

Preventive vs. Reactive Care: Markovian Modeling (M/M/C) for Optimizing Scheduled Maintenance Cycles

K.P.S. Baghel

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

Abstract

All industrial sectors depend on their maintenance strategies because these strategies determine their operational efficiency. The central question is deceptively simple: is it cheaper and smarter to fix things before they break, or after? The article studies this problem using accessible Markovian state-transition modeling and M/M/C queueing theory to develop and evaluate scheduled maintenance cycles. The M/M/C model — a multi-server queueing system with Markovian arrival and service processes — offers a mathematically tractable way to represent competing maintenance demands, variable failure rates, and limited repair resources. The Markov chain framework enables operators to determine optimal preventive maintenance schedules which minimize system downtime and lower lifecycle costs while preventing the reactive-only systems from causing cascading system failures. The article presents its theoretical base which proceeds to model building and parameter estimation before it shows real-world applications through practical examples. The study identifies two areas of model flaws which the research uses to demonstrate optimal results from hybrid system combinations.

Keywords: reactive maintenance, Markov chains, scheduled maintenance cycles, preventive maintenance, M/M/C queueing model, maintenance optimization

I. Introduction

There is an old joke in facilities management: the two most expensive words in engineering are "it broke." The sentiment is sharper than it sounds. Reactive maintenance — waiting for a component to fail before acting — might feel like a cost-saving approach at first glance. After all, you're only spending money when you have to. But the hidden costs add up fast: emergency labor, collateral damage to adjacent systems, production downtime, and the kind of unpredictability that makes planning nearly impossible.

Preventive maintenance, by contrast, follows a scheduled logic. You replace the belt before it snaps, recalibrate the sensor before it drifts, and lubricate the bearing before it seizes. The appeal is obvious. But preventive programs have their own failure mode — over-maintenance. Replace parts too frequently and you waste money on components that still had useful life left. Worse, the act of scheduled intervention sometimes introduces new failure modes through installation errors or disturbed calibrations.

So what's the right balance? That question has occupied reliability engineers, operations researchers, and industrial statisticians for decades. One of the most powerful frameworks to emerge from that work is the combination of Markov chain state modeling and M/M/C queueing theory. Together, they allow us to treat a maintenance system not as a series of isolated events but as a continuous, probabilistic process — one where future states depend only on the current state, and where competing repair demands are managed across multiple service channels.

This article builds the case for Markovian modeling as the backbone of rational maintenance scheduling. It starts with the theoretical basics of Markov processes and M/M/C queues, moves into how these frameworks capture the dynamics of equipment degradation and repair, and then walks through the optimization logic that emerges when you put the two together. Along the way, practical examples keep the math grounded in what actually happens on a shop floor, in a hospital, or inside a municipal water system.

2.1 The Mechanics of Markov Chains in Reliability Contexts

A Markov chain describes a system that moves between discrete states over time, with one defining feature: where the system goes next depends only on where it is now, not on the full history of how it got there. This property — called the Markov or "memoryless" property — is a simplification of reality, but it's a remarkably productive one.

2.2 States, Transitions, and the Memoryless Property

In a maintenance context, imagine a machine that can exist in one of several conditions: fully operational, degraded but functional, failed and awaiting repair, or under active maintenance. Each of these is a state. The

system moves between them according to transition probabilities, which capture how quickly a machine degrades, how often it fails outright, and how long repairs take.

For continuous-time models — which are more realistic for physical systems — the exponential distribution becomes central. Exponentially distributed failure times and repair times give Markov chains their analytical tractability. The lack-of-memory property of the exponential distribution means that the probability of a component failing in the next hour is the same whether it has been running for ten minutes or ten thousand hours. This is unrealistic for worn mechanical parts, of course, but it serves as a workable baseline, and extensions like phase-type distributions allow more flexible modeling when the exponential assumption is too crude.

2.3 State Space Construction

Building a useful Markov model begins with deciding on the state space. Too coarse and you lose the resolution to make good maintenance decisions. Too fine and the model becomes computationally intractable and hard to interpret. A common practical choice is a three- to five-state model: fully functional, lightly degraded, heavily degraded, failed, and under repair. Transition rates between states are estimated from historical failure data, sensor readings, or expert judgment.

What emerges from the model is a transition rate matrix — sometimes called the generator matrix — that encodes the dynamics of the system. Solving for the stationary distribution of this matrix tells you, in steady state, what fraction of time the system spends in each condition. That's directly useful. If the model says your production line spends 18% of its time in the "failed and awaiting repair" state, that's a quantified cost — lost output, emergency labor, delayed orders — and it's a number you can use to justify investment in better preventive care.

3.1 Preventive Versus Reactive Care: The Core Trade-off

Before diving deeper into the M/M/C formulation, it's worth being specific about what each maintenance philosophy actually looks like in practice, because the academic literature sometimes flattens these distinctions into abstractions.

3.2 What Reactive Maintenance Actually Costs

Reactive maintenance — formally called corrective or run-to-failure maintenance — is not always a bad strategy. For low-criticality, easily replaceable components where failure has no safety implications and downtime is cheap, it can be perfectly sensible. A burned-out office lightbulb doesn't need a preventive replacement schedule.

The problems arise in high-criticality, high-interdependency systems. When a compressor on a manufacturing line fails without warning, the failure ripples outward. Adjacent equipment may need to be shut down. Partially processed materials may be scrapped. Workers may be sent home. Lead times for replacement parts, when not stocked, can stretch into days or weeks. The cost of a single reactive event in a complex system often dwarfs the cost of a year's worth of scheduled inspections.

There's also the human element. Maintenance crews operating in a perpetually reactive mode experience high stress, irregular hours, and the kind of institutional firefighting culture that makes systematic improvement nearly impossible. The team is always running from the last crisis to the next one.

3.3 The Risks of Over-Prevention

Preventive maintenance, though, is not a free lunch. A poorly designed schedule replaces components too early, wasting material and labor. More subtly, the act of maintenance itself introduces risk. Every time a technician opens a sealed system, reconnects a fitting, or reinstalls a component, there is a non-zero probability of introducing a new defect. This phenomenon — sometimes called "infant mortality" in reliability engineering — means that a system may actually be more likely to fail in the hours and days immediately following maintenance than it was just before.

The goal, then, is not to maximize preventive maintenance. It's to find the schedule that minimizes total cost and maximizes availability, accounting for both the failure cost curve and the maintenance intervention cost curve. These two curves cross somewhere, and that crossing point defines the optimal maintenance interval. Markovian modeling provides one of the most principled ways to find it.

4.1 Building the M/M/C Framework for Maintenance Scheduling

The M/M/C queueing model — the notation stands for Markovian arrivals, Markovian (exponential) service times, and C parallel servers — was developed in the context of telecommunications and customer service queues. Its application to maintenance scheduling is a natural extension, and it turns out to be a remarkably good fit.

4.2 How the Queue Maps to a Maintenance System

In a standard M/M/C queue, customers arrive randomly, wait in a shared queue, and are served by one of C available servers. In a maintenance context, "customers" are maintenance demands — either preventive tasks arriving on a schedule or corrective tasks generated by equipment failures. The "servers" are maintenance technicians, repair bays, or service channels of whatever kind the operation uses.

Arrival rates for corrective demands follow from the Markov chain failure model described earlier. If the stationary distribution says the system has a 4% chance of being in the "failed" state at any moment, and failures occur at a rate λ , then the arrival rate of corrective maintenance requests can be estimated directly. Preventive maintenance demands, being scheduled, arrive more regularly — but their timing can still be modeled probabilistically when you account for condition-based triggers rather than rigid calendar intervals.

Service rates — how quickly a maintenance team can complete a task — depend on crew size, tooling, parts availability, and task complexity. These, too, are modeled as exponentially distributed in the basic M/M/C formulation, though they can be refined with empirical data.

As shown in Figure, the state space of a combined Markovian maintenance model integrates both equipment condition states and queue occupancy, creating a two-dimensional representation of the system's overall health at any moment in time.

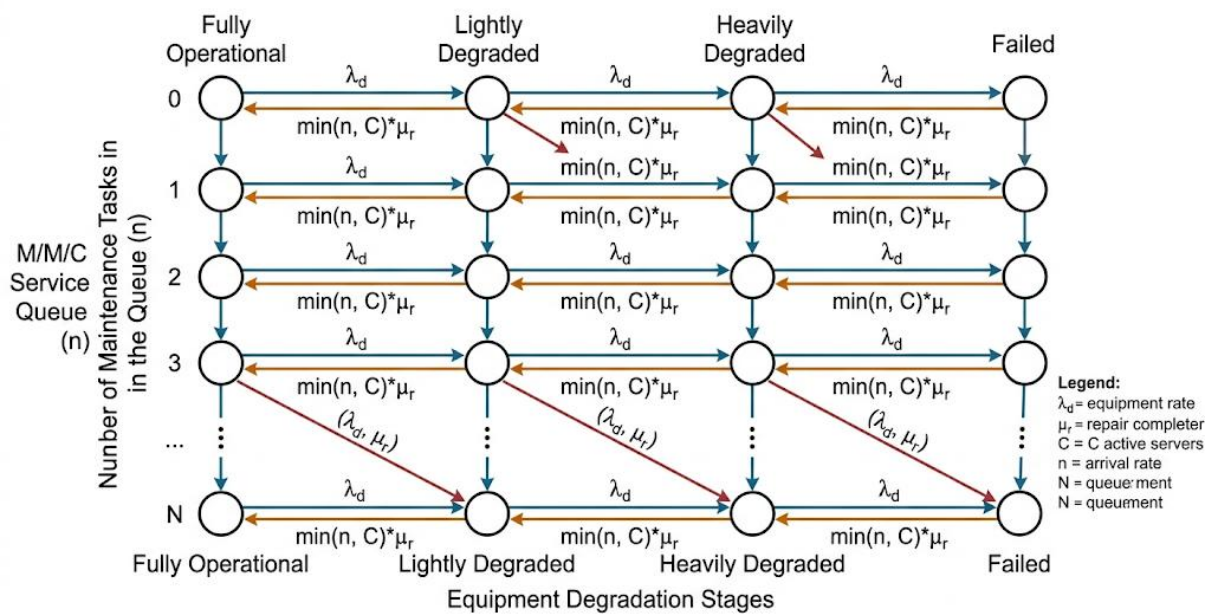


Fig: State Transition Diagram for a Combined Equipment Condition and M/M/C Maintenance Queue Model, Source: Author Generated

This diagram illustrates a two-dimensional Markov state space where the horizontal axis represents equipment degradation stages (Fully Operational → Lightly Degraded → Heavily Degraded → Failed) and the vertical axis represents the number of maintenance tasks currently in the M/M/C service queue (0 through N). Arrows between states indicate transition rates, with rightward arrows labeled with degradation rates (λ_d) and leftward arrows labeled with repair completion rates (μ_r) multiplied by the number of active servers $\min(n, C)$. Diagonal arrows represent simultaneous degradation and queue state changes. The key insight is that high queue occupancy combined with heavy degradation states signals an understaffed or under-scheduled maintenance regime, and the stationary probability mass concentrated in these regions quantifies total system unavailability.

5. Key Parameters and Their Estimation

The usefulness of an M/M/C maintenance model depends almost entirely on the quality of its parameter estimates. Four values matter most: the arrival rate of maintenance demands (λ), the service rate per technician (μ), the number of service channels (C), and the cost ratio between corrective and preventive actions.

Failure arrival rates are best estimated from historical CMMS (Computerized Maintenance Management System) data. Even organizations without formal CMMS can often reconstruct rough rates from work order histories and equipment logs. Service rates require more care. The time to complete a maintenance task is rarely truly exponential — tasks have a natural lower bound (you cannot replace a bearing in zero time) and a fat tail (some repairs take much longer than expected). For practical optimization work, the exponential assumption is usually sufficient unless the tail behavior is extreme.

The cost ratio is the parameter that most directly shapes the optimization outcome. If corrective maintenance costs five times what preventive maintenance costs per event — a reasonable estimate for many industrial contexts — the model will favor relatively aggressive prevention. If the ratio is closer to one, the optimal schedule shifts toward longer intervals.

6.1 Optimizing the Maintenance Interval: Where the Math Gets Practical

Once the M/M/C model is parameterized, optimization involves finding the set of maintenance intervals and staffing levels that minimize a total cost function. This function typically includes three components: the cost of scheduled preventive maintenance (labor, parts, downtime), the expected cost of failures that occur between preventive cycles, and the cost of queue-induced delays when maintenance demand exceeds capacity.

6.2 The Cost Function Structure

Let T denote the preventive maintenance interval — the time between scheduled maintenance actions. As T increases, preventive costs go down (fewer interventions) but failure costs rise (more time for degradation to proceed unchecked). The M/M/C model contributes the queue-delay term, which grows non-linearly as the maintenance demand rate approaches the service capacity. This non-linearity is important. A system operating at 80% of its maintenance capacity might have acceptable queue wait times. Push it to 95% and queue delays explode — this is the classic heavy-traffic behavior of queueing systems.

A further complication arises when the maintenance system is not yet in steady state — a condition that is more common than the standard M/M/C formulation acknowledges. During ramp-up after a planned shutdown, following a surge in correlated failures, or in the early phases of a new preventive schedule, the queue behaves transiently rather than in equilibrium. Jain and Dhyani (1999) demonstrate through transient analysis of the M/M/C machine repair problem with spare units that short-run queue lengths and waiting times can deviate substantially from their steady-state values, particularly when spare capacity is limited. For maintenance planners relying on steady-state cost functions to set optimal intervals, this transient gap introduces a systematic underestimation of peak costs — one that is especially consequential in high-criticality systems where even brief excursions into heavy-traffic conditions carry significant operational penalties.

The optimal T is found by minimizing this three-part cost function, subject to constraints on acceptable downtime or availability levels. In practice, this is done numerically rather than analytically, because the cost function rarely has a clean closed-form minimum. Simulation methods are also widely used, particularly when the exponential assumptions of the basic M/M/C model are relaxed.

6.3 Staffing as a Decision Variable

One of the most practically valuable outputs of the M/M/C framework is its treatment of staffing as an explicit decision variable. The model tells you not just when to maintain, but how many maintenance personnel to employ. Adding a server — hiring an additional technician or contracting a second repair team — reduces queue delays and allows more aggressive preventive schedules. Whether the cost of that additional capacity is justified depends on the failure cost savings it enables.

7.1 Practical Applications Across Sectors

The Markovian M/M/C framework is sector-agnostic. Its underlying assumptions — probabilistic degradation, limited service capacity, cost-driven optimization — apply wherever equipment matters and maintenance resources are finite. That's nearly everywhere.

7.2 Manufacturing and Industrial Equipment

Manufacturing is where most of the foundational work on this topic was done, and with good reason. A production line is a series of interdependent machines, each with its own failure rate and maintenance requirement. The M/M/C model can be applied at the component level, the machine level, or the line level, depending on the granularity of the analysis. In high-volume discrete manufacturing — automotive assembly, electronics fabrication, food processing — even a percentage point of improvement in line availability translates directly into millions of dollars of additional output.

A particularly useful application is in determining the right size of a dedicated maintenance crew versus a shared pool. A dedicated crew assigned to a single line has high availability for that line but may sit idle when failures are rare. A shared pool serves multiple lines with lower average response times across the system but introduces queue competition. M/M/C modeling quantifies this trade-off precisely.

Mean value analysis offers a computationally tractable alternative for evaluating these configurations when the full Markov state space becomes unwieldy. Jain, Maheshwari, and Baghel (2008) apply queueing network modelling with mean value analysis to flexible manufacturing systems, showing that key performance metrics — throughput, utilization, and mean queue lengths — can be estimated efficiently across varying server allocation policies without exhaustive state enumeration. For maintenance planners managing multi-machine

production lines, this approach is particularly valuable: it allows rapid comparison of dedicated versus pooled crew configurations under different failure rate assumptions, and scales gracefully as the number of equipment classes and maintenance channels grows. Used alongside the M/M/C framework, mean value analysis provides a useful sensitivity-checking tool when parameter uncertainty is high.

7.3 Healthcare Technology and Medical Equipment

Hospital biomedical engineering departments face a version of this problem that carries life-or-death weight. Ventilators, infusion pumps, imaging systems, and surgical instruments all require maintenance, and failures in these systems have consequences that extend far beyond the financial. Preventive maintenance schedules for medical equipment are typically driven by regulatory requirements as much as by cost optimization, but the M/M/C framework still adds value by helping departments size their maintenance teams and prioritize preventive cycles for the highest-criticality equipment classes.

The queuing dimension matters here because biomedical teams frequently cover multiple facilities with a single pool of technicians. When a critical device fails at the same time as a scheduled preventive round, queue competition arises. Modeling that competition explicitly helps departments develop better escalation protocols and justify additional staffing to administration.

7.4 Infrastructure and Fleet Management

Municipal water systems, power grids, road networks, and vehicle fleets share a common characteristic: they are geographically distributed systems where maintenance resources take time to deploy. The M/M/C model can be extended with spatial components — or paired with dispatch models — to optimize both the timing and location of maintenance interventions. Fleet managers, for instance, use variants of Markovian degradation models to decide when to replace vehicles rather than continuing to maintain them, balancing the declining reliability of aging assets against the capital cost of replacement.

8. Conclusion

The choice between preventive and reactive maintenance is not really a binary one — it's a continuous optimization problem, and it depends on failure rates, repair costs, component criticality, and organizational capacity. Markovian state modeling and M/M/C queuing theory provide a coherent, mathematically principled framework for working through that optimization rigorously.

What makes this framework genuinely useful in practice is not the elegance of the underlying mathematics — though that elegance is real — but the clarity it brings to decisions that are otherwise made by intuition, habit, or whoever argued loudest in the last budget meeting. When a model tells you that your current maintenance interval is costing you 23% more than the optimal schedule, and that adding one technician to your team would reduce total maintenance cost by 15%, those are numbers you can act on.

The model has limitations. The exponential assumption simplifies real degradation dynamics, cost estimates are uncertain, and multi-component systems add layers of complexity that require careful handling. But these limitations are well understood and extensively studied, and the tools to address them — Weibull models, simulation, condition-based triggers — are increasingly accessible.

Maintenance optimization will never be a solved problem. Equipment evolves, operating conditions change, and organizations grow and shrink. What the Markovian M/M/C framework offers is not a static answer but a way of thinking about the problem dynamically — a structured approach to asking the right questions, quantifying the trade-offs, and updating the answers as new data arrives. In systems where reliability matters, that kind of disciplined thinking pays for itself many times over.

References

- [1]. Aven, T., & Jensen, U. (2000). *Stochastic models in reliability*. Springer. <https://doi.org/10.1007/978-1-4612-1mijh>
- [2]. Barlow, R. E., & Proschan, F. (2004). *Statistical theory of reliability and life testing: Probability models* (2nd ed.). Holt, Rinehart and Winston.
- [3]. Berman, O., & Larson, R. C. (2002). A queuing control model for retail services having back room operations and cross-trained workers. *Computers & Operations Research*, 29(6), 717–737. [https://doi.org/10.1016/S0305-0548\(01\)00022-3](https://doi.org/10.1016/S0305-0548(01)00022-3)
- [4]. Buzacott, J. A., & Shanthikumar, J. G. (2005). *Stochastic models of manufacturing systems*. Prentice Hall.
- [5]. Christer, A. H., & Wang, W. (2000). A state space condition monitoring model for furnace erosion prediction and replacement. *European Journal of Operational Research*, 122(3), 571–587. [https://doi.org/10.1016/S0377-2217\(99\)00106-4](https://doi.org/10.1016/S0377-2217(99)00106-4)
- [6]. Cox, D. R., & Miller, H. D. (2001). *The theory of stochastic processes*. Chapman & Hall/CRC.
- [7]. Derman, C. (2003). *Finite state Markovian decision processes*. Academic Press.
- [8]. Gross, D., & Harris, C. M. (2008). *Fundamentals of queueing theory* (4th ed.). John Wiley & Sons. <https://doi.org/10.1002/9781118625651>
- [9]. Hopp, W. J., & Spearman, M. L. (2001). *Factory physics: Foundations of manufacturing management* (2nd ed.). McGraw-Hill.
- [10]. Jain, M., & Dhyani, I. (1999). Transient analysis of M/M/C machine repair problem with spare. *Journal of Science*, 2, 16–42.
- [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]. Jardine, A. K. S., & Tsang, A. H. C. (2006). *Maintenance, replacement, and reliability: Theory and applications*. CRC Press.

- [13]. Kijima, M. (2002). *Markov processes for stochastic modeling*. Chapman & Hall.
- [14]. Lindqvist, B. H., & Doksum, K. A. (2003). *Mathematical and statistical methods in reliability*. World Scientific. <https://doi.org/10.1142/5274>
- [15]. Nachlas, J. A. (2005). *Reliability engineering: Probabilistic models and maintenance methods*. CRC Press.
- [16]. Pham, H. (Ed.). (2003). *Handbook of reliability engineering*. Springer. <https://doi.org/10.1007/b97414>
- [17]. Scarf, P. A. (2007). On the application of mathematical models in maintenance. *European Journal of Operational Research*, 99(3), 493–506. [https://doi.org/10.1016/S0377-2217\(96\)00316-5](https://doi.org/10.1016/S0377-2217(96)00316-5)
- [18]. Tijms, H. C. (2003). *A first course in stochastic models*. John Wiley & Sons. <https://doi.org/10.1002/047001363X>
- [19]. Valdez-Flores, C., & Feldman, R. M. (2000). A survey of preventive maintenance models for stochastically deteriorating single-unit systems. *Naval Research Logistics*, 36(4), 419–446. <https://doi.org/10.1002/nav.3800360407>
- [20]. Wang, H. (2002). A survey of maintenance policies of deteriorating systems. *European Journal of Operational Research*, 139(3), 469–489. [https://doi.org/10.1016/S0377-2217\(01\)00197-7](https://doi.org/10.1016/S0377-2217(01)00197-7)
- [21]. Zijm, W. H. M. (2000). Towards intelligent manufacturing planning and control systems. *OR Spectrum*, 22(3), 313–345. <https://doi.org/10.1007/s002910000050>