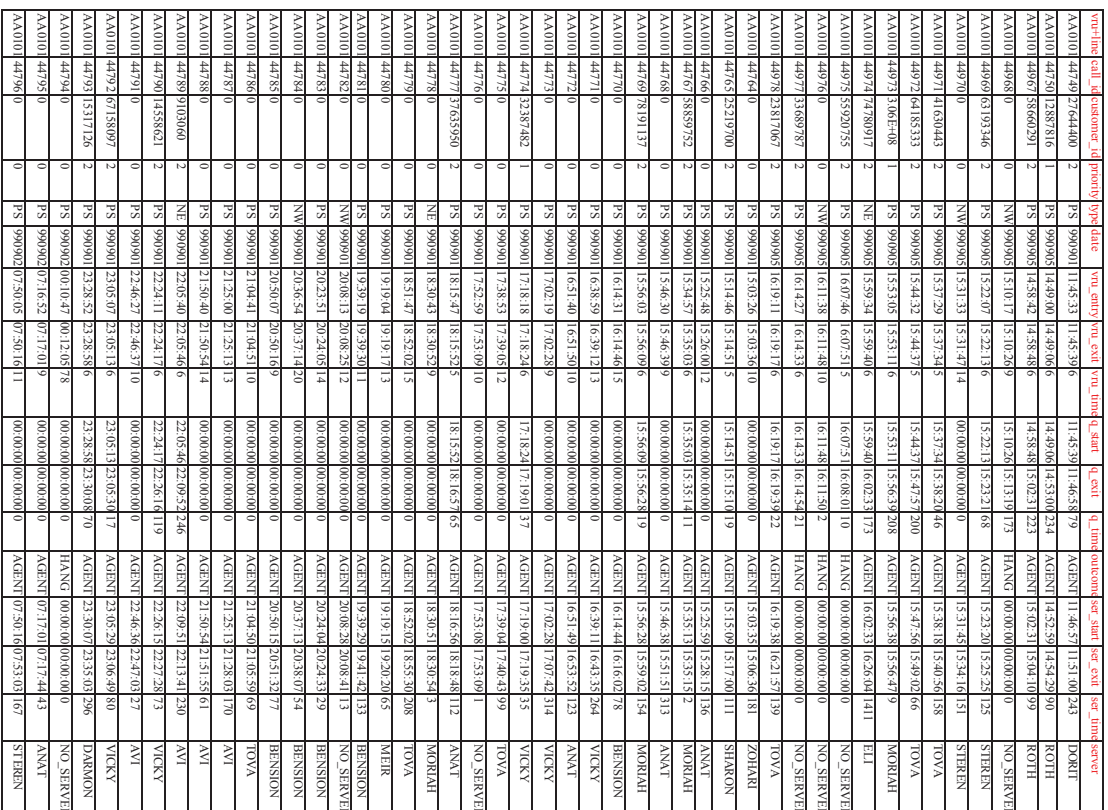
Our Databases: Operations (vs. Marketing, Surveys, . . .)

- Face-to-Face data (branch banking) – recitations; QUESTA
- Telephone data (small banking call center) – homework; JASA
- **DataMOCCA** (large cc’s: repository; interface) – class/research; Website

Future Research:

Healthcare, Multimedia, Field-Support; Operation+Marketing,

Measurements: Telephone Services Log-File of Call-by-Call Data



Measurements:
Prevalent Averages (ACD Data)

Command Center Intraday Report

Date		Updated Through: All Day									
06/13 - Tue		Recvd	Answw	Abn %	ASA	AHT	Occ %	On	On Prod	Sch	Open Sch Avail %
Total:		129,960	126,321	2.8%	31	318	90.9%	88.4%	1631.7	1585.0	96.6%
INQ	Charlotte	20,577	19,860	3.5%	30	307	95.1%	85.4%	222.7	234.6	95.0%
INQ	Columbus MSC	7,973	7,773	2.5%	36	314	94.9%	89.8%	89.2	94.5	94.4%
INQ	Phoenix	17,102	16,757	2.0%	31	298	92.7%	91.8%	187.3	194.8	96.2%
INQ	Scranton	1,257	1,254	0.2%	6	515	78.6%	28.9%	28.5	35.1	81.2%
INQ	Tampa	9,174	8,899	3.4%	42	366	91.5%	93.6%	123.1	125.9	97.8%
CEN	Boulevard	6,070	5,937	2.2%	33	362	86.7%	90.2%	86.0	88.4	97.3%
CEN	Bristol	10,667	10,505	1.5%	25	355	95.1%	93.1%	136.3	139.6	97.6%
CEN	Columbus Claims	5,258	5,153	2.0%	27	293	86.7%	89.8%	60.5	62.2	97.3%
STH	Atlanta	7,514	7,338	2.3%	40	318	82.1%	89.5%	98.6	99.8	98.8%
STH	Sherman	19,669	18,833	4.3%	46	252	90.8%	90.6%	175.5	174.9	100.4%
STH	Wilmington	10,422	9,888	5.1%	21	285	89.9%	92.1%	108.7	114.6	94.8%
WST	Visalia	14,277	14,164	0.8%	10	382	87.2%	85.0%	215.2	220.6	97.6%

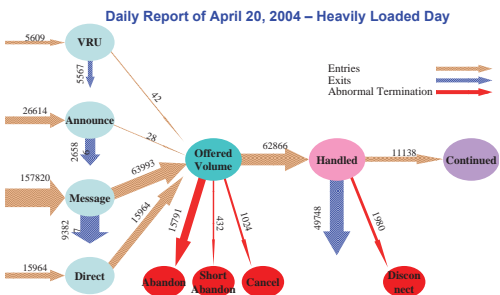
12 cc's

- Center

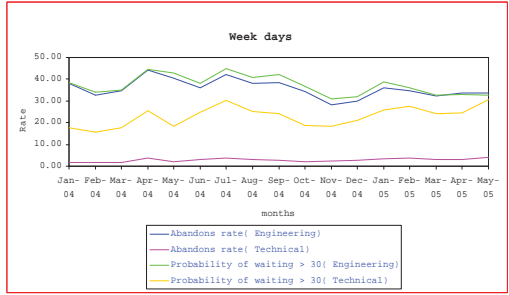
Time	Recvd	Answw	Abn %	ASA	AHT	Occ %	Prod%	On	On Prod	Sch	Open Sch	Avail %
0	20,577	19,860	3.5%	30	307	95.1%	85.4%	222.7	234.6	95.0%		
8:00	332	308	7.2%	27	302	87.1%	79.5%	59.3	68.9	88.5%		
8:30	653	615	5.8%	58	283	96.1%	81.1%	104.1	111.7	93.2%		
9:00	886	796	8.1%	63	308	97.1%	84.7%	140.4	145.3	96.6%		
9:30	1,152	1,138	1.2%	28	303	90.8%	81.5%	211.1	221.3	95.4%		
10:00	1,330	1,286	3.3%	22	307	98.4%	84.3%	223.1	229.0	97.4%		
10:30	1,384	1,338	1.9%	33	296	99.0%	84.1%	222.5	227.9	97.6%		
11:00	1,380	1,280	7.2%	34	306	98.2%	84.0%	222.0	223.9	99.2%		
11:30	1,272	1,247	2.0%	44	298	94.6%	82.8%	218.0	233.2	93.5%		
12:00	1,179	1,177	0.2%	1	306	91.6%	88.6%	218.3	222.5	98.1%		
12:30	1,174	1,160	1.2%	10	302	95.5%	93.5%	203.8	208.8	97.4%		
13:00	1,018	999	1.9%	9	314	95.4%	91.2%	182.9	187.0	97.8%		
13:30	1,061	981	9.4%	67	306	100.0%	88.8%	163.4	182.5	89.5%		
14:00	1,173	1,082	7.8%	78	313	99.5%	85.7%	188.9	213.0	98.3%		
14:30	1,212	1,179	2.7%	23	304	96.8%	86.0%	206.1	220.9	93.3%		
15:00	1,137	1,122	1.3%	15	320	96.9%	83.5%	205.8	222.1	92.1%		
15:30	1,169	1,137	2.7%	17	311	97.1%	84.5%	202.2	207.0	97.7%		
16:00	1,107	1,059	4.3%	46	315	99.2%	79.4%	187.1	192.9	97.0%		
16:30	914	892	2.4%	22	307	95.2%	81.8%	160.0	172.3	92.6%		
17:00	615	615	0.0%	2	328	83.0%	93.5%	135.0	146.2	92.2%		
17:30	420	420	0.0%	0	328	73.8%	95.4%	103.5	116.1	89.2%		
18:00	49	49	0.0%	14	180	84.2%	89.1%	5.8	1.4	416.2%		

DataMOCCA

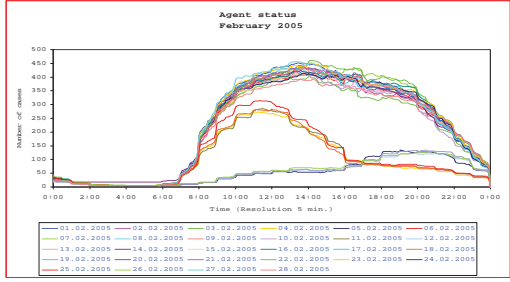
Daily Report



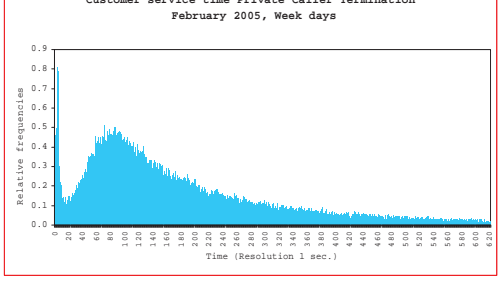
Time Series



Cross Tabulation

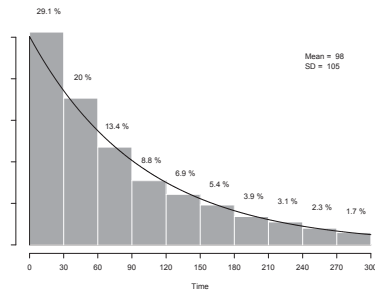


Histogram

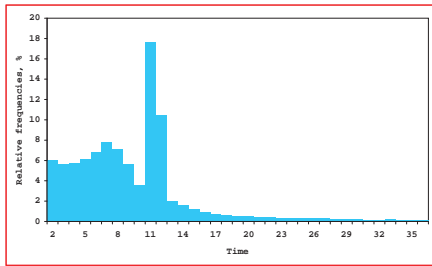


Beyond Averages: Waiting Times in a Call Center

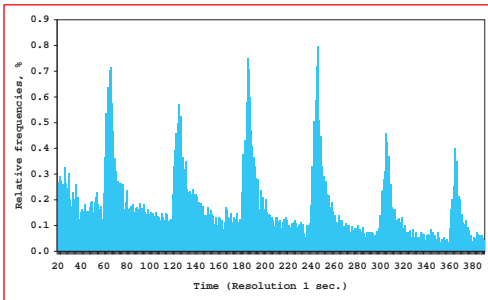
Small Israeli Bank



Large U.S. Bank



Medium Israeli Bank



8

The Second Prerequisite: (Operational) Models

Empirical Models

- Conceptual
 - Service-Process **Data** = **Flow** Network
 - **Service Networks** = **Queueing Networks**
- Descriptive
 - QC-Tools: Pareto, Gantt, Fishbone Diagrams,...
 - Histograms, Hazard-Rates, ...
 - Data-MOCCA: Repository + Interface

Explanatory

- Nonparametric: Comparative Statistics, Regression,...
- Parametric: Log-Normal Services, (Doubly) Poisson Arrivals, Exponential (Im)Patience

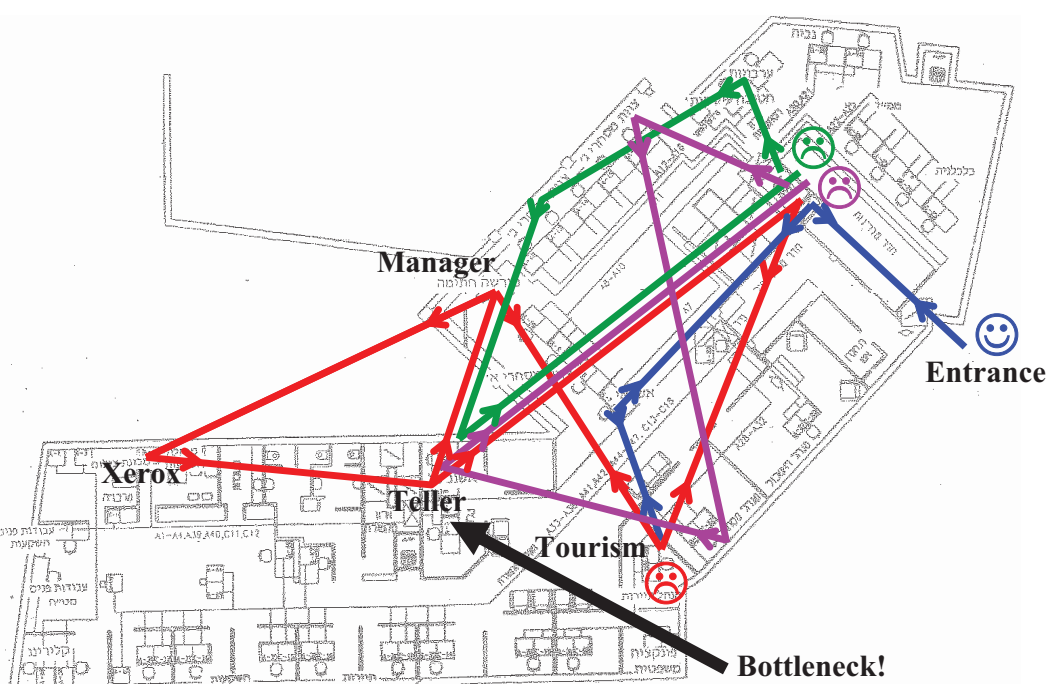
Analytical Models

- Fluid (Deterministic) Models
- Stochastic Models (Birth & Death, $G/G/n$, Jackson,...)

Conceptual Model: Service Networks = Queueing Networks

- People, **waiting for service**: teller, repairman, ATM
- Telephone-calls, to be answered: busy, music, info.
- Forms, to be sent, processed, printed; **for a partner**
- Projects, to be developed, approved, implemented
- Justice, to be made: pre-trial, hearing, retrial
- Ships, for a pilot, berth, unloading crew
- Patients, for an ambulance, emergency room, operation
- Cars, in rush hour, for parking
- Checks, waiting to be processed, cashed
- **Queues** **Scarce Resources, Synchronization Gaps**
Costly, but here to stay
 - Face-to-face Nets (Chat) (min.)
 - Tele-to-tele Nets (Telephone) (sec.)
 - Administrative Nets (Letter-to-Letter) (days)
 - Fax, e.mail (hours)
 - Face-to-ATM, Tele-to-IVR
 - Mixed Networks (Contact Centers)

Conceptual Model: Bank Branch = Queueing Network



Bank Branch: A Queuing Network

Transition Frequencies Between Units in The Private and Business Sections:

From Unit	To Unit	Private Banking			Business				Exit
		Bankers	Authorized Personal	Compensations	Tellers	Tellers	Overdrafts	Authorized Personal	Full Service
Private	Bankers		1%	1%	4%	4%	0%	0%	90%
	Authorized Personal	12%		5%	4%	6%	0%	0%	73%
Banking	Compensations	7%	4%		18%	6%	0%	0%	64%
	Tellers	6%	0%	1%		1%	0%	0%	90%
Services	Tellers	1%	0%	0%	0%		1%	0%	94%
	Overdrafts	2%	0%	1%	1%	19%		5%	64%
	Authorized Personal	2%	1%	0%	1%	11%	5%		69%
	Full Service	1%	0%	0%	0%	8%	1%	2%	88%
	Entrance								
		13%	0%	3%	10%	58%	2%	14%	0%

Legend: 0%-5% 5%-10% 10%-15% >15%

Dominant Paths - Business:

Unit	Station 1	Station 2	Total
Parameter	Tourism	Teller	Dominant Path
Service Time	12.7	4.8	17.5
Waiting Time	8.2	6.9	15.1
Total Time	20.9	11.7	32.6
Service Index	0.61	0.41	0.53

Dominant Paths - Private:

Unit	Station 1	Station 2	Total
Parameter	Banker	Teller	Dominant Path
Service Time	12.1	3.9	16.0
Waiting Time	6.5	5.7	12.2
Total Time	18.6	9.6	28.2
Service Index	0.65	0.40	0.56

Service Index = % time being served

Mapping the Offered Load (Bank Branch)

Department	Business Services		Private Banking		Banking Services	
	Tourism	Teller	Teller	Teller	Comprehensive	
Time						
8:30 – 9:00						
9:00 – 9:30						
9:30 – 10:00						
10:00 – 10:30						
10:30 – 11:00						
11:00 – 11:30						
11:30 – 12:00						
12:00 – 12:30						
Break						
16:00 – 16:30						
16:30 – 17:00						
17:00 – 17:30						
17:30 – 18:00						

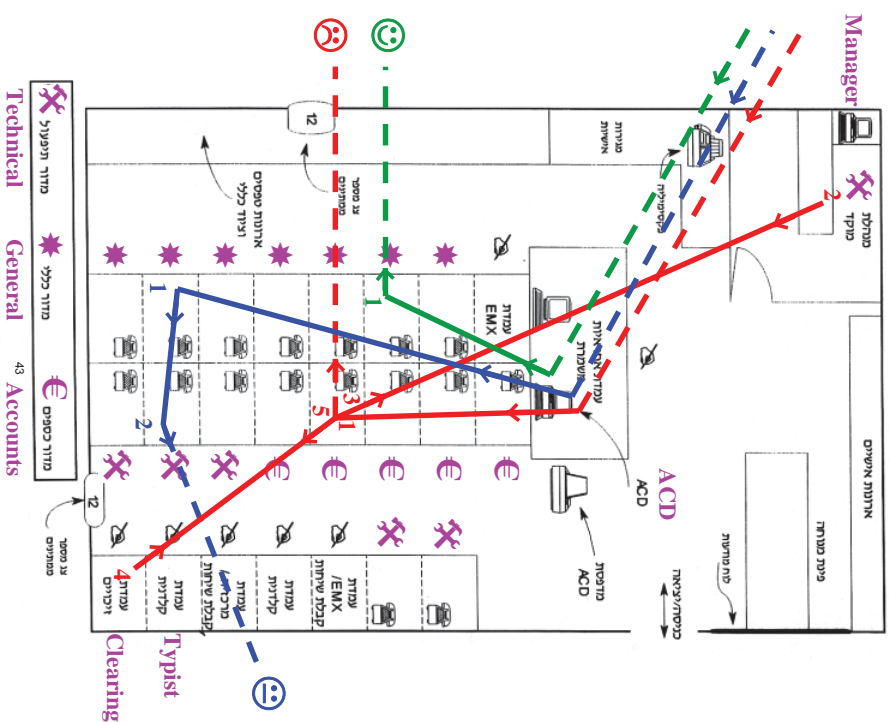
Legend: Not Busy Busy Very Busy

Note: What can / should be done at 11:00 ?

Conclusion: Models are not always necessary but measurements are !

Conceptual Model: Call-Center Network

Schematic Chart – Pelephone Call-Center 1994
= Tele Net = Queueing Network



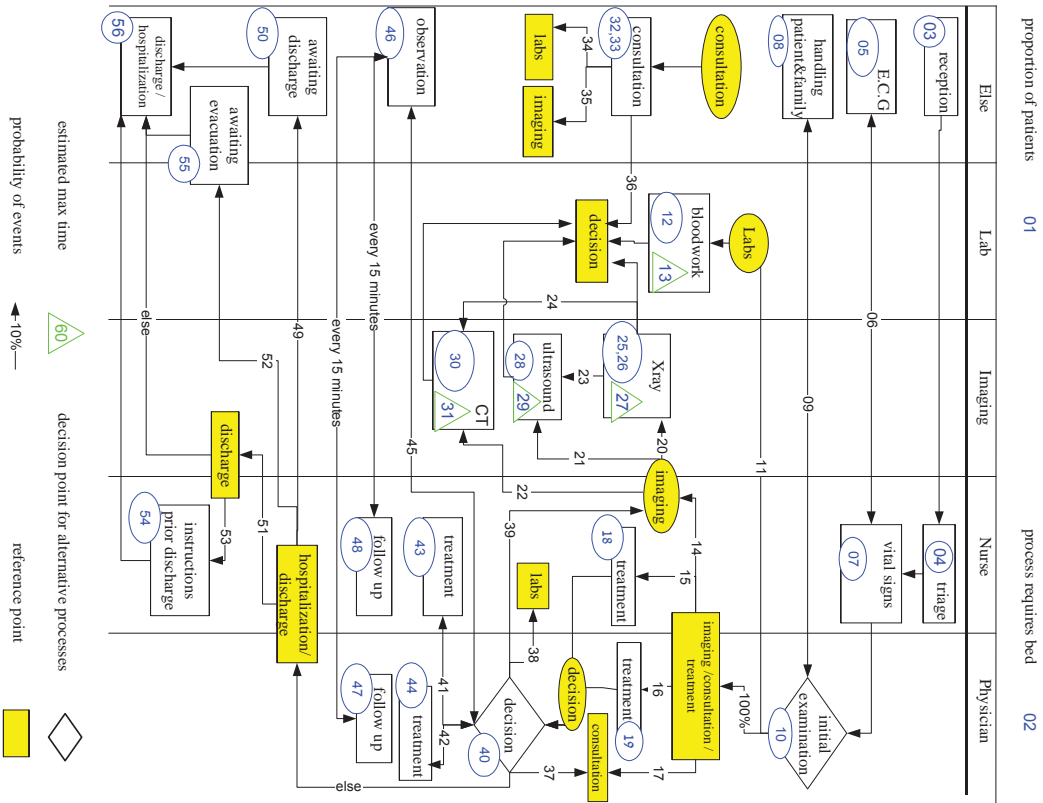
Conceptual Model: Call-Center Network

Current Status - Analysis

	Accounts Center	General Center	Technical Center
Peak days in a week	Sun, Fri	Sun	Sun
Peak days in a month	12	8-14, 2-3	10-20
Avg. applications no. in a day	4136	2476	1762
Avg. applications no. in an hour - λ_{avg}	253.6	193	167
Peak hours in a day	11:00-12:00	10:00-11:00	9:00-10:00
Avg. applications no. in peak hours - λ_{max}	422	313	230
Avg. waiting time (secs.)	10.9	20.0	55.9
Avg. service time (secs.)	83.5	131.3	143.2
Service index	0.88	0.87	0.72
Abandonment percentage	2.7	5.6	11.2
Avg. waiting time before abandonment (secs.)	9.7	16.8	43.2
Avg. staffing level	9.7	10.3	5.2
Target waiting time	12	25	-

Conceptual Model: Hospital Network

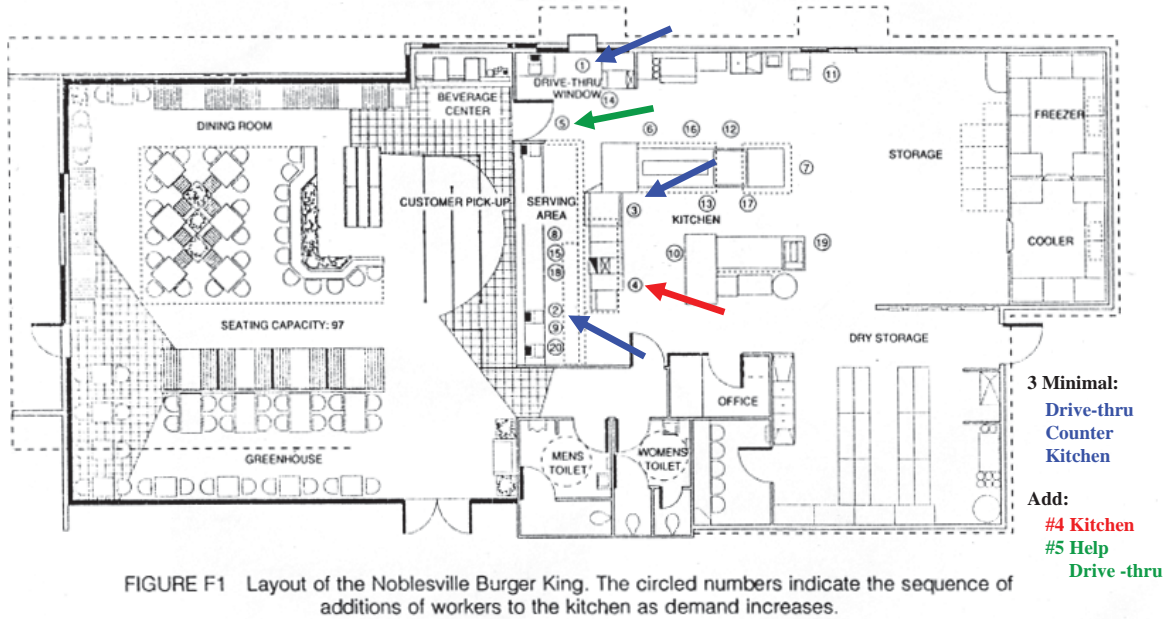
Emergency Department: Generic Flow



Conceptual Model: Burger King Bottlenecks

Bottleneck Analysis: Short – Run Approximations
Time – State Dependent Q-Net

TOUR F / A WORKER-PACED LINE FLOW PROCESS AND A SERVICE FACTORY 155



Analytical Models: Little's Law, or The First Law of Congestion

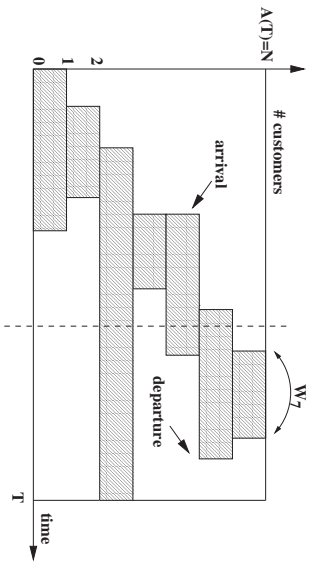


- λ = average arrival rate;
- L = average **number** within system;
- W = average **time** within system.

Little's Law

$$L = \lambda W$$

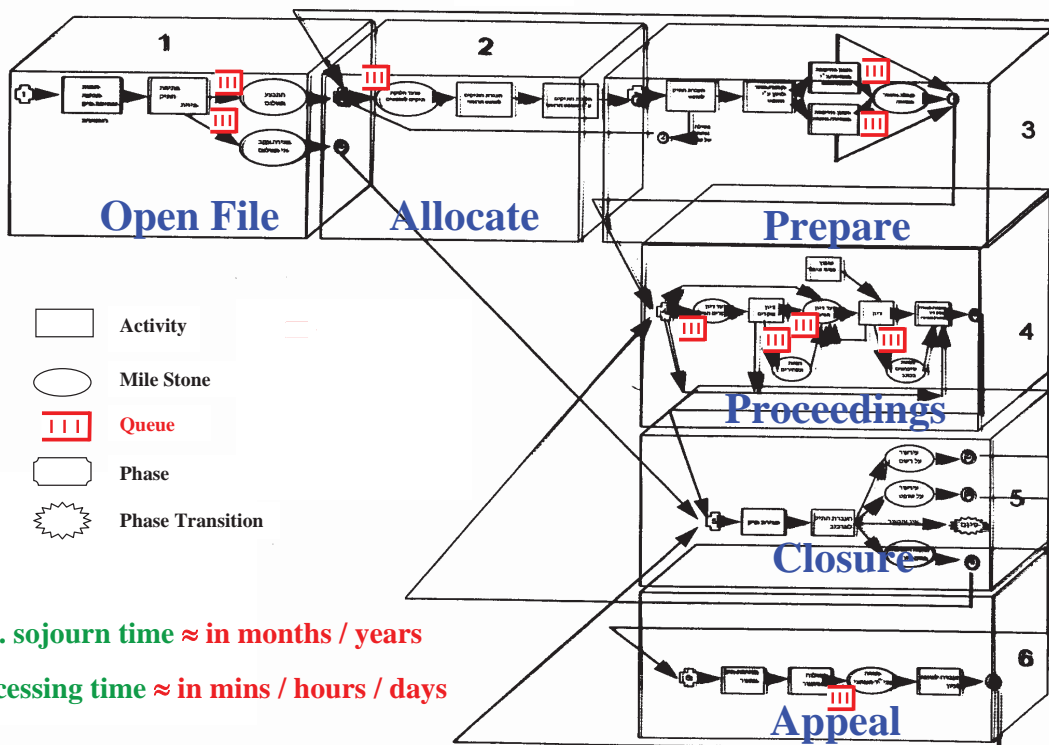
Finite-Horizon Version



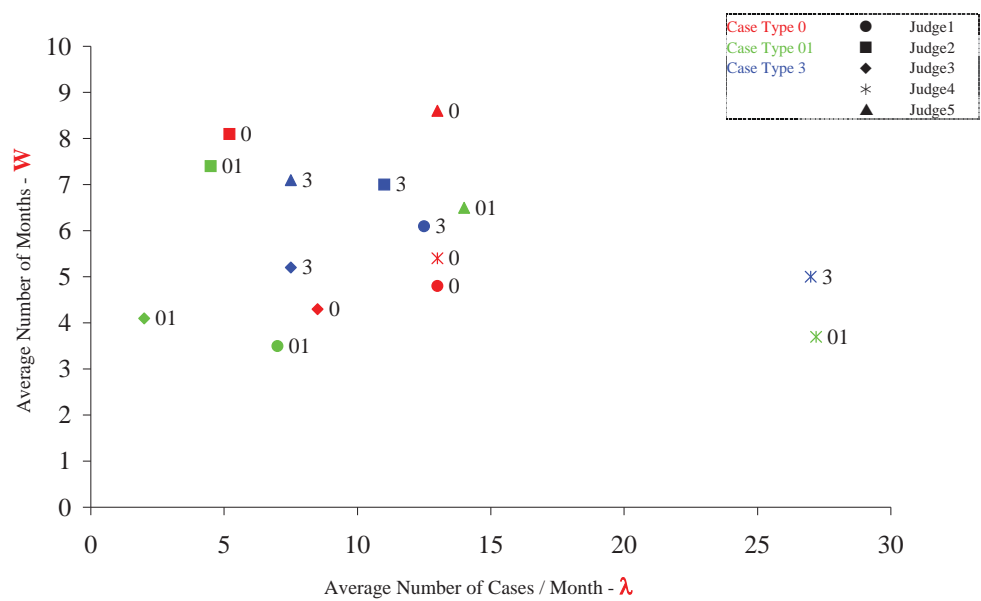
Long-Run (Stochastic) Example

$$M/M/1: L = \frac{\rho}{1 - \rho} = \frac{\lambda}{\mu - \lambda}, \quad W = \frac{1}{\mu - \lambda} = \frac{1}{\mu} \frac{1}{1 - \rho}.$$

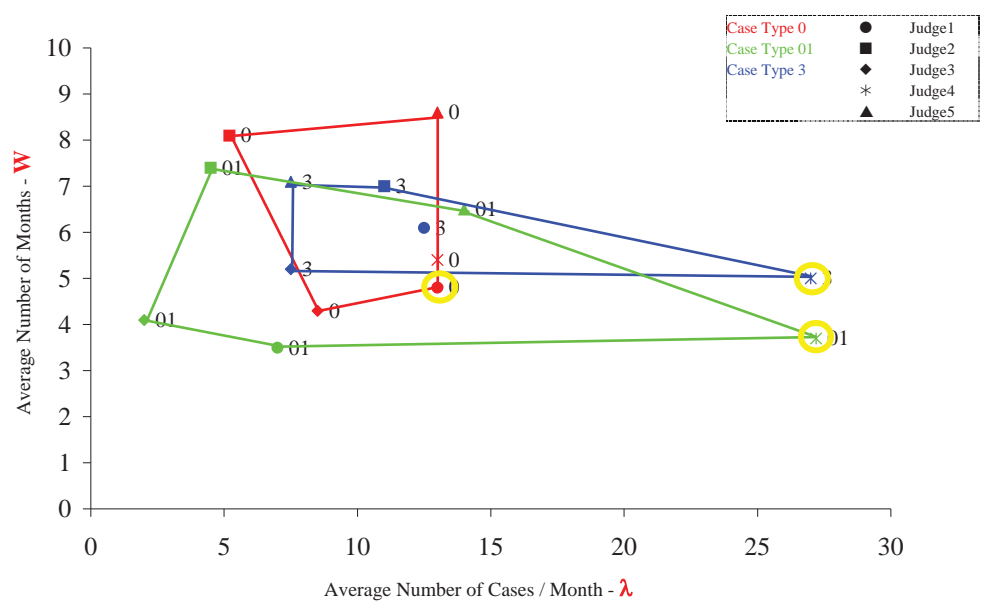
Conceptual Model: The Justice Network, or The Production of Justice



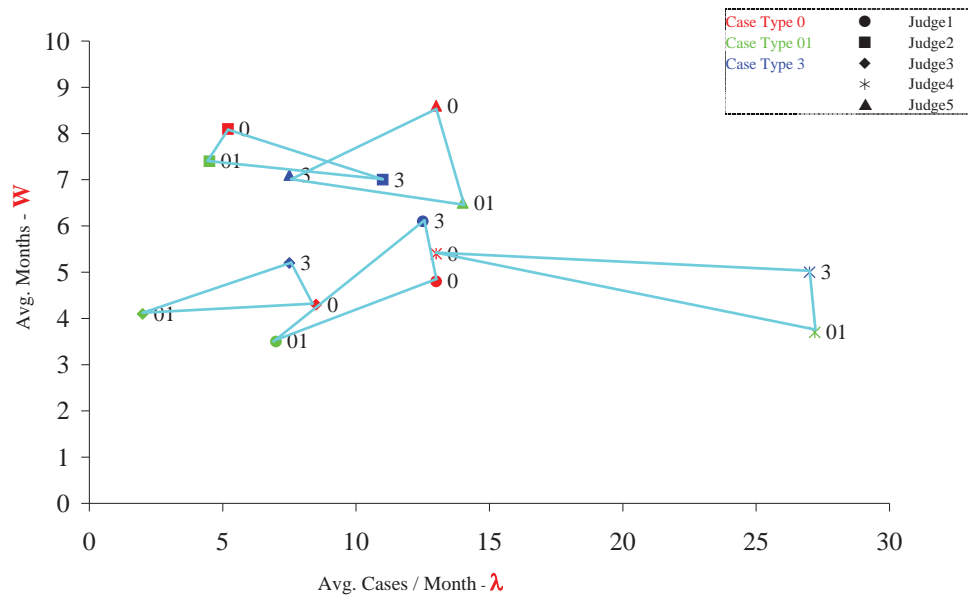
Judges: Operational Performance - Base case



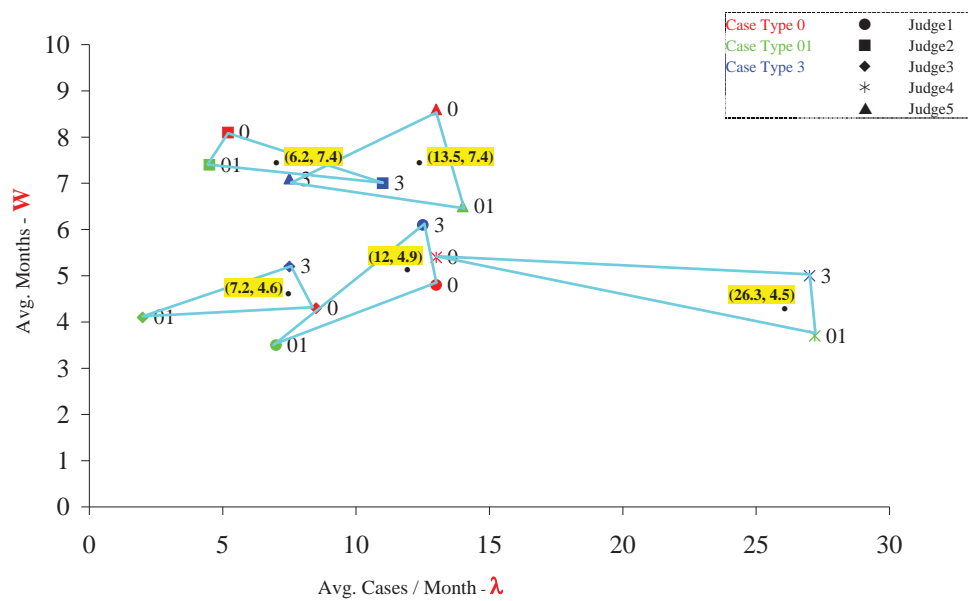
3 Case-Types: Performance by 5 Judges



5 Judges: Performance by 3 Case-Types



Judges: Performance Analysis



Conceptual Fluid Model

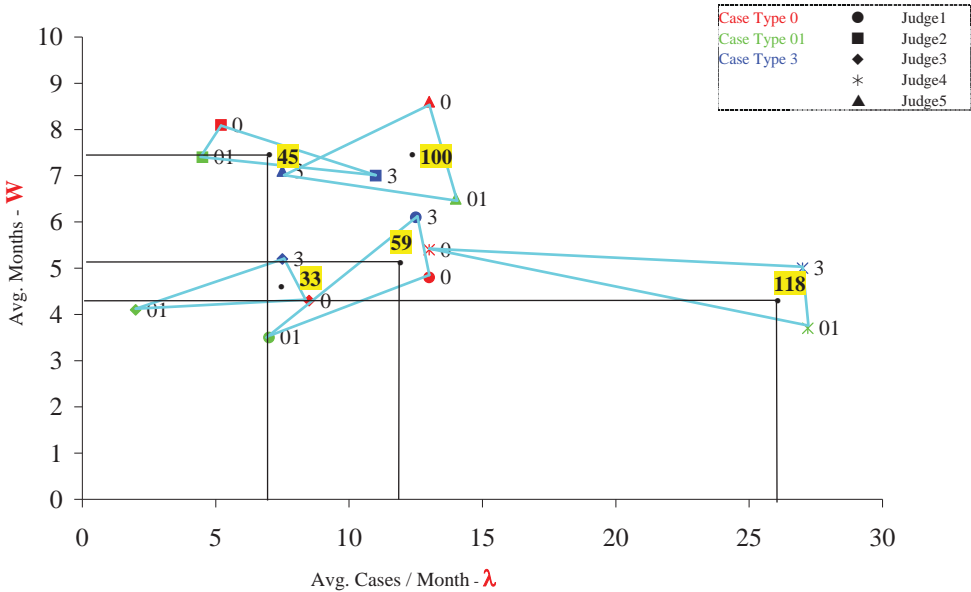
Customers/units are modeled by **fluid** (continuous) flow.

Labor-day Queueing at Niagara Falls

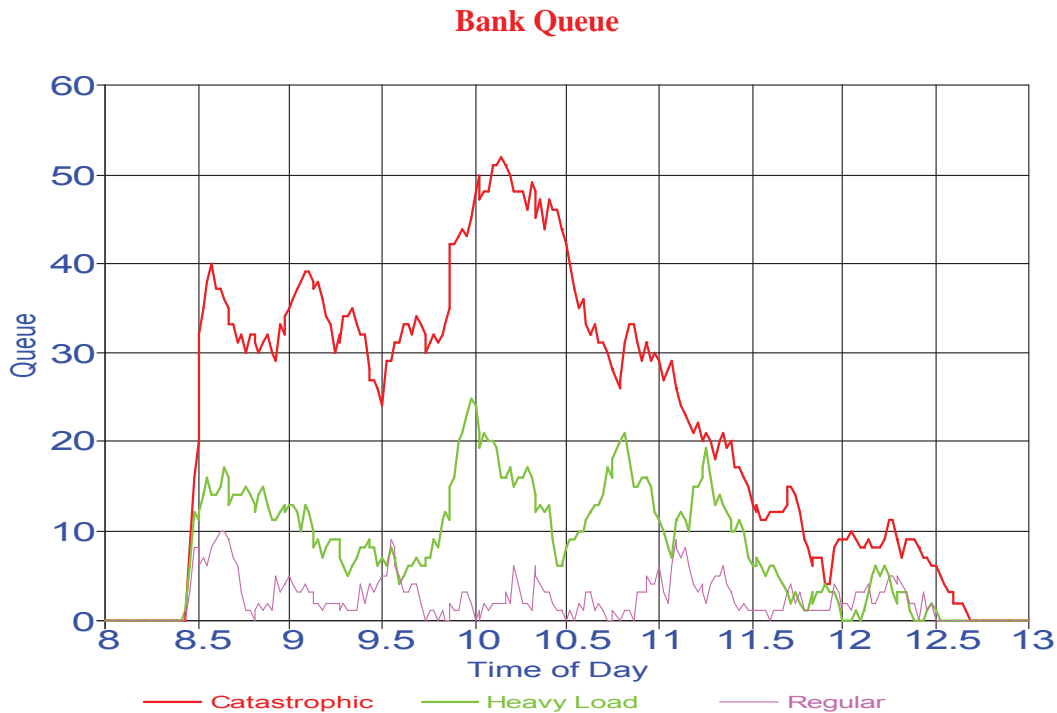


- Appropriate when **predictable variability** prevalent;
- Useful **first-order** models/approximations, often **suffice**;
- Rigorously justifiable via Functional Strong Laws of Large Numbers.

Judges: Best/Worst Performance



Empirical Fluid Model: Queue-Length at a Catastrophic/Heavy/Regular Day



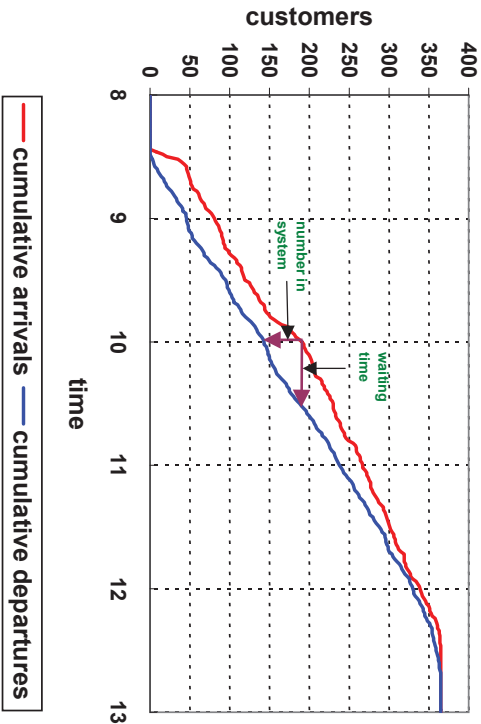
31

Empirical Models: Fluid, Flow

Derived directly from event-based (call-by-call) measurements. For example, an isolated service-station:

- $A(t)$ = **cumulative** # arrivals from time 0 to time t ;
- $D(t)$ = **cumulative** # departures from system during $[0, t]$;
- $L(t) = A(t) - D(t)$ = # customers in system at t .

Arrivals and Departures from a Bank Branch Face-to-Face Service



When is it possible to calculate waiting time in this way?

32

Mathematical Fluid Models

Differential Equations:

- $\lambda(t)$ – arrival rate at time $t \in [0, T]$.
- $c(t)$ – maximal potential processing rate.
- $\delta(t)$ – effective processing (departure) rate.
- $Q(t)$ – total amount in the system.

Then $Q(t)$ is a solution of

$$\dot{Q}(t) = \lambda(t) - \delta(t); \quad Q(0) = q_0, \quad t \in [0, T].$$

In a Call Center Setting (no abandonment)

$N(t)$ statistically-identical servers, each with service rate μ .

$c(t) = \mu N(t)$: maximal potential processing rate.

$\delta(t) = \mu \cdot \min(N(t), Q(t))$: processing rate.

$$\dot{Q}(t) = \lambda(t) - \mu \cdot \min(N(t), Q(t)), \quad Q(0) = q_0, \quad t \in [0, T].$$

How to actually solve? Mathematics (theory, numerical), or simply: Start with $t_0 = 0$, $Q(t_0) = q_0$.

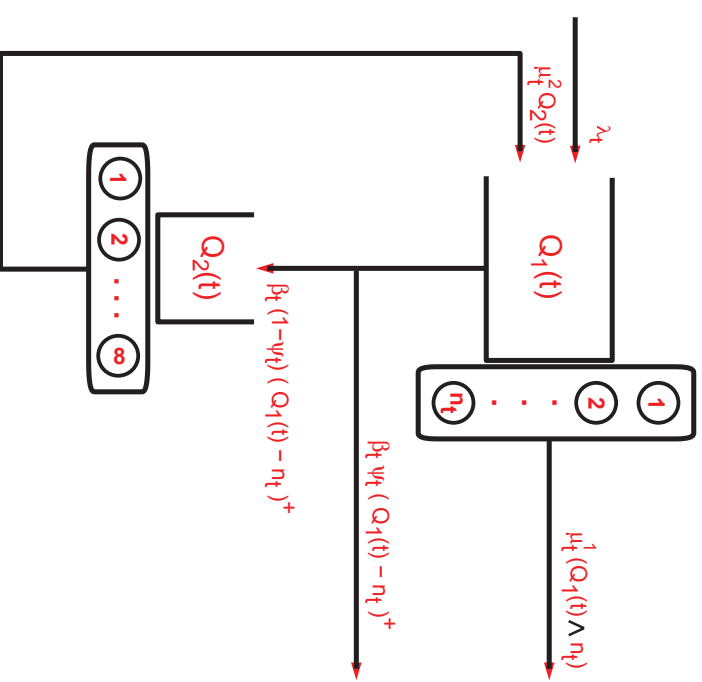
Then, for $t_n = t_{n-1} + \Delta t$:

$$Q(t_n) = Q(t_{n-1}) + \lambda(t_{n-1}) \cdot \Delta t - \mu \min(N(t_{n-1}), Q(t_{n-1})) \cdot \Delta t.$$

Time-Varying Queues with Abandonment and Retrials

Based on three paper with Massey, Reiman, Rider and Stolyar.

Call Center: a Multiserver Queue with Abandonment and Retrials



Primitives: Time-Varying Predictability

- λ_t exogenous arrival rate;
e.g., continuously changing, sudden peak.
- μ_t^1 service rate;
e.g., change in nature of work or fatigue.
- n_t number of servers;
e.g., in response to predictably varying workload.
- $Q_1(t)$ number of customers in call center
(queue+service).
- β_t abandonment rate while waiting;
e.g., in response to IVR discouragement
at predictable overloading.
- ψ_t probability of no retrial.
- μ_t^2 retrial rate;
if constant, $1/\mu^2$ – average time to retry.
- $Q_2(t)$ number of customers that will retry.

In our examples, we vary λ_t only, other primitives are constant.

Fluid Model

Replacing random processes by their rates yields

$$Q^{(0)}(t) = (Q_1^{(0)}(t), Q_2^{(0)}(t))$$

Solution to nonlinear differential balance equations

$$\begin{aligned} \frac{d}{dt} Q_1^{(0)}(t) &= \lambda_t - \mu_t^1 (Q_1^{(0)}(t) \wedge n_t) \\ &\quad + \mu_t^2 Q_2^{(0)}(t) - \beta_t (Q_1^{(0)}(t) - n_t)^+ \\ \frac{d}{dt} Q_2^{(0)}(t) &= \beta_1 (1 - \psi_t) (Q_1^{(0)}(t) - n_t)^+ \\ &\quad - \mu_t^2 Q_2^{(0)}(t) \end{aligned}$$

Justification: **Functional Strong Law of Large Numbers**,

with $\lambda_t \rightarrow \eta \lambda_t$, $n_t \rightarrow \eta n_t$.

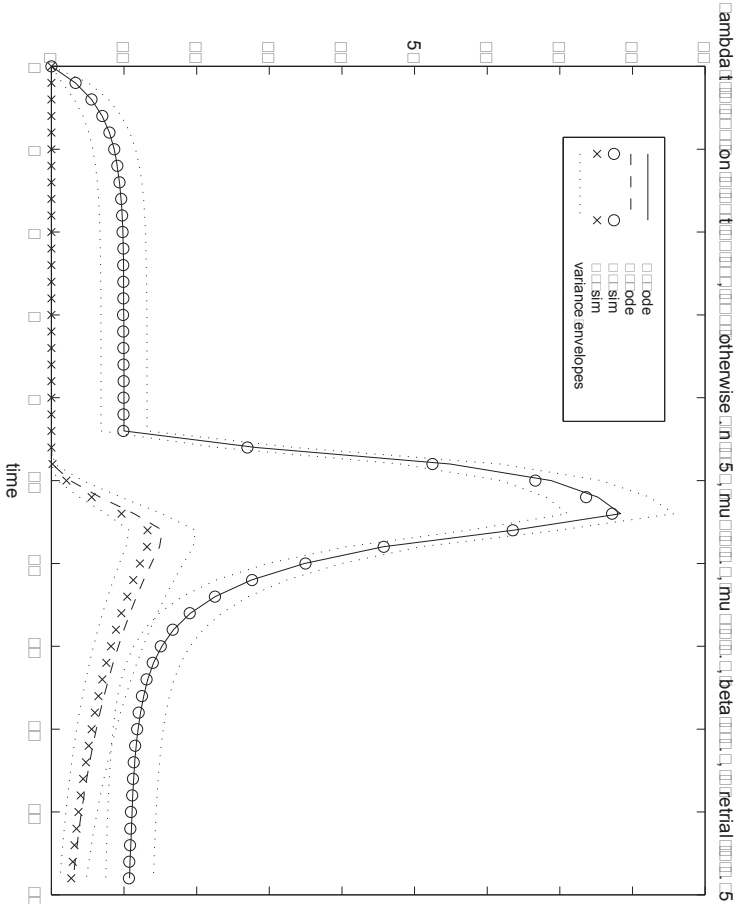
As $\eta \uparrow \infty$,

$$\frac{1}{\eta} Q^\eta(t) \rightarrow Q^{(0)}(t), \quad \text{uniformly on compacts, a.s.}$$

given convergence at $t = 0$

Sudden Rush Hour

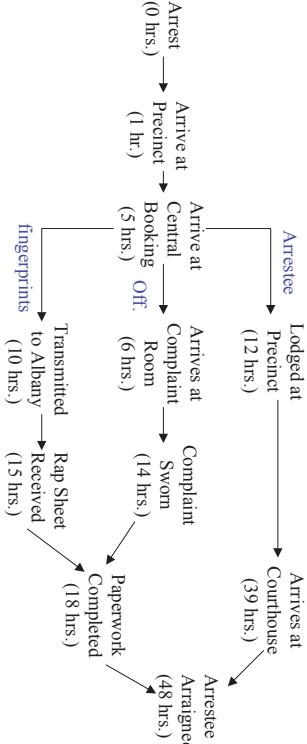
$$\begin{aligned}
 n &= 50 \text{ servers} & \mu &= 1 \\
 \lambda_t &= 110 & \text{for } 9 \leq t \leq 11, & \lambda_t = 10 \text{ otherwise}
 \end{aligned}$$



Stochastic Framework: DS PERT/CPM

DS = Dynamic Stochastic (Fork-Join, Split-Match)
PERT = Program **E**valuation and **R**eview **T**echnique
CPM = Critical Path Method
 Operations Research in Project Management. Standard Successful.

New-York Arrest-to-Arraignment System (Larson et al., 1993)



- CRM – task times are deterministic/averages (standard).
 S-PERT (Stochastic PERT) – task times random variables.
 DS-PERT/CPM – multi-project (dynamic) environment, with tasks processed at dedicated service stations.
- **Capacity analysis:** Can we do it? (LP)
 - **Response-time analysis:** How long will it take? (S-Nets)
 - **What if:** Can we do better? (Sensitivity, Parametric)
 - **Optimality:** What is the best one can do?

Stochastic Model of a Basic Service Station

Building blocks:

- Arrivals
- Service durations (times)
- Customers' (im)patience.
- Customers' returns (during service process, after service)

First **study** these building blocks one-by-one:

- Empirical analysis, which motivates
- Theoretical model(s).

Then **integrate** building blocks, via protocols, into (Basic) Models:

- Erlang-B/C (Arrivals, Services)
- Erlang-A (+ Abandonment), Erlang-R (+ Returns).

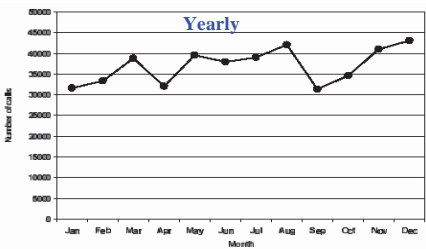
The models support, for example,

- Staffing Workforce, for which Basic Models are already useful; and beyond:
- Routing Customers
- Scheduling Servers
- Matching Customers-Needs with Servers-Skills (SBR).

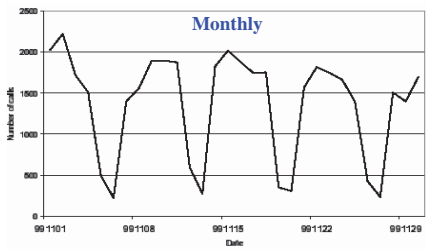
Arrivals to Service

Arrivals to a Call Center (1999): Time Scale

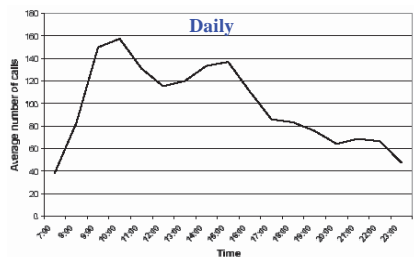
Strategic



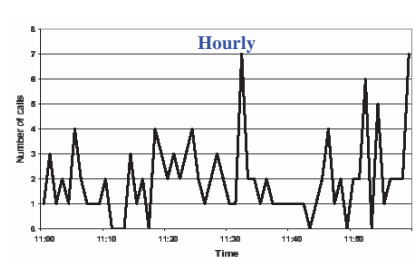
Tactical



Operational



Stochastic



Arrivals Process, in 1976

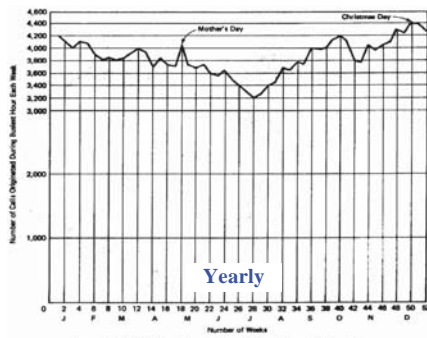


Figure 1 Typical distribution of calls during the busiest hour for each week during a year.

(E. S. Buffa, M. J. Cosgrove, and B. J. Luce, "An Integrated Work Shift Scheduling System")

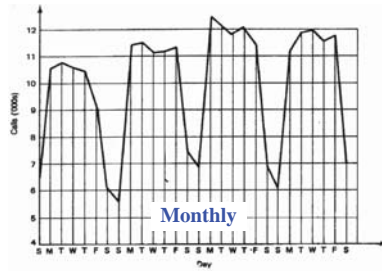


Figure 2 Daily call load for Long Beach, January 1972.

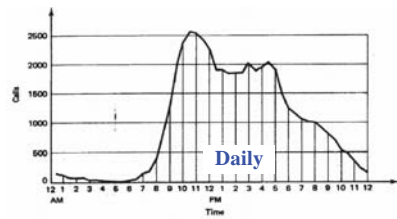


Figure 3 Typical half-hourly call distribution (Bundy O A).

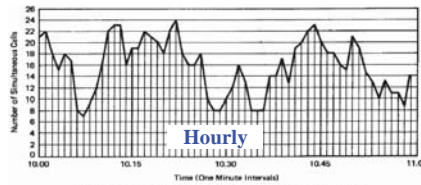


Figure 4 Typical intrahour distribution of calls, 10:00-11:00 A.M.

Q-Science: Predictable Variability

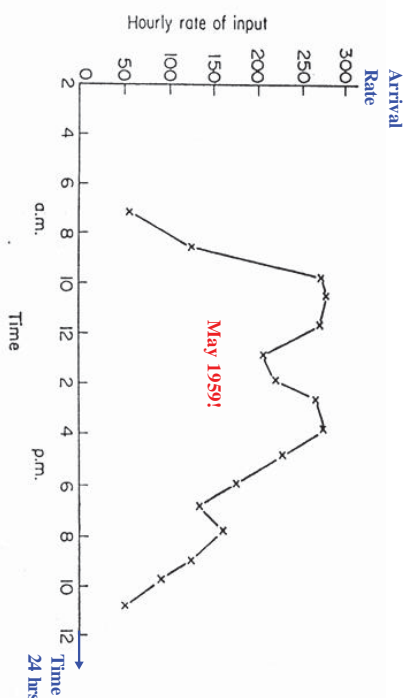
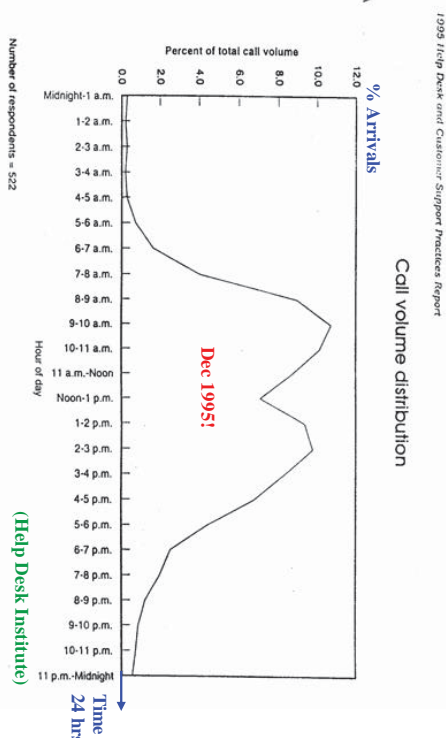
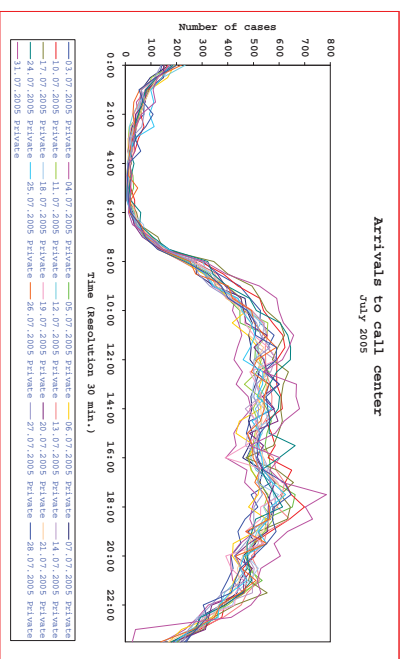


Fig. 15-1 The variation in the hourly input rates of reservations calls during a typical day (in May 1959) (Lee A.M., Applied Q-Th)



Arrivals to Service: Poisson Processes

Weekday Arrival Rates (Israeli CC, MOCCA)



- Arrivals over short (but not too short) intervals (15, 30 min) are close to homogeneous **Poisson**, with **over-dispersion**.
 - Arrivals over the day are (over-dispersed) **non-homogeneous Poisson**.
- Practice: model as **Poisson with piecewise-constant arrival rates**.
- Poisson Phenomena:**
- **PASTA** = Poisson Arrivals See Time Averages;
 - **Biased sampling:** Why is the service time we encounter upon arrival longer than a “typical” service time?

Arrivals to Service: Forecasting

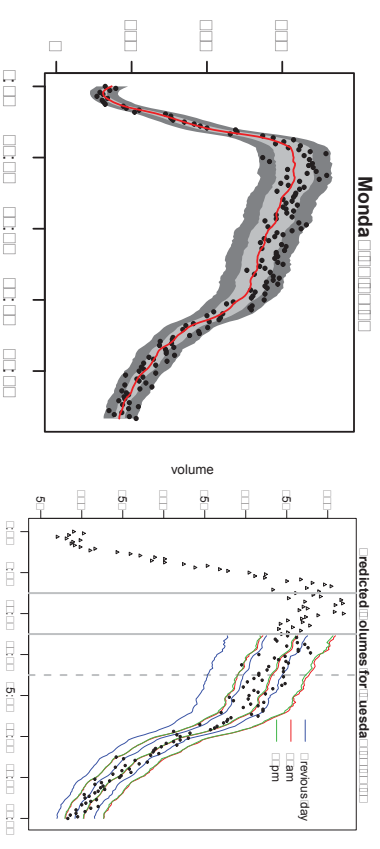
How to **predict** Poisson arrival rates? **Time Series** models. Days are divided into **time intervals** over which arrival rates are assumed **constant**.

Standard Resolutions: 15 min, 30 min, 1 hour.

N_{jk} = number of arrivals on day j during interval k . Assume K time intervals and J days overall.

- **One-day-ahead** prediction: N_1, \dots, N_{J-1} , known. Predict N_{J1}, \dots, N_{JK} .
- **Several days (weeks) ahead** prediction.
- **Within-day** prediction.

Forecast Accuracy (U.S. Bank, Weinberg)



Service Times (Durations)

<http://few3.technion.ac.il/serveng/Lectures/ServiceFull.pdf>

Why Significant? +1 second of 1000 agents costs \$500K yearly.

Why Interesting?

Must accurately **Model, Estimate, Predict, Analyze**:

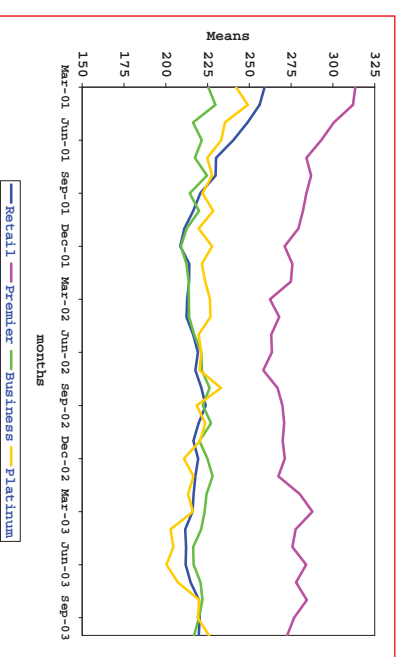
- Resolution: Sec's (phone)? min's (email)? hr's (hospital)
- Parameter, Distribution (Static) or Process (Dynamic)?
- Does it include after-call work?
- Does it include interruptions?
 - Whisper time, hold time, phones during face-to-face,...
- Does it account for return services?

How affected by **covariates**?

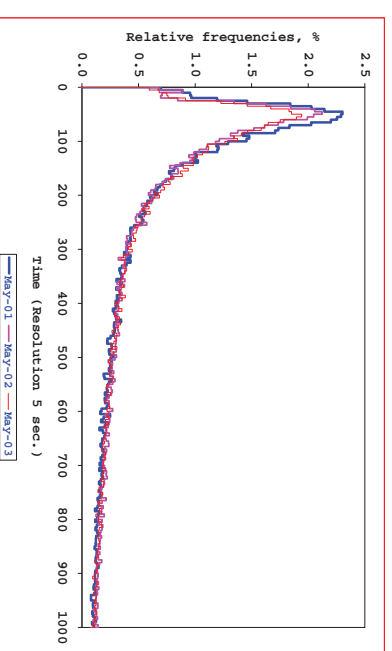
- Experience and Skill of agents (Learning Curve)
- Type of Customer: Service Type, VIP Status
- Time-of-Day: Congestion-Level
- Human Factor: Incentives, pending workload, fatigue

Service Times: Trends and Stability

Average Customer Service Time, Weekdays (MOCCA)



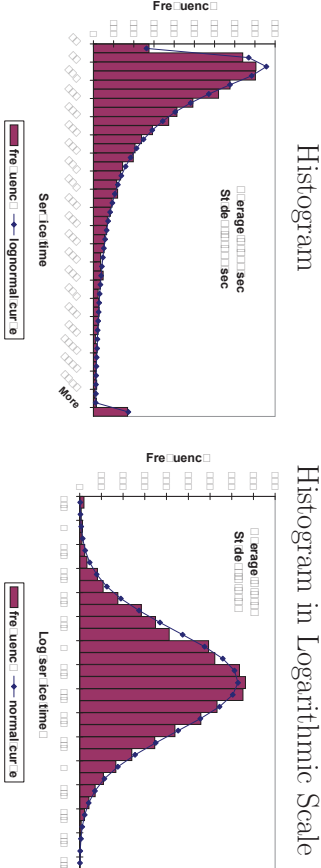
USBank Service-Time Histograms for Telesales (MOCCA)



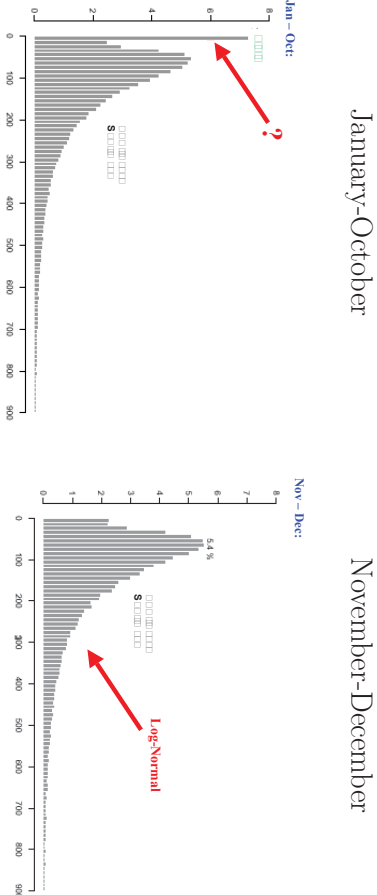
Service Times: Static Models, or
Averages Do Not Tell the Whole Story

Distributions: Parametric (Exponential, Lognormal),
Semi-Parametric (Phase-Type), Non-Parametric (Empirical).

Lognormal Service Times in an Israeli Bank

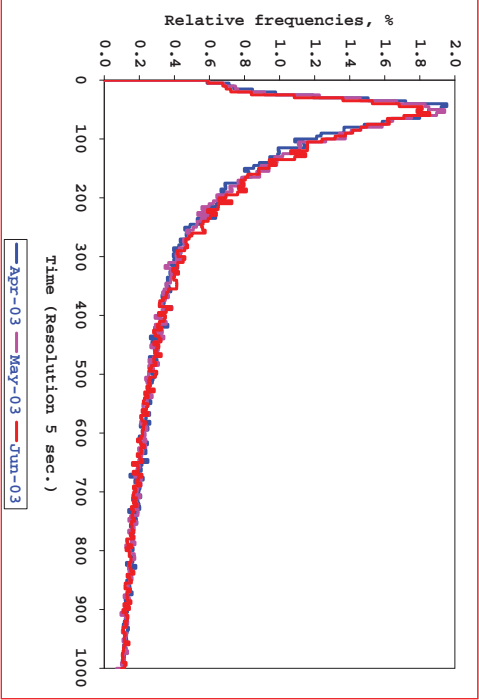
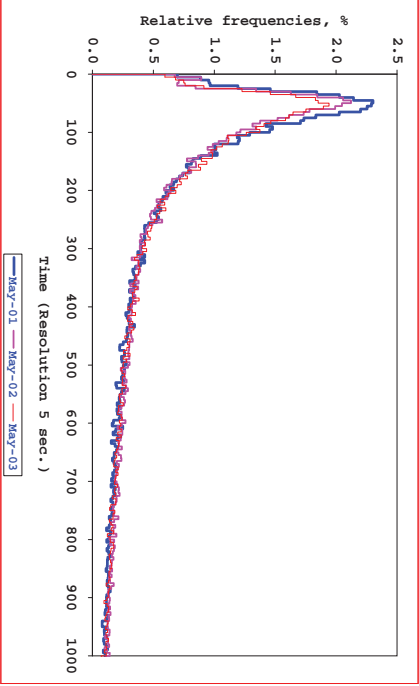


A Typical Call Center?



Service Times: 5 Sec's Resolution

USBank. Service-Time Histograms for Telesales (MOCCA)



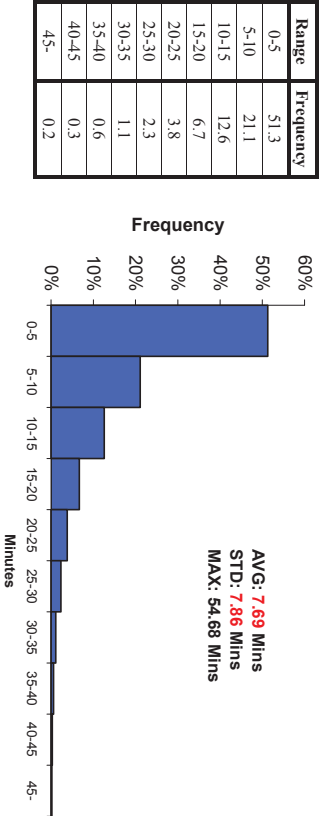
Local Municipalities

Department	Station No.	Total Customers	Avg. Arrival Rate (1/Hr)	Avg. Service Time (Mins)	STD (Mins)	Maximal Service Time (Mins)	Utilization	Avg. Waiting Time (Mins)
Water	N/A	187	1.8 ± 0.2	8.87 ± 1.0	8.15	54.68	13.3%	4.76
Tellers	N/A	1328	12.6 ± 0.5	8.82 ± 0.4	8.55	49.37	30.8%	7.73
Cashier	N/A	757	7.2 ± 0.4	6.64 ± 0.4	6.94	29.95	79.7%	3.89
Manager	N/A	190	1.8 ± 0.2	7.99 ± 1.0	8.44	38.97	24.1%	9.16
Discounts	N/A	317	3.0 ± 0.3	4.59 ± 0.4	4.54	36.72	23.1%	3.65

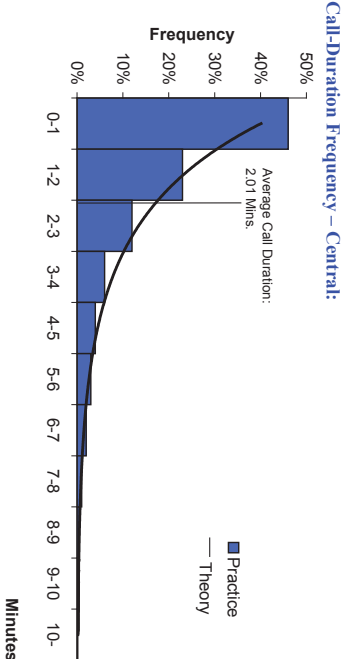
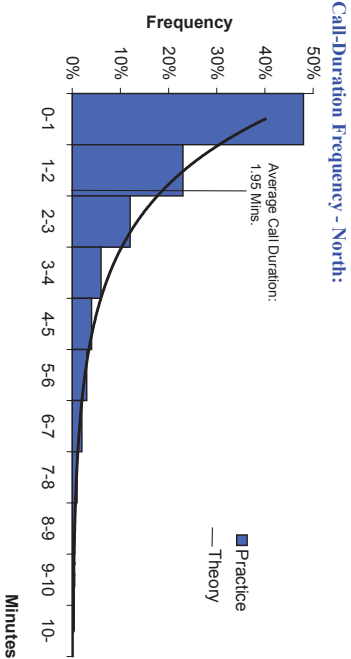
Water	1	57	N/A	7.80 ± 1.70	7.61	31.28	6.5%	N/A
	2	130	N/A	9.34 ± 1.20	8.37	54.68	19.3%	N/A
	3	336	N/A	9.04 ± 0.80	8.93	49.05	48.2%	N/A
	4	208	N/A	9.93 ± 1.00	8.82	49.12	33.0%	N/A
Tellers	5	417	N/A	8.97 ± 0.70	8.55	49.37	59.4%	N/A
	6	144	N/A	9.53 ± 1.20	8.75	41.70	21.8%	N/A
	7	156	N/A	8.03 ± 1.10	7.96	35.27	19.8%	N/A
	8	67	N/A	3.74 ± 0.70	3.58	21.03	4.0%	N/A
Cashier	9	757	N/A	6.64 ± 0.40	6.94	29.95	79.7%	N/A
Manager	10	190	N/A	1.99 ± 1.00	8.44	38.97	24.1%	N/A
Discounts	11	317	N/A	4.59 ± 0.40	4.54	36.72	23.1%	N/A

*Service time ranges given with 90% confidence.

Service Time Histogram – Overall:



Service Times: Exponential (Phone Calls)

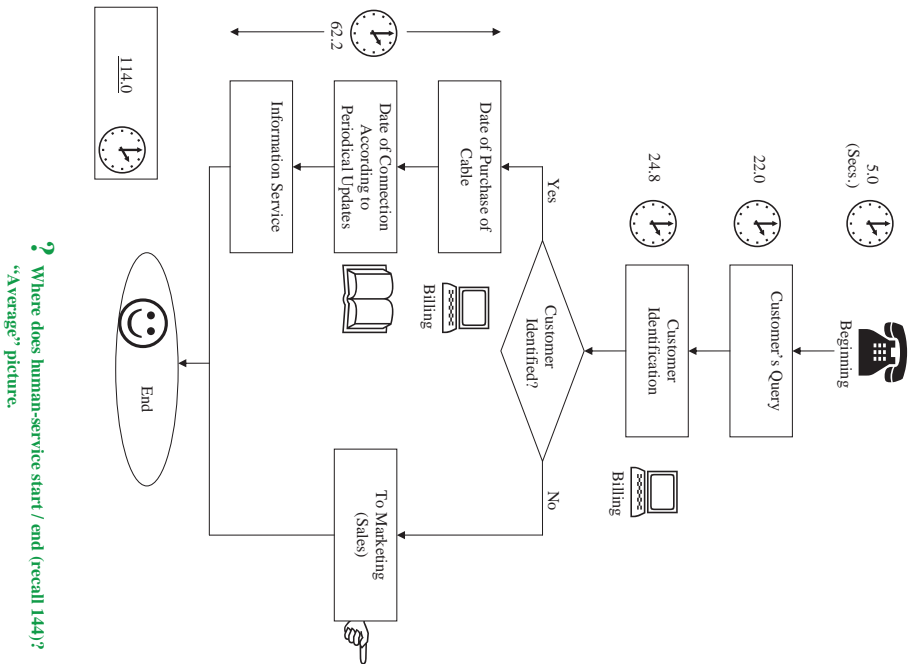


Q. How to recognize “Exponential” when you “see” one?

A. Geometric Approximation.

Service Times: Phase-Type Model

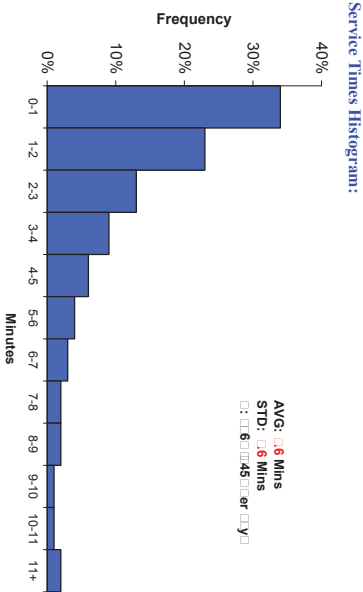
Late Connections



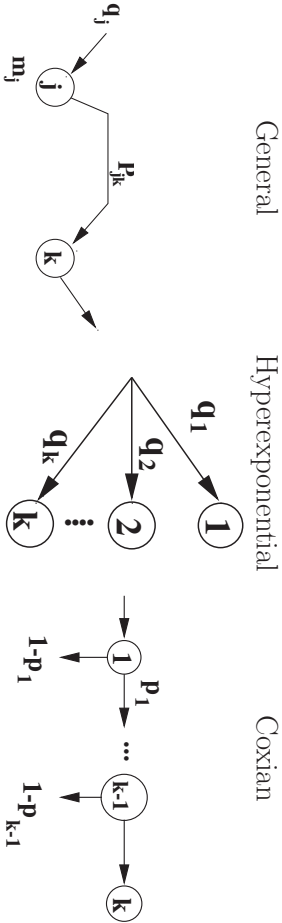
Service Times: Exponential, Phase-Type

Static Model: Exponential Duration

Face-to-Face Services in a Government Office

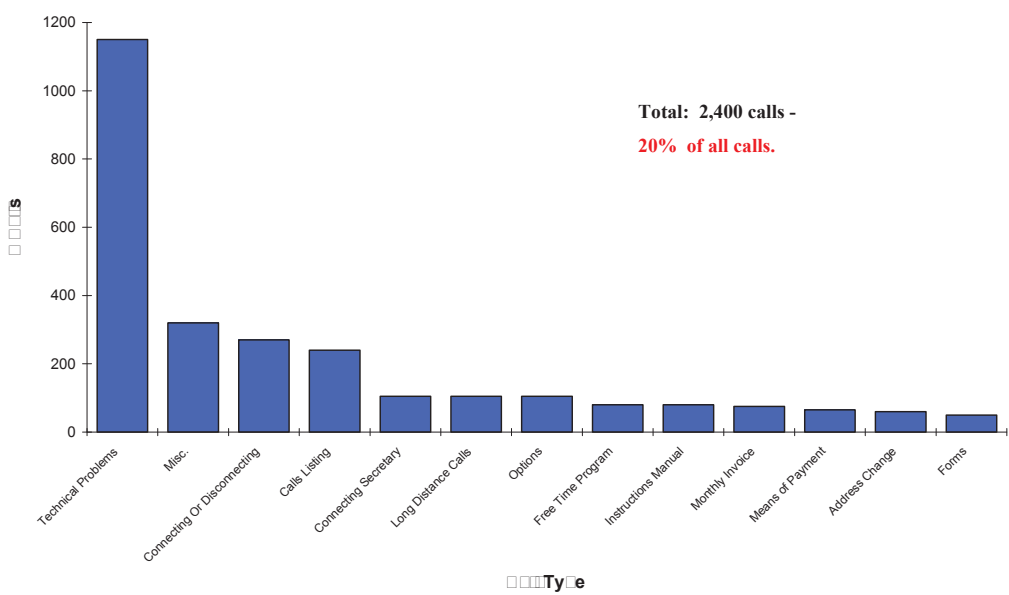


Dynamic Model: Phase-Type Duration



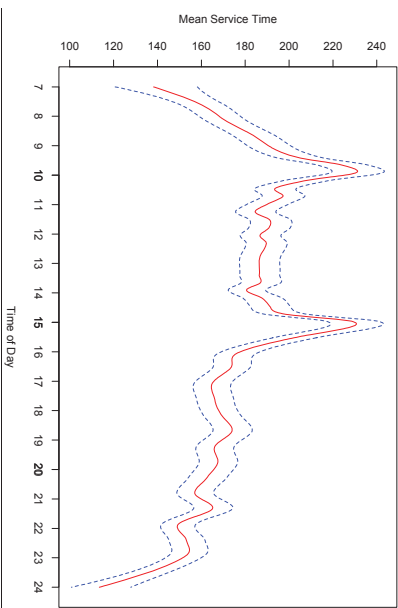
Service Times: Returns

Bank Classification of “Continued – Calls”

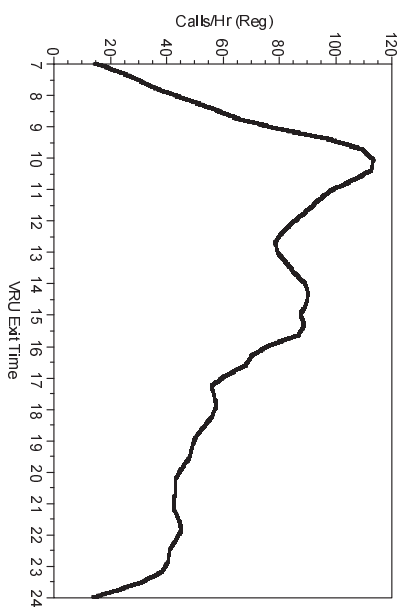


Service Times: The Human Factor, or Why Longest During Peak Loads?

Mean-Service-Time (Regular) vs. Time-of-Day (95% CI) (n=42613)



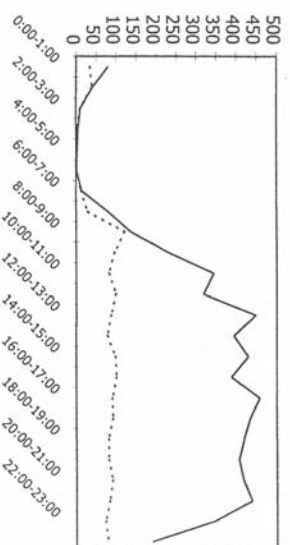
Arrivals to Queue or Service - Regular Calls (Inhomogeneous Poisson)



Customers' (Im)Patience

Marketing Campaign at a Call Center

Average wait 376 sec, 24% calls answered



Abandonment **Important and Interesting**

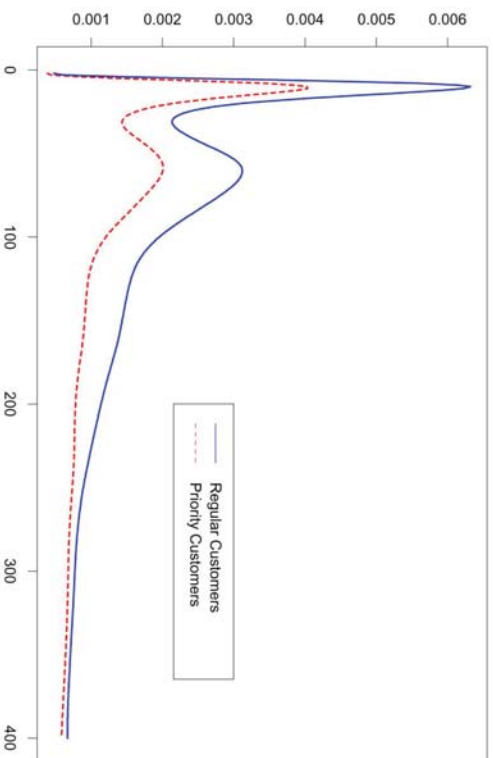
- One of two customer-subjective performance measures (2nd=Redials)
- Poor service level (future losses)
- Lost business (present losses)
- 1-800 costs (present gains; out-of-pocket vs. alternative)
- Self-selection: the “fittest survive” and wait less (much less)
- Accurate Robust models (vs. distorted instability-prone)
- Beyond Operations/OR: Psychology, Marketing, Statistics
- Beyond Telephony: VRU/IVR (Opt-Out-Rates), Internet (over 60%), Hospitals ED (LWBS).

Understanding (Im)Patience

- **Observing** (Im)Patience – Heterogeneity:
Under a single roof, the fraction abandoning varies from 6% to 40%, depending on the type of service/customer.
- **Describing** (Im)Patience Dynamically:
Irritation proportional to Hazard Rate (Palm's Law).
- **Managing** (Im)Patience:
 - VIP vs. Regulars: who is more “Patient”?
 - What are we actually measuring?
 - (Im)Patience Index:
“How long *Expect* to wait” relative to
“How long *Willing* to wait”.
- **Estimating** (Im)Patience: Censored Sampling.
- **Modeling** (Im)Patience:
 - The “Wait” Cycle:
Expecting, Willing, Required, Actual, Perceived, etc.
The case of the *Experienced & Rational* customer.
 - (Nash) Equilibrium Models.

Palm's Law of Irritation (1943-53): \propto Hazard-Rate of (Im)Patience Distribution

Small Israeli Bank (1999):
 Regular over **Priority (VIP)** Customers



Hazard-Rate function of $\tau \geq 0$ (absolutely continuous):

$$h(t) = \frac{g(t)}{1 - G(t)},$$

g = Density function of τ ,

G = Distribution function of τ .

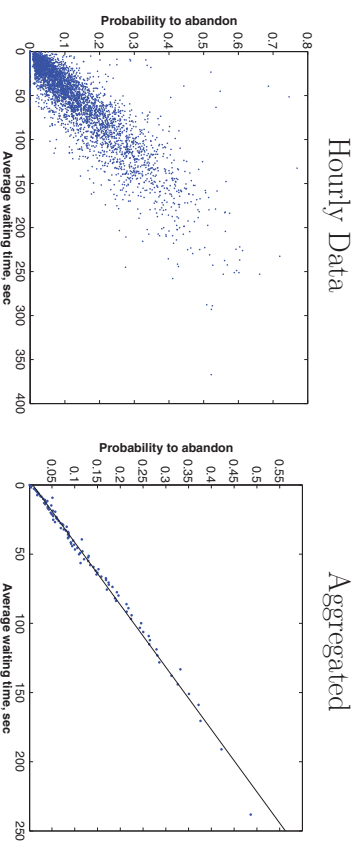
Intuition: $P\{\tau \leq t + \Delta | \tau > t\} \approx h(t) \cdot \Delta$.

$$P\{Ab\} \propto E[W_q]$$

Claim: (Im)Patience that is $\exp(\theta)$ implies

$$P\{Ab\} = \theta \cdot E[W_q].$$

Small Israeli Bank: 1999 Data



The graphs are based on 4158 hour intervals.

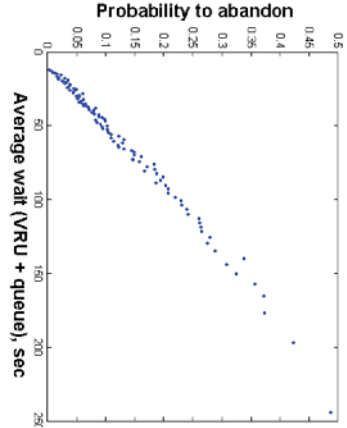
Regression \Rightarrow average patience $(1/\theta) \approx \frac{250}{0.56} \approx 446$ sec.

But (im)patience at this bank is **not** exponential ! ?

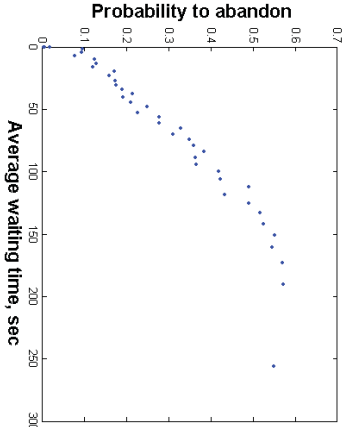
Moreover,

Queueing Science: Human Behavior

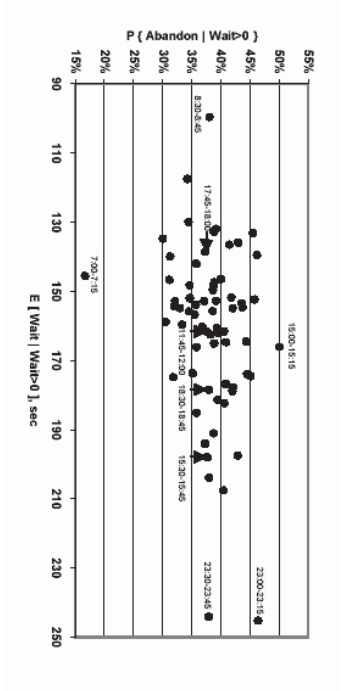
Delayed Abandons (IVR)



Balking (New Customers)

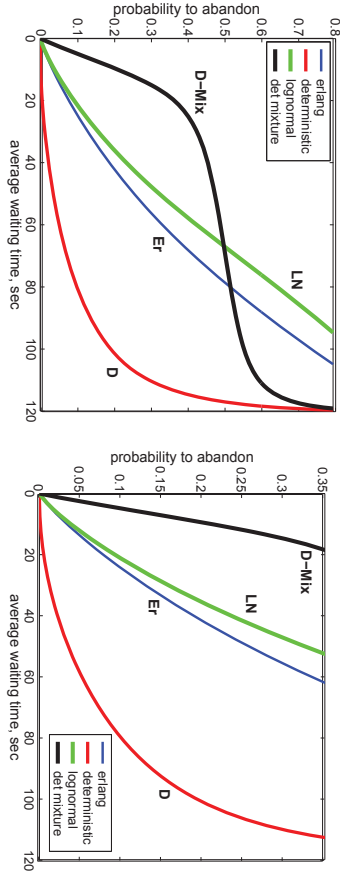


Learning (Internet Customers)



Examples of non-linear relations

moderate loads



Patience distributions:

- D: Deterministic: 2 minutes exactly;
- Er: Erlang with two exp(mean=1) phases;
- LN: Lognormal, both average and standard deviation equal to 2;
- D-Mix: 50-50% mixture of two constants: 0.2 and 3.8.

A Patience Index

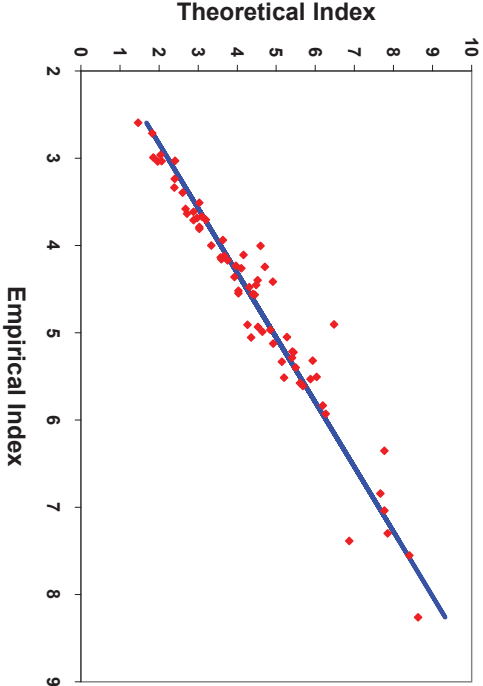
How to quantify (im)patience?

Theoretical Patience Index = $\frac{\text{Willing to Wait}}{\text{Expected to Wait}}$.

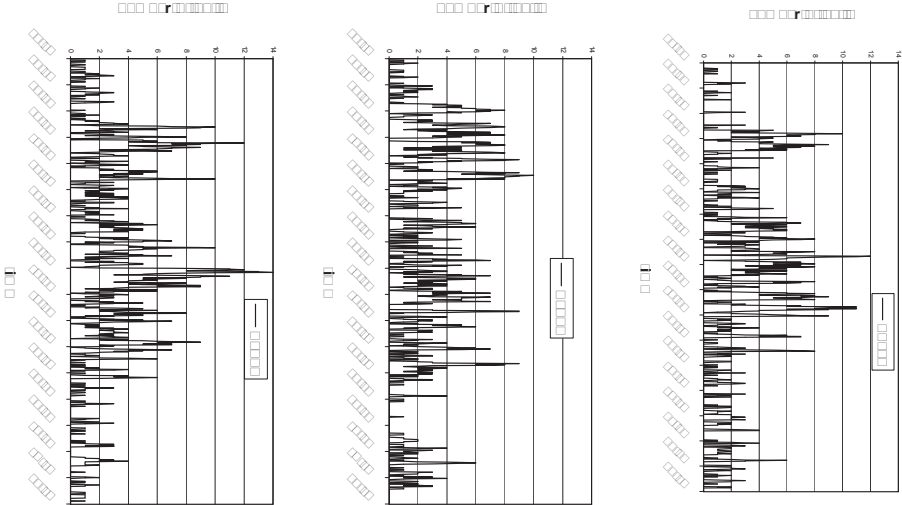
How to measure? Calculate? Assume Experienced customers. Then, a simple (but not too simple) model suggests the easy-to-measure:

Empirical Patience Index $\triangleq \frac{\% \text{ Served}}{\% \text{ Abandoned}}$.

Patience index – Empirical vs. Theoretical



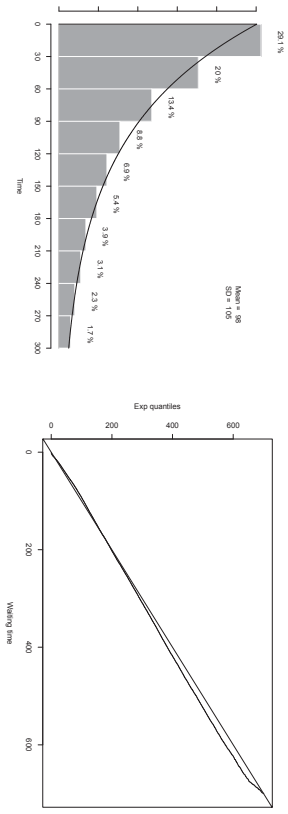
Queues = Integrating the Building Blocks



Delays = Integrating the Building Blocks

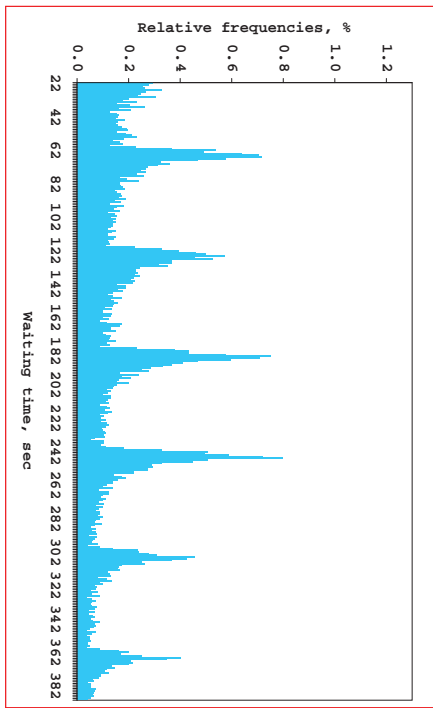
Exponential Delays:

Small Call Center of an Israeli Bank (1999)



Delays:

Medium-Size Call Center of an Israeli Bank (2006)

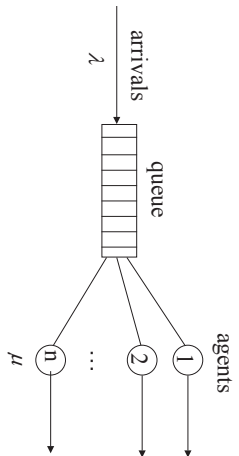


Basic (Markovian) Queueing Models of a Basic Service Station

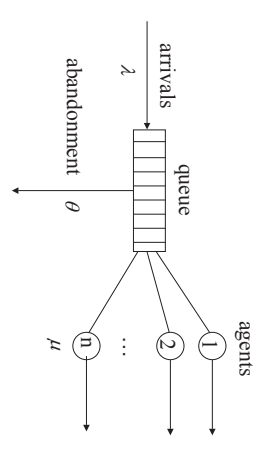
Poisson arrivals, Exponential service times, Exponential (im) patience.

Mathematical Framework: Markov Jump-Processes (Birth&Death).

M/M/n (Erlang-C) Queue



M/M/n+M (Palm/Erlang-A) Queue



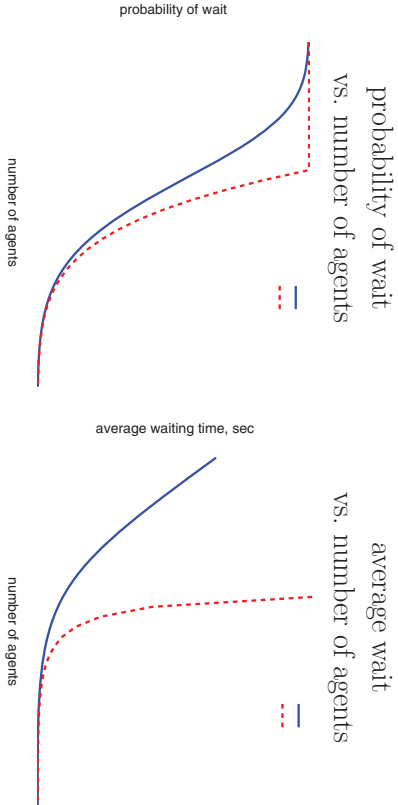
Additional Markovian Models: Balking; Trunks; Retrials.

Applications: Performance Analysis, Design (EOS), Staffing.

”The Fittest Survive” and Wait Less - Much Less!

Erlang-A vs. Erlang-C

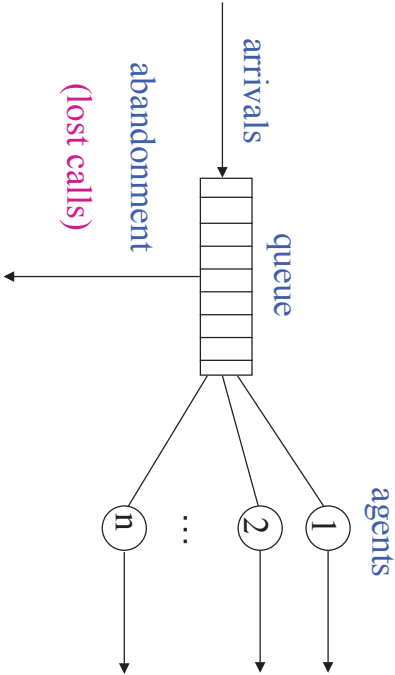
48 calls per min, 1 min average service time,
 2 min average patience



If 50 agents:

	M/M/n	M/M/n+M	M/M/n, $\lambda \downarrow 3.1\%$
Fraction abandoning	–	3.1%	–
Average waiting time	20.8 sec	3.7 sec	8.8 sec
Waiting time's 90-th percentile	58.1 sec	12.5 sec	28.2 sec
Average queue length	17	3	7
Agents' utilization	96%	93%	93%

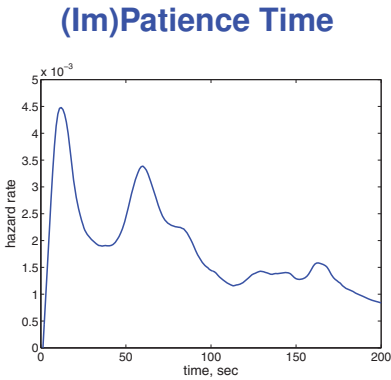
Modelling (Im)Patience: Time Willing vs. Time Required to Wait



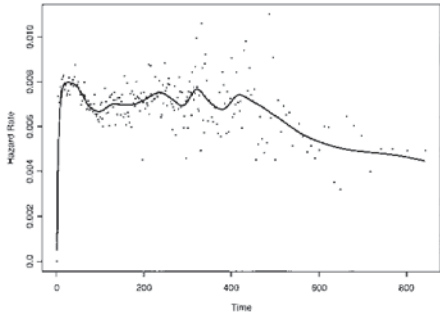
- **(Im)Patience Time $\tau \sim G$:**
 Time a customer **willing to wait** for service.
- **Offered Wait V :**
 Time a customer **required to wait** for service;
 in other words, waiting-time of an infinitely-patient customer.
- If $\tau \leq V$, customer **Abandons**;
 otherwise, customer **Served**.
- **Actual wait $W = \min(\tau, V)$.**

Call Center Data: Hazard Rates (Un-Censored)

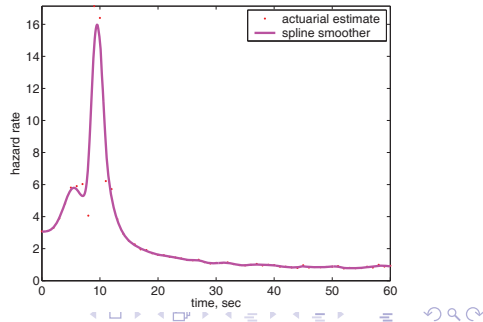
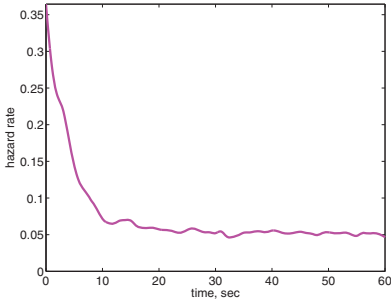
Israel



Required/Offered Wait



U.S.



30

Predicting Performance

Model **Primitives** (eg. Erlang-A):

- Arrivals to service (eg. Poisson)
- (Im)Patience while waiting τ (eg. Exp)
- Service times (eg. Exp)
- Number of Agents.

Model **Output**: **Offered-Wait** V

Operational Performance Measure calculable in terms of (τ, V) .

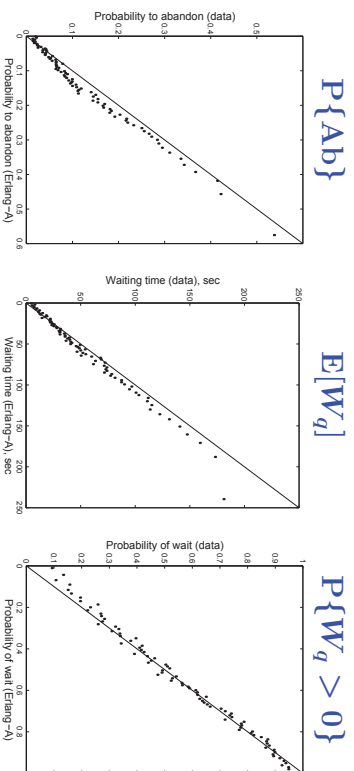
- eg. Average Wait = $E[\min\{\tau, V\}]$
- eg. % Abandonment = $P\{\tau < V\}$

Applications:

- **Performance Analysis**
- **Design, Phenomena** (Pooling, Economies of Scale)
- **Staffing – How Many Agents** (FTE's = Full-Time-Equivalent's)
- Note: Control requires model-refinements - later, in SBR.

Erlang-A: A Simple Model at the Service of Complex Realities

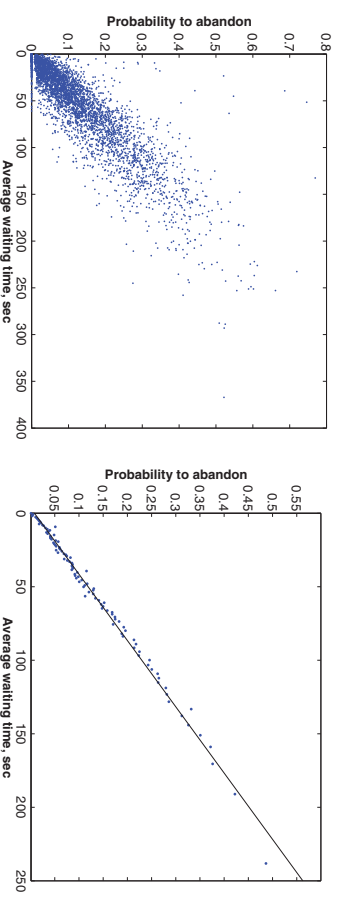
- Small Israeli bank (10 agents);
- Data-Based Estimation of Patience ($P\{Ab\}/E[W_q]$);
- Graph: Actual Performance vs. Erlang-A Predictions (aggregation of 40 similar hours).



- **Question:** Why Erlang-A works? indeed, all its underlying assumptions fail (Arrivals, Services, Impatience)
- **Towards a Theoretical Answer:** Robustness and Limitations, via Asymptotic (QED) Analysis.
- **Practical Significance:** Asymptotic results applicable in small systems (eg. healthcare).

Queueing Science: In Support of Erlang-A

Israeli Bank: Yearly Data
Hourly Data Aggregated



Data: $P\{Ab\} \propto E[W_q]$.

Theory: $P\{Ab\} = \theta \cdot E[W_q]$, **if** (Im)Patience = $\text{Exp}(\theta)$.

Proof: Let λ = Arrival Rate. Then, by Conservation & Little:

$$\lambda \cdot P\{Ab\} = \theta \cdot E[L_q] = \theta \cdot \lambda \cdot E[W_q], \quad \text{q.e.d.}$$

Recipe: Use Erlang-A, with $\hat{\theta} = P\{\widehat{Ab}\}/E[\widehat{W}_q]$ (slope above).

But (Im)Patience is **not** Exponentially distributed !?

Queueing Science: via Data & Theory, Linearity Robust.

Service Engineering: via Theory & Simulations, often-enough,

- Reality $\approx M/G/n + G \approx \text{Erlang-A}$, in which $\theta = g(0)$;
- $P\{Ab\} \approx g(0) \cdot E[W_q]$, hence **recipe prevails, often enough**.

4CallCenters: Personal Tool for Workforce Management

Calculations based on the M.Sc. thesis of Ofer Garnett.

Is extensively used in Service Engineering.

Install at

<http://ie.technion.ac.il/serveng/4CallCenters/Downloads.htm>

4CallCenters: Output Example

4CallCenters v2.01

File Table Settings Help

Performance Profiler Starting Query Advanced Profiling Advanced Queries What-if Analysis

Advanced Queries
center's parameters - pressing 'Compute' will find the value(s) of this parameter for which all your goals are met.

Compute Add to Table Delete Rows Clear All Export Graph Settings

Goals	Query	Input	Target	Number of Agents	Average Handling Time	Calls per Interval	Average Patience	Agents' Occupancy	%Abandon	Average Time in Queue	%Answer within Target
Multi-Value		00:20	04:00	Range	05:00			3%		80%	
Upper		00:20.0	10.0	04:00.0	100.0	05:00.0		65.3%	2.0%	00:06.0	90.1%
1		00:20.0	13.0	04:00.0	150.0	05:00.0		74.7%	2.9%	00:08.7	85.0%
2		00:20.0	17.0	04:00.0	200.0	05:00.0		76.7%	2.3%	00:06.8	87.4%
3		00:20.0	20.0	04:00.0	250.0	05:00.0		81.0%	2.8%	00:08.3	84.2%
4		00:20.0	24.0	04:00.0	300.0	05:00.0		81.5%	2.2%	00:06.6	86.8%
5		00:20.0	27.0	04:00.0	350.0	05:00.0		84.2%	2.5%	00:07.6	84.5%
6		00:20.0	30.0	04:00.0	400.0	05:00.0		86.3%	2.9%	00:08.6	82.4%
7		00:20.0	34.0	04:00.0	450.0	05:00.0		86.2%	2.3%	00:07.0	85.2%
8		00:20.0	37.0	04:00.0	500.0	05:00.0		87.8%	2.6%	00:07.8	83.5%
9		00:20.0	40.0	04:00.0	550.0	05:00.0		89.1%	2.8%	00:08.5	81.9%
10		00:20.0	44.0	04:00.0	600.0	05:00.0		88.8%	2.4%	00:07.1	84.5%
11		00:20.0	47.0	04:00.0	650.0	05:00.0		89.8%	2.6%	00:07.7	83.1%
12		00:20.0	50.0	04:00.0	700.0	05:00.0		90.9%	2.9%	00:08.3	81.5%

Settings
Parameters
Indicators

4CallCenters: Congestion Curves

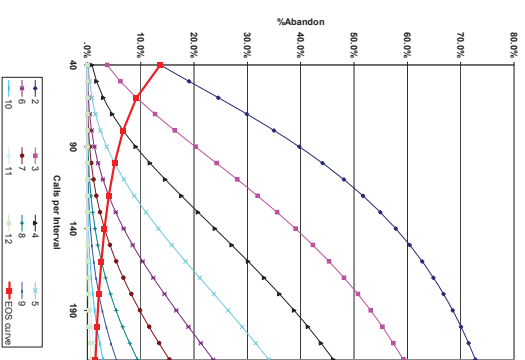
Vary input parameters of Erlang-A and display output (performance measures) in a table or graphically.

Example: $1/\mu = 2$ minutes, $1/\theta = 3$ minutes;

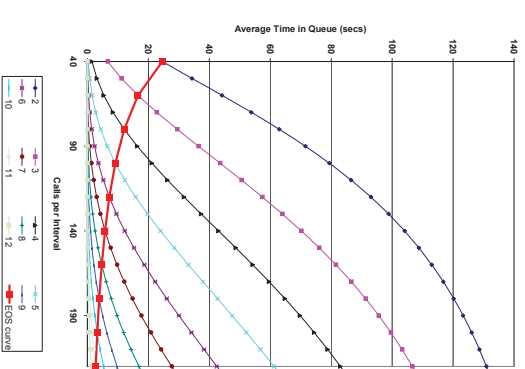
λ varies from 40 to 230 calls per hour, in steps of 10;

n varies from 2 to 12.

Probability to abandon



Average wait



Red curve: offered load per server fixed.

EOS (Economies-Of-Scale) observed.

Why the two graphs are similar?

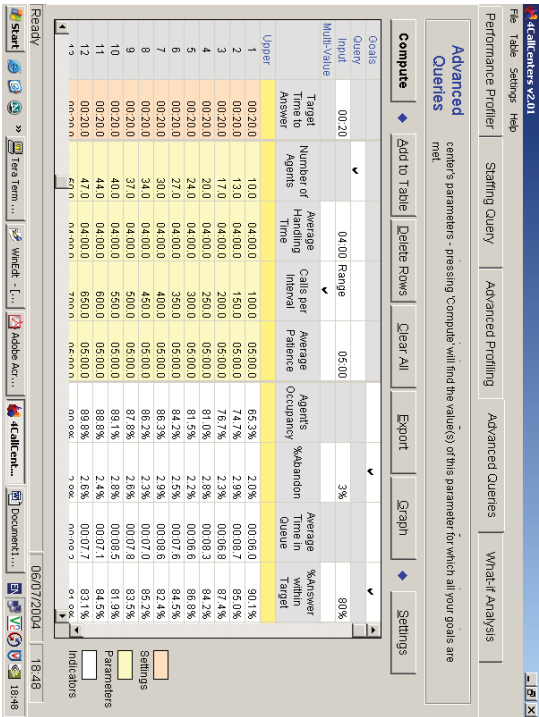
4CallCenters: Advanced Staffing Queries

Set multiple performance goals.

Example: $1/\mu = 4$ minutes, $1/\theta = 5$ minutes;
 λ varies from 100 to 1200, in steps of 50.

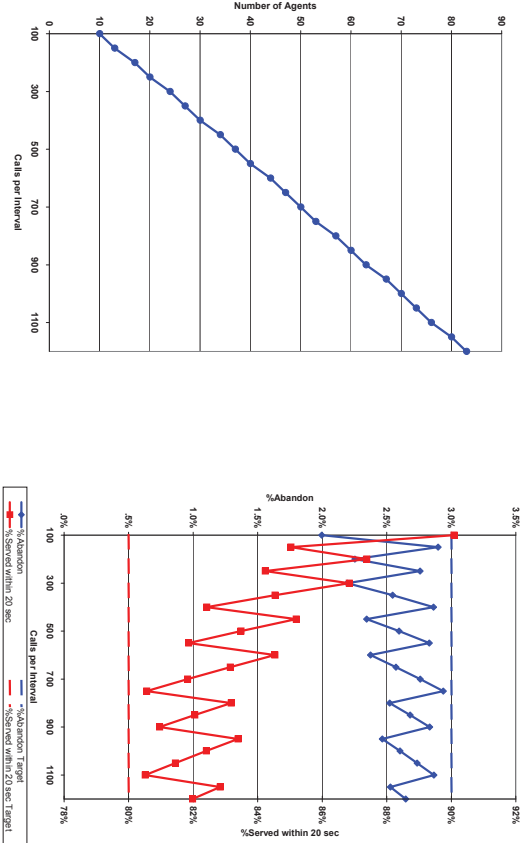
Performance targets:
 $P\{Ab\} \leq 3\%;$ $P\{W_q < 20 \text{ sec}; Sr\} \geq 0.8.$

4CallCenters output



Advanced Staffing Queries II

Recommended staffing level Target performance measures



EOS: 10 agents needed for 100 calls per hour but only 83 for 1200 calls per hour.

Call Centers: Hierarchical Operational View

Forecasting Customers: Statistics, Time-Series

Agents : HRM (Hire, Train; Incentives, Careers)

Staffing: Queueing Theory

Service Level, Costs

FTE's (Seats)
per unit of time

Shifts: IP, Combinatorial Optimization; LP

Union constraints, Costs

Shift structure

Rostering: Heuristics, AI (Complex)

Individual constraints

Agents Assignments

Skills-based Routing: Stochastic Control

Operational Regimes in Many-Server Queues

The **Quality-Efficiency Tradeoff** in services (call centers).

Offered Load: $R = \lambda \times E[S]$ Erlangs, namely minutes of work (= service) that arrive per minute.

Efficiency-Driven (ED):

$$n \approx R - \gamma R, \quad \gamma > 0.$$

Understaffing with respect to the offered load.

Quality-Driven (QD):

$$n \approx R + \delta R, \quad \delta > 0.$$

Overstaffing with respect to the offered load.

Quality and Efficiency-Driven (QED):

$$n \approx R + \beta \sqrt{R}, \quad -\infty < \beta < \infty.$$

The **Square-Root Staffing Rule**:

- Introduced by **Erlang**, already in 1924!
- Rigorized by **Halfin-Whitt**, only in 1981 (Erlang-C);
- Above version: with Garnett, Reiman, Zeltyn (Erlang-A/G).

Operational Regimes: Rules-of-Thumb

Assume that **offered load** R is not small ($\lambda \rightarrow \infty$).

ED regime:

$$n \approx R - \gamma R, \quad 0.1 \leq \gamma \leq 0.25.$$

- Essentially **all** customers delayed prior to service;
- %Abandoned $\approx \gamma$ (10-25%);
- Average wait ≈ 30 seconds - 2 minutes.

QD regime:

$$n \approx R + \delta R, \quad 0.1 \leq \delta \leq 0.25.$$

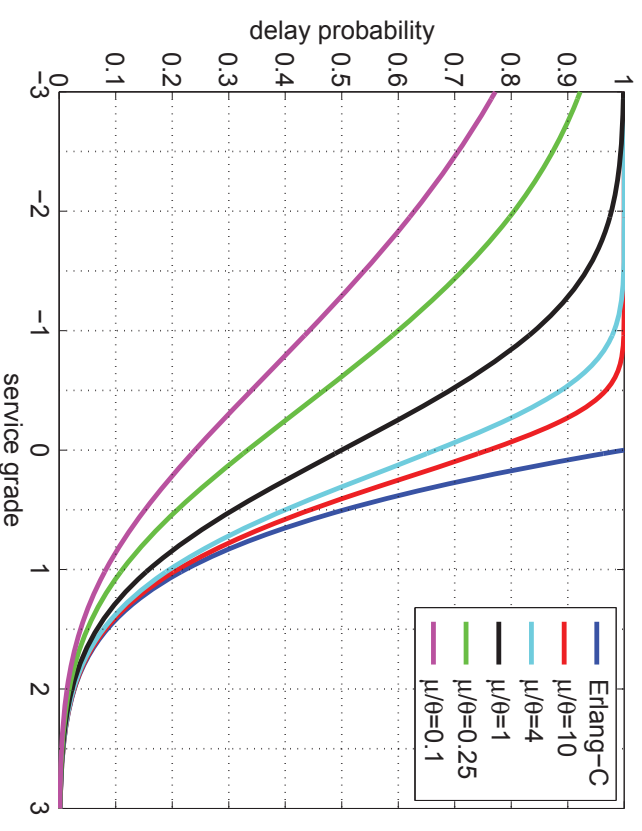
Essentially **no** delays.

QED regime:

$$n \approx R + \beta \sqrt{R}, \quad -1 \leq \beta \leq 1.$$

- %Delayed between **25% and 75%**;
- %Abandoned is 1-5%;
- Average wait is one-order less than average service-time (seconds vs. minutes).

The QED Regime in Erlang-A: Delay Probability



Note. Erlang-C is the limit of Erlang-A, as patience increases indefinitely.

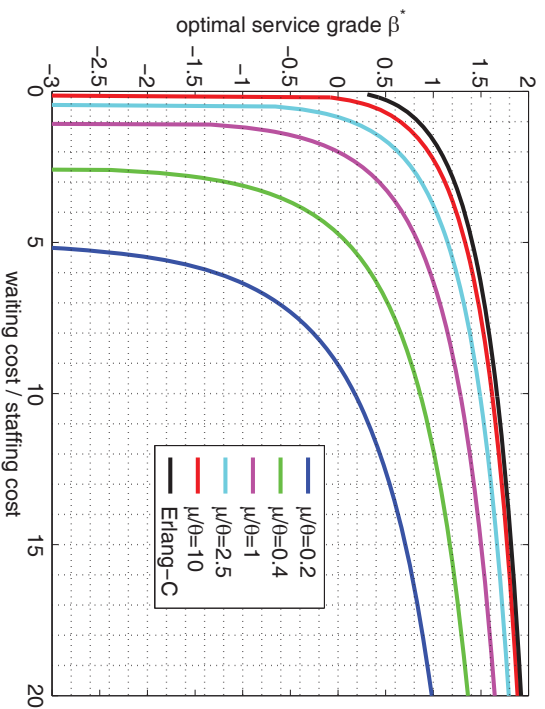
Dimensioning Erlang-A: Optimal QoS

Cost $= c \cdot n + d \cdot \lambda E[W_q]$.

(Abandonment cost can be accommodated via $P\{Ab\} = \theta E[W_q]$.)

Optimal staffing level:

$$n^* \approx R + \beta^*(r; s) \sqrt{R}, \quad r = d/c, \quad s = \sqrt{\mu/\theta},$$



- $r < \theta/\mu$ implies that “close-the-gate” is optimal.
- $r \leq 20 \Rightarrow \beta^* < 2$; $r \leq 500 \Rightarrow \beta^* < 3$!
- **Remarkable** accuracy and robustness, via numerical tests.

Non-Parametric Queueing Models: A Basic Service Station

Assumptions:

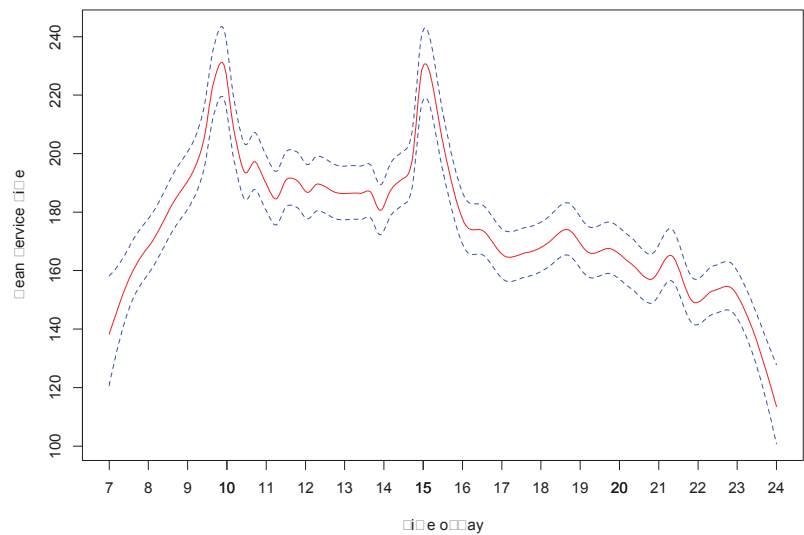
- Non-Poisson (Renewal) Arrivals;
- Non-Exponential i.i.d. Service Times;
- Non-Exponential i.i.d. (Im)Patience.

Analysis:

- Intractable Models, hence resort to Approximations;
- Single- and Moderately-Few Servers in Heavy-Traffic;
(Many-Server Models with General Service Times is still a Theory in the Making);
- Steady-State Analysis;
- Two-Moment Theory: Means and Coefficients-of-Variations;
- Priorities;
- Optimal Scheduling of Customer Classes: The $c\mu$ -Rule, and Relatives.

Interdependence of the Building Blocks

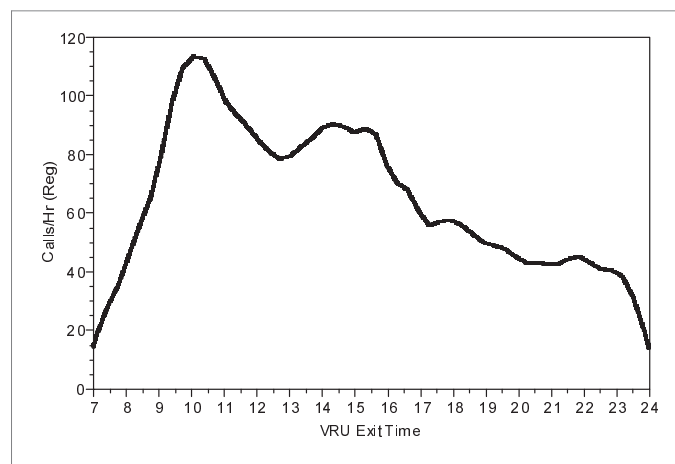
Figure 12: Mean Service Time (Regular) vs. Time-of-day (95% CI) ($n = 42613$)



Arrival Rates: Longest Services at Peak Loads

Arrivals: Inhomogeneous Poisson

Figure 1: Arrivals (to queue or service) – “Regular” Calls



Service Times: Short and Long

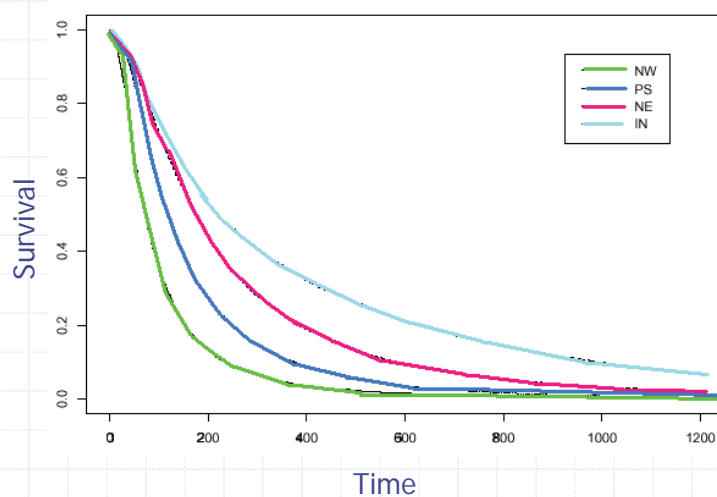
Service Time

	Overall	Regular service	New customers	Internet	Stock
Mean	188	181	111	381	269
SD	240	207	154	485	320
Med	114	117	64	196	169

Service Times: Stochastically Ordered

Service Time

Survival curve, by Types



Means (In Seconds)

NW (New) = 111

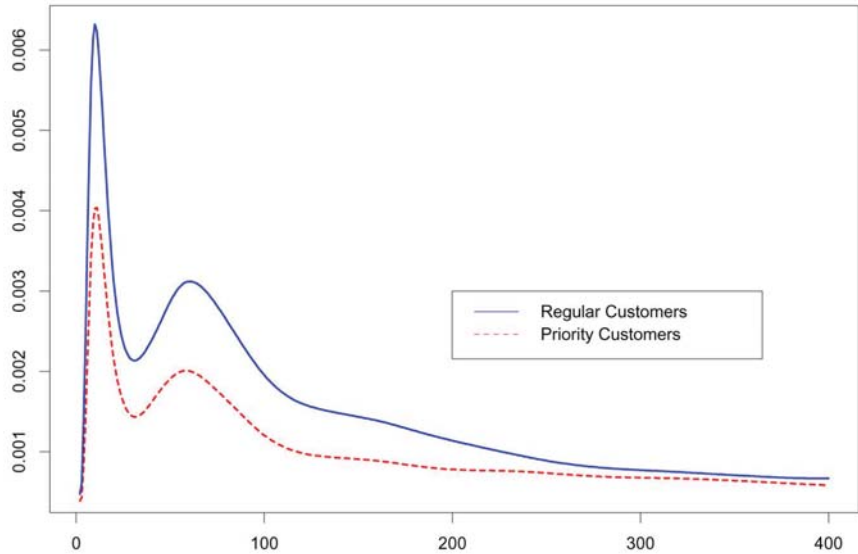
PS (Regular) = 181

NE (Stocks) = 269

IN (Internet) = 381

(Im)Patience: Regulars vs. VIP

Hazard Rate: Empirical (Im)Patience



BONUS SUPPLEMENT: E-TAILING'S FUTURE 

www.businessweek.com

BusinessWeek

OCTOBER 23, 2000 A PUBLICATION OF THE MCGRAW-HILL COMPANIES

Mutual Funds
How to avoid a big tax bill

Wall Street
Will tech's slide keep spreading?

Dot-coms
The search for new business models

Managed Care
Employers seek a new solution

WHY SERVICE STINKS
Companies know just how good a customer you are—and unless you're a high roller, they would rather lose you than fix your problem

Page 118

#BXNDZLN*****CR-RT SORF**0833
#00323865631763#3010201 019469
52/INDUSTRIAL 0830
ENGINEERING LIBRARY 103
PO BOX 830657
BIRMINGHAM AL 35283-0657
AOL Keyword: BW

Customer Relationships Management

NationsBank's Design of the Service Encounter

Examples of Specifications:
Assignable Grade Of Service

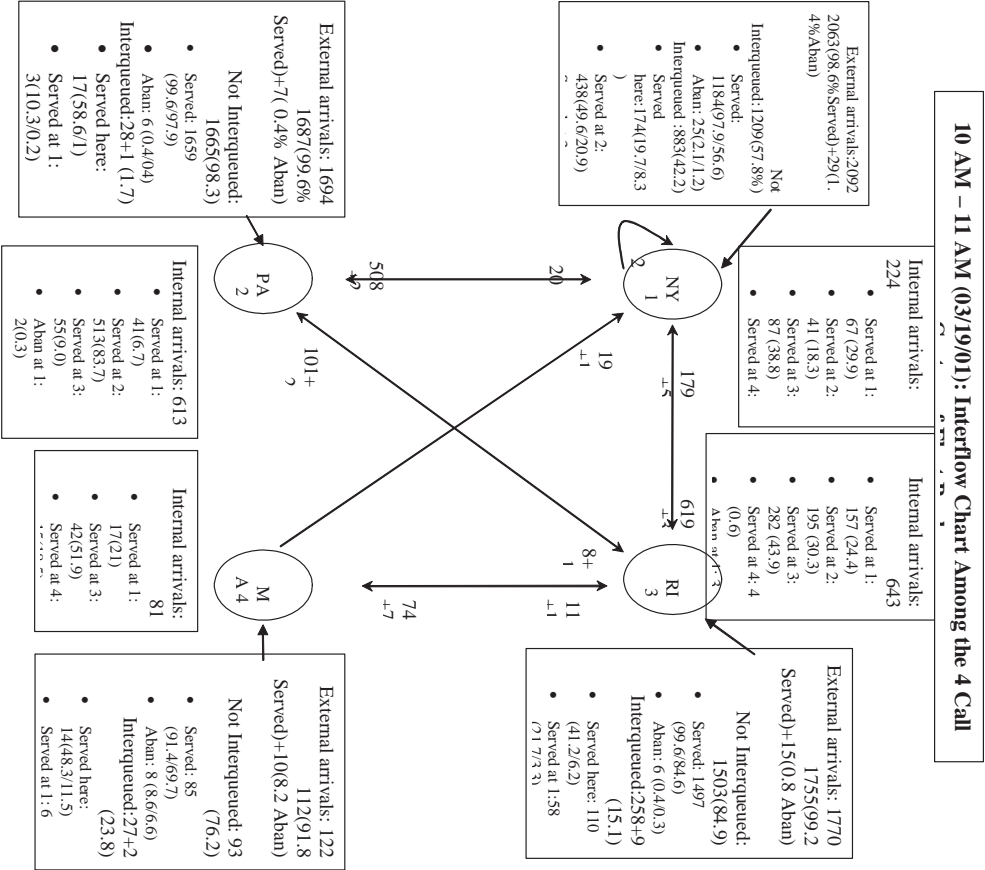
	RG1	RG2	RG3
VRU Target	70% of calls	85% of calls	90% of calls
Abandonment rate	< 1 %	< 5%	< 9%
Speed of Answer	100% in 2 rings	80% in 20 seconds	50% in 20 seconds
Average Talk Time	no limit	4 min. average	2 min. average
Rep. Training	universal	product experts	basic product
Rep. Personalization	request rep / callback	FCFS	FCFS
Trans. Confirmation	call / fax	call / mail	mail
Problem Resolution	during call	within 2 business days	within 8 business days

NationsBank CRM: Relationship Groups:

- RG1: high-value customers;
- RG2: marginally profitable customers (with potential);
- RG3: unprofitable customer.

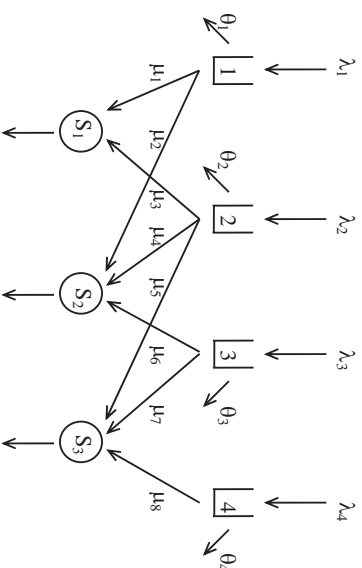
CRM = Customer **Revenue** Management

Distributed Call Center (U.S. Bank)



Skills-Base Routing: Operational Complexities

Multi-queue parallel-server system = schematic depiction of a **telephone call-center**:



Here the λ 's designate arrival rates, the μ 's service rates, the θ 's abandonment rates, and the S 's are the number of servers in each server-pool.

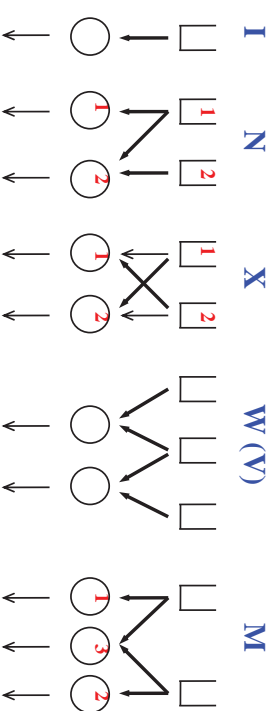
Skills-Based Design:

- **Queue**: "customer-type" requiring a specific type of service;
- **Server-Pool**: "skills" defining the service-types it can perform;
- **Arrow**: leading into a server-pool define its skills / constituency.

For example, a server with skill 2 (**S2**) can serve customers of type 3 (**C3**) at rate μ_6 customers/hour.

Customers of type 3 arrive randomly at rate λ_3 customers/hour, equipped with an impatience rate of θ_3 .

Some Canonical Designs - Animation



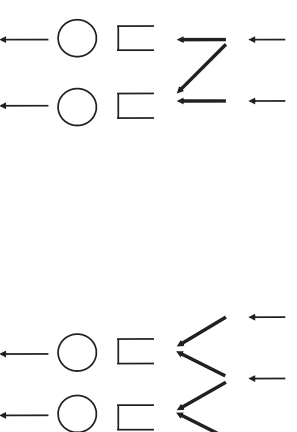
I – dedicated (specialized) agents

N: for example,

- C1 = VIP, then S2 are serving C1 to improve service level.
- C2 = VIP, then S2 serve C1 to improve efficiency.
- S2 = Bilingual.

X: for example, S1 has C1 as Primary and C2 as Secondary Types.

V: Pure Scheduling; **Upside-down V**: Pure Routing.



Major Design / Engineering Decisions

1. Classifying customers into **types** (**Marketing**):
Tech. support vs. Billing, VIP vs. Members vs. New
2. Determining server **skills, incentives, numbers** (**HRM, OM, OR**)
Universal vs. Specialist, Experienced / Novice, Uni- / Multi-lingual;
Staffing: how many servers?
3. Prerequisite Infrastructure - MIS / IT / Data-Bases (**CS, Statistics**)
CTI, ERP, Data-Mining

Major Control Decisions

4. Matching customers and agents (**OR**)
 - **Customer Routing**: Whenever an agent turns idle and there are queued customers, which customer (if any) should be routed to this agent.
 - **Agent Scheduling**: Whenever a customer arrives and there are idle agents, which agent (if any) should serve this customer.
5. **Load Balancing**
 - Routing of customers to distributed call centers (eg. nation-wide)

SBR: Where are We?

Still a **challenge**, both theoretically and practically.

- “Exact” analysis of Markovian models (but mostly “queue-less”), by Koole et al.
- The ED-regime is relatively-well covered, in conventional heavy-traffic a-la Stolyar’s (control) and the fluid-models of Harrison et al (staffing + control, accommodating also non-parametric models with “time-varying randomness”).
- Control in the QED-regime is “theoretically-covered” by Atar et al. (exponential service-times).
- Staffing + Control in the QED-regime covers special cases: Gurvich, Armony; Dai, Tezcan; Gurvich, Whitt; ...

Still plenty to do.

Interesting and Significant Additional Topics

- Stochastic Service **Networks**:
 - Classical Markovian: Jackson and Gordon-Newell, Kelly/BCMP Networks;
 - Non-Parametric Network Approximations (QNA, SBR).
- Service **Quality** (Psychology, Marketing);
- Additional Significant Service Sectors: **Healthcare**, Hospital-ity, Retail, Professional Services (Consulting), \dots ; e-health, e-retail, e- \dots ;
- Convergence of Services and Manufacturing: After-Sale or Field Support (life-time customer-value);
- Service **Supply-Chains**;
- **New-Service Development** (or Service-Engineering in Germany);
- Design and Management of the **Customer-System Interface**: Multi-Media Channels; Appointments; Pricing; \dots
- Revenue Management (Finite Horizon, Call Centers, \dots)

Call Centers = Q 's w/ Impatient Customers 15 Years History, or “A Modelling Gallery”

1. Kella, Meilijson: Practice \Rightarrow Abandonment important
2. Shimkin, Zohar: No data \Rightarrow Rational patience in Equilibrium
3. Carnon, Zakay: Cost of waiting \Rightarrow Psychological models
4. Garnett, Reiman; Zeltyn: Palm/Erlang-A to replace Erlang-C/B as the standard **Steady-state** model
5. Massey, Reiman, Rider, Stolyar: Predictable variability \Rightarrow **Fluid** models, **Diffusion** refinements
6. Ritov; Sakov, Zeltyn: Finally Data \Rightarrow **Empirical** models
7. Brown, Gans, Haipeng, Zhao: **Statistics** \Rightarrow Queueing Science
8. Atar, Reiman, Shaikhet: Skills-based routing \Rightarrow **Control** models
9. Nakibly, Meilijson, Pollatchek: Prediction of waiting \Rightarrow **Online Models** and Real-Time Simulation
10. Garnett: Practice \Rightarrow **4CallCenters.com**
11. Zeltyn: Queueing Science \Rightarrow **Empirically-Based Theory**
12. Borst, Reiman; Zeltyn: Dimensioning $M/M/N+G$
13. Momcilovic: **Non-Parametric** (G/GI/N+GI) QED Q 's
14. Jennings; Feldman, Massey, Whitt: Time-stable performance (ISA)