

# Skills-Based Routing and its Operational Complexities

## Service Engineering

Wharton's Call Center Forum

May 9, 2003

**e.mail : [avim@tx.technion.ac.il](mailto:avim@tx.technion.ac.il)**

**Website: [http://ie.technion.ac.il/serveng\(2003\)](http://ie.technion.ac.il/serveng(2003))**

**Tool : <http://4CallCenters.com> (register & use)**

## Supporting Material (in Web-Site)

Gans, Koole, and M.: “Telephone Call Centers: [Tutorial, Review and Research Prospects.](#)” Invited review to *MSOM*, 2002

Garnett and M.: "An [Introduction](#) to Skills-Based Routing and its Operational Complexities", Teaching Note, 2000; under revision.

Borst, M. and Reiman.: “Dimensioning Large Telephone Call Centers.” Accepted to *OR*, 2002.

Atar, M. and Reiman: “Scheduling a Multi-Class Queue with Many Exponential Servers: Asymptotic Optimality in Heavy-Traffic.” Submitted to *Annals Appl Prob*, 2002.

M. and Stolyar: “Scheduling Flexible Servers with Convex Delay Costs: Heavy-Traffic Optimality of the Generalized  $c\mu$ -Rule.” Under revision to *OR*, 2002.

Yahalom and M.: "Optimal Scheduling for a Multi-server Multi-class Non-preemptive Queueing System", in preparation.

Armony and M.: "Optimal Routing, Staffing and Networking in the QED (Halfin-Whitt) Regime: Homogenous Customers and Heterogeneous Servers", in preparation.

# Contents

1. Introduction to **Skills-Based-Routing (SBR)**:  
Examples: CRM, Distributed Call Centers  
Truly a Multi-Disciplinary Challenge
2. Focus: **Agent Scheduling, Customer Routing** and some  
Workforce Management (**Staffing**)
3. **Operational Regimes**: Quality-Driven, Efficiency-Driven  
**QED** (Quality and Efficiency Driven; Haflin-Whitt)
4. **Square-Root Staffing** Laws: Scale Economies (Erlang)
5. Efficiency-Driven SBR: Index strategies in the General Case
6. QED SBR: Only special Cases (I, N, V, Upside-Down V)
7. **Beyond** Conventional Queueing Theory:  
Abandonment, Retrials; Time-Varying Queues

## NationsBank CRM:

### What are the relationship groups?

---

- The groups
  - RG1 : high-value customers
  - RG2 : marginally profitable customers (with potential)
  - RG3 : unprofitable customer
- What does it mean for a customer in each group to be profitable? Customer Revenue Management

---

3

Wharton

## NationsBank's Design of the Service Encounter

---

### Examples of Specifications: Assignable Grade Of Service (AGOS)

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

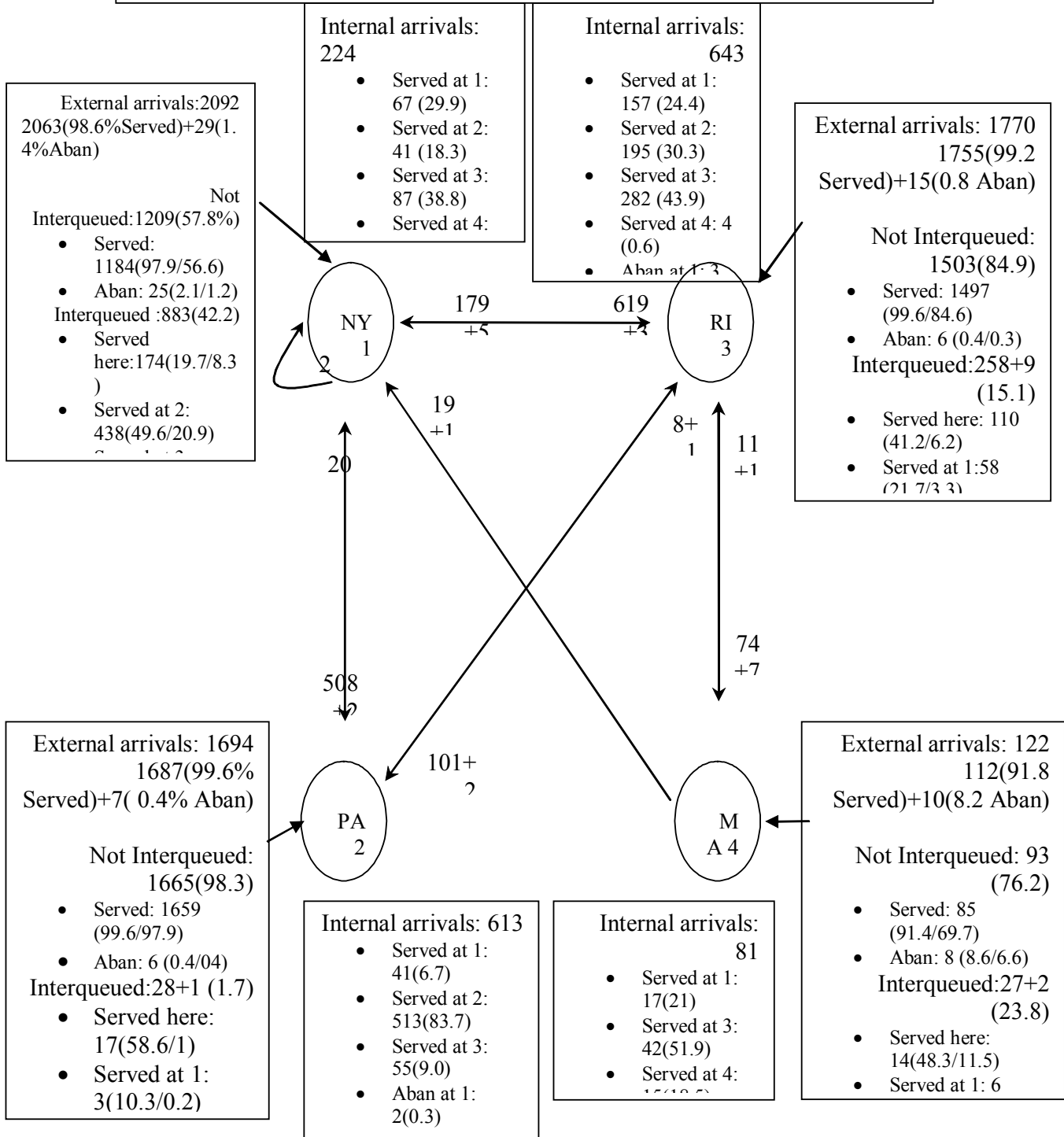
---

5

Wharton

# Distributed Call Center: Member1

**10 AM – 11 AM (03/19/01): Interflow Chart Among the 4 Call**



# Workforce Management: Hierarchical Operational View

**Forecasting** Customers: Statistics, Time-Series  
Agents : HRM (Hire, Train; Incentives, Careers)

**Staffing:** Queueing Theory

Service Level, Costs

# FTE's (Seats)  
per unit of time

**Shifts:** IP, Combinatorial Optimization; LP

Union constraints, Costs

Shift structure

**Rostering:** Heuristics, AI (Complex)

Individual constraints

Agents Assignments

**Skills-based Routing:** Stochastic Control

# An Introduction to Skills-Based Routing and its Operational Complexities

**By Ofer Garnett and Avishai Mandelbaum**

**Technion, ISRAEL**

( **Full** Version )

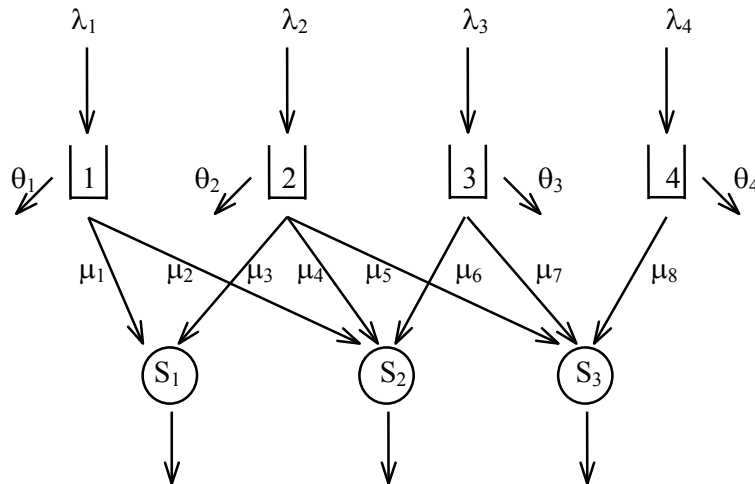
Contents:

- 1. Introduction**
- 2. N-design with single servers**
- 3. X-design with multi-server pools and impatient customers**
- 4. Technical Appendix: Simulations – the computational effort**

Acknowledgement: This teaching-note was written with the financial support of the Fraunhofer IAO Institute in Stuttgart, Germany. The authors are grateful to Dr. Thomas Meiren and Prof. Klaus-Peter Fährnich of the IAO for their assistance and encouragement.

## Introduction

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



Here the  $\lambda$ 's designate arrival rates, the  $\mu$ 's service rates, the  $\theta$ 's abandonment rates, and the  $S$ 's are the number of servers in each server-pool.

## **Skills-Based Design:**

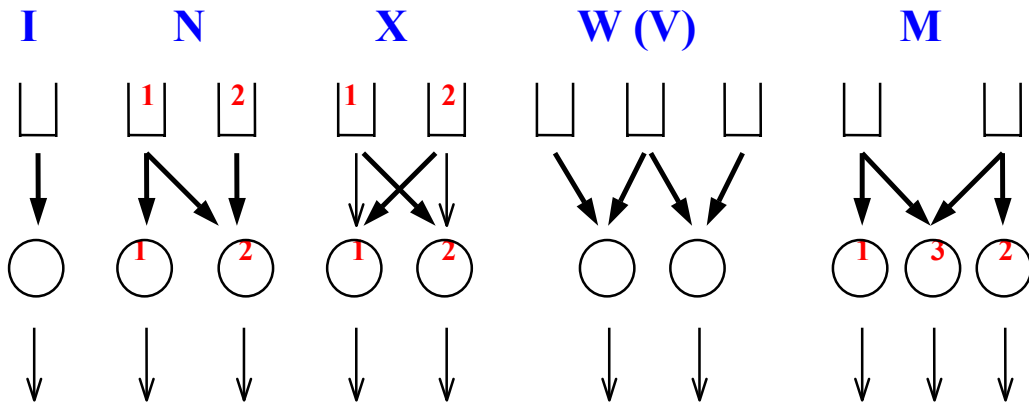
- **Queue:** "customer-type" requiring a specific type of service;
- **Server-Pool:** "skills" defining the service-types it can perform;
- **Arrow:** leading into a server-pool define its skills / constituency.

For example, a server with skill 2 (**S2**) can serve customers of type 3 (**C3**) at rate  $\mu_6$  customers/hour.

Customers of type 3 arrive randomly at rate  $\lambda_3$  customers/hour, equipped with an impatience rate of  $\theta_3$ .



## Some Canonical Designs - Animation



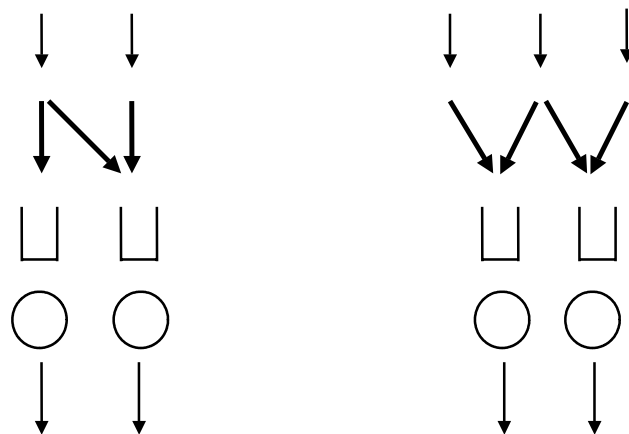
**I** – dedicated (specialized) agents

**N**: for example,

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

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

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



## Major **Design / Engineering** Decisions

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

## Major **Control** Decisions

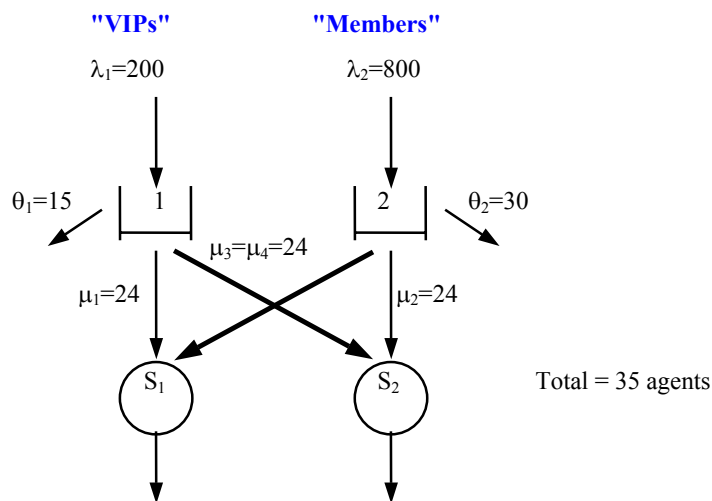
4. Matching customers and agents (**OR**)
  - **Agent Scheduling**: Whenever an agent turns idle and there are queued customers, which customer (if any) should be routed to this agent.
  - **Customer Routing**: Whenever a customer arrives and there are idle agents, which agent (if any) should serve this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

## **Multidisciplinary** Challenging Research

**Skills-Based Routing:** protocol for online matching of S's and C's.

- **Prevalent:** Static Priorities of customer types and agent skills
- **Index-based:** Dynamic Priorities via continuous review
- **Threshold-based:** Dynamic Management by Exception
- **Others:** discrete review, credit schemes (SLA), scripts; call backs

Example: **Scripts** for Staffing, Scheduling, Routing



**Setup A :** (X-design)

"VIP" servers :  $S_1 = 20$

- If "VIP" queue not empty serve the "VIP" queue + all "Members" waiting more than **40** seconds, as a single FIFO queue.
- If "VIP" queue is empty, serve the first in the "Member" queue.

"Member" servers :  $S_2 = 15$

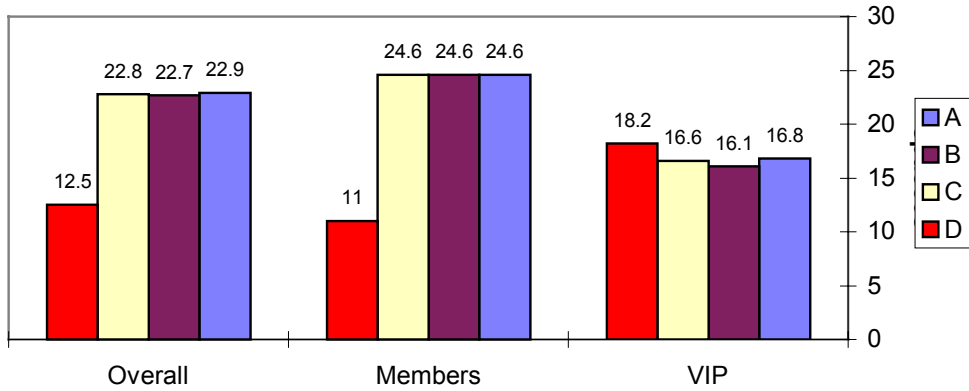
- If "Member" queue not empty serve the "Member" queue + all "VIPs" waiting more than **6** seconds, as a single FIFO queue.
- If "Member" queue is empty, serve the first in the "VIP" queue.

**Setup C :** (V-design; feasible since servers are assumed equally skilled.)

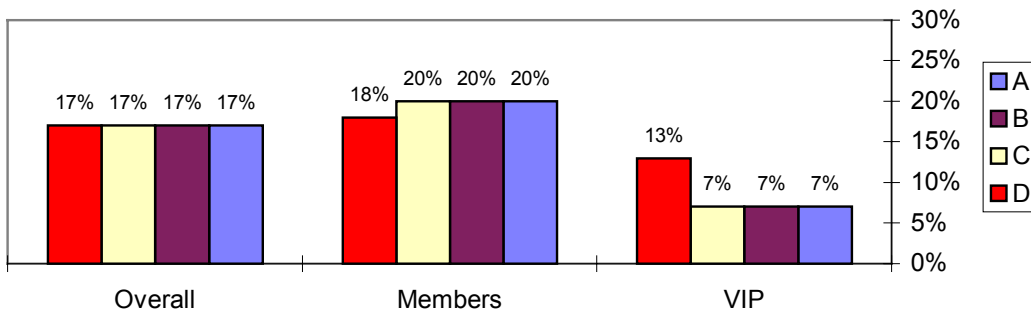
Total servers: 35

- Serve as a FIFO queue, but "VIPs" enter the queue with a virtual **15** second wait (i.e. as if they had joined the queue 15 seconds earlier).

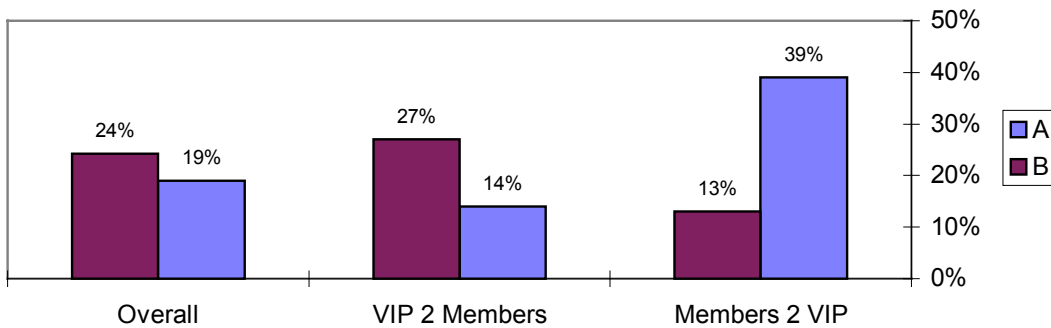
**Chart 2 : 1000 Calls/hour - ASA**



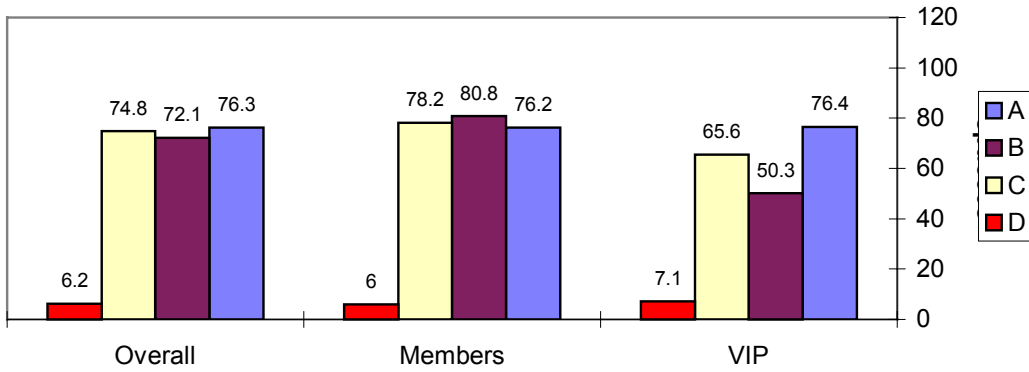
**Chart 3 : 1000 Calls - Abandonment**



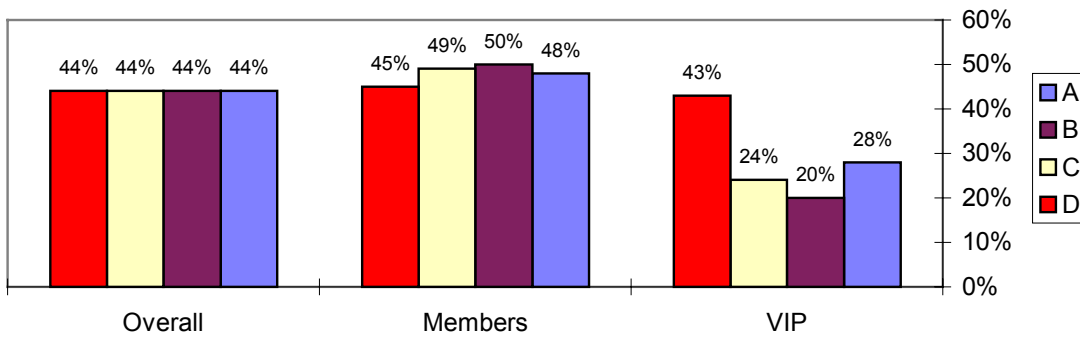
**Chart 4 : 1000 Calls - Overflows**



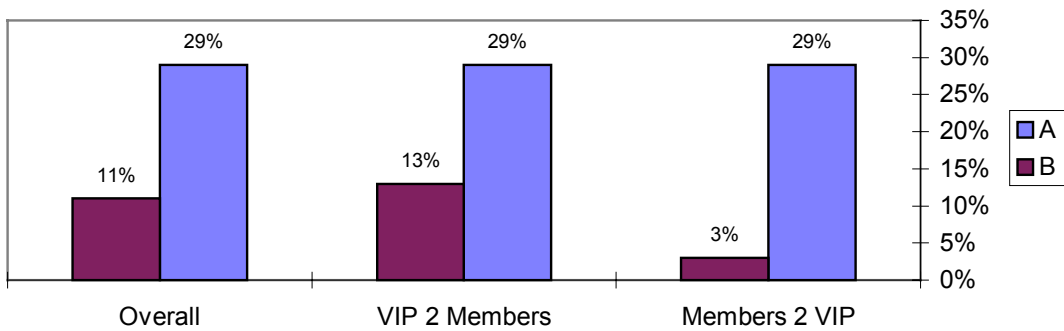
## WHAT IF : 1500 Calls/hour - ASA



## Chart 7 : 1500 Calls - Abandonment



## Chart 8 : 1500 Calls - Overflows



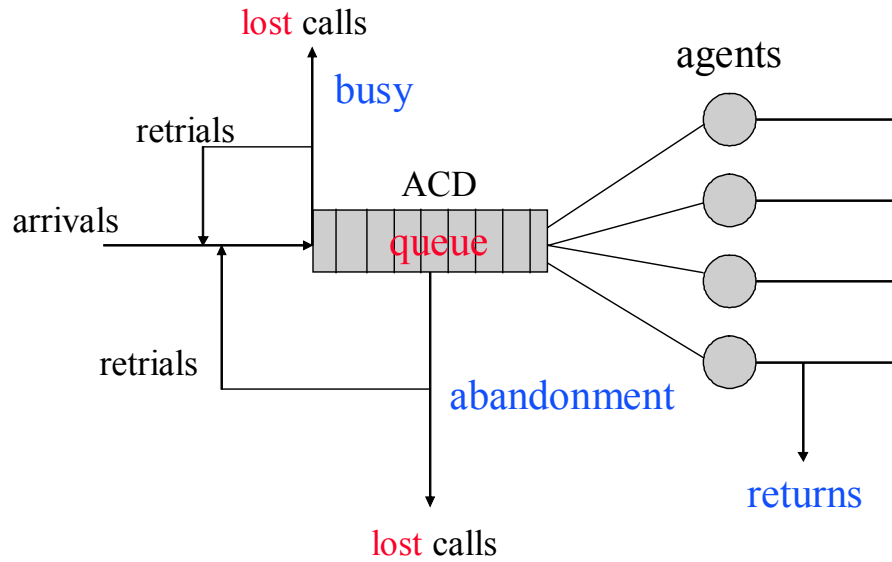
## Reality

- Technology enables smart systems
- Reality becomes increasingly complex
- Solutions are urgently needed
- Theory lags significantly behind needs
- Ad-hoc methods: heuristics, simulation-based

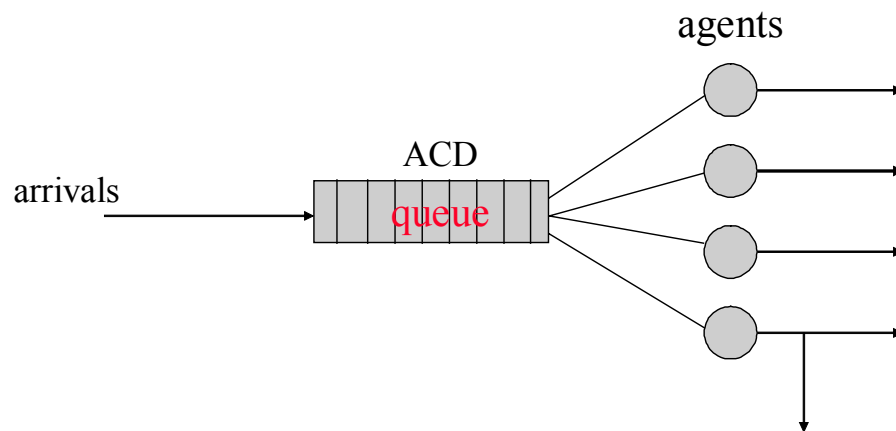
## Research Status

- Efficiency-driven SBR well understood and solved
- QED SBR is challenging and advancing
- **Small yet significant models for theoretical insight**
- Principles/Guidelines for design, staffing, control
- Implementation: fine-tuning of parameters, scale-up

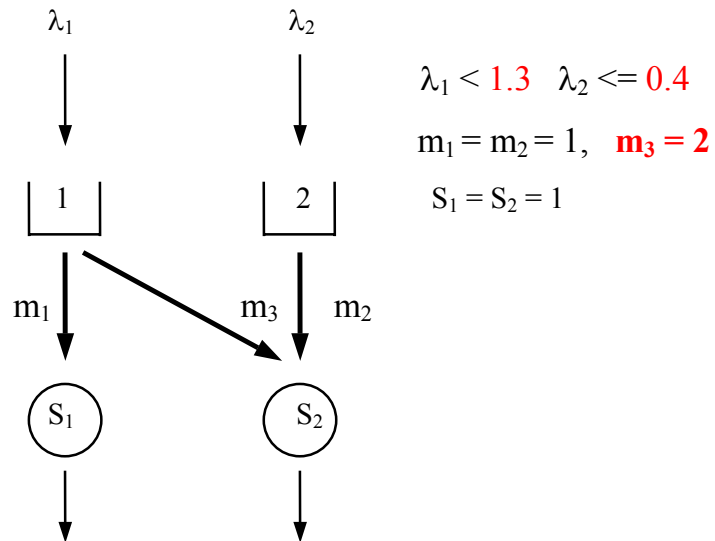
# The Basic Call Center



$$\text{Erlang-C} = M/M/N$$

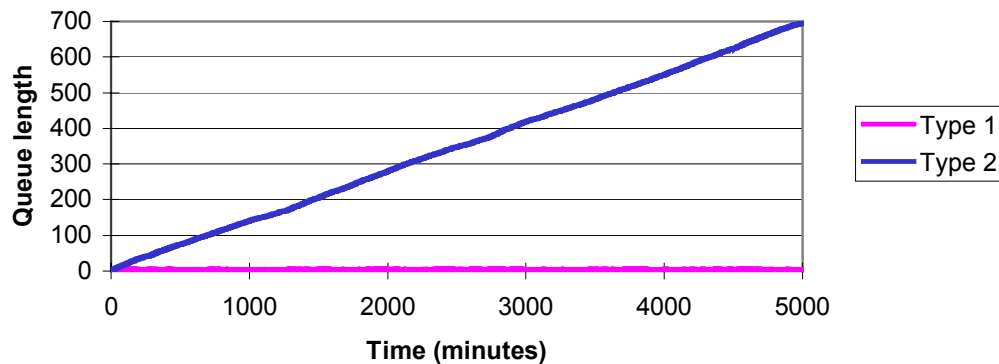


## Static Priorities (Cross-Training): Some Subtleties



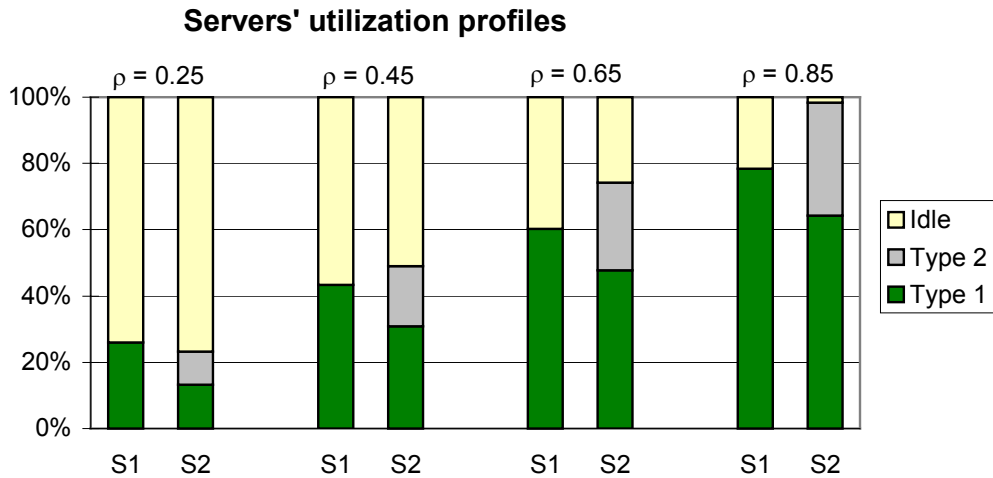
- C1 are **VIP**, hence S2 **helps** S1 by giving priority to C1 over C2.
- If both servers are idle - **C<sub>i</sub>** customers are routed to server **S<sub>i</sub>**

Queue length: S2 helps with VIP C1, Heavy Loading -



Q2 "explodes", while Q1 is negligibly small – why ?

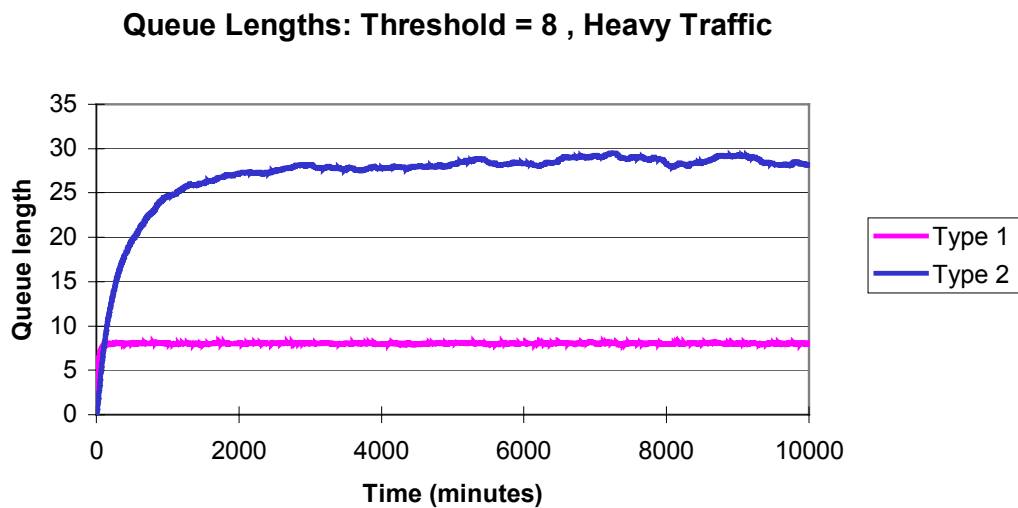




Instability: S2 **overworked** serving C1 and neglecting C2, while S1 is **20%** idle.

To avoid "overzealous help", apply **Threshold control**:

S2 assists S1 **only when Q1 is at or above a certain threshold**



Both Q1 and Q2 are stable.

Now fine-tuning of the threshold value

## Efficiency-Driven SBR - the "EASY" Case

Examples: Scarce agents, hence must be well utilized.

Email-dominance, hence can delay response.

Classical **special** case: **V**-design

- **Agent Scheduling**: upon service completion, if

1. Same mean service times: serve the costliest queue (largest **c**)
2. Same delay costs: serve the shortest service (smallest **m**)
3. Generally: serve the largest **c/m** (= index).

**General** (N, X, W, M, ... ) solution: **Index Control** is optimal

- **Customer Routing**: irrelevant, since essentially all customers wait.

- **Agent Scheduling**: upon service completion, the server chooses the queue with the largest index and serves its "oldest" customer.

- **Index**: marginal waiting-cost per unit of average service-time  
(Example: actual "waiting-time" of the "oldest" customer in queue)

**However**: well-managed telephone services are **not**

(at least should not be) Efficiency-Driven !?



## Rough Performance Analysis

**Peak** 10:00 – 10:30 a.m., with 100 agents  
400 calls  
3:45 minutes average service time

Offered load  $R = \lambda \times E(S)$   
 $= 400 \times 3:45 = 1500 \text{ min./30 min.}$   
 $= 50 \text{ Erlangs}$

Occupancy  $\rho = R/N$   
 $= 50/100 = 50\%$

$\Rightarrow$  **Quality-Driven Operation** (Light-Traffic)

Above:  $R = 50$ ,  $N = R + 50$ ,  $\approx$  all served immediately.

Rule of Thumb:  $N = \lceil R + \delta R \rceil$ ,  $\delta > 0$  service-grade.

**Quality-driven:** 100 agents, 50% utilization

⇒ **Can** increase offered load - but **by how much?**

**Erlang-C      N=100      E(S) = 3:45 min.**

$\lambda/hr$	$\rho$	$E(W_q) = \text{ASA}$	% Wait = 0
800	50%	0	100%
1000	62.5%	0	100%
1200	75%	0	99.7%
1400	87.5%	0:02 min.	88%
1500	93.8%	0:15 min.	60%
1550	96.9%	0:48 min.	35%
1580	98.8%	2:34 min.	15%
1585	<b>99.1%</b>	<b>3:34 min.</b>	12%

⇒ **Efficiency-driven Operation (Heavy Traffic)**

Above:  $R = 99$ ,  $N = R + 1$ ,  $\approx$  all delayed.

Rule of Thumb:  $N = \lceil R + \gamma \rceil$ ,  $\gamma > 0$  service grade.

## Changing N (Staffing)

$$E(S) = 3:45$$

$\lambda$ /hr	$N$	OCC	ASA	% Wait = 0
1585	100	99.1%	3:34	12%
1599	<b>100</b>	99.9%	<b>59:33</b>	0%
1599	<b>100+1</b>	98.9%	<b>3:06</b>	13%
1599	102	98.0%	1:24	24%
1599	105	<b>95.2%</b>	<b>0:23</b>	<b>50%</b>

⇒ **New Rationalized Operation**

**Heavy traffic**, in the sense that  $OCC > 95\%$ ;

**Light traffic**, 50% answered immediately

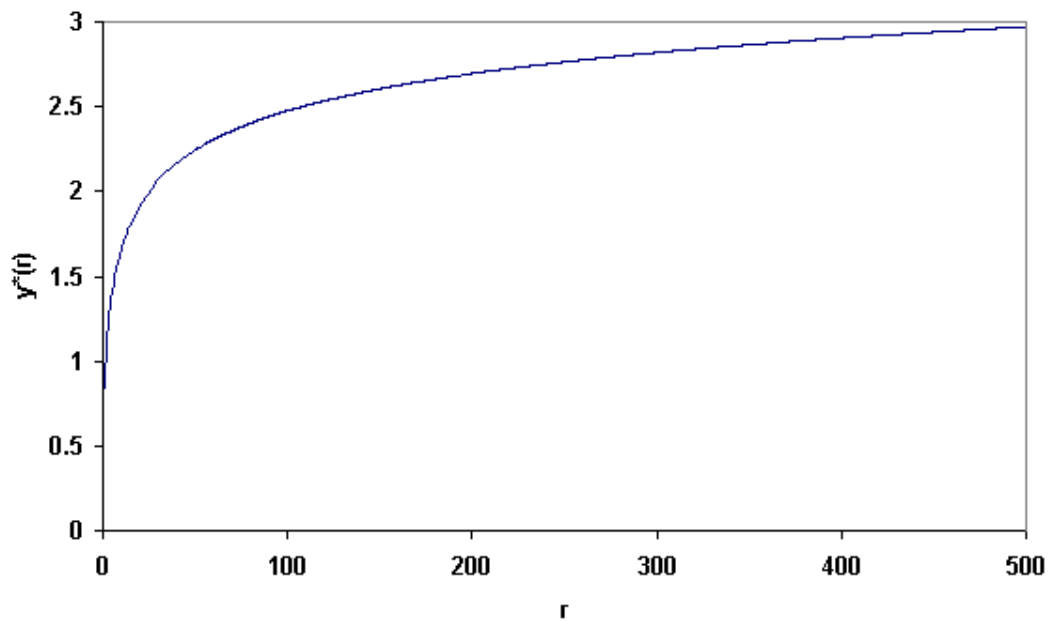
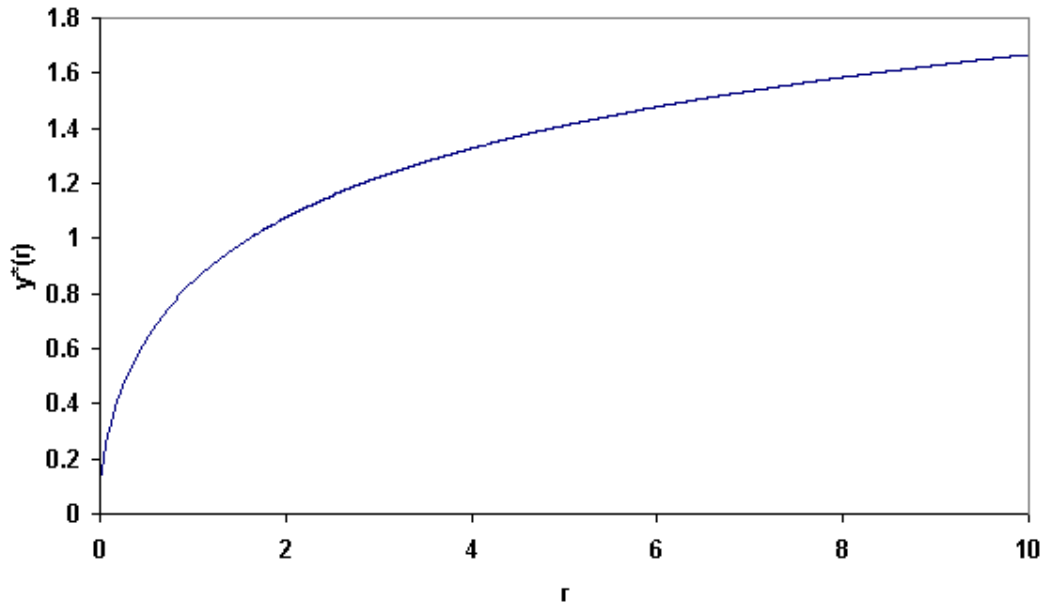
**QED Regime = Quality- and Efficiency-Driven Regime**

Above:  $R = 100$ ,  $N = R + 5$ , **50% delayed.**

**√ Safety-Staffing**  $N = \lceil R + \beta \sqrt{R} \rceil$ ,  $\beta > 0$  .

Square-Root Safety Staffing:  $N = R + y^*(r)\sqrt{R}$

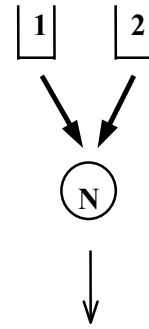
$r = \text{cost of delay (1-800)} / \text{cost of staffing (salary)}$



## V-Design: Pure Scheduling

N agents, fully flexible

C1 = VIP



## Optimal Scheduling: Agent Reservation

- C1(=VIP) always served, if possible;
- C2 served only if # of idle agents exceeds a threshold.

**QED** regime:  $\sqrt{\cdot}$  Safety-Staffing, as before.

**Threshold Size** (relative to N) determines Service Levels:

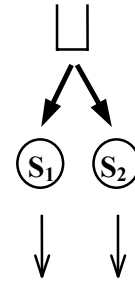
- Large: C1 is Q-served, C2 is E-served
- Small: C1 and C2 indistinguishable QED
- Moderate: C1 is Q-served, C2 is QED



# Upside-Down-V Design: Pure Routing

Homogeneous Customers

Heterogeneous Agents: **S1 = Faster**



Optimal Routing: **"Slow-Server"** phenomenon

- S1(=Fast) always employed, if possible;
- S2(= Slow) employed if # in queue exceeds a threshold.

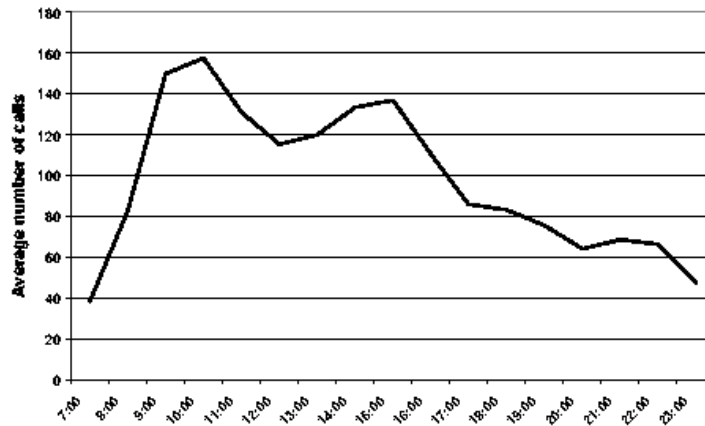
**QED** regime:  $\sqrt{\cdot}$  Safety-Staffing for S1+S2.

- No threshold needed: just have all servers work when possible, ensuring that the "fast" get the priority.
- Can do also detailed staffing: how many S1 and S2.
- Distributed call centers: similar

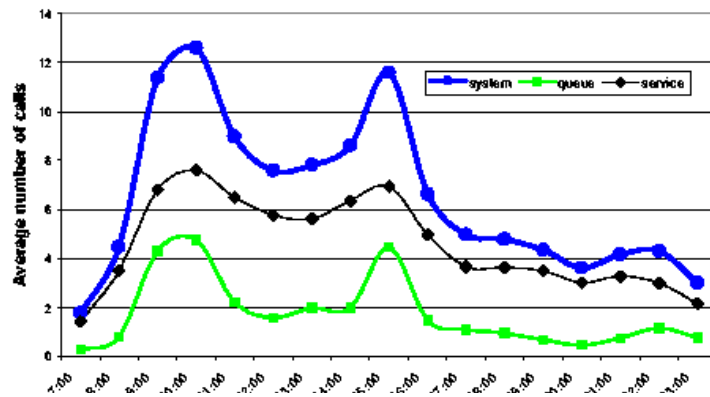
But N-Design active challenging research

# Beyond Erlang-C: Predictable Variability

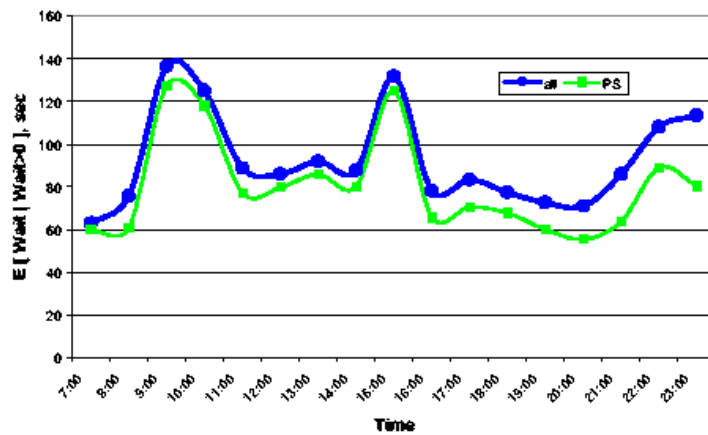
## Arrivals



## Queues



## Waiting



# Erlang-A: (Im)Patience

## Common Performance

### BCMS SKILL REPORT

Switch Name: FDC/HAMPDEN

Date: 7:00 pm WED MAR 10, 1999

Skill: 37

Skill Name: !BA AUTH1

Acceptable Service Level: 30

DAY	ACD		AVG		AVG		AVG		TOTAL		TOTAL		% IN
	CALLS	ANS	SPEED	ABAND	ABAND	TALK	AFTER	FLOW	FLOW	AUX/	AVG	SERV	
3/04/99	637	0:19	219	0:26	1:57	92:05	0	0	4310:06	8.7	66		
3/05/99	849	0:06	135	0:06	1:35	179:58	0	0	4299:43	11.3	85		
3/06/99	1330	0:11	363	0:13	1:42	280:22	0	0	5592:29	13.2	73		
3/07/99	1213	0:12	358	0:18	1:46	226:20	0	0	4830:15	11.5	72		
3/08/99	631	0:26	382	0:33	1:57	150:50	0	0	3743:04	7.9	49		
3/09/99	570	0:40	487	0:43	1:52	148:41	0	0	3979:04	6.7	38		
3/10/99	512	0:29	292	0:28	1:41	243:06	0	0	3046:00	7.9	50		
SUMMARY	<u>5742</u>	0:18	<u>2236</u>	0:26	1:46	1321:22	0	0	****:**	9.6	63		

Arrivals

Abandons 40%

Switch Name: FDC/HAMPDEN

Date: 7:00 pm WED MAR 10, 1999

Skill: 46

Skill Name: !BA AUTHORIZATION

Acceptable Service Level: 30

DAY	ACD		AVG		AVG		AVG		TOTAL		TOTAL		% IN
	CALLS	ANS	SPEED	ABAND	ABAND	TALK	AFTER	FLOW	FLOW	AUX/	AVG	SERV	
3/04/99	1185	0:22	479	0:31	2:08	190:16	0	0	4213:22	8.4	61		
3/05/99	1805	0:05	308	0:04	1:38	337:20	0	0	4299:43	11.3	84		
3/06/99	2437	0:12	642	0:12	1:51	444:03	0	0	5592:29	13.2	73		
3/07/99	2260	0:13	558	0:14	1:46	326:33	0	0	4830:14	11.5	74		
3/08/99	1260	0:35	676	0:28	2:06	308:19	0	0	3743:04	7.9	48		
3/09/99	1126	0:40	653	0:34	2:10	250:40	0	0	3979:04	6.7	44		
3/10/99	890	0:30	472	0:32	2:16	162:13	0	0	3046:00	7.9	51		
SUMMARY	<u>10963</u>	0:19	<u>3788</u>	0:22	1:55	2019:24	0	0	****:**	9.6	65		

30%

### BCMS SKILL REPORT

Switch Name: FDC/HAMPDEN

Date: 7:01 pm WED MAR 10, 1999

Skill: 33

Skill Name: GA Authorization

Acceptable Service Level: 30

DAY	ACD		AVG		AVG		AVG		TOTAL		TOTAL		% IN
	CALLS	ANS	SPEED	ABAND	ABAND	TALK	AFTER	FLOW	FLOW	AUX/	AVG	SERV	
3/04/99	1248	0:27	61	0:42	1:57	330:04	0	0	4390:04	9.5	72		
3/05/99	1521	0:14	37	0:20	1:58	353:48	0	0	6035:35	13.0	85		
3/06/99	2388	0:20	130	0:34	2:10	550:16	0	0	6369:58	14.4	76		
3/07/99	1748	0:14	66	0:30	2:08	432:16	0	0	4616:11	11.7	82		
3/08/99	925	0:18	50	1:00	1:53	191:06	0	0	3835:19	8.4	81		
3/09/99	856	0:26	57	0:53	1:54	125:16	0	0	4388:02	8.1	73		
3/10/99	959	1:15	125	1:55	1:48	186:44	0	0	4198:39	8.9	53		
SUMMARY	<u>9645</u>	0:25	<u>526</u>	0:57	2:02	2169:30	0	0	****:**	10.6	76		

6%

### BCMS SKILL REPORT

Switch Name: FDC/HAMPDEN

Date: 7:02 pm WED MAR 10, 1999