



Achieving consistency where distributed transactions have failed.

BY MARTIN KLEPPMANN, ALASTAIR R. BERESFORD,
AND BOERGE SVINGEN

Online Event Processing

FOR ALMOST HALF a century, ACID transactions (satisfying the properties of atomicity, consistency, isolation, and durability) have been the abstraction of choice for ensuring *consistency* in data-storage systems. The well-known *atomicity* property ensures that either all or none of a transaction's writes take

effect in the case of a failure; *isolation* prevents interference from concurrently running transactions; and *durability* ensures that writes made by committed transactions are not lost in the case of a failure.

While transactions work well within the scope of a single database product, transactions that span several different data-storage products from distinct vendors have been problematic: many storage systems do not support them, and those that do often perform poorly. Today, large-scale applications are often implemented by combining several distinct data-storage technologies that are optimized for different access patterns. Distributed transactions have failed to gain adoption in most such settings, and most large-scale applications instead rely on ad hoc, unreliable approaches for maintaining the consistency of their data systems.

In recent years, however, there has been an increase in the use of event

logs as a data-management mechanism in large-scale applications. This trend includes the event-sourcing approach to data modeling, the use of change data capture systems, and the increasing popularity of log-based publish/subscribe systems such as Apache Kafka. Although many databases use logs internally (for example, write-ahead logs or replication logs), this new generation of log-based systems is different: rather than using logs as an implementation detail, they raise them to the level of the application-programming model.

Since this approach uses application-defined events to solve problems that traditionally fall in the transaction-processing domain, we name it OLEP (online *event* processing) to contrast with OLTP (online *transaction* processing) and OLAP (online *analytical* processing). This article explains the reasons for the emergence of OLEP and shows how it allows ap-

lications to guarantee strong consistency properties across heterogeneous data systems, without resorting to atomic commit protocols or distributed locking. The architecture of OLEP systems allows them to achieve consistent high performance, fault tolerance, and scalability.

Application Architecture Today: Polyglot Persistence

Different data-storage systems are designed for different access patterns, and there is no single one-size-fits-all storage technology that is able to serve all possible uses of data efficiently. Consequently, many applications today use a combination of several different storage technologies, an approach sometimes known as *polyglot persistence*.


For example:

► *Full-text search*. When users need to perform a keyword search on a dataset (for example, a product catalog), a full-text search index is required. Although some relational databases, such as PostgreSQL, include a basic full-text indexing feature, more advanced uses generally require a dedicated search server such as Elasticsearch. To improve the indexing or search result ranking algorithms, the search engine's indexes may need to be rebuilt from time to time.


► *Data warehousing*. Most enterprises export operational data from their OLTP databases and load it into a data warehouse for business analytics. The storage layouts that perform well for such analytic workloads, such as column-oriented encoding, are very different from those of OLTP storage engines, necessitating the use of distinct systems.

► *Stream processing*. Message brokers allow an application to subscribe to a stream of events as they happen (for example, representing the actions of users on a website), and stream processors provide infrastructure for interpreting and reacting to those streams (for example, detecting patterns of fraud or abuse).

► *Application-level caching*. To improve the performance of read-only requests, applications often maintain caches of frequently accessed objects (for example, in memcached). When the underlying



The architecture of online event processing systems allows them to achieve consistent high performance, fault tolerance, and scalability.



data changes, applications employ custom logic to update the affected cache entries accordingly.

Note these storage systems are not fully independent of each other. Rather, it is common for one system to hold a copy or materialized view of data in another system. Thus, when data in one system is updated, it often needs to be updated in another, as illustrated in Figure 1.

OLTP transactions are predefined and short. In the traditional view, as implemented by most relational database products today, a transaction is an interactive session in which a client's queries and data modification commands are interleaved with arbitrary processing and business logic on the client. Moreover, there is no time limit for the duration of a transaction, since the session traditionally may have included human interaction.

However, reality today looks different. Most OLTP database transactions are triggered by a user request made via HTTP to a Web application or Web service. In the vast majority of applications, the span of a transaction extends no longer than the handling of a single HTTP request. This means that by the time the service sends its response to the user, any transactions on the underlying databases have already been committed or aborted. In a user workflow that spans several HTTP requests (for example, adding an item to a cart, going to checkout, confirming the shipping address, entering payment details, and giving a final confirmation), no one transaction spans the entire user workflow; there are only short, noninteractive transactions to handle single steps of the workflow.

Moreover, an OLTP system generally executes a fairly small set of known transaction patterns. On this basis, some database systems encapsulate the business logic of transactions as *stored procedures* that are registered ahead of time by the application. To execute a transaction, a stored procedure is invoked with certain input parameters, and the procedure then runs to completion on a single execution thread without communicating with any nodes outside of the database.

Heterogeneous distributed transactions are problematic. It is important to distinguish between two types of distributed transactions:

- ▶ Homogeneous distributed transactions are those in which the participating nodes are all running the same database software. For example, Google's Cloud Spanner and VoltDB are recent database systems that support homogeneous distributed transactions.

- ▶ Heterogeneous distributed transactions span several different storage technologies by distinct vendors. For example, the X/Open XA (extended architecture) standard defines a transaction model for performing 2PC (two-phase commit) across heterogeneous systems, and the JTA (Java Transaction API) makes XA available to Java applications.

While some homogeneous transaction implementations have proved successful, heterogeneous transactions continue to be problematic. By their nature, they can only rely on a lowest common denominator of participating systems. For example, XA transactions block execution if the application process fails during the *prepare* phase; moreover, XA provides no deadlock detection and no support for optimistic concurrency-control schemes.³

Many of the systems listed here, such as search indexes, do not support XA or any other heterogeneous transaction model. Thus, ensuring the atomicity of writes across different storage technologies remains a challenging problem for applications.

Building Upon Event Logs

Figure 1 shows an example of polyglot persistence: an application that needs to maintain records in two separate storage systems such as an OLTP database (for example, an RDBMS) and a full-text search server. If heterogeneous distributed transactions are available, the system can ensure atomicity of writes across the two systems. Most search servers do not support distributed transactions, however, leaving the system vulnerable to these potential inconsistencies:

- ▶ *Non-atomic writes.* If a failure occurs, a record may be written to one of the systems but not the other, leaving them inconsistent with each other.

- ▶ *Different order of writes.* If there are two concurrent update requests A and B for the same record, one system may process them in the order A, B while the other system processes them in the order B, A. Thus, the systems may disagree on which write was the latest, leaving them inconsistent.

Figure 2 presents a simple solution to these problems: when the application wants to update a record, rather than performing direct writes to the two storage systems, it appends an *update event* to a log. The database and the search index each subscribe to this log and write updates to their storage in the order they appear in the log.⁴ By sequencing updates through a log, the database and the search index apply the same set of writes in the same order, keeping them consistent with each other. In effect, the database and the search index are *materialized views* onto the sequence of events in the log. This approach solves both of the aforementioned problems as follows:

- ▶ Appending a single event to a log is atomic; thus, either both subscribers see an event, or neither does. If a subscriber fails and recovers, it resumes processing any events that it has not processed previously. Thus, if an update is written to the log, it will eventually be processed by all subscribers.

- ▶ All subscribers of the log see its events in the same order. Thus, each of the storage systems will write records in the same serial order.

In this example, the log serializes writes only, but the application may read from the storage systems at any time. Since the log subscribers are asynchronous, reading the index may return a record that does not yet exist in the database, or vice versa; such transient inconsistencies are not a problem for many applications. For those applications that require it, reads can also be serialized through the log; an example of this is presented later.

The log abstraction. There are sev-

Figure 1. Record written to a database and to search index.

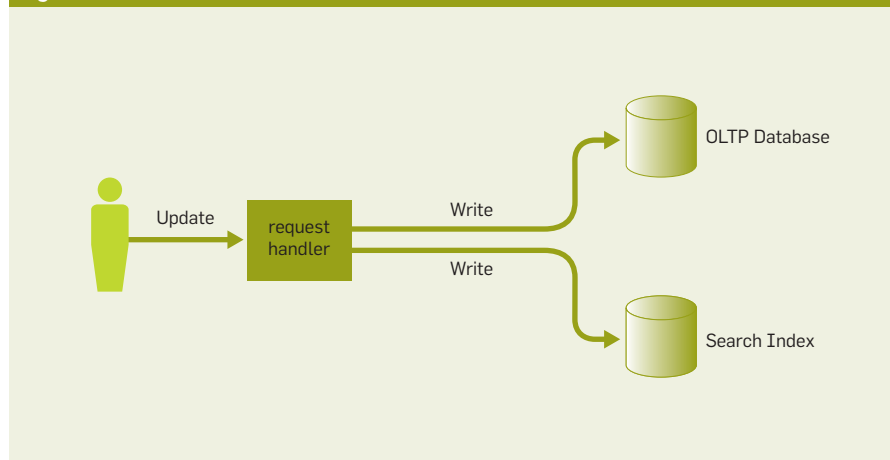
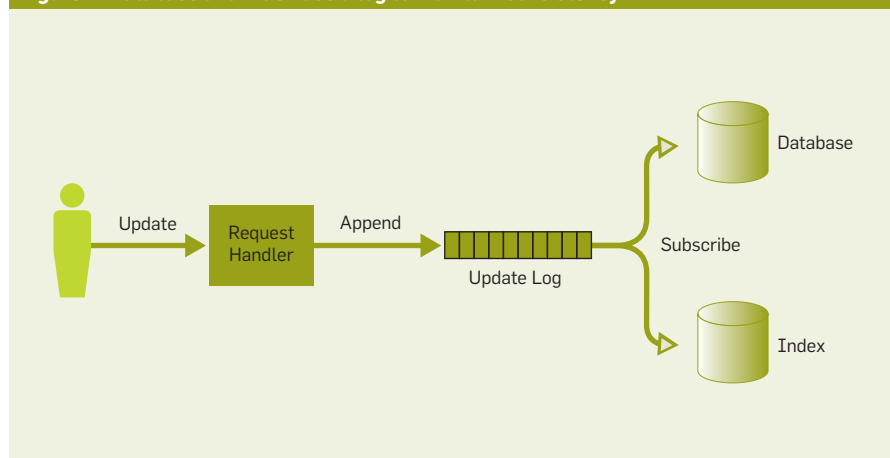


Figure 2. Database and Index use a log to maintain consistency.



eral log implementations that can serve this role, including Apache Kafka, CORFU (from Microsoft Research), Apache Pulsar, and Facebook’s LogDevice. The required log abstraction has the following properties:

- ▶ *Durable*. The log is written to disk and replicated to several nodes, ensuring that no events are lost in a failure.

- ▶ *Append-only*. New events can be added to the log only by appending them at the end. Besides appending, the log may allow old events to be discarded (for example, by truncating log segments older than some retention period or by performing key-based log compaction).

- ▶ *Sequential reads*. All subscribers of the log see the same events in the same order. Each event is assigned a monotonically increasing LSN (log sequence number). A subscriber reads the log by starting from a specified LSN and then receiving all subsequent events in log order.

- ▶ *Fault-tolerant*. The log remains highly available for reads and writes in the presence of failures.

- ▶ *Partitioned*. An individual log may have a maximum throughput it can support (for example, the throughput of a single network interface or a single disk). The system can be assumed to scale linearly, however, by having

many *partitions*—that is, many independent logs that can be distributed across many machines—and to have no ordering guarantee across different log partitions. Multiple logical logs may be multiplexed into a single physical log partition.

The following assumptions are made about subscribers of a log:

- ▶ A subscriber may maintain state (for example, a database) that is read and updated based on the events in the log, and that survives crashes. Moreover, a subscriber may append further events to any log (including its own input).

- ▶ A subscriber periodically checkpoints the latest LSN it has processed to stable storage. When a subscriber crashes, upon recovery it resumes processing from the latest checkpointed LSN. Thus, a subscriber may process some events twice (those processed between the last checkpoint and the crash), but it never skips any events. Events in the log are processed at least once by each subscriber.

- ▶ The events in a single log partition are processed sequentially on a single thread, using deterministic logic. Thus, if a subscriber crashes and restarts, it may append duplicate events to other logs.

These assumptions are satisfied by existing log-based stream-processing

frameworks such as Apache Kafka Streams and Apache Samza. Updating state deterministically based on an ordered log corresponds to the classic *state machine replication* principle.⁵ Since it is possible for an event to be processed more than once when recovering from a failure, state updates must also be *idempotent*.

Aside: Exactly-once semantics.

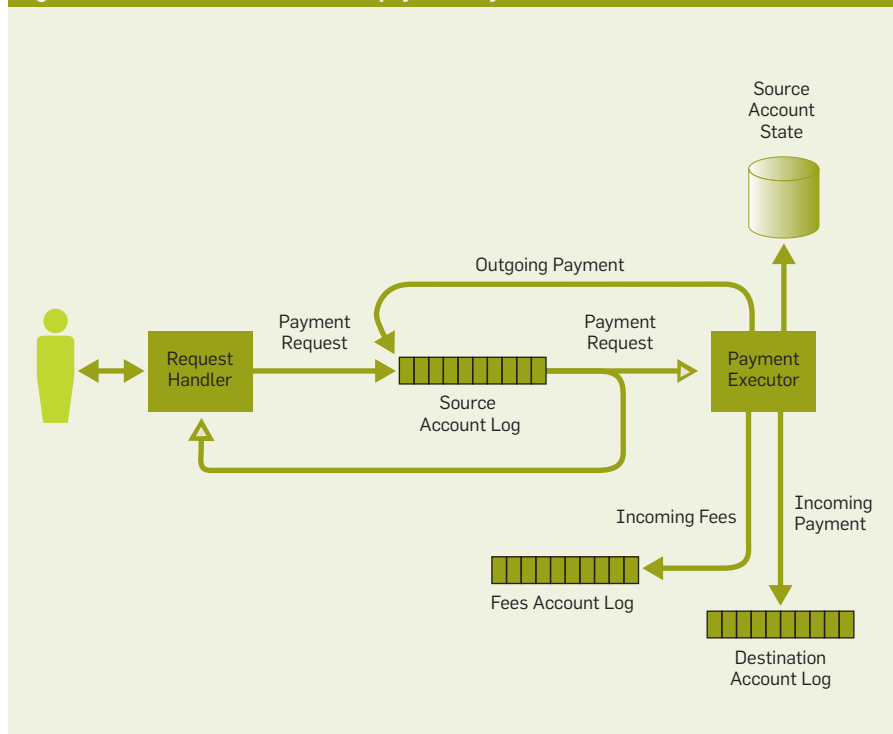
Some log-based stream processors such as Apache Flink support so-called *exactly-once semantics*, which means that even though an event may be processed more than once, the effect of the processing will be the same as if it had been processed exactly once. This behavior is implemented by managing side effects within the processing framework and atomically committing these side effects together with the checkpoint that marks a section of the log as processed.

When a log consumer writes to external storage systems, however, as in Figure 2, exactly-once semantics cannot be ensured, since doing so would require a heterogeneous atomic commit protocol across the stream processor and the storage system, which is not available on many storage systems, such as full-text search indexes. Thus, frameworks with exactly-once semantics still exhibit at-least-once processing when interacting with external storage and rely on idempotence to eliminate the effects of duplicate processing.

Atomicity and enforcing constraints. A classic example where atomicity is required is in a banking/payments system, where a transfer of funds from one account to another account must happen atomically, even if the two accounts are stored on different nodes. Moreover, such a system typically needs to maintain consistency properties or invariants (for example, an account cannot be overdrawn by more than some set limit). Figure 3 shows how such a payments application can be implemented using the OLEP approach instead of distributed transactions. Arrows with solid heads denote appending an event to a log, while arrows with hollow heads denote subscribing to the events in a log. It works as follows:

1. When a user wishes to transfer funds from a source account to a destination account, he or she first appends

Figure 3. Flow of events in a financial payments system.



a *payment request* event to the log of the source account. This event merely indicates the intention to transfer funds; it does not imply that the transfer has been successful. The event carries a unique ID to identify the request.

2. A single-threaded *payment executor* process subscribes to the source-account log. It maintains a database containing transactions on the source account and the current balance. This process deterministically checks whether the payment request should be allowed, based on the current balance and perhaps other factors. This log consumer is very similar to the execution of a stored procedure.


3. If the executor decides to grant the payment request, it writes that fact to its local database and appends events to several different logs: as a minimum, an outgoing payment event to the source account log and an incoming payment event to the log for the destination account. If a fee is due for this payment (for example, because of an overdrawn account or currency conversion), an additional outgoing payment event for the fees may be appended to the source-account log, and a corresponding incoming payment event may be appended to the log of a fees account. The original event ID is included in all of these generated events so that their origin can be traced.

4. Since the executor subscribes to the source-account log, the outgoing payment event will be delivered back to the executor. It uses the unique event ID to determine that it has already processed this payment and recorded it in its database.


5. The payment events on other accounts, such as the incoming payment on the destination account, are similarly processed by single-threaded executors, with a separate executor per account. The event processing is made idempotent by suppressing duplicates based on the original event ID.

6. The server handling the user's request may also subscribe to the source-account log and thus be notified when the payment request has been processed. This status information can be returned to the user.

If the payment executor crashes and restarts, it may reprocess some payment requests that were partially processed



Heterogeneous transactions continue to be problematic. By their very nature, they can only rely on a lowest common denominator of participating systems.



before the crash. Since the executor is deterministic, upon recovery it will make the same decisions to approve or decline requests, and thus potentially append duplicate payment events to the source, destination, and fees logs. Based on the ID in the events, however, it is easy for downstream processes to detect and ignore such duplicates.

Multipartition processing. In this payment example, each account has a separate log and thus may be stored on a different node. Moreover, each payment executor only needs to subscribe to events from a single account, and different executors handle different accounts. These factors allow the system to scale linearly to an arbitrary number of accounts.

In this example, the decision of whether to allow the payment request is conditional only on the balance of the source account; you can assume that the payment into the destination account always succeeds, since its balance can only increase. For this reason, the payment executor needs to serialize the payment request only with respect to other events in the source account. If other log partitions need to contribute to the decision, the approval of the payment request can be performed as a multistage process in which each stage serializes the request with respect to a particular log.

Splitting a “transaction” into a multistage pipeline of stream processors allows each stage to make progress based only on local data; it ensures that one partition is never blocked waiting for communication or coordination with another partition. Unlike multipartition transactions, which often impose a scalability bottleneck in distributed transaction implementations, this pipelined design allows OLEP systems to scale linearly.

Advantages of event processing. Besides this scalability advantage, developing applications in an OLEP style has several further advantages:

- ▶ Since every log can support many independent subscribers, it is easy to create new derived views or services based on an event log. For example, in the payment scenario of Figure 3, a new account log subscriber could send a push notification to a customer's smartphone if a certain spending limit on the customer's credit card is

reached. A new search index or view over an existing dataset can be built simply by consuming the event log from beginning to end.³


► If an application bug causes bad events to be appended to a log, it is fairly easy to recover: subscribers can be programmed to ignore the incorrect events, and any views derived from the events can be recomputed. In contrast, in a database that supports arbitrary insertions, updates, and deletes, it is much harder to recover from incorrect writes, potentially requiring the database to be restored from a backup.

► Similarly, debugging is much easier with an append-only log than a mutable database, because events can be replayed in order to diagnose what happened in a particular situation.


► For data-modeling purposes, an append-only event log is increasingly preferred over freeform database mutations; this approach is known in the domain-driven design community as *event sourcing*.² The rationale is that events capture state transitions and business processes more accurately than insert/update/delete operations on tables, and those state updates are better described as side effects resulting from processing an event. For example, the event “*student cancelled course enrollment*” clearly expresses intent, whereas the side effects “*one row was deleted from the enrollments table*” and “*one cancellation reason was added to the student feedback table*” are much less clear.

► From a data analysis point of view, an event log is more valuable than the state in a database. For example, in an e-commerce setting, it is valuable for business analysts to see not only the final state of the cart at checkout, but also the full sequence of items added to and removed from the cart, since the removed items carry information, too (for example, one product is a substitute for another, or the customer may return to buy a certain item on a later occasion).

► With a distributed transaction, if any one of the participating nodes is unavailable, the whole transaction must abort, so failures are amplified. In contrast, if a log has multiple subscribers, they make progress independently from each other: if one subscriber fails,



Debugging is much easier with an append-only log than a mutable database because events can be replayed in order to diagnose what happened in a particular situation.



that does not impede the operation of the publisher or other subscribers, so faults are contained.

Disadvantages of the OLEP approach. In the previous examples, log consumers update the state in data stores (the database and search index in Figure 2; the account balances and account statements in Figure 3). While the OLEP approach ensures every event in the log will eventually be processed by every consumer, even in the face of crashes, there is no upper bound on the time until an event is processed.

This means if a client reads from two different data stores that are updated by two different consumers or log partitions, then the values read by the client may be inconsistent with each other. For example, reading the source and destination accounts of a payment may return the source account after the payment has been processed, but the destination account before it has been processed. Thus, even though the accounts will eventually converge toward a consistent state, they may be inconsistent when read at one particular point in time.

Note that in an ACID context, preventing this anomaly falls under the heading of *isolation*, not *atomicity*; a system with atomicity alone does not guarantee that two accounts will be read in a consistent state. A database transaction running at “read committed” isolation level—the default isolation level in many systems including PostgreSQL, Oracle DB, and SQL Server—may experience the same anomaly when reading from two accounts.³ Preventing this anomaly requires a stronger isolation level: “repeatable read,” snapshot isolation, or serializability.

At present, the OLEP approach does not provide isolation for read requests that are sent directly to data stores (rather than being serialized through the log). Hopefully, future research will enable stronger isolation levels such as snapshot isolation across data stores that are updated from a log.

Case Study: The New York Times

The *New York Times* maintains all textual content published since the newspaper’s founding in 1851 in a single log

partition in Apache Kafka.⁶ Image files are stored in a separate system, but URLs and captions of images are also stored as log events.

Whenever a piece of content (known as an *asset*) is published or updated, an event is appended to this log. Several systems subscribe to this log: for example, the full text of each article is written to an indexing service for full-text search; various cached pages (for example, the list of articles with a particular tag, or all pieces by a particular author) need to be updated; and personalization systems notify readers who may be interested in a new article.

Each asset is given a unique identifier, and an event may create or update an asset with a given ID. Moreover, an event may reference the identifiers of other assets—much like a normalized schema in a relational database, where one record may reference the primary key of another record. For example, an image (with caption and other metadata) is an asset that may be referenced by one or more articles.

The order of events in the log satisfies two rules:

- ▶ Whenever one asset references another, the event that publishes the referenced asset appears in the log before the referencing asset.

- ▶ When an asset is updated, the latest version is the one published by the latest event in the log.

For example, an editor might publish an image and then update an article to reference the image. Every consumer of the log then passes through three states in sequence:

1. The old version of the article (not referencing the image) exists.
2. The image also exists but is not yet referenced by any article.
3. The article and image both exist, with the article referencing the image.

Different log consumers will pass through these three states at different times but in the same order. The log order ensures that no consumer is ever in a state where the article references an image that does not yet exist, ensuring referential integrity.

Moreover, whenever an image or caption is updated, all articles referencing that image need to be updated in caches and search indexes. This can easily be achieved with a log con-

sumer that uses a database to keep track of references between articles and images. This consistency model lends itself very easily to a log, and it provides most of the benefits of distributed transactions without the performance costs.

Further details on the *New York Times's* approach appear in a blog post.⁶


Conclusion

Support for distributed transactions across heterogeneous storage technologies is either nonexistent or suffers from poor operational and performance characteristics. In contrast, OLEP is increasingly used to provide good performance and strong consistency guarantees in such settings.

In data systems it is very common for logs (for example, write-ahead logs) to be used as internal implementation details. The OLEP approach is different: it uses event logs, rather than transactions, as the *primary application programming model* for data management. Traditional databases are still used, but their writes come from a log rather than directly from the application. This approach has been explored by several influential figures in industry, such as Jay Kreps,⁴ Martin Fowler,² and Greg Young under names such as event sourcing and CQRS (Command/Query Responsibility Segregation).^{1,7}

The use of OLEP is not simply pragmatism on the part of developers, but rather it offers a number of advantages. These include linear scalability; a means of effectively managing polyglot persistence; support for incremental development where new application features or storage technologies are added or removed iteratively; excellent support for debugging via direct access to the event log; and improved availability (because running nodes can continue to make progress when other nodes have failed).

Consequently, OLEP is expected to be increasingly used to provide strong consistency in large-scale systems that use heterogeneous storage technologies.

Acknowledgments. This work was supported by a grant from The Boeing Company. Thanks to Pat Helland for feedback on a draft of this article. 

Related articles on queue.acm.org

Consistently Eventual

Pat Helland

<https://queue.acm.org/detail.cfm?id=3226077>

Evolution and Practice: Low-latency Distributed Applications in Finance

Andrew Brook

<https://queue.acm.org/detail.cfm?id=2770868>

It Isn't Your Father's Real Time Anymore

Phillip Laplante

<https://queue.acm.org/detail.cfm?id=1117409>

References

1. Betts, D., Domínguez, J., Melnik, G., Simonazzi, F. and Subramanian, M. *Exploring CQRS and Event Sourcing*. Microsoft Patterns & Practices, 2012; <http://aka.ms/cqrs>.
2. Fowler, M. Event sourcing, 2005; <https://www.martinfowler.com/eaaDev/EventSourcing.html>.
3. Kleppmann, M. *Designing Data-intensive Applications*. O'Reilly Media, 2017.
4. Kreps, J. The log: What every software engineer should know about real-time data's unifying abstraction. LinkedIn Engineering, 2013; <https://bit.ly/199IMwY>.
5. Schneider, F.B. Implementing fault-tolerant services using the state machine approach: A tutorial. *ACM Computing Surveys* 22, 4 (1990), 299–319; <https://dl.acm.org/citation.cfm?doi=98163.98167>.
6. Svingen, B. Publishing with Apache Kafka at the *New York Times*. (Sept. 5 2017); <https://open.nytimes.com/publishing-with-apache-kafka-at-the-new-york-times-7f0e3b7d2077>.
7. Vernon, V. *Implementing Domain-driven Design*. Addison-Wesley, 2013.

Martin Kleppmann is a distributed-systems researcher at the University of Cambridge and author of *Designing Data-Intensive Applications* (<http://dataintensive.net/>). Previously he was a software engineer, cofounding two startups and working on large-scale data infrastructure at LinkedIn.

Alastair R. Beresford is a reader in computer security at the University of Cambridge. His work examines the security and privacy of large-scale distributed computer systems, with a particular focus on networked mobile devices.

Boerge Svingen is a director of engineering at the *New York Times*. He was a founder of Fast Search & Transfer (alltheweb.com, FAST ESP) as well as a founder and CTO of Open AdExchange.