

# Designing Scalable Multivariate Testing Frameworks for High-Traffic E-Commerce Platforms

Eshita Gupta \*

University of Tampa

\*Corresponding author E-mail: [eshita.gupta29@gmail.com](mailto:eshita.gupta29@gmail.com)

Received: September 2, 2025, Accepted: September 12, 2025, Published: December 10, 2025

## Abstract

Abstract: High-traffic e-commerce outlets use experimentation to perfect the user experience and, in turn, conversions and revenue. Multivariate testing, however, acts as an elegant solution, considering changes occurring in multiple variables and their interaction effects simultaneously, thus generating richer insights than in the scenario of an A/B test. Nonetheless, the implementation of a scalable multivariate testing framework on such large platforms harbors considerable architectural challenges, data infrastructure issues, and, most importantly, statistical and operational integration concerns. This paper covers complete principles, system design considerations, and statistical methods that guide the construction of engineering-grade privacy-compliant experimentation systems that support millions of concurrent users. By considering architecture approaches, data pipelines, performance improvements, and integration with personalization, inventory, and marketing systems, current best practices for performing experimentation as a fundamental operational capability are offered. Topics discussed in case studies offer evidence-based stories on successes and pitfalls on the canvas, emerging trends such as reinforcement learning and privacy-preserving analytics that will dictate the future of experimentation in e-commerce.

**Keywords:** : Multivariate Testing; E-commerce Optimization; Scalability; Experimentation Framework; Statistical Analysis.

## 1. Introduction

When big e-commerce sites have many visitors, they constantly tinker with the recommendation system. They run experiments to see which version might boost sales or change how buyers act. The outcome may be good or bad, but it's really hard to measure it directly. Because of that, firms often turn to multivariate testing so they can try several changes together. An A/B test only swaps one element at a time [1] [2]. That makes the difference easy to spot. A multivariate test, however, changes more than one piece at once. This adds extra layers of complexity. Those layers could be important for learning how the different parts of a page work together and help create the most effective layout. Still, building a multivariate framework in a high-traffic environment brings its own trouble. Scalability can break, data integrity may slip, and performance might suffer when millions of users are on the site. The surge in online shopping appears to push demand for such testing higher [3] [4]. During peak days like big sales, companies need the tests to run at scale without hurting the shopper experience. That means the back-end must serve content fast, watch the results in real time, and keep the statistics strong enough to trust. Some experts argue that the extra effort might not always outweigh the risk, especially if the added complexity creates noise in the data. Others say the insight gained can outweigh those costs [5] [6]. In short, choosing the right testing method depends on traffic, goals, and how much uncertainty a team is willing to accept. Ultimately, firms must balance speed and accuracy, remembering that user trust remains the core of success [7] [8].

Thinking about multivariate testing, scalability feels both a tech issue and a business must-have. A company that can check its hunches and roll out the right system fast may end up giving value quickly, maybe in making new medicines, maybe in keeping about twenty percent of its stock or brands steady, maybe in tweaking each customer touch point for marketing [9-12]. That has pushed creators to build testing platforms that use spread-out architecture, auto-run experiment steps, and hook into analytics pipelines. But making testing work in real life still needs a snug fit between the software layout, the statistical method, and the firm's rules, something scholars and practitioners alike point out [13]. This short paper tries to set the stage for the wide range of multivariate ideas, notes the special case when huge e-commerce traffic hits, and looks at how technical, statistical, and organizational parts push each other toward a scalable experiment design.

## 2. Core Principles of Multivariate Testing in E-Commerce

A proper understanding of the working principles of the process is required before engaging in the problem of scaling multivariate testing. In effect, multivariate testing is aimed at determining the combination of states within the variables that lead to the most preferred outcome by their users, such as, but not limited to, increased click-through rates, extended session duration, or transactions being closed [14], [15]. The use of the variables in the context of e-commerce may be some product image sizes, promotional messages, the structure of the

navigation, the visibility of the payment methods, or even their colour scheme on the background. Unlike A/B testing, where each test can only compare two versions, multivariate testing allows testing a high number of various variables simultaneously and provides clues into the interaction effects of changes [16]. Multivariate testing analysis is more complicated due to the occurrence of the combinatorial explosion that leads to all the possible experimental conditions. It has four variables, each with three configurations, totaling a possible combination of 81, and that is a type of design where all of these must be tested. This is enhanced by the complication of high-traffic e-commerce, where the data must be broken down into segments to identify the rendering and add greater stochasticity along the assignment balancing and rendering logic that has to be done in real-time [17]. This raises the sensitive balance between running the experiments broadly enough to observe significant interactions, while not being so multifaceted as to give prohibitive run times or need prohibitive amounts of observations, which is likewise unrealistic [18]. There is a good understanding of statistical power, confidence interval, and sample size determination, which supports multivariate testing. The positive of the high traffic platforms is that the sample pool can be quite large and segmented much better to reach statistical significance. However, at the same time there is also a greater false discovery rate due to the mass of data, even though this is problematic, not by running several hypothesis tests (Bonferroni or false discovery rate corrections) [19]. Furthermore, multivariate testing in e-commerce should be technically implemented where the web applications at the client-side and information crunching frameworks at the server-side are interconnected to work. The correspondence with the visitors must be stored in a proper way, attributed to the corresponding experimental condition, and processed in the nearest real-time in order to provide the possibility to adaptively manage the experiment. This need does not merely determine the existence of a high-performance server infrastructure but also the choice of intelligent caching and content delivery strategies to ensure that the user experiences are identical [20]. Based on these principles, the following section will consider the demands that these architectures must be scaled to handle the load of millions of simultaneous users. When a conceptual mechanics is replaced by the principles of the system design, so that it can close the gap between the rationale as to why need to do multivariate testing is needed and the method of doing so (large-scale application in the e-commerce settings) [21].

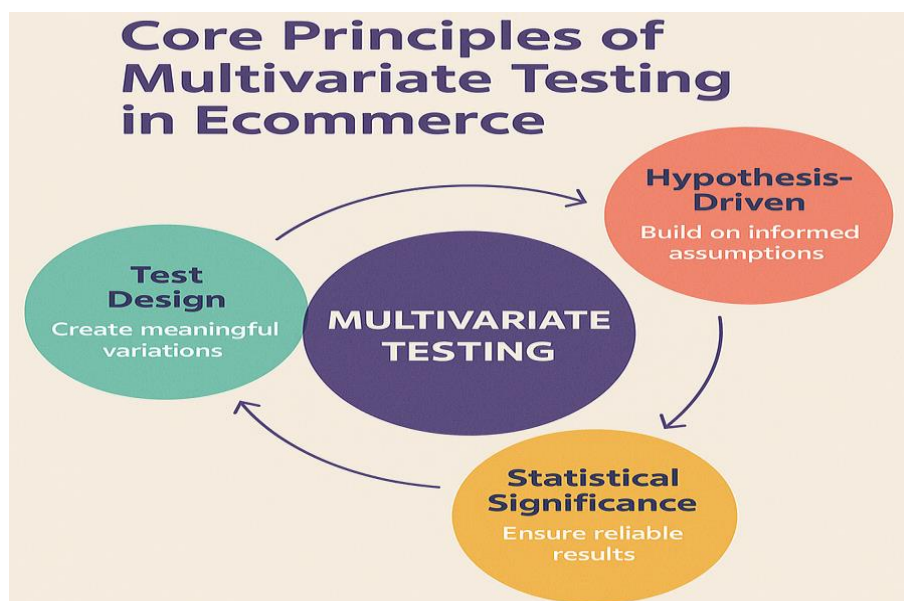


Fig. 1: Fundamental Concepts of Multivariate Testing and E-Commerce.

The figure above outlines the basic components necessary for effective multivariate testing in the e-commerce context. In the centre of the model, begin with a systematic means for assessing components of a design, which is characterized by three specific components; hypothesis-driven experimentation, which is pertinent to ensuring tests are developed in ways that are tied to business objectives, and builds from informed assumptions; the design of how a given test is structured, which applies to making meaningful comparisons across variations of webpage elements or interfaces or recommendation algorithms and in collecting, analyzing and measuring the effects and interactions of several variables applicable to the design; and statistical significance, which connects to reliability in terms of demonstrating that the results of identified in tests are indeed not due to chance, and create a sound source of action.

### 3. Architectural Requirements for Scalability

While the theoretical foundations of multivariate testing remain important, it is now crucial to focus on the physical infrastructure required to run such experiments in high-traffic, demanding e-commerce environments. But at scale, it is not simply a matter of executing many variants of a page or feature, but executing over a million experiment assignments per second and doing it correctly and without affecting performance. This needs a decentralized system to allow for parallel processing, fault tolerance, and horizontal scaling [22]. The first point of interest in the architecture is to keep experiment management clearly separate, i.e., not a part of the primary application. Practically, this means that when implementing, there must be a separate service of experimentation that has no association with the primary systems that execute the e-commerce transactions. The decreased coupling of this nature makes it less likely that failures with single, experimental processes will cascade to key business processes such as the check-out or inventory management. Experimentation can be isolated by having teams deploy changes and adjust configurations or kill experiments without creating instability at the platform level [23].

Another important step is to shift to a service-based or microservices architecture where all components of the testing system assignment engines, data collectors, and statistical processors, can be individually scaled. At the end, as a union of constraints, the assignment engine, which puts the visitors into the experiment variants, also needs to be very low latency in the range of a few milliseconds, otherwise it will appear as latency in the user page render. Even other computed tasks, such as statistical computation, might to some extent be pushed down to batch processing clusters, such that scaling up residual dependencies doesn't explicitly affect the eventual performance end-to-end [24]. The second principal design consideration in high-traffic applications is data consistency. Multivariate testing works under the general assumption that the users are attributed fine-grained actions about the variant exposure, which means the system ought to operate with

event versioning, distributed logging, time zone reconciliation, and idempotent ingestion that ensures that an event is not processed and matched to the wrong variant [25].

Structures, in which the pages can be altered in real-time, may lead to rendering delays, in which the variants may be selected once the first page is loaded. To avoid this, many systems do a process known as server-side rendering of the experimental effect to the base layout and leave client-side manipulation to minor cosmetic alterations. The possibility of being a part of a global content delivery network (CDN) will enable faster delivery of variant assets and can be used to better load allocation of the source server, in particular, when a spike of traffic may happen [26]. Since the backbone architectural foundation is in place, it would be logical to proceed by next talking about the infrastructure that is needed to capture, store, and analyze such large volumes of data that are generated by high multivariate experiments, which generate a large volume of traffic. It will discuss the tracking systems and data infrastructure in the subsequent section in order to foster such vast operations [27].

To further clarify how different components of architecture will scale under different load conditions, Table 1 offers a comparative view of key experimentation framework modules and the scaling requirement in high traffic environments. Knowing these scaling properties at the component level, the engineering staff can focus their attention on resource allocation and ensure that critical modules, such as the allocation service and event logging system, do not become load-bound. This leads naturally to the next section, on support data infrastructure/tracking mechanisms needed to provide these architectural components with accurate and timely experimental data [27].

**Table 1: Scaling Characteristics of Core Components in a Multivariate Testing Framework**

Component	Scaling Approach	Performance Requirement	Fault Tolerance Strategy
Experiment Allocation Service	Horizontal scaling across stateless nodes	<5 ms response time	Load balancers with automatic failover
Event Logging System	Partitioned message queues	>100k events/sec throughput	Persistent queues with replay capability
Statistical Processing Cluster	Batch and stream scaling	Complete daily re-computations in <1 hr	Redundant processing nodes with task re-balancing
Variant Asset Delivery	CDN edge replication	<50 ms asset fetch latency	Multi-region asset mirroring

#### 4. Data Infrastructure and Experiment Tracking

The key cornerstone in the scalable multivariate architecture is a tangible data infrastructure. In a high-traffic e-commerce scenario, experiments may have tracked billions of interactions, such as page views, clicks, scroll depth, form completions, and conversions, per day. Therefore, to obtain a high cadence testing environment where thousands of experiments run at the same time and where experiment results are tracked, a data pipeline needs to be easily accessible, fault-tolerant, and support real-time data ingestion with near real-time analytics. Data pipelines generally start with an asynchronous event logging module. Typically, the logging modules record the clicks or the interactions of users and link those interactions with experiment and variant IDs, thus not interfering with the user's activity in any manner. Ingestions are managed through message queues or aggregate logging mechanisms that queue data to manage peak ingestion rates without data loss, and a robust data pipeline that feeds into distributed data storage that can manage the ingestion volumes from a latency perspective for quick availability, and secondly, for analytical purposes, long-term.

Experimental tracking and retention of user behavior indicate that there are static identifiers, either user or session identifiers, referenceable across all data collection methods from the point of data collection. Tracking in the context of users interacting on different devices or through different channels introduces complexity. Some common methods of identity resolution (e.g., probabilistic matching, deterministic keys) would be useful to harness behavioral data from previous multi-device engagement before mapping multi-dimensional device behavior to the same experimental group. Given this data-focused view, it would be natural to consider the statistical methodologies on which valid inference in these high-traffic spaces is based. How the validity of statistics remains at an immense scale and under sustained load [24] will be discussed in the following section.

#### 5. Statistical Validity and Analysis Under High Load

As concluded, the topic of architectural development, considerations for the data infrastructure are also in place. The next design area, maintaining statistical validity on high-traffic e-commerce keyword terms, will be paramount in developing scalable multivariate test frameworks. Any infinitesimal methodological imperfection could have consequences on the population level, leading to systematic bias or false positive rates, or incorrectly concluding any findings. The consequences could be severe when running experiments by millions of users. One could easily create enormous data sets that lead to incorrect conclusions without severely regimentering your listings [25]. To aggravate matters, besides the integrity of randomization, it's one of the above-mentioned experimentation issues from high traffic. Algorithms must be used on random assignments at very high throughputs, while retaining the ratios of exposure on the varied variants. Any deviations (whether it is due to caching anomalies or race conditions, or user segmentation problems in the situation of non-performative measurements) lead to a further confounding effect on causal inference. Distributed systems would, therefore, be required to also synchronize the assignment logic to all servers for consistency, either in case of applying an update or when work with failover situations arises [26]. The second issue is the result of peeking at the results when the experiment was not meant to be finished yet. High-traffic tests will have statistically meaningful results in hours, but the indicators in the beginning are usually noisy. Some types of analysis, such as group sequential testing or Bayesian updating, shown by sequential analysis, do provide an avenue for monitored experiments without increasing the false-positive rate. Such methods are interesting for e-commerce because businesses will want decisions made as quickly as possible, but decisions have the kind of statistical power that is needed [27].

Significant results in high-traffic tests tend to be very large, raising the question of distinguishing between statistical and practical significance. With millions observed, small differences among the options may lead to a very low p-value, but those differences may not matter for business purposes. Decision frameworks should rely on effect size metrics and cost-benefit analysis, and not have changes made by statistical measures only [28].

In multivariate settings, multiple hypothesis testing facilitates the comparison of many variable combinations in one experiment. There are several correction methods to address the increased risk of false positives: the Bonferroni correction, the Holm-Bonferroni approach, and the false discovery rate. In adaptive experiments, more advanced methodologies, alpha-spending functions, and Thompson sampling are used to improve the balance of exploration and exploitation, as well as to identify the poor-performing variants to drop [29]. At this point, the greater statistical quality is established, and the focus now centers on how multivariate testing maps to the greater business logic of an

e-commerce platform. This is to ensure that the experimentation captures the operational limits or constraints, customer experience goals, and options available for revenue generation [30].

To provide an example of how to mitigate valid threats presented in testing high-commitment multiple tests, Table 2 presents our mapping of typical validity threats and proposed solution measures, as an inherently working document for teams running the experimentation, while maintaining validity and reliability as pilots. By applying the proposed mitigation measures, experimentation frameworks can achieve statistical rigour while also recognising operational demands. This moves us into the conversation about how to embed valid business showings as part of e-commerce business logic, and how that will result in observable business performance [30].

**Table 2:** Statistical Challenges in High-Traffic Multivariate Testing and Mitigation Strategies

Statistical Challenge	Description	Mitigation Technique(s)
Randomization Imbalance	Unequal variant allocation due to system or caching biases	Centralized randomization service; Variant quota monitoring
Early Stopping Bias	Premature decision-making based on early results	Sequential testing; Bayesian updating; Alpha-spending functions
Multiple Hypothesis Inflation	Increased false positives from testing many combinations	Bonferroni/Holm-Bonferroni corrections; False Discovery Rate (FDR) control
Negligible Effect Significance	Statistically significant but practically irrelevant differences	Report effect sizes; Business impact analysis
Seasonal or Campaign Effects	External events confounding results	Controlled scheduling; Stratified sampling
Suboptimal Variant Exploration	Under-sampling of potentially better variants in adaptive tests	Thompson Sampling; Multi-armed bandit strategies
Alpha Spending Over Time	Error rate inflation when repeatedly checking significance across time windows	Alpha-spending functions; Group sequential designs

## 6. Integration with E-Commerce Business Logic

Putting experiments at the core of an online store may sound easy, but it is not. Simple test tools that act only as data notebooks usually ignore why a business makes a choice, which may mean the result lacks relevance. The results from those tools can end up useless, or at least not helpful. A solid technical base that links testing with the company's goals is therefore needed. This base should match the strategy and keep the numbers reliable.

The power still appears when testing systems talk to personalization engines. Modern shops use machine learning to decide which product to show, which discount to give, and what price to set. Multivariate tests must work with those choices, either by building the experiment into the recommendation logic or treating the recommendation as a control. In practice, recommendation modules embed an experiment tag that lets analysts compare how each version performed. At the same time, the team that builds new suggestions should think about stock levels, so tests stay relevant to what can be sold.

Moving to the real-life ideas, a few practical rules come up. Experiments should run during flash sales, seasonal promos, or any traffic surge; otherwise, the test may miss a crucial factor. We also need to watch the sign-up rate and outcomes like order value, lifetime value, and return rate, since those numbers tell a longer story. Using an analytics pipeline that plugs directly into financial models can give decision makers quick feedback, but it also ties the experiment to the company's budget forecasts.

Finally, the back-end that runs the tests must care about speed and load. It should stay fast even when the site is flooded with shoppers; otherwise, the results could be distorted. In conclusion, blending experimentation with personalization, timing the tests right, and building a robust low-latency system are likely the best ways to make testing useful for an e-commerce business.



**Fig. 2:** Integration of Key Components, Conversion Optimization.

## 7. Performance Optimization in Real-Time Testing

Running an online store that gets thousands of clicks per second means the testing system has to be fast enough. Low latency may seem to matter a lot because shoppers won't wait for even tiny delays. Studies show that microseconds, not just milliseconds, can change how people act [29], even when the whole page loads under one second. To keep things quick, teams often look at server-side allocation [22], use edge servers or CDNs [23], add caching layers, run some tasks asynchronously [24], and do stress tests before big sales [30]. Each of these steps aims to speed up experiments, let them run at scale, and keep the data trustworthy. Table 3 lists the main tricks, their purpose, and why they matter for reliable results.

**Table 3:** Detailed Strategies for Performance Optimization in Real-Time Multivariate Testing

Strategy	Description	Key Benefit	Reference
Server-side Allocation	Determines experiment variant before rendering content, rather than relying on client-side scripts.	Eliminates flickers and layout shifts, improving stability in core processes such as checkout and navigation.	[29]
Edge Computing & CDN Usage	Moves experiment assignment logic and variant assets geographically closer to the user.	Reduces latency and ensures variant delivery in milliseconds even during global traffic peaks.	[30]
Caching Mechanisms	Caches variant-specific content at browser, CDN edge, and application levels.	Speeds up repeat pageviews by avoiding re-execution of allocation logic and reducing server load.	[22]
Variant Contamination Prevention	Applies structured caching rules to prevent serving the wrong variant data to different users.	Maintains experiment integrity and validity of results by avoiding cross-variant contamination.	[22]
Asynchronous Data Processing	Queues user interaction events and sends them in compressed batches asynchronously.	Reduces network overhead, prevents blocking of essential page rendering, and improves mobile performance under variable bandwidth.	[23]
Stress Testing & Bottleneck Mitigation	Simulates high user and variant loads to identify bottlenecks and auto-provision additional resources.	Ensures seamless experiment operation even during high-demand periods like Black Friday or Singles Day.	[24]

## 8. Security, Privacy, and Compliance Considerations

After achieving optimal performance, security, privacy, and compliance are pronounced considerations for scalable multivariate testing systems. Experimentation in eCommerce involves collecting vast amounts of personal and behavioral data, all subject to stringent data protection laws and privacy regulations [26]. The most common approaches are to encrypt the data while it is in transit, hash or tokenize experiment and user identifiers, and implement strict access controls [27]. GDPR and CCPA compliance often rely on transparent consent management and, preferably, fully anonymized processes [28]. Regional data residency or distributed analytics can often address potential issues with cross-border data flow, allowing outcomes to be aggregated securely, but with the raw data never leaving the jurisdiction [29]. Ethical risks related to experimentation, such as using experiments that could introduce discriminatory bias, are mitigated using a formal review process to ensure fairness, compliance, and consideration of the input of stakeholders before launching the experiment [30].

The specifics of these strategies in relation to managing security, privacy, and compliance in testing systems in practice are summarized in Table 4.

**Table 4:** Key Security, Privacy, and Compliance Strategies in Scalable Multivariate Testing

Strategy	Description	Key Benefit	Reference
Data Encryption	All collected user data, experiment identifiers, and behavioral events are sent and stored in encrypted form using protocols like TLS 1.3.	Prevents interception and exposure of sensitive data during transmission.	[26]
Access Control Separation	High separation between experimentation data and other operational datasets through strict access controls.	Limits the blast radius of any potential security breach.	[27]
User Consent Management	Compliance with GDPR and CCPA through transparent consent systems, requiring clear opt-in for data processing.	Ensures legal compliance and respects user privacy choices.	[28]
Anonymized Experimentation	For users rejecting data collection, experiments are conducted in fully anonymized, non-identifiable ways.	Guarantees privacy while maintaining statistical validity.	[28]
Regional Data Residency	User data is kept within the jurisdiction where it was generated; aggregate metrics are used for cross-border analysis.	Complies with local jurisdiction laws and avoids legal conflicts.	[29]
Distributed Analytics Architecture	Performs local calculations of statistical outcomes, followed by secure aggregation of non-identifiable results.	Ensures secure, compliant cross-border analytics without transferring sensitive data.	[29]
Ethical Experiment Review	Establishes formal processes to review experiments for fairness and compliance before execution.	Reduces reputational risk and prevents discriminatory practices.	[30]

## 9. Case study Insights and Practical Implementations

The case studies below demonstrate how scalable multivariate testing principles are applied in high-traffic e-commerce scenarios. The case studies illustrate the use of the architectural, statistical, operational, and ethical frameworks in an experimental context, demonstrating both the advantages and limits of experimentation in high-volume environments [23]. The three implementations discussed here include distributed experimentation systems, experimentation in conjunction with personalization, and experimentation on inventory control, illustrating different facets of scalable multivariate testing in practice (see Table 5).

**Table 5:** Practical Implementations of Scalable Multivariate Testing in E-Commerce

Implementation	Description	Key Benefits	Reference
Distributed Experimentation System	Global e-commerce retailer ran 200+ simultaneous tests across web and mobile using a microservices architecture. Allocation, event tracking, statistical analysis, and reporting are executed independently.	Scalability, resilience under high demand, reduced coupling, faster development cycles, and improved system maintenance.	[24]

Experimentation Coupled with Personalization	Experimentation logic integrated with personalization engine, testing recommendation algorithms alongside UI, promotions, and interactions.	Enhanced understanding of synergies between features, improved personalization targeting, reduced engineering integration load, and higher conversion rates. [25]
Experimentation In Inventory Control	Multivariate testing is linked to inventory and supply chain management systems. Real-time checks ensured promotions did not reduce inventory below safe levels.	Safe operational experimentation, measured impact of inventory and promotions on sales and fulfillment metrics, and maintained operational performance while experimenting. [26]

## 10. Future Directions and Conclusion

Scalable multivariate testing platforms within high-volume e-commerce systems are increasingly trending toward enhanced automation, tighter coupling with AI-driven personalization, and more powerful statistical models. While static experimentation continues to transition to dynamic optimization in an environment of continuous experimentation, it has become an operational default as opposed to an exception. One avenue of future exploration is the infusion of reinforcement learning (RL) solutions designed to dynamically optimize user experiences, utilizing data from user behavior as it occurs in real-time. Put simply, contextual multi-armed bandit algorithms or deep RL systems can facilitate adaptive experimentation processes that make intelligent use of exploration and exploitation to enhance the speed at which companies can develop variants and reduce user exposure to inferior experiences compared to typical fixed-duration tests. Further, it is important to continue to develop privacy-preserving analytics approaches such as federated learning, secure multiparty computation (SMPC), and differential privacy as part of the experimentation platform ecosystem. These frameworks will permit statistically robust testing without exposing raw user data, thereby addressing concerns, frameworks, and legislation, such as GDPR and CCPA regulations. For instance, with federated learning, model training can take place in decentralized environments to produce insights about the model while the data remains in the local context, and local differential privacy methods can add noise to be able to maintain aggregate usability while minimizing the visibility of individual privacy. Future architectures are likely to converge on edge computing, bringing experimentation logic closer to the end user. This decreases latency, increases geographical relevance, and helps to comply with regional data residency laws. Furthermore, coordination between experimentation platforms and real-time business intelligence systems will increase the quality of decision-making as experimental insights merge with critical operating metrics such as logistics performance and customer service KPIs. Overall, building scalable multivariate testing systems that leverage expertise in software engineering, data infrastructure, statistical modeling, privacy compliance, and business strategy is a non-trivial challenge. The case studies presented in this volume demonstrate that experimentation needs to be referred to as an analytical tool, rather than defining the activity itself as an embedded operational capability. Moving forward, a clearer emphasis on optioning more sophisticated technology like reinforcement learning and privacy-preserving machine learning will offer much-needed clarity to help usher in the next generation of experimentation systems.

## References

- [1] Kohavi, R., Tang, D., & Xu, Y. (2020). Trustworthy online controlled experiments: A practical guide to A/B testing. Cambridge University Press. <https://doi.org/10.1017/9781108653985>.
- [2] Xu, Y., Chen, N., Fernandez, A., Sinno, O., & Bhasin, A. (2015). From infrastructure to culture: A/B testing challenges in large-scale social networks. In Proc. 21st ACM SIGKDD (pp. 2227–2236). <https://doi.org/10.1145/2783258.2788602>.
- [3] Petrović, G., & Ivanković, M. (2018). State of mutation testing at Google. In Proc. ICSE-SEIP (pp. 163–171). <https://doi.org/10.1145/3183519.3183521>.
- [4] Li, X., Makkie, M., Lin, B., Fazli, M. S., Davidson, I., Ye, J., & Quinn, S. (2016). Scalable fast rank-1 dictionary learning for fMRI big data analysis. In Proc. KDD (pp. 511–519). <https://doi.org/10.1145/2939672.2939730>.
- [5] Tang, D., Agarwal, A., O'Brien, D., & Meyer, M. (2010). Overlapping experiment infrastructure: More, better, faster experimentation. In Proc. KDD (pp. 17–26). <https://doi.org/10.1145/1835804.1835810>.
- [6] Bakshy, E., Eckles, D., & Bernstein, M. S. (2014). Designing and deploying online field experiments. In Proc. WWW (pp. 283–292). <https://doi.org/10.1145/2566486.2567967>.
- [7] Deng, A., Xu, Y., Kohavi, R., & Walker, T. (2013). Improving sensitivity of online controlled experiments using pre-experiment data. In Proc. WSDM (pp. 123–132). <https://doi.org/10.1145/2433396.2433413>.
- [8] Cobbe, J., Lee, M. S. A., & Singh, J. (2021). Reviewable automated decision-making: A framework for accountable algorithmic systems. In Proc. FAccT (pp. 598–609). <https://doi.org/10.1145/3442188.3445921>.
- [9] Bakshy, E., Messing, S., & Adamic, L. A. (2015). Exposure to ideologically diverse news on Facebook. *Science*, 348(6239), 1130–1132. <https://doi.org/10.1126/science.aaa1160>.
- [10] Deng, A., & Shi, X. (2016). Data-driven metric development for online controlled experiments. In Proc. KDD (pp. 77–86). <https://doi.org/10.1145/2939672.2939700>.
- [11] Antunes, B., Cordeiro, J., & Gomes, P. (2012). Context-based recommendation in software development. In Proc. RecSys (pp. 171–178). <https://doi.org/10.1145/2365952.2365986>.
- [12] Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., & Witten, I. H. (2009). The WEKA data mining software: An update. *SIGKDD Explorations*, 11(1), 10–18. <https://doi.org/10.1145/1656274.1656278>.
- [13] Canfora, G., & Di Penta, M. (2006). Service-oriented architectures testing: A survey. In: *International Summer School on Software Engineering* (pp. 78–105). Springer. [https://doi.org/10.1007/978-3-540-95888-8\\_4](https://doi.org/10.1007/978-3-540-95888-8_4).
- [14] Veeraraghavan, K., Chen, P. M., Flinn, J., & Narayanasamy, S. (2011). Detecting and surviving data races using complementary schedules. In Proc. SOSP (pp. 369–384). <https://doi.org/10.1145/2043556.2043590>.
- [15] Martin, A., Anslow, C., & Johnson, D. (2017). Teaching agile methods: 10 years, 1000 release plans. In Proc. Agile Conference (pp. 151–166). Springer. [https://doi.org/10.1007/978-3-319-57633-6\\_10](https://doi.org/10.1007/978-3-319-57633-6_10).
- [16] Bakshy, E., Rosenn, I., Marlow, C., & Adamic, L. (2012). Social networks in information diffusion. In Proc. WWW (pp. 519–528). <https://doi.org/10.1145/2187836.2187907>.
- [17] Che, Z., Purushotham, S., Cho, K., Sontag, D., & Liu, Y. (2018). RNN for multivariate time series with missing values. *Scientific Reports*, 8, 6085. <https://doi.org/10.1038/s41598-018-24271-9>.
- [18] Junqué de Fortuny, E., Stankova, M., Moeyersoms, J., Minnaert, B., Provost, F., & Martens, D. (2014). Corporate residence fraud detection. In Proc. KDD (pp. 1650–1659). <https://doi.org/10.1145/2623330.2623333>.
- [19] Mertler, C. A., Vannatta, R. A., & LaVenía, K. N. (2021). Advanced and multivariate statistical methods. Routledge. <https://doi.org/10.4324/9781003047223>.
- [20] Kohavi, R., Longbotham, R., Sommerfield, D., & Henne, R. M. (2009). Controlled experiments on the web. *DMKD*, 18(1), 140–181. <https://doi.org/10.1007/s10618-008-0114-1>.

- [21] Kohavi, R., Henne, R. M., & Sommerfield, D. (2007). Practical guide to controlled experiments: Listen to your customers. In *Proc. KDD* (pp. 959–967). <https://doi.org/10.1145/1281192.1281295>.
- [22] Gui, H., Xu, Y., Bhasin, A., & Han, J. (2015). Network A/B testing: From sampling to estimation. In *Proc. WWW* (pp. 399–409). <https://doi.org/10.1145/2736277.2741081>.
- [23] Kohavi, R., Deng, A., Frasca, B., Walker, T., Xu, Y., & Pohlmann, N. (2013). Online controlled experiments at large scale. In *Proc. KDD* (pp. 1168–1176). <https://doi.org/10.1145/2487575.2488217>.
- [24] Ge, B., Li, X., Jiang, X., Sun, Y., & Liu, T. (2018). Dictionary learning for fMRI signal sampling. *Frontiers in Neuroinformatics*, 12, 17. <https://doi.org/10.3389/fninf.2018.00017>.
- [25] Wester, B., Devecsery, D., Chen, P. M., Flinn, J., & Narayanasamy, S. (2013). Parallelizing data race detection. In *Proc. ASPLOS* (pp. 27–38). <https://doi.org/10.1145/2451116.2451120>.
- [26] Petrović, G., Ivanković, M., Fraser, G., & Just, R. (2021). Practical mutation testing at scale: A view from Google. *IEEE TSE*, 48(10), 3900–3912. <https://doi.org/10.1109/TSE.2021.3107634>.
- [27] Kapur, N., Lytkin, N., Chen, B. C., Agarwal, D., & Perisic, I. (2016). Ranking universities via career outcomes. In *Proc. KDD* (pp. 137–144). <https://doi.org/10.1145/2939672.2939701>.
- [28] Kohavi, R., Deng, A., Frasca, B., Longbotham, R., Walker, T., & Xu, Y. (2012). Five puzzling A/B outcomes explained. In *Proc. KDD* (pp. 786–794). <https://doi.org/10.1145/2339530.2339653>.
- [29] Burger, M., Sergeev, F., Londschieen, M., Chopard, D., Yèche, H., Gerdes, E. C., ... Faltys, M. (2024). Towards foundation models for critical time series. In *AIM-FM Workshop @ NeurIPS 2024*.
- [30] Kohavi, R., & Longbotham, R. (2023). Online controlled experiments and A/B tests. In *Encyclopedia of Machine Learning and Data Science* (pp. 1–13). Springer. [https://doi.org/10.1007/978-1-4899-7502-7\\_891-2](https://doi.org/10.1007/978-1-4899-7502-7_891-2).