

Service Engineering

Class 10

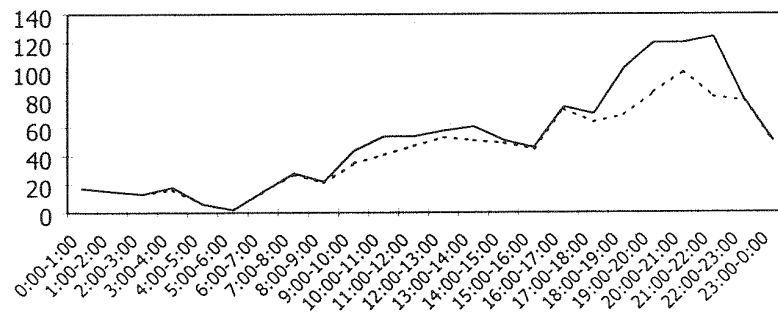
Stochastic Markovian Service Station in Steady State - Part II: The Palm/Erlang-A Queue

- Reviewing Abandonment and (Im)Patience.
- Definition of the Erlang-A Queue.
- Comparison with the Erlang-C Queue.
- Steady-State Distribution and Performance Measures.
- Probability to Abandon vs. Average Wait: $P\{\text{Ab}\} = \theta \cdot E[W_q]$.
- Estimating the (Im)Patience Parameter.
- General (Im)Patience Distribution: M/M/n+G Queue.
- Erlang-A: Fitting a Simple Model to a Complex Reality.

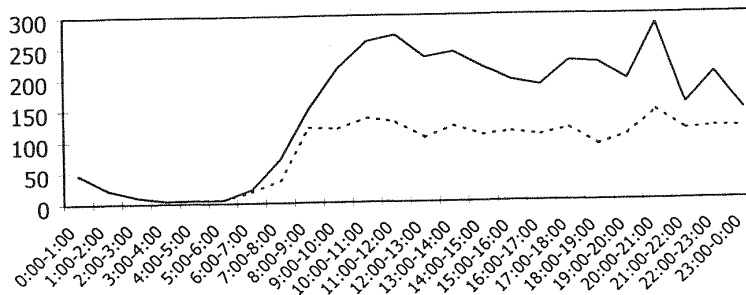
Example: How Bad Can It Get?

Call Center of a Long-Distance Service Provider. Daily Reports.

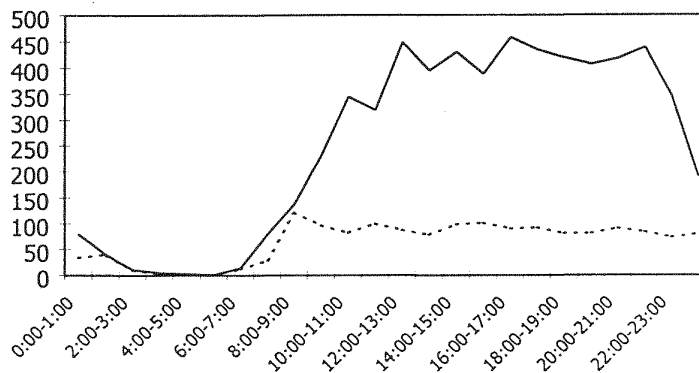
Average wait 72 sec, 81% calls answered (Saturday)



Average wait 217 sec, 53% calls answered (Thursday)



Average wait 376 sec, 24% calls answered (Sunday)



Example: How Good Can It Get?

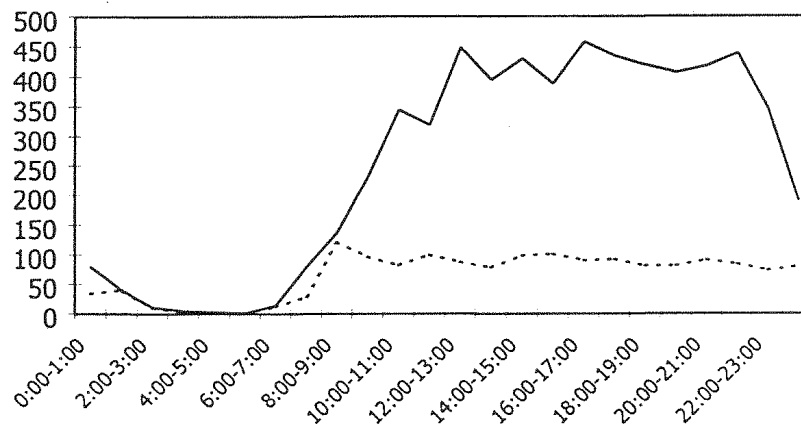
Call Center of a Health Insurance Provider. ACD Report.

Time	Calls	Answered	Abandoned%	ASA	AHT	Occ%	# of agents
Total	20,577	19,860	3.5%	30	307	95.1%	
8:00	332	308	7.2%	27	302	87.1%	59.3
8:30	653	615	5.8%	58	293	96.1%	104.1
9:00	866	796	8.1%	63	308	97.1%	140.4
9:30	1,152	1,138	1.2%	28	303	90.8%	211.1
10:00	1,330	1,286	3.3%	22	307	98.4%	223.1
10:30	1,364	1,338	1.9%	33	296	99.0%	222.5
11:00	1,380	1,280	7.2%	34	306	98.2%	222.0
11:30	1,272	1,247	2.0%	44	298	94.6%	218.0
12:00	1,179	1,177	0.2%	1	306	91.6%	218.3
12:30	1,174	1,160	1.2%	10	302	95.5%	203.8
13:00	1,018	999	1.9%	9	314	95.4%	182.9
13:30	1,061	961	9.4%	67	306	100.0%	163.4
14:00	1,173	1,082	7.8%	78	313	99.5%	188.9
14:30	1,212	1,179	2.7%	23	304	96.6%	206.1
15:00	1,137	1,122	1.3%	15	320	96.9%	205.8
15:30	1,169	1,137	2.7%	17	311	97.1%	202.2
16:00	1,107	1,059	4.3%	46	315	99.2%	187.1
16:30	914	892	2.4%	22	307	95.2%	160.0
17:00	615	615	0.0%	2	328	83.0%	135.0
17:30	420	420	0.0%	0	328	73.8%	103.5
18:00	49	49	0.0%	14	180	84.2%	5.8

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** operational performance measures (Second one is 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, unstable, sensitive)
- **Beyond Operations/OR**: Psychology, Marketing, Statistics
- **Beyond Telephony**: VRU/IVR (Opt-Out-Rates), Internet (over 60%), Hospitals ED (LWBS).

Understanding (Im)Patience

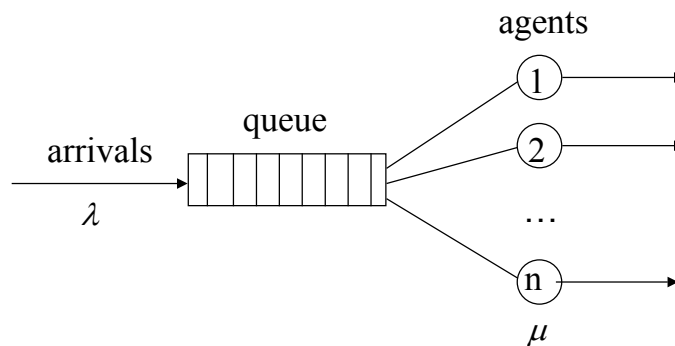
- **Observing** (Im)Patience – Heterogeneity:
Under a single roof, the fraction abandoning varies from 6% to 40%, depending on the type of service/customer.
- **Describing** (Im)Patience Dynamically:
Irritation proportional to Hazard Rate (Palm's Law).
- **Managing** (Im)Patience:
 - VIP vs. Regulars: 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.

Basic (Markovian) Queueing Models of a Basic Service Station

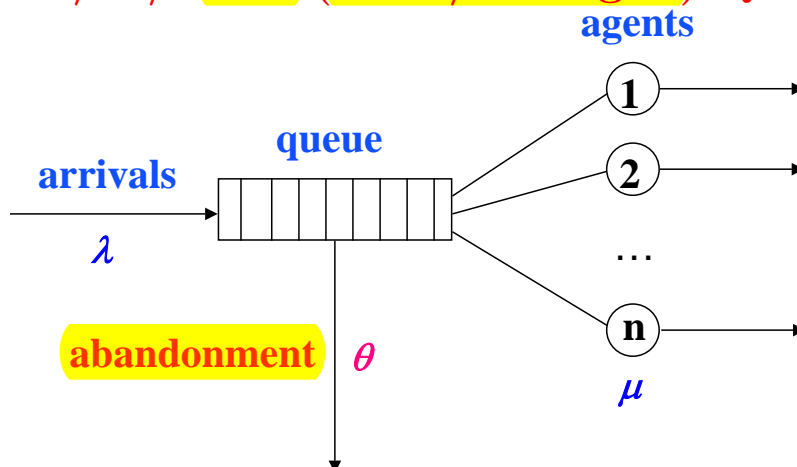
Poisson arrivals, **Exponential** service times, **Exponential** (im)patience.

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

M/M/n (Erlang-C) Queue



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



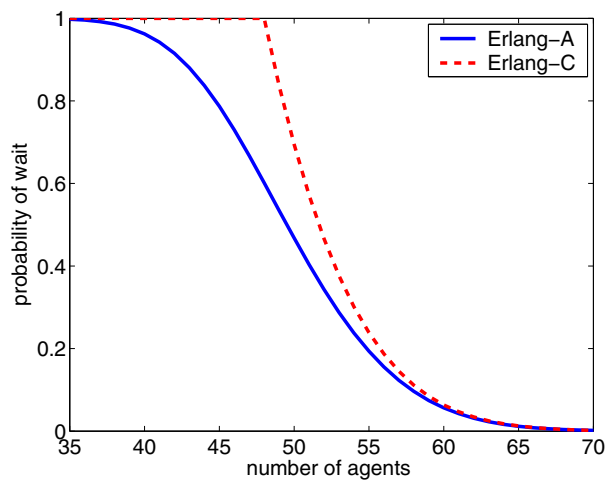
Additional Markovian Models: Balking, Trunks; Retrials.

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

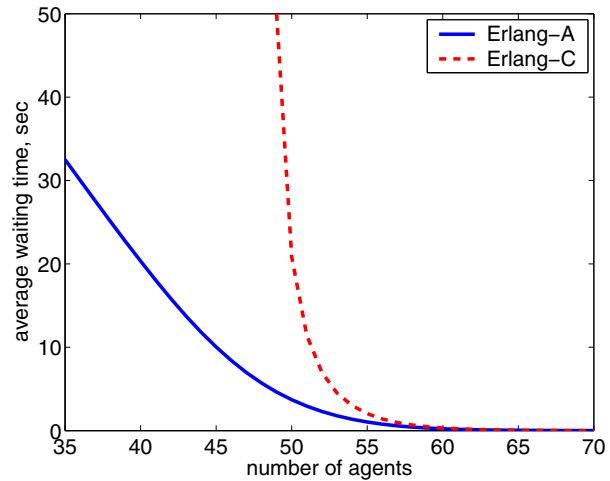
Erlang-A vs. Erlang-C

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

probability of wait
vs. number of agents



average wait
vs. number of agents



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%

“The fittest survive” and wait less - much less.

Abandonment reduces workload when needed – at high-congestion periods.

Modelling (Im)Patience: Time-to-Abandon and Offered-Wait, or Time-Willing vs. Time-Required to Wait

- **(Im)Patience time** $\tau \stackrel{d}{=} \exp(\theta)$:
time a customer is **willing to wait** for service.
- **Offered wait** V :
time a customer is **required to wait** for service; in other words, waiting time of a (virtual) customer with infinite patience.
- If $\tau \leq V$, customer **abandons**;
otherwise, **gets service**;
- **Actual wait** $W = \min\{\tau, V\}$ (sometimes W_q).

Predicting (Operational) Performance

Model **Primitives** (Building Blocks):

- Arrivals to service (eg. Poisson);
- (Im)Patience while waiting (eg. Exponential);
- Service times (eg. Exponential);
- Servers (eg. i.i.d.).

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

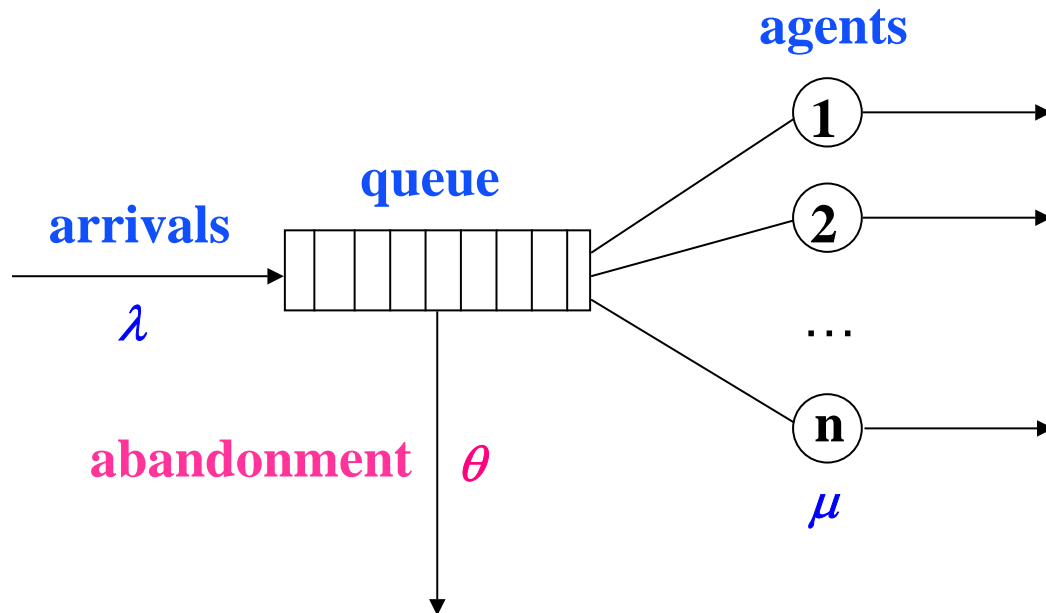
Operational Performance Measure calculable in terms of (τ, V) .

- eg. % Abandonment = $P\{\tau < V\}$ (or $P\{5 \text{ sec} < \tau < V\}$)
- eg. Average Wait = $E[\min\{\tau, V\}]$ (or $E[\tau | \tau < V]$)

Applications:

- **Performance Analysis**
- **Design, Phenomena** (Pooling, Economies of Scale)
- **Staffing – How Many Agents** (FTE's = Full-Time-Equivalent's)
Note: Within the Basic Model of heterogeneous customers and servers (vs. priorities, SBR - later).

Erlang-A (Palm, M/M/n+M; M-M/M/n)



Simplest model with abandonment, used by well-run call centers.

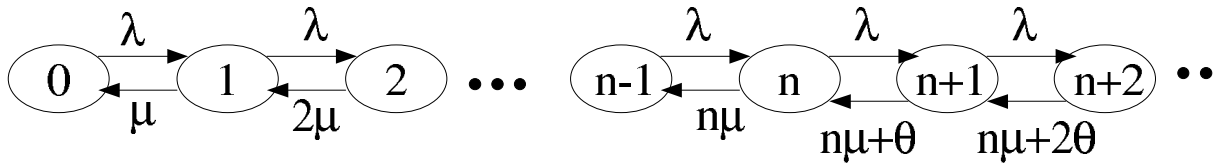
Parameters:

- λ – **Poisson** arrival rate.
- μ – **Exponential** service rate.
- n – number of service agents.
- θ – **Exponential** individual abandonment rate.

Erlang-A = Birth-and-Death Process

$L(t)$ – number-in-system at time t (served plus queued);
 $L = \{L(t), t \geq 0\}$ – Markov Birth-and-Death process.

Transition-rate diagram



Steady-state equations:

$$\begin{cases} \lambda \pi_j = (j+1) \cdot \mu \pi_{j+1}, & 0 \leq j \leq n-1 \\ \lambda \pi_j = (n\mu + (j+1-n)\theta) \cdot \pi_{j+1}, & j \geq n. \end{cases}$$

Steady-state distribution:

$$\pi_j = \begin{cases} \frac{(\lambda/\mu)^j}{j!} \pi_0, & 0 \leq j \leq n \\ \prod_{k=n+1}^j \left(\frac{\lambda}{n\mu + (k-n)\theta} \right) \frac{(\lambda/\mu)^n}{n!} \pi_0, & j \geq n+1, \end{cases}$$

where

$$\pi_0 = \left[\sum_{j=0}^n \frac{(\lambda/\mu)^j}{j!} + \sum_{j=n+1}^{\infty} \prod_{k=n+1}^j \left(\frac{\lambda}{n\mu + (k-n)\theta} \right) \frac{(\lambda/\mu)^n}{n!} \right]^{-1}.$$

Numerical drawback: infinite sums.

Erlang-A: Stability

Claim: Erlang-A is always stable.

Proof:

$$\begin{aligned}\pi_0^{-1} &= \sum_{j=0}^n \frac{(\lambda/\mu)^j}{j!} + \sum_{j=n+1}^{\infty} \prod_{k=n+1}^j \left(\frac{\lambda}{n\mu + (k-n)\theta} \right) \frac{(\lambda/\mu)^n}{n!} \\ &\leq \sum_{j=0}^{\infty} \frac{(\lambda/\min(\mu, \theta))^j}{j!} = e^{-\lambda/\min(\mu, \theta)}.\end{aligned}$$

(Used the inequality $n\mu + (k-n)\theta \geq k \min(\mu, \theta)$, for all $k \geq n$.)

Remark: Let d_j = death-rate in state j , $0 < j < \infty$.

Then, in fact,

$$j \cdot \min(\mu, \theta) \leq d_j \leq j \cdot \max(\mu, \theta).$$

Now observe that the bounds are death-rates of M/M/ ∞ queues, with service rates $\min(\mu, \theta)$ and $\max(\mu, \theta)$.

This implies that Erlang-A is sandwiched (stochastically) between two M/M/ ∞ queues.

\Rightarrow The stationary (limiting) distribution is sandwiched (stochastically) between Poisson distributions.

Special case: $\mu = \theta \Rightarrow$ Erlang-A $\stackrel{d}{=} \text{M/M}/\infty$.
 \Rightarrow Square-Root Staffing
(via Poisson \approx Normal; more on that later).

Steady-State Distribution via Special Functions (Palm)

Gamma function:

$$\Gamma(x) \triangleq \int_0^\infty t^{x-1} e^{-t} dt, \quad x > 0.$$

Incomplete Gamma function:

$$\gamma(x, y) \triangleq \int_0^y t^{x-1} e^{-t} dt, \quad x > 0, y \geq 0.$$

$$A(x, y) \triangleq \frac{x e^y}{y^x} \cdot \gamma(x, y) = 1 + \sum_{j=1}^{\infty} \frac{y^j}{\prod_{k=1}^j (x + k)}, \quad x > 0, y \geq 0.$$

Recall $E_{1,n}$ = *blocking probability* in Erlang-B (M/M/n/n):

$$E_{1,n} = \frac{\frac{(\lambda/\mu)^n}{n!}}{\sum_{j=0}^n \frac{(\lambda/\mu)^j}{j!}} = \frac{(\lambda/\mu)^n}{e^{\lambda/\mu}} \cdot \frac{1}{\Gamma(n+1) - \gamma(n+1, \lambda/\mu)}.$$

(Can be efficiently calculated via recursion.)

Then

$$\pi_j = \begin{cases} \pi_n \cdot \frac{n!}{j! \cdot \left(\frac{\lambda}{\mu}\right)^{n-j}}, & 0 \leq j \leq n, \\ \pi_n \cdot \frac{\left(\frac{\lambda}{\theta}\right)^{j-n}}{\prod_{k=1}^{j-n} \left(\frac{n\mu}{\theta} + k\right)}, & j \geq n+1, \end{cases}$$

where

$$\pi_n = \frac{E_{1,n}}{1 + \left[A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) - 1 \right] \cdot E_{1,n}}.$$

Operational Performance Measures

The most prevalent performance measure is $P\{W_q \leq T; \text{Sr}\}$ (or “worse” $P\{W_q \leq T \mid \text{Sr}\}$).

We recommend:

- $P\{W_q \leq T; \text{Sr}\}$ - fraction of **well-served**;
- $P\{\text{Ab}\}$ - fraction of **poorly-served**.

with **T** determined via “*Waiting less than T is Well-Served*”.

Or even a four-dimensional refinement:

- $P\{W_q \leq T; \text{Sr}\}$ - fraction of **well-served**;
- $P\{W_q > T; \text{Sr}\}$ - fraction of **served**, with **potential for improvement** (say, a higher priority on next visit);
- $P\{W_q > \epsilon; \text{Ab}\}$ - fraction of **poorly-served**;
- $P\{W_q \leq \epsilon; \text{Ab}\}$ - fraction of those whose **service-level** is **undetermined**.

with **ε**: “*Abandoning before ε is Harmless*”.

Properties of $P\{Ab\}$

- $P\{Ab\}$ increases monotonically in θ, λ ;
 $P\{Ab\}$ decreases monotonically in n, μ
(Bhattacharya and Ephremides, 1991);
- $P\{Ab\} \leq P\{Block\}$ in Erlang-B (Boxma and de Waal, 1994)
(think zero-patience).
- Note: In $M/M/n+G$, with $E[\tau]$ fixed, deterministic patience minimizes $P\{Ab\}$ but maximizes $E[W_q]$ (Zeltyn's PhD, 2004).

Additional Useful Performance Measures

- **ASA** (Average Speed of Answer) – used extensively in call centers; usually taken to be $E[W_q | \text{Sr}]$ (could be misleading);
- Average Wait $E[W_q]$;
- Delay Probability $P\{W_q > 0\}$ – important (later), yet unused;
- Agents' **Occupancy** $\rho = \frac{\lambda \cdot (1 - P\{\text{Ab}\})}{n\mu}$;
- Average Queue-Length $E[L_q]$.

Operational Performance Measures: Calculation via 4CallCenters

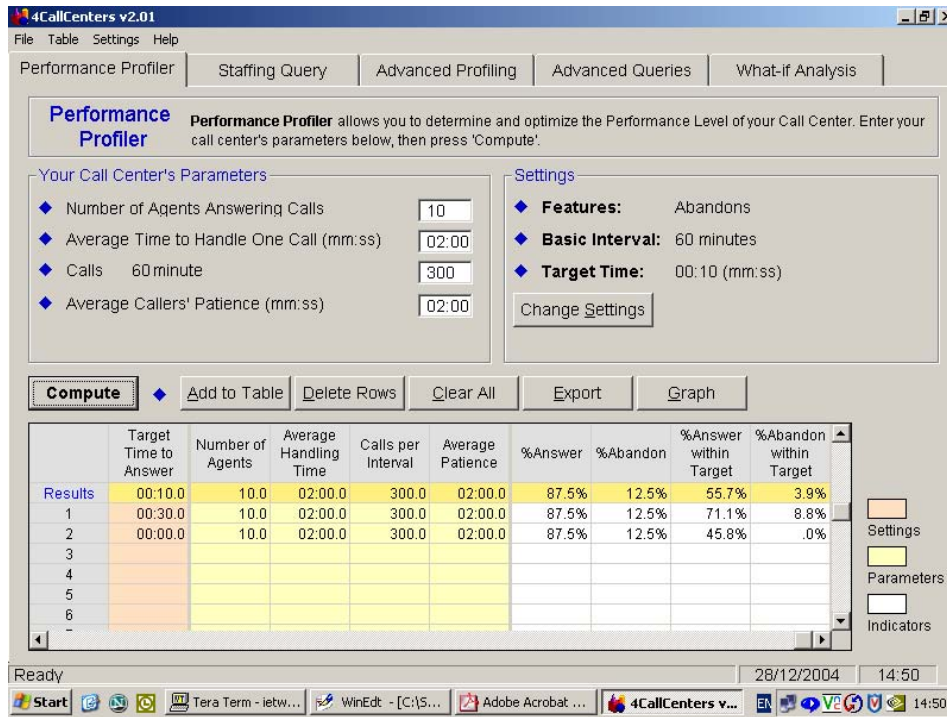
- Performance measures of the form $E[f(V, \tau)]$.
- Calculable, by numerically stable algorithms.

For example,

$f(v, \tau)$	$E[f(V, \tau)]$
$1_{\{v > \tau\}}$	$P\{V > \tau\} = P\{\text{Ab}\}$
$1_{(t, \infty)}(v \wedge \tau)$	$P\{W_q > t\}$
$1_{(t, \infty)}(v \wedge \tau) 1_{\{v > \tau\}}$	$P\{W_q > t; \text{Ab}\}$
$(v \wedge \tau) 1_{\{v > \tau\}}$	$E\{W_q; \text{Ab}\}$
$g(v \wedge \tau)$	$E[g(W_q)]$

From these, one derives additional measures, eg. $E[W_q | \text{Ab}]$.

Operational Performance Measures: Calculation via 4CallCenters



Erlang-A parameters:

$\lambda = 300$ calls/hour, $1/\mu = 2$ min, $n = 10$, $1/\theta = 2$ min.

Target times $T = 30$ sec, $\epsilon = 10$ sec.

- $P\{W_q \leq T; Sr\} = 71.1\%$;
- $P\{W_q > T; Sr\} = 87.5\% - 71.1\% = 16.4\%$;
- $P\{W_q > \epsilon; Ab\} = 12.5\% - 3.9\% = 8.6\%$;
- $P\{W_q \leq \epsilon; Ab\} = 3.9\%$.
- Delay probability $P\{W_q > 0\} = 100\% - 45.8\% = 54.2\%$.

Additional Performance Measures: Calculation via 4CallCenters

4CallCenters v2.01
File Table Settings Help

Performance Profiler | Staffing Query | Advanced Profiling | Advanced Queries | What-if Analysis

Performance Profiler Performance Profiler allows you to determine and optimize the Performance Level of your Call Center. Enter your call center's parameters below, then press 'Compute'.

Your Call Center's Parameters

- Number of Agents Answering Calls: 10
- Average Time to Handle One Call (mm:ss): 02:00
- Calls 60 minute: 300
- Average Callers' Patience (mm:ss): 02:00

Settings

- Features: Abandons
- Basic Interval: 60 minutes
- Target Time: 00:10 (mm:ss)

Change Settings

Compute | Add to Table | Delete Rows | Clear All | Export | Graph

	Agent's Occupancy	Agent's Availability	%Answer	%Abandon	Average Speed of Answer	Average Time in Queue	%Answer within Target	%Abandon within Target	Average Queue Length
Results	87.5%	12.5%	87.5%	12.5%	00:13.8	00:15.0	55.7%	3.9%	1.3
1									
2									
3									
4									
5									
6									
-									

Ready 28/12/2004 14:57

Start | Tera Te... | WinEdt... | Adobe ... | 4CallCe... | fourCC... | Docume... | 14:57

- Average Time in Queue = $E[W_q] = 15$ sec;
- ASA = $E[W_q | Sr] = 13.8$ sec;
- Agents' Occupancy $\rho = 87.5\%$;
- Average Queue Length $E[L_q] = 1.3$.

Operational Performance Measures: Calculation via Special Functions

For example,

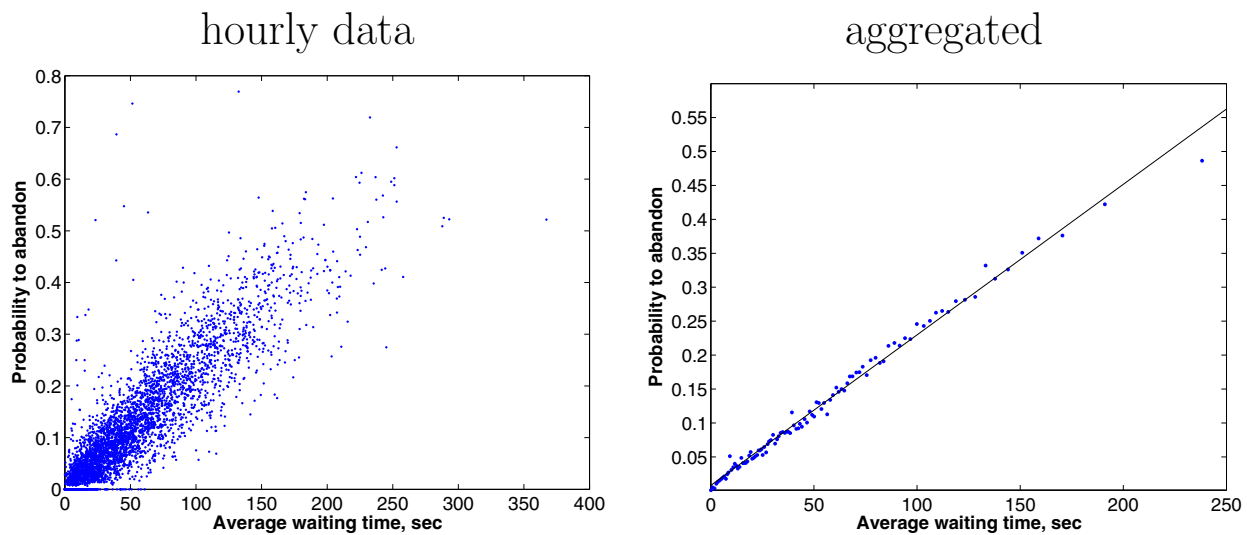
$$\begin{aligned} P\{\mathbf{W}_q > 0\} &= \sum_{j=n}^{\infty} \pi_j = \frac{A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) \cdot E_{1,n}}{1 + \left(A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) - 1\right) \cdot E_{1,n}}, \\ P[\text{Ab} | \mathbf{W}_q > 0] &= \frac{1}{\rho A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right)} + 1 - \frac{1}{\rho}, \\ E[\mathbf{W}_q | \mathbf{W}_q > 0] &= \frac{1}{\theta} \cdot \left[\frac{1}{\rho A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right)} + 1 - \frac{1}{\rho} \right]. \end{aligned}$$

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

Recall. In a queueing model with patience that is $\exp(\theta)$:

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

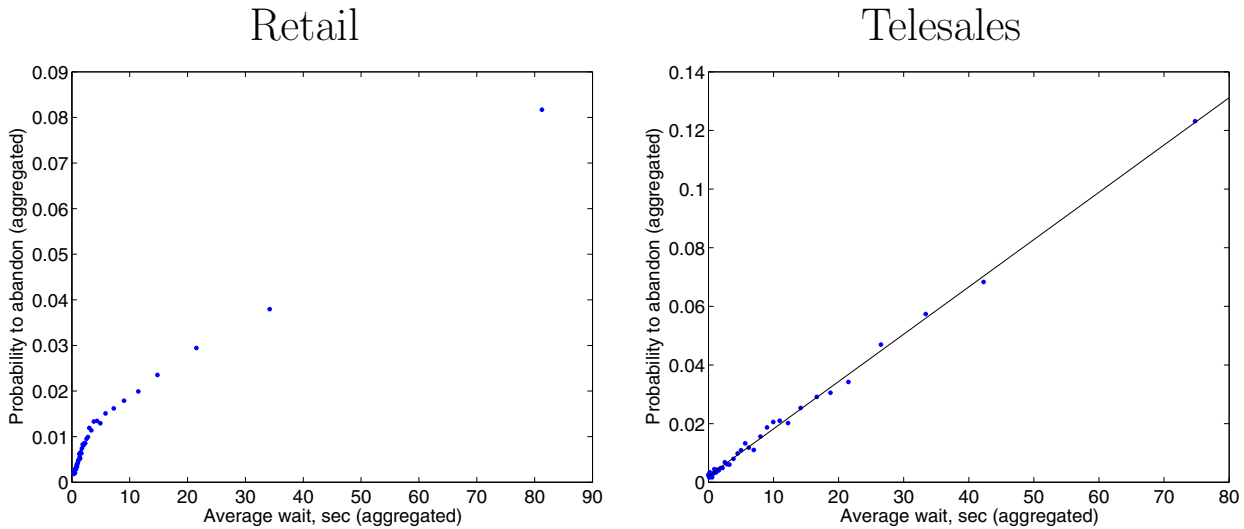
Israeli Bank: Yearly Data



The graphs are based on 4158 hour-intervals.

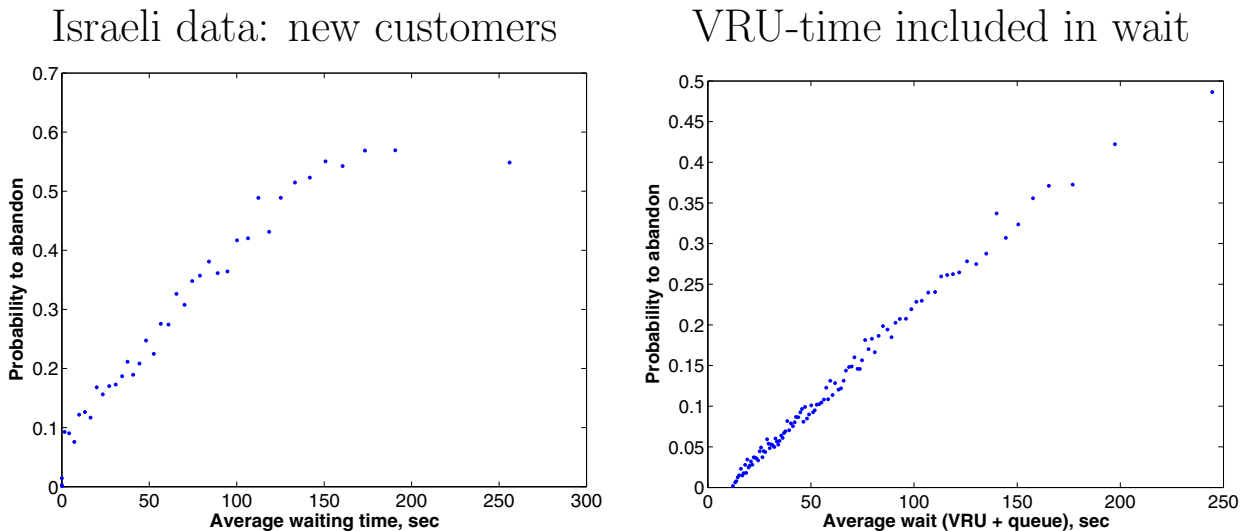
Regression \Rightarrow Average Patience $(1/\theta) \approx \frac{250}{0.56} \approx 446$ sec.

U.S. Bank



Retail – significant abandonment during first seconds of wait.

Linear patterns with non-zero intercepts



Left-hand plot \approx exp patience with **balking**:

0 with probability p , $\exp(\theta)$ with probability $(1 - p)$.

Right-hand plot \approx delayed patience: $c + \exp(\theta)$, $c > 0$.

Parameter Estimation and Prediction I; 4CallCenters, Erlang-A, and beyond

Estimation: Inference from historical data (e.g. Exp, LogNormal), with parameters assumed fixed over time-periods (overall).

Prediction: Forecast behavior beyond the available data.

Arrivals (λ)

- Poisson arrivals, time-varying but assumed with constant rate at 15/30/60 min. scale;
- Significant uncertainty concerning future rates \Rightarrow prediction;
- Helpful: Predict separately *daily volumes* and *fraction* of arrivals per time interval.

Services (μ , or $E(S)$)

- Typically stable from day to day \Rightarrow estimation;
- Can vary, depending on time-of-day;
- Typically, service time \neq talk time, and the former is needed.

First approach:

Service Time = talk time + wrap-up time (after-call work) + ...;

Second Estimation Approach:

$$E(\widehat{S}) = \frac{\text{Total Working Time} - \text{Total Accessible (Idle) Time}}{\# \text{ Served Customers}}.$$

Parameter Estimation and Prediction II

Number of Agents (n)

- Obtaining accurate historical data on n can be hard.
- Output of WFM software (given λ , μ , θ , and performance goals). One gets, in fact, the number of FTE's (Full Time Equivalent positions).
- Agents on Schedule = FTE's \times RSF (Rostered Staff Factor) (RSF > 1). Reasons: absenteeism, unscheduled breaks, ...

(Im)Patience (θ)

- Observations are **censored!** (typically heavy censoring):
 - Customer abandons \Rightarrow patience τ known;
 - Customer served \Rightarrow offered-wait V known ($\Rightarrow \tau > V$).
- Estimate via

$$\hat{\theta} = \frac{\# \text{ Abandoning}}{\text{Total Waiting Time (Abandoning + Served)}};$$

or via slope of the Regression of $P\{\text{Ab}\}$ over $E[\mathbf{W}_q]$, as before;
or both.

Estimating (Im)Patience Distribution I

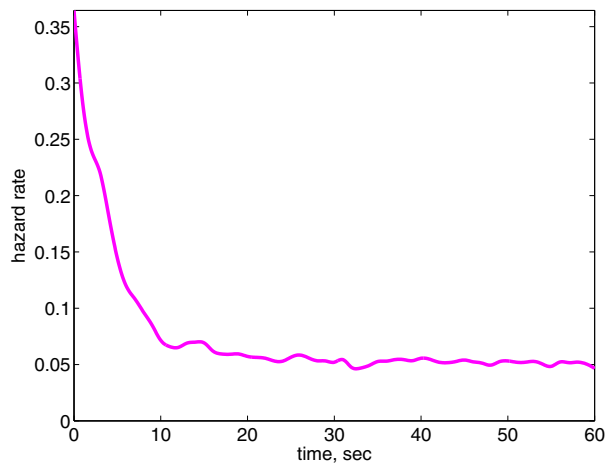
Are patience times really exponential?

To “uncensor data”, use the Kaplan-Meier estimator (standard).

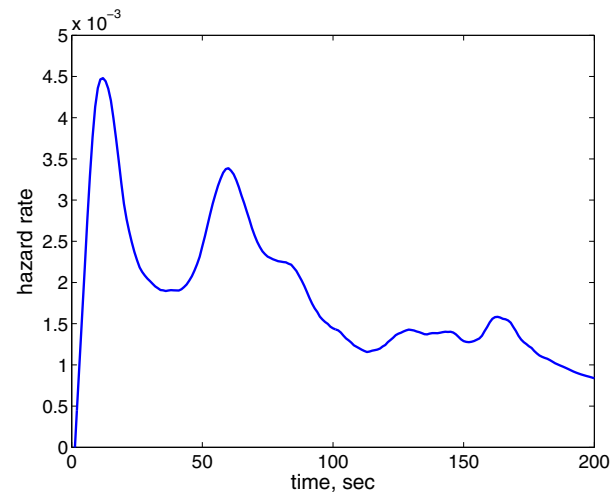
Output: Estimates of survival function and hazard-rate function.

Empirical Hazard Rates of (Im)Patience

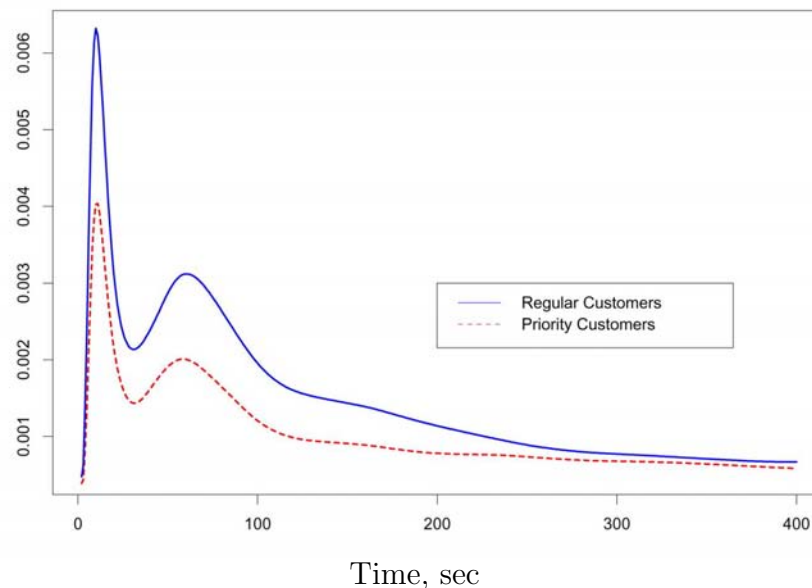
U.S. Bank



Israeli Bank

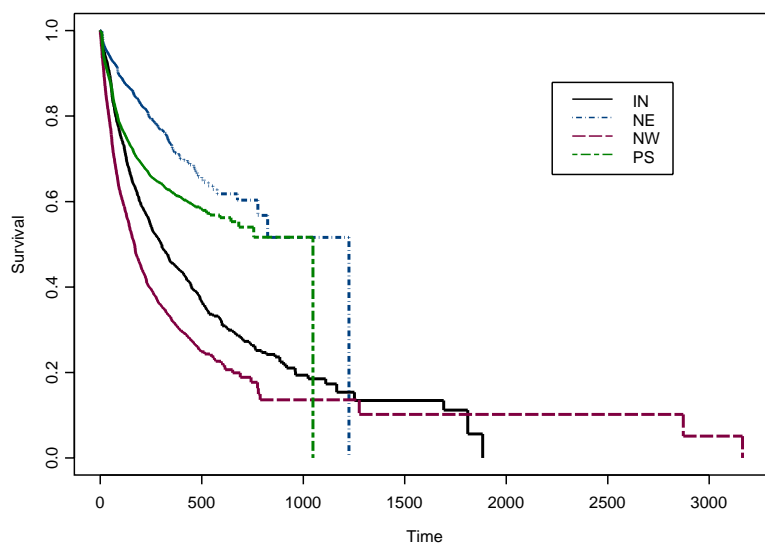


Israeli Bank: Regulars vs. VIP's



Estimating (Im)Patience Distribution II

Israeli Bank: Service Types



IN – Internet; NE – Stocks; NW – New; PS – Regulars

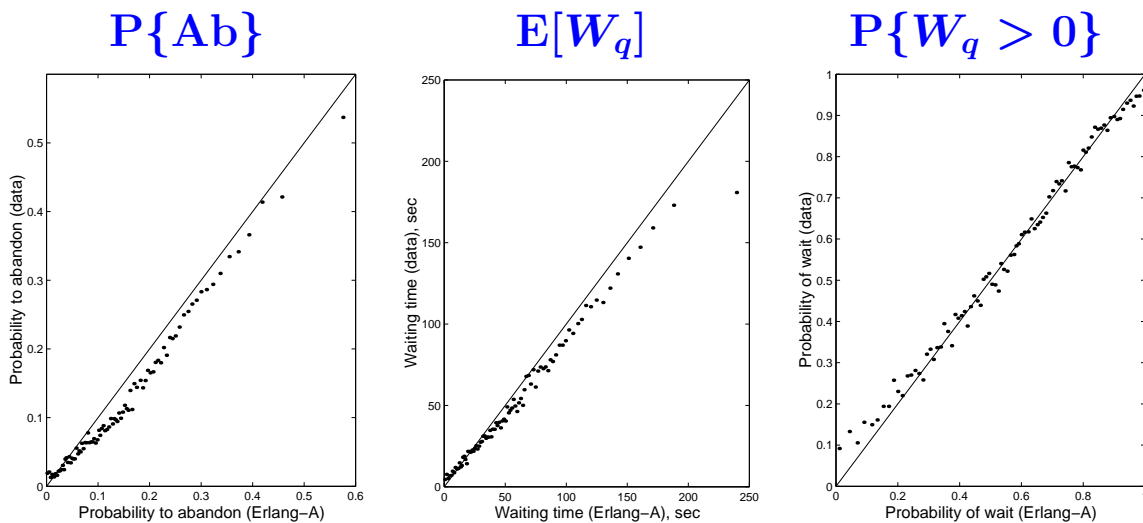
Conclusions:

- Patience time are, in general, non-exponential;
- Tele-customers are (perhaps surprisingly) **very** patient;
- Hazard-Rates very informative concerning dynamic *qualitative* evolution of (im)patience (peaks, IFR, DFR). (Palm: proportional to irritation);
- Survival functions useful for (stochastic) comparisons;
- Kaplan-Meier often problematic for estimating *quantitative* characteristics (mean, variance, median). (Eg. $E[\widehat{\tau}] = \int_0^\infty \widehat{S}(x)dx$.)

Question: Can Erlang-A be applied with non-exponential (im)patience?

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

- Small Israeli bank (10 agents);
- Data-Based Estimation of $\hat{\theta} = \frac{\# \text{ Abandoning}}{\text{Total Waiting Time}}$;
- Graph: Actual Performance vs. Erlang-A Predictions (aggregation of 40 similar hours): Model provides tight upper bounds.

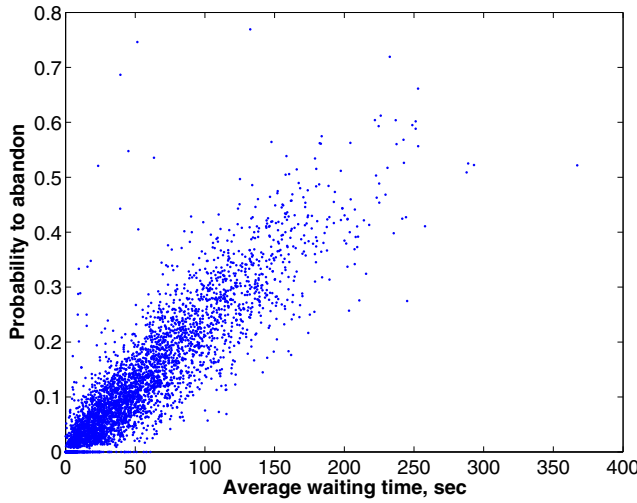


- **Question:** Why does Erlang-A works? indeed, **all** its underlying assumptions fail (Arrivals, Services, Impatience).
- **Towards a Theoretical Answer:** Robustness and Limitations, via Asymptotic (QED/QD) Analysis - later.
- **Practical Significance:** Asymptotic results applicable in small systems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

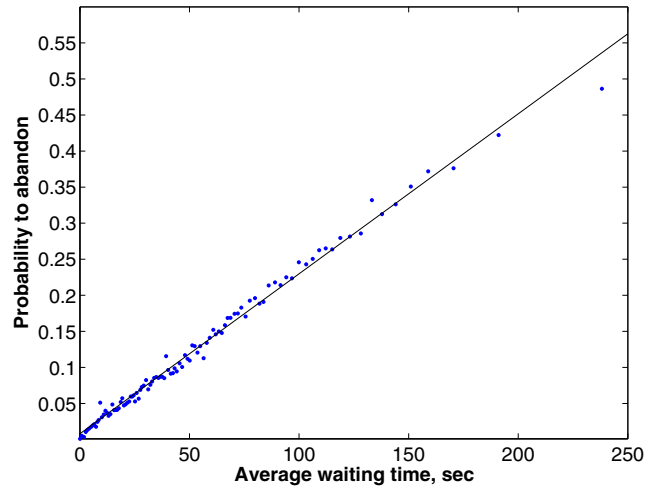
Queueing Science: In Support of Erlang-A

Israeli Bank: Yearly Data

Hourly Data



Aggregated



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

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

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

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

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

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

Queueing Science: via Data & Theory, Linearity Robust.

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

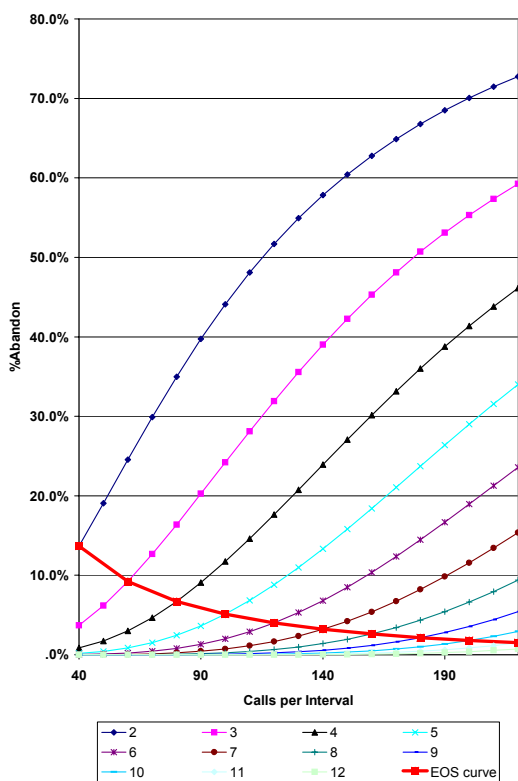
- Reality $\approx M/G/n + G \approx \text{Erlang-A}$, in which $\theta = g(0)$;
- $P\{\text{Ab}\} \approx g(0) \cdot E[W_q]$, hence **recipe prevails, often enough**.

4CallCenters: Congestion Curves

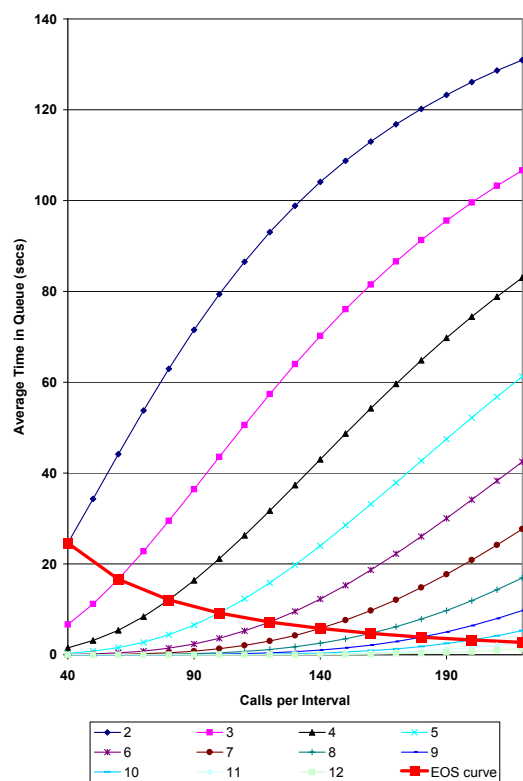
Vary input parameters of Erlang-A and display output (performance measures) in a table or graphically.

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 are the two graphs similar?

4CallCenters: Advanced Staffing Queries I

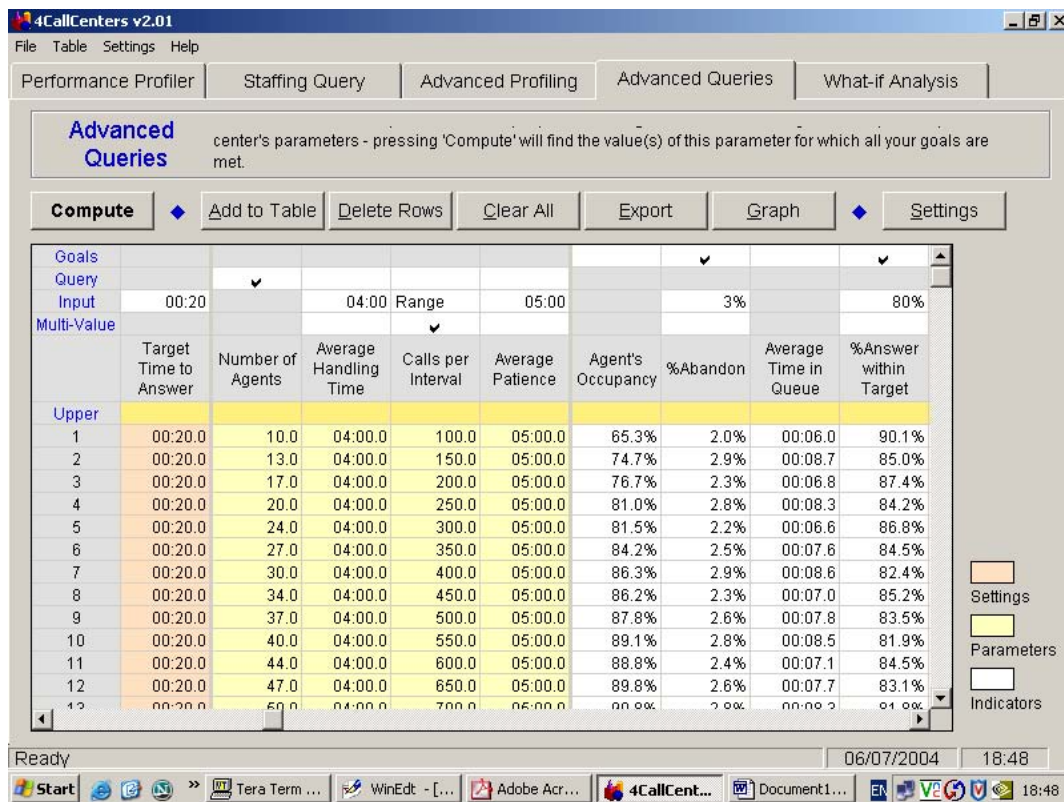
Set multiple performance goals.

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

Performance targets:

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

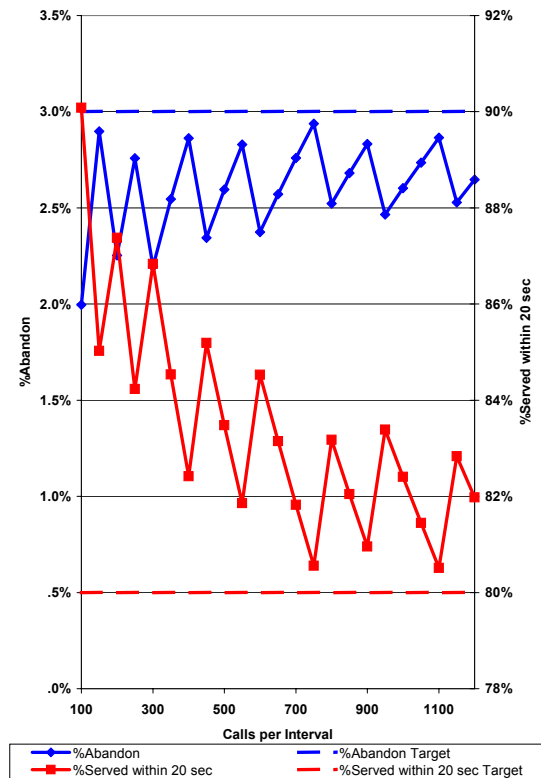
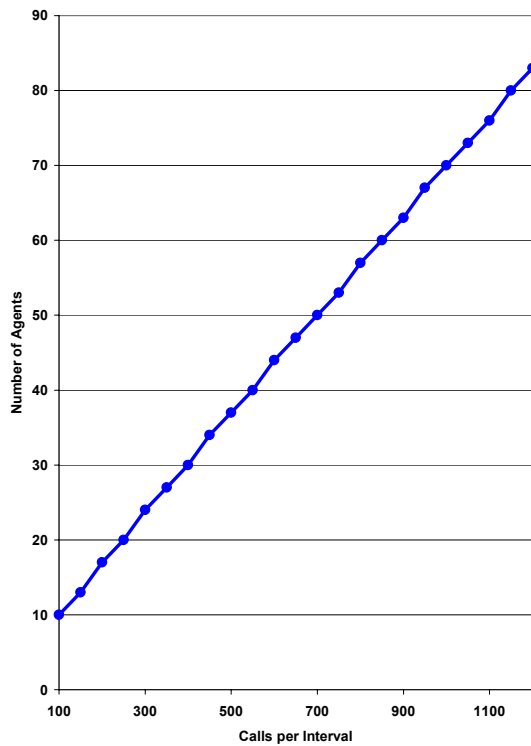
4CallCenters output



4CallCenters: 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.

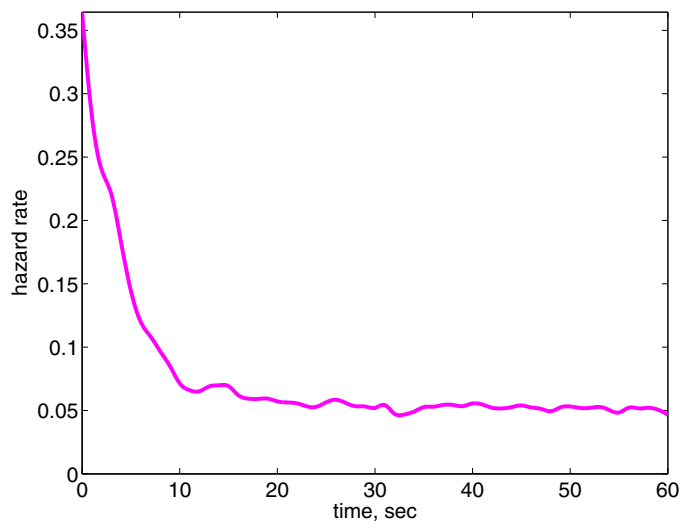
Back to General (Im)Patience: Empirical Patience Distributions

Are patience times Exponential?

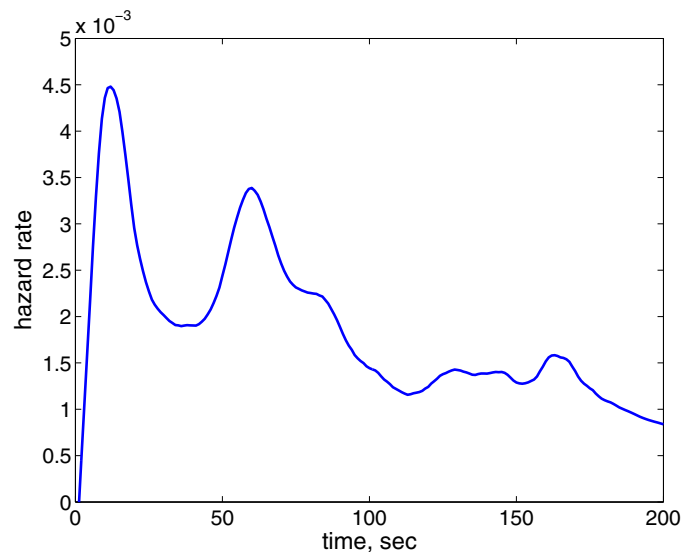
In the call centers that we studied, **they are not!**

Empirical hazard rates of patience times

U.S. bank

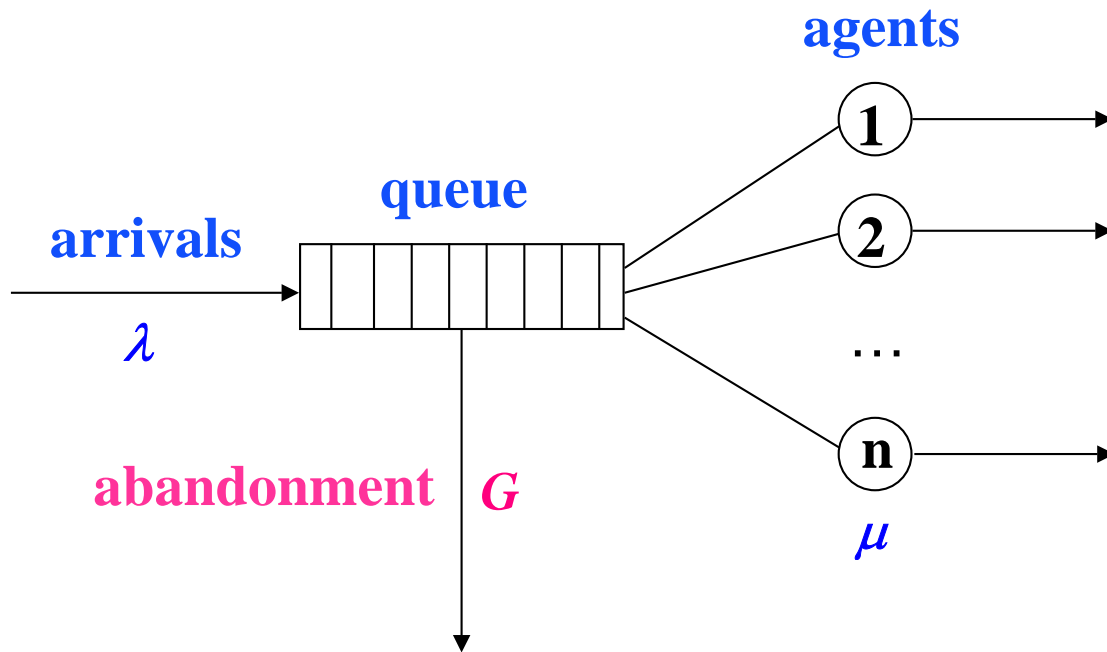


Israeli bank



To “uncensor data” use Kaplan-Meier (product-limit) estimator.
Output: estimates of **survival function** and **hazard rate**.

The M/M/n+G Queue



Patience times $\stackrel{d}{=} G(\text{eneral})$, i.i.d., independent of all else.

Performance measures can be computed, but calculations are cumbersome.

M/M/ n +G: Building Blocks, for calculating Performance Measures

Reference (Support Material in website): with Zeltyn, prepared for Bank of America.

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

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

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

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

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

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

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

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

Finally,

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

M/M/ n +G: Performance Measures

$\{\text{Ab}\} = \{\text{Abandonment}\}$, $\{\text{Sr}\} = \{\text{Served}\}$,

W – waiting time, V – offered wait,

Q – queue length.

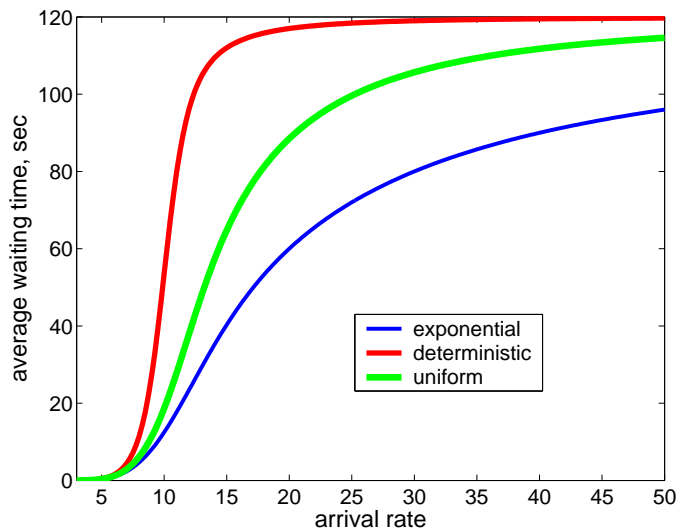
$$\begin{aligned}
 P\{V > 0\} &= \frac{\lambda J}{\mathcal{E} + \lambda J}, \\
 P\{W > 0\} &= \frac{\lambda J}{\mathcal{E} + \lambda J} \cdot \bar{G}(0), \\
 P\{\text{Ab}\} &= \frac{1 + (\lambda - n\mu)J}{\mathcal{E} + \lambda J}, \\
 P\{\text{Sr}\} &= \frac{\mathcal{E} + n\mu J - 1}{\mathcal{E} + \lambda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

M/M/n+G: Impact on Performance of Patience-Distribution

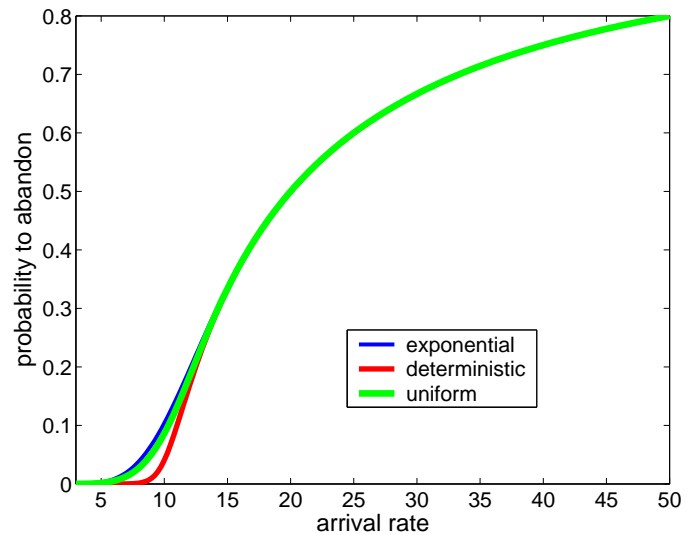
Parameters: 1 min average service time, 2 min average patience, 10 agents, arrival rate varies from 3 to 50 per minute.

G = Exponential, Deterministic, Uniform (mean = 2 min)

Average Wait

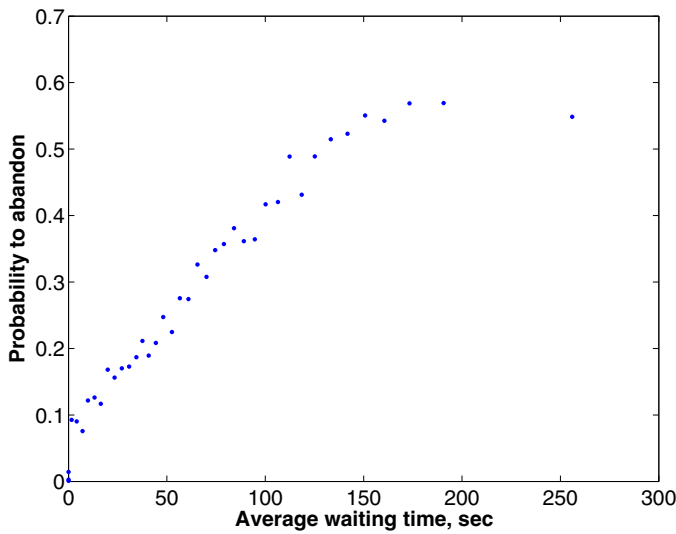


Probability to Abandon

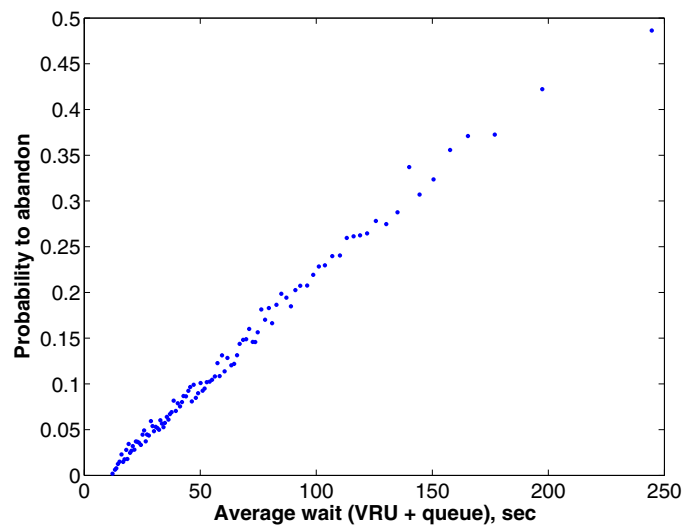


Applications of M/M/ n +G Model: Linear Patterns of $P\{Ab\}/E[W_q]$ with Non-Zero Intercepts

Israeli data: new customers



VRU-time included in wait



Left-hand plot \approx exp patience with **balking**:
0 with probability p ,
 $\exp(\theta)$ with probability $(1 - p)$.

Right-hand plot \approx **delayed patience**:
 $c + \exp(\theta)$, $c > 0$.

Simple Models at the Service of Complex Realities: **A Patience Index**

How to quantify (im)patience? Assuming experienced customer,

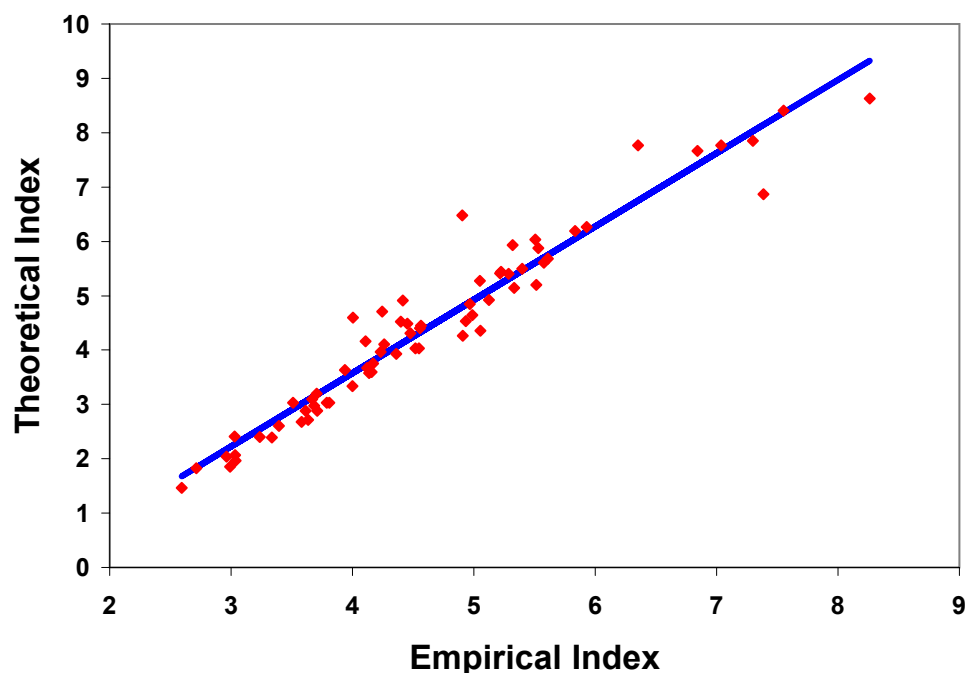
$$\begin{aligned}\text{Theoretical Patience Index} &\triangleq \frac{\text{time willing to wait}}{\text{time required to wait}} \\ &= \frac{\text{average patience}}{\text{average offered wait}}. \quad (1)\end{aligned}$$

Demanding calculations. Hence, “assume” τ and V Exponential:

$$\text{Empirical Patience Index} \triangleq \frac{\% \text{ served}}{\% \text{ abandoned}}.$$

Easily calculable from ACD reports.

Patience Index – Empirical vs. Theoretical



PATIENCE INDEX

- How to Define? Measure? Manage?

<u>Statistics</u>	<u>Time Till</u>	<u>Interpretation</u>
360K served (80%)	2 min.	? must = expect
90K abandon (20%)	1 min.	? willing to wait

“Time willing to wait” of served is **censored** by their “wait”.

“Uncensoring” (simplified)

Willing to wait $1 + 2 \times \frac{360K}{90K} = 1 + 2 \times 4 = \mathbf{9}$ min.

Expect to wait $2 + 1 \times \frac{90K}{360K} = 2 + 1 \times \frac{1}{4} = \mathbf{2.25}$ min.

Patience Index = $\frac{\text{time willing}}{\text{time expect}} = 4 = \frac{\# \text{ served/wait} > 0}{\# \text{ abandon/wait} > 0}$

\uparrow definition \uparrow measure

Customer-Focused Queueing Theory

Waiting experience of experienced customer often cycles through:

1. Time that a customer *expects* to wait;
2. Time that a customer is *willing* to wait (τ , patience or need);
3. Time that a customer *required* wait (V , offered wait);
4. Time that a customer *actually* waits ($W_q = \min(\tau, V)$);
5. Time that a customer *perceives* waiting.

Experienced customers $\Rightarrow 1=3$.

Rational customers $\Rightarrow 4=5$.

Thus left with (τ, V) , as in Erlang-A.

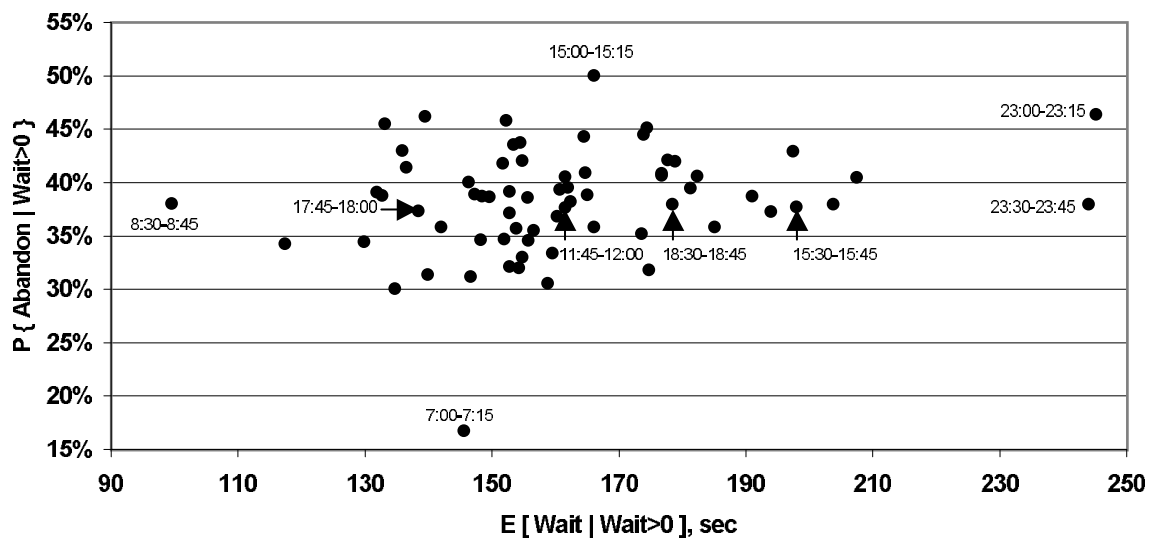
Eg. 200 abandonment in Direct-Banking: Perceived vs. Actual.

Reason to Abandon	Actual Abandon Time (sec)	Perceived Abandon Time (sec)	Perception Ratio
Fed up waiting (77%)	70	164	2.34
Not urgent (10%)	81	128	1.6
Forced to (4%)	31	35	1.1
Something came up (6%)	56	53	0.95
Expected call-back (3%)	13	25	1.9

Adaptive Behavior of (Im)patient Customers

Question: Do customers adapt their patience to system performance (offered wait)?

Israeli Bank: Internet-Support Customers



Supporting theory in “Rational abandonment from invisible queues”, with Shimkin & Zohar.

Queues with Impatient Customers

M/M/N+M (G)

Palm/Erlang-A

e.mail: avim@tx.technion.ac.il

Website: <http://ie.technion.ac.il/serveng>

Supporting Material (in Website)

Gans, Koole, and M.: "Telephone Call Centers: **Tutorial**, **Review** and Research Prospects." *Review of State-of-the-Art Research*.

Brown, Gans, M., Sakov, Shen, Zeltyn, Zhao: "**Statistical** Analysis of a Telephone Call Center: A Queueing-Science Perspective." *Analysis of Arrivals, Services and Patience*.

Garnett, M. and Reiman: "**Designing** a Call Center with Impatient Customers." *Erlang-A, based on Garnett's MSc thesis*.

M. and Zeltyn: "The Impact of Customer Patience on Delay and Abandonment: Some **Empirically-Driven Experiments** with the M/M/N+G Queue." *On the relation between $P(Ab)$ and $E(Wait)$* .

Zeltyn: Ph.D. thesis, on M/M/N+G.

Palm: "Intensitatsschwankungen im fernsprechverkehr," (In English) Ericsson Technics, 1943.

Palm: "Methods of judging the annoyance caused by congestion." Tele, 1953: **Recommended**.

Bacelli and Hebuterne: "On queues with impatient customers." In Performance '81, ed. Gelenbe, 1981.

The Palm/Erlang-A Queue, with Applications to Call Centers*

Avishai Mandelbaum and Sergey Zeltyn

Faculty of Industrial Engineering & Management
Technion,
Haifa 32000, ISRAEL

emails: avim@tx.technion.ac.il, zeltyn@ie.technion.ac.il

December 28, 2004

Contents

1	Introduction	1
2	Significance of abandonment in practice and modelling	3
3	Birth-and-death process representation; Steady-state	6
4	Operational measures of performance	9
4.1	Practical measures: Waiting Time	9
4.2	Practical measures: accounting for Abandonment	10
4.3	Calculations: the 4CallCenters software	12
4.4	Delay probability $P\{W>0\}$	13
4.5	Fraction abandoning $P\{Ab\}$	13
4.6	Theoretical relations among $P\{Ab\}$, $E(W)$, $E(Q)$	14
4.7	A general approach for computing operational performance measures	15
4.8	Empirical relations between $E(W)$ and $P\{Ab\}$	15
5	Parameter estimation and prediction in a call center environment	17
6	Approximations	19

*Parts of the text are adapted from [11], [19], [22], [28] and [40]

7	Applications in call centers	23
7.1	Erlang-A performance measures: comparison against real data	23
7.2	Erlang-A approximations: comparison against real data	24
8	Human behavior	25
8.1	Balking and delayed impatience	25
8.2	Examples of the patience-time hazard rate	26
8.3	Adaptive behavior of impatient customers	27
8.4	Patience index	29
9	Advanced features of the 4CallCenters software	30
10	Some open research topics	33
10.1	Dimensioning the Erlang-A queue	33
10.2	Uncertainty in parameter values	34
10.3	Additional topics	35
A	Derivation of some Erlang-A performance measures	39

4 Operational measures of performance

In order to understand and apply the Erlang-A model, one must first define its measures of performance, and then be able to calculate them. Moreover, since call centers can get very large (thousands of agents), the implementation of these calculations must be both fast and numerically stable.

4.1 Practical measures: Waiting Time

The most popular measure of operational (positive) performance is the fraction of served customers that have been waiting less than some given time, or formally $P\{W \leq T, \text{Sr}\}$, where W is the (random) waiting time in steady-state, $\{\text{Sr}\}$ is the event “customer gets service” and T is a target time that is determined by Management/Marketing. For example, in a call center that caters to emergency calls, $T = 0$ (or T very small) would be appropriate. A common rule of thumb (without any theoretical backing, as far as we know) is the goal that at least 80% of the customers be served within 20 seconds; formally, $P\{W \leq 20, \text{Sr}\} \geq 0.8$. To this, one sometimes adds $E[W]$, or $E[W|W > 0]$, as some measure of an average (negative) experience for those who waited.

An important measure that is rarely used in practice is $P\{W > 0\}$, the fraction of customers who encounter a delay. This is a useful stable measure of congestion. Its importance stems from the fact that it identifies an organization's operational focus, in the following sense:

- $P\{W > 0\}$ close to 0 indicates a Quality-Driven operation, where the focus is on *service quality*;
- $P\{W > 0\}$ close to 1 indicates an Efficiency-Driven operation, where the focus is on *servers' efficiency* (in the sense of high servers' utilization);
- $P\{W > 0\}$ strictly between 0 and 1 (for example 0.5) indicates a careful *balance* between Quality and Efficiency, which we abbreviate to **QED = Quality & Efficiency Driven** operational regime.

The above three-regime dichotomy is rather delicate. For example, consider a system in which customers' average patience is close to the average service duration (for example, let both be equal to one minute), and assume that its offered load λ/μ is 100 Erlangs. Then, staffing of 100 servers would lead to the QED regime, with high levels of both service and efficiency that are balanced as follows: about 50% of the customers are served immediately upon arrival, the average wait is 2.3 seconds, 4% of the customers abandon due to their impatience, and servers' utilization levels are 96%. The QED regime still prevails at staffing levels between 95 and 105. With 90 servers, the system is efficiency-driven: 11% of the customers abandon, only 15% are served immediately, and utilization is over 99%. With 110 agents, it is quality-driven: abandonment is less than 1%, and 83% are served immediately.

In Section 6, we shall add details about the three operational regimes. This will be done in the context of describing regime-specific approximations for performance measures. However, there is much more to say about this important subject, and readers are referred to [19, 9] and Section 4 in the review [17] for details.

4.2 Practical measures: accounting for Abandonment

In a quality-driven service, $P\{W > 0\}$ seems the “right” measure of operational performance. We thus turn to alternative modes of operations and consider hereafter services in which $P\{W > 0\}$ is not close to vanishing.

As explained before, performance measures must take into account those customers who abandon. Indeed, if forced into choosing a *single* number as a proxy for operational performance, we recommend the probability to abandon $\mathbf{P}\{\mathbf{Ab}\}$, the fraction of customers who explicitly declare that the service offered is not worth its wait. Some managers actually opt for a refinement that excludes those who abandon within a very short time, formally $P\{W > \epsilon; \mathbf{Ab}\}$, for some small $\epsilon > 0$, for example $\epsilon = 3$ seconds. The justification is that those who abandon within 3 seconds can not be characterized as poorly served. There is also a practical rational that arises from physical limitations, specifically that such “immediate” abandonment could in fact be a malfunction or an inaccuracy of the measurement devices.

The single abandonment measure $P\{\mathbf{Ab}\}$ can be in fact refined to account explicitly for those customers who were or were not well-served. Thus, we propose:

- $P\{W \leq T; \mathbf{Sr}\}$ - fraction of well-served;
- $P\{\mathbf{Ab}\}$ - fraction of poorly-served.

A further refinement, that yields a four-dimensional service measure, could be:

- $P\{W \leq T; \mathbf{Sr}\}$ - fraction of well-served;
- $P\{W > T; \mathbf{Sr}\}$ - fraction of served, with a potential for improvement (say, a higher priority on their next visit);
- $P\{W > \epsilon; \mathbf{Ab}\}$ - fraction of poorly-served;
- $P\{W \leq \epsilon; \mathbf{Ab}\}$ - fraction of those whose service-level is undetermined - see the above for an elaboration.

Remark 4.1 4CallCenters [16] calculates, for a given target time, both $P\{W \leq T; \mathbf{Sr}\}$, the fraction of customers who are served within target, and $P\{W \leq \epsilon; \mathbf{Ab}\}$, those who abandon within target. To calculate the other two measures, it suffices to have $P\{\mathbf{Ab}\}$, also calculated by 4CallCenters. Indeed,

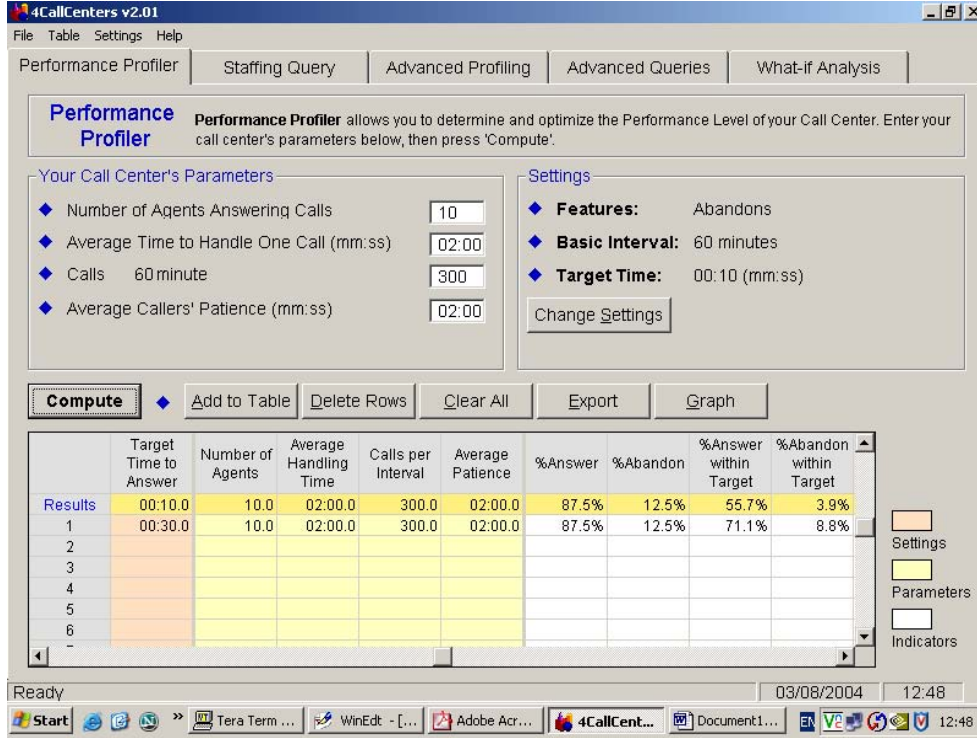
$$\begin{aligned} P\{W > T; \mathbf{Sr}\} &= 1 - P\{\mathbf{Ab}\} - P\{W \leq T; \mathbf{Sr}\}, \\ P\{W > \epsilon; \mathbf{Ab}\} &= P\{\mathbf{Ab}\} - P\{W \leq \epsilon; \mathbf{Ab}\}. \end{aligned}$$

Since a single target must be used ($T = \epsilon$ above), one must apply the program twice if different targets are required.

4.3 Calculations: the 4CallCenters software

Black-box Erlang-A calculations, as well as many other useful features, are provided by the free-to-use software 4CallCenters [16]. (This software is being regularly debugged and upgraded.) The calculation methods are described in Appendix B of [19]; they were developed in the Technician's M.Sc. thesis of the first author, Ofer Garnett.

Figure 5: 4Callcenters. Example of output.



These calculations are in fact for measures of the form $E[f(V, \tau)]$, for various functions f (Table 3 in [19]). For example,

$$E[W] = E[\min\{V, \tau\}] , \quad P\{\text{Abandon}\} = E[1_{\{\tau < V\}}] .$$

Figure 5 displays a 4CallCenters output and demonstrates how to calculate the four-dimensional service measure, introduced in Subsection 4.2.

The values of the four Erlang-A parameters are displayed in the middle of the upper half of the screen: $n = 10$, $1/\mu = 2$ minutes, $\lambda = 300$ calls per hour, $1/\theta = 2$ minutes. Let $T = 30$

seconds and $\epsilon = 10$ seconds. Then one should perform computations twice: with *Target Time* 30 and 10 seconds. (Both computations appear in Figure 5.) We get:

- $P\{W \leq T; \text{Sr}\}$ - fraction of well-served is equal to 71.1%;
- $P\{W > T; \text{Sr}\}$ - fraction of served, with a potential for improvement, is 16.4% (87.5% – 71.1%);
- $P\{W > \epsilon; \text{Ab}\}$ - fraction of poorly-served is 8.6% (12.5% – 3.9%);
- $P\{W \leq \epsilon; \text{Ab}\}$ - fraction of those whose service-level is undetermined is 3.9%.

Note that the 4CallCenters output includes many more performance measures than those displayed in Figure 5: one could scroll the screen to values of agents' occupancy, average waiting time, average queue length, etc.

In Section 9 we describe several examples of the more advanced capabilities of 4CallCenters.

4.4 Delay probability $P\{W>0\}$

In this note, we content ourselves with few representative insightful calculations, based on conditioning and the incomplete gamma function introduced above. We start with the *delay probability* $P\{W > 0\}$, which represents the fraction of customers who are forced to actually wait for service. (The others are served immediately upon calling.) Recall that this measure identifies operational regimes of performance.

Following Palm [31], we show in the Appendix that the representations (3.5) and (3.7) immediately imply

$$P\{W > 0\} = \sum_{j=n}^{\infty} \pi_j = \frac{A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) \cdot E_{1,n}}{1 + \left(A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) - 1\right) \cdot E_{1,n}}; \quad (4.1)$$

here, the first equality in (4.1) follows from PASTA.

4.5 Fraction abandoning $P\{\text{Ab}\}$

We proceed with calculating the probability to abandon, which represents the fraction abandoning. Define $P_j\{\text{Sr}\}$ to be the probability of ultimately getting served, for a customer that encounters all servers busy and j customers in queue, upon arrival (equivalently, $n + j$ in the system). "Competition among exponentials" now implies that

$$P_0\{\text{Sr}\} = \frac{n\mu}{n\mu + \theta}.$$

Then,

$$P_1\{\text{Sr}\} = \frac{n\mu + \theta}{n\mu + 2\theta} \cdot P_0\{\text{Sr}\} = \frac{n\mu}{n\mu + 2\theta},$$

where we conditioned on the first event, after an arrival that encounters all servers busy and a single customer in queue; this event is either a service completion (with probability $\frac{n\mu + \theta}{n\mu + 2\theta}$) or an abandonment. More generally, via induction:

$$P_j\{\text{Sr}\} = \frac{n\mu + j\theta}{n\mu + (j+1)\theta} \cdot P_{j-1}\{\text{Sr}\} = \frac{n\mu}{n\mu + (j+1)\theta}, \quad j \geq 1.$$

The probability to abandon service, given all servers busy and j customers in the queue upon arrival, finally equals

$$P_j\{\text{Ab}\} = 1 - P_j\{\text{Sr}\} = \frac{(j+1)\theta}{n\mu + (j+1)\theta}, \quad j \geq 0. \quad (4.2)$$

It follows that

$$P[\text{Ab}|W > 0] = \sum_{j=n}^{\infty} \pi_j P_j\{\text{Ab}\} / P\{W > 0\} = \frac{1}{\rho A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right)} + 1 - \frac{1}{\rho}. \quad (4.3)$$

The first equality in (4.3) is a consequence of PASTA, and the second is derived in the Appendix. The fraction abandoning, $P\{\text{Ab}\}$, is simply the product $P[\text{Ab}|W > 0] \times P\{W > 0\}$.

4.6 Theoretical relations among $P\{\text{Ab}\}$, $E(W)$, $E(Q)$

A remarkable property of Erlang-A, which in fact generalizes to other models with patience that is $\exp(\theta)$, is the following linear relation between the fraction abandoning $P\{\text{Ab}\}$ and average wait $E[W]$:

$$P\{\text{Ab}\} = \theta \cdot E[W]. \quad (4.4)$$

Proof: The proof is based on the balance equation

$$\theta \cdot E[Q] = \lambda \cdot P\{\text{Ab}\}, \quad (4.5)$$

and on Little's formula

$$E[Q] = \lambda \cdot E[W], \quad (4.6)$$

where Q is the steady-state queue length. The balance equation (4.5) is a steady-state equality between the rate that customers abandon the queue (left hand side) and the rate that abandoning customers (i.e. - customers who eventually abandon) enter the system. Substituting Little's formula (4.6) into (4.5) yields formula (4.4). ■

Observe that (4.4) is equivalent to

$$P[\text{Ab}|W > 0] = \theta \cdot E[W|W > 0]. \quad (4.7)$$

Then, the average waiting time of delayed customers is computed via (4.3) and (4.7):

$$E[W|W > 0] = \frac{1}{\theta} \cdot \left[\frac{1}{\rho A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right)} + 1 - \frac{1}{\rho} \right]. \quad (4.8)$$

The unconditional average wait $E[W]$ equals the product of (4.1) with (4.8).

4.7 A general approach for computing operational performance measures

Expressions for additional performance measures of Erlang-A are derived in Riordan [32]. However, we recommend to use more general M/M/n+G formulae, as the main alternative to the 4CallCenters software. Indeed, M/M/n+G is a generalization of Erlang-A, in which patience times are generally distributed. A comprehensive list of M/M/n+G formulae, as well as guidance for their application, appears in Mandelbaum and Zeltyn [30]. The preparation of [30] was triggered by a request from a large U.S. bank. Consequently, this bank has been routinely applying Erlang-A in the workforce management of its 10,000 telephone agents, who handle close to 150 millions calls yearly. (In fact, Erlang-A replaced a simulation tool that had been used before.)

The handout [30] also explains how to adapt the M/M/n+G formulae to Erlang-A, in which patience is exponentially distributed:

$$G(x) = 1 - e^{-\theta x}, \quad \theta > 0.$$

Specifically, see Sections 1,2 and 5 of [30].

Finally, we explain how to calculate the four service measures from Section 4.2. The list on page 4 of [30] contains formulae for $P\{\text{Ab}\}$, $P\{W > T\}$ and $P\{\text{Ab}|W > T\}$. The product of the last two provides us with $P\{W > T; \text{Ab}\}$. The other three service measures are easily derived. For example,

$$P\{W > T; \text{Sr}\} = P\{W > T\} - P\{W > T; \text{Ab}\}.$$

4.8 Empirical relations between $E(W)$ and $P\{\text{Ab}\}$

Figure 6 illustrates the relation (4.4). It was plotted using yearly data of an Israeli bank call center [12]. (See also Brown et al. [11] for statistical analysis of this call center data.) First,

Appendix

A Derivation of some Erlang-A performance measures

Steady-state distribution. Using formulae (3.4), (3.5) and definition (3.6) one gets

$$\begin{aligned}\pi_0^{-1} &= \sum_{j=0}^n \frac{(\lambda/\mu)^j}{j!} + \frac{(\lambda/\mu)^n}{n!} \cdot \sum_{j=n+1}^{\infty} \prod_{k=n+1}^j \left(\frac{\lambda}{n\mu + (k-n)\theta} \right) \\ &= \frac{(\lambda/\mu)^n}{n!} \cdot \left[\frac{1}{E_{1,n}} + \sum_{j=1}^{\infty} \frac{(\lambda/\theta)^j}{\prod_{k=1}^j (n\mu/\theta + k)} \right] = \frac{(\lambda/\mu)^n}{n!} \cdot \left[\frac{1}{E_{1,n}} + A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) - 1 \right].\end{aligned}$$

Hence

$$\pi_0 = \frac{E_{1,n}}{1 + \left[A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) - 1 \right] \cdot E_{1,n}} \cdot \frac{n!}{(\lambda/\mu)^n}.$$

For $1 \leq j \leq n$

$$\pi_j = \pi_0 \cdot \frac{(\lambda/\mu)^j}{j!} = \frac{E_{1,n}}{1 + \left[A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) - 1 \right] \cdot E_{1,n}} \cdot \frac{n!}{j! \cdot (\lambda/\mu)^{n-j}}.$$

Specifically,

$$\pi_n = \frac{E_{1,n}}{1 + \left[A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) - 1 \right] \cdot E_{1,n}}. \quad (\text{A.1})$$

Finally, for $j > n$,

$$\pi_j = \pi_n \cdot \frac{\lambda^{j-n}}{\prod_{k=1}^{j-n} (n\mu + k\theta)} = \frac{E_{1,n}}{1 + \left(A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) - 1 \right) \cdot E_{1,n}} \cdot \frac{\left(\frac{\lambda}{\theta}\right)^{j-n}}{\prod_{k=1}^{j-n} \left(\frac{n\mu}{\theta} + k\right)}. \quad (\text{A.2})$$

Probability of wait. From PASTA, (A.1) and (A.2), the delay probability is equal to

$$\begin{aligned} \mathbb{P}\{W > 0\} &= \sum_{j=n}^{\infty} \pi_j = \frac{E_{1,n}}{1 + \left[A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) - 1 \right] \cdot E_{1,n}} \cdot \left[1 + \sum_{j=n+1}^{\infty} \frac{(\lambda/\theta)^{j-n}}{\prod_{k=1}^{j-n} \left(\frac{n\mu}{\theta} + k\right)} \right] \\ &= \frac{A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) \cdot E_{1,n}}{1 + \left[A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) - 1 \right] \cdot E_{1,n}}. \end{aligned} \quad (\text{A.3})$$

Probability to abandon. First, we need to perform some preliminary calculations. Differentiating (3.5), we get

$$\frac{\partial}{\partial y} A(x, y) = \frac{\partial}{\partial y} \left[\frac{x e^y}{y^x} \gamma(x, y) \right] = \frac{x}{y} + \left(1 - \frac{x}{y} \right) \cdot A(x, y).$$

Then, for $x > 0$, $y > 0$,

$$\begin{aligned} \sum_{j=0}^{\infty} \frac{(j+1)y^j}{\prod_{k=1}^{j+1} (x+k)} &= \frac{\partial}{\partial y} \left[\sum_{j=1}^{\infty} \frac{y^j}{\prod_{k=1}^j (x+k)} \right] \\ &= \frac{\partial}{\partial y} [A(x, y) - 1] = \frac{\partial}{\partial y} A(x, y) = \frac{x}{y} + \left(1 - \frac{x}{y} \right) \cdot A(x, y). \end{aligned} \quad (\text{A.4})$$

Using (A.3) and (4.2), the conditional probability to abandon is equal to

$$\begin{aligned} \mathbb{P}\{\text{Ab}|W > 0\} &= \frac{\sum_{j=n}^{\infty} \pi_j \cdot \mathbb{P}_{j-n}\{\text{Ab}\}}{\mathbb{P}\{W > 0\}} \\ &= \frac{1}{A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right)} \cdot \sum_{j=n}^{\infty} \frac{\left(\frac{\lambda}{\theta}\right)^{j-n}}{\prod_{k=1}^{j-n} \left(\frac{n\mu}{\theta} + k\right)} \cdot \frac{\theta(j+1-n)}{n\mu + \theta(j+1-n)} \end{aligned}$$

(by convention, $\prod_{k=1}^0 \left(\frac{n\mu}{\theta} + k\right) \triangleq 1$)

$$\begin{aligned} &= \frac{1}{A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right)} \cdot \sum_{j=0}^{\infty} \frac{\left(\frac{\lambda}{\theta}\right)^j \cdot (j+1)}{\prod_{k=1}^{j+1} \left(\frac{n\mu}{\theta} + k\right)} = \frac{1}{A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right)} \cdot \frac{\partial}{\partial y} \left[A\left(\frac{n\mu}{\theta}, y\right) \right]_{y=\lambda/\theta} \\ &= \frac{1}{A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right)} \cdot \left[\frac{n\mu}{\lambda} + \left(1 - \frac{n\mu}{\lambda} \right) A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) \right] = \frac{1}{\rho A\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right)} + 1 - \frac{1}{\rho}, \end{aligned}$$

where the last line follows from (A.4).

Designing a Call Center with Impatient Customers

O. Garnett* A. Mandelbaum*[†] M. Reiman[‡]

March 26, 2002

ABSTRACT. The most common model to support workforce management of telephone call centers is the $M/M/N/B$ model, in particular its special cases $M/M/N$ (Erlang C, which models out busy-signals) and $M/M/N/N$ (Erlang B, disallowing waiting). All of these models lack a central prevalent feature, namely that impatient customers might decide to leave (abandon) before their service begins.

In this paper we analyze the simplest abandonment model, in which customers' patience is exponentially distributed and the system's waiting capacity is unlimited ($M/M/N + M$). Such a model is both rich and analyzable enough to provide information that is practically important for call center managers. We first outline a method for exact analysis of the $M/M/N + M$ model, that while numerically tractable is not very insightful. We then proceed with an asymptotic analysis of the $M/M/N + M$ model, in a regime that is appropriate for large call centers (many agents, high efficiency, high service level). Guided by the asymptotic behavior, we derive approximations for performance measures and propose "rules of thumb" for the design of large call centers. We thus add support to the growing acknowledgment that insights from diffusion approximations are directly applicable to management practice.

*Davidson Faculty of Industrial Engineering and Management, Technion, Haifa 32000, ISRAEL.

[†]Research supported by the fund for the promotion of research at the Technion, by the Technion V.P.R. funds - Smoler Research Fund, and B. and G. Greenberg Research Fund (Ottawa), and by the Israel Science Foundation (grant no. 388/99).

[‡]Bell Laboratories, Murray Hill, NJ 07974, USA.

Appendix B:

Calculating $E[f(V, X)]$ in an $M/M/N/B + M$ Model

To calculate $E[f(V, X)]$, we start with the following decomposition:

$$\begin{aligned} E[f(V, X)] &= E[f(V, X) \cdot 1_{(0, \infty)}(V)] + E[f(V, X) \cdot 1_{\{0\}}(V)] \\ &= E[f(V, X) \cdot 1_{(0, \infty)}(V)] + E[f(0, X)] \cdot (\pi_B + \sum_{k=0}^{N-1} \pi_k) . \end{aligned} \quad (4)$$

Here we use π to denote the stationary distribution of the queue-length process $Q(t)$, namely

$$\lim_{t \rightarrow \infty} P\{Q(t) = n\} = \pi_n, \quad n = 0, 1, 2, \dots, B.$$

A general expression for these probabilities is given by

$$\pi_k = \begin{cases} \frac{(\lambda/\mu)^k}{k!} \pi_0, & 0 \leq k \leq N \\ \prod_{j=N+1}^k \left(\frac{\lambda}{N\mu + (j-N)\theta} \right) \frac{(\lambda/\mu)^N}{N!} \pi_0, & N < k \leq B \end{cases}$$

where

$$\pi_0 = \left[\sum_{k=0}^N \frac{(\lambda/\mu)^k}{k!} + \sum_{k=N+1}^B \prod_{j=N+1}^k \left(\frac{\lambda}{N\mu + (j-N)\theta} \right) \frac{(\lambda/\mu)^N}{N!} \right]^{-1}.$$

Remark:

For a blocked customer (i.e the queue was full upon his arrival) the convention $V = 0$ is introduced.

For all functions f which seem of interest in our case, $E[f(0, X)]$ evaluates to 0 or 1. Therefore we proceed to calculate the first expression. We present three different methods for performing this calculation, each with its own virtues and drawbacks.

Our calculations require the distribution function of V . Recall that V is the *potential waiting time* of a typical customer. What is meant by a “typical” customer? Consider the sequence $\{w_n, n \in \mathbb{N}\}$, where w_n is the potential waiting time of the n -th customer. Let F_w be the stationary distribution of this sequence. Quoting from Baccelli and Hebuterne [2], F_w is also the stationary distribution of the process $\nu(t)$ - the *virtual waiting time* at time t (i.e. the time spent waiting in queue of a hypothetical infinitely-patient customer arriving at time t). Therefore a typical customer’s potential waiting time, V , has distribution function F_w .

Similarly we are interested in V_n , which is a random variable whose distribution is that of V given n customers in queue upon arrival, and all agents busy, $n = 0, 1, \dots$; V_n has distribution function F_n .

The distribution of V is not given beforehand, and is derived through analysis of the model. On the other hand, V_n can be expressed as the sum of $n + 1$ independent exponential random variables with parameters $N\mu$, $N\mu + \theta$, \dots , $N\mu + n\theta$, the i -th of these representing the period of time the customer spent in the i -th place in queue, before advancing to the $(i - 1)$ -th (due to end of service or abandonment from the queue in front of him).

Method A: Conditioning on the number of customers in the queue upon arrival, and substituting the explicit expression given by Riordan [35] (equation (83) on page 111) for $\bar{F}_n(t) = 1 - F_n(t)$, we have

$$E[f(V, X)1_{(0,\infty)}(V)] = c\pi_N \sum_{k=0}^{B-N-1} (-1)^k \frac{(\lambda/\theta)^k}{k!} I(k) \sum_{n=k}^{B-N-1} \frac{(\lambda/\theta)^{n-k}}{(n-k)!}, \quad (5)$$

where

$$I(k) = \theta^2 c \int_0^\infty \int_0^\infty f(t, x) e^{-(c+k)\theta t} e^{-\theta x} dt dx \quad \text{and} \quad c = N\mu/\theta$$

Calculating the values of $I(k)$ is usually a simple task. The main drawback of this method are the alternating signs in the first sum, which cause it to be numerically unstable. Therefore we present the next method, which avoids this problem.

Method B: Starting similarly to Method A, and using the relation

$$\sum_{k=0}^n \binom{n}{k} (-e^{-\theta t})^k = (1 - e^{-\theta t})^n$$

to eliminate one sum, we arrive at

$$E[f(V, X)1_{(0,\infty)}(V)] = \theta^2 c \pi_N \sum_{n=0}^{B-N-1} \frac{(\lambda/\theta)^n}{n!} J(n), \quad (6)$$

where

$$J(n) = \int_0^\infty \int_0^\infty f(t, x) e^{-(x+ct)\theta} (1 - e^{-\theta t})^n dx dt. \quad (7)$$

Here calculating the values of $J(n)$ tends to be more costly since the integrals must usually be solved numerically.

These methods lose some of their attractiveness when dealing with infinite buffers ($B = \infty$). Then sums appearing in both methods become infinite, and must be truncated

at some point for implementation (the alternating signs in Method A can be problematic in the aspect of truncation too). Since this case forces us to consider the issue of precision tolerance, we present the third method, which is a straightforward numerical integration.

Method C: Following through Riordan [35], and solving the more general case of any buffer size B , we arrive at the function f_V^+ , where $\frac{f_V^+}{P\{V>0\}}$ is a density function, given by

$$f_V^+(t) = N\mu\pi_N \left[1 - \frac{\gamma(B-N, \frac{\lambda}{\theta}(1-e^{-\theta t}))}{\Gamma(B-N)} \right] \cdot \exp \left\{ \frac{\lambda}{\theta}(1-e^{-\theta t}) - N\mu t \right\}, \quad t > 0. \quad (8)$$

Here Γ and γ denote the gamma and incomplete gamma functions respectively, defined by

$$\Gamma(x) = \int_0^\infty t^{x-1} \exp(-t) dt \quad \text{and} \quad \gamma(x, y) = \int_0^y t^{x-1} \exp(-t) dt, \quad y > 0.$$

Now we are left with the evaluation of the double integral

$$E[f(V, X) \cdot 1_{(0, \infty)}(V)] = \int_0^\infty \int_0^\infty f(t, x) \theta e^{-x\theta} f_V^+(t) dx dt. \quad (9)$$

The integral with respect to x is usually solved analytically and rather easily (depending on f), leaving us to perform one numerical integration (with respect to t).

Some additional remarks concerning the infinite buffer case:

Remarks:

1. When the system's buffer is unlimited, solving the stationary distribution equations involves an infinite sum. A solution is given by Palm [32], expressing the stationary distribution as a function of the easily calculated blocking probability in an $M/M/N/N$ system (denoted here $P\{Bl\}$), with the same arrival and service rates:

$$\pi_n = \begin{cases} \frac{P\{Bl\}}{1 + (A(\frac{\lambda}{N\mu}, \frac{N\mu}{\theta}) - 1)P\{Bl\}} \cdot \frac{N!}{n! \left(\frac{\lambda}{\mu}\right)^{N-n}}, & n < N \\ \frac{P\{Bl\}}{1 + (A(\frac{\lambda}{N\mu}, \frac{N\mu}{\theta}) - 1)P\{Bl\}} \cdot \frac{\left(\frac{\lambda}{\theta}\right)^{n-N}}{\left(\frac{N\mu}{\theta} + 1\right) \cdots \left(\frac{N\mu}{\theta} + (n - N)\right)}, & n \geq N \end{cases}$$

where

$$A(x, y) = \frac{ye^{xy}}{(xy)^y} \cdot \gamma(y, xy).$$

2. For $B = \infty$ the density function f_V^+ given here becomes a special case of the result by Baccelli and Hebuterne [2] for an $M/M/N + G$ model with patience distribution F , namely:

$$f_V^+(t) = N\mu\pi_N \exp \left\{ \lambda \int_0^t (1 - F(u)) du - N\mu t \right\}, \quad t > 0.$$

The M/M/n+G Queue: Summary of Performance Measures

Avishai Mandelbaum and Sergey Zeltyn

Faculty of Industrial Engineering & Management
Technion
Haifa 32000, ISRAEL

emails: avim@tx.technion.ac.il, zeltyn@ie.technion.ac.il

May 10, 2004

Contents

1	M/M/n+G: primitives and building blocks	1
1.1	Special case. Deterministic patience (M/M/n+D).	2
1.2	Special case. Exponential patience (M/M/n+M, Erlang-A).	3
2	Performance measures, exact formulae	4
3	Performance measures, QED approximations	5
4	Performance measures, efficiency-driven approximations	6
5	Guidelines for applications	7
5.1	Exact formulae: numerical calculations	7
5.2	QED approximation	7
5.3	Efficiency-driven approximation	7

1 M/M/n+G: primitives and building blocks

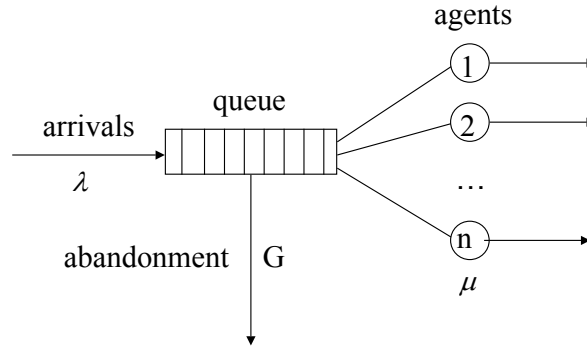
Primitives:

λ – arrival rate,

μ – service rate (= reciprocal of average service time),

n – number of servers,

G – patience distribution ($\bar{G} = 1 - G$: survival function).



Building blocks.

Define

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

Let

$$\begin{aligned} J &\triangleq \int_0^\infty \exp \{ \lambda H(x) - n \mu x \} dx , \\ J_1 &\triangleq \int_0^\infty x \cdot \exp \{ \lambda H(x) - n \mu x \} dx , \\ J_H &\triangleq \int_0^\infty H(x) \cdot \exp \{ \lambda H(x) - n \mu x \} dx . \end{aligned}$$

In addition, let

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

and

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

Finally, introduce

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

1.1 Special case. Deterministic patience (**M/M/n+D**).

Patience times equal to a constant D . Then

$$H(x) = \begin{cases} x, & 0 \leq x \leq D \\ D, & x > D \end{cases}.$$

If $\lambda - n\mu \neq 0$,

$$\begin{aligned} J &= \frac{1}{n\mu - \lambda} - \frac{\lambda}{n\mu(n\mu - \lambda)} \cdot e^{-(n\mu - \lambda)D}, \\ J(t) &= \begin{cases} \frac{1}{n\mu - \lambda} \cdot e^{-(n\mu - \lambda)t} - \frac{\lambda}{n\mu(n\mu - \lambda)} \cdot e^{-(n\mu - \lambda)D}, & t < D \\ \frac{1}{n\mu} \cdot e^{\lambda D - n\mu t}, & t \geq D \end{cases} \\ J_1 &= \frac{1}{(n\mu - \lambda)^2} - \left[\frac{1}{(n\mu - \lambda)^2} - \frac{1}{(n\mu)^2} + \frac{\lambda D}{n\mu(n\mu - \lambda)} \right] \cdot e^{-(n\mu - \lambda)D}, \\ J_H &= \frac{1}{(n\mu - \lambda)^2} \cdot [1 - e^{-(n\mu - \lambda)D}] - \frac{\lambda D}{n\mu(n\mu - \lambda)} \cdot e^{-(n\mu - \lambda)D}, \\ J_H(t) &= \begin{cases} \frac{1}{(n\mu - \lambda)^2} \cdot [e^{-(n\mu - \lambda)t} - e^{-(n\mu - \lambda)D}] + \frac{t}{n\mu - \lambda} \cdot e^{-(n\mu - \lambda)t} - \frac{\lambda D}{n\mu(n\mu - \lambda)} \cdot e^{-(n\mu - \lambda)D}, & t < D \\ \frac{D}{n\mu} \cdot e^{\lambda D - n\mu t}, & t \geq D \end{cases} \end{aligned}$$

If $\lambda - n\mu = 0$,

$$\begin{aligned} J &= D + \frac{1}{n\mu}, \\ J(t) &= \begin{cases} D - t + \frac{1}{n\mu}, & t < D \\ \frac{1}{n\mu} \cdot e^{\lambda D - n\mu t}, & t \geq D \end{cases} \\ J_1 &= \frac{D^2}{2} + \frac{D}{n\mu} + \frac{1}{(n\mu)^2}, \\ J_H &= \frac{D^2}{2} + \frac{D}{n\mu}, \\ J_H(t) &= \begin{cases} \frac{D^2 - t^2}{2} + \frac{D}{n\mu}, & t < D \\ \frac{D}{n\mu} \cdot e^{\lambda D - n\mu t}, & t \geq D \end{cases} \end{aligned}$$

1.2 Special case. Exponential patience (M/M/n+M, Erlang-A).

Patience times are iid $\exp(\theta)$. Then

$$H(x) = \frac{1}{\theta} \cdot (1 - e^{-\theta x}).$$

Define the *incomplete Gamma function*

$$\gamma(x, y) \triangleq \int_0^y t^{x-1} e^{-t} dt, \quad x > 0, \quad y \geq 0.$$

($\gamma(x, y)$ can be calculated in Matlab.) Then

$$\begin{aligned} J &= \frac{\exp\left\{\frac{\lambda}{\theta}\right\}}{\theta} \cdot \left(\frac{\theta}{\lambda}\right)^{\frac{n\mu}{\theta}} \cdot \gamma\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta}\right) \\ J(t) &= \frac{\exp\left\{\frac{\lambda}{\theta}\right\}}{\theta} \cdot \left(\frac{\theta}{\lambda}\right)^{\frac{n\mu}{\theta}} \cdot \gamma\left(\frac{n\mu}{\theta}, \frac{\lambda}{\theta} e^{-\theta t}\right) \\ J_H &= \frac{J}{\theta} - \frac{\exp\left\{\frac{\lambda}{\theta}\right\}}{\theta^2} \cdot \left(\frac{\theta}{\lambda}\right)^{\frac{n\mu}{\theta}+1} \cdot \gamma\left(\frac{n\mu}{\theta} + 1, \frac{\lambda}{\theta}\right) \\ J_H(t) &= \frac{J(t)}{\theta} - \frac{\exp\left\{\frac{\lambda}{\theta}\right\}}{\theta^2} \cdot \left(\frac{\theta}{\lambda}\right)^{\frac{n\mu}{\theta}+1} \cdot \gamma\left(\frac{n\mu}{\theta} + 1, \frac{\lambda}{\theta} e^{-\theta t}\right) \end{aligned}$$

Remark. J_1 cannot be expressed via the incomplete Gamma function. Consequently, formulae that involve J_1 (see the next page), must be calculated either numerically, or by approximations, as discussed in the sequel.

2 Performance measures, **exact** formulae

Many important performance measures of the M/M/ n +G queue can be conveniently expressed via the building blocks above. Define

$P\{\text{Ab}\}$ – probability to abandon,

$P\{\text{Sr}\}$ – probability to be served,

Q – queue length,

W – waiting time,

V – offered wait (time that a customer with infinite patience would wait).

Then

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