Sergey Zeltyn                              February 2005
zeltyn@ie.technion.ac.il

**STAT 991. Service Engineering.**
**The Wharton School. University of Pennsylvania.**

# Abandonment and Customers' Patience in Tele-Queues. The Palm/Erlang-A Model.

Based on:

- Mandelbaum A. and Zeltyn S.
  The Palm/Erlang-A Queue, with Applications to Call Centers.
  Lecture note to *Service Engineering* course.
  http://iew3.technion.ac.il/serveng/References/Erlang_A_Dec04.pdf

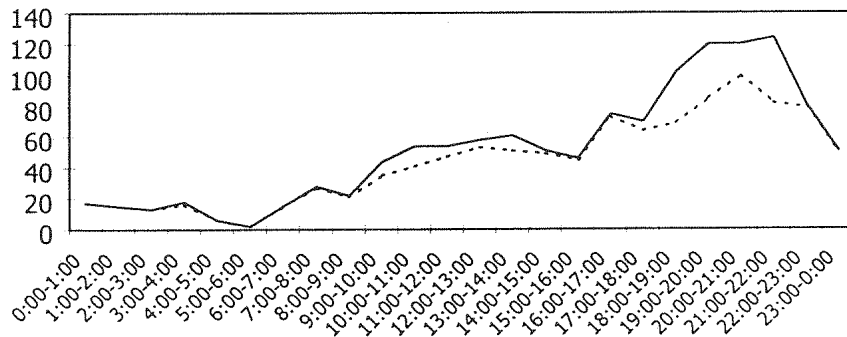- Mandelbaum A. *Service Engineering* course, Technion.
  http://iew3.technion.ac.il/serveng2005W

No abandonment in models of the previous lecture.

**However, abandonment takes place and can be very significant and very important.**
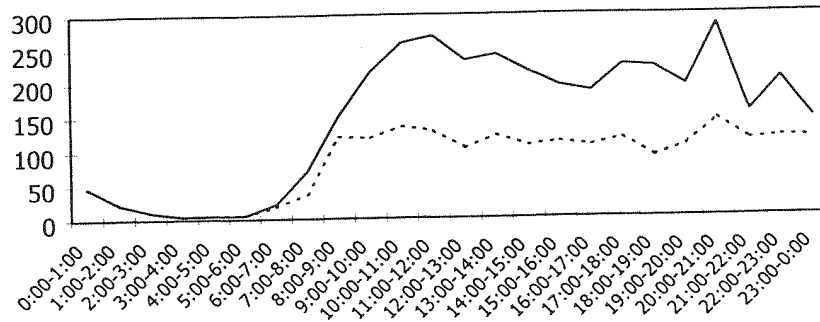
# Example 1. "Catastrophic situation".
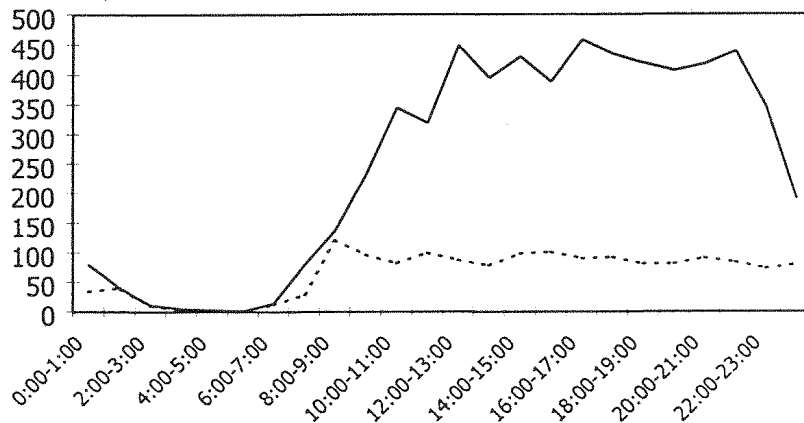## Call center of telephone company.

Average wait 72 sec, 81% calls answered (Saturday, 06/11/99)



Average wait 217 sec, 53% calls answered (Thursday, 25/11/99)



Average wait 376 sec, 24% calls answered (Sunday, 21/11/99)

# Example 2. Moderate abandonment.
## (Moderate, but still very important.)
## Health Insurance. Charlotte – Center. 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 |

## Abandonment Important Practically

- One of two customer-subjective performance measures ($2^{nd}$=Redials);

- Lost business (now);

- Poor service level (future losses);

- 1-800 costs (out-of-pocket vs. alternative);

- Self-selection: the "fittest survive" and wait less;

- Must account for (carefully) in models and measures. Otherwise, wrong picture of reality: misleading performance measures, hence staffing.
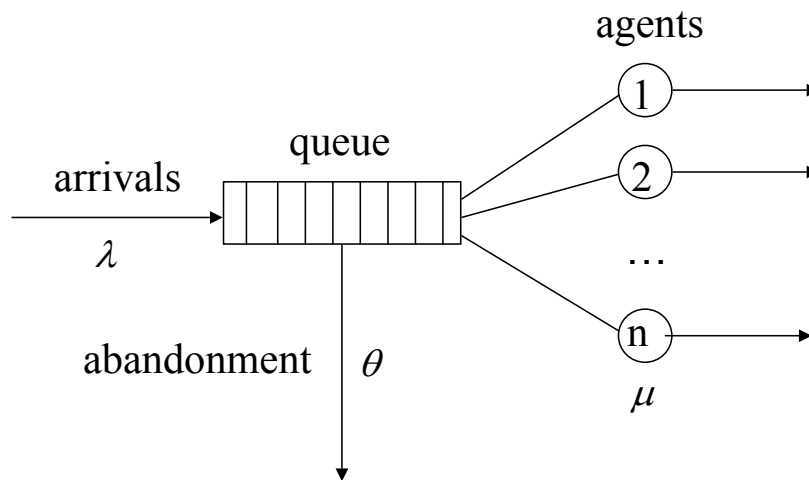
- Unstable models (vs. robustness).


## Abandonment also Interesting Theoretically

- Queueing Science
  (Paradigm: experiment, measure, model, validate);

- Research: OR + Psychology + Marketing
  (Modelling: steady-state, transient, equilibrium);

- Wide Scope of Applications: in addition to Phone,

    - VRU/IVR: opt-out-rates;
    - Internet: business-drivers (60% and more).

# The Erlang-A (Palm, M/M/n+M) Model

Simplest model with abandonment, used by well-run call centers.
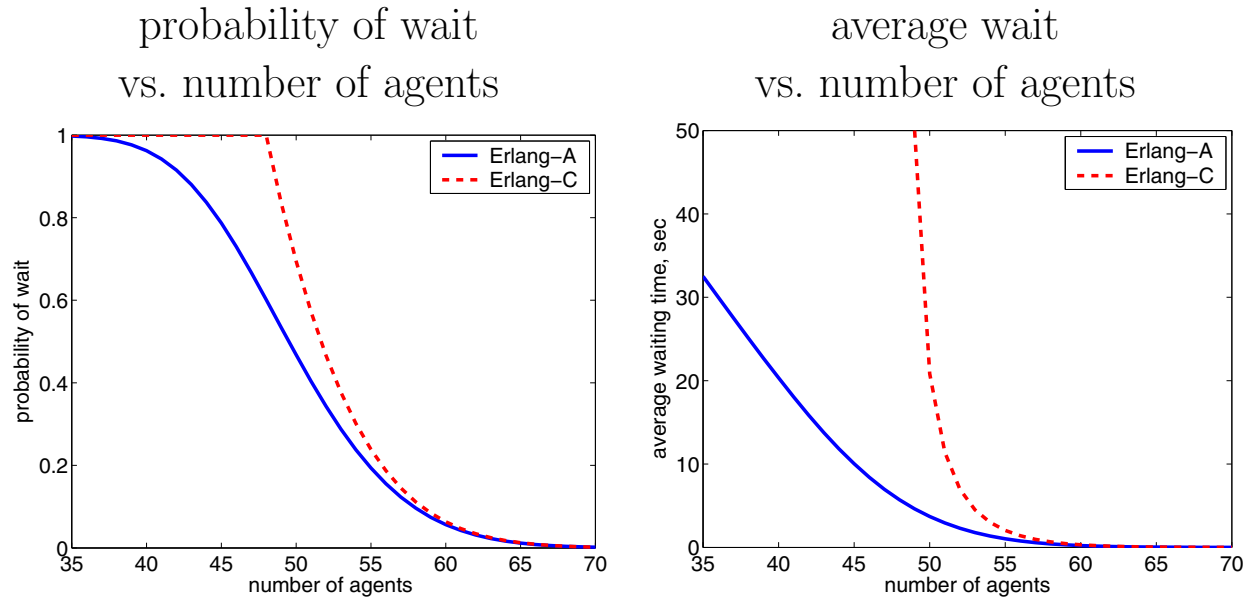
## Building blocks:

- $\lambda$ – Poisson arrival rate.

- $\mu$ – Exponential service rate.

- $n$ – number of service agents.

- $\theta$ – individual abandonment rate.



- **Patience time** $\tau \sim \exp(\theta)$:
  time a customer is willing to wait for service;

- **Offered wait** $V$:
  waiting time of a customer with infinite patience;

- If $\tau \leq V$, customer abandons; otherwise, gets service;

- **Actual wait** $W_q = \min(\tau, V)$.

# 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.

# 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:

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

## Steady-state distribution:

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

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.

# Stability

Erlang-A is always stable!

$d_j$ – death-rate in state $j$, $\ 0 < j < \infty$:

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

Bounds are death rates of M/M/$\infty$ queues with service rates $\min(\mu, \theta)$ and $\max(\mu, \theta)$.

Proof of stability:

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

$$ \le \ \sum_{j=0}^{\infty} \frac{(\lambda/\min(\mu, \theta))^j}{j!} \ = \ e^{-\lambda/\min(\mu, \theta)} \,. $$

(Use that $n\mu + (k-n)\theta \ge k \min(\mu, \theta)$.)

**Steady-state distribution via special functions (Palm):**

*Gamma function:*

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

*Incomplete Gamma function:*

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

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

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

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

Can be calculated also via recursion.

Then one can show:

$$\pi_j \ =\ \begin{cases} \pi_n\cdot\dfrac{n!}{j!\cdot\left(\frac{\lambda}{\mu}\right)^{n-j}}\,, & 0 \leq j \leq n\,, \\[4mm] \pi_n\cdot\dfrac{\left(\frac{\lambda}{\theta}\right)^{j-n}}{\Pi_{k=1}^{j-n}\left(\frac{n\mu}{\theta}+k\right)}\,, & j \geq n+1\,, \end{cases}$$

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 popular performance measure is $P\{W_q \leq T; \mathrm{Sr}\}$ (or even worse $P\{W_q \leq T \mid \mathrm{Sr}\}$).

We recommend:

- $P\{W_q \leq T; \mathrm{Sr}\}$ - fraction of well-served;

- $P\{\mathrm{Ab}\}$ - fraction of poorly-served.

or a four-dimensional refinement:

- $P\{W_q \leq T; \mathrm{Sr}\}$ - fraction of well-served;

- $P\{W_q > T; \mathrm{Sr}\}$ - fraction of served, with potential for improvement (say, a higher priority on next visit);

- $P\{W_q > \epsilon; \mathrm{Ab}\}$ - fraction of poorly-served;

- $P\{W_q \leq \epsilon; \mathrm{Ab}\}$ - fraction of those whose service-level is undetermined.

## Properties of P{Ab}:

- P{Ab} increases monotonically in $\theta$, $\lambda$;
  P{Ab} decreases monotonically in $n$, $\mu$
  (Bhattacharya and Ephremides (1991))

- M/M/$n$+G: if $\mathrm{E}[\tau]$ is fixed, deterministic patience minimizes P{Ab} (Mandelbaum and Zeltyn (2004))

## Additional important performance measures:

- Delay probability $P\{W_q > 0\}$;

- Average wait $E[W_q]$;

- ASA (Average Speed of Answer) – used extensively in call centers; usually defined as $E[W_q|\mathrm{Sr}]$;

- Agents' occupancy $\rho = \dfrac{\lambda \cdot (1 - P\{\mathrm{Ab}\})}{n\mu}$.

- Average queue-length $E[L_q]$.

# Operational Performance Measures: calculation via 4CallCenters

Performance measures of the form $E[f(V, \tau)]$.
Calculable, in numerically stable procedures.
For example,

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

# Operational Performance Measures: calculation via 4CallCenters



Erlang-A parameters:

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

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

- $\mathrm{P}\{W_q \leq T; \mathrm{Sr}\} = 71.1\%$;

- $\mathrm{P}\{W_q > T; \mathrm{Sr}\} = 87.5\% - 71.1\% = 16.4\%$;

- $\mathrm{P}\{W_q > \epsilon; \mathrm{Ab}\} = 12.5\% - 3.9\% = 8.6\%$;

- $\mathrm{P}\{W_q \leq \epsilon; \mathrm{Ab}\} = 3.9\%$.

- Delay probability $\mathrm{P}\{W_q > 0\} = 100\% - 45.8\% = 54.2\%$.

# Additional performance measures



- Average Time in Queue = $\mathrm{E}[W_q] = 15$ sec;

- ASA = $\mathrm{E}[W_q | \mathrm{Sr}] = 13.8$ sec;

- Agents' Occupancy $\rho = 87.5\%$;

- Average Queue Length $\mathrm{E}[L_q] = 1.3$.

# Operational Performance Measures: calculation via special functions

For example,

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

# Operational Performance Measures: calculation via M/M/n+G formulae

M/M/$n$+G – generalization of Erlang-A, patience times distributed with cdf $G(\cdot)$. See

`http://iew3.technion.ac.il/serveng/References/references.html`

- Mandelbaum A. and Zeltyn S. (2004) M/M/$n$+G queue. Summary of performance measures;

- Zeltyn S. (2004) Call centers with impatient customers: exact analysis and many-server asymptotics of the M/M/$n$+G queue, Ph.D. Thesis.

Explained how to adapt M/M/$n$+G to Erlang-A:

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

# The relation P{Ab}/E[$W_q$]

**Theoretical:** In Erlang-A (and other queueing models with $\exp(\theta)$ patience):

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

**Proof.** Balance equation:

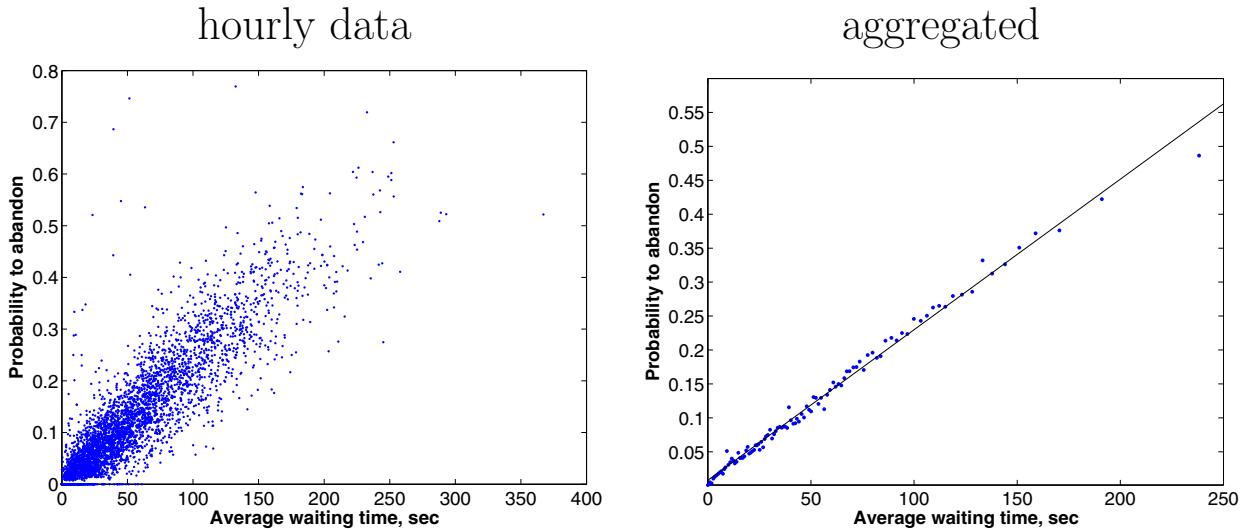$$\theta \cdot E[L_q] \ = \ \lambda \cdot P\{Ab\} \,. \tag{1}$$

Little's formula:

$$E[L_q] \ = \ \lambda \cdot E[W_q] \,. \tag{2}$$

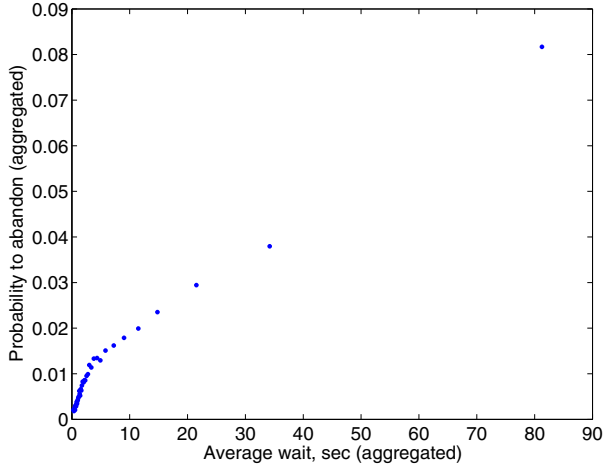Substitute (2) into (1). ∎

# Empirical relations

## Israeli bank: yearly data

hourly data                                          aggregated



The graphs are based on 4158 hour intervals.

# U.S. bank

### Retail



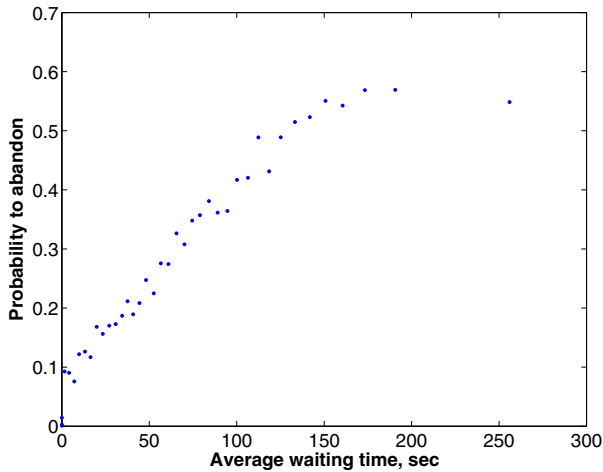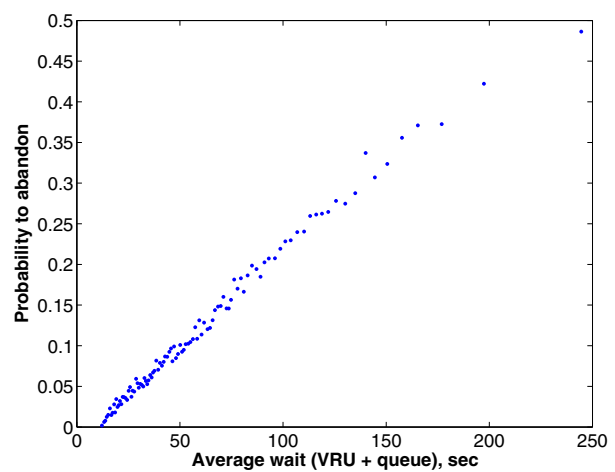### Telesales



Retail – significant abandonment during first seconds of wait.

## Linear patterns 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: $c + \exp(\theta), \ \ c > 0$.

# Erlang-A: parameter estimation and prediction

**Estimation:** inference from historical data (e.g. exp, normal) were parameters assumed fixed over time.

**Prediction:** forecast behavior of sample outside of original data set.

## Arrivals ($\lambda$)

- Typically Poisson, time-varying rates, constant at 15/30/60 min scale;

- Significant uncertainty concerning future rates $\Rightarrow$ prediction;

- Predict separately *daily volumes* and *fraction* of arrivals per time interval.

## Services ($\mu$)

- Typically stable from day to day $\Rightarrow$ estimation;

- Can change depending on time-of-day;

- Typically, service time $\neq$ talk time.

**First approach:**
service time = talk time + wrap-up time (after-call work) + ...;

**Second approach:**

$$\text{service time} = \frac{\text{Total Working Time} - \text{Total Idle Time}}{\text{Number of Served Customers}}.$$
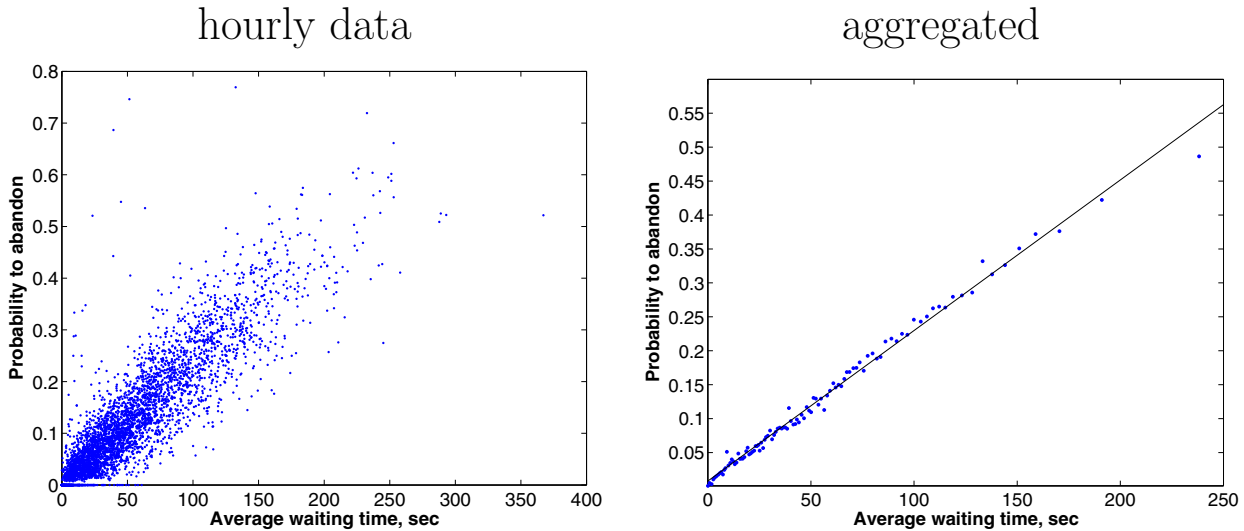
# Number of agents (n)

- Output of WFM software given $\lambda$, $\mu$, $\theta$, performance goals. One gets number of FTE's (Full Time Equivalent positions).

- Agents on schedule = FTE's $\cdot$ RSF (Rostered Staff Factor) (RSF > 1). Reasons: absenteeism, unscheduled breaks, . . .

- Obtaining historical data on $n$ can be hard.

# Patience ($\boldsymbol{\theta}$)

Observations are **censored**! (heavily)

- Customer abandoned $\Rightarrow$ patience $\tau$ known;

- Customer served $\Rightarrow$ offered wait $V$ known $\Rightarrow \tau > V$.

Avoiding direct "uncensoring": use $\ \mathrm{P}\{\mathrm{Ab}\} \ = \ \theta \cdot \mathrm{E}[W_q]\,.$

hourly data                          aggregated



Regression $\Rightarrow$ average patience $(1/\theta) \approx \dfrac{250}{0.56} \approx 446$ sec.
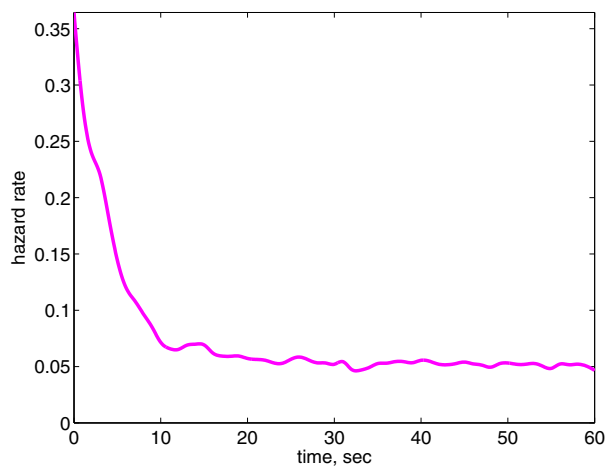
# Estimating patience distribution

Are patience times really exponential?
To "uncensor data" use Kaplan-Meier (product-limit) estimator.
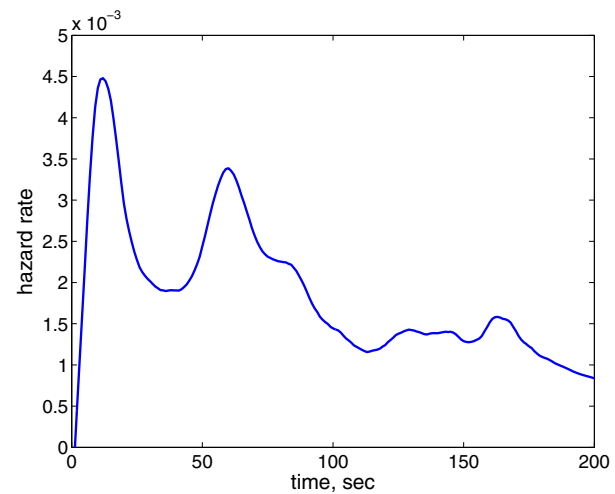Output: estimates of survival function and hazard rate.

## Empirical hazard rates of patience times

U.S. bank                                    Israeli bank



## Israeli bank: regular vs. priority customers



19

## Israeli bank: service types



IN – Internet Assistance;   NE – Stock Transactions;
NW – New Customers;   PS – Regular

## Conclusions:

- Patience times are, in general, non-exponential;

- Most tele-customers are **very** patient;

- Kaplan-Meier is very informative concerning patience
  *qualitative* patterns (abandonment peaks, comparisons, ... );

- Kaplan-Meier can be problematic concerning estimation of
  *quantitative* characteristics (mean, variance, median).
  $\widehat{E[\tau]} = \int_0^\infty \widehat{S(x)}dx$, where $S(x)$ - survival function of patience.
  However, $\widehat{S(x)}$ not reliable for large $x$.

**Question:** can we apply Erlang-A with non-exponential patience?

# Fitting a simple model to a complex reality

## Erlang-A Formulae vs. Data Averages (Israeli Bank)



P{Ab}

E[$W_q$]

P{$W_q > 0$}

# Conclusions:

- Points: hourly data vs. Erlang-A output;

- Formulae with continuous $n$ used;

- Patience estimated via $P\{Ab\}/E[W_q]$ relation;

- Erlang-A estimates – close upper bounds.

# Fitting a simple model to a complex reality: Patience index

How to define (im)patience?

$$\text{Theoretical Patience Index} \;\triangleq\; \frac{\text{time willing to wait}}{\text{time required to wait}}$$
$$= \frac{\text{average patience}}{\text{average offered wait}}. \qquad (3)$$
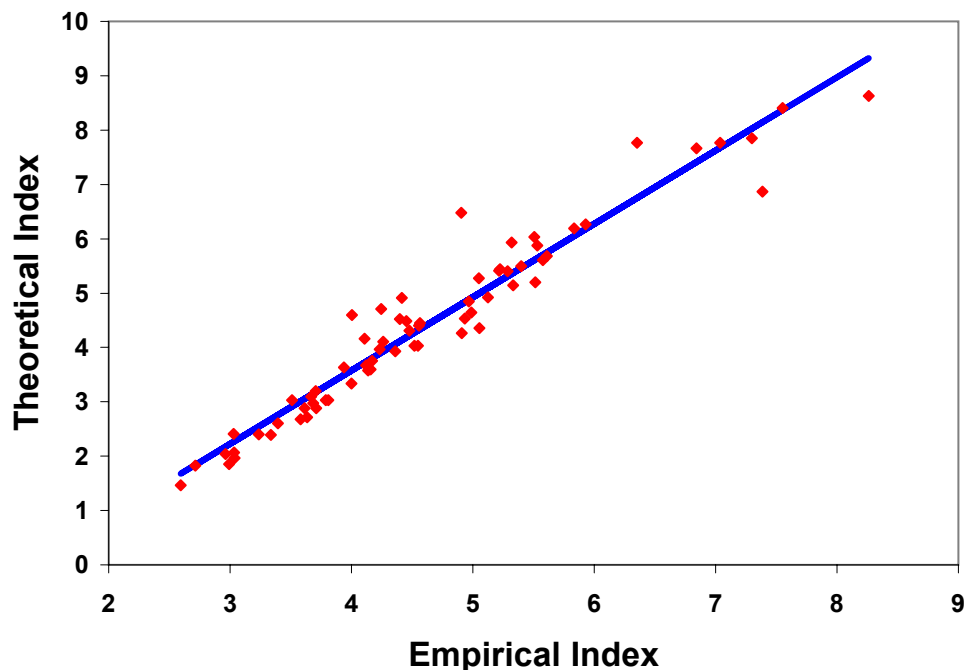
Calculation can be difficult.

$$\text{Empirical Patience Index} \;\triangleq\; \frac{\%\ \text{served}}{\%\ \text{abandoned}}. \qquad (4)$$

Easily calculable from ACD reports.

If $\tau$ and $V$ exponentially distributed, (4) is MLE of (3).

## 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**  $1 + 2 \times \dfrac{360K}{90K} = 1 + 2 \times 4 = \mathbf{9}$ min.

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

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

$\qquad\qquad\qquad\qquad\qquad\uparrow \qquad\qquad\qquad\qquad\qquad \uparrow$

$\qquad\qquad\qquad\qquad\text{definition} \qquad\qquad\qquad \text{measure}$

# Customer-Focused Queueing Theory

Waiting experience can be summarized by:

1. Time that a customer *expects* to wait;

2. Time that a customer is *willing* to wait ($\tau$, patience or need);

3. Time that a customer *must* 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;
Then left with $(\tau, V, W_q)$, as introduced before.

200 abandonment in Direct-Banking: perceived vs. actual waiting.

| 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 impatient customers

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

## Israeli bank: Internet-support customers



Rational abandonment from invisible queues: Mandelbaum, Shimkin, Zohar.

# Advanced features of 4CallCenters

## Advanced profiling

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

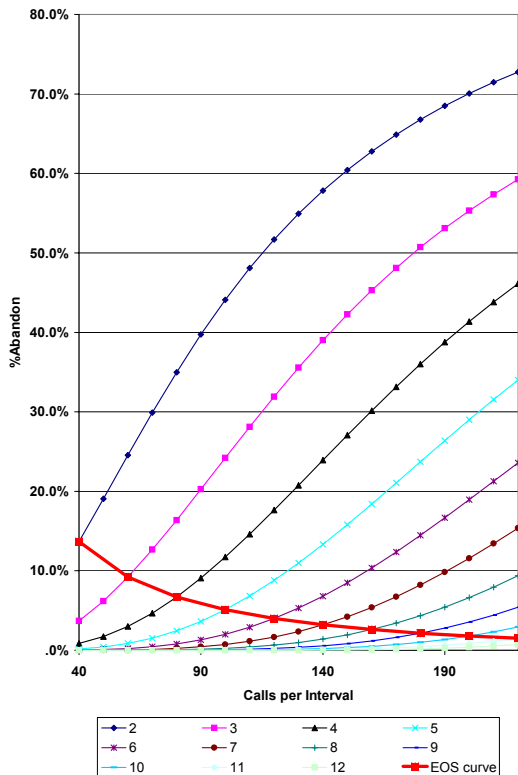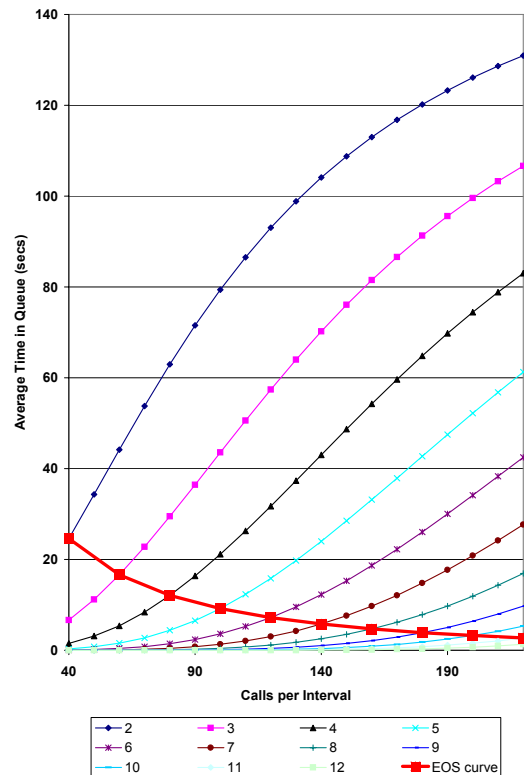**Example:** $1/\mu = 2$ minutes, $1/\theta = 3$ minutes;
$\lambda$ varies from 40 to 230 calls per hour, in steps of 10;
$n$ varies from 2 to 12.

Probability to abandon          Average wait



Red curve: EOS (Economies-Of-Scale).
Why the two graphs are similar?
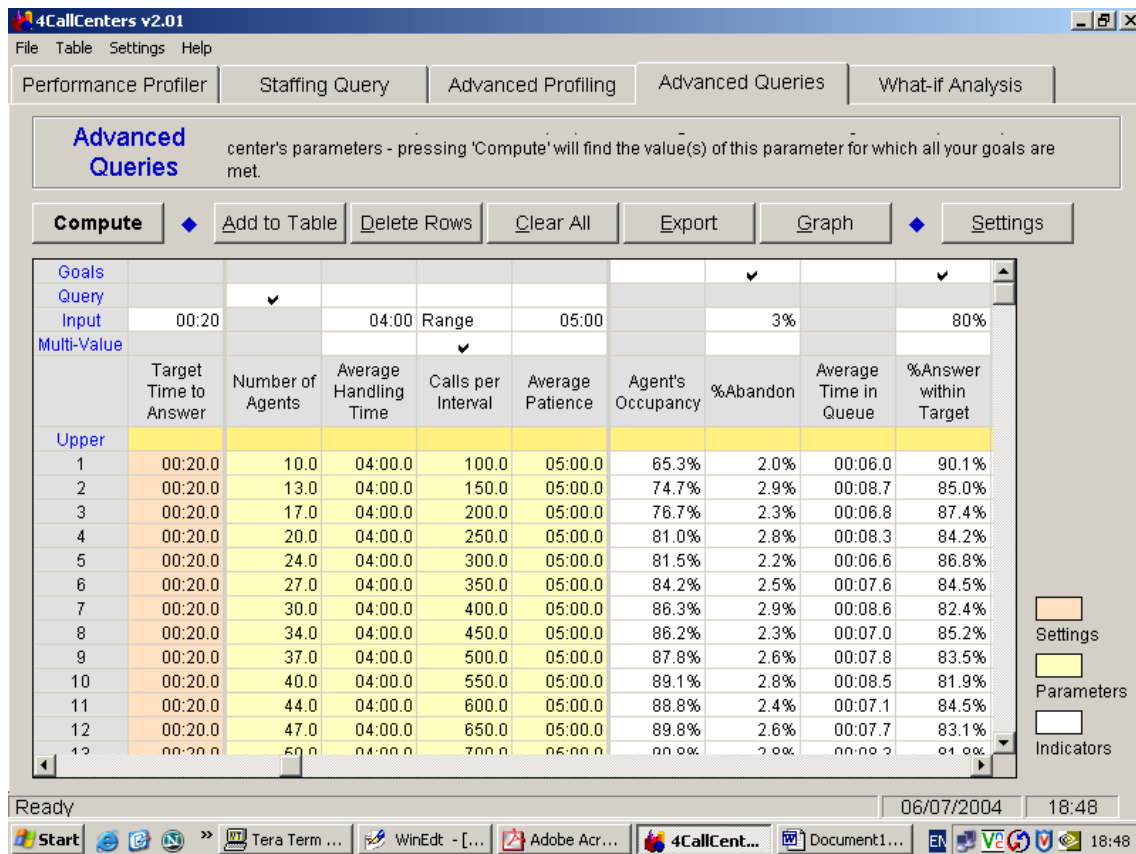
# Advanced staffing queries

Multiple performance goals.

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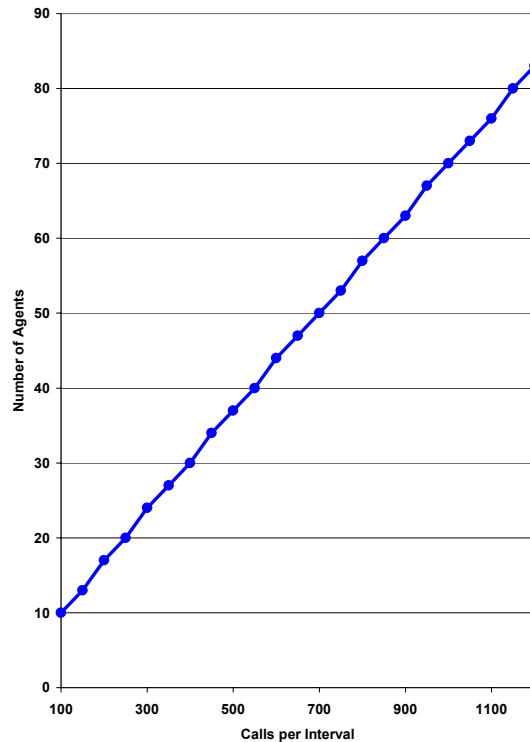$\lambda$ varies from 100 to 1200, in steps of 50.

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

## 4CallCenters output



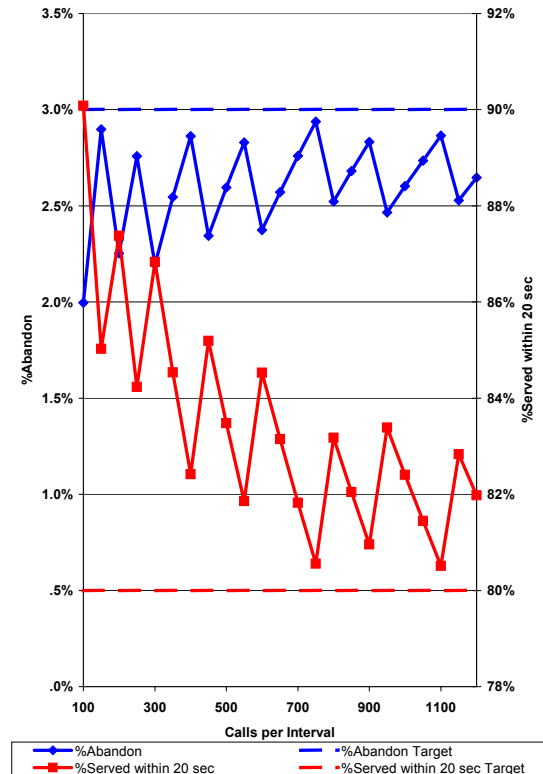| | Target Time to Answer | Number of Agents | Average Handling Time | Calls per Interval | Average Patience | Agent's Occupancy | %Abandon | Average Time in Queue | %Answer within Target |
|---|---|---|---|---|---|---|---|---|---|
| Upper | | | | | | | | | |
| 1 | 00:20.0 | 10.0 | 04:00.0 | 100.0 | 05:00.0 | 65.3% | 2.0% | 00:06.0 | 90.1% |
| 2 | 00:20.0 | 13.0 | 04:00.0 | 150.0 | 05:00.0 | 74.7% | 2.9% | 00:08.7 | 85.0% |
| 3 | 00:20.0 | 17.0 | 04:00.0 | 200.0 | 05:00.0 | 76.7% | 2.3% | 00:06.8 | 87.4% |
| 4 | 00:20.0 | 20.0 | 04:00.0 | 250.0 | 05:00.0 | 81.0% | 2.8% | 00:08.3 | 84.2% |
| 5 | 00:20.0 | 24.0 | 04:00.0 | 300.0 | 05:00.0 | 81.5% | 2.2% | 00:06.6 | 86.8% |
| 6 | 00:20.0 | 27.0 | 04:00.0 | 350.0 | 05:00.0 | 84.2% | 2.5% | 00:07.6 | 84.5% |
| 7 | 00:20.0 | 30.0 | 04:00.0 | 400.0 | 05:00.0 | 86.3% | 2.9% | 00:08.6 | 82.4% |
| 8 | 00:20.0 | 34.0 | 04:00.0 | 450.0 | 05:00.0 | 86.2% | 2.3% | 00:07.0 | 85.2% |
| 9 | 00:20.0 | 37.0 | 04:00.0 | 500.0 | 05:00.0 | 87.8% | 2.6% | 00:07.8 | 83.5% |
| 10 | 00:20.0 | 40.0 | 04:00.0 | 550.0 | 05:00.0 | 89.1% | 2.8% | 00:08.5 | 81.9% |
| 11 | 00:20.0 | 44.0 | 04:00.0 | 600.0 | 05:00.0 | 88.8% | 2.4% | 00:07.1 | 84.5% |
| 12 | 00:20.0 | 47.0 | 04:00.0 | 650.0 | 05:00.0 | 89.8% | 2.6% | 00:07.7 | 83.1% |
| 13 | 00:20.0 | 50.0 | 04:00.0 | 700.0 | 05:00.0 | 90.8% | 2.8% | 00:08.3 | 81.8% |

# 4CallCenters. Advanced staffing queries.
## Dynamics of staffing level and performance.

Recommended staffing level

Target performance measures



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