# QED Q's

(Service-Science & -Engineering of)

# **Quality- and Efficiency-Driven Queues**

(Call/Contact Centers)

#### Avishai Mandelbaum

Technion, Haifa, Israel

http://ie.technion.ac.il/serveng

SP&Applications, August 2007

Based on joint work with Sergey Zeltyn, ...

Technion SEE Center / Lab: Paul Feigin, Valery Trofimov, RA's, ...

◆ロト ◆卸 ▶ ◆ 重 ▶ ◆ 重 ・ 釣 ♀ (

- ▶ Data-Based Introduction (Service Engineering, Call Centers)
  - ▶ "Production" of Health, Justice, Banking-Services, Tele-Services
  - Simple Models at the Service of Complex Realities

- Data-Based Introduction (Service Engineering, Call Centers)
  - "Production" of Health, Justice, Banking-Services, Tele-Services
  - Simple Models at the Service of Complex Realities
- ► The Basic Call-Center Model: Erlang-A (M/M/N+M)
- Validating Erlang-A? All Assumptions Violated

- Data-Based Introduction (Service Engineering, Call Centers)
  - "Production" of Health, Justice, Banking-Services, Tele-Services
  - Simple Models at the Service of Complex Realities
- ► The Basic Call-Center Model: Erlang-A (M/M/N+M)
- Validating Erlang-A? All Assumptions Violated
- ▶ But Erlang-A Works! Why? Framework - Asymptotic Regimes: QED, ED, ED+QED

- Data-Based Introduction (Service Engineering, Call Centers)
  - "Production" of Health, Justice, Banking-Services, Tele-Services
  - Simple Models at the Service of Complex Realities
- ► The Basic Call-Center Model: Erlang-A (M/M/N+M)
- Validating Erlang-A? All Assumptions Violated
- But Erlang-A Works! Why?
  Framework Asymptotic Regimes: QED, ED, ED+QED
  - General Patience: M/M/N+G
  - General Services: M/G/N+G
  - ► Time-Varying Arrivals: M<sub>t</sub>/G/N<sub>t</sub>+G
  - Heterogenous Customers: M<sup>J</sup>/M/N+M
  - ► Current Research: M<sup>J</sup>/G/N<sub>t</sub>+G

- Data-Based Introduction (Service Engineering, Call Centers)
  - "Production" of Health, Justice, Banking-Services, Tele-Services
  - Simple Models at the Service of Complex Realities
- ► The Basic Call-Center Model: Erlang-A (M/M/N+M)
- Validating Erlang-A? All Assumptions Violated
- But Erlang-A Works! Why?
  Framework Asymptotic Regimes: QED, ED, ED+QED
  - General Patience: M/M/N+G
  - General Services: M/G/N+G
  - ► Time-Varying Arrivals: M<sub>t</sub>/G/N<sub>t</sub>+G
  - ► Heterogenous Customers: M<sup>J</sup>/M/N+M
  - ► Current Research: M<sup>J</sup>/G/N<sub>t</sub>+G
- Explain Practice: "Right Answers for the Wrong Reasons"



- Data-Based Introduction (Service Engineering, Call Centers)
  - "Production" of Health, Justice, Banking-Services, Tele-Services
  - Simple Models at the Service of Complex Realities
- ► The Basic Call-Center Model: Erlang-A (M/M/N+M)
- Validating Erlang-A? All Assumptions Violated
- But Erlang-A Works! Why?
  Framework Asymptotic Regimes: QED, ED, ED+QED
  - General Patience: M/M/N+G
  - General Services: M/G/N+G
  - ► Time-Varying Arrivals: M<sub>t</sub>/G/N<sub>t</sub>+G
  - ► Heterogenous Customers: M<sup>J</sup>/M/N+M
  - ► Current Research: M<sup>J</sup>/G/N<sub>t</sub>+G
- Explain Practice: "Right Answers for the Wrong Reasons"
- ► The Technion SEE Center / Laboratory: DataMOCCA



1. Simple Useful Models at the Service of Complex Realities.

1. Simple Useful Models at the Service of Complex Realities.

Note: Useful must be Simple; Simple often rooted in deep analysis.

1. Simple Useful Models at the Service of Complex Realities.

Note: Useful must be Simple; Simple often rooted in deep analysis.

2. Data-Based Research & Teaching is a Must & Fun.

Supported by **DataMOCCA** = Data **MO**dels for **Call Center Analysis**. Initiated with Wharton, developed at Technion, available for adoption.

1. Simple Useful Models at the Service of Complex Realities.

Note: Useful must be Simple; Simple often rooted in deep analysis.

2. Data-Based Research & Teaching is a Must & Fun.

Supported by **DataMOCCA** = Data **MO**dels for **Call Center Analysis**. Initiated with Wharton, developed at Technion, available for adoption.

3. Back to the **Basic-Research Paradigm** (Physics, Biology, ...): **Measure, Model, Experiment, Validate, Refine, etc.** 

1. Simple Useful Models at the Service of Complex Realities.

Note: Useful must be Simple; Simple often rooted in deep analysis.

2. Data-Based Research & Teaching is a Must & Fun.

Supported by **DataMOCCA** = Data **MO**dels for **Call Center Analysis**. Initiated with Wharton, developed at Technion, available for adoption.

- 3. Back to the **Basic-Research Paradigm** (Physics, Biology, ...): **Measure, Model, Experiment, Validate, Refine, etc.**
- **4. Ancestors** & **Practitioners** often knew/apply the "**right answer**": simply did/do not have our tools/desire/need to prove it so.

Supported by Erlang (1915), Palm (1945),..., thoughtful managers.

1. Simple Useful Models at the Service of Complex Realities.

Note: Useful must be Simple; Simple often rooted in deep analysis.

2. Data-Based Research & Teaching is a Must & Fun.

Supported by **DataMOCCA** = Data **MO**dels for **Call Center Analysis**. Initiated with Wharton, developed at Technion, available for adoption.

- 3. Back to the **Basic-Research Paradigm** (Physics, Biology, ...): **Measure, Model, Experiment, Validate, Refine, etc.**
- **4. Ancestors** & **Practitioners** often knew/apply the "**right answer**": simply did/do not have our tools/desire/need to prove it so. Supported by **Erlang (1915), Palm (1945),...**, thoughtful managers.
- **5. Scientifically-based design principles and tools (software)**, that support the balance of service **quality**, process **efficiency** and business **profitability**, from the (often-conflicting) views of **customers**, **servers**, **managers**: **Service Engineering**.

# **Background Material (Downloadable)**

► Technion's "Service-Engineering" Course (≥ 1995): http://ie.technion.ac.il/serveng

# **Background Material (Downloadable)**

- ► Technion's "Service-Engineering" Course (≥ 1995): http://ie.technion.ac.il/serveng
- 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."
   Technion Report, 2004-2006.
- ► M. "Call Centers: Research Bibliography with Abstracts." Version 7, December 2006.



# Present Focus: Call Centers, but Expanding

#### **U.S. Statistics (Relevant Elsewhere)**

- Over 60% of annual business volume via the telephone
- ► 100,000 200,000 call centers
- ► 3 6 million employees (2% 4% workforce)
- ▶ 1000's agents in a "single" call center = 70 % costs.
- 20% annual growth rate
- \$200 \$300 billion annual expenditures

# **Present Focus: Call Centers, but Expanding**

#### **U.S. Statistics (Relevant Elsewhere)**

- Over 60% of annual business volume via the telephone
- ► 100,000 200,000 call centers
- 3 − 6 million employees (2% − 4% workforce)
- ▶ 1000's agents in a "single" call center = 70 % costs.
- 20% annual growth rate
- \$200 \$300 billion annual expenditures

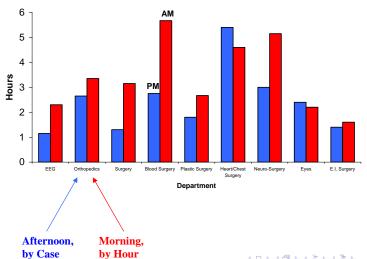
#### Expanding, eg. Healthcare:

- Similar Challenges: Scarce transactional data, natural queueing-network view, human-operations interface (7% LWBS), nurse-staffing (over 3 millions), . . .
- ▶ **Unique** Challenges: Risk, economies vs. dis-economies of scale, synchronization gaps, . . . ,



# The Human Factor, or Even "Doctors" Can Manage

Operations Time - Morning (by Hour) vs. Afternoon (by Case):



# **Prerequisite: Data**

#### **Averages Prevalent.**

But I need data at the level of the **Individual Transaction**: For each service transaction (during a phone-service in a call center, or a patient's stay in a hospital), its **operational history** = time-stamps of events.

# **Prerequisite: Data**

#### **Averages Prevalent.**

But I need data at the level of the **Individual Transaction**: For each service transaction (during a phone-service in a call center, or a patient's stay in a hospital), its **operational history** = time-stamps of events.

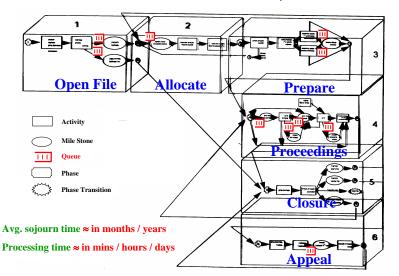
Sources: "Service-floor" (vs. Industry-level, Surveys, ...)

- Administrative (Court, via "paper analysis")
- ► Face-to-Face (Bank, via bar-code readers)
- ► Telephone (Call Centers, via ACD / CTI)
- Future:
  - Hospitals (via RFID)
  - ► IVR (VRU), internet, chat (multi-media)
  - Operational + Financial + Marketing / Clinical history



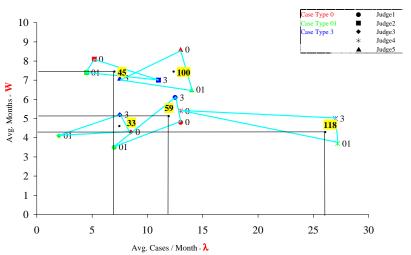
# "Production of Justice" (Administrative) Network

#### The Labor-Court Process in Haifa, Israel



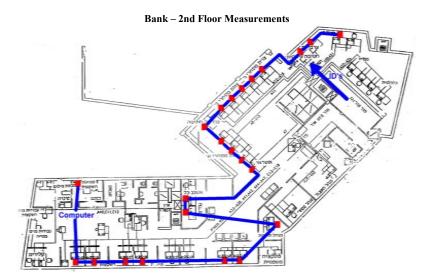
# **Little's Law in Court (Creative Averaging)**

## Judges: The Best/Worst (Operational) Performer



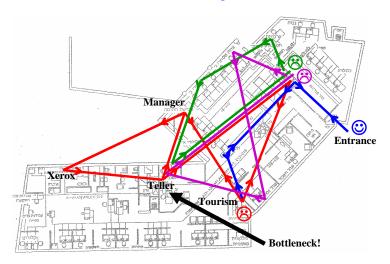
# **Measurements: Face-to-Face Services**

#### 23 Bar-Code Readers at a Bank Branch



#### "Face-to-Face Services" Network

#### **Bank Branch = Queueing Network**



#### **Transition Probabilities of a Jackson Network**

#### Transition Frequencies Between Units in The Private and Business Sections:

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

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

#### Dominant Paths - Business:

Unit Parameter	Station 1 Tourism	Station 2 Teller	Total 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

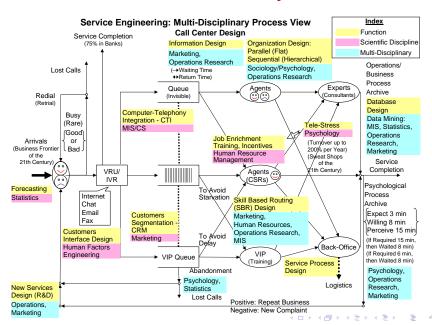
### **Call-Center Environment: Service Network**



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



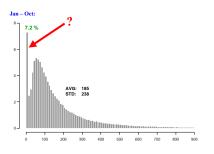
## **Call-Center Network: Gallery of Models**



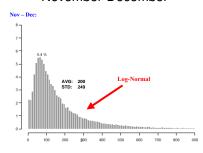
# **Beyond Averages: Service Times in a Call Center**

#### Histogram of Service Times in an Israeli Call Center

#### January-October



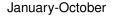
#### November-December



▶ 7.2% Short-Services:

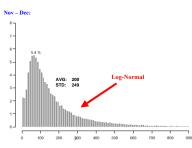
# **Beyond Averages: Service Times in a Call Center**

#### Histogram of Service Times in an Israeli Call Center



# Jan - Oct: 7,2 % AVG: 185 STD: 238

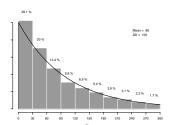
#### November-December



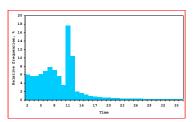
- ▶ **7.2% Short-Services:** Agents' "Abandon" (improve bonus, rest)
- Distributions, not only Averages, must be measured.
- ▶ Lognormal service times prevalent in call centers

# **Beyond Averages: Waiting Times in a Call Center**

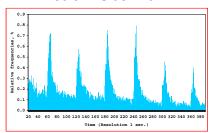
#### **Small Israeli Bank**



Large U.S. Bank



#### Medium Israeli Bank



# The "Phases of Waiting" for Service

#### Common Experience:

- Expected to wait 5 minutes, Required to 10
- ► Felt like 20, Actually waited 10 (hence Willing ≥ 10)

# The "Phases of Waiting" for Service

#### Common Experience:

- Expected to wait 5 minutes, Required to 10
- ► Felt like 20, Actually waited 10 (hence Willing ≥ 10)

#### An attempt at "Modeling the Experience":

```
1. Time that a customer 2. willing to wait ((Im)Patience: \tau)
3. required to wait ((Im)Patience: \tau)
4. actually waits (W_q = min(\tau, V))
5. perceives waiting.
```

# The "Phases of Waiting" for Service

#### Common Experience:

- Expected to wait 5 minutes, Required to 10
- ► Felt like 20, Actually waited 10 (hence Willing ≥ 10)

#### An attempt at "Modeling the Experience":

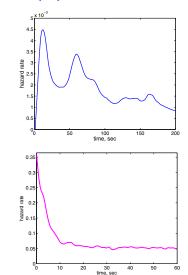
```
1. Time that a customer \begin{array}{cccc} & & & & & & & \\ & & & & & & & \\ 2. & & & & & & \\ 3. & & & & & & \\ 4. & & & & & & \\ 5. & & & & & & \\ \end{array} expects to wait \begin{array}{cccc} & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ \end{array} ((Im)Patience: \tau) (Offered Wait: V) (Offered Wait: V) (W_q = \min(\tau, V)) perceives waiting.
```

```
Experienced customers ⇒ Expected = Required  
"Rational" customers ⇒ Perceived = Actual.
```

Then left with  $(\tau, V)$ .

# **Call Center Data: Hazard Rates (Un-Censored)**

(Im)Patience Time  $\tau$ 



Israel

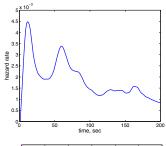
U.S.

# **Call Center Data: Hazard Rates (Un-Censored)**

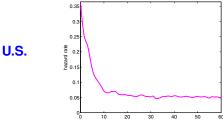
(Im)Patience Time  $\tau$ 

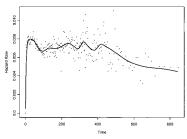
Required/Offered Wait V

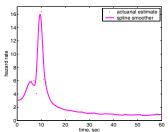
Israel



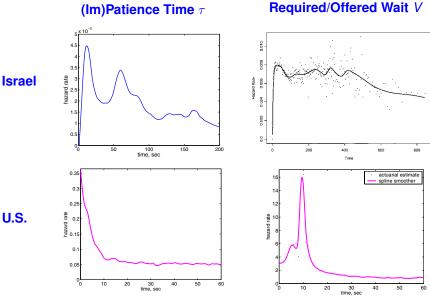
time, sec







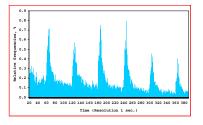
# **Call Center Data: Hazard Rates (Un-Censored)**



Note: 5% abandoning ⇒ 95% (im)patience-observations censored!



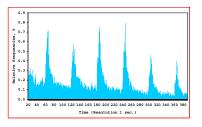
# A "Waiting-Times" Puzzle at a Large Israeli Bank



#### Peaks Every 60 Seconds. Why?

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

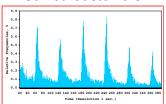
# A "Waiting-Times" Puzzle at a Large Israeli Bank



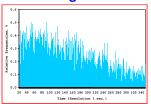
### Peaks Every 60 Seconds. Why?

- ► Human: Voice-announcement every 60 seconds.
- System: Priority-upgrade (unrevealed) every 60 sec's (Theory?)

#### **Served Customers**



#### **Abandoning Customers**



# **Models for Performance Analysis**

- ▶ (Im)Patience: r.v.  $\tau$  = Time a customer is willing to wait
- ▶ Offered-Wait: r.v. V = Time a customer is required to wait (= Waiting time of a customer with infinite patience).
- ▶ Abandonment =  $\{\tau \le V\}$
- **Service** = {*τ* > *V*}
- ▶ Actual Wait  $W_q = \min\{\tau, V\}$ .

21

# **Models for Performance Analysis**

- ▶ (Im)Patience: r.v.  $\tau$  = Time a customer is willing to wait
- ▶ Offered-Wait: r.v. V = Time a customer is required to wait (= Waiting time of a customer with infinite patience).
- ▶ Abandonment =  $\{\tau \le V\}$
- **Service** = {*τ* > *V*}
- ▶ Actual Wait  $W_q = \min\{\tau, V\}$ .

Modeling:  $\tau = \text{input}$  to model, V = output.

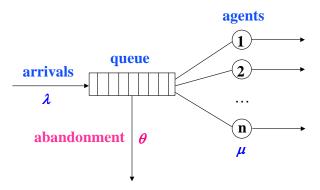
Operational Performance-Measure calculable in terms of  $(\tau, V)$ :

- eg. Avg. Wait =  $E[min\{\tau, V\}]$  (  $E[W_q|Served] = E[V|\tau > V]$  )
- ▶ eg. % Abandon =  $P\{\tau \le V\}$  (  $P\{5 \sec < \tau \le V\}$  )

Application: Staffing - How Many Agents? (When? Who?)



## The Basic Staffing Model: Erlang-A (M/M/N +M)



#### Erlang-A (Palm 1940's): Birth & Death Q, with parameters:

- $\lambda$  **Arrival** rate (Poisson)
- ▶  $\mu$  **Service** rate (Exponential)
- $\bullet$   $\theta$  Impatience rate (Exponential)
- ▶ n Number of Service-Agents.



# **Testing the Erlang-A Primitives**

Arrivals: Poisson?

Service-durations: Exponential?

▶ (Im)Patience: Exponential?

# **Testing the Erlang-A Primitives**

Arrivals: Poisson?

Service-durations: Exponential?

(Im)Patience: Exponential?

Primitives independent?

Customers / Servers Heterogeneous?

Service discipline FCFS?

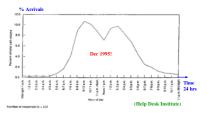
**...?** 

Validation: Support? Refute?

# **Arrivals to Service: only Poisson-Relatives**

#### **Arrival Rate to Three Call Centers**

Dec. 1995 (U.S. 700 Helpdesks)





Time

p.m.

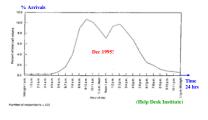
24 hrs

c.m.

## **Arrivals to Service: only Poisson-Relatives**

#### **Arrival Rate to Three Call Centers**

Dec. **1995** (U.S. 700 Helpdesks)



May 1959 (England)



#### November 1999 (Israel)

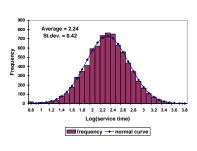


#### Observation:

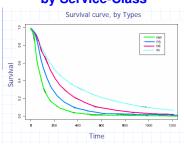
Peak Loads at 10:00 & 15:00

# **Service Durations: LogNormal Prevalent**

### Israeli Bank Log-Histogram



# Survival-Functions by Service-Class



- New Customers: 2 min (NW);
- Regulars: 3 min (PS);

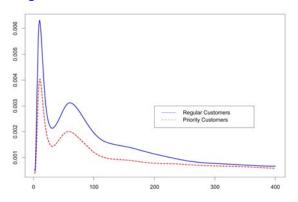
- ► Stock: 4.5 min (NE);
- Tech-Support: 6.5 min (IN).

Observation: VIP require longer service times.



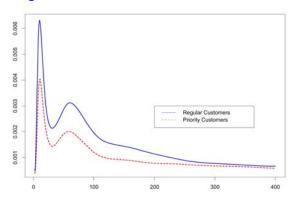
# (Im)Patience while Waiting (Palm 1943-53)

# Irritation ∝ Hazard Rate of (Im)Patience Distribution Regular over VIP Customers – Israeli Bank



# (Im)Patience while Waiting (Palm 1943-53)

# Irritation Hazard Rate of (Im)Patience Distribution Regular over VIP Customers − Israeli Bank



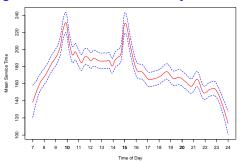
- Peaks of abandonment at times of Announcements
- ► Call-by-Call Data (DataMOCCA) required (& Un-Censoring).

Observation: VIP are more patient (Needy)



# A "Service-Time" Puzzle at a Small Israeli Bank Inter-related Primitives

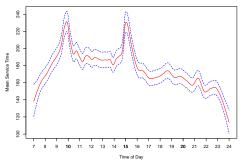
#### Average Service Time over the Day – Israeli Bank



Prevalent: Longest services at peak-loads (10:00, 15:00). Why?

# A "Service-Time" Puzzle at a Small Israeli Bank Inter-related Primitives

Average Service Time over the Day – Israeli Bank

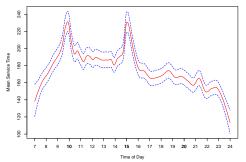


Prevalent: Longest services at peak-loads (10:00, 15:00). Why? Explanations:

► Common: Service protocol different (longer) during peak times.

# A "Service-Time" Puzzle at a Small Israeli Bank Inter-related Primitives

#### Average Service Time over the Day – Israeli Bank



Prevalent: Longest services at peak-loads (10:00, 15:00). Why? Explanations:

- Common: Service protocol different (longer) during peak times.
- Operational: The needy abandon less during peak times; hence the VIP remain on line, with their long service times.

# **Erlang-A: Practical Relevance?**

#### **Experience:**

- ► Arrival process **not pure Poisson** (time-varying,  $\sigma^2$  too large)
- Service times not Exponential (typically close to LogNormal)
- ▶ Patience times **not Exponential** (various patterns observed).
- Building Blocks need not be independent (eg. long wait possibly implies long service)
- Customers and Servers not homogeneous (classes, skills)
- Customers return for service (after busy, abandonment)
- ▶ · · · , and more.

# **Erlang-A: Practical Relevance?**

#### **Experience:**

- ► Arrival process **not pure Poisson** (time-varying,  $\sigma^2$  too large)
- Service times not Exponential (typically close to LogNormal)
- ▶ Patience times **not Exponential** (various patterns observed).
- Building Blocks need not be independent (eg. long wait possibly implies long service)
- Customers and Servers not homogeneous (classes, skills)
- Customers return for service (after busy, abandonment)
- ▶ · · · , and more.

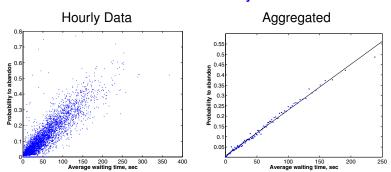
Question: Is Erlang-A Practically Relevant?



# Estimating (Im)Patience: via $P{Ab} \propto E[W_q]$

Assume  $Exp(\theta)$  (im)patience. Then,  $P\{Ab\} = \theta \cdot E[W_q]$ .

#### Israeli Bank: Yearly Data

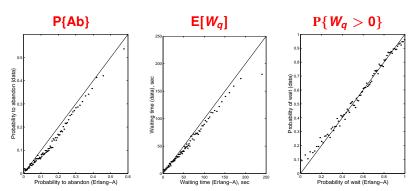


Graphs based on 4158 hour intervals.

Estimate of mean (im)patience:  $250/0.55 \approx 450$  seconds.

# Erlang-A: Fitting a Simple Model to a Complex Reality

- Small Israeli Banking Call-Center (10 agents)
- ▶ (Im)Patience ( $\theta$ ) estimated via P{Ab} / E[ $W_q$ ]
- Graphs: Hourly Performance vs. Erlang-A Predictions, during 1 year (aggregating groups with 40 similar hours).



# **Erlang-A: Simple, but Not Too Simple**

#### **Further Natural Questions:**

- 1. Why does Erlang-A practically work? justify robustness.
- 2. When does it fail? chart boundaries.
- 3. Generalize: time-variation, SBR, networks, uncertainty, ...

# **Erlang-A: Simple, but Not Too Simple**

#### **Further Natural Questions:**

- 1. Why does Erlang-A practically work? justify robustness.
- 2. When does it fail? chart boundaries.
- 3. Generalize: time-variation, SBR, networks, uncertainty, ...

**Answers** via **Asymptotic Analysis**, as load- and staffing-levels increase, which reveals model-essentials:

- ► Efficiency-Driven (ED) regime: Fluid models (deterministic)
- Quality- and Efficiency-Driven (QED): Diffusion refinements.



# **Erlang-A: Simple, but Not Too Simple**

#### **Further Natural Questions:**

- 1. Why does Erlang-A practically work? justify robustness.
- 2. When does it fail? chart boundaries.
- 3. Generalize: time-variation, SBR, networks, uncertainty, ...

**Answers** via **Asymptotic Analysis**, as load- and staffing-levels increase, which reveals model-essentials:

- ► Efficiency-Driven (ED) regime: Fluid models (deterministic)
- Quality- and Efficiency-Driven (QED): Diffusion refinements.

**Motivation**: Moderate-to-large service systems (**100's - 1000's** servers), notably **call-centers**.

Results turn out **accurate** enough to also cover **10-20** servers. Important – relevant to **hospitals** (nurse-staffing: de Véricourt & Jennings, 2006), ...

# **Operational Regimes: Conceptual Framework**

Assume: Offered Load  $R = \frac{\lambda}{\mu}$  (=  $\lambda \times E[S]$ ) not too small.

QD Regime:  $N \approx R + \delta R$   $[(N - R)/R \rightarrow \delta, \text{ as } N, \lambda \uparrow \infty]$ 

▶ Essentially **no** delays:  $[P\{W_q > 0\} \rightarrow 0]$ .

ED Regime:  $N \approx R - \gamma R$ 

- Garnett, M. & Reiman 2003
- Essentially all customers are delayed
- ▶ Wait same order as service-time;  $\gamma$ % Abandon (10-25%).



# **Operational Regimes: Conceptual Framework**

Assume: Offered Load  $R = \frac{\lambda}{\mu}$  (=  $\lambda \times E[S]$ ) not too small.

**QD Regime:**  $N \approx R + \delta R$   $[(N - R)/R \rightarrow \delta, \text{ as } N, \lambda \uparrow \infty]$ 

▶ Essentially no delays:  $[P\{W_q > 0\} \rightarrow 0]$ .

## ED Regime: $N \approx R - \gamma R$

- Garnett, M. & Reiman 2003
- Essentially all customers are delayed
- Wait same order as service-time; γ% Abandon (10-25%).

# QED Regime: $N \approx R + \beta \sqrt{R}$

- Erlang 1924, Halfin & Whitt 1981
- ▶ %Delayed between 25% and 75%
- ▶ Wait one-order below service-time (sec vs. min); 1-5% Abandon.



# **Operational Regimes: Conceptual Framework**

Assume: Offered Load  $R = \frac{\lambda}{\mu}$  (=  $\lambda \times E[S]$ ) not too small.

QD Regime: 
$$N \approx R + \delta R$$
  $[(N - R)/R \rightarrow \delta, \text{ as } N, \lambda \uparrow \infty]$ 

▶ Essentially **no** delays:  $[P\{W_q > 0\} \rightarrow 0]$ .

## ED Regime: $N \approx R - \gamma R$

- ► Garnett, M. & Reiman 2003
- Essentially all customers are delayed
- Wait same order as service-time;  $\gamma$ % Abandon (10-25%).

# QED Regime: $N \approx R + \beta \sqrt{R}$

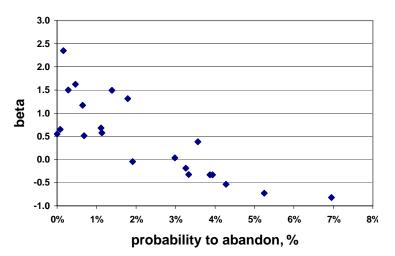
- Erlang 1924, Halfin & Whitt 1981
- %Delayed between 25% and 75%
- ▶ Wait one-order below service-time (sec vs. min); 1-5% Abandon.

# QED+ED: $N \approx (1 - \gamma)R + \beta\sqrt{R}$

- Zeltyn & M. 2006
- ▶ QED refining ED to accommodate "timely-delays":  $P\{W_q > T\}$ .

### **QED: Practical Support**

QOS parameter  $\beta = (N - R)/\sqrt{R}$  vs. %Abandonment



## QED: Theoretical Support (Garnett, M., Reiman '02; Zeltyn '03)

Consider a sequence of M/M/N+G models, N=1,2,3,...

Then the following **points of view** are equivalent:

$$%{Wait > 0} \approx \alpha,$$

$$0 < \alpha < 1$$
;

• Customers 
$$% \{Abandon\} \approx \frac{\gamma}{\sqrt{N}},$$

$$0 < \gamma$$
;

$$OCC \approx 1 - \frac{\beta + \gamma}{\sqrt{N}}$$

OCC 
$$\approx 1 - \frac{\beta + \gamma}{\sqrt{N}}$$
  $-\infty < \beta < \infty$ ;

$$N \approx R + \beta \sqrt{R}$$

• Managers 
$$N \approx R + \beta \sqrt{R}$$
,  $R = \lambda \times E(S)$  not small;

QED performance (ASA, ...) is easily computable, all in terms of  $\beta$  (the square-root safety staffing level) – see later.

# QED Approximations (Zeltyn, M. '06)

G – patience distribution,

 $g_0$  – patience density at origin  $(g_0 = \theta, \text{ if } \exp(\theta)).$ 

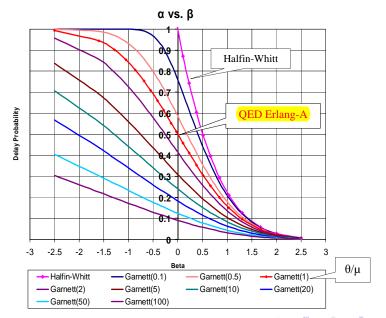
$$\begin{split} \boldsymbol{N} \; &= \; \frac{\lambda}{\mu} + \beta \sqrt{\frac{\lambda}{\mu}} + o(\sqrt{\lambda}) \;, \qquad -\infty < \beta < \infty \;. \\ & \qquad \qquad \mathrm{P}\{\mathrm{Ab}\} \; \approx \; \frac{1}{\sqrt{N}} \cdot \left[h(\hat{\beta}) - \hat{\beta}\right] \cdot \left[\sqrt{\frac{\mu}{g_0}} + \frac{h(\hat{\beta})}{h(-\beta)}\right]^{-1} \;, \\ & \qquad \qquad \mathrm{P}\left\{W > \frac{T}{\sqrt{N}}\right\} \; \approx \; \left[1 + \sqrt{\frac{g_0}{\mu}} \cdot \frac{h(\hat{\beta})}{h(-\beta)}\right]^{-1} \cdot \frac{\bar{\Phi}\left(\hat{\beta} + \sqrt{g_0\mu} \cdot T\right)}{\bar{\Phi}(\hat{\beta})} \;, \\ & \qquad \mathrm{P}\left\{\mathrm{Ab} \; \middle| \; W > \frac{T}{\sqrt{N}}\right\} \; \approx \; \frac{1}{\sqrt{N}} \cdot \sqrt{\frac{g_0}{\mu}} \cdot \left[h\left(\hat{\beta} + \sqrt{g_0\mu} \cdot T\right) - \hat{\beta}\right] \;. \end{split}$$

Here

$$\begin{array}{rcl} \widehat{\beta} & = & \beta\sqrt{\frac{\mu}{g_0}} \\ \bar{\Phi}(x) & = & 1-\Phi(x)\,, \\ h(x) & = & \phi(x)/\bar{\Phi}(x)\,, \ \ \mbox{hazard rate of } N(0,1). \end{array}$$



# Garnett / Halfin-Whitt Functions: $P\{W_q > 0\}$

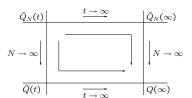


# **Process Limits (Queueing, Waiting)**

•  $\hat{Q}_N = \{\hat{Q}_N(t), t \geq 0\}$  : stochastic process obtained by centering and rescaling:

$$\hat{Q}_N = \frac{Q_N - N}{\sqrt{N}}$$

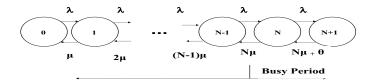
- $\hat{Q}_N(\infty)$  : stationary distribution of  $\hat{Q}_N$
- $\hat{Q} = {\hat{Q}(t), t \geq 0}$ : process defined by:  $\hat{Q}_N(t) \stackrel{d}{\to} \hat{Q}(t)$ .



Approximating (Virtual) Waiting Time

$$\hat{V}_N = \sqrt{N} \ V_N \Rightarrow \hat{V} = \left[ \frac{1}{\mu} \ \hat{Q} \right]^+$$
 (Puhalskii, 1994)

# **QED Intuition via Excursions: Busy/Idle Periods**



Q(0) = N: all servers busy, no queue.

Let 
$$T_{N,N-1}=$$
 Busy Period (down-crossing  $N\downarrow N-1$  )

$$T_{N-1,N} =$$
 Idle Period (up-crossing  $N-1 \uparrow N$ )

Then 
$$P(Wait > 0) = \frac{T_{N,N-1}}{T_{N,N-1} + T_{N-1,N}} = \left[1 + \frac{T_{N-1,N}}{T_{N,N-1}}\right]^{-1}$$



38

# **QED Intuition via Excursions: Asymptotics**

$$\begin{array}{ll} \text{Calculate} & T_{N-1,N} = \frac{1}{\lambda_N E_{1,N-1}} \sim \frac{1}{N\mu \times h(-\beta)/\sqrt{N}} \sim \frac{1}{\sqrt{N}} \cdot \frac{1/\mu}{h(-\beta)} \\ & T_{N,N-1} = \frac{1}{N\mu\pi_+(0)} \sim \frac{1}{\sqrt{N}} \cdot \frac{\beta/\mu}{h(\delta)/\delta}, \quad \delta = \beta\sqrt{\mu/\theta} \\ & \text{Both apply as} & \sqrt{N} \left(1-\rho_N\right) \to \beta, \; -\infty < \beta < \infty. \end{array}$$
 Hence, 
$$P(Wait > 0) \sim \left[1 + \frac{h(\delta)/\delta}{h(-\beta)/\beta}\right]^{-1}.$$

# **QED Intuition via Excursions: Asymptotics**

$$\begin{array}{ll} \text{Calculate} & T_{N-1,N} = \frac{1}{\lambda_N E_{1,N-1}} \sim \frac{1}{N\mu \times h(-\beta)/\sqrt{N}} \sim \frac{1}{\sqrt{N}} \cdot \frac{1/\mu}{h(-\beta)} \\ & T_{N,N-1} = \frac{1}{N\mu\pi_+(0)} \sim \frac{1}{\sqrt{N}} \cdot \frac{\beta/\mu}{h(\delta)/\delta}, \quad \delta = \beta\sqrt{\mu/\theta} \\ & \text{Both apply as} \quad \sqrt{N} \left(1-\rho_N\right) \to \beta, \ -\infty < \beta < \infty. \end{array}$$

Special case: 
$$\mu = \theta$$
:  
Then  $\mathbf{Q} \stackrel{d}{=} \mathbf{M}/\mathbf{M}/\infty$ , since sojourn-time is  $\exp(\mu = \theta)$ ; and  $\mathbf{P}\{\mathbf{Wait} > \mathbf{0}\} \approx \mathbf{1/2}$ , since  $\delta = \beta$ .



# **Dimensioning a Service System**

Operational Regimes provide a conceptual framework.

# **Dimensioning a Service System**

Operational Regimes provide a conceptual framework.

#### Questions:

- 1. How accurate are QD/ED/QED approximations?
- 2. How to determine the regime? QOS parameters?
- 3. Is there a regime robust enough to cover the others?

#### **Dimensioning a Service System**

Operational Regimes provide a **conceptual framework**.

#### Questions:

- 1. How accurate are QD/ED/QED approximations?
- 2. How to determine the regime? QOS parameters?
- 3. Is there a regime robust enough to cover the others?

**Answers**, via many-server **Asymptotic Analysis** (w/ Borst & Reiman, 2004; Zeltyn, 2006):

1. Approximations are **extremely accurate**.

#### **Dimensioning a Service System**

Operational Regimes provide a **conceptual framework**.

#### Questions:

- 1. How accurate are QD/ED/QED approximations?
- 2. How to determine the regime? QOS parameters?
- 3. Is there a regime robust enough to cover the others?

**Answers**, via many-server **Asymptotic Analysis** (w/ Borst & Reiman, 2004; Zeltyn, 2006):

- 1. Approximations are extremely accurate.
- 2. Dimensioning:
  - ► Constraint Satisfaction: eg. Min. *n* s.t. QOS constraints.
  - Cost / Profit Optimization: eg. Min costs of Staffing + Congestion.

#### **Dimensioning a Service System**

Operational Regimes provide a **conceptual framework**.

#### Questions:

- How accurate are QD/ED/QED approximations?
- 2. How to determine the regime? QOS parameters?
- 3. Is there a regime robust enough to cover the others?

**Answers**, via many-server **Asymptotic Analysis** (w/ Borst & Reiman, 2004; Zeltyn, 2006):

- 1. Approximations are **extremely accurate**.
- 2. Dimensioning:
  - ► Constraint Satisfaction: eg. Min. *n* s.t. QOS constraints.
  - Cost / Profit Optimization: eg. Min costs of Staffing + Congestion.
- 3. Robustness depends:
  - Without Abandonment: QED covers all, at amazing accuracy.
  - With Abandonment: ED, QED, ED+QED all have a role.



#### **Operational Regimes: Rules-of-Thumb**

Constraint	P{Ab}		$\mathrm{E}[W]$		$P\{W > T\}$	
	Tight	Loose	Tight	Loose	Tight	Loose
	1-10%	≥ 10%	$\leq 10\% \mathrm{E}[\tau]$	$\geq 10\% \mathrm{E}[\tau]$	$0 \le T \le 10\% \mathrm{E}[\tau]$	$T \geq 10\% \mathrm{E}[\tau]$
Offered Load					$5\% \le \alpha \le 50\%$	$5\% \le \alpha \le 50\%$
Small (10's)	QED	QED	QED	QED	QED	QED
Moderate-to-Large	QED	ED,	QED	ED,	QED	ED+QED
(100's-1000's)		QED		QED if $\tau \stackrel{d}{=} \exp$		

ED: 
$$N \approx R - \gamma R$$
 (0.1  $\leq \gamma \leq$  0.25).

QD: 
$$N \approx R + \delta R$$
 (0.1  $\leq \delta \leq$  0.25).

**QED:** 
$$N \approx R + \beta \sqrt{R}$$
  $(-1 \le \beta \le 1)$ .

**ED+QED:** 
$$N \approx (1 - \gamma)R + \beta \sqrt{R}$$
 ( $\gamma, \beta$  as above).

41

## The ED Regime: M/M/n+G

 ${f ED}-{f E}$ fficiency- ${f D}$ riven.

Assume  $G(x) = \gamma$  has a unique solution  $x^*$  and  $g(x^*) > 0$ .

Staffing:  $n = R \cdot (1 - \gamma) + o(\sqrt{R})$ ,  $0 < \gamma < 1$ .

## The ED Regime: M/M/n+G

**ED** – **E**fficiency-**D**riven.

Assume  $G(x) = \gamma$  has a unique solution  $x^*$  and  $g(x^*) > 0$ .

Staffing:  $n = R \cdot (1 - \gamma) + o(\sqrt{R})$ ,  $0 < \gamma < 1$ .

#### Performance Measures

- $P\{W_q = 0\}$  decreases exponentially in n.
- Probability to abandon converges to:

$$P{Ab} \approx \gamma \approx 1 - \frac{1}{\rho}.$$

• Offered wait converges to  $x^*$ :

$$\mathrm{E}[V] \; \approx \; x^* \,, \qquad V \; \xrightarrow{p} \; x^* \,.$$

• Distribution  $G^*$  of  $\min(x^*, \tau)$ :

$$G^*(x) = \begin{cases} G(x)/\gamma, & x \le x^* \\ 1, & x > x^* \end{cases}$$

Asymptotic distribution of wait:

$$W_q \xrightarrow{w} G^*, \quad E[W_q] \rightarrow E[\min(x^*, \tau)].$$

## ED+QED Regime: Motivation

Min n, s.t.

 $P\{W_q > 0\} \le \alpha$  – use **QED** staffing.

 $E[W_q] \le T$  – use **ED** staffing.

What about  $\mathbf{P}\{W_q > T\} \leq \alpha$ , T > 0? (Most prevalent SL constraint in call centers.)

## ED+QED Regime: Motivation

Min n, s.t.

$$P\{W_q > 0\} \le \alpha$$
 – use **QED** staffing.

$$E[W_q] \le T$$
 – use **ED** staffing.

What about  $\mathbf{P}\{W_q > T\} \leq \alpha$ , T > 0? (Most prevalent SL constraint in call centers.)

## ED approximation:

$$P\{W_q > T\} \approx \begin{cases} \bar{G}(T), & T < G^{-1}(\gamma), \\ 0, & T > G^{-1}(\gamma). \end{cases}$$

or (as a function of "staffing"):

$$\mathrm{P}\{W_q > T\} \; \approx \; \begin{cases} \bar{G}(T) \,, \; \gamma > G(T) \,, \\ 0 \,, \qquad \gamma < G(T) \,, \end{cases}$$

**Too crude** to capture  $\alpha$  exactly.

**Solution:** Refine around  $\gamma = G(T)$ ,

$$n = (1 - G(T)) \cdot R + \beta \sqrt{R}, \quad -\infty < \beta < \infty.$$

## ED+QED Performance Measures

**Theorem.** The following statements are equivalent:

- 1. Staffing level:  $n = (1 \gamma)R + \beta\sqrt{R} + o(\sqrt{R})$ ;
- 2. Tail probability:  $P\{W_q > T\} = \alpha + o(1)$ ;
- 3. Probability to abandon:

$$P\{Ab\} = \gamma - \frac{\beta}{\sqrt{R}} + o\left(\frac{1}{\sqrt{R}}\right);$$

4. Average wait:

$$E[W_q] = \int_0^T \bar{G}(u)du - \frac{\beta}{\sqrt{R}} \cdot \frac{1}{h_G(T)} + o\left(\frac{1}{\sqrt{R}}\right).$$

## ED+QED Performance Measures

**Theorem.** The following statements are equivalent:

- 1. Staffing level:  $n = (1 \gamma)R + \beta\sqrt{R} + o(\sqrt{R})$ ;
- 2. Tail probability:  $P\{W_q > T\} = \alpha + o(1)$ ;
- 3. Probability to abandon:

$$P\{Ab\} = \gamma - \frac{\beta}{\sqrt{R}} + o\left(\frac{1}{\sqrt{R}}\right);$$

4. Average wait:

$$E[W_q] = \int_0^T \bar{G}(u)du - \frac{\beta}{\sqrt{R}} \cdot \frac{1}{h_G(T)} + o\left(\frac{1}{\sqrt{R}}\right).$$

Here  $0 < \alpha < \bar{G}(T)$ ,  $\gamma = G(T)$ ,  $h_G(T)$  = patience hazard-rate at T and

$$\beta \ = \ \bar{\Phi}^{-1} \left( \frac{\alpha}{\bar{G}(T)} \right) \cdot \sqrt{\frac{g(T)}{\mu}} \, .$$

Corollary. Approximation for the tail probability:

Note: If  $\alpha \geq \bar{G}(T)$  then n = 0 satisfies  $P\{W_q > T\} \leq \alpha$ .

# M/M/n+G Performance Measures: Building Blocks

$$H(x) \triangleq \int_0^x \bar{G}(u) du$$

where  $\bar{G}(\cdot) = 1 - G(\cdot)$ , the survival-function of patience.

$$J \triangleq \int_0^\infty \exp\left\{\lambda H(x) - n\mu x\right\} dx,$$

$$J_1 \triangleq \int_0^\infty x \cdot \exp\left\{\lambda H(x) - n\mu x\right\} dx,$$

$$J_H \triangleq \int_0^\infty H(x) \cdot \exp\left\{\lambda H(x) - n\mu x\right\} dx,$$

$$J(t) \triangleq \int_t^\infty \exp\left\{\lambda H(x) - n\mu x\right\} dx.$$

$$J_1(t) \triangleq \int_t^\infty x \cdot \exp\left\{\lambda H(x) - n\mu x\right\} dx,$$

$$J_H(t) \triangleq \int_t^\infty H(x) \cdot \exp\left\{\lambda H(x) - n\mu x\right\} dx.$$

Finally,

$$\mathcal{E} \triangleq \frac{\sum_{j=0}^{n-1} \frac{1}{j!} \left(\frac{\lambda}{\mu}\right)^j}{\frac{1}{(n-1)!} \left(\frac{\lambda}{\mu}\right)^{n-1}} = \int_0^\infty e^{-t} \left(1 + \frac{t\mu}{\lambda}\right)^{n-1} dt.$$

## M/M/n+G Performance Measures

 $P{Ab}$  = probability to abandon,  $W_q$  = waiting time, V = offered wait, Q = queue length.

$$\begin{split} & \text{P}\{V > t\} \, = \, \frac{\lambda J(t)}{\mathcal{E} + \lambda J}, \; (\textbf{Baccelli \& Hebuterne}, \, 1981) \\ & \text{P}\{W_q > 0\} \, = \, \frac{\lambda J}{\mathcal{E} + \lambda J} \cdot \bar{G}(0) \, , \\ & \text{P}\{\text{Ab}\} \, = \, \frac{1 + (\lambda - n\mu)J}{\mathcal{E} + \lambda J} \, , \\ & \text{E}[V] \, = \, \frac{\lambda J_1}{\mathcal{E} + \lambda J} \, , \\ & \text{E}[W_q] \, = \, \frac{\lambda^2 J_H}{\mathcal{E} + \lambda J} \, , \\ & \text{E}[Q] \, = \, \frac{\lambda^2 J_H}{\mathcal{E} + \lambda J} \, , \\ & \text{E}[W_q \mid \text{Ab}] \, = \, \frac{J + \lambda J_H - n\mu J_1}{(\lambda - n\mu)J + 1} \, , \\ & \text{P}\{W_q > t\} \, = \, \frac{\lambda \bar{G}(t)J(t)}{\mathcal{E} + \lambda J} \, , \\ & \text{E}[W_q \mid W_q > t] \, = \, \frac{J_H(t) - (H(t) - t\bar{G}(t)) \cdot J(t)}{\bar{G}(t)J(t)} \, , \\ & \text{P}\{\text{Ab} \mid W_q > t\} \, = \, \frac{\lambda - n\mu - G(t)}{\lambda \bar{G}(t)} + \frac{\exp\{\lambda H(t) - n\mu t\}}{\lambda \bar{G}(t)J(t)} \, . \end{split}$$

## M/M/n+G: Laplace Method

Asymptotic calculation of integrals:

- 1. Show that the integral (mass) is concentrated near a **certain point**.
- 2. Use **Taylor expansion** to approximate integrand near this point.

Apply to **Building Blocks** and **Performance Measures** above.

### **Examples:**

QED regime: 
$$n = \frac{\lambda}{\mu} + \beta \sqrt{\frac{\lambda}{\mu}} + o(\sqrt{\lambda}).$$

$$J = \frac{1}{\sqrt{n}} \cdot \frac{1}{\sqrt{\mu g_0}} \cdot \frac{1}{h_G(\hat{\beta})} + o\left(\frac{1}{\sqrt{n}}\right),\,$$

where

$$\hat{\beta} \triangleq \beta \sqrt{\frac{\mu}{g_0}}.$$

$$ext{ED+QED regime:} \ n = ar{G}(T) \cdot rac{\lambda}{\mu} + eta iggl(rac{\lambda}{\mu} + o(\sqrt{\lambda}).$$

$$J \sim \exp\{\lambda H(T) - n\mu T\} \cdot \exp\left\{\frac{\beta^2 \mu}{2g(T)}\right\} \cdot \sqrt{\frac{2\pi}{\lambda g(T)}}.$$

#### Back to "Why does Erlang-A Work?"

Theoretical Answer:  $M_t^J/G/N_t + G \approx (M/M/N + M)_t, t \geq 0.$ 

- General Patience: Behavior at the origin is all that matters.
- ► General Services: Empirical insensitivity beyond the mean.
- ► Time-Varying Arrivals: Modified Offered-Load approximations.
- ► Heterogeneous Customers: 1-D state collapse.



#### Back to "Why does Erlang-A Work?"

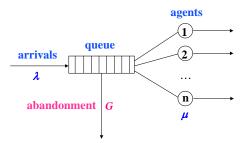
Theoretical Answer:  $M_t^J/G/N_t + G \approx (M/M/N + M)_t, t \geq 0.$ 

- General Patience: Behavior at the origin is all that matters.
- ► General Services: Empirical insensitivity beyond the mean.
- ► Time-Varying Arrivals: Modified Offered-Load approximations.
- ► Heterogeneous Customers: 1-D state collapse.

**Practically**: Why do (stochastic-ignorant) Call Centers work? "The right answer for the wrong reason"



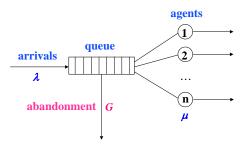
#### "Why does Erlang-A Work?" General Patience



(Im)Patience times Generally Distributed: M/M/n+G

**Exact** analysis in steady-state (Baccelli & Hebuterne, 1981): solve Kolmogorov's PDE's (semi-Markov) for the offered-wait *V*.

#### "Why does Erlang-A Work?" General Patience



(Im)Patience times **Generally Distributed**: M/M/n+G

**Exact** analysis in steady-state (Baccelli & Hebuterne, 1981): solve Kolmogorov's PDE's (semi-Markov) for the offered-wait **V**.

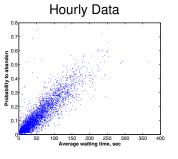
**QED** analysis (w/ Zeltyn, 2006):  $\mathbf{n} \approx \mathbf{R} + \beta \sqrt{\mathbf{R}}$ .

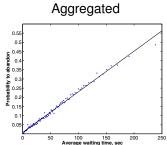
- Assume (Im)Patience density g(0) > 0.
- ▶ V asymptotics  $(\lambda \uparrow \infty)$ : Laplace Method.
- ▶ QED Approximations: Use Erlang-A, with  $\theta \leftrightarrow g(0)$ .



#### **General Patience: Fitting Erlang-A**

#### Israeli Bank: Yearly Data





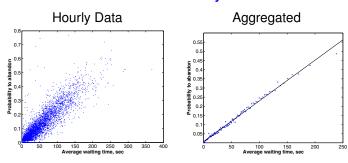
#### Theory:

**Erlang-A:** 
$$P\{Ab\} = \theta \cdot E[W_q];$$

$$M/M/N+G$$
: P{Ab}  $\approx g(0) \cdot E[W_q]$ .

#### **General Patience: Fitting Erlang-A**

#### Israeli Bank: Yearly Data



#### Theory:

**Erlang-A:**  $P{Ab} = \theta \cdot E[W_q];$ 

M/M/N+G: P{Ab}  $\approx g(0) \cdot E[W_q]$ .

#### Recipe:

In both cases, use Erlang-A, with  $\hat{\theta} = \widehat{P}\{Ab\}/\widehat{E}[W_q]$  (slope above).



Established:  $M/M/N+G \approx M/M/N+M$  ( $\theta = g(0)$ ).

47

Established:  $M/M/N+G \approx M/M/N+M$  ( $\theta = g(0)$ ).

**Now:**  $M/G/N+G \approx M/M/N+G$  (E[S] same in both).

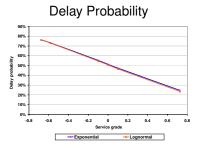
Established:  $M/M/N+G \approx M/M/N+M$   $(\theta = g(0))$ .

**Now:**  $M/G/N+G \approx M/M/N+G$  (E[S] same in both).

**Numerical Experiments:** Whitt (2004), Rosenshmidt (2006) demonstrate a **useful fit** for typical call-center parameters.

Lognormal (CV=1) vs. Exponential Service Times, QED Regime; 100 agents, average patience = average service

## 

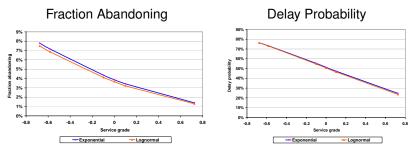


Established:  $M/M/N+G \approx M/M/N+M$  ( $\theta = g(0)$ ).

**Now:**  $M/G/N+G \approx M/M/N+G$  (E[S] same in both).

**Numerical Experiments:** Whitt (2004), Rosenshmidt (2006) demonstrate a **useful fit** for typical call-center parameters.

Lognormal (CV=1) vs. Exponential Service Times, QED Regime; 100 agents, average patience = average service



**QED G-Services**: G/GI/N (Reed, 2007), G/D<sub>K</sub>/N+G (w/ Momčilović),

4 미 > 〈@ > 〈 본 > 〈 経 > 〈 본 > 〈 経 > 〈 본 > 〈 본 > 〈 経 >

Established:  $M/G/N+G \approx M/M/N+M$   $(\theta = g(0))$ .

Established:  $M/G/N+G \approx M/M/N+M$   $(\theta = g(0))$ .

**Now:**  $M_t/G/N_t+G\approx (M/G/N+G)_t$   $(N_t,\lambda \text{ well chosen}).$ 

48

Established:  $M/G/N+G \approx M/M/N+M$  ( $\theta = g(0)$ ).

**Now:**  $M_t/G/N_t + G \approx (M/G/N + G)_t$   $(N_t, \lambda \text{ well chosen}).$ 

Two steps (Feldman, M., Massey & Whitt, 2006):

- Modified Offered-Load: λ
  - ▶ Consider  $M_t/G/N_t + G$  with arrival rate  $\lambda(t)$ ,  $t \ge 0$ .
  - Approximate its **time-varying** performance **at time** t with a **stationary**  $M/G/N_t + G$ , in which  $\lambda = E\lambda(t S_e)$ . ( $S_e \stackrel{d}{=}$  residual-service: congestion-lag behind peak-load.)

Established:  $M/G/N+G \approx M/M/N+M$  ( $\theta = g(0)$ ).

**Now:**  $M_t/G/N_t + G \approx (M/G/N + G)_t$   $(N_t, \lambda \text{ well chosen}).$ 

Two steps (Feldman, M., Massey & Whitt, 2006):

- Modified Offered-Load: λ
  - ▶ Consider  $M_t/G/N_t + G$  with arrival rate  $\lambda(t)$ ,  $t \ge 0$ .
  - Approximate its **time-varying** performance at **time** t with a stationary  $M/G/N_t + G$ , in which  $\lambda = E\lambda(t S_e)$ . ( $S_e \stackrel{d}{=}$  residual-service: congestion-lag behind peak-load.)
- 2. Square-Root Staffing: Nt
  - Let  $R_t = E\lambda(t S_e) \times ES$  be the Offered-Load at time t ( $R_t = Number-in-system in a corresponding <math>M_t/G/\infty$ .)
  - Staff  $N_t = R_t + \beta \sqrt{R_t}$ .



48

Established:  $M/G/N+G \approx M/M/N+M$  ( $\theta = g(0)$ ).

**Now:**  $M_t/G/N_t + G \approx (M/G/N + G)_t$   $(N_t, \lambda \text{ well chosen}).$ 

Two steps (Feldman, M., Massey & Whitt, 2006):

- Modified Offered-Load: λ
  - ▶ Consider  $M_t/G/N_t + G$  with arrival rate  $\lambda(t)$ ,  $t \ge 0$ .
  - Approximate its **time-varying** performance at **time** t with a stationary  $M/G/N_t + G$ , in which  $\lambda = E\lambda(t S_e)$ . ( $S_e \stackrel{d}{=}$  residual-service: congestion-lag behind peak-load.)
- 2. Square-Root Staffing: N<sub>t</sub>
  - Let  $R_t = E\lambda(t S_e) \times ES$  be the **Offered-Load** at time t ( $R_t = Number-in-system in a corresponding <math>M_t/G/\infty$ .)
  - Staff  $N_t = R_t + \beta \sqrt{R_t}$ .

**Serendipity: Time-stable** performance, supported by **ISA** = Iterative Staffing Algorithm, and QED diffusion limits  $(M_t/M/N + M, \mu = \theta)$ .



## Time-Varying Arrivals: √ Safety-Staffing

Model  $M_t/M/N_t + M$ 

Parameters  $\lambda(t)$   $\mu$  ?  $\theta$ 

 $\mathbf{N_t} = \mathbf{R_t} + \beta \sqrt{\mathbf{R_t}}$ 

## Time-Varying Arrivals: √ Safety-Staffing

Model 
$$M_t/M/N_t + M$$

Parameters  $\lambda(t)$   $\mu$  ?  $\theta$ 

$$\mathbf{N}_{t} = \mathbf{R}_{t} + \beta \sqrt{\mathbf{R}_{t}}$$

$$\mu = \theta$$
:  $L_t \stackrel{d}{=} Poisson(R_t) \stackrel{d}{\approx} N(R_t, R_t)$ , since  $M_t / M / \infty$ 

$$R_t = E\lambda(t-S) \cdot E(S) = E \int_{t-S}^{t} \lambda(u) du$$
 offered load

## Time-Varying Arrivals: √. Safety-Staffing

Model 
$$M_t/M/N_t + M$$

Parameters  $\lambda(t)$   $\mu$  ?  $\theta$ 

$$\mathbf{N}_{t} = \mathbf{R}_{t} + \beta \sqrt{\mathbf{R}_{t}}$$

$$\mu = \theta$$
:  $L_t \stackrel{d}{=} Poisson(R_t) \stackrel{d}{\approx} N(R_t, R_t)$ , since  $M_t / M / \infty$ 

$$R_t = E\lambda(t-S) \cdot E(S) = E \int_{t-S}^{t} \lambda(u) du$$
 offered load

Given 
$$L_t \approx R_t + Z\sqrt{R_t}$$
,  $Z \stackrel{d}{=} N(0,1)$ 

choose 
$$N_t = R_t + \beta \sqrt{R_t}$$

$$\Rightarrow \quad \alpha = P(W_t > 0) \approx P(L_t \ge N_t) = P(Z \ge \beta) = 1 - \phi(\beta)$$

$$\Rightarrow \beta = \phi^{-1} (1 - \alpha)$$
 time-stable  $\alpha \equiv P(W_t > 0)$ ?

## Time-Varying Arrivals: √ Safety-Staffing

Model 
$$M_t/M/N_t + M$$

Parameters  $\lambda(t)$   $\mu$  ?  $\theta$ 

$$\mathbf{N}_{t} = \mathbf{R}_{t} + \beta \sqrt{\mathbf{R}_{t}}$$

$$\mu = \theta$$
:  $L_t \stackrel{d}{=} Poisson(R_t) \stackrel{d}{\approx} N(R_t, R_t)$ , since  $M_t / M / \infty$ 

$$R_t = E\lambda(t-S) \cdot E(S) = E \int_{t-S}^{t} \lambda(u) du$$
 offered load

Given 
$$L_t \approx R_t + Z\sqrt{R_t}$$
,  $Z \stackrel{d}{=} N(0,1)$ 

choose 
$$N_t = R_t + \beta \sqrt{R_t}$$

$$\Rightarrow \quad \alpha = P(W_t > 0) \underset{\text{PASTA}}{\approx} P(L_t \ge N_t) = P(Z \ge \beta) = 1 - \phi(\beta)$$

$$\Rightarrow \beta = \phi^{-1} (1 - \alpha)$$
 time-stable  $\alpha \equiv P(W_t > 0)$ ?

Indeed, but in fact TIME-STABLE PERFORMANCE

## Time-Varying Arrivals: √ Safety-Staffing

Model 
$$M_t/M/N_t + M$$

Parameters  $\lambda(t)$   $\mu$  ?  $\theta$ 

$$\mathbf{N}_{t} = \mathbf{R}_{t} + \beta \sqrt{\mathbf{R}_{t}}$$

$$\mu = \theta$$
:  $L_t \stackrel{d}{=} Poisson(R_t) \stackrel{d}{\approx} N(R_t, R_t)$ , since  $M_t / M / \infty$ 

$$R_t = E\lambda(t-S) \cdot E(S) = E \int_{t-S}^{t} \lambda(u) du$$
 offered load

Given 
$$L_t \approx R_t + Z\sqrt{R_t}$$
,  $Z \stackrel{d}{=} N(0,1)$ 

choose 
$$N_t = R_t + \beta \sqrt{R_t}$$

$$\Rightarrow \quad \alpha = P(W_t > 0) \underset{\text{PASTA}}{\approx} P(L_t \ge N_t) = P(Z \ge \beta) = 1 - \phi(\beta)$$

$$\Rightarrow \beta = \phi^{-1} (1 - \alpha)$$
 time-stable  $\alpha \equiv P(W_t > 0)$ ?

Indeed, but in fact

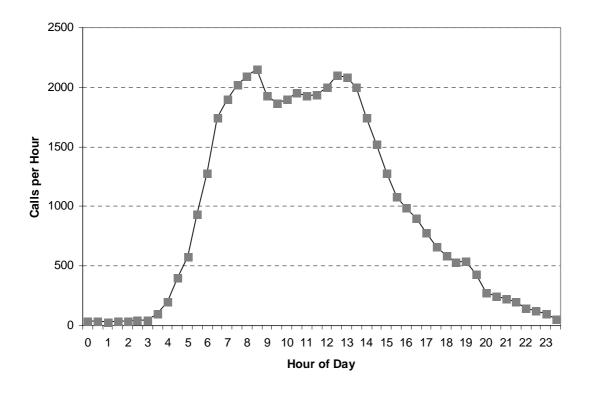
TIME-STABLE PERFORMANCE

 $(\mu \neq \theta, \text{ or generally : Iterative Simulation-Based Algorithm})$ 

## Example: "Real" Call Center

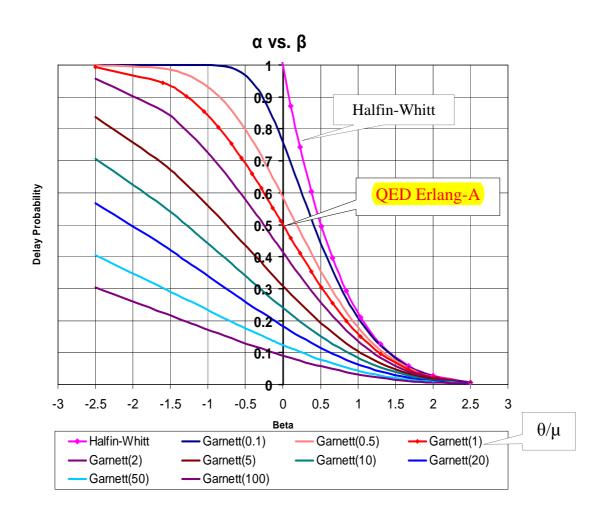
(The "Right Answer" for the "Wrong Reasons")

Time-Varying (two-hump) arrival functions common (Adapted from Green L., Kolesar P., Soares J. for benchmarking.)



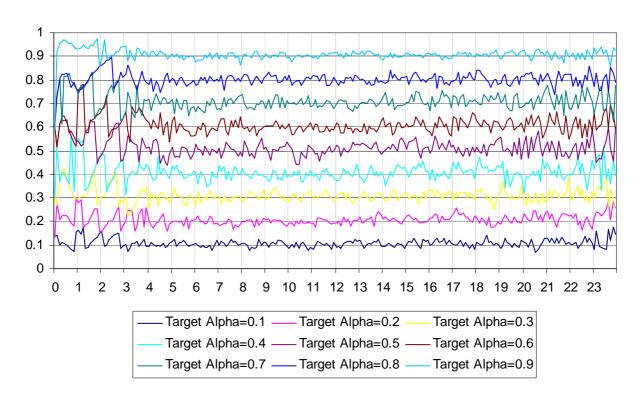
Assume: Service and abandonment times are both Exponential, with mean 0.1 (6 min.)

## **HW/GMR** Delay Functions



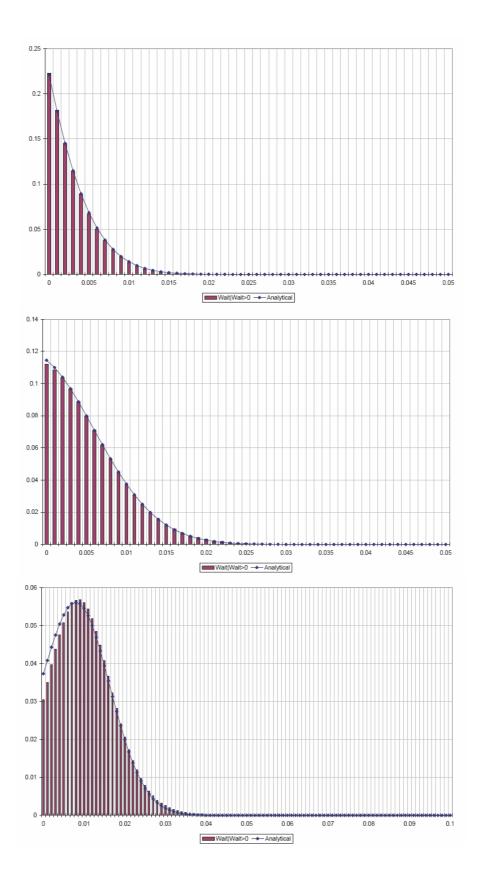
## Delay Probability α

#### **Delay Probability**



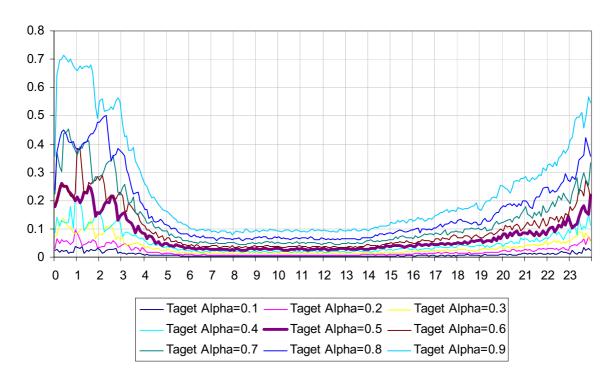
Real Call Center: Empirical waiting time, given positive wait

(1)  $\alpha$ =0.1 (QD) (2)  $\alpha$ =0.5 (QED) (3)  $\alpha$ =0.9 (ED)

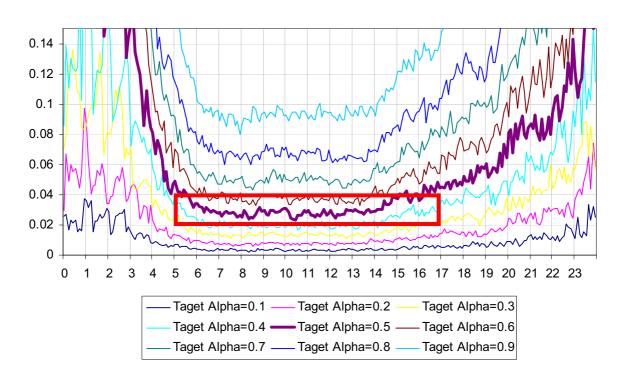


# **Abandon Probability**

## **Abandon Probability**



## **Abandon Probability**



# The "Right Answer" (for the "Wrong Reasons")

**Prevalent Practice** 

$$N_t = \lceil \lambda(t) \cdot E(S) \rceil$$
 (PSA)

"Right Answer" 
$$N_t \approx R_t + \beta \cdot \sqrt{R_t}$$
 (MOL)

$$R_t = E\lambda(t - S) \cdot E(S)$$

# The "Right Answer" (for the "Wrong Reasons")

Prevalent Practice 
$$N_t = \lceil \lambda(t) \cdot E(S) \rceil$$
 (PSA)

"Right Answer" 
$$N_t \approx R_t + \beta \cdot \sqrt{R_t}$$
 (MOL)

$$R_t = E\lambda(t - S) \cdot E(S)$$

Practice 
$$\approx$$
 "Right"  $\beta \approx 0$  (QED)

and  $\lambda(t) \approx$  stable over service-durations

Practice Improved  $N_t = \lceil \lambda [t - E(S)] \cdot E(S) \rceil$ 

# The "Right Answer" (for the "Wrong Reasons")

Prevalent Practice 
$$N_t = \lceil \lambda(t) \cdot E(S) \rceil$$
 (PSA)

"Right Answer" 
$$N_t \approx R_t + \beta \cdot \sqrt{R_t}$$
 (MOL) 
$$R_t = E\lambda(t-S) \cdot E(S)$$

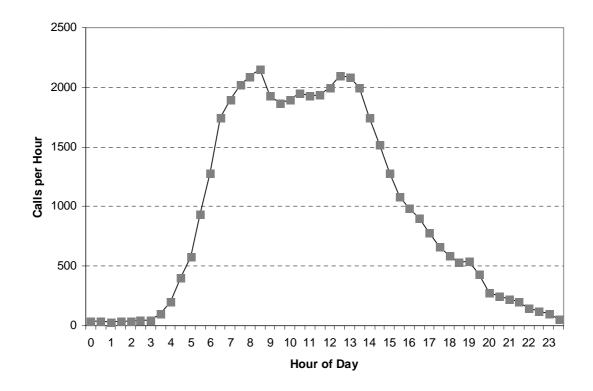
Practice 
$$\approx$$
 "Right"  $\beta \approx 0$  (QED)

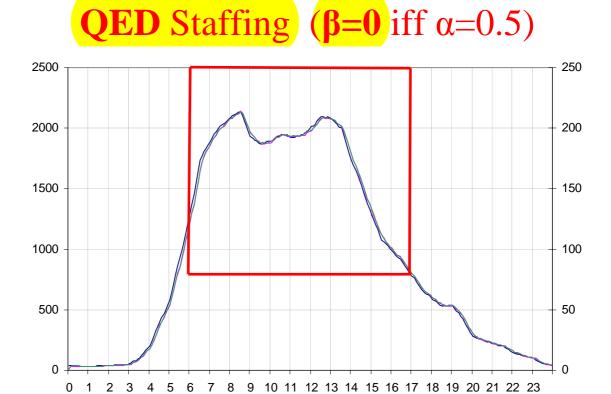
and  $\lambda(t) \approx$  stable over service-durations

Practice Improved 
$$N_t = \lceil \lambda [t - E(S)] \cdot E(S) \rceil$$

When Optimal? for moderately-patient customers:

- 1. Satisfization ⇔ At least 50% to be serve immediately
- 2. Optimization  $\Leftrightarrow$  Customer-Time = 2 x Agent-Salary





Staffing

Arrived

Offered Load

**Now:**  $M_t^J/G/N_t + G \approx (M^J/G/N + G)_t$  (well staffed & controlled).

**Service Levels**: Class 1 = VIP, ..., Class J = best-effort.

**Now:**  $M_t^J/G/N_t + G \approx (M^J/G/N + G)_t$  (well staffed & controlled).

**Service Levels**: Class 1 = VIP, ..., Class J = best-effort.

Staffing, Control (w/ Gurvich & Armony 2005; Feldman & Gurvich):

- ▶ Consider  $M_t^J/G/N_t + G$  with arrival rates  $\lambda_i(t), t \geq 0$ .
- Assume i.i.d. servers.
- ▶ Let  $R_t = E \sum_i \lambda_i (t S_e) \times ES$  be the **Offered-Load** at time t.

**Now:**  $M_t^J/G/N_t + G \approx (M^J/G/N + G)_t$  (well staffed & controlled).

**Service Levels**: Class 1 = VIP, ..., Class J = best-effort.

#### Staffing, Control (w/ Gurvich & Armony 2005; Feldman & Gurvich):

- ▶ Consider  $M_t^J/G/N_t + G$  with arrival rates  $\lambda_i(t), t \geq 0$ .
- Assume i.i.d. servers.
- ▶ Let  $R_t = E \sum_i \lambda_i (t S_e) \times ES$  be the **Offered-Load** at time t.
- ▶ **Staff**  $N_t = R_t + \beta \sqrt{R_t}$ , with  $\beta$  determined by a desired QED performance for the lowest-priority class J.
- ► Control via threshold priorities, where the thresholds are determined by ISA according to desired service levels.
- Approximate time-varying performance at time t with a stationary threshold-controlled  $M^J/G/N_t+G$ , in which  $\lambda_i = \mathbb{E}\lambda_i(t-S_e)$ .



**Now:**  $M_t^J/G/N_t + G \approx (M^J/G/N + G)_t$  (well staffed & controlled).

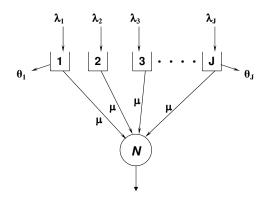
**Service Levels**: Class 1 = VIP, ..., Class J = best-effort.

#### Staffing, Control (w/ Gurvich & Armony 2005; Feldman & Gurvich):

- ▶ Consider  $M_t^J/G/N_t + G$  with arrival rates  $\lambda_i(t), t \geq 0$ .
- Assume i.i.d. servers.
- ▶ Let  $R_t = E \sum_i \lambda_j (t S_e) \times ES$  be the **Offered-Load** at time t.
- ▶ **Staff**  $N_t = R_t + \beta \sqrt{R_t}$ , with  $\beta$  determined by a desired QED performance for the lowest-priority class J.
- Control via threshold priorities, where the thresholds are determined by ISA according to desired service levels.
- Approximate time-varying performance at time t with a stationary threshold-controlled  $M^J/G/N_t+G$ , in which  $\lambda_j = \mathbb{E}\lambda_j(t-S_e)$ .

**Serendipity: Multi-Class Multi-Skill**, w/ **class-dependent** services. Support: ISA, QED diffusion limits (Atar, M. & Shaikhet, 2007).

# The V-Design



- J customer classes: **Arrivals** Poisson $(\lambda_j)$ , **Patience**  $Exp(\theta_j)$ .
- N iid servers: service durations Exp(μ).
- Delay probabilities  $\alpha_1 < \alpha_2 \ldots < \alpha_J$

Objective: Min N, subject to service level differentiation:

$$P\{W_j > 0\} \le \alpha_j, \ j = 1, \dots, J.$$

## **Proposed solution:**

- Static priorities  $1>2>\dots$  with thresholds  $S_1>S_2>\dots$  i.e. a class-j customer served if it is of the present highest-priority and the number of idle servers is  $S_j$  or more.
- Performance analysis in steady-state with no abandonments (Schaack & Larson 1986).

# The V-Design in the QED regime

Gurvich, Armony and Mandelbaum ('06)

Consider a sequence indexed by N = 1, 2, ...

Thresholds: 
$$0 = S_1^N \le S_2^N \le ... \le S_J^N \le N$$

Assume: 
$$S_J^N = o(\sqrt{N})$$

Assume: 
$$\lambda_j^N/\lambda^N\mu \to 
ho_j > 0 \,, \forall j$$
 (all classes non-negligible)

Then the following **points of view** are equivalent:

• QED: 
$$\lim_{N\to\infty} P_N\{W_J>0\} = \alpha_J, \quad 0<\alpha_J<1;$$

• Customer: 
$$\lim_{N\to\infty} \sqrt{N} P_N \{ Ab_J \} = \Delta, \quad 0 < \Delta < \infty;$$

• Server: 
$$\lim_{N\to\infty}\sqrt{N}\ (1-\rho_N)=\beta, \quad -\infty<\beta<\infty;$$

• Manager: 
$$N = R + \beta \sqrt{R} + o(\sqrt{R})$$
,  $R = \lambda/\mu$  large.

Here,  $\alpha_J$  and  $\Delta$  are determined by the QED M/M/N + M.

## Erlang-A is all that is needed:

$$\limsup_{N\to\infty} P\{\overline{W_j^N}>0\} \leq \underline{\alpha_j}, \text{ if } \lim\inf_{N\to\infty} S_{j+1}^N - S_j^N \geq \frac{\ln\alpha_j/\alpha_{j+1}}{\ln\sum_{j=1}^j \rho_i}$$

# The V-Design in the QED regime

Armony, Gurvich and Mandelbaum ('06)

## High priorities (j < J):

- $W_j|W_j > 0$  as in  $M(\lambda_1)/M(N\mu)/1 + M(\theta_1)$ .
- $W_j|W_j>0\stackrel{d}{=}\Theta\left(\frac{1}{N}\right)$ , but
- $W_j \stackrel{d}{=} o\left(\frac{1}{N}\right)$  for  $S_j = O(\log(N))$ .

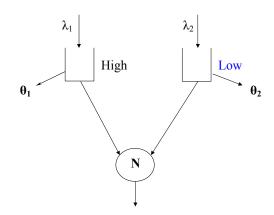
Low priority  $(W_J)$  unchanged as long as  $S_J^N = o(\sqrt{N})$ .

## Comparison to static priorities:

- Static: All classes share same delay probability  $P\{N \text{ servers busy}\}$
- Thresholds:  $P\{W_j > 0\}$  differentiated through thresholds.

 $o(\sqrt{N})$  thresholds guarantee QD waiting times for high-priority and, at the same time, are not hurting the low-priority (who are still QED-served).

# What's Going On?



Thresholds:  $0 = S_1^N \le S_2^N \le N$ 

1. # High Priority  $Q_1^N$ : "See" an  $M(\lambda_1)/M(N\mu)/1 + M(\theta_1)$  queue in light traffic.

$$Q_1^N = o(\sqrt{N}).$$

Only  $Q_2^N$  and  $\theta_2$  play a role in asymptotics.

2. **# Busy Servers**  $B^N$ : Only high priority enter service when  $B^N \geq N - S_2^N$ .  $(B^N - S_2^N)^+$  is in light traffic:

$$(B^N - S_2^N)^+ = o(\sqrt{N})$$

## **Conclusion:**

$$B^{N} + Q_{1}^{N} + Q_{2}^{N} \approx M/M/(N - S_{2}^{N}) + M(\theta_{2})$$
  
  $\approx M/M/N + M(\theta_{2}) \quad (S_{2}^{N} = o(\sqrt{N}))$ 

Asymptotic equivalence to M/M/N + M

# Additional Simple (QED) Models of Complex Realities: Exponential Services; i.i.d. Customers, i.i.d. Servers

#### Performance Analysis:

- Khudiakova, Feigin, M. (Semi-Open): Call-Center + IVR/VRU;
- De Véricourt, Jennings (Closed + Delay), then w/ Yom-Tov (Semi-Open): Nurse staffing (ratios), bed sizing;
- Randhawa, Kumar (Closed + Loss): Subscriber queues.
- ▶ Optimal Staffing: Accurate to within 1, even with very small *n*'s, for both constraint-satisfaction and cost/revenue optimization (staffing, abandonment and waiting costs).
  - Armony, Maglaras: (M<sub>x</sub>/M/N) Delay information (Equilibrium);
  - ▶ Borst, M., Reiman (M/M/N): Asymptotic framework;
  - Zeltyn, M. (M/M/N+G): Optimization still ongoing.
- ► Time-Varying Queues, via 2 approaches:
  - Jennings, M., Massey, Whitt, then w/ Feldman: Time-Stable Performance (ISA, leading to Modified Offered Load);
  - M., Massey, Reiman, Rider, Stolyar: Unavoidable Time-Varying Performance (Fluid & Diffusion models, via Uniform Acceleration).

## Less-Simple (QED) Models: General Service-Times

The Challenge: Must keep track of the state of n individual servers, as  $n \uparrow \infty$ . (Recall Kiefer & Wolfowitz).

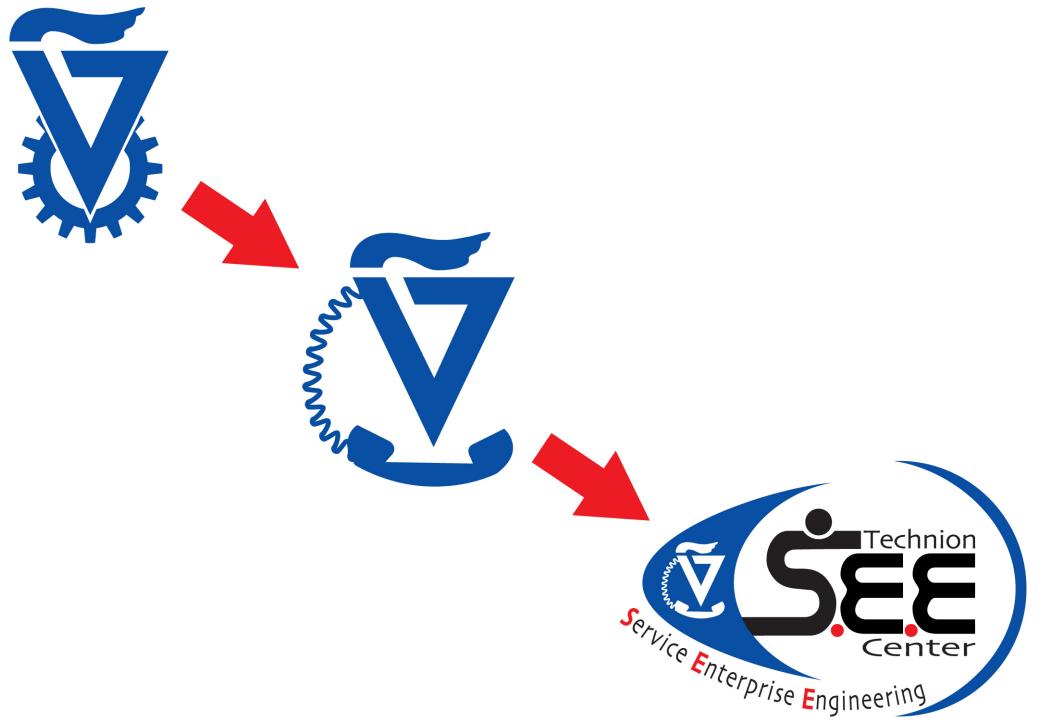
- ► Shwartz, M. (M/G/N), Rosenshmidt, M. (M/G/N+G): Simulations; LogNormal better then Exp, 2-valued same as D.
- Whitt (GI/M+0/N): Covering CV > 1;
- Puhalskii, Reiman (GI/PH/N): Markovian process-limits (no steady-state); also priorities;
- ▶ Jelencović, M., Momčilović (GI/D/N): steady-state (via round-robin); then M., Momčilović (G/D<sub>K</sub>/N): process-limits, via "Lindley-Trees"; G/D<sub>K</sub>/N+G ongoing.
- ► Kaspi, Ramanan (G/G/N): Fluid, next Diffusion (measure-valued ages, following Kiefer & Wolfowitz);
- ▶ Reed (GI/GI/N): Fluid, Diffusion (Skorohod-Like Mapping).

# Complex (QED) Models: Skills-Based Routing (Heterogeneous Customers or/and Servers - Theory)

- V-Model: Harrison, Zeevi; Atar, M., Reiman; Gurvich, M., Armony; then Class-dependent services: Atar, M., Shaikhet;
- Reversed-V: Armony, M.;
   then Pool-dependent services: Dai, Tezcan; Gurvich, Whitt (G-cμ); Atar, M., Shaikhet (Abandonment);
- General: Atar, then w/ Shaikhet (Null-controllability, Throughput-suboptimality); Gurvich, Whitt (FQR);
- ▶ Distributed Networks: Tezcan;
- ▶ Random Service Rates: Atar (Fastest or longest-idle server).

## The Technion SEE Center / Laboratory





- ▶ Technion: P. Feigin, V. Trofimov, Statistics / SEE Laboratory.
- ▶ Wharton: L. Brown, N. Gans, H. Shen (UNC).
- ▶ industry:
  - U.S. Bank: 2.5 years, 220M calls, 40M by 1000 agents.
  - Israeli Cellular: 2.5 years, 110M calls, 25M calls by 750 agents; ongoing.

- ▶ Technion: P. Feigin, V. Trofimov, Statistics / SEE Laboratory.
- ▶ Wharton: L. Brown, N. Gans, H. Shen (UNC).
- industry:
  - U.S. Bank: 2.5 years, 220M calls, 40M by 1000 agents.
  - Israeli Cellular: 2.5 years, 110M calls, 25M calls by 750 agents; ongoing.

**Project Goal:** Designing and Implementing a (universal) data-base/data-repository and interface for storing, retrieving, analyzing and displaying **Call-by-Call-based Data** / **Information**.

- ▶ Technion: P. Feigin, V. Trofimov, Statistics / SEE Laboratory.
- ▶ Wharton: L. Brown, N. Gans, H. Shen (UNC).
- ▶ industry:
  - U.S. Bank: 2.5 years, 220M calls, 40M by 1000 agents.
  - Israeli Cellular: 2.5 years, 110M calls, 25M calls by 750 agents; ongoing.

**Project Goal:** Designing and Implementing a (universal) data-base/data-repository and interface for storing, retrieving, analyzing and displaying **Call-by-Call-based Data / Information**.

#### System Components:

- ► Clean Databases: operational-data of individual calls / agents.
- ► **Graphical Online Interface**: easily generates graphs and tables, at varying resolutions (seconds, minutes, hours, days, months).

- ▶ Technion: P. Feigin, V. Trofimov, Statistics / SEE Laboratory.
- ▶ Wharton: L. Brown, N. Gans, H. Shen (UNC).
- industry:
  - ► U.S. Bank: 2.5 years, 220M calls, 40M by 1000 agents.
  - Israeli Cellular: 2.5 years, 110M calls, 25M calls by 750 agents; ongoing.

**Project Goal:** Designing and Implementing a (universal) data-base/data-repository and interface for storing, retrieving, analyzing and displaying **Call-by-Call-based Data / Information**.

#### System Components:

- ► Clean Databases: operational-data of individual calls / agents.
- Graphical Online Interface: easily generates graphs and tables, at varying resolutions (seconds, minutes, hours, days, months).

Free for academic adoption: Mini version available on a DVD; working version 7GB tables, or 20GB raw zipped, for each call center – ask for a DVD, or my mini-HD.