

Empirical Adventures in Call Centers Emergency Departments, ...

Avishai Mandelbaum Technion IE&M

with Graduate Students, Research Partners
Technion SEE Center, IBM+Rambam+Technion OCR

WITOR, September 2009

The SEE Center - Project DataMOCCA





DataMOCCA

Data **MO**dels for **C**all **C**enter **A**nalysis

Project Collaborators:

Technion: Paul Feigin, Avi Mandelbaum

Technion SEElab: Valery Trofimov, Ella Nadjharov, Igor Gavako, Katya Kutsy, Polyna Khudyakov, Shimrit Maman, Pablo Liberman

Students (PhD, MSc, BSc), RAs

Wharton: Larry Brown, Noah Gans, Haipeng Shen (N. Carolina),

Students, Wharton Financial Institutions Center

Companies: U.S. Bank, Israeli Telecom, 2 Israeli Banks,

Israeli Hospitals, ...

The SEE Center - Project DataMOCCA

Goal: Designing and Implementing a (universal) data-base/data-repository and interface for storing, retrieving, analyzing, displaying and interacting with transaction-based data.



The SEE Center - Project DataMOCCA

Goal: Designing and Implementing a (universal) data-base/data-repository and interface for storing, retrieving, analyzing, displaying and interacting with transaction-based data.



Enable the Study of:

- Customers (Callers, Patients)

- Servers (Agents, Nurses)

- Managers (System)

Waiting, Abandonment, Returns

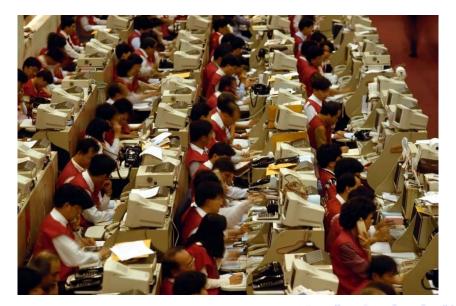
Service Duration, Activity Profile

Loads, Queue Lengths, Trends

Call-Center: Hidden Complex Service Network



Call-Centers: "Sweat-Shops of the 21st Century"



A "Good" Hospital in Beijing



DataMOCCA: System Components

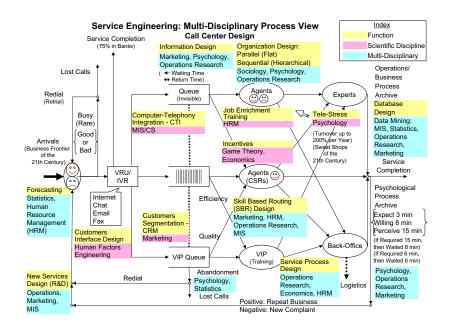
- Clean Databases: Operational histories of individual customers and servers (mostly with IDs).
 - In Call Centers: from IVR to Exit;
 - In Hospitals: from ED to Exit (or just ED).
- 2. SEEStat: Online GUI (friendly, flexible, powerful)
 - Queueing-Science perspective;
 - Operational data (vs. financial, contents or clinical);
 - Flexible customization (e.g. seconds to months);

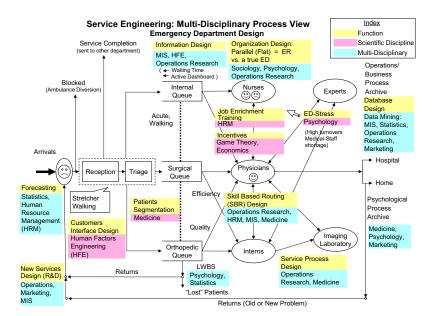
3. Tools:

- Online statistics (survival analysis, mixtures, smoothing);
- Dynamic Graphs (flow-charts, work-flows)
- Simulators (CC, ED; data-driven).

Current Databases

- **1.** U.S. <u>Bank</u> (**PUBLIC**): 220M calls, 40M agent-calls, 1000 agents, 2.5 years, 7-40GB.
- 2. Israeli Banks:
 - Small (PUBLIC): 350K calls, 15 agents, 1 year. Started it all in 1999 (JASA), now "romancing" again (Medium, with 300 agents);
 - Large (ongoing): 500 agents, 1.5 years, 3-8GB.
- 3. Israeli <u>Telecom</u> (ongoing): 800 agents, 3.5 years; 5-55GB.
- 4. Israeli Hospitals:
 - Six ED's (to be made PUBLIC);
 - Large (ongoing): 1000 beds, 45 medical units, 75,000 patients hospitalized yearly, 4 years, 7GB.
- 5. Website (pilot).



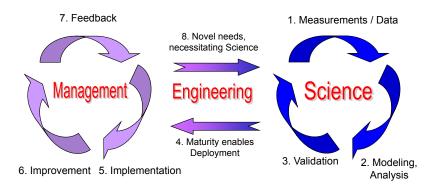


Expanding the Scientific Paradigm (OCR)

- Physics, Biology, ...: Measure, Model, Experiment, Validate, Refine.
- Human-complexity triggered the above in Transportation, Economics.

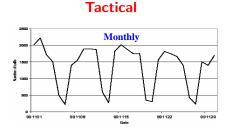
Expanding the Scientific Paradigm (OCR)

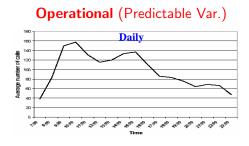
- Physics, Biology, ...: Measure, Model, Experiment, Validate, Refine.
- Human-complexity triggered the above in Transportation, Economics.
- Expand to:



Arrivals to a Call Center (Israel, 1999): Time Scales



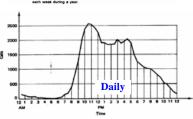






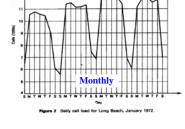
Arrivals to a Call Center (U.S., 1976): Queueing Science





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

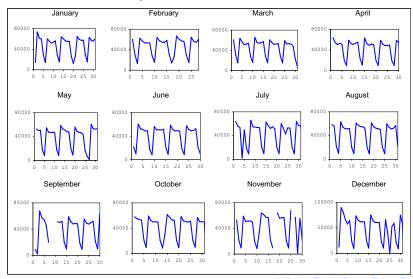
(E. S. Buffa, M. J. Cosgrove, and B. J. Luce, "An Integrated Work Shift Scheduling System")





Monthly Arrivals to Service

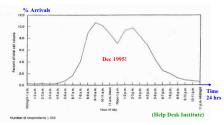
U.S. Bank: Daily Arrival-Rates, over a Month, 2002



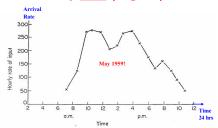
Daily Arrivals to Service: Time-Inhomogeneous (Poisson?)

Intraday Arrival-Rates (per hour) to Call Centers





May 1959 (England)



November 1999 (Israel)

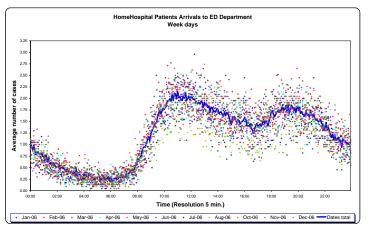


Observation:

Peak Loads at 10:00 & 15:00

Arrivals to an Emergency Department (ED)

Large Israeli ED, 2006

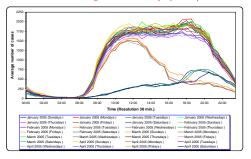


- Second peak at 19:00 (vs. 15:00 in call centers).
- How much stochastic variability?

Intraday Arrival Rates: Does a Day have a Shape?

Arrival Patterns, Israeli Telecom

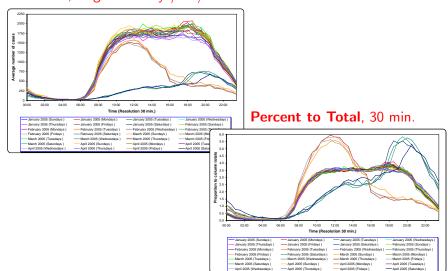
Arrivals, Avg. Weekdays/1-4/2005



Intraday Arrival Rates: Does a Day have a Shape?

Arrival Patterns, Israeli Telecom

Arrivals, Avg. Weekdays/1-4/2005



A (Common) Model for Call Arrivals

Whitt (99'), Brown et. al. (05'), Gans et. al. (09'), and others:

Doubly-stochastic (Cox, Mixed) Poisson with instantaneous rate

$$\Lambda(t) = \lambda(t) \cdot X ,$$

where $\int_0^T \lambda(t) dt = 1$.

• $\lambda(t)$ = "Shape" of weekday

[Predictable variability]

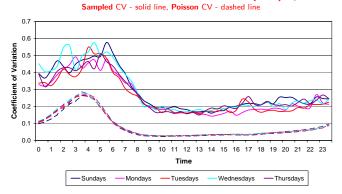
X = Total # arrivals

[Unpredictable variability]

w/ Maman & Zeltyn (09'): Above assumes "too-much" stochastic variability!

Over-Dispersion (Relative to Poisson), Maman et al. ('09)

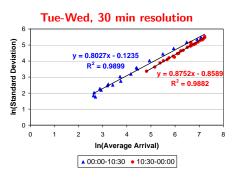
Israeli-Bank Call-Center Arrival Counts - Coefficient of Variation (CV), per 30 min.

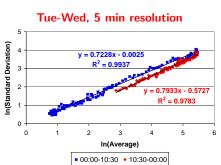


- 263 regular days, 4/2007 3/2008.
- Poisson CV = $1/\sqrt{\text{mean arrival-rate}}$.
- Sampled CV's ≫ Poisson CV's ⇒ Over-Dispersion.

Over-Dispersion: Fitting a Regression Model







Significant linear relations (Aldor & Feigin):

$$ln(STD) = c \cdot ln(AVG) + a$$

Over-Dispersion: Random Arrival-Rate Model

The **linear relation** between ln(STD) and ln(AVG) motivates the following model:

Arrivals distributed Poisson with a Random Rate

$$\Lambda = \lambda + \lambda^{c} \cdot X, \quad 0 < c < 1;$$

- X is a random-variable with E[X] = 0, capturing the magnitude of **stochastic deviation** from mean arrival-rate.
- *c* determines **scale-order** of the over-dispersion:
 - c=1, proportional to λ ;
 - c=0, Poisson-level, same as $0 \le c \le 1/2$.

In call centers, over-dispersion (per 30 min.) is of order λ^{c} , $c \approx 0.8 - 0.85$.

Over-Dispersion: Distribution of X?

- Fitting a **Gamma Poisson** mixture model to the data: Assume a (conjugate) prior Gamma distribution for the arrival rate $\Lambda \stackrel{d}{=} Gamma(a, b)$. Then, $Y \stackrel{d}{=} Poiss(\Lambda)$ is Negative Binomial.
- Very good fit of the Gamma Poisson mixture model, to data of the Israeli Call Center, for the majority of time intervals.
- Relation between our c-based model and Gamma-Poisson mixture is established.
- Distribution of X derived, under the Gamma prior assumption: X is asymptotically normal, as $\lambda \to \infty$.

Over-Dispersion: The QED-c Regime

QED-c Staffing: Under offered-load $R = \lambda \cdot E[S]$,

$$n = R + \beta \cdot R^c$$
, $0.5 < c < 1$

Performance measures:

a. Delay probability:
$$P\{W_q > 0\} \sim 1 - F(\beta)$$

b. Abandonment probability:
$$P\{Ab\} \sim \frac{E[X-\beta]_+}{n^{1-c}}$$

c. Average offered wait:
$$E[V] \sim \frac{E[X - \beta]_+}{n^{1-c} \cdot g_0}$$

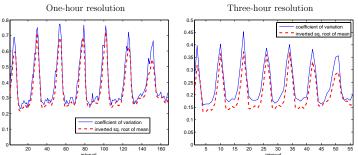
d. Average actual wait:
$$E_{\Lambda,n}[W] \sim E_{\Lambda,n}[V]$$



Over-Dispersion: The Case of ED's

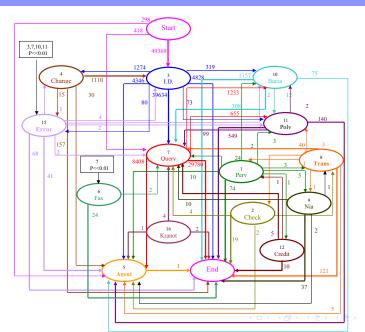
Israeli-Hospital Emergency-Department

Arrival Counts - Coefficient of Variation, per 1-hr. & 3-hr.



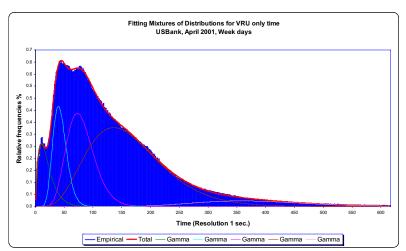
- 194 **weeks**, 1/2004 10/2007 (excluding 5 weeks war in 2006).
- Moderate over-dispersion: c = 0.5 reasonable for hourly resolution.
- ED beds in conventional QED (Less var. than call centers!?).

Call Transitions in the IVR - Phase Type



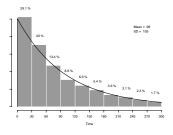
Service Times: Fitting Distribution

Fitting Mixture of 5 Gamma Components

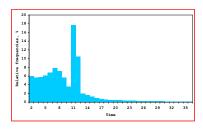


Beyond Averages: Waiting Times in a Call Center

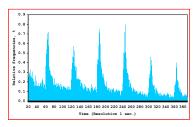
Small Israeli Bank



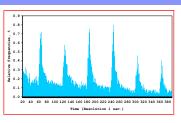
Large U.S. Bank



Medium Israeli Bank



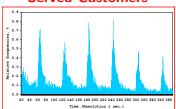
"Waiting-Times" Puzzle at a Large Israeli Bank



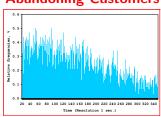
Peaks Every 60 Seconds. Why?

- Human: Voice-announcement every 60 seconds.
- System: Priority-upgrade (unrevealed) every 60 secs (Theory?)

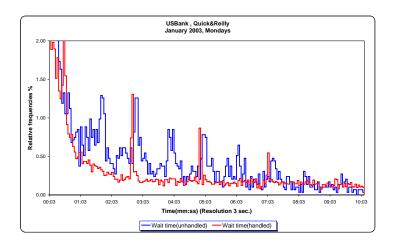
Served Customers



Abandoning Customers



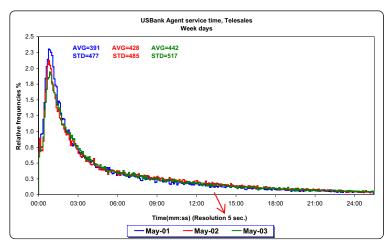
Still a Puzzle at a US Bank



- Different cycles of peaks in the waiting times of both served (protocol?) and abandoning (psychology?) customers.
- No theory for periodic updates of either priorities or information.

Service Times: Service Science

US Bank: Service Time Histograms for Telesales, 2001-3

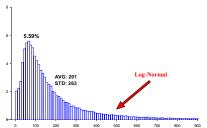


Service Times: Distribution and Psychology

Histogram of Service Times in a Small Israeli Bank



November-December

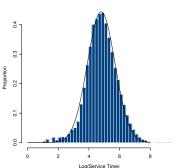


- Lognormal service times prevalent in call centers
- 6.8% Short-Services: Agents' "Abandon" (improve bonus, rest)
- Distributions, not only Averages, must be measured.

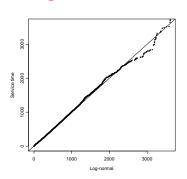
Validating LogNormality of Service Times

Israeli Call Center, Nov-Dec, 1999



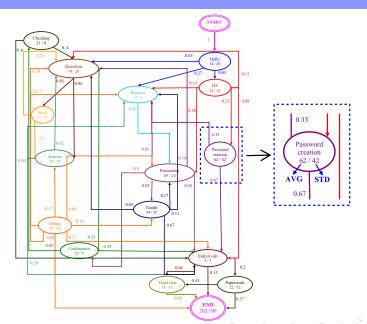


LogNormal QQPlot



Call Transitions in the Service

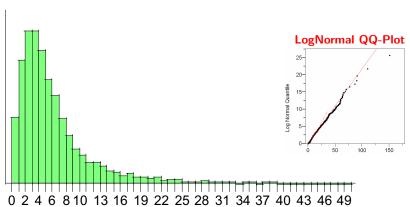
Israeli Bank, Retail Service



Length of Stay: Resolution Dependence

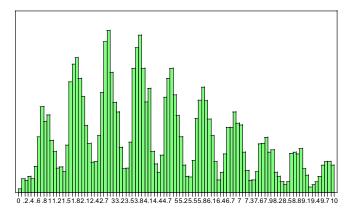
Israeli Large Hospital: LOS in IW

Days Resolution



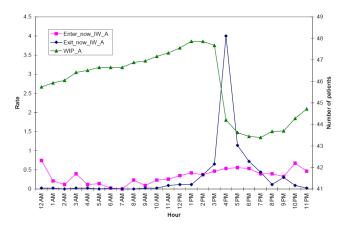
Length of Stay: Resolution Dependence

Hours Resolution: LN = Normal Mixture



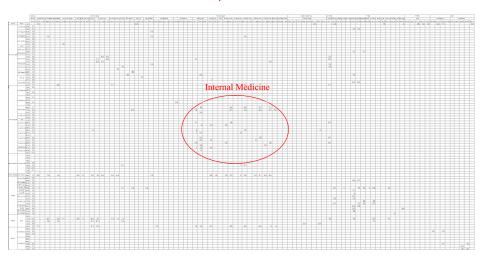
Length of Stay: Resolution Dependence

Internal Ward A: Arrivals / Departures / # Patients , by hour



Ongoing: Empirical Analysis of an ED, IW and Everything In Between, w/ Y. Marmor, Y. Tseytlin, G. Yom-Tov, M. Armony.

90×90 Matrix, Sub-Ward Resolution



8 × 8 Matrix, Division Resolution

Including Arrivals and Releases

| | Home | Surgery | Internal | Psychology | Intensive Care | Pediatrics | Emergency Dep. | Gynecology |
|----------------|------|---------|----------|------------|----------------|------------|----------------|------------|
| Home | | 8.4 | 3.2 | 0.1 | | 18.3 | 60.3 | 9.7 |
| Surgery | 90 | 7.9 | 1.3 | | 0.7 | 0.1 | | |
| Internal | 84.4 | 1.9 | 13 | 0.1 | 0.5 | | | 0.1 |
| Psychology | 94.3 | 1.9 | 3.8 | | | | | |
| Intensive Care | 17.2 | 40.9 | 38.4 | | | 0.9 | | 2.6 |
| Pediatrics | 78.8 | 0.6 | | | | 20.6 | | |
| Emergency Dep. | 69.9 | 8.9 | 19.2 | 0.2 | 0.3 | 1 | | 0.5 |
| Gynecology | 55.3 | 0.3 | 0.2 | | 0.1 | | | 44.1 |

Transitions Inside the Hospital

| | Surgery | Internal | Psychology | Intensive Care | Pediatrics | Emergency Dep. | Gynecology |
|----------------|---------|----------|------------|----------------|------------|----------------|------------|
| Surgery | 78.3 | 12.7 | 0.2 | 7 | 1.4 | | 0.4 |
| Internal | 12 | 83.3 | 0.6 | 3.4 | 0.2 | | 0.5 |
| Psychology | 33.3 | 66.7 | | | | | |
| Intensive Care | 49.5 | 46.4 | | | 1 | | 3.1 |
| Pediatrics | 2.6 | 0.2 | | 0.1 | 96.9 | | 0.2 |
| Emergency Dep. | 29.7 | 63.7 | 0.6 | 0.9 | 3.4 | | 1.7 |
| Gynecology | 0.7 | 0.4 | | 0.1 | | | 98.8 |

- About 50% of transitions between ED and internal wards.
- Most transitions are inside the specific hospitalized unit.

"ED-to-IW" Routing

IW Operational Measures, or Efficiency vs. Fairness Israeli Large Hospital (1/5/06 to 30/10/08, excluding 1-3/07)

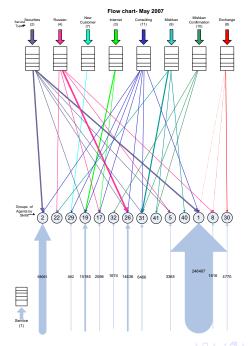
| | Ward A | Ward B | Ward C | Ward D |
|---------------------------|--------|--------|--------|--------|
| ALOS (days) | 6.37 | 4.47 | 5.36 | 5.56 |
| Avg Occupancy Rate | 97% | 95% | 86% | 92% |
| Avg # Patients per Month | 206 | 187 | 210 | 210 |
| Standard capacity | 45 | 30 | 44 | 42 |
| Avg # Patients /Bed/Month | 4.57 | 6.25 | 4.77 | 4.77 |
| Return Rate | 15.4% | 15.6% | 16.2% | 14.8% |

"ED-to-IW" Routing

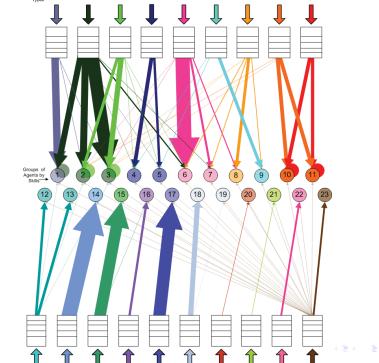
IW Operational Measures, or Efficiency vs. Fairness Israeli Large Hospital (1/5/06 to 30/10/08, excluding 1-3/07)

| | Ward A | Ward B | Ward C | Ward D |
|---------------------------|--------|--------|--------|--------|
| ALOS (days) | 6.37 | 4.47 | 5.36 | 5.56 |
| Avg Occupancy Rate | 97% | 95% | 86% | 92% |
| Avg # Patients per Month | 206 | 187 | 210 | 210 |
| Standard capacity | 45 | 30 | 44 | 42 |
| Avg # Patients /Bed/Month | 4.57 | 6.25 | 4.77 | 4.77 |
| Return Rate | 15.4% | 15.6% | 16.2% | 14.8% |

- The "fastest" + smallest Ward B subject to highest workload:
 occupancy, flux: unfair.
- Calls for ED-to-IW routing, which is both efficient and fair (w/ Tseytlin (MSc), Tseytlin & Momcilovic, Tseytlin & Zviran): exact analysis, QED approximation (natural hours wait for days service), partial bed-information.



Skills Groups- May 2003 Service Types Telesales EBO Premier Business Retail Platinum Subanco cco loads Banking Quality Service (12) (13) (5) (10) Groups of Agents by Skills 46 (10 (5 (20) (30 (32) 28 (16) (33)(35)(40) (42) (43) (1) (44) (33)(42)(36)(20)(46) (1) 94% Costumer 97% Priority 100% Telesales 87% Premier 100% Quick& 100% Retail Loans Service 136 Agents 13% Retail Reilly 233 Agents 6% Retail 3% Case Quality 152 Agents 70 Agents 100 Agents 82 Agents (30)(10)(35)84% Business (28)68% Retail 96% Online (43) (48)11% Retail 91% Business 30% EBO Banking 100% BPS 100% AST 4% Platinum 9% Retail 2% Business 4% Retail 68 Agents 19 Agents 1% Telesales 40 Agents 72 Agents 122 Agents 27 Agents (40)(5) (32)(16)99% Case 99% Retail 74% Business (44)81% Retail Quality 0.75% Business 17% Platinum 100% CCO 18% Subanco 1% Priority 0.25% Premier 9% Retail 136 Agents Service 34 Agents 400 Agents 18 Agents 31 Agents



System Design: Simplification via State-Space Collapse

Service Rate: Class or Pool Dependent?

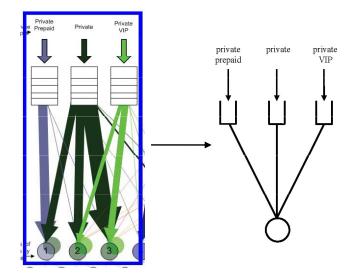
| Agents Group\Service Class | Private Prepaid | Private | Private VIP |
|----------------------------|--------------------|---------|-------------|
| Private Prepaid | 163.1 | 236.1 | |
| Private - Private VIP (1) | | 243.5 | 195.1 |
| Private - Private VIP (2) | | 244 | 201.4 |

⇒ Class-dependent service rate

| Agents Group\Service Class | Business | Business VIP | Business Preservation |
|----------------------------|----------|--------------|--------------------------|
| Business (1) | 276.9 | 261.5 | |
| Business (2) | 336.7 | 334.5 | |
| Business VIP | 315.9 | 280.5 | |
| Business Preservation | | 386.2 | 634.1 |

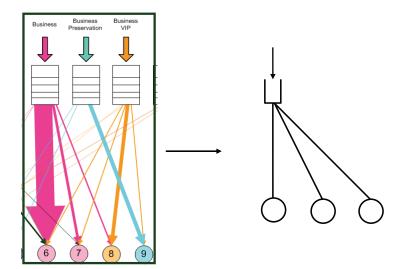
Many-Server Approximations: State-Space Collapse

$\textbf{Class-Dependent} \approx \textit{V-}\textbf{Model}$



Many-Server Approximations: State-Space Collapse

Pool-Dependent $\approx \Lambda$ -Model



Unpredictable Variability: The Multi-Class Case

Unpredictable variability: $X = (X_1, ..., X_l)$

Pairs: $(X_{Retail}, X_{Business})$ and $(X_{Business}, X_{Platinum})$

US Bank: Correlations, 600 weekdays (Gurvich et al., '09)

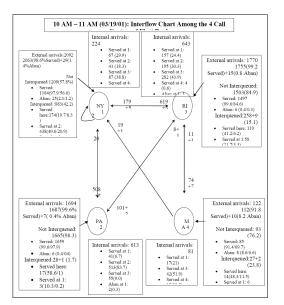




- Positive correlation (not independent)
- Research: Impact on design and control decisions ?

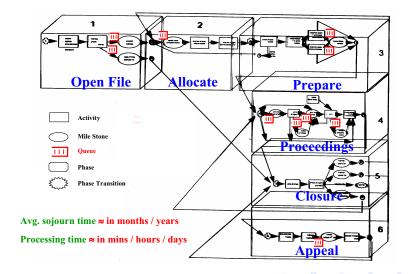
System Design: Inter-queue Model

US Bank



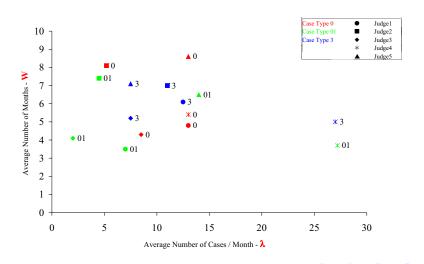
Conceptual Model: The "Production of Justice"

The Labor-Court Process in Haifa, Israel



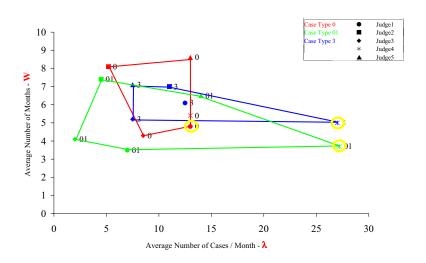
Analytical Model: Little's Law in Court (I)

Judges: Operational Performance - Base Case



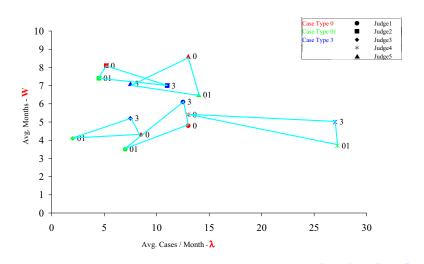
Analytical Model: Little's Law in Court (II)

Judges: Performance by Case-Type



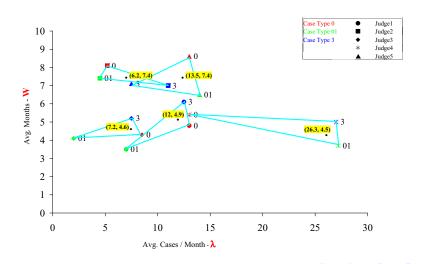
Analytical Model: Little's Law in Court (III)

Judges: Performance Analysis



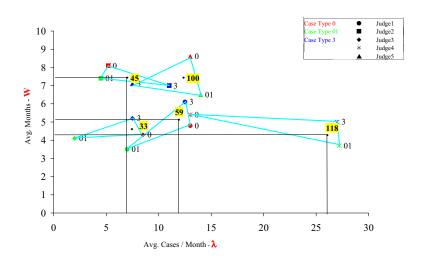
Analytical Model: Little's Law in Court (IV)

Judges: Performance Analysis



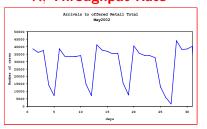
Analytical Model: Little's Law in Court (V)

Judges: Performance Analysis

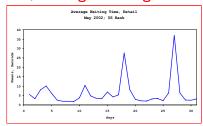


US Bank: Retail calls, May 2002

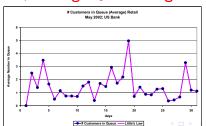
λ , Throughput Rate



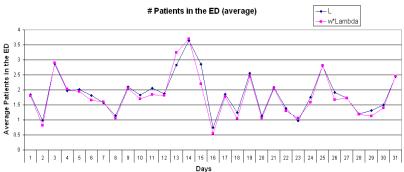
W, Average Waiting Time



L, Average Queue Length



Israeli ED, October 1999, Day Resolution

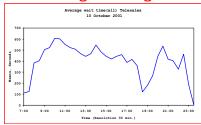


US Bank: Telesales Calls, October 10, 2001

λ , Throughput Rate



W, Average Waiting Time

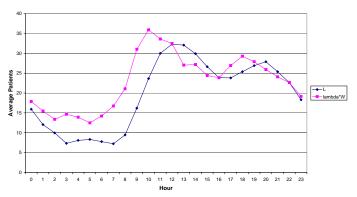


L, Average Queue Length



Israeli ED, Hour Resolution

Patients in the ED (average)



Workload and Offered-Load

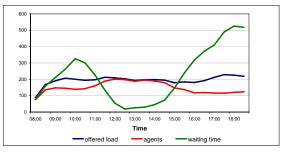
- Workload: Stochastic process, representing the amount of work present at time t, under the assumptions of infinitely many resources (service commences immediately upon arrival).
- Offered-Load: Function of time $t \ge 0$, representing the average of the workload at time t.

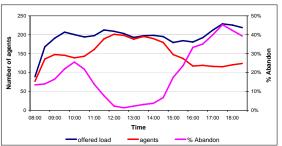
The Offered-Load, R(t), determines staffing level via c-staffing (c=0.5 is conventional square-root staffing):

$$N(t) = R(t) + \beta \cdot [R(t)]^{c}$$

Offered-Load vs. # Agents

Israeli Cable Company, Retail Service, January 2009





Offered-Load Representations (or Time-Varying Little)

For the $M_t/GI/N_t+GI$ queue, the **offered-load** $R=\{R(t),\ t\geq 0\}$, has the following representations:

$$R(t) = E[L(t)] = \int_{-\infty}^{t} \lambda(u) \cdot P(S > t - u) du = E\left[A(t) - A(t - S)\right] =$$

$$= E\left[\int_{t-S}^{t} \lambda(u) du\right] = E[\lambda(t - S_e)] \cdot E[S],$$

where

 $A = \{A(t), t \ge 0\}$ is the Arrival process;

S is a generic service time;

 S_e is a generic excess (residual) service.

In stationary models, where $\lambda(t) \equiv \lambda$, the offered-load R(t) is the familiar $\lambda \cdot E[S]$ (or λ/μ), measured in Erlangs.

Imputing Service Times of Abandoning Customers

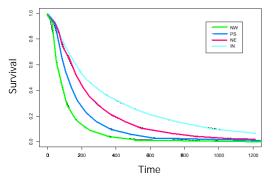
In calculating the offered-load, one must account for service-times of abandoning customers.

A prevalent assumptions is that service times and (im)patience times are independent. Experience suggests that this assumption is often violated.

For example, it is not unreasonable that customers who anticipate longer service times, will be willing to wait more for service before abandoning.

Service Times: Stochastic Order

Small Israeli Bank: Survival Functions by Type



Service times stochastic order: $S_{NW} \stackrel{st}{<} S_{PS} \stackrel{st}{<} S_{NE} \stackrel{st}{<} S_{IN}$

Patience times stochastic order: $au_{\rm NW} \stackrel{\rm st}{<} au_{\rm IN} \stackrel{\rm st}{<} au_{\rm PS} \stackrel{\rm st}{<} au_{\rm NE}$



Relationship Between Service-Time and (Im)Patience

Ongoing research (w/ M. Reich, Y. Ritov) develops a procedure for calculating the function $E(S|\tau=w)$:

1. Introduce $g(w) = E(S|\tau > W = w)$, which is the mean service time of those who waited exactly w units of time and were served. Then calculate g via the non-linear regression:

$$S_i = g(W_i) + \varepsilon_i$$
,

where *i* indexing served customers.

2. Calculate $E(S|\tau=w)$ via the (established) relation

$$E(S|\tau=w)=g(w)-\frac{g'(w)}{h_{\tau}(w)},$$

where $h_{\tau}(w)$ is the hazard-rate function of (im)patience, to be estimated via un-censoring.

Finally, extend the above to calculate the distribution of S, given w, which is then used to impute service-times for calculating the offered-load.

143