



Architecture Overview

Engineering, Security, and
Operations

January 2021



In this guide, you'll learn:

How Canvas has been architected as a native cloud application, providing unmatched availability, scalability, and reliability for our customers.

Table of Contents

Canvas Architecture	3
Hosting Regions	3
Product Security	3
Separation of Tenant Data	4
Architecture and Data Flow	4
Scaling, Backup, Recovery, and Redundancy	5
Load Balancers	5
Application Servers	6
Cache Servers	6
Database Servers	6
Distributed File Storage	7

Canvas Architecture

Canvas is a dynamic Ruby on Rails web application built on cloud-native, multi-tenant architecture capable of automatically scaling to serve tens of millions of users.

HOSTING REGIONS

For US customers, Instructure uses two Amazon Web Services (AWS) regions, ensuring that client data is not stored outside of the United States:

- US East (Northern Virginia) Region with 3 EC2 Availability Zones
- US West (Oregon) Region with 3 EC2 Availability Zones

For international clients, Instructure uses the following AWS regions:

- Canada Central (Montreal) Region with 2 EC2 Availability Zones
- EU West (Ireland) with 3 EC2 Availability Zones
- EU Central (Germany) Region with 2 EC2 Availability Zones
- Asia Pacific (Sydney) Region with 2 EC2 Availability Zones
- Asia Pacific (Singapore) Region with 2 EC2 Availability Zones.

PRODUCT SECURITY

Instructure produces, on an annual basis, a SOC2 Type II report for Canvas covering the following principles: Security, Availability, Confidentiality, Processing Integrity, and Privacy.

As one of the benefits of utilizing AWS cloud infrastructure, we also benefit from the following security certifications:

- SOC 1 Type II (ISAE 3402), SOC 2 Type II, and SOC 3 Type II reports
- ISO 9001, 27001 (CSA Star Level 2), 27017, and 27018 certified
- Level 1 PCI-DSS service provider
- FISMA-Moderate operation level
- GDPR ready, FERPA compliant (shared responsibility model)
- Cyber Essentials PLUS certification

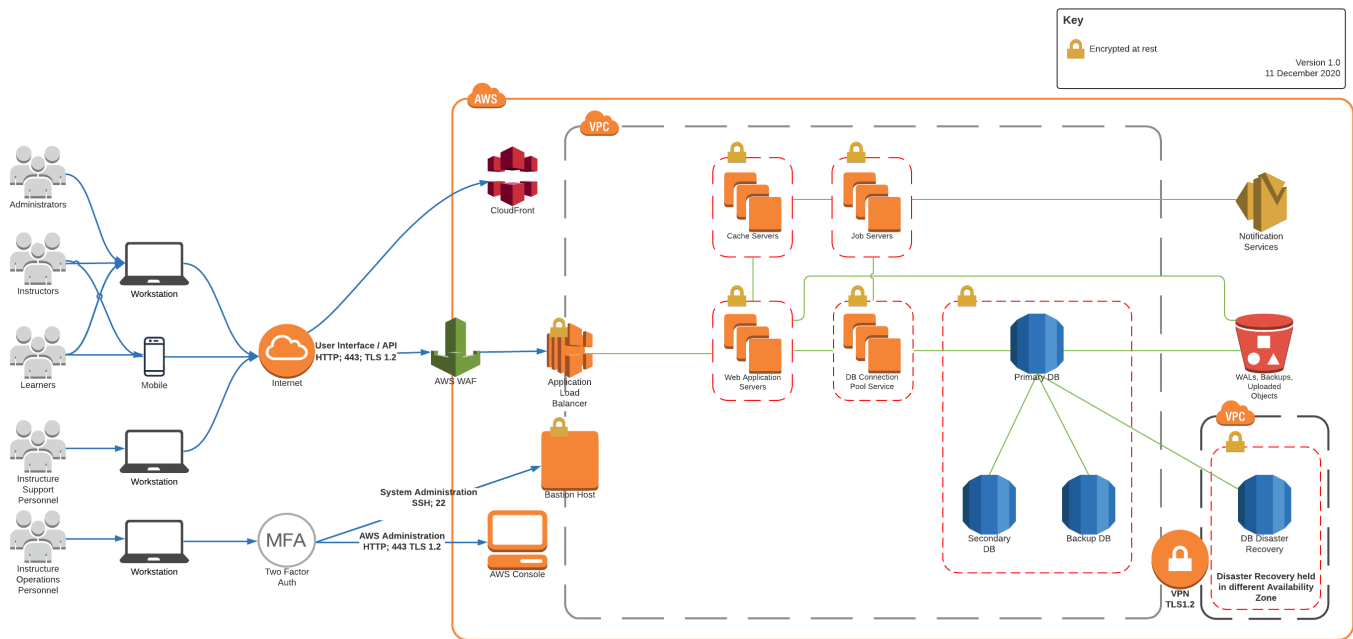


SEPARATION OF TENANT DATA

Separation of tenants is accomplished in AWS via logical separation in natively multi-tenant software. Customer data is segregated via database sharding (horizontal partitioning). Horizontal partitioning is a design principle whereby rows of a database table are held separately, rather than splitting by columns (as for normalization). Each partition forms part of a shard. The advantage is the number of rows in each table is reduced, reducing index size, and improving performance.

Sharding is based on real-world aspect of the data (e.g., segmented by customer) and data cannot leak from one shard to another, nor can clients gain access to data in another shard as the method of inferring the client shard is accomplished after authentication. As client credentials are only valid for a single account, and therefore shard, user authentication is intrinsically tied to the shard identity. Validation of segregated client data occurs during weekly disaster recovery testing.

ARCHITECTURE AND DATA FLOW



Scaling, Backup, Recovery, and Redundancy

AWS data center electrical and network systems are designed to be fully redundant and maintainable without impact to operations, 24 hours a day, seven days a week. Uninterruptible Power Supply (UPS) units are available in the event of an electrical failure for critical and essential loads in the facility. Data centers use generators to provide backup power for the entire facility.

The Canvas architecture replicates data in near real-time and data is backed up on a daily basis. Instructure creates daily offsite database backups of Canvas data and content including course content, student submissions, student-created content, analytics, rubrics, learning outcomes, and metadata. Data is stored redundantly in multiple data centers and multiple geographic locations through Amazon S3.

The Canvas architecture is horizontally scalable and uses a mix of in-house developed and AWS-provided technologies, enabling it to respond to usage spikes in real-time and accommodate expanded, long-term usage. Through automatic scaling and automated provisioning technology, Canvas adjusts cloud resources to handle large usage loads before they cause slowdowns. When concurrent user numbers grow, Canvas automatically adds resources, so users don't experience outages or slowdown.

Assuring the recovery and redundancy of the Canvas platform, we take advantage of multiple geographically separate sites and Availability Zones which provide resilience in the face of most failure modes including natural disasters or system failures. The Canvas application is designed to make full use of the real-time redundancy and capacity capabilities offered by AWS, running across multiple availability zones in regions throughout the world. Primary storage is provided by Amazon S3, which is designed for durability exceeding 99.99999999%.

The Canvas architecture is also resilient to failure and capable of rapid recovery from component failure. The Canvas application, its media and file storage, and its databases are each independently redundant. If an application hosting node were to fail, all traffic would transfer to living nodes. If load increases, an automated provisioning system ensures that more hosting nodes are made available to handle the traffic—either in response to increased load or in predictive anticipation of future workloads. The database and file stores are also horizontally scalable, adding capacity for both additional storage and load as needed.

LOAD BALANCERS

AWS Elastic Load Balancers are deployed in a highly available active/active configuration, which handles incoming requests and dispatches the underlying connections evenly to available application servers. The load balancer maintains a dynamic list of available application servers for dispatch. The load balancer sends regular heartbeats—a simple network message—to verify the application server is healthy, available, and capable of receiving additional work. The load balancer will not dispatch work to unresponsive application servers. Additional capacity is automatically added to the load balancing pool as traffic and demand increases.



APPLICATION SERVERS

Application servers process incoming requests from the load balancers. They are responsible for executing the business logic, rendering HTML, and returning some static assets to the Canvas user's web browser. Additionally, these servers are balanced across multiple availability zones to ensure maximum fault tolerance.

Application servers are constantly monitored individually for load and capacity information. When all application servers reach a certain load threshold, a new application server is automatically provisioned and deployed. Instructure's in-house automation can dynamically and intelligently schedule new application servers in anticipation of high load times, such as during the beginning and end of semesters.

CACHE SERVERS

The caching layer provides performance optimization. A healthy cache means the application servers need to make fewer trips to the database which speeds up response times. The caching layer is made up of numerous machines running Redis. Data is spread out evenly across all machines. Additionally, Amazon Cloudfront (a caching CDN) is used to quickly deliver static assets to Canvas users. These CDN endpoints are globally distributed, thereby making the network path for these requests as efficient as possible.

Cache servers are constantly monitored. When a cache server fails, a new one is provisioned and deployed to take its place. When a cache server fails, the data that would have been stored on it, is simply retrieved from the database instead.

Cache servers are completely memory based. Memory usage is monitored continuously. When the cache hit rates falls below an acceptable threshold, new cache servers are provisioned and deployed.

DATABASE SERVERS

Course and user data are stored in relational databases. The databases are partitioned by client institution for performance and data isolation purposes. Each institution utilizes a pair of databases: A Primary database and a Secondary database in a separate availability zone.

There is also a third Backup server in each region and (if available) in a separate availability zone. All database changes are streamed in real-time to each other, and to a durable data layer (S3). This means Canvas database information for U.S. customers is stored in three separate geographically separated locations. Canadian customers benefit from two separate geographically separated locations. Additionally, database backups (a different form of data redundancy for different purposes) are tested weekly.

If the Primary database fails, the Secondary will be promoted to Primary and a new Secondary database provisioned and deployed. Upon failure of the Secondary database, a new Secondary database is provisioned and deployed. In the unlikely event of simultaneous component failure or data corruption, the standby backup server can be used to create a new database pair.

Databases are constantly monitored for resource usage and response time. If either database approaches peak load, individual customers will be relocated to clusters with available capacity.



DISTRIBUTED FILE STORAGE

Course media, including videos, image files, audio recordings, etc. and student-uploaded files, like assignments, documents, and learning artifacts are stored outside the database in a separate and scalable Amazon Simple Storage Service (S3) bucket that is designed for durability exceeding 99.99999999%. All objects within the S3 buckets are encrypted and replicated between geographically separate sites and have version control enabled so previous versions of an object can be restored with minimal effort.



© 2021 Instructure Inc. All rights reserved.