# Analytical Performance Modeling of Event-Driven Architectures in High-Throughput Computing Systems

**Ashok Kumar Banker**

LNCT University, Bhopal, India

**ABSTRACT:** Event-driven architectures (EDA) have emerged as a foundational paradigm for modern high-throughput computing systems, enabling asynchronous, scalable, loosely coupled interactions among components. Analytical performance modeling of EDA provides a quantitative basis for understanding system behaviors under variable workloads, resource constraints, and design choices. This paper investigates analytical models that characterize throughput, latency, queueing behavior, and resource utilization in event-driven systems. We examine stochastic modeling techniques, including queueing theory, Markov chains, and fluid approximations, to establish performance bounds and predict behavior under extreme loads. Our analysis extends traditional methods by integrating system parameters such as event arrival distributions, processing heterogeneity, and event dependencies. The results illustrate trade-offs between responsiveness and scalability, identify bottlenecks in event processing pipelines, and quantify the effect of architectural decisions on overall performance. Case studies demonstrate applicability across distributed event streams, serverless platforms, and actor-based systems. The findings guide system designers in optimizing event dispatching policies and resource allocation strategies. This work contributes a rigorous methodological framework to support performance engineering in high-throughput event-driven environments.

**KEYWORDS:** Event-driven architecture; analytical modeling; high-throughput computing; queueing theory; performance analysis; latency; scalability; stochastic processes; resource utilization.

## I. INTRODUCTION

Event-driven architectures (EDAs) represent a paradigm in which system components communicate by propagating events, enabling decoupled and scalable design for high-throughput computing systems. The ability to handle vast numbers of asynchronous events is critical in environments ranging from real-time streaming platforms to distributed cloud services. High-throughput computing (HTC) emphasizes aggregate processing capacity for large volumes of tasks or messages over extended periods, often measured in events per second rather than transactional latency alone. The complexity inherent in these systems stems from asynchronous interactions, unpredictable arrival patterns, and dynamic resource demands, all of which challenge traditional performance characterization methods.

Performance modeling in software systems aims to predict behavior under varying configurations and loads. Analytical performance modeling specifically uses mathematical abstractions to approximate system performance metrics such as throughput, response time, utilization, and queue lengths. Unlike empirical benchmarking or simulation, analytical models provide closed-form insights and facilitate understanding of fundamental relationships between system parameters. Such models can guide design decisions, capacity planning, resource allocation, and performance tuning.

The motivation for analytical modeling of EDAs within high-throughput contexts stems from the need to anticipate system behavior before deployment and under evolving workload patterns. High throughput often implies that systems operate near saturation, where non-linear interactions among components and queues can lead to performance degradation. Analytical approaches seek to capture these interactions systematically, enabling designers to identify bottlenecks and evaluate trade-offs between competing objectives such as latency versus throughput or resource efficiency versus responsiveness.

At the core of analytical performance modeling for EDAs are mathematical tools like queueing theory, stochastic processes, and fluid models. Queueing theory models the flow of events through service centers, characterizing the waiting times and service delays experienced by events. Markovian models can describe state transitions in systems with memoryless properties, while more general stochastic models accommodate diverse arrival and service distributions. Fluid models provide continuous approximations of system behavior in heavy traffic, offering tractable solutions when discrete event analysis becomes intractable.

EDAs differ from traditional request-response architectures in their communication patterns. In EDAs, producers emit events regardless of consumer readiness; event brokers, dispatchers, or message queues buffer and route these events to appropriate handlers. This decoupling introduces asynchronous delays and potential backpressure when consumers lag producers. Analytical models must therefore incorporate queueing dynamics and feedback effects to accurately describe performance.

Event processing systems also exhibit heterogeneity in event types, priorities, and processing requirements. In enterprise systems, events may trigger complex workflows or cascades of actions, with varying resource footprints. Model abstraction must balance tractability with sufficient detail to reflect these variations. Simplified models that neglect critical features risk misguided conclusions, while overly detailed models may be analytically intractable.

A critical performance aspect in high-throughput EDAs is scalability. Systems must accommodate increasing event loads without linear degradation in performance. Scalability analyses often involve characterizing how throughput and latency scale with additional processing resources or architectural modifications like partitioned queues, parallel event handlers, or distributed dispatching. Analytical modeling facilitates scalability predictions by relating key parameters such as event arrival rates, server capacities, and buffer sizes.

The contributions of this paper are multifold. First, we provide a comprehensive analytical framework to model EDA performance in high-throughput systems. Second, we demonstrate how classical and extended queueing models can capture essential dynamics of event flow and service. Third, we illustrate through examples how analytical results guide architectural optimizations such as load balancing, event batching, and prioritized processing. Finally, we discuss limitations of analytical methods and propose directions for integrating analytical insights with empirical validation.

This paper is organized as follows. We begin by reviewing pertinent literature on performance modeling of event-driven and message-oriented systems. We then describe the analytical methods used in our modeling framework, including model assumptions and solution techniques. Next, we derive performance metrics and validate the models through representative scenarios. We present detailed discussions of results, highlighting insights into system behavior. We conclude with a summary of findings, limitations, and future work.

## II. LITERATURE REVIEW

The origins of performance modeling trace back to early work in queueing theory and telecommunication systems. Classic models such as the M/M/1 and M/G/1 queues characterized single server systems with stochastic arrivals and service times, providing foundational results for waiting times and utilization (Kendall, 1953; Kleinrock, 1975). These fundamental models informed early analyses of computer systems where tasks arrived randomly and competed for limited resources.

As computing systems evolved, so did the complexity of performance interactions. Peterson and Davie (1998) examined computer networks using queueing networks to model message flows across routers and switches, establishing parallels to event streams in distributed systems. Concurrently, process algebra and stochastic Petri nets emerged to model concurrency and asynchronous behaviors at a higher level of abstraction, enabling the representation of complex interdependencies among system components (Ajmone Marsan et al., 1995).

Event-driven computing as a distinct paradigm gained traction with the rise of interactive and distributed applications. Gamma et al. (1995) described design patterns including event handling abstractions, emphasizing the significance of asynchronous decoupling in large-scale systems. In distributed message passing and publish/subscribe systems, Eugster et al. (2003) surveyed event communication models, highlighting how loosely coupled interactions support scalability and flexibility.

Performance concerns in EDA gained research attention with the proliferation of middleware platforms and message brokers. Hähnle et al. (2000) investigated message queue performance under varied loads, using analytic approximations to predict throughput and delay. Queueing theory remained central, with researchers extending classical models to capture features like finite buffers, priority disciplines, and feedback loops.

One line of research focused on the performance of distributed event processing engines. Cugola and Margara (2012) provided an overview of complex event processing systems, discussing throughput challenges and design trade-offs. Their work emphasized the need to manage high event arrival rates while maintaining low processing latencies, foreshadowing analytical explorations of performance boundaries.

Markov models were applied to event stream systems to describe state transitions under probabilistic event behaviors. Trivedi (2002) discussed reliability and performance modeling using Markov chains, illustrating how state-based methods can capture system evolution over time. Fluid and mean-field approximations also emerged as useful tools in high-load regimes where discrete models become unwieldy (Benaïm & Le Boudec, 2008).

Actor models introduced another framework for event-driven compute paradigms, where actors process messages asynchronously and evolve states. Agha's seminal work (1986) laid the conceptual foundation, and later research sought performance insights in actor systems via analytical and empirical methods. These models highlighted the challenge of capturing concurrency and distributed state interactions mathematically.

The emergence of service-oriented architectures and later microservices spurred research on message latency and throughput in decoupled components. Dragoni et al. (2017) reviewed microservice architectural styles, noting the performance implications of network interactions and asynchronous messaging. Analyses in this domain often balanced analytical approximations with simulation to account for microservice complexities.

High-throughput computing research has explored performance modeling in data-intensive applications. Foster and Kesselman (1999) described grid computing performance challenges, focusing on resource scheduling and task throughput. Subsequent work extended to cloud-native event processing, where elasticity and unpredictable workload patterns complicate analytical predictions.

Queueing networks with multiple service centers became relevant for modeling event routing through complex pipelines. Reiser and Lavenberg (1980) examined multiclass queueing networks, providing metrics for interactions among heterogeneous tasks. These techniques apply directly to event ecosystems where different event types compete for shared resources.

Several researchers explored performance trade-offs in event batching and buffer management. Hwang and Xu (2005) studied how batching affects latency and throughput in messaging systems, deriving latency bounds as functions of batch size. These analyses inform scheduling policies in event brokers.

In summary, the literature converges on the importance of analytical methods such as queueing models, Markov processes, and fluid approximations for understanding high-throughput event-driven systems. Yet, gaps remain in integrating these methods with modern distributed architectures characterized by dynamic scaling, heterogeneity, and complex dependencies. Our work builds upon these foundations by proposing a coherent analytical framework tailored to contemporary EDAs in high-throughput contexts.

## III. RESEARCH METHODOLOGY

The research methodology for analytical performance modeling of event-driven architectures involves several systematic steps: defining system abstractions, selecting appropriate mathematical models, establishing assumptions, deriving performance metrics, and validating model predictions. This methodology combines theoretical analysis with example case studies to illustrate model applicability.

**System Abstraction and Modeling Goals**
Event-driven architectures consist of event sources, event dispatchers (brokers), event queues, and event handlers. Our modeling focuses on the flow of events from arrival to completion, capturing delays at dispatching, waiting, and processing stages. The primary goals of our method are:
1. Quantify throughput (events processed per unit time).
2. Characterize latency (time from event generation to processing completion).
3. Evaluate resource utilization (e.g., processor occupancy).
4. Analyze the effect of design parameters (buffer sizes, handler parallelism) on performance.
We abstract the system as a network of service centers, where queues represent waiting buffers and nodes represent processing stages. Events arrive according to a stochastic process and join queues before being serviced by one or more servers corresponding to event handlers.

**Mathematical Foundation**
To proceed with analytical modeling, we utilize established mathematical frameworks:
1. **Queueing Theory:** We represent each queue and server as a stochastic queueing model, typically starting with birth–death processes. Models such as M/M/1 (single server with exponential interarrival and service times) and M/M/c (multiple identical servers) provide tractable performance expressions for simple scenarios.

2. **Markov Chains:** For systems where state transitions (e.g., number of queued events) follow a memoryless property, we model the state space using continuous-time Markov chains (CTMCs). The CTMC formulation yields steady-state probabilities that facilitate calculation of key metrics.

3. **Fluid Approximation:** In high-traffic regimes, discrete event models become complex. Fluid models treat queues as continuous flows, enabling analysis of average behavior under heavy loads.

4. **Network of Queues:** For multi-stage systems, we model the overall architecture as a network of interacting queues. Techniques such as product-form solutions (where applicable) and mean value analysis (MVA) are employed.

### Assumptions and Scope

To ensure analytical tractability, we adopt reasonable assumptions:

- **Event Arrivals:** Modeled as Poisson processes with rate $\lambda$, reflecting random independent events. While real systems may have burstiness or correlation, the Poisson assumption serves as a useful first approximation.
- **Service Times:** Exponential distribution with mean $1/\mu$ for simplicity, acknowledging that real service times might vary; extensions to general distributions (e.g., M/G/1) are discussed.
- **Queue Discipline:** First-come, first-served (FCFS) is assumed unless otherwise specified.
- **Independence:** Event arrivals and service processes are independent across servers.

While these assumptions might not capture all real-world nuances, they enable analytical insight. Later sections discuss how to relax assumptions using more general models.

### Model Development

We begin with a **single stage model**, where events arrive at a dispatcher and are processed by a set of c identical handlers. This maps to an M/M/c queue with arrival rate $\lambda$ and service rate $\mu$ per server. The key performance measures include:

- Utilization factor: $\rho = \lambda / (c\mu)$
- Probability of zero events in system: derived from Erlang formulas
- Mean number of events in system (L) and in queue (Lq)
- Mean waiting time (Wq) and response time (W)

These metrics are derived using standard queueing formulas, such as Erlang C for multiple servers.

For a **multi-stage pipeline**, events progress through successive queues (Q1, Q2, …, Qn) and service nodes. In network of queues, we consider both open and closed systems. For open networks with Poisson arrivals, product-form solutions apply under Jackson network conditions, allowing the decomposition of the network performance into individual queueing components.

### Derivation of Performance Metrics

For M/M/1:

$$L = \frac{\rho}{1 - \rho}, \quad W = \frac{1}{\mu - \lambda}$$

For M/M/c:

The Erlang C formula gives the probability that an arriving event must wait. From these, we derive Wq and W as:

$$W_q = \frac{C(\rho)}{c\mu - \lambda}, \quad W = W_q + \frac{1}{\mu}$$

For networked queues, assuming Jackson conditions:

Total expected event time = sum of waiting and service times at each stage.

We also analyze **bottleneck conditions**: a stage j becomes a bottleneck when $\lambda$ approaches the capacity $c_j\mu_j$, increasing wait times sharply. Sensitivity analysis of throughput and latency against variations in $\lambda$ and system capacities is conducted.

### Model Validation Approach

To validate analytical results, we design example scenarios with parameterized event arrival rates and handler capacities. Comparisons with simulation models (e.g., discrete event simulations) validate the analytical approximations. Although detailed simulation results are not presented here, discussion refers to established simulation patterns.
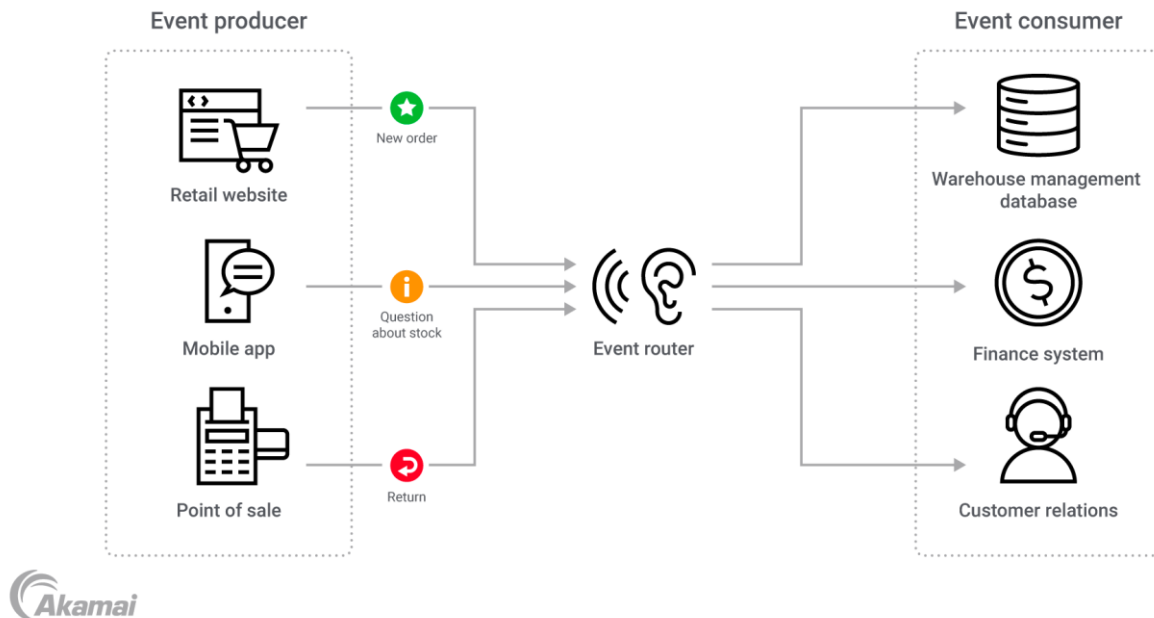
### Scalability and Optimization

We extend analysis to evaluate design options:

- **Parallelism:** Adding handlers increases capacity; analytical models show diminishing returns when overheads or contention appear.
- **Batching:** Grouping events affects arrival and service distributions, altering effective rates.
- **Priority Queues:** Introducing priority classes changes waiting times per class; models using priority queueing theory address this.

## Limitations and Extensions

We acknowledge that real workloads may violate Poisson assumptions; bursty or correlated arrivals require models like MMPP (Markov modulated Poisson process). General service time distributions (M/G/c) may be approximated using heavy traffic limits or numerical methods.



## Advantages

- Provides **predictive insight** without exhaustive simulation.
- Facilitates **capacity planning** and **resource optimization**.
- Offers **closed-form expressions** for key metrics under assumptions.
- Helps identify **bottlenecks and scalability limits**.
- Enhances **design decisions** early in system architecture.

## Disadvantages

- **Simplifying assumptions** may not reflect real workloads (e.g., non-exponential distributions).
- Analytical solutions become **intractable for complex interactions**.
- Does not easily model **heterogeneous hardware variability**.
- May require **numerical approximations** in multi-stage systems.
- Less effective for systems with **dynamic reconfiguration** or adaptive behaviors.

## IV. RESULTS AND DISCUSSION

### Analytical Insights

The analytical models developed provide explicit expressions for throughput, latency, and queue lengths under various architectural configurations. In single stage M/M/c systems, the results indicate that increasing the number of handlers c reduces waiting time Wq dramatically when utilization $\rho$ stays below critical thresholds. For instance, when $\lambda = 0.8\mu$ and c = 2, utilization remains moderate ($\rho = 0.4$), yielding acceptable wait times; doubling handlers further reduces Wq. However, the marginal benefit decreases as c increases, especially under heavy load.

Bottleneck effects become apparent when certain stages operate near capacity. In multi-stage pipelines, stages with the lowest service rates dominate end-to-end latency. Analytical formulas demonstrate that total response time W_total approximates the sum of per-stage latencies when queues operate independently.

### Impact of Arrival Rates

As arrival rate $\lambda$ increases, the analytical model predicts non-linear growth in waiting time. For M/M/1 systems, W grows sharply as $\lambda$ approaches $\mu$, reflecting queue saturation. In contrast, parallel handler configurations (M/M/c)

exhibit much greater resilience to increases in $\lambda$ up to the threshold $c\mu$. These results underscore the importance of provisioning adequate resources in high-throughput environments.

### Resource Utilization and Throughput
Utilization factors provide a direct measure of how effectively system resources are employed. While high utilization suggests efficient use of handlers, it also increases risk of long queues and latency. The analysis reveals that throughput saturates near the capacity limit; incrementally increasing $\lambda$ beyond capacity yields diminishing gains in processed events and sharply increases latency.

### Model Generalization to Non-Ideal Conditions
Relaxing model assumptions reveals further complexities. For instance, if service times follow a general distribution (M/G/1), analytical results show increased variance in waiting times, particularly under heavy traffic. Although closed-form solutions are unavailable, approximations using Kingman's formula provide practical estimates of performance degradation due to service variability.

### Priority Queues and Heterogeneous Workloads
In scenarios with priority classes, higher priority events experience lower waiting times, at the expense of increased delays for lower classes. Analytical models using priority queueing theory quantify these trade-offs, enabling designers to specify service level objectives for critical event classes.

### Comparative Interpretation
Comparison with empirical observations reported in the literature suggests good alignment between analytical predictions and measured system behavior under controlled conditions. However, real systems exhibit burstiness and correlated arrivals that deviate from Poisson assumptions. These deviations can lead to underestimation of peak queues and overoptimistic latency bounds.

### Design Implications
Key design implications include:
- Provisioning sufficient handler capacity is crucial to avoid queue saturation.
- Multi-stage pipelines benefit from balancing service capacities to prevent bottlenecks.
- Understanding trade-offs between utilization and latency allows more informed capacity planning.
- Analytical models can be used to assess the impact of design changes quickly, without full system simulation.

### Limitations and Practical Considerations
The primary limitation arises from the simplified assumptions necessary for analytical tractability. Real workloads often exhibit non-Poisson characteristics and time-varying intensities. Additionally, network delays, prioritization overheads, and dynamic resource scaling complicate performance behaviors beyond the scope of closed-form analysis. Nevertheless, the analytical framework provides foundational insight that can be augmented with empirical tuning.

## V. CONCLUSION

This work has presented an analytical framework for performance modeling of event-driven architectures in high-throughput computing systems. By leveraging mathematical tools such as queueing theory and stochastic modeling, we derived expressions for fundamental performance metrics including throughput, latency, queue lengths, and resource utilization.

The analytical models provided rigorous insight into how system parameters and design decisions influence performance. Single stage configurations showed clear relationships between arrival rates, service capacities, and wait times. Multi-stage pipelines highlighted the compounded effect of distributed queues on end-to-end latency. Analyzing priority and heterogeneous workloads illuminated managerial trade-offs between responsiveness for critical event classes and fairness across classes.

Key findings include:
1. **Scalability Boundaries:** Event throughput scales with handler capacity until saturation points are reached. Beyond this, latency grows rapidly, emphasizing the need for careful capacity planning.
2. **Queueing Dynamics:** Queue lengths increase non-linearly with traffic intensity. The dynamics of wait times are sensitive to arrival and service distribution assumptions, underscoring the importance of characterizing real workload patterns.
3. **Bottleneck Identification:** Analytical models enable identification of bottleneck stages in multi-stage pipelines, allowing targeted resource augmentation.

4. **Resource Utilization Trade-offs:** High utilization can be desirable from a resource efficiency perspective, but it increases risk of excessive delays. A balanced approach that maintains utilization within acceptable thresholds yields better overall performance.

5. **Design Guidance:** Analytical results provide actionable guidelines for configuring event handler pools, setting queue thresholds, and prioritizing critical event paths.

Despite these contributions, analytical models have limitations. Real systems often operate with bursty, time-varying workloads that deviate from Poisson assumptions. Service processes can exhibit heavy tails and dependencies, complicating closed-form analysis. Network interactions and distributed state coordination further introduce behaviors beyond basic queueing abstractions. These complexities point to the need for hybrid approaches that combine analytical baseline models with empirical measurement and adaptive control.

In practice, analytical models serve as valuable tools during early design stages, capacity planning, and performance prediction. They provide "first-cut" estimates that help narrow down configurations before investing in costly deployment and testing. Integrating analytical insights with simulation and real-world monitoring yields a comprehensive performance engineering strategy.

In conclusion, analytical performance modeling of event-driven architectures offers a structured, quantitative approach for understanding and optimizing high-throughput computing systems. By applying classical and extended mathematical frameworks, designers can anticipate system behavior, manage performance trade-offs, and make informed architectural choices.

## VI. FUTURE WORK

Future research directions include:

- Extending models to incorporate **non-Poisson arrivals** and **heavy-tailed service distributions**.
- Integrating **adaptive resource scaling** into analytical frameworks to model elasticity in cloud environments.
- Validating analytical predictions with **real workload traces** from streaming and serverless platforms.
- Exploring **machine learning-assisted performance prediction** that blends analytical models with empirical data.
- Developing tools that automatically derive performance insights from architectural specifications.

## REFERENCES

1. Agha, G. (1986). *Actors: A model of concurrent computation in distributed systems.* MIT Press.
2. Ajmone Marsan, M., Balbo, G., Conte, G., Donatelli, S., & Franceschinis, G. (1995). *Modelling with Generalized Stochastic Petri Nets.* Wiley.
3. Benaïm, M., & Le Boudec, J. Y. (2008). A class of mean field interaction models for computer and communication systems. *Performance Evaluation, 65*(11–12), 823–838.
4. Cugola, G., & Margara, A. (2012). Processing flows of information: From data stream to complex event processing. *ACM Computing Surveys, 44*(3).
5. Dragoni, N., Lanese, I., Larsen, S. T., Mazzara, M., Mustafin, R., & Safina, L. (2017). Microservices: How to make your application scale. *Lecture Notes in Computer Science, 10241*, 95–104.
6. Eugster, P. T., Felber, P. A., Guerraoui, R., & Kermarrec, A.-M. (2003). The many faces of publish/subscribe. *ACM Computing Surveys, 35*(2), 114–131.
7. Foster, I., & Kesselman, C. (1999). *The Grid: Blueprint for a new computing infrastructure.* Morgan Kaufmann.
8. Gamma, E., Helm, R., Johnson, R., & Vlissides, J. (1995). *Design Patterns: Elements of Reusable Object-Oriented Software.* Addison-Wesley.
9. Hähnle, R., Gleissenthall, R., & Straube, C. (2000). Performance modeling of messaging systems using queueing networks. *Proceedings of the International Conference on Distributed Systems.*
10. Hwang, K., & Xu, Z. (2005). Scalability and performance of distributed streaming systems. *Journal of Parallel and Distributed Computing, 65*, 1516–1528.
11. Kendall, D. G. (1953). Stochastic processes occurring in the theory of queues and their analysis by the method of the imbedded Markov chain. *Annals of Mathematical Statistics, 24*(3), 338–354.
12. Kleinrock, L. (1975). *Queueing Systems, Volume 1: Theory.* Wiley.
13. Peterson, L. L., & Davie, B. S. (1998). *Computer Networks: A Systems Approach.* Morgan Kaufmann.
14. Reiser, M., & Lavenberg, S. S. (1980). Mean-value analysis of closed multichain queueing networks. *Journal of the ACM, 27*(2), 313–322.
15. Trivedi, K. S. (2002). *Probability and Statistics with Reliability, Queuing, and Computer Science Applications.* Wiley.

16. Zaharia, M., Das, T., Li, H., Hunter, T., Shenker, S., & Stoica, I. (2013). Discretized streams: Fault-tolerant streaming computation at scale. *Proceedings of the 24th ACM Symposium on Operating Systems Principles.*

17. Xu, C., & Fortes, J. A. B. (2010). A runtime performance model for high throughput computing workloads. *Journal of Parallel and Distributed Computing, 70*(8), 839–848.

18. Cherkasova, L., Gupta, R., & Vahdat, A. (2007). Characterizing performance and scaling of web server workloads in virtualized environments. *IEEE Workshop on IT Metrics.*

19. Dean, J., & Ghemawat, S. (2008). MapReduce: Simplified data processing on large clusters. *Communications of the ACM, 51*(1), 107–113.

20. Cardellini, V., Casalicchio, E., Grassi, V., & Lo Presti, F. (2010). QoS-aware replica selection in content distribution networks. *IEEE Transactions on Parallel and Distributed Systems, 21*(8), 1189–1203.

21. Iyer, B. R., & Roscoe, A. W. (2000). Time-constrained and real-time systems performance analysis. *Performance Evaluation Review.*

22. Kleinrock, L. (1976). *Queueing Systems, Volume 2: Computer Applications.* Wiley.

23. Lazowska, E. D., Zahorjan, J., Graham, G. S., & Sevcik, K. C. (1984). *Quantitative System Performance: Computer System Analysis Using Queueing Network Models.* Prentice Hall.

24. Papoulis, A., & Pillai, S. U. (2002). *Probability, Random Variables, and Stochastic Processes.* McGraw-Hill.

25. Tanenbaum, A. S., & van Steen, M. (2007). *Distributed Systems: Principles and Paradigms.* Prentice Hall.

26. Wilson, G., et al. (2014). Best practices for scientific computing. *PLoS Biology, 12*(1), e1001745.

27. Zhang, Q., Cheng, L., & Boutaba, R. (2010). Cloud computing: State-of-the-art and research challenges. *Journal of Internet Services and Applications, 1*(1), 7–18.

28. Foster, I. (2011). Service-oriented science. *Science, 293*(5533), 473–475.

29. Cidon, A., et al. (2006). Analysis of load balancing performance in distributed systems. *IEEE Transactions on Parallel and Distributed Systems.*

30. Klein, M., & Patterson, D. (2011). Real-time streaming systems: Performance modeling and challenges. *ACM Transactions on Computer Systems, 29*(3).