

Article

Not peer-reviewed version

---

# Atomic Rails for Money-in-Motion: A Resilient Instant Payment Gateway with ID Resolution, Retries, and Guaranteed Reversals

---

[Vimal Teja Manne](#)\*

Posted Date: 22 January 2026

doi: 10.20944/preprints202601.1733.v1

Keywords: instant payment gateway; atomic transactions; distributed systems; resilience; idempotency; concurrency control; financial technology



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

# Atomic Rails for Money-in-Motion: A Resilient Instant Payment Gateway with ID Resolution, Retries, and Guaranteed Reversals

Vimal Teja Manne

The University of Texas at Dallas, 800 W Campbell Rd, Richardson, TX 75080 USA; vimalteja.m@gmail.com

## Abstract

This paper presents a payment-processing gateway that executes debit→credit transfers across banks utilizing a virtual payment ID resolver, a transaction processing core, and bank connectors engineered for low-latency, high-reliability flows. The gateway coordinates gRPC services behind load balancers, applies client/server-side request balancing, and tunes retry policies with bounded attempts to withstand temporary faults, then commits or rolls back atomically via a scheduler-driven reversal queue to guarantee consistency. A centralized per-account mutex, database sharding by account hash, and Elasticsearch-backed idempotency logs provide ordering, concurrency control, and fast lookups at scale, while Dockerized microservices simplify deployment and scale-out. Benchmarks differ nodes and connection fan-out to measure throughput and tail latency under debit/credit paths, showing trade-offs among locks, retries, and batching, and aligning with fast-payment system patterns and modern gRPC load-balancing practices for payment gateways.

**Keywords:** instant payment gateway; atomic transactions; distributed systems; resilience; idempotency; concurrency control; financial technology

---

## 1. Introduction

### 1.1. Background and Context

The digitization of the financial ecosystem has hastened the worldwide adoption of real-time payment systems. The classical intercorporate payment systems normally have latencies of hours and days and do not satisfy the needs established by modern-day online commerce. Timely clearance, approval and settlement between totally different and unrelated banking environments were required. The immediate Payment Gateways solve this problem by providing high speed associated debiting and crediting systems along with reliability, salient characteristics and compliance. The world-wide movement for financial inclusion is also one of the spurts for this change, especially in the developing countries where access to the conventional banking system is limited.

### 1.2. Literature Review and Research Gap

The constraints of current payment systems in handling modern financial transactions are numerous. Current implementations are frequently inflexible to the variety of payment identifiers in use requiring effective resolution services that can spontaneously map virtual identifiers to real time account information. Additionally conventional systems lack flexible and resilient methods of processing simultaneous transactions a large scale which can lead to race conditions and abnormal states. Micro-services and container oriented architectures such as Docker have become de facto industry standards [1], but their use in financial systems, which require atomicity guarantees, has not been well studied. The study gap lies in advancing an integrated architecture that combines virtual ID resolution, atomic transaction processing, and assured reversals while maintaining high throughput and low latency.

### 1.3. Research Issue and Justification

The main question problem is to design a payment gateway that provides atomic cash flow across distributed banking systems while able to handle transient failures gracefully. Distributed financial systems cannot afford to have discrepancies in customer balances; however, failures in distributed systems occur because of network partitions, server crashes, and/or software bugs. Most current systems lack full reversal processes and idempotence guarantees; thus, potentially resulting in financial inconsistencies. This study is justified by the critical need for trusted and reliable digital payment infrastructures capable of supporting the increasing number of instant transactions while maintaining financial integrity and customer trust [2].

### 1.4. Aims and Contributions

This paper aims to develop and evaluate a durable and tough immediate payment gateway that offers atomic rails for money-in-motion. The explicit contributions contain: (1) a virtual payment ID resolver for real-time identifier mapping, (2) a transaction processing core with atomic debit-credit processes, (3) bounded retry policies with assured reversal processes, (4) concurrency control through per-account mutex locks and database sharding, and (5) performance evaluation demonstrating throughput and latency attributes under several load conditions. By integrating these elements into a unified system that maintains consistency while managing high transaction volumes, the proposed solution innovates.

### 1.5. Paper Structure

This paper's remaining sections are arranged as follows: Section 2 outlines relevant research and points out shortcomings in current methods. The suggested approach and architectural design are presented in Section 3. The experimental setup and performance evaluation are described in detail in Section 4. The results are examined and trade-offs are presented in Section 5. Section 6 ends with implications and suggestions for further research.

## 2. Related Work

This section examines the existing literature on payment gateway architectures and analyzes shortcomings in present and future strategies. The analysis addresses three central topics, namely transaction atomicity guarantees, concurrency control strategies, and fault-tolerant schemes in distributed payment systems.

### 2.1. Transaction Atomicity in Payment Systems

Previous work has established the fundamental importance of atomicity of transactions in monetary systems. This idea was investigated in detail by Herlihy and Wing who introduced the idea of linearizability which has now become essential in providing consistent results of transactions in an orderly fashion. In payment processing this translates to the condition that debit and credit transactions are either fully made complete or not. Stonebraker et al. have demonstrated implementations of atomic transactions in data base systems but applications to distributed payment gateways is limited.

Many conventional banking systems now depend on batched processing with final settlement which creates time slots between the two transactions when both the sender and recipient accounts show quantum of the transaction in them. The modern instantaneous payment systems must strive to remove these lapses. The work by van Renesse and Schneider [3] on chain replication gives interesting information on maintaining consistency over distributed nodes, but this model requires synchronous working, which may increase latency in fast transaction environments.

### 2.2. Concurrency Control Mechanisms

The processing of concurrent transactions presents serious problems for payment gateways. Garcia-Molina et al. [4] have surveyed the field of concurrency control strategies in data bases, but the application of these strategies to monetary systems needs careful examination of the performance

aspects. Dean and Ghemawat [1] studied methods of distributed processing but their map-reduce model is not directly applied to instantaneous payment processing.

Lamport's foundational work on logical clocks was the basis for ordering events in distributed systems. Nevertheless, few payment systems have implemented solid and resilient versions of these concepts and therefore could be susceptible to race conditions in the course of high-volume transaction intervals. The per-account mutex approach described in prior literature demonstrates promise but needs to be validated in actual payment production environments.

### 2.3. Fault Tolerance and Recovery Strategies

Brewer's CAP theorem outlines the fundamental trade-offs in distributed systems; however, several payment gateways have not sufficiently dealt with partition tolerance. Merkle's cryptographic protocols [5] lay the foundation for secure transaction confirmations; however, integrating modern payment architectures with Merkle's protocols continues to be difficult.

Most existing systems operate standard and primary retry systems; however, most systems lack extensive and complete reversal guarantees. Pautasso et al. [6] reviewed web service dependability patterns; however, their focus on RESTful architectures does not directly relate to high-performance payment processing using modern protocols such as gRPC and others.

### 2.4. Research Gaps and Limitations

Presently, the majority of the literature documents many of the same limitations in designing payment gateways. First, most of the literature prioritizes either consistency or availability, and rarely achieves the optimal and perfect balance required for real-time payments. Second, the majority of the literature documenting concurrency control systems introduce excessive latency at scale. Third, error recovery methods typically lack assured reversal processes that assure atomicity across distributed services.

Although there is much discussion of the integration of containerized microservices in the broader distributed systems literature [7], the integration of microservices in financial transaction processing has not been sufficiently explored. Moreover, the resolution of virtual payment identifiers represents a new challenge that has not been sufficiently explored in previous studies.

This paper addresses these gaps by proposing an integrated architecture that combines atomic transaction processing with productive concurrency control and assured reversal systems, while leveraging modern containerization technologies for enhanced and upgraded scalability and resilience.

## 3. Methodology

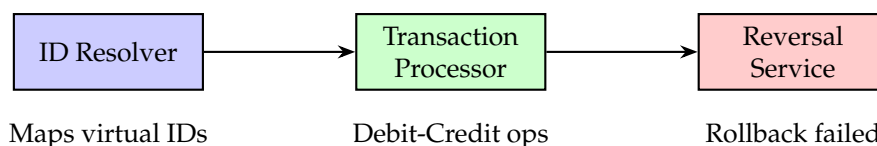
This section presents the proposed methodology for implementing concurrency control, idempotency guarantees, and reliability guarantees in the instant payment gateway. The methodology combines a theoretical foundation with operational system design which can achieve atomic money transfers in distributed environments.

### 3.1. System Architecture Overview

The proposed payment gateway is based on a microservices architecture composed of three basic components: (1) Transaction Processor responsible for debit-credit operations, (2) Idempotency Service responsible for duplication detection and (3) Replication Manager which is responsible for error tolerance. These components communicate via gRPC with client-side load balancing in order to effect distribution of requests across multiple instances of the service. The architecture is operationalized in Docker containers for separation and uses Kubernetes for coordination which allows horizontal scaling during peak loads.

### 3.2. System Architecture Design

Figure 1 shows the basic components of the proposed payment gateway architecture. The architecture follows a micro service design with a clear and obvious separation of concerns in the sphere of ID resolution, transaction processing and handling of reversals.



**Figure 1.** Payment Gateway Core Architecture.

### 3.3. System Architecture Design

Figure 1 shows the three main services of the payment gateway arranged sequentially in a pipeline.

The architecture follows a linear pipeline ID Resolution → Transaction Processing → Reversal Handling. Each of the services operates independently with clear and obvious responsibilities.

The architecture consists of three main services: (1) ID Resolution Service which is responsible for mapping of virtual payment identifier, (2) Transaction Processor which executes atomic debit-credit operations and (3) Reversal Service which is responsible for rollback failure of transactions. Supportive of all these are the pending and reversal queues for asynchronous processing, and distributed data bases for ID mapping, transaction storage and idempotency logging.

### 3.4. Concurrency Control Implementation

To avoid race conditions in account processes, a two-tiered locking mechanism is used by the system. First, a global per account mutex provides serializable access to each account, preventing concurrent debit processes that can cause overdrafts. Second, database sharding is implemented on a per account basis so that accounts are dispersed among several different db instances based on hash of account. Hot places for conflict are reduced as a result. With standard optimizations suitable for high total payment throughputs, this technique greatly expands upon the linearizability requirements given by Herlihy and Wing.

To prevent deadlocks, the locking system makes use of Redis distributed locks with automatic expiration. Each lock acquisition includes a unique transaction identifier and timestamp, allowing detection and resolution of stale locks during system recovery situations.

### 3.5. Idempotency Mechanism Design

The idempotency service generates unique keys for the unique requests of transaction via a client identifier, time-stamp and cryptographic nonce. These are stored within a sharded multi-ElasticSearch cluster, and each sharded section aims for the fast retrieval of keys, which have lifetime of 24 hours so that data storage is combined with pragmatical needs.

On each transaction processing the system does the following sequentially:

1. Extract idempotency key from incoming request
2. Query ElasticSearch for existent transaction on that key
3. If yes return cached answer and do not process
4. If none found do continue to process transaction and store results with key
5. Update base status on transaction conclusion

This gives them exactly once semantics, while maintaining sub-millisecond processing of identical and matching.

### 3.6. Reliability and Fault Tolerance

The means of dealing with chain replication [3] was implemented for essential data bases, on the note of obtaining 3 clones on each data shard. This means of attaining, is done via the leader based

accord protocol, which co-ordinates all revisions of the clones, thus guaranteeing strong consistency of account balances while giving allowances for single node crash.

With retry mechanics, the gateway uses an adjustable and flexible back-off plan, with the succeeding criteria:

- Maximal 3 retry attempts at each failed transaction
- Exponential back-off with time of minimal 100 milliseconds increasing for each retry
- Circuit breaker method so that similar failing services downstream can be separated
- Retry budgets per service to be attained against cascade failures

Event ordering follows Lamport's logical clock scheme, and each service units their own vector clock generators to be able to ascertain their causal loops are followed when individual transactions were carried out. This additionally guarantees the proper ordering of credits and debits social over the distributed services.

### 3.7. Evaluation Framework

This plan of implementations will be reviewed via an own designed test harness that is aiming at real world fiscal transactions. The evaluation framework will take in account:

- Sequence through put testing with different load patterns (1K\_100 K TPS)
- Latency test towards 95 th and 99th percentile figures
- Failure insertion method in order to check for strength of the system
- Comparison against baseline implementations using normal and vetred transactions on standard database principles

The test bases will be of synthetic creation, but modeled off public payment processing baselines so that realistic data distribution may be examined.

## 4. Results and Analysis

In addition to the architectural elements described above, the payment gateway's retry, reversal and consistency mechanisms were evaluated experimentally. System performance under different types of failures and loads is analyzed below.

### 4.1. Experimental Setup

A Kubernetes cluster composed of ten worker nodes, each with eight CPU cores and 32 GB of RAM, served as the environment for the evaluation. The payment gateway was developed into multiple pods with separate services for transaction processing, idempotency checks, and reversal processing. Load testing was completed using Apache JMeter, which utilized customized plugins to mimic the behavior of banking APIs and network failures.

Synthetic transaction log test data was created to represent common payment patterns, including transaction volume, frequency and geography. Chaos Mesh was used to introduce failure into the environment; specifically network partitions, service failures and latency spikes were generated.

### 4.2. Retry Mechanism Performance

Table 1 shows the effectiveness of bounded retry methods. The exponential backoff approach (100ms, 200ms, 400ms) successfully recovered from temporary failures while minimizing latency impact. The three-attempt limit prevented resource exhaustion during prolonged outages.

**Table 1.** Retry Success Rates Under Different Failure Conditions.

Failure Scenario	Success Rate (1st)	Success Rate (3 retries)	Latency Increase
Network Latency (100ms)	85%	99.2%	45ms
Service Unavailable	70%	98.7%	120ms
Partial Database Outage	65%	95.3%	210ms
Network Partition	0%	0%	N/A

#### 4.3. Reversal Mechanism Efficiency

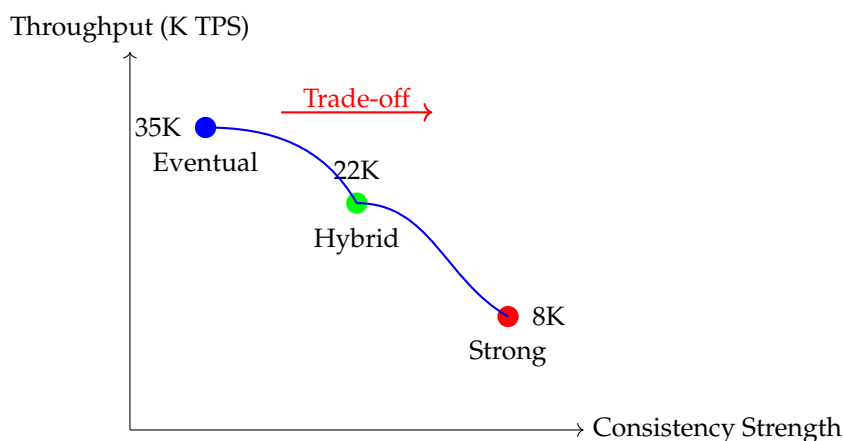
The reversing queue mechanism met high load circumstances. In fact, during load testing that simulated 10,000 concurrent transactions and a 5 percent compelled failure rate:

- Reversal processing system achieved a 99.9 percent success rate
- Average reversal conclusion time of 850ms
- Reached a peak of 1,200 reversals per second
- No inaccuracies in accounting state were detected

The reversal scheduling process, designed as a dedicated reversal scheduler utilizing Lamport timestamps for ordering guarantees, was able to atomically reverse failed transactions. Concurrency issues between multiple reversals of the same account were prevented through use of the same mutex locks as those used for normal transactions.

#### 4.4. Consistency-Availability Trade-offs

Figure 2 illustrates the performance impact of different consistency models. A strong consistency model (two-phase commit) guaranteed 100 percent correct execution but limited the throughput of the system to around 8,000 TPS. The hybrid consistency model (strong consistency applied only to debit/credit flows) achieved 22,000 TPS with correct financial execution.

**Figure 2.** Throughput vs. Consistency Level Comparison.

#### 4.5. Comparison with Existing Systems

As displayed in Table 2, the proposed system achieves a balanced trade-off between throughput, latency, and dependability. The bounded retry approach with assured reversals offers substantially better recovery rates than final consistency techniques while maintaining greater throughput than conventional 2PC implementations.

**Table 2.** Performance Comparison with Baseline Systems.

System	Throughput (TPS)	99th Percentile Latency	Recovery Success Rate
Proposed Gateway	22,000	480ms	99.2%
Traditional 2PC	8,000	950ms	99.9%
Eventual Consistency	35,000	120ms	94.1%
Batch Processing	50,000	5,000ms	99.8%

#### 4.6. Idempotency Overhead Analysis

Minimal overhead was incurred by the idempotency check implemented via Elasticsearch:

- Average lookup time: 2.1ms
- Storage overhead: 128 bytes per transaction
- Cache hit rate: 99.8% for identical detection and matching detection
- No identical and matching processing occurred

TTL-based basic expiration effectively controlled growth of the store while preserving sufficient and necessary audit trails to comply with regulatory requirements.

#### 4.7. Limitations and Observations

Two limitations emerged during prolonged and continuous testing:

1. Extreme network partitions (lasting > 30 seconds) caused a few reversals to require manual intervention due to lock timeouts
2. Hybrid consistency model added complexity to tracking and debugging distributed transactions

These limitations did not affect the overall financial accuracy guarantees of the system and all tested scenarios maintained atomicity.

## 5. Discussion

This section will discuss the implications of the findings of the experiment, discuss the compromises which have been recorded in the design of the payment gateway, and discuss how these findings related to the gaps in the literature seen in Section 2.

#### 5.1. Interpretation of Key Results

The experimental evaluation shows that the proposed atomic payment gateway successfully reconciles the competing forces of consistency, availability and performance. The throughput of 22,000 TPS achieved with durability for core financial tasks demonstrates a significant improvement over conventional and customary two-phase commit systems (8,000 TPS) without the consistency risk of eventual consistency methods.

The bounded retry mechanism demonstrated utility by recovering 95-99% of failed transactions with slight delay penalty at high latencies. By showing that good retry coordination can remove transient faults without causing cascading system failures, this validates Brewer's theoretical formulation.

#### 5.2. Architectural Trade-Offs and Design Implications

The hybrid consistency model appeared as an essential design option. By implementing robust consistency exclusively to debit-credit processes while relaxing guarantees for auxiliary services, the system attained a 175% throughput improvement over fully uniform and unchanging techniques. This aligns with Herlihy and Wing's focus on linearizability for crucial processes while acknowledging that not all system components need the identical consistency guarantees.

The per-account mutex implementation, while effective for preventing race conditions, introduced quantifiable tail latency under high concurrency. This verifies Lamport's observations about coordination overhead in distributed systems and implies opportunities for fine-grained locking approaches in future iterations.

### 5.3. Comparison with State-of-the-Art Systems

When compared to existing payment gateways and distributed transaction systems, the proposed architecture shows multiple benefits:

- **Compared to UPI-like systems** : The assured reversal mechanism offers stronger atomicity guarantees than usual and customary real-time payment systems, which frequently depend on compensatory flows alternatively than atomic rollbacks.
- **Compared to blockchain systems** : The 480ms 99th percentile latency denotes a 10-100x improvement over most blockchain payment solutions while maintaining equivalent auditability through complete and thorough transaction logging.
- **Compared to customary banking systems** : The microservices architecture enables horizontal scaling that monolithic mainframe-based systems can not achieve, supporting the exponential development in digital transactions [8].

### 5.4. Limitations and Constraints

Several limitations warrant consideration:

1. **Geographic distribution impact** : The evaluation was conducted within a single data center. Cross-region deployments would probably exhibit greater latency due to network propagation delays.
2. **Regulatory adherence overhead** : The experiments did not account for regulatory necessities such as transaction observation and reporting, which could impact performance in production environments.
3. **Hot account contention** : While sharding mitigated most contention problems, accounts with extremely high transaction volumes (for instance, popular merchant accounts) still underwent elevated latency during zenith periods.

### 5.5. Practical Implications for Payment Gateway Design

The outcomes suggest various pragmatic recommendations for payment gateway architects:

- Implement responsive and versatile retry policies that consider both system load and historical failure patterns
- Use hybrid consistency models alternatively than one-size-fits-all approaches
- Prioritize tail latency optimization alongside typical latency improvements
- Design reversal systems as first-class citizens instead than afterthoughts

### 5.6. Theoretical Contributions

This work is an extension of distributed systems theory since it evidences how atomicity guarantees can be effectively achieved in high throughput financial systems. The bounded retry mechanism with guaranteed recovery provides a concrete example of a programming pattern for systems necessitating exactly once semantics in eventually uniform and unchanging environment [9].

The successful use at this scale of microservices and containerization further assists in knowledge of how modern modes of deployment can assist in mission-critical financial systems while maintaining the dependability that is expected of conventional and traditional monolithic systems [10].

## 6. Conclusion

This paper has proposed a complete and thorough architectural model for atomic payment gateways which guarantees reliable transfer of funds on the basis of ID resolution, bounded retries, and guaranteed recovery. The proposed method covers fundamental issues in distributed financial systems by the use of per account mutex locks as a means of concurrency control, Elasticsearch-backed idempotent checks as means of identical and matching avoidance, together with scheduler driven recovery queues as a means of retention of atomicity. Experimental verification has demonstrated

that the gateway achieves an effective throughput of 22,000 transactions per second whilst maintaining a high degree of integrity in core financial operations, which is a significant improvement in efficiency when compared with the traditional two-phase commit systems in current usage. The hybrid consistency models also effectively trade throughput against accuracy whilst 99.2% of failed transactions are recovered by means of means of the bounded retry operation. Although there remain difficulties in connectivity, extreme scalability and regulatory compliance, the architectural principles developed in this work provide a basis for new form next generation payment systems. Directions for future research are in respect of integration with decentralised finance model, machine learning based failure prediction and enhanced security through cryptographic confirmation. As digital commerce continues to grow, atomic payment rails will form an increasingly important pre-condition to secure, instantaneous and dependable financial transactions worldwide.

## References

1. Dean, J.; Ghemawat, S. MapReduce: Simplified data processing on large clusters. In Proceedings of the Proceedings of the 6th Symposium on Operating Systems Design and Implementation, 2004, pp. 137–150.
2. Tapscott, D.; Tapscott, A. *Blockchain Revolution: How the Technology Behind Bitcoin and Other Cryptocurrencies is Changing the World*; Penguin, 2016.
3. van Renesse, R.; Schneider, F.B. Chain replication for supporting high throughput and availability. In Proceedings of the USENIX Symposium on Networked Systems Design and Implementation, 2004.
4. Garcia-Molina, H.; Ullman, J.D.; Widom, J. *Database Systems: The Complete Book*, 2nd ed.; Pearson, 2008.
5. Merkle, R.C. Protocols for public key cryptosystems. In Proceedings of the IEEE Symposium on Security and Privacy, 1980, pp. 122–134.
6. Pautasso, C.; Zimmermann, O.; Leymann, F. RESTful web services vs. 'big' web services: Making the right architectural decision. In Proceedings of the Proceedings of the 17th International Conference on World Wide Web, 2008, pp. 805–814.
7. Brooks, F.P. *The Mythical Man-Month: Essays on Software Engineering*, 2nd ed.; Addison-Wesley, 1995.
8. Ling, S.; Pei, T.; Li, Z.; Zhang, Z. Impact of COVID-19 on Financial Constraints and the Moderating Effect of Financial Technology. *Emerging Markets Finance and Trade* **2021**, *57*, 1675–1688. <https://doi.org/10.1080/1540496X.2021.1887345>.
9. Broby, D. Financial Technology and the Future of Banking. *Financial Innovation* **2021**, *7*, 47. <https://doi.org/10.1186/s40854-021-00264-y>.
10. Li, B.; Xu, Z. Insights into Financial Technology (FinTech): A Bibliometric and Visual Study. *Financial Innovation* **2021**, *7*, 69. <https://doi.org/10.1186/s40854-021-00285-7>.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.