

FI

FI

FI, built on Spanner, is “...a rewrite of Google’s advertising backend. FI uses five replicas spread across the United States. Most other applications will probably replicate their data across 3 to 5 datacenters in one geographic region, but with relatively independent failure modes. That is, most applications will choose lower latency over higher availability, as long as they can survive 1 or 2 datacenter failures.” (Corbett et al. 2013, 8:2).

“[Adwords] was originally based on a MySQL database that was manually sharded many ways. The uncompressed dataset is tens of terabytes, which is small compared to many NoSQL instances, but was large enough to cause difficulties with sharded MySQL.” (Corbett et al. 2013, 8:17).

“The MySQL sharding scheme assigned each customer and all related data to a fixed shard. This layout enabled the use of indexes and complex query processing on a per-customer basis, but required some knowledge of the sharding in application business logic. Resharding this revenue-critical database as it grew in the number of customers and the size of their data was extremely costly. The last resharding took over two years of intense effort, and involved coordination and testing across dozens of teams to minimize risk. This operation was too complex to do regularly: as a result, the team had to limit growth on the MySQL database by storing some data in external Bigtables, which compromised transactional behavior and the ability to query across all data.” (Corbett et al. 2013, 8:17).

FI

FI team chose to use Spanner for several reasons:

- Spanner removes the need to manually reshard
- Spanner provides synchronous replication and automatic failover.
 - With MySQL, failover was difficult and risked data-loss and downtime.
- FI requires strong transactional semantics (transactions across arbitrary data, and consistent reads): not practical in other NoSQL systems.

FI team also needed secondary indexing on their data. Spanner doesn't yet provide automatic support for secondary indexes, but FI team could implement their own consistent global indexes using Spanner transactions.

(Corbett et al. 2013, 8:17).

F1: A Distributed SQL Database That Scales

Jeff Shute
Chad Whipkey
David Menestrina

Radek Vingralek
Eric Rollins
Stephan Ellner
Traian Stancescu

Bart Samwel
Mircea Oancea
John Cieslewicz
Himani Apte

Ben Handy
Kyle Littlefield
Ian Rae*

Google, Inc.

*University of Wisconsin-Madison

ABSTRACT

F1 is a distributed relational database system built at Google to support the AdWords business. F1 is a hybrid database that combines high availability, the scalability of NoSQL systems like Bigtable, and the consistency and usability of traditional SQL databases. F1 is built on Spanner, which provides synchronous cross-datacenter replication and strong consistency. Synchronous replication implies higher commit latency, but we mitigate that latency by using a hierarchical schema model with structured data types and through smart application design. F1 also includes a fully functional distributed SQL query engine and automatic change tracking and publishing.

1. INTRODUCTION

F1¹ is a fault-tolerant globally-distributed OLTP and OLAP database built at Google as the new storage system for Google's AdWords system. It was designed to replace a sharded MySQL implementation that was not able to meet our growing scalability and reliability requirements.

The key goals of F1's design are:

1. **Scalability:** The system must be able to scale up, trivially and transparently, just by adding resources.

consistent and correct data.

Designing applications to cope with concurrency anomalies in their data is very error-prone, time-consuming, and ultimately not worth the performance gains.

4. **Usability:** The system must provide full SQL query support and other functionality users expect from a SQL database. Features like indexes and ad hoc query are not just nice to have, but absolute requirements for our business.

Recent publications have suggested that these design goals are mutually exclusive [5, 11, 23]. A key contribution of this paper is to show how we achieved all of these goals in F1's design, and where we made trade-offs and sacrifices. The name F1 comes from genetics, where a *Filial 1 hybrid* is the first generation offspring resulting from a cross mating of distinctly different parental types. The F1 database system is indeed such a hybrid, combining the best aspects of traditional relational databases and scalable NoSQL systems like Bigtable [6].

F1 is built on top of Spanner [7], which provides extremely scalable data storage, synchronous replication, and strong consistency and ordering properties. F1 inherits those fea-

FI

“FI is a fault-tolerant globally distributed OLTP and OLAP database built at Google... It was designed to replace a sharded MySQL implementation that was not able to meet our growing scalability and reliability requirements”. (Shute et al. 2013, p.1).

FI

OnLine Transaction Processing

“FI is a fault-tolerant globally distributed OLTP and OLAP database built at Google... It was designed to replace a sharded MySQL implementation that was not able to meet our growing scalability and reliability requirements”. (Shute et al. 2013, p.1).

FI

OnLine
Analytical
Processing

“FI is a fault-tolerant globally distributed OLTP and OLAP database built at Google... It was designed to replace a sharded MySQL implementation that was not able to meet our growing scalability and reliability requirements”. (Shute et al. 2013, p.1).

FI

“FI is a fault-tolerant globally distributed OLTP and OLAP database built at Google... It was designed to replace a sharded MySQL implementation that was not able to meet our growing scalability and reliability requirements”. (Shute et al. 2013, p.1).

Design Goals:

FI

“FI is a fault-tolerant globally distributed OLTP and OLAP database built at Google... It was designed to replace a sharded MySQL implementation that was not able to meet our growing scalability and reliability requirements”. (Shute et al. 2013, p.1).

Design Goals:

- **Scalability:** scale up, trivially and transparently, just by adding resources.

FI

“FI is a fault-tolerant globally distributed OLTP and OLAP database built at Google... It was designed to replace a sharded MySQL implementation that was not able to meet our growing scalability and reliability requirements”. (Shute et al. 2013, p.1).

Design Goals:

- Scalability: scale up, trivially and transparently, just by adding resources.
- **Availability:** system must never go down.

FI

“FI is a fault-tolerant globally distributed OLTP and OLAP database built at Google... It was designed to replace a sharded MySQL implementation that was not able to meet our growing scalability and reliability requirements”. (Shute et al. 2013, p.1).

Design Goals:

- Scalability: scale up, trivially and transparently, just by adding resources.
- Availability: system must never go down.
- **Consistency**: provide ACID transactions.

FI

“FI is a fault-tolerant globally consistent and OLAP database built at Google... It was designed to replace MySQL implementation that was not able to meet our growth and reliability requirements”. (Shute et al. 2013, p.1).

**Attomic
Consistent
Isolated
Durable
(Lecture 10)**

Design Goals:

- Scalability: scale up, trivially and transparently, just by adding resources.
- Availability: system must never go down.
- Consistency: provide ACID transactions.

FI

“FI is a fault-tolerant globally distributed OLTP and OLAP database built at Google... It was designed to replace a sharded MySQL implementation that was not able to meet our growing scalability and reliability requirements”. (Shute et al. 2013, p.1).

Design Goals:

- Scalability: scale up, trivially and transparently, just by adding resources.
- Availability: system must never go down.
- Consistency: provide ACID transactions.
- Usability: full SQL support.

FI

“FI is a fault-tolerant globally distributed OLTP and OLAP database built at Google... It was designed to replace a sharded MySQL implementation that was not able to meet our growing scalability and reliability requirements”. (Shute et al. 2013, p.1).

Design Goals:

- Scalability: scale up, trivially and transparently, just by adding resources.
- Availability: system must never go down.
- Consistency: provide ACID transactions.
- Usability: full SQL support.

FI running Adwords since early 2012: 100s of applications; 1000s of users; DB >100Tb; serves >100,000s requests/s; scans >10,000,000,000,000 rows.

F1 Basic Architecture

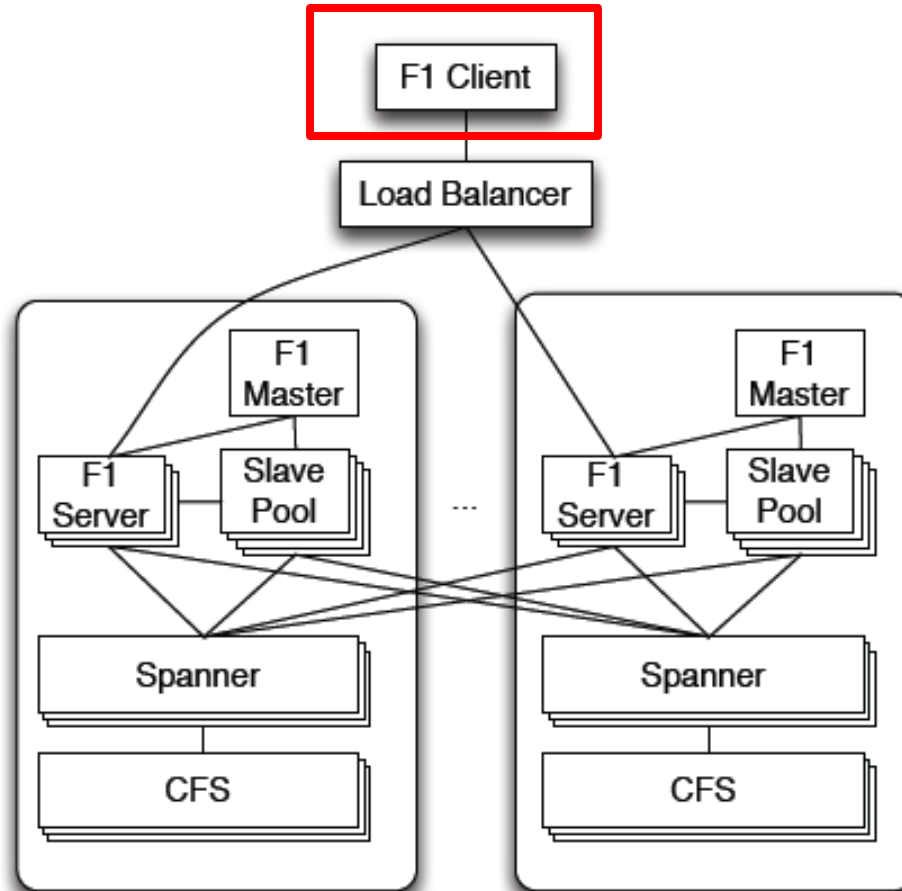


Figure 1: The basic architecture of the F1 system, with servers in two datacenters.

F1 Basic Architecture

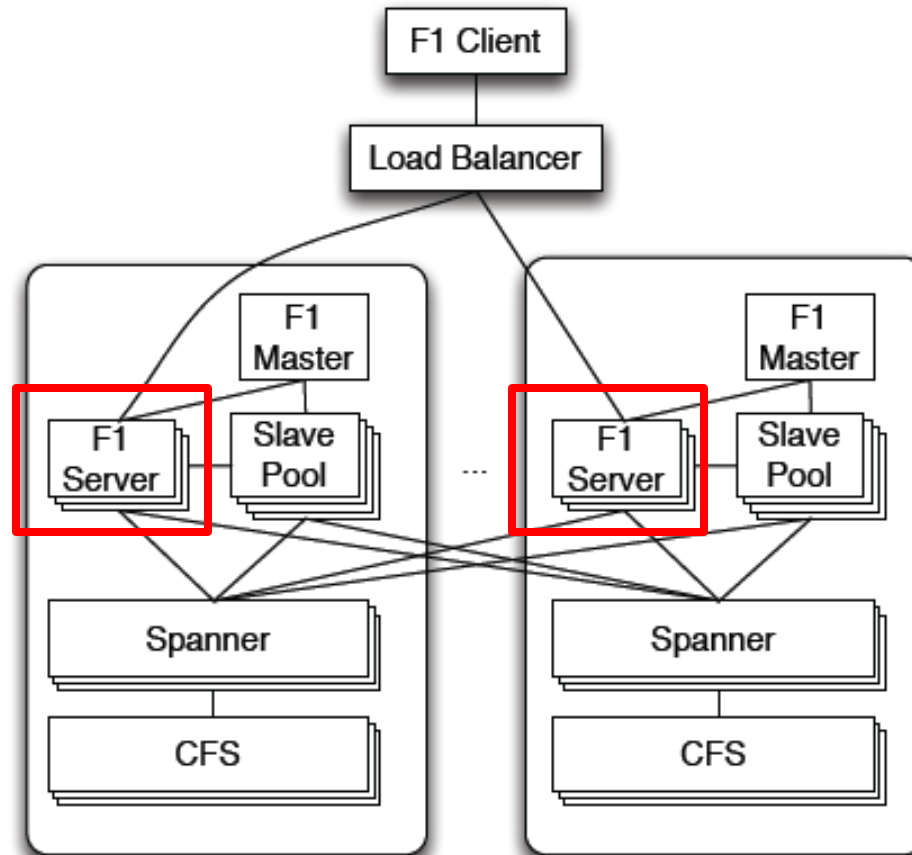


Figure 1: The basic architecture of the F1 system, with servers in two datacenters.

F1 Basic Architecture

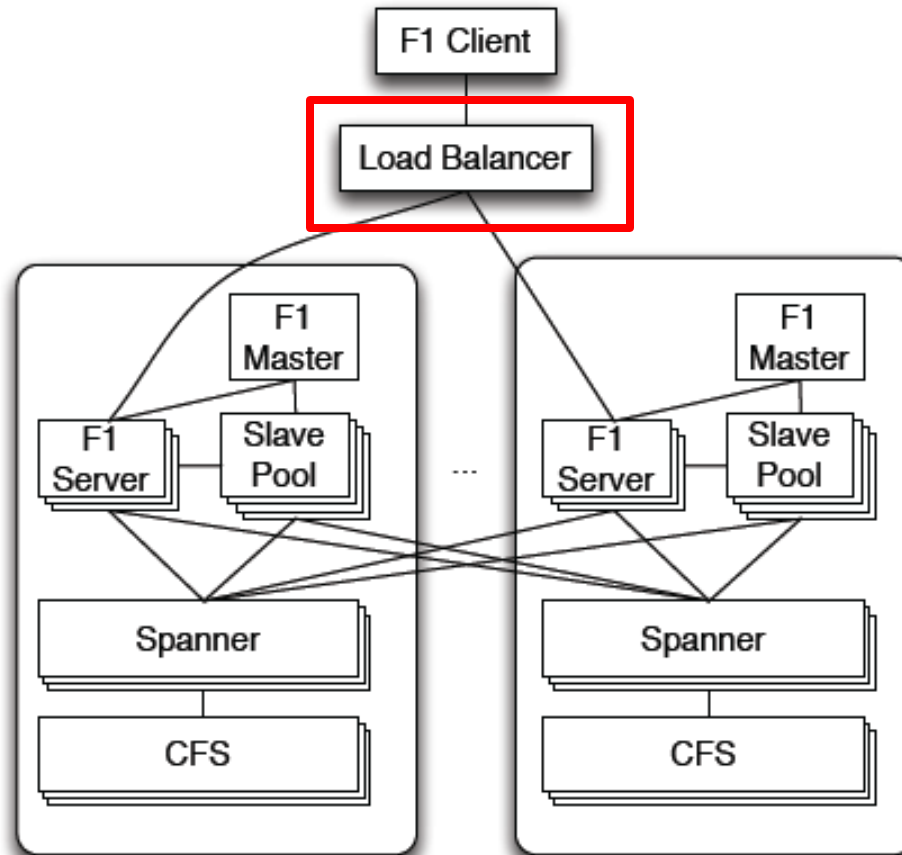


Figure 1: The basic architecture of the F1 system, with servers in two datacenters.

F1 Basic Architecture

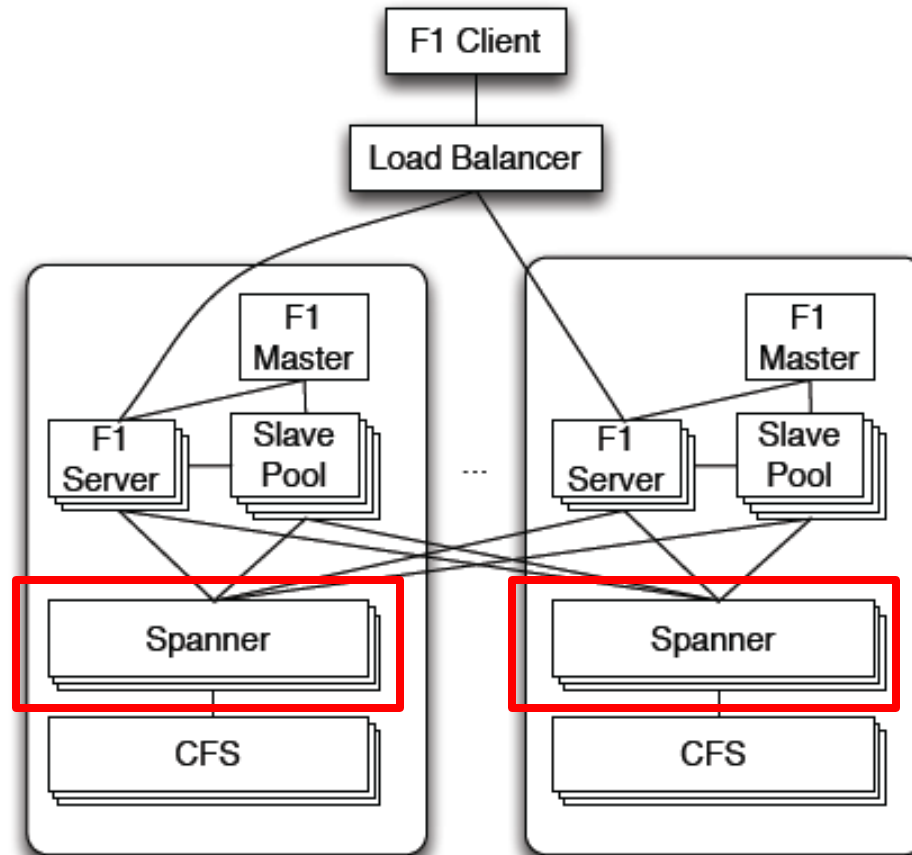


Figure 1: The basic architecture of the F1 system, with servers in two datacenters.

F1 Basic Architecture

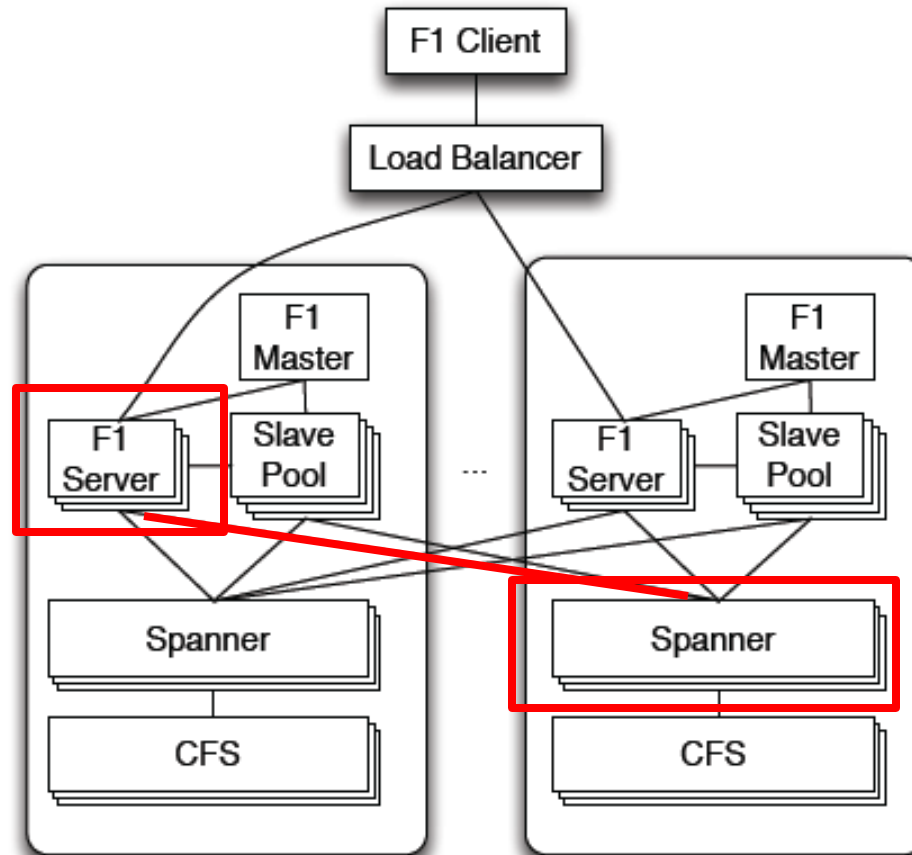


Figure 1: The basic architecture of the F1 system, with servers in two datacenters.

F1 Basic Architecture

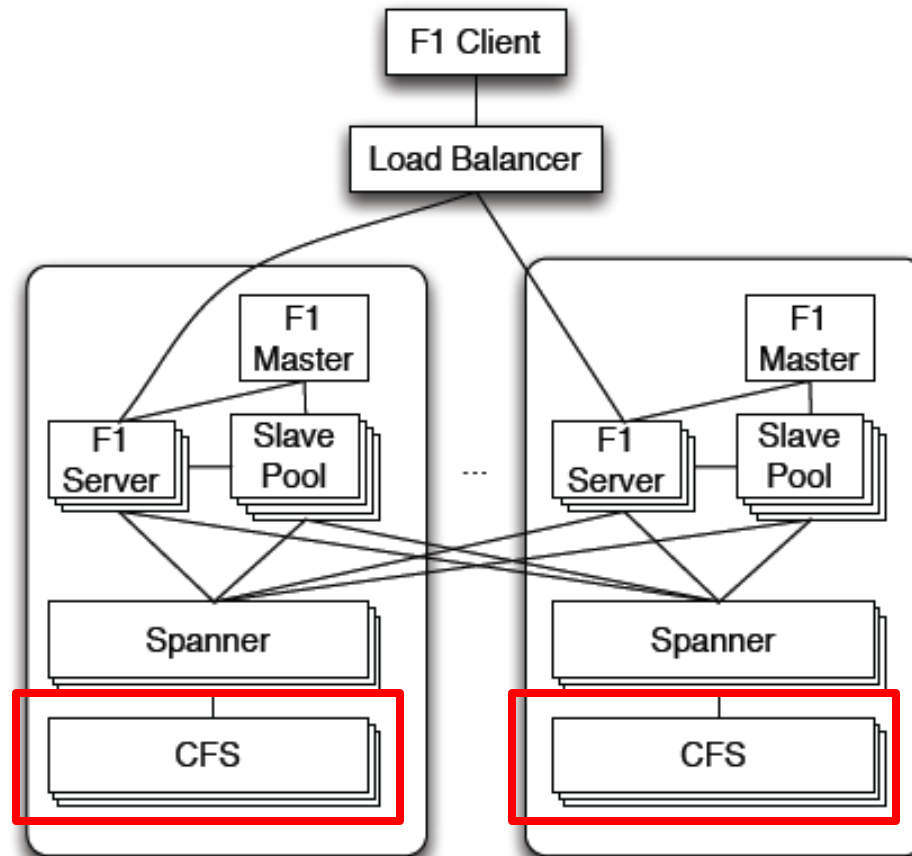


Figure 1: The basic architecture of the F1 system, with servers in two datacenters.

F1 Basic Architecture

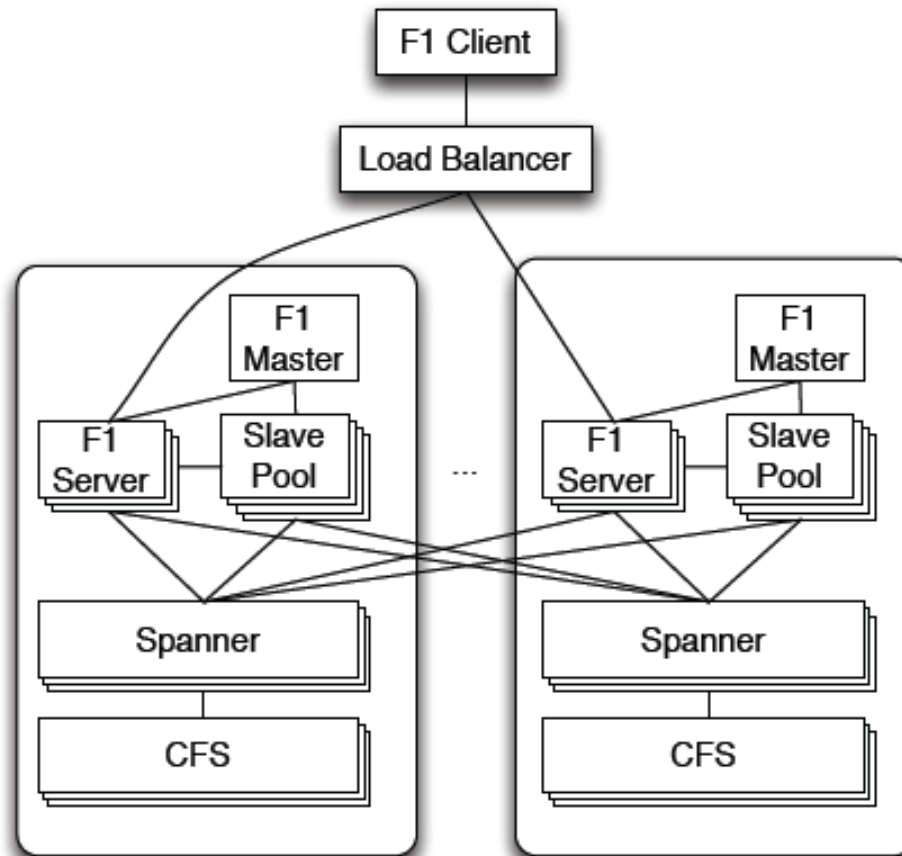


Figure 1: The basic architecture of the F1 system, with servers in two datacenters.

F1 Basic Architecture

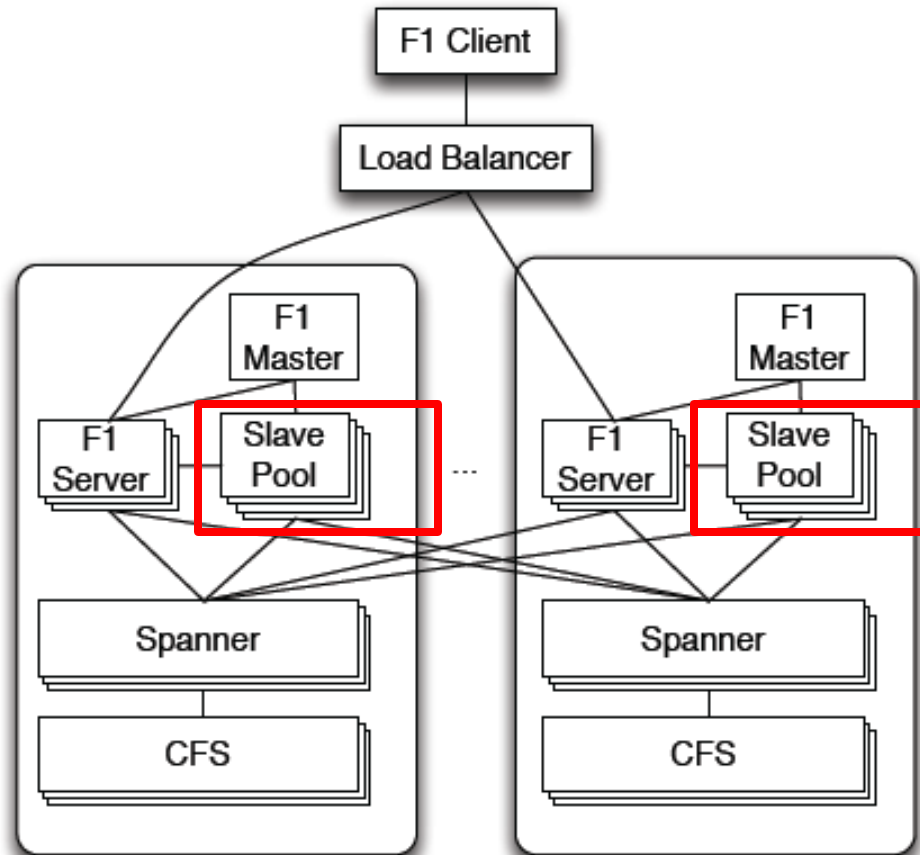


Figure 1: The basic architecture of the F1 system, with servers in two datacenters.

F1 Basic Architecture

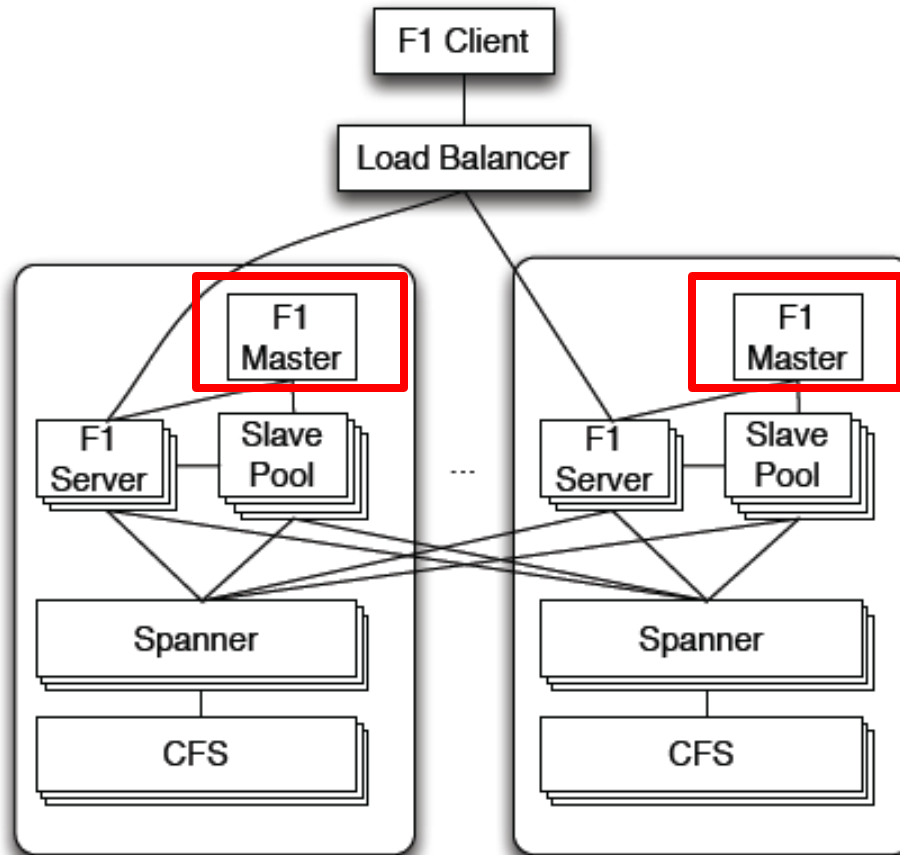


Figure 1: The basic architecture of the F1 system, with servers in two datacenters.

F1 Basic Architecture

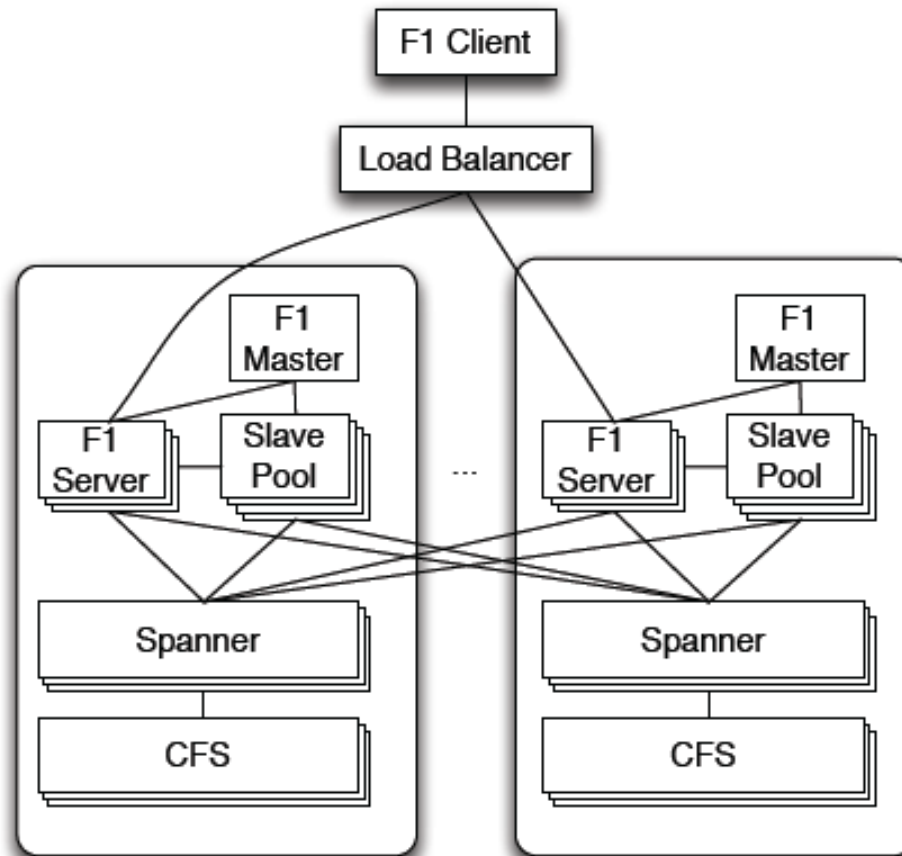


Figure 1: The basic architecture of the F1 system, with servers in two datacenters.

F1 Data Model

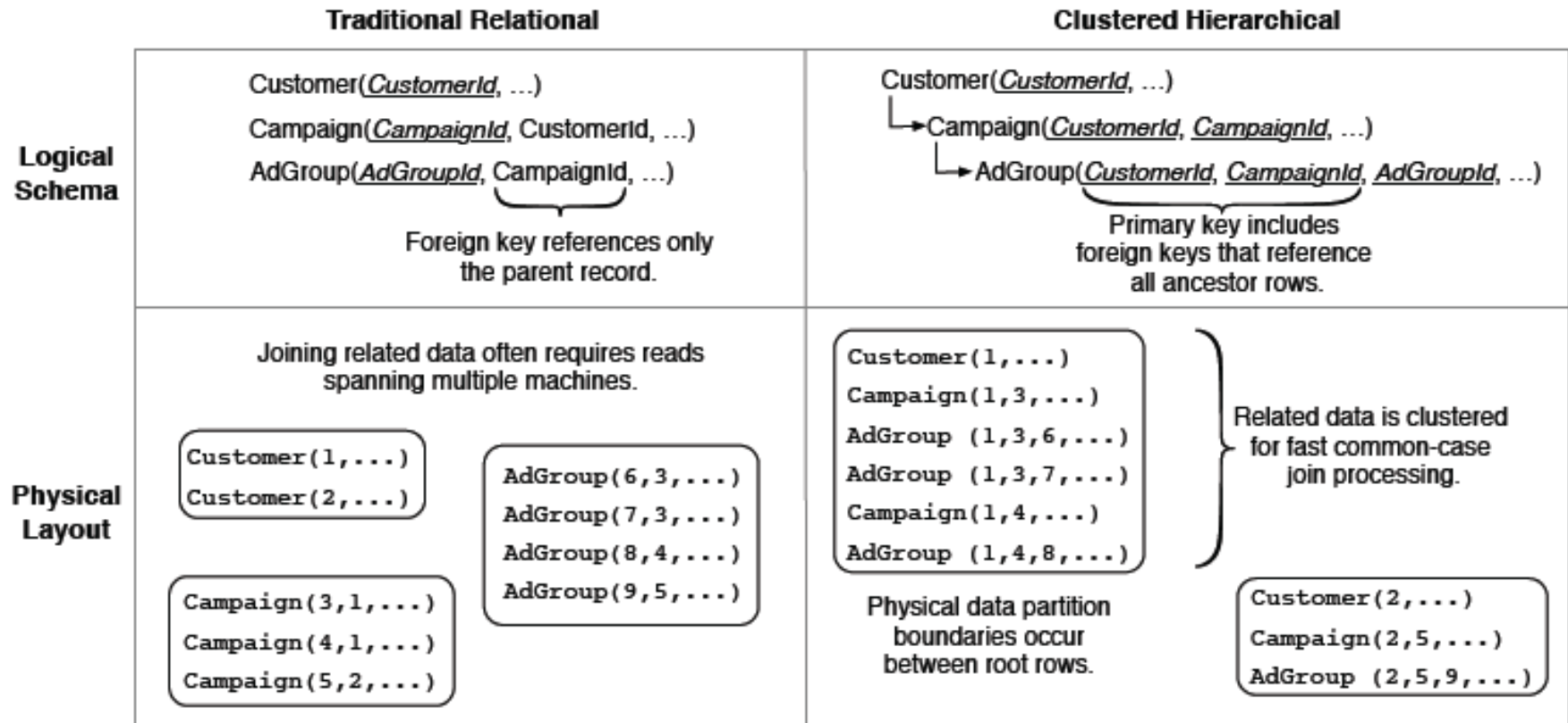


Figure 2: The logical and physical properties of data storage in a traditional normalized relational schema compared with a clustered hierarchical schema used in an F1 database.

F1 Data Model

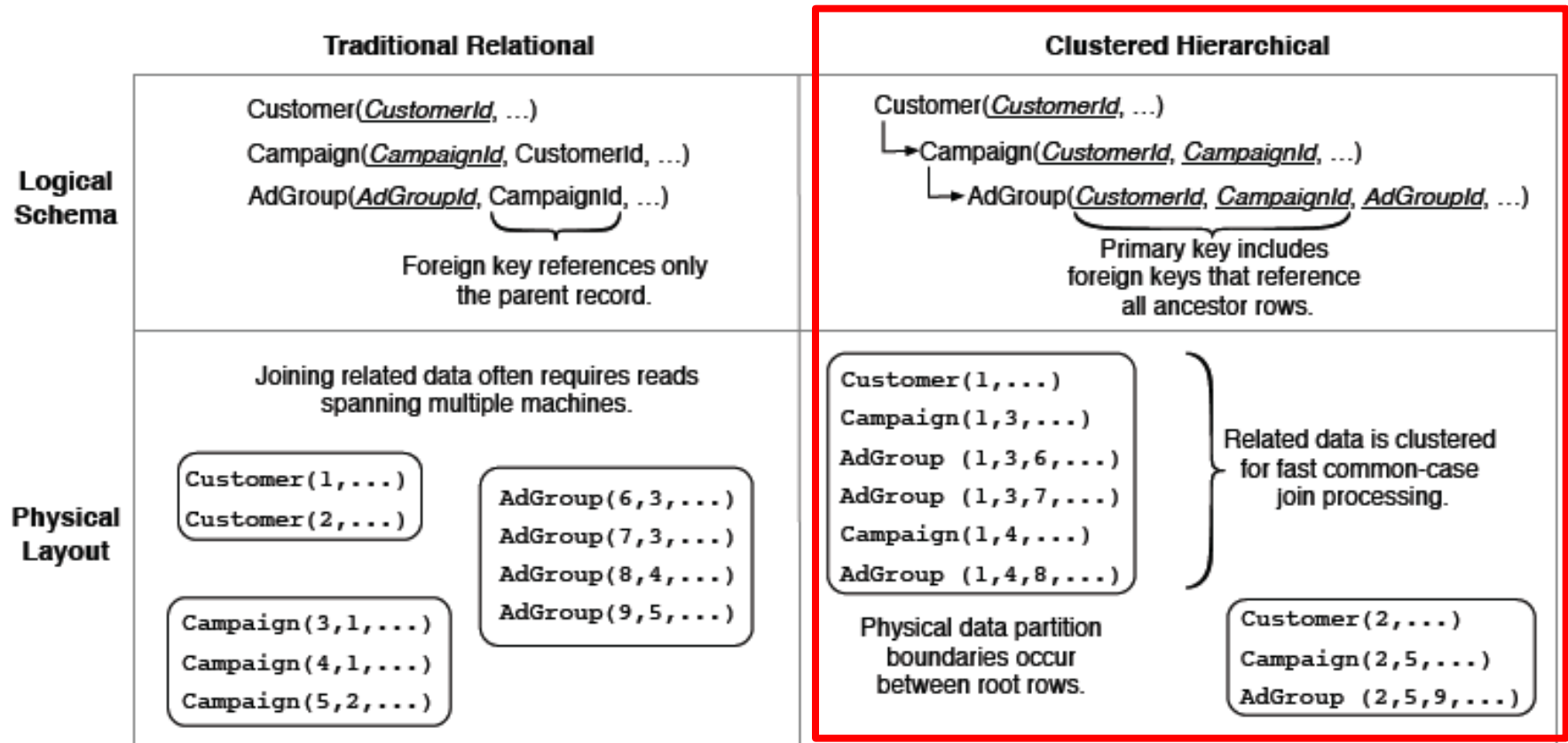


Figure 2: The logical and physical properties of data storage in a traditional normalized relational schema compared with a clustered hierarchical schema used in an F1 database.

F1 Data Model

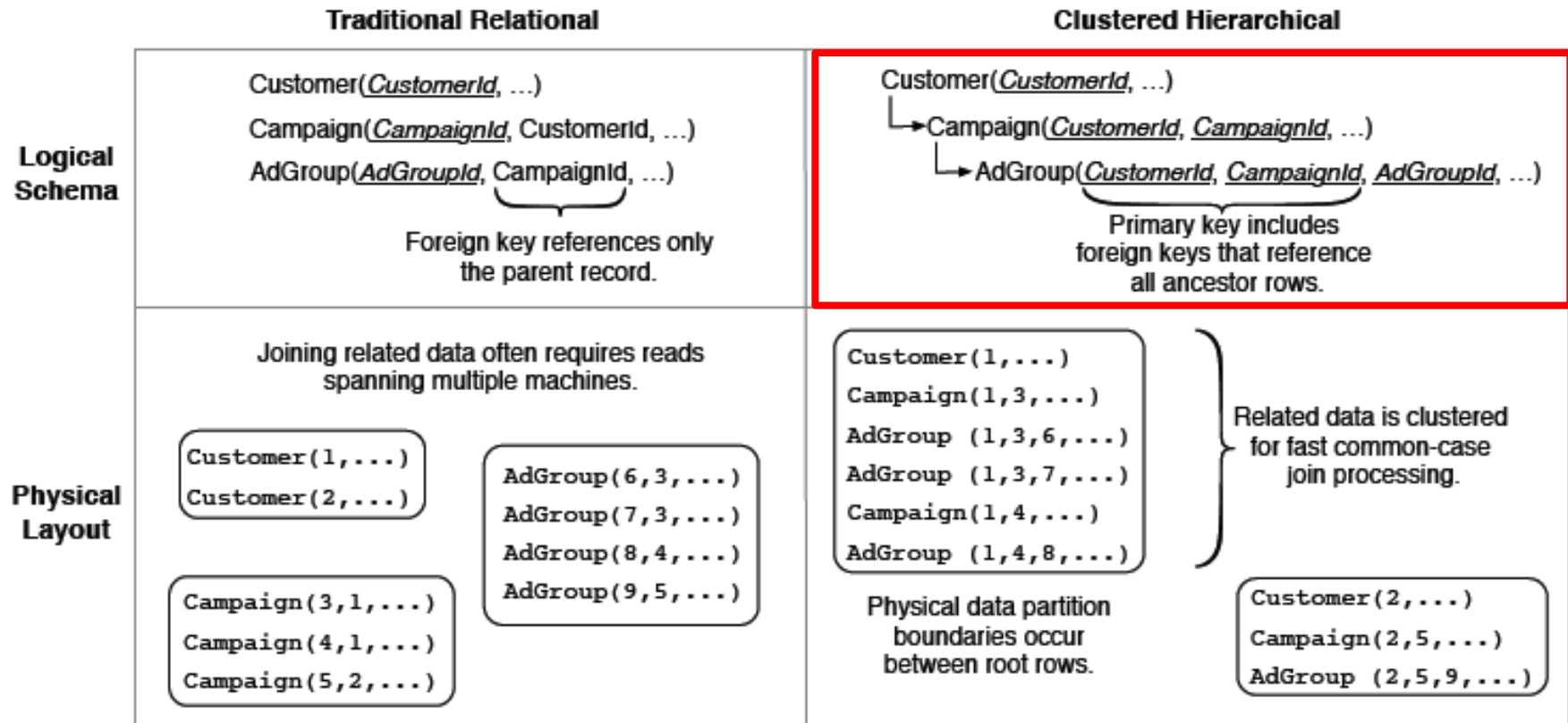


Figure 2: The logical and physical properties of data storage in a traditional normalized relational schema compared with a clustered hierarchical schema used in an F1 database.

F1 Data Model

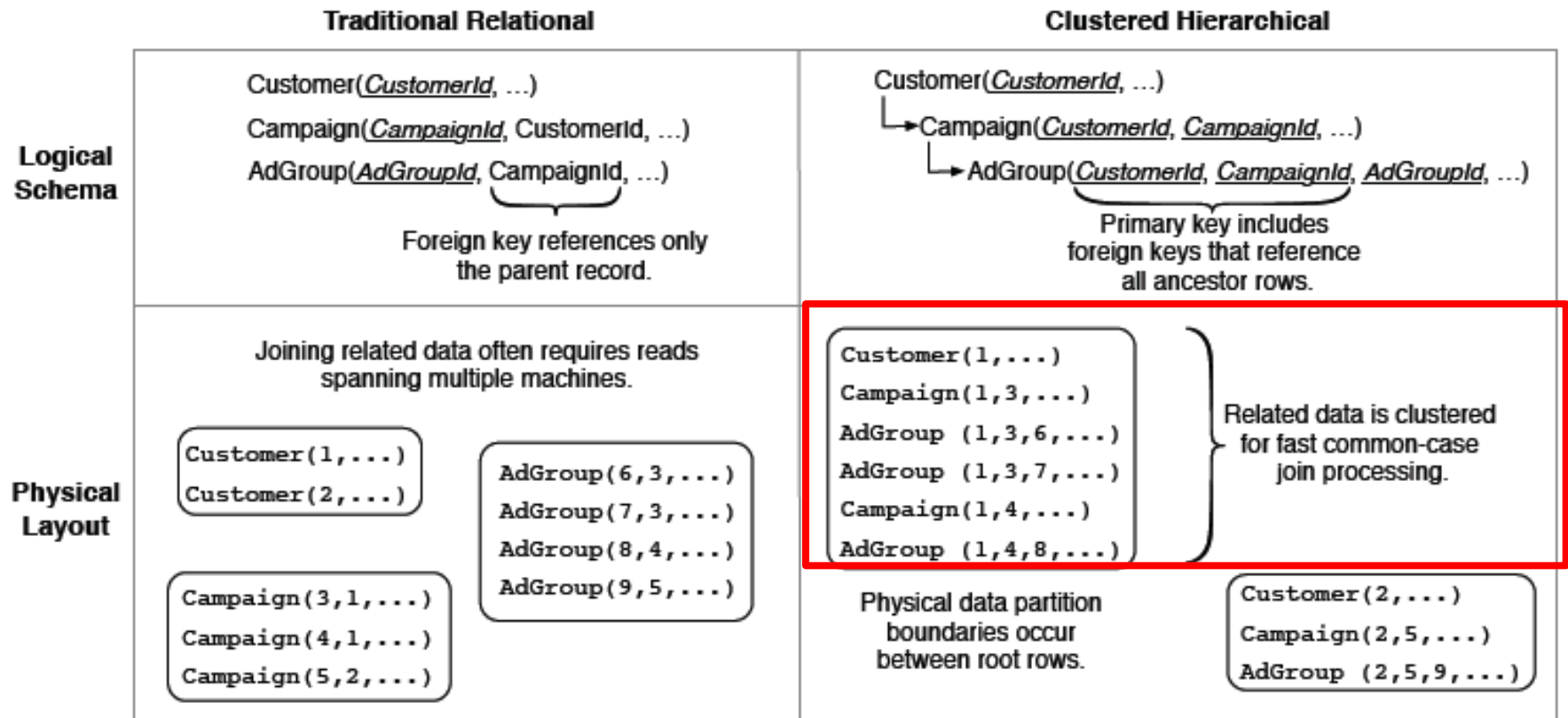


Figure 2: The logical and physical properties of data storage in a traditional normalized relational schema compared with a clustered hierarchical schema used in an F1 database.

F1 Data Model

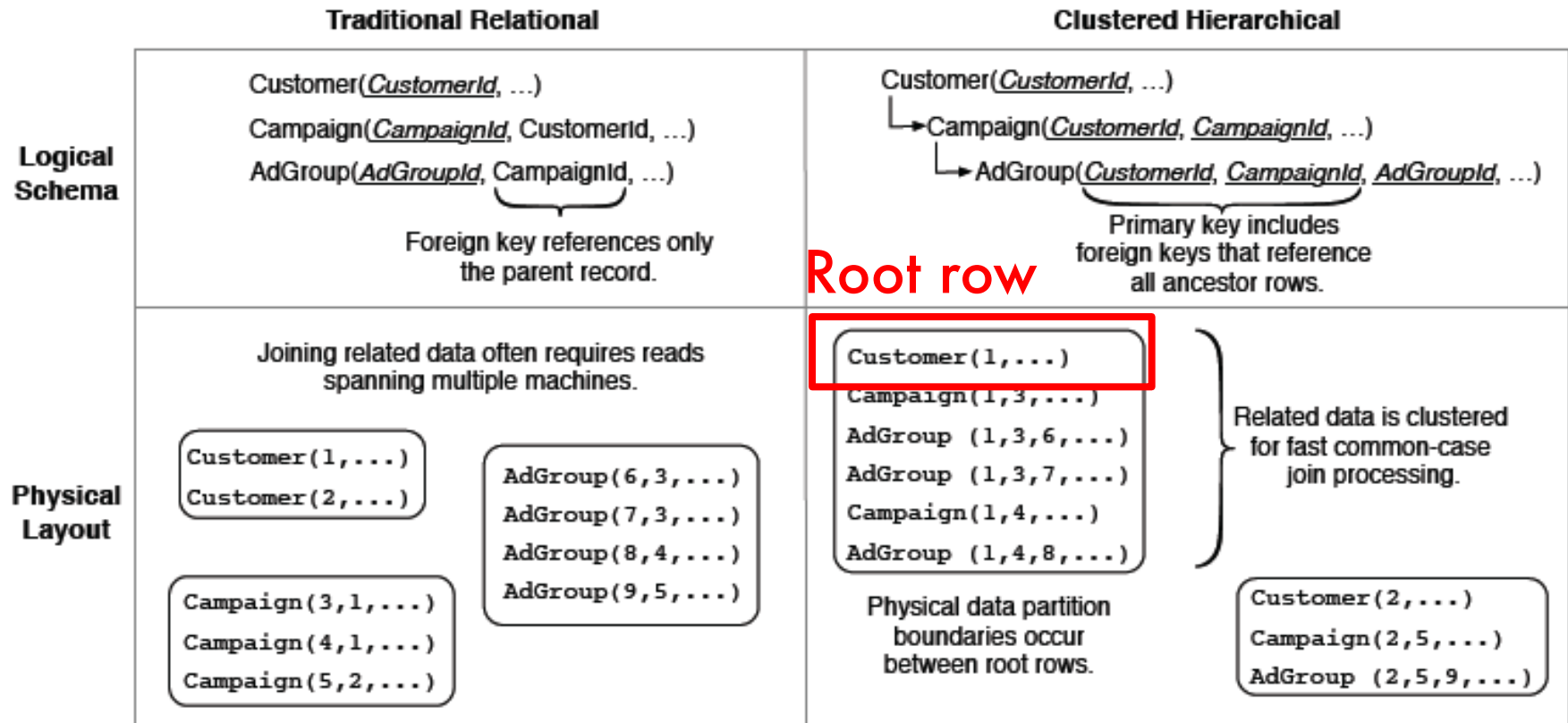


Figure 2: The logical and physical properties of data storage in a traditional normalized relational schema compared with a clustered hierarchical schema used in an F1 database.

F1 Data Model

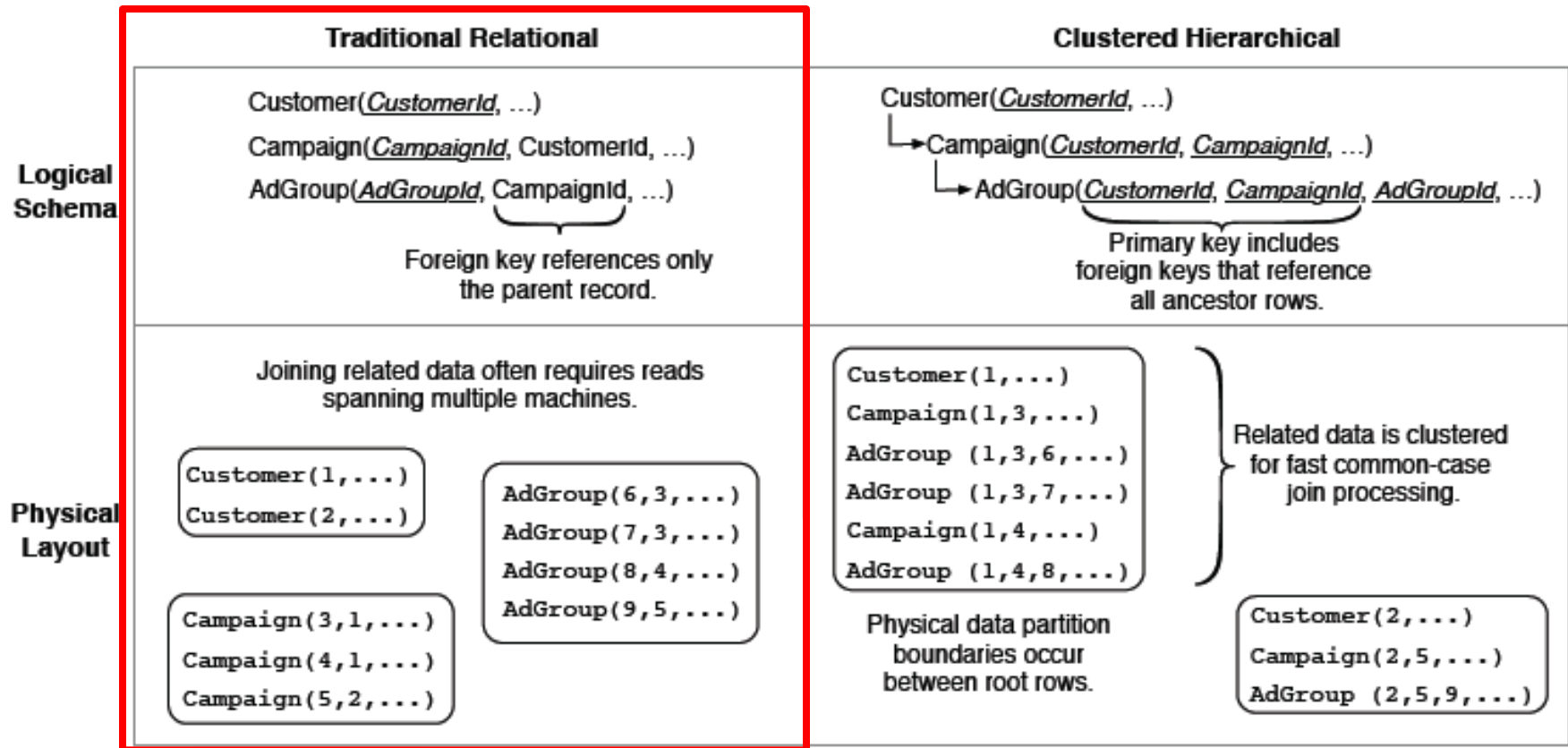


Figure 2: The logical and physical properties of data storage in a traditional normalized relational schema compared with a clustered hierarchical schema used in an F1 database.

F1 Data Model

	Traditional Relational	Clustered Hierarchical
Logical Schema	<p> <code>Customer(<u>CustomerId</u>, ...)</code> <code>Campaign(<u>CampaignId</u>, CustomerId, ...)</code> <code>AdGroup(<u>AdGroupId</u>, CampaignId, ...)</code> </p> <p>Foreign key references only the parent record.</p>	<p> <code>Customer(<u>CustomerId</u>, ...)</code> ↳ <code>Campaign(<u>CustomerId</u>, <u>CampaignId</u>, ...)</code> ↳ <code>AdGroup(<u>CustomerId</u>, <u>CampaignId</u>, <u>AdGroupId</u>, ...)</code> </p> <p>Primary key includes foreign keys that reference all ancestor rows.</p>
Physical Layout	<p>Joining related data often requires reads spanning multiple machines.</p> <div> <div> <code>Customer(1,...)</code> <code>Customer(2,...)</code> </div> <div> <code>AdGroup(6,3,...)</code> <code>AdGroup(7,3,...)</code> <code>AdGroup(8,4,...)</code> <code>AdGroup(9,5,...)</code> </div> </div> <div> <code>Campaign(3,1,...)</code> <code>Campaign(4,1,...)</code> <code>Campaign(5,2,...)</code> </div>	<div> <code>Customer(1,...)</code> <code>Campaign(1,3,...)</code> <code>AdGroup (1,3,6,...)</code> <code>AdGroup (1,3,7,...)</code> <code>Campaign(1,4,...)</code> <code>AdGroup (1,4,8,...)</code> </div> <p>Physical data partition boundaries occur between root rows.</p> <div> <code>Customer(2,...)</code> <code>Campaign(2,5,...)</code> <code>AdGroup (2,5,9,...)</code> </div> <p>Related data is clustered for fast common-case join processing.</p>

Figure 2: The logical and physical properties of data storage in a traditional normalized relational schema compared with a clustered hierarchical schema used in an F1 database.

F1 Data Model

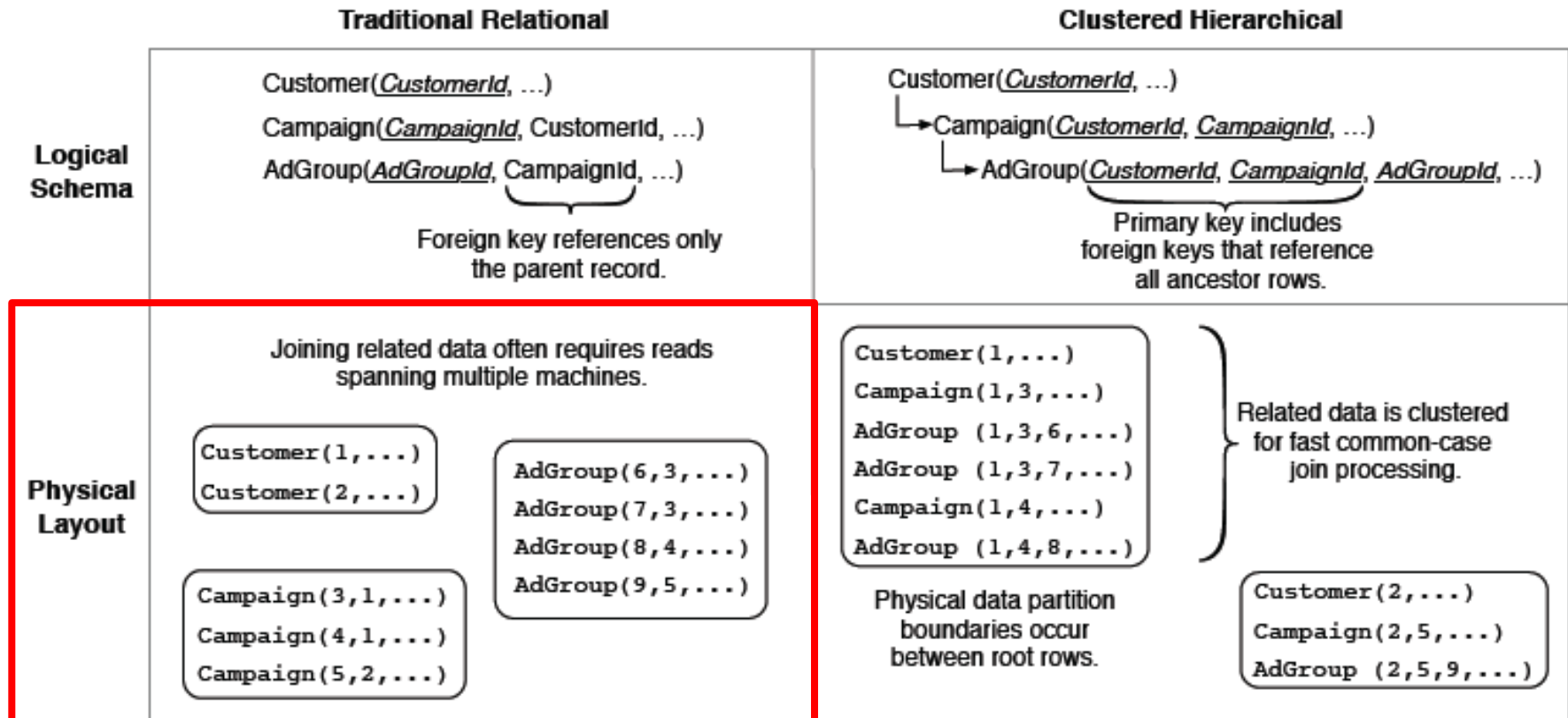


Figure 2: The logical and physical properties of data storage in a traditional normalized relational schema compared with a clustered hierarchical schema used in an F1 database.

Google Core Technology: Summary

In this lecture we've done a very rapid tour through...

- MapReduce
- GFS
- BigTable
- Spanner
- FI

In Lecture 15 we look at the open source equivalents of MapReduce, GFS, and BigTable, all developed by Apache Software Foundation.

It seems highly unlikely that Apache will ever set up projects directly aimed at developing open-source versions of Spanner, at least until there is widespread adoption of TrueTime and the associated GLPS/atomic-clock technology.

And, because FI is built on Spanner (+ TrueTime), direct open-source version of FI is unlikely, but other “NewSQL” database technologies are under development.