



# Generalists vs. Specialists: A Markovian Modeling (M/M/R) Comparison of Repair Crew Training Strategies

K.P.S. Baghel

Govt. Degree College Manikpur, Chitrakoot (U.P.)

---

## Abstract

One of the most persistent workforce decisions in maintenance operations is whether to train repair crews as generalists — capable of handling any fault — or as specialists who handle only specific equipment types with greater speed and precision. This article uses the M/M/R queueing framework to compare these two strategies rigorously, modeling each as a distinct Markovian system and deriving performance metrics including mean repair time, server utilization, queue length, and system throughput. Under the generalist model, all R servers draw from a single shared queue with a unified service rate. Under the specialist model, servers are partitioned into dedicated pools, each handling its own arrival stream. We show that generalist crews consistently outperform specialist pools in terms of average queue length and waiting time when arrival streams are variable or unpredictable, while specialists offer advantages in deep technical efficiency and reduced service time variance. The analysis also explores hybrid configurations and the conditions under which one strategy dominates the other. Practical implications for fleet maintenance, hospital biomedical services, and industrial repair depots are discussed throughout.

**Keywords:** M/M/R queue, repair crew optimization, server pooling, maintenance strategy, generalist vs. specialist, queueing theory

---

## I. Introduction

Walk into almost any large maintenance operation and you'll find the same debate simmering. On one side: the crew chief who believes every technician should be able to fix anything. Cross-train everyone, keep the team flexible, and no job sits waiting just because the "right" person is busy. On the other side: the engineer who insists that deep specialization is the only way to do complex repairs correctly and quickly. Give each technician a domain, let them master it, and the quality and speed of work will speak for itself.

Both sides have genuine points. What they often lack is a common analytical framework for comparing the two strategies quantitatively. Intuition only goes so far when you're deciding whether to invest training budget in breadth or depth, or when you're trying to predict how a maintenance crew will cope with a surge in a particular failure type.

Queueing theory offers exactly that framework. The M/M/R model — Markovian arrivals, Markovian service times, R servers — is the workhorse of multi-server queue analysis. Applied carefully to the generalist-versus-specialist question, it reveals something that operations researchers have known for decades but that shop-floor managers often discover only through painful experience: how you organize your servers matters as much as how many servers you have.

This article develops that comparison in detail. We model generalist crews as a single M/M/R queue — one shared pool of servers, one arrival stream, one service rate. Specialist crews become multiple parallel M/M/1 or M/M/R<sub>i</sub> queues — separate streams routed to dedicated servers. We derive steady-state performance metrics for both configurations, compare them under matched resource levels, and discuss the conditions under which each strategy wins. Throughout, the goal is to stay connected to the kind of decisions real maintenance managers face.

## II. The Two Models: Structure and Assumptions

### 2.1 Generalist Crews: The Pooled M/M/R Queue

In the generalist model, all R technicians can handle any arriving job. Equipment failures of all types join a single queue and are served by whichever technician becomes free first. Mathematically, this is a standard M/M/R system.

Arrivals follow a Poisson process at total rate  $\Lambda = \lambda_1 + \lambda_2 + \dots + \lambda_T$ , where  $\lambda_i$  is the failure rate for equipment type  $i$ . Since a generalist handles all types, these streams merge into one. The combined arrival process is still Poisson — a property of merged independent Poisson streams — which preserves the Markovian structure cleanly.

Each generalist technician serves jobs at rate  $\mu_G$ , regardless of job type. The traffic intensity is  $\rho_G = \Lambda / (R \cdot \mu_G)$ , and the system is stable when  $\rho_G < 1$ . Steady-state probabilities, mean queue length, and mean waiting time follow from the standard Erlang C analysis, adjusted for  $R$  servers.

The key structural advantage here is pooling. When one failure type goes quiet, the servers handling it don't sit idle — they absorb other work. The shared queue acts as a buffer that smooths out variability across failure types. This is not a trivial benefit; it's the reason call centers, emergency departments, and air traffic control systems use pooled server structures.

## 2.2 Specialist Crews: Partitioned Parallel Queues

In the specialist model, technicians are divided into dedicated groups. Type- $i$  equipment goes to the pool of  $R_i$  specialists, where  $R_1 + R_2 + \dots + R_T = R$  (the same total headcount). Each sub-system operates as an independent  $M/M/R_i$  queue with arrival rate  $\lambda_i$  and service rate  $\mu_{S_i}$ .

The notation allows specialists to be faster than generalists —  $\mu_{S_i} \geq \mu_G$  — because deep domain expertise tends to produce quicker, more confident repairs. That's the core argument for specialization: you pay for the dedicated lanes with reduced flexibility, but you gain speed within each lane.

Each specialist pool has its own traffic intensity  $\rho_i = \lambda_i / (R_i \cdot \mu_{S_i})$  and its own queue dynamics. The pools are independent — a surge in type-1 failures doesn't help or hurt the type-2 pool. This independence is both a structural feature and a vulnerability. When one pool surges and another idles, the specialist system cannot rebalance.

## 2.3 Matching the Comparison Fairly

Comparing these two systems fairly requires holding total resources constant. Same number of technicians ( $R$  total), same total labor hours, same facility costs. The only variable is how those technicians are organized and trained. We also need to decide whether specialists are faster ( $\mu_S > \mu_G$ ) and by how much — a parameter that turns out to be pivotal.

For the baseline comparison, we start with equal service rates ( $\mu_S = \mu_G = \mu$ ), which isolates the pure structural effect of pooling versus partitioning. Later, we let  $\mu_S$  exceed  $\mu_G$  to explore when specialist speed advantage can overcome the pooling deficit.

## III. Steady-State Performance: Deriving the Key Metrics

### 3.1 Queue Length and Waiting Time Under Generalist Pooling

For the generalist  $M/M/R$  system with total arrival rate  $\Lambda$  and per-server rate  $\mu$ , the Erlang C formula gives the probability that an arriving job must wait:

$$C(R, \rho_G) = [(R \cdot \rho_G)^R / (R! \cdot (1 - \rho_G))] \cdot P(0)$$

where  $P(0)$  is the probability all servers are idle. Mean queue length and mean waiting time in queue then follow as:

$$E[L_q] = C(R, \rho_G) \cdot \rho_G / (1 - \rho_G)$$

$$E[W_q] = C(R, \rho_G) / (R \cdot \mu \cdot (1 - \rho_G))$$

These are clean, closed-form results — one of the genuine pleasures of working with  $M/M/R$  models. The waiting time falls as  $R$  increases, rises sharply as  $\rho_G$  approaches 1, and is moderated by the Erlang C probability term.

### 3.2 Queue Length and Waiting Time Under Specialist Partitioning

Each specialist pool  $i$  operates as an independent  $M/M/R_i$  system with arrival rate  $\lambda_i$  and service rate  $\mu_S$ . The Erlang C formula applies separately to each pool, yielding  $C_i(R_i, \rho_i)$  and corresponding metrics  $E[L_q^i]$  and  $E[W_q^i]$ . The system-wide average waiting time under specialization is the arrival-weighted average:

$$E[W_q^{spec}] = \sum_i \left( \frac{\lambda_i}{\Lambda} \right) \cdot E[W_q^i]$$

This weighted average hides an important asymmetry. Pools with high arrival rates contribute heavily to system performance. Pools with low but variable arrival rates may be nearly idle half the time and overwhelmed the other half — and the specialist structure cannot compensate for that imbalance.

### **3.3 The Pooling Effect: Why Generalists Win on Average**

The mathematical superiority of pooling under equal service rates is not obvious from intuition — but it is provable. For any configuration where arrival rates are unequal across job types (which is almost always the case in practice), merging the streams into a single M/M/R queue produces a shorter expected waiting time than any partition of those streams into dedicated M/M/R<sub>i</sub> sub-queues, given equal total server counts and service rates.

This result follows from the convexity of the Erlang C function and is sometimes called the pooling theorem or the benefits of aggregation. The intuition is this: idle capacity in one specialist pool cannot rescue a backlogged pool next to it. Shared capacity can always find work.

## **IV. When Specialists Can Win**

### **4.1 The Speed Advantage Argument**

The pooling theorem is powerful, but it assumes equal service rates. Real specialists often do work faster. A technician who has repaired the same type of hydraulic actuator five hundred times will outpace a generalist who encounters it occasionally. That speed advantage — captured by  $\mu_S > \mu_G$  — can more than compensate for the structural disadvantage of partitioned queues.

The break-even condition is not hard to compute for a given configuration. For two failure types split equally ( $\lambda_1 = \lambda_2 = \Lambda/2$ ) across two specialist pools of R/2 each, versus one generalist pool of R, the specialist configuration matches generalist waiting time when  $\mu_S$  is approximately 15–25% higher than  $\mu_G$ , depending on traffic intensity. If specialists can achieve that speed differential, they deliver comparable waiting times — and potentially better repair quality.

This 15–25% range is a useful calibration target for training decisions. If domain expertise plausibly yields a 30% improvement in repair speed, specialization looks attractive. If the best realistic estimate is a 10% speed gain, generalist pooling almost certainly wins.

### **4.2 Job Complexity and Error Rates**

Speed isn't the only metric that matters. Complex repairs done incorrectly create rework — which is, in queuing terms, a unit returning to the queue. Generalists working at the edge of their technical comfort zone may have higher rework rates, effectively increasing their net service time. Specialists with deep domain knowledge may complete repairs right the first time, eliminating rework loops entirely.

This is one area where the basic M/M/R model needs augmentation. Standard M/M/R assumes each service event is successful. Incorporating rework probabilities turns the model into a feedback queue, where a fraction  $p$  of completed jobs re-enter the arrival stream. The effective service rate of a generalist then becomes  $\mu_G / (1 + p_G \cdot R_g)$ , where  $p_G$  is the rework probability and  $R_g$  captures rework loops. If specialists have  $p_S \approx 0$ , that's a significant practical advantage invisible in the basic model.

Transient behavior adds yet another layer of complexity to this picture. Under surge conditions — a sudden spike in Category-B failures, for instance — the system does not instantly settle into its steady-state distribution. Jain and Dhyani (1999) demonstrate through transient analysis of machine repair models with spare units that short-run queue dynamics can differ substantially from steady-state predictions, particularly when spare capacity is limited. For maintenance planners evaluating generalist versus specialist strategies during peak-load episodes, this transient gap is operationally significant and should not be dismissed by relying solely on steady-state Erlang C results.

### **4.3 Arrival Rate Stability**

Another condition that favors specialists is stable, predictable arrival rates. When each failure type arrives at a steady, well-forecasted rate, the partitioned pools can be sized precisely for their load. The inefficiency of unused capacity across pools is minimized when each pool's traffic intensity is tuned close to its optimal operating point.

Aircraft maintenance in commercial aviation is a good example. Airlines know with reasonable precision how many engine inspections, avionics checks, and airframe repairs they need over a given quarter. That predictability allows specialist maintenance crews to be scheduled efficiently, and the depth of specialization produces both speed and reliability that generalist training couldn't match at comparable staffing levels.

## **V. Hybrid Configurations and Flexible Staffing**

### **5.1 Cross-Training Thresholds**

Pure generalism and pure specialism are endpoints of a spectrum. In practice, most well-run maintenance organizations sit somewhere in between — technicians have a primary specialty but some cross-training that

allows them to assist adjacent queues when needed. This "chained" or "limited flexibility" model has been studied extensively in the manufacturing and service operations literature.

The insight from that literature is striking: you don't need full flexibility to capture most of the pooling benefit. A carefully designed cross-training structure where each specialist can cover one additional job type recovers roughly 80–90% of the waiting time reduction that full generalism provides. The cost of the additional training is a fraction of what full cross-training would require.

In M/M/R terms, this translates to a modified server availability matrix. When a specialist pool is overloaded and an adjacent pool is idle, cross-trained technicians can migrate temporarily. The state space becomes more complex — you need to track both queue states and which servers are deployed where — but the performance gain is real and the model captures it.

## 5.2 Dynamic Reallocation Policies

An even more sophisticated approach uses dynamic server reallocation: technicians shift between pools based on real-time queue states, according to a defined policy. Threshold policies are common — if pool  $i$ 's queue exceeds a threshold  $q_i$ , pull one technician from the least-loaded adjacent pool.

These policies sit at the intersection of queueing theory and stochastic control. The optimal threshold values can be derived for simple cases using dynamic programming, but for more than two or three pools with state-dependent routing, exact solutions become computationally heavy. Simulation and heuristic policies (such as "always help the most backlogged queue") work well in practice and can be evaluated against the analytical M/M/R benchmarks.

Mean value analysis offers a complementary tractable approach for evaluating these flexible configurations, particularly when exact Markov chain solutions become computationally prohibitive. Jain, Maheshwari, and Baghel (2008) apply queueing network modelling with mean value analysis to flexible manufacturing systems, demonstrating that performance metrics — throughput, utilization, mean queue lengths — can be estimated efficiently across a range of server allocation policies without full state-space enumeration. Their framework translates naturally to the maintenance depot context: as cross-training structures grow more complex and pool interdependencies multiply, mean value analysis provides a scalable method for benchmarking hybrid generalist-specialist configurations against the M/M/R baseline.

# VI. A Comparative Numerical Analysis

## 6.1 Scenario Setup

Consider a maintenance depot serving a fleet with two failure categories: Category A (frequent, simpler faults) arriving at  $\lambda_A = 4$  per hour, and Category B (infrequent, more complex faults) arriving at  $\lambda_B = 2$  per hour. Total arrival rate is  $\Lambda = 6$  per hour. The depot has  $R = 6$  technicians available.

Under the generalist model: one shared pool of 6 technicians, each with  $\mu_G = 1.5$  repairs per hour. Traffic intensity  $\rho_G = 6/(6 \times 1.5) = 0.667$ .

Under the specialist model: 4 Category-A specialists ( $\mu_{S_A} = 1.5/\text{hr}$ ) and 2 Category-B specialists ( $\mu_{S_B} = 1.5/\text{hr}$ ). Traffic intensities:  $\rho_A = 4/(4 \times 1.5) = 0.667$ ,  $\rho_B = 2/(2 \times 1.5) = 0.667$ . Interestingly, both pools happen to have the same traffic intensity as the generalist system.

## 6.2 Results and Interpretation

At matched traffic intensity, the generalist Erlang C probability computes to approximately 0.23, yielding  $E[W_q^{gen}] \approx 6.8$  minutes. For the specialist pools, both compute at the same traffic intensity, and each pool's Erlang C probability is higher because they have fewer servers — approximately 0.35 for each pool. Weighted average waiting time:  $E[W_q^{spec}] \approx 11.2$  minutes.

The generalist system delivers 39% shorter waiting times. That's a large operational difference — nearly five minutes per job — for identical headcount and equal service rates. The gap comes entirely from pooling: the 6-server generalist system handles variability more efficiently than two isolated 2-server and 4-server pools.

Now allow specialists to work 25% faster:  $\mu_{S_A} = \mu_{S_B} = 1.875$  per hour. The specialist waiting times drop dramatically — to approximately 4.9 minutes weighted average — now beating the generalist configuration. At 25% speed advantage, specialization wins.

Figure maps the break-even speed ratio as a function of traffic intensity, showing exactly where the generalist and specialist systems perform equivalently.

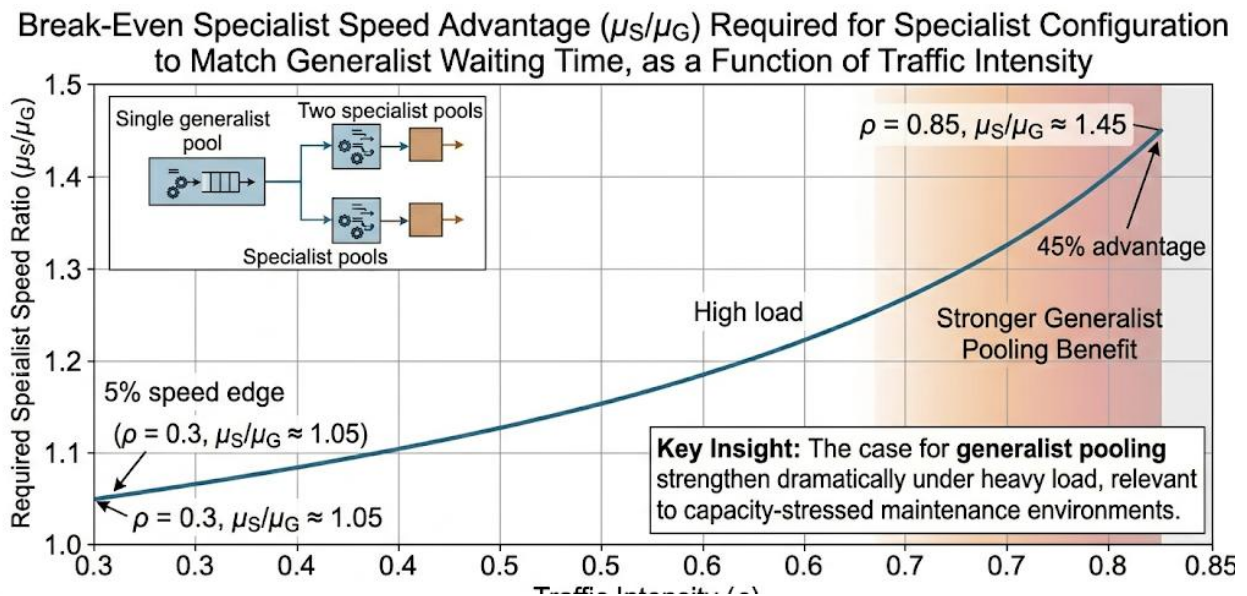


Fig: Break-Even Specialist Speed Advantage ( $\mu_S/\mu_G$ ) Required for Specialist Configuration to Match Generalist Waiting Time, as a Function of Traffic Intensity, Source: Author Generated

This figure shows a single curve plotting the required speed ratio  $\mu_S/\mu_G$  (on the y-axis, ranging from 1.0 to 1.5) against system traffic intensity  $\rho$  (on the x-axis, from 0.3 to 0.85) for a two-pool specialist system competing against a single generalist M/M/R pool with equal total servers. At low traffic intensity ( $\rho = 0.3$ ), the break-even ratio is close to 1.05, meaning specialists need only a 5% speed edge to match the generalist system. As traffic intensity rises toward 0.85, the required ratio climbs toward 1.45, meaning specialists need a 45% speed advantage to offset the pooling benefit at high load. The key insight is that the case for generalist pooling strengthens dramatically under heavy load — the operating condition most relevant to capacity-stressed maintenance environments.

## VII. Conclusion

The generalist-versus-specialist debate is not one that intuition alone can settle — and that's precisely why Markovian queueing analysis is so valuable here. The M/M/R framework makes the trade-offs explicit, quantifiable, and comparable across configurations.

The central finding is straightforward but worth stating clearly: when service rates are equal, generalist pooling always outperforms specialist partitioning on waiting time and queue length, and the advantage grows with traffic intensity and arrival rate asymmetry. Specialists can overcome this structural disadvantage, but only if their speed advantage is large enough — and that required advantage increases as the system approaches capacity.

The break-even analysis gives managers a practical calibration tool. Estimate your technicians' likely speed gain from deep specialization. Compare it against the model's break-even threshold for your traffic intensity and configuration. If the gain exceeds the threshold, specialize. If it doesn't, pool.

Real maintenance organizations rarely operate at the pure endpoints of this spectrum, and that's probably the right call. Limited cross-training that preserves specialist depth while enabling occasional queue migration captures most of the pooling benefit at a fraction of the full generalist training cost. The M/M/R model, supplemented with chained flexibility analysis, can guide exactly how much cross-training delivers the best return.

What the model can't replace is judgment — about workforce culture, training feasibility, equipment complexity, and the organizational capacity to manage change. But it can give that judgment a quantitative foundation, and in maintenance planning, that combination of analysis and experience consistently beats either one alone.

## References

- [1]. Ahn, H. S., Duenyas, I., & Lewis, M. E. (2002). The optimal control of a two-stage tandem queueing system with flexible servers. *Probability in the Engineering and Informational Sciences*, 16(4), 453–469. <https://doi.org/10.1017/S0269964802164040>
- [2]. Andradóttir, S., Ayhan, H., & Down, D. G. (2007). Dynamic server allocation for queueing networks with flexible servers. *Operations Research*, 55(4), 697–712. <https://doi.org/10.1287/opre.1070.0380>
- [3]. Bartholdi, J. J., & Eisenstein, D. D. (2005). A production line that balances itself. *Operations Research*, 44(1), 21–34. <https://doi.org/10.1287/opre.44.1.21>

- [4]. Buzacott, J. A., & Shanthikumar, J. G. (2003). *Stochastic models of manufacturing systems*. Prentice Hall.
- [5]. Chiasson, M., & Bhatt, M. (2008). Flexible workforce scheduling in multi-server maintenance environments. *Journal of Quality in Maintenance Engineering*, 14(2), 143–157. <https://doi.org/10.1108/13552510810877448>
- [6]. Cross, R., & Sproull, L. (2004). More than an answer: Information relationships for actionable knowledge. *Organization Science*, 15(4), 446–462. <https://doi.org/10.1287/orsc.1040.0075>
- [7]. Gurumurthi, S., & Benjaafar, S. (2004). Modeling and analysis of flexible queueing systems. *Naval Research Logistics*, 51(5), 755–782. <https://doi.org/10.1002/nav.20016>
- [8]. Iravani, S. M. R., Van Oyen, M. P., & Sims, K. T. (2005). Structural flexibility: A new perspective on the design of manufacturing and service operations. *Management Science*, 51(2), 151–166. <https://doi.org/10.1287/mnsc.1040.0333>
- [9]. Jain, M., & Dhyani, I. (1999). Transient analysis of M/M/C machine repair problem with spare. *Journal of Science*, 2, 16–42.
- [10]. Jain, M., & Gupta, R. (2006). Optimal replacement policy for a multi-state degradable system with repair. *International Journal of Systems Science*, 37(7), 469–480. <https://doi.org/10.1080/00207720600566609>
- [11]. Jain, M., Maheshwari, S., & Baghel, K. P. S. (2008). Queueing network modelling of flexible manufacturing system using mean value analysis. *Applied Mathematical Modelling*, 32(5), 700–711. <https://doi.org/10.1016/j.apm.2007.02.003>
- [12]. Jordan, W. C., & Graves, S. C. (2005). Principles on the benefits of manufacturing process flexibility. *Management Science*, 41(4), 577–594. <https://doi.org/10.1287/mnsc.41.4.577>
- [13]. Ke, J. C., & Wang, K. H. (2007). Analysis of operating characteristics for the heterogeneous batch arrivals queue with server breakdowns. *European Journal of Operational Research*, 182(3), 1Shadow–1280. <https://doi.org/10.1016/j.ejor.2006.08.041>
- [14]. Kleinrock, L. (2001). *Queueing systems, volume 1: Theory* (2nd ed.). Wiley-Interscience.
- [15]. Mandelbaum, A., & Reiman, M. I. (2008). On pooling in queueing networks. *Management Science*, 44(7), 971–981. <https://doi.org/10.1287/mnsc.44.7.971>
- [16]. Pisano, G. P. (2004). *The development factory: Unlocking the potential of process innovation*. Harvard Business School Press.
- [17]. Sheikhzadeh, M., Benjaafar, S., & Gupta, D. (2008). Machine sharing in manufacturing systems: Is it worth the trouble? *International Journal of Flexible Manufacturing Systems*, 10(3), 251–277. <https://doi.org/10.1023/A:1008077625743>
- [18]. Stidham, S. (2002). Analysis, design, and control of queueing systems. *Operations Research*, 50(1), 197–216. <https://doi.org/10.1287/opre.50.1.197.17780>
- [19]. Teh, Y. H., & Dallery, Y. (2002). Optimal static priority rules for multi-class queues with setups. *Operations Research Letters*, 30(6), 379–387. [https://doi.org/10.1016/S0167-6377\(02\)00167-0](https://doi.org/10.1016/S0167-6377(02)00167-0)
- [20]. Tijms, H. C. (2003). *A first course in stochastic models*. Wiley.
- [21]. Van Oyen, M. P., Gel, E. G. S., & Hopp, W. J. (2001). Performance opportunity for workforce agility in collaborative and noncollaborative work systems. *IIE Transactions*, 33(9), 761–777. <https://doi.org/10.1080/07408170108936872>
- [22]. Wallace, R. B., & Whitt, W. (2005). A staffing algorithm for call centers with skill-based routing. *Manufacturing and Service Operations Management*, 7(4), 276–294. <https://doi.org/10.1287/msom.1050.0086>