
Virtual Waiting Room on AWS

Implementation Guide

Virtual Waiting Room on AWS: Implementation Guide

Copyright © 2023 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

Amazon's trademarks and trade dress may not be used in connection with any product or service that is not Amazon's, in any manner that is likely to cause confusion among customers, or in any manner that disparages or discredits Amazon. All other trademarks not owned by Amazon are the property of their respective owners, who may or may not be affiliated with, connected to, or sponsored by Amazon.

Table of Contents

Welcome	1
Cost	2
Daily cost for maintaining the solution without any events	2
Cost for 50,000 waiting room users during 2-hour event	2
Cost for 100,000 waiting room users during 2-hour event	3
Architecture overview	4
How the solution works	6
Solution components	8
Waiting room public and private APIs	8
Authorizers	10
OpenID adapter	10
Sample inlet strategies	11
Sample waiting room	12
Security	14
Monitoring	14
IAM roles	14
Amazon CloudFront	15
Security groups	15
Design considerations	16
Deployment options	16
Supported protocols	16
Waiting room inlet strategies	16
MaxSize	16
Periodic	17
Customizing and extending the solution	17
Quotas	17
Regional deployments	18
AWS CloudFormation templates	19
Automated deployment	21
Prerequisites	21
Deployment overview	21
Step 1. Launch the getting-started stack	21
Step 2. (Optional) Test the waiting room	23
Generate AWS keys to call the IAM secured APIs	23
Open up the sample waiting room's control panel	23
Test the sample waiting room	23
Deploying separate stacks	25
1. Launch the core stack	25
2. (Optional) Launch the Authorizers stack	26
3. (Optional) Launch the OpenID stack	27
4. (Optional) Launch the sample inlet strategy stack	28
5. (Optional) Launch the sample waiting room stack	29
Updating the stack from a previous version	31
Performance data	32
Findings	32
Troubleshooting	33
Additional resources	34
Uninstall the solution	35
Using the AWS Management Console	35
Using AWS Command Line Interface	35
Deleting the Amazon S3 buckets	35
Source code	36
Revisions	37
Contributors	38

Notices	39
AWS glossary	40

Absorb large bursts of traffic to your website with the Virtual Waiting Room on AWS

Publication date: November 2021 ([last update \(p. 37\)](#): April 2023)

The Virtual Waiting Room on AWS solution helps control incoming user requests to your website during large bursts of traffic. It creates a cloud infrastructure designed to temporarily offload incoming traffic to your website, and it provides options to customize and integrate a virtual waiting room. This solution can be integrated with either new or existing websites to seamlessly scale to handle sudden surges in traffic.

Examples of large-scale events which could produce a surge in website traffic include:

- Start of sale for concert or sporting event tickets
- Fire sale or other large retail sale, such as Black Friday
- New product launch with broad marketing announcements
- Exam access and class attendance for online testing and lessons
- Release of medical appointment slots
- Launch of a new direct-to-customer service that requires account creation and payments

The solution acts as a holding area for visitors to your website and allows traffic to pass through when there is enough capacity. The client software used by visitors can be configured to transparently allow traffic through the waiting room until the website is at maximum capacity; at which point the waiting room holds back visitors. When your website has capacity for more traffic, the solution generates [JSON Web Tokens](#) (JWT) that allow users to access the website. For example, if you have an event that lasts for two hours and your website can process 50 users per second, but you expect volume of 250 per second, then you can use this solution to regulate the traffic while allowing users to keep their position in the queue.

This solution provides the following key features:

- Structured queuing of users into your website
- Scalability to control traffic for very large event sizes
- JSON web token generation to allow entry to the target site
- All functionality is controlled through REST APIs
- Turnkey API Gateway authorizer for client solutions
- Standalone integration or use with OpenID

This implementation guide describes architectural considerations and configuration steps for deploying Virtual Waiting Room on AWS in the Amazon Web Services (AWS) Cloud. It includes links to [AWS CloudFormation](#) templates that launch and configure the AWS services required to deploy this solution using AWS best practices for security and availability.

The guide is intended for IT architects, developers, DevOps staff, data analysts, and marketing technology professionals who have practical experience architecting in the AWS Cloud.

Cost

You are responsible for the cost of the AWS services used while running this solution. As of this revision, the cost for running this solution with the default settings in the US East (N. Virginia) Region is approximately **\$10.00/day per stack plus charges for API requests and data traffic relative to the size of the event.**

Daily cost for maintaining the solution without any events

AWS service	Requests/Time	Cost
Amazon API Gateway	0	\$0.00
Amazon CloudFront	0	\$0.00
Amazon CloudWatch	0	\$0.00
Amazon DynamoDB	0	\$0.00
Amazon ElastiCache	Compute node hours (Redis)	~\$6.00
AWS Lambda	Free tier*	\$0.00
AWS Secrets Manager	Free tier*	\$0.00
Amazon Simple Storage Service (Amazon S3)	Free tier*	\$0.00
Amazon Virtual Private Cloud (Amazon VPC)	VPC endpoint hours	~\$4.00
TOTAL:		~\$10.00

*The cost estimate is based upon a clean environment. If you are using this AWS service outside of this solution, you might exceed the free tier quota.

The following tables show estimated costs for a **50,000-user** and a **100,000-user waiting room** with an **event duration ranging 2-4 hours** with 500 users/second incoming and 1,000 users/min outgoing. Prices are subject to change. For full details, refer to the pricing webpage for each AWS service used in this solution.

Estimated cost for 50,000 waiting room users during 2-hour event

AWS service	Dimensions	Cost
-------------	------------	------

Amazon API Gateway	Requests	\$2.00
CloudFront	Requests, Bandwidth	\$75.00
CloudWatch	Metrics, Alarms, Storage	\$1.00
Amazon CloudWatch Events	Events	\$1.00
DynamoDB	Read/Write units, Storage	\$1.00
ElastiCache	Node hours	\$8.00
Lambda	Requests, Compute time	\$1.00
AWS Secrets Manager	Secrets, Requests	\$1.00
Amazon S3	Requests, Storage	\$1.00
Amazon VPC	Data transfer, Endpoint time	\$2.00
TOTAL		\$94.00

Estimated cost for 100,000 waiting room users during 2-hour event

AWS service	Dimensions	Cost
Amazon API Gateway	Requests	\$4.00
CloudFront	Requests, Bandwidth	\$296.00
CloudWatch	Metrics, Alarms, Storage	\$1.00
CloudWatch Events	Events	\$1.00
DynamoDB	Read/Write units, Storage	\$4.00
ElastiCache	Node hours	\$32.00
Lambda	Requests, Compute time	\$1.00
AWS Secrets Manager	Secrets, Requests	\$1.00
Amazon Simple Queue Service (Amazon SQS)	Requests	\$1.00
Amazon S3	Requests, Storage	\$1.00
Amazon VPC	Data transfer, Endpoint time	\$6.00
TOTAL		\$348.00

Architecture overview

Deploying this solution with the required and optional templates, using default parameters, builds the following environment in the AWS Cloud.

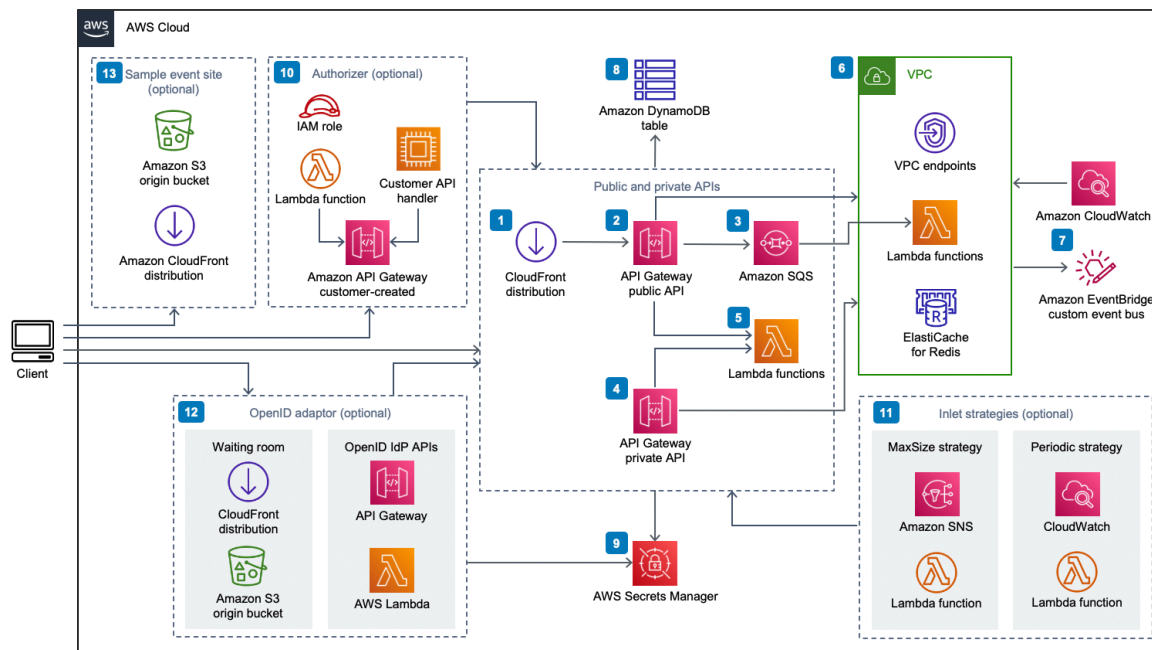


Figure 1: Virtual Waiting Room on AWS architecture

The AWS CloudFormation templates deploy the following infrastructure:

1. An [Amazon CloudFront](#) distribution to deliver public API calls for the client.
2. [Amazon API Gateway](#) public API resources to process queue requests from the virtual waiting room, track the queue position, and support validation of tokens that allow access to the target website.
3. An [Amazon Simple Queue Service](#) (Amazon SQS) queue to regulate traffic to the [AWS Lambda](#) function that processes the queue messages. Instead of invoking the Lambda function for each request, the SQS queue batches the incoming bursts of requests.
4. API Gateway private API resources to support administrative functions.
5. Lambda functions to validate and process public and private API requests, and return the appropriate responses.
6. [Amazon Virtual Private Cloud](#) (VPC) to host the Lambda functions that interact directly with the [Amazon ElastiCache for Redis](#) cluster. VPC endpoints allow Lambda functions in the VPC to communicate with services within the solution.
7. An [Amazon CloudWatch](#) rule to invoke a Lambda function that works with a custom [Amazon EventBridge](#) bus to periodically broadcast status updates.
8. [Amazon DynamoDB](#) tables to store token, queue position, and serving counter data.
9. [AWS Secrets Manager](#) to store keys for token operations and other sensitive data.
- 10(Optional) Authorizer component consisting of an [AWS Identity and Access Management](#) (IAM) role and a Lambda authorizer function for use with API Gateway.
- 11(Optional) [Amazon Simple Notification Service](#) (Amazon SNS), CloudWatch, and Lambda functions to support two inlet strategies.

- 12(Optional) OpenID adaptor component with API Gateway and Lambda functions to allow an OpenID provider to authenticate users to your website. CloudFront distribution with an [Amazon Simple Storage Service](#) (Amazon S3) bucket for the waiting room page for this component.
- 13(Optional) CloudFront distribution with Amazon S3 origin bucket for the sample waiting room web application.

How the solution works

This section describes the steps in an AWS Virtual Waiting Room workflow at a high level. Refer to the [Developer Guide on GitHub](#) for details about building, customizing, and integrating a waiting room for your website.

The waiting room's public API can be located behind your site perimeter security or it can be available without any authorization. Depending on which approach you use for integrating the waiting room with the website, the user might be required to first authenticate to the website before being allowed to navigate to the waiting room and obtain a position in the queue.

Client software must have the Event ID to enter the waiting room and make other requests. An Event ID is a unique ID required for most requests against the public and private APIs. The Event ID is set during installation of the core API stack. During operation, the Event ID can be provided as a URL parameter or cookie via the waiting room page; it can be provided as part of authentication token claims or it can be distributed to clients through a different data path.

There are cases where the client needs both the Event ID and Request ID to make certain API calls. The Request ID is a unique ID issued from the waiting room representing a specific client in line.

The following steps describe the flow of API requests for queue entry, waiting for the queue to progress, and exiting the waiting room with an access token for the website.

User enters the waiting room:

1. The user is presented with a screen or page that represents the waiting room entry point. They choose to enter the queue and the client software (browser, mobile, device) calls the `assign_queue_num` public API to request a queue position.
2. The API request is immediately delivered to the Amazon SQS queue by API Gateway.
3. The `assign_queue_num` API call returns when the request is placed into the queue. The client receives a unique Request ID that can be used later to retrieve the queue position, time of the request, and an access token.
4. The `AssignQueueNum` Lambda function receives batches of up to ten requests from the SQS queue. Lambda service fans out invocations to process multiple batches of requests.
5. The `AssignQueueNum` Lambda function validates each message in its batch, increments the queue counter in ElastiCache for Redis, and stores each request in ElastiCache for Redis with its associated queue position.
6. Each message is deleted as it is successfully processed. Messages involved in an error condition are reprocessed once in a later batch. After a second failure, they are sent to a dead-letter-queue connected to a [CloudWatch alarm](#).
7. The client can begin polling the `queue_num` API after it receives the Request ID from the `assign_queue_num` call. The client sends the Event ID and Request ID to the `queue_num` API and receives a numeric queue position or a response indicating the request hasn't been processed yet. The client might need to make this call more than once during large events. The `GetQueueNum` Lambda function is invoked by API Gateway and returns the client's numeric position in the queue from DynamoDB.

User waits in the waiting room:

8. After the client has its position in the queue, it can begin polling the `serving_num` API at a regular interval. The `serving_num` API is called with the Event ID and returns the current serving position of the queue. The response from the `serving_num` API tells the client when they can move from the

waiting room to the actual target site where the final transaction can occur. The `GetServingNum` Lambda function returns the current serving position of the waiting room.

9. When the serving position is equal to or greater than the client's queue (request) position, the client can request a JSON Web Token (JWT) from the public API. The token can be used with the target site to finalize the transaction. The `generate_token` API is called with the Event ID and Request ID. API Gateway invokes the `GenerateToken` Lambda function with the parameters.
10. The `GenerateToken` Lambda function validates the request and checks if this token has been previously generated. The Lambda function queries the DynamoDB table for a matching token. If found, that token is returned to the caller and it is not regenerated. This process prevents a single Request ID from being used to generate multiple, different tokens with new expiration times.
11. If the token is not found in DynamoDB, the Lambda function retrieves keys to create the token and saves the token in DynamoDB with the Event ID and the client's Request ID. The Lambda function writes an event to EventBridge to signal that a new token was generated. The Lambda function increments an ElastiCache for Redis counter that keeps track of the number of tokens generated for the event.
12. If `queue_pos_expiry` is turned on, the client can query the remaining time before its expiry by calling the `queue_pos_expiry` API which invokes the `GetQueuePositionExpiryTime` Lambda function.

User leaves the waiting room:

13. When the client receives its token, it enters the target site to begin its transaction. Depending on how your infrastructure supports an integration with JWT, the client may need to present the token in a request header, a cookie, or another way. The authorizer for API Gateway can be used to validate the token included in a client's request. Any commercial or open-source libraries for validating and managing JWTs can be used with Virtual Waiting Room on AWS tokens. If the token is valid, the client is allowed to continue their transaction.
14. After the client completes their transaction, a private API is called to update the status of the client's token and is completed in DynamoDB.

Queue position expiry:

15. When this feature is activated, the Request ID corresponding to a particular queue position is eligible to generate a token only for a specified time interval.

Increment serving counter on queue position expiry:

16. When this feature is activated, the serving counter is automatically incremented based on expired queue positions that were unable to generate tokens.

Solution components

Waiting room public and private APIs

The Virtual Waiting Room on AWS solution's primary purpose is to control the generation of JSON Web Tokens (JWT) for clients in a controlled way to avoid bursts of new users that might overwhelm the destination website. The JWTs can be used for site protection, preventing access to web pages until the waiting room token is obtained, and also for API access authorization.

The core template installs a public API and private (IAM-authorized) API used for most Virtual Waiting Room on AWS operations. The public API is configured with a CloudFront distribution with multiple caching policies based on the API path. A DynamoDB table and EventBridge event bus are created. The template adds a new VPC with two Availability Zones (AZs), an ElastiCache for Redis cluster in both AZs, and several Lambda functions. Lambda functions that interact with ElastiCache for Redis have network interfaces within the VPC and all other Lambda functions have default network connectivity. The core APIs are the lowest layer of interaction with the solution. Other Lambda functions, Amazon Elastic Compute Cloud (Amazon EC2) instance, and containers can act as extensions and call the core APIs to build waiting rooms, control inlet traffic, and react to events generated from the solution.

In addition, the core stack creates an alarm for all of its Lambda function errors and throttle conditions, as well as alarms for each API Gateway deployment for 4XX and 5XX status codes.

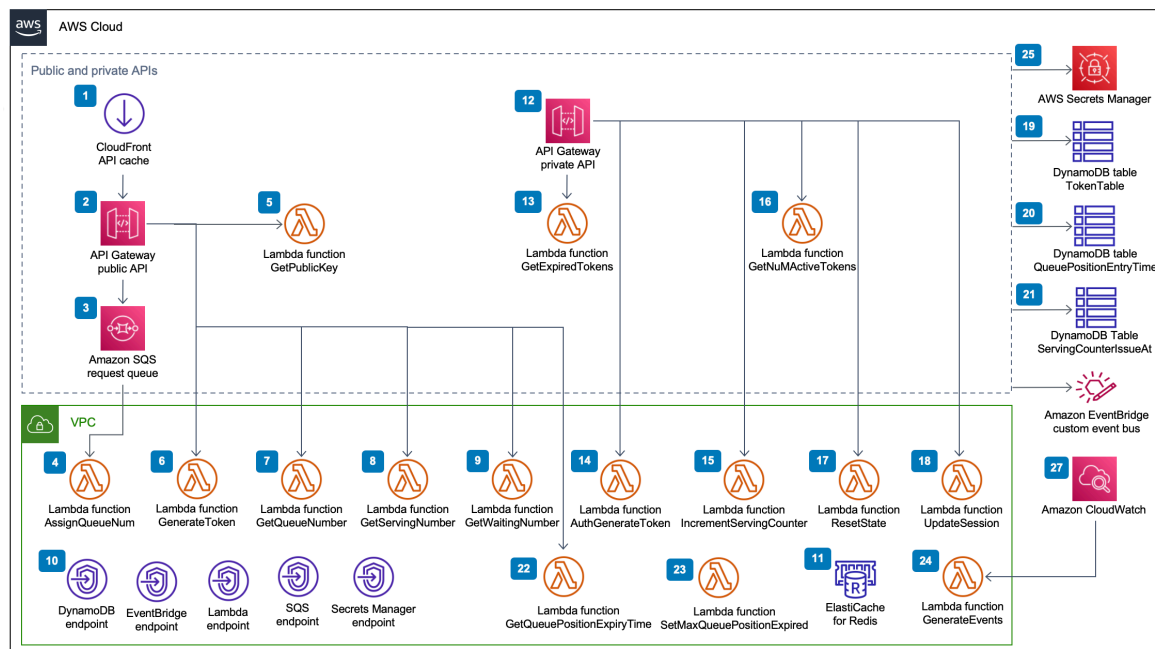


Figure 2: Virtual Waiting Room on AWS Public and private APIs component

1. CloudFront distribution delivers public API calls for the client and caches result where appropriate.
2. Amazon API Gateway public API process queue requests from the virtual waiting room, track the queue position, and support validation of tokens that allow access to the target website.
3. SQS queue regulates traffic to the AWS Lambda function that processes the queue messages.

4. The `AssignQueueNum` Lambda function validates each message in its batch received, increments the queue counter in `ElastiCache` for Redis, and stores each request in `ElastiCache` for Redis with its associated queue position.
5. The `GetPublicKey` Lambda function retrieves the public key value from `Secrets Manager`.
6. The `GenerateToken` Lambda function generates a JWT for a valid request that has been allowed to complete its transaction at the target site. It writes an event to the waiting room's custom event bus that a token has been generated. If a token has been previously generated for this request, no new token is generated.
7. The `GetQueueNumber` Lambda function retrieves and returns the client's numeric position in the queue from `ElastiCache` for Redis.
8. The `GetServingNumber` Lambda function retrieves and returns the number currently being served by the waiting room from `ElastiCache` for Redis.
9. The `GetWaitingNum` Lambda function returns the number currently queued in the waiting room and have yet to be issued a token.
10. VPC endpoints allow Lambda functions in the VPC to communicate with services within the solution.
11. `ElastiCache` for Redis cluster stores all requests to enter the waiting room with a valid Event ID. It also stores several counters like number of requests enqueued, number currently being served, number of tokens generated, number of sessions completed, and number of sessions abandoned.
12. API Gateway private API resources to support administrative functions. The private APIs are AWS IAM authenticated.
13. The `GetExpiredTokens` Lambda function returns a list of request IDs with expired tokens.
14. The `AuthGenerateToken` Lambda function generates a token for a valid request that has been allowed to complete its transaction at the target site. The issuer and the validity period of a token initially set during the core stack deployment can be overridden. It writes an event to the waiting room's custom event bus that a token has been generated. If token has been previously generated for this request, no new token is generated.
15. The `IncrementServingCounter` Lambda function increments the waiting room's serving counter stored in `ElastiCache` for Redis given an increment by value.
16. The `GetNumActiveTokens` Lambda function queries `DynamoDB` for the number of tokens that have yet to expire, have not been used to complete its transaction, and have not been marked abandoned.
17. The `ResetState` Lambda function resets all counters stored in `ElastiCache` for Redis. It also deletes and recreates the `TokenTable`, `QueuePositionEntryTime`, and `ServingCounterIssuedAt` `DynamoDB` tables.
18. The `UpdateSession` Lambda function updates the status of a session (token) stored in the `TokenTable` `DynamoDB` table. Session status is denoted by an integer. Sessions set to a status of 1 indicates completed, and -1 indicates abandoned. It writes an event to the waiting room's custom event bus that a session has been updated.
19. The `TokenTable` `DynamoDB` table stores token data.
20. The `QueuePositionEntryTime` `DynamoDB` table stores queue position and entry time data.
21. The `ServingCounterIssuedAt` `DynamoDB` table stores updates to the serving counter.
22. The `GetQueuePositionExpireTime` Lambda function is invoked when the client requests the remaining queue position expiry time.
23. The `SetMaxQueuePositionExpired` Lambda function sets the maximum queue position that has expired corresponding to the `ServingCounterIssuedAt` table values. It runs every minute if the `IncrSvcOnQueuePositionExpiry` parameter is set to `true` during core stack deployment.
24. The `GenerateEvents` Lambda function writes various waiting room metrics to the waiting room's custom event bus. It gets run every minute if the `Enable Events Generation` parameter is set to `true` during core stack deployment.
25. AWS `Secrets Manager` stores keys for token operations and other sensitive data.
26. Amazon `EventBridge` custom event bus receives an event every time a token is generated and a session is updated in the `TokenTable` `DynamoDB` table. It also receives events when the serving

counter is moved in the `SetMaxQueuePositionExpired` Lambda. It gets written to with various waiting room metrics, if activated during core stack deployment.

27 Amazon CloudWatch event rule is created if the `Enable Events Generation` parameter is set to true during core stack deployment. This event rule initiates the `GenerateEvents` Lambda function every minute.

Authorizers

The solution includes an API Gateway Lambda authorizers stack. The stack consists of one IAM role and a Lambda function. The `APIGatewayAuthorizer` Lambda function is an authorizer for API Gateway that can validate the signature and claims of a token issued by the Virtual Waiting Room on AWS API. The Lambda function supplied with the stack can be used to protect cloud APIs until a user has progressed through the waiting room and receives an access token. The authorizer automatically retrieves and caches the public key and configuration from the core API for token verification. It can be used without modification and can be installed in any AWS Region that supports AWS Lambda.

OpenID adapter

The [OpenID adapter](#) stack deploys an API Gateway and Lambda functions that act as an OpenID identity provider. The OpenID adapter provides a set of OIDC-compatible APIs that can be used with existing web hosting software that support OIDC identity providers, such as AWS Elastic Load Balancers, WordPress, or as a federated identity provider for Amazon Cognito or similar service. The adapter allows a customer to use the waiting room in the AuthN/AuthZ flow when using off-the-shelf web hosting software with limited integration options. The stack also installs a CloudFront distribution with one Amazon S3 bucket as an origin and another S3 bucket for logging requests. The OpenID adapter serves a sample waiting room page, similar to the one provided in the sample waiting room stack, but designed for an OpenID authentication flow. The process of becoming authenticated involves getting a position in the waiting room queue and waiting until the serving position is equal or greater than the client's queue position. The OpenID waiting room page redirects back to the target site, which uses the OpenID API to complete the token acquisition and session configuration for the client. This solution's API endpoints map directly to the official OpenID Connect 1.0 flow specification, name-for-name. Refer to [OpenID Connect Core 1.0 Authentication](#) for details.

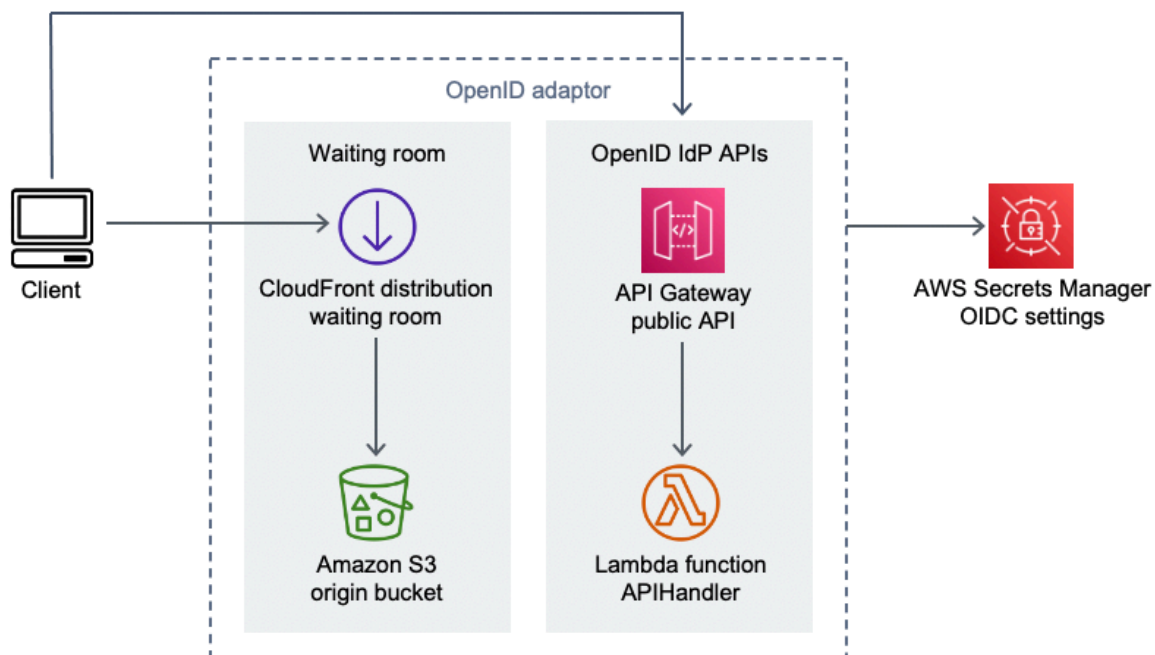


Figure 3: Virtual Waiting Room on AWS OpenID adaptor component

1. CloudFront distribution serves the S3 bucket's content to the user.
2. S3 bucket hosts sample waiting room pages.
3. Amazon API Gateway API provides a set of OIDC-compatible APIs that can be used with existing web hosting software that support OIDC identity provider's Lambda authorize function.
4. The APIHandler Lambda function handles requests for all API Gateway resource paths. Different Python functions within the same module are mapped to each API path. For example, the `/authorize` resource path in API Gateway invokes `authorize()` within the Lambda Function.
5. OIDC settings are stored in Secrets Manager.

Sample inlet strategies

Inlet strategies determine when the solution's serving counter should move forward to accommodate more users in the target site. For more conceptual information about waiting room inlet strategies, refer to [Design considerations \(p. 16\)](#).

There are two sample inlet strategies provided by the solution: *MaxSize* and *Periodic*.



Figure 4: Virtual Waiting Room on AWS Inlet strategies component

Max Size inlet strategy option:

1. A client issues an Amazon SNS notification that invokes the MaxSizeInlet Lambda function to increase the serving counter based on the message payload.
2. The MaxSizeInlet Lambda function expects to receive a message that it uses it determine how much to increment the serving counter.

Periodic inlet strategy option:

3. A CloudWatch rule invokes a Lambda function every minute to increase the serving counter by a fixed quantity.
4. The PeriodicInlet Lambda function increments the serving counter by the given size if the time is between the start and end time provided. Optionally, it checks a CloudWatch alarm and, if the alarm is in OK state, performs the increment, otherwise skips it.

Sample waiting room

The sample waiting room integrates with the public and private APIs in addition to the custom authorizer to demonstrate a minimal end-to-end waiting room solution. The main web page is stored in an S3 bucket and used as an origin to CloudFront. It takes the user through the following steps:

1. Get in line at the waiting room for entry into the site.
2. Obtain the client's position in line.
3. Obtain the serving position of the waiting room.
4. Obtain a token set once the serving position is equal or greater to the client's position.
5. Use the token to call an API protected by the Lambda authorizer.

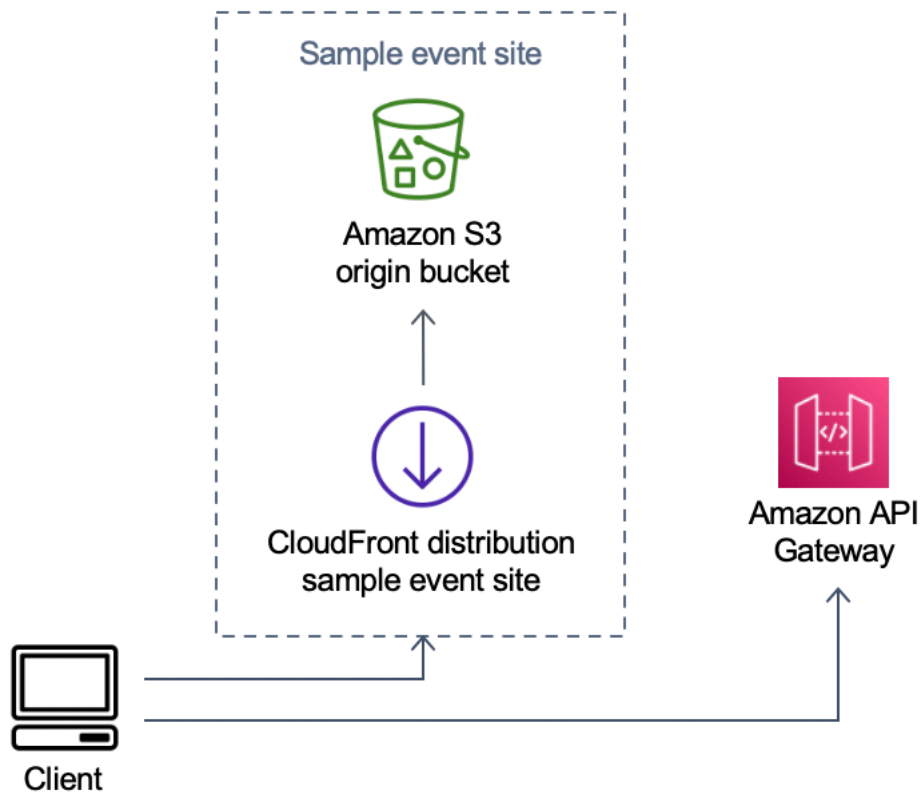


Figure 5: Virtual Waiting Room on AWS Sample event site component

1. S3 bucket hosts the sample content for the waiting room and control panel.
2. CloudFront distribution serves the S3 bucket content to user.
3. Sample API Gateway deployment with shopping-like resource paths like `/search` and `/checkout`. This API is installed by the stack and configured with the token authorizer. It is intended as an example of a simple way to protect an API with the waiting room. Requests that present a valid token are forwarded to the Lambda, otherwise an error is returned. There is no functionality to the API other than the response from the Lambda function attached.

Security

When you build systems on AWS infrastructure, security responsibilities are shared between you and AWS. This [shared model](#) reduces your operational burden because AWS operates, manages, and controls the components including the host operating system, the virtualization layer, and the physical security of the facilities in which the services operate. For more information about AWS security, visit [AWS Cloud Security](#).

ElastiCache for Redis is assigned a network interface inside the private VPC. The Lambda functions that interact with ElastiCache for Redis are also assigned network interfaces within a VPC. All other resources have network connectivity in the shared AWS network space. Lambda functions with VPC interfaces that interact with other AWS services use VPC endpoints to connect to these services.

The public and private keys used for creating and validating JSON web tokens are generated at deployment time and stored in Secrets Manager. The password used to connect to ElastiCache for Redis is also generated at deployment time and stored in Secrets Manager. The private key and ElastiCache for Redis password are not accessible via any solution API.

The public API must be accessed through CloudFront. The solution generates an API key for API Gateway, which is used as the value of a custom header, `x-api-key`, in CloudFront. CloudFront includes this header when making origin requests. For additional details, refer to [Adding custom headers to origin requests](#) in the *Amazon CloudFront Developer Guide*.

The private APIs are configured to require AWS IAM authorization for invocation. The solution creates the ProtectedAPIGroup IAM user group with the appropriate permissions to invoke the private APIs. An IAM user added to this group are authorized to invoke the private APIs.

IAM policies used in roles and permissions that are attached to various resources created by the solution grant only the permissions required to perform the necessary tasks.

For resources such as S3 buckets, SQS queues, and SNS topics generated by the solution, encryption at rest and during transit is activated wherever possible.

Monitoring

The core API stack includes several CloudWatch alarms that can be monitored to detect problems while the solution is operating. The stack creates an alarm for Lambda function errors and throttle conditions, and changes the state of the alarm from OK to ALARM if an error or throttle condition occurs in a one-minute period.

The stack also creates alarms for each API Gateway deployment for 4XX and 5XX status codes. The alarm changes state from OK to ALARM if a 4XX or 5XX status code is returned from the API within a one-minute period.

These alarms return to an OK state after one minute of no errors or throttles.

IAM roles

AWS Identity and Access Management (IAM) roles allow customers to assign granular access policies and permissions to services and users on the AWS Cloud. This solution creates IAM roles that grant the solution's AWS Lambda functions access to create Regional resources.

Amazon CloudFront

The `virtual-waiting-room-on-aws.template` CloudFormation template, which creates the core public and private APIs of the waiting room, also deploys a CloudFront distribution for the public API. CloudFront caches the responses from the public API, thereby reducing load on API Gateway and the Lambda functions performing work.

This solution also has an optional sample waiting room template that deploys a simple web application [hosted](#) in an Amazon Simple Storage Service (Amazon S3) bucket. To help reduce latency and improve security, an Amazon CloudFront distribution is deployed with an origin access identity, which is a CloudFront user that provides public access to the solution's website bucket contents. For more information, refer to [Restricting Access to Amazon S3 Content by Using an Origin Access Identity](#) in the *Amazon CloudFront Developer Guide*.

Security groups

The [VPC security groups](#) created in this solution are designed to control and isolate network traffic to the ElastiCache for Redis. Lambdas that need to communicate with the ElastiCache for Redis are placed in the same Security Group as that of the ElastiCache for Redis. We recommend that you review the security groups and further restrict access as needed once the deployment is up and running.

Design considerations

Deployment options

If this is the first time installing, or you are not sure what to install, deploy the `virtual-waiting-room-on-aws-getting-started.template` nested CloudFormation template, which installs the core, authorizers, and sample waiting room templates. This provides you a minimal waiting room with a simple flow.

Supported protocols

The Virtual Waiting Room on AWS solution can be integrated with the following:

- JSON Web Token verification libraries and tools
- Existing API Gateway deployments
- REST API clients
- OpenID clients and providers

Waiting room inlet strategies

Inlet strategies encapsulate the logic and data needed to move clients from the waiting room to the website. An inlet strategy can be implemented as a Lambda function, container, Amazon EC2 instance, or any other compute resource. It does not need to be a cloud resource as long as it can call the waiting room public and private APIs. The inlet strategy receives events about the waiting room, the website, or other outside indicators that help it decide when more clients can have tokens issued and enter the site. There are several approaches to inlet strategies. Which one you adopt depends on the resources available to you and constraints in the design of the website being protected.

The primary action taken by the inlet strategy is to call the `increment_serving_num` Amazon API Gateway private API with a relative value that indicates how many more clients can enter the site. This section describes two sample inlet strategies. These can be used as-is, customized, or you can employ a completely different approach.

MaxSize

Using the MaxSize strategy, the `MaxSizeInlet` Lambda function is configured with the maximum number of clients that can use the website concurrently. This is a fixed value. A client issues an Amazon SNS notification that invokes the `MaxSizeInlet` Lambda function to increase the serving counter based on the message payload. The source of the SNS message can come from anywhere, including code on the website or a monitoring tool that observes the site's utilization level.

The `MaxSizeInlet` Lambda function expects to receive a message that can include:

- `exited` : number of transactions that have completed
- list of request IDs to be marked completed

- list of request IDs to be marked abandoned

This data is used to determine how much to increment the serving counter. There can be cases where there is no additional capacity to increment the counter, based on the current number of clients.

Periodic

When using the periodic strategy, a CloudWatch rule invokes the `PeriodicInlet` Lambda function every minute to increase the serving counter by a fixed quantity. The periodic inlet is parameterized with the event start time, end time, and increment amount. Optionally, this strategy also checks a CloudWatch alarm and, if the alarm is in OK state, it performs the increment, otherwise it skips it. The site integrators can connect a utilization metric to an alarm, and use that alarm to pause the periodic inlet. This strategy only changes the serving position while the current time is between the start and end times, and optionally, the specified alarm is in the OK state.

Customizing and extending the solution

Your organization's site administrator must decide on the integration methods to use with the waiting room. There are two options:

1. Basic integration directly using APIs and API Gateway authorizers.
2. OpenID integration through an identity provider.

In addition to the integration above, you might be required to configure the domain name redirection. You are also responsible for deploying a customized waiting room site page.

The Virtual Waiting Room on AWS solution is designed for extension through two mechanisms: EventBridge for unidirectional event notification and REST APIs for bidirectional communication.

Quotas

The primary scale limitation for Virtual Waiting Room on AWS is the Lambda throttle limit for the installed AWS Region. When installed into an AWS account with the default Lambda concurrent run quota, the Virtual Waiting Room on AWS solution can handle up to 500 clients per second requesting a position in the queue. The 500 client per second rate is based on the solution having all Lambda function concurrent quota limits available exclusively. If the Region in the account is shared with other solutions that invoke Lambda functions, the Virtual Waiting Room on AWS solution should have at least 1,000 concurrent invocations available. You can use CloudWatch metrics to chart the Lambda concurrent invocations in your account over time to make a determination. You can use the [Service Quotas console](#) to request increases. Increasing the Lambda throttle limit only raises the monthly account charges if additional invocations actually occur.

For each additional 500 clients per second, increase your throttle limit by 1,000.

Incoming users per second expected	Recommended concurrent execution quota
0-500	1,000 (default)
501-1,000	2,000
1,001-1,500	3,000

Lambda has a fixed burst limit of 3,000 concurrent invocations. For more information, refer to [Lambda function scaling](#). Client code should expect and retry some API calls if an error code is returned indicating a temporary throttle situation. The sample waiting room client includes this code as an example of how to design clients used in high capacity and high burst events.

This solution is also compatible with Lambda *reserved* and *provisioned* concurrency with custom configuration steps. For details, refer to [Managing Lambda reserved concurrency](#).

The upper limit of users that can enter the waiting room, receive a token, and continue to a transaction is limited by the upper limit of ElastiCache for Redis counters. The counters are used for the waiting room serving position and tracking summary state of the solution. The counters used in ElastiCache for Redis have an upper limit of 9,223,372,036,854,775,807. A DynamoDB table is used to store a copy of each token issued to a waiting room user. DynamoDB has no practical limit on a table's size.

Regional deployments

Services used by this solution are supported in all AWS Regions. For the most current availability of AWS services by Region, refer to the [AWS Regional Services List](#).

AWS CloudFormation templates

To automate deployment, this solution uses the following AWS CloudFormation templates, which you can download before deployment.

If this is the first time installing, or you are not sure what to install, deploy the `virtual-waiting-room-on-aws-getting-started.template` AWS CloudFormation template, which installs the core, authorizers, and sample waiting room code templates. This allows you to test a working waiting room with a simple flow.



View template

`virtual-waiting-room-on-aws-api-gateway-cw-logs-role.template`: Use this template to add a default role ARN to API Gateway at the account-level for CloudWatch logging permissions. Refer to [Prerequisites \(p. 21\)](#) for details on whether your account requires the deployment of this template or not.



View template

`virtual-waiting-room-on-aws-getting-started.template`: Use this nested template to install the core, authorizers, and sample waiting room stacks.



View template

`virtual-waiting-room-on-aws.template`: Use this core template to install the core public and private REST APIs and cloud services for creating waiting room events. Install this template in the account and Region where you need the waiting room REST APIs, ElastiCache for Redis, and DynamoDB table.



View template

`virtual-waiting-room-on-aws-authorizers.template`: Use this template to install the Lambda authorizer designed for verifying waiting room-issued tokens and intended for protecting end-users' APIs. Requires the core stack. Some outputs from the core stack are needed as parameters to deploy this stack. This is an optional template.



View template

`virtual-waiting-room-on-aws-openid.template`: Use this template to install an OpenID identity provider for waiting room integration with authorization

interfaces. Requires the core stack. Some outputs from the core stack are needed to deploy this stack. This is an optional template.

[View template](#)

virtual-waiting-room-on-aws-sample-inlet-strategy.template: Use this template to install sample inlet strategies intended for use between a target site and the waiting room. Inlet strategies help encapsulate logic to determine when to allow more users into the target site. Requires the core stack. Outputs from the core stack are needed to deploy this stack. This is an optional template.

[View template](#)

virtual-waiting-room-on-aws-sample.template: Use this template to install a sample minimal web and API Gateway configuration for a waiting room and target site. Requires the core and authorizers stacks. Outputs from the core and authorizers stacks are required as parameters to deploy this stack. This is an optional template.

Automated deployment

Before you launch the solution, review the cost, architecture, network security, and other considerations discussed in this guide. Follow the step-by-step instructions in this section to configure and deploy the solution into your account.

Time to deploy: Approximately 30 minutes (getting-started stack only)

Prerequisites

- AWS account console permissions equivalent to [Administrator Access](#).
- Activate CloudWatch logging from API Gateway:
 - Sign in to the [API Gateway console](#) and select the Region where you plan to install the stacks.

If you have existing APIs defined in this Region:

1. Select any API.
 2. From the left navigation, select **Settings**.
 3. Check for a value in the **CloudWatch log role ARN** field.
- If there is no ARN, install the [virtual-waiting-room-on-aws-api-gateway-cw-logs-role.template](#).
 - If there is an ARN, start with [launching the getting-started stack \(p. 21\)](#).

If there are no existing APIs defined in this Region, install the [virtual-waiting-room-on-aws-api-gateway-cw-logs-role.template](#).

- Knowledge of the architecture and implementation details of the target site to protect.

Deployment overview

Use the following steps to deploy this solution on AWS. For detailed instructions, follow the links for each step.

[Step 1. Launch the getting-started stack \(p. 21\)](#)

- Launch the AWS CloudFormation template into your AWS account.
- Review the templates parameters and enter or adjust the default values as needed.

[Step 2. \(Optional\) Test the waiting room \(p. 23\)](#)

- Generate AWS keys to call the IAM secured APIs.
- Open up the sample waiting room's control panel.
- Test the sample waiting room.

Step 1. Launch the getting-started stack

This automated AWS CloudFormation template deploys the core, authorizers, and sample waiting room templates which allows you to view and test a working waiting room. You must read and understand the Prerequisites before launching the stack.

Note

You are responsible for the cost of the AWS services used while running this solution. For more details, visit the [Cost \(p. 2\)](#) section in this guide, and refer to the pricing webpage for each AWS service used in this solution.

1. Sign in to the AWS Management Console and select the button to launch the `virtual-waiting-room-on-aws-getting-started.template` AWS CloudFormation template.



Alternatively, you can [download the template](#) as a starting point for your own implementation.

2. The template launches in the US East (N. Virginia) Region by default. To launch the solution in a different AWS Region, use the Region selector in the console navigation bar.
3. On the **Create stack** page, verify that the correct template URL is in the **Amazon S3 URL** text box and choose **Next**.
4. On the **Specify stack details** page, assign a name to your solution stack. For information about naming character limitations, refer to [IAM and STS Limits](#) in the *AWS Identity and Access Management User Guide*.
5. Under **Parameters**, review the parameters for this solution template and modify them as necessary. This solution uses the following default values.

Parameter	Default	Description
Event ID	Sample	Unique ID for this instance of the waiting room, GUID format suggested.
Validity Period	3600	Token validity period in seconds.
Enable Events Generation	false	If set to true, metrics related to the Waiting Room are written to its event bus every minute
Redis Port	1785	The port number to use for connecting to ElastiCache for Redis server. Recommended not to use the default ElastiCache for Redis port of 6379.
EnableQueuePositionExpiry	true	If set to false, queue position expiry period is not applied.
QueuePositionExpiryPeriod	900	It is the time interval in seconds beyond which a queue position is ineligible to generate a token.
IncrSvcOnQueuePositionExpiry	false	If set to true, the serving counter is automatically advanced based on expired

Parameter	Default	Description
		queue positions that did not successfully generate tokens.

6. Choose **Next**.
7. On the **Configure stack options** page, choose **Next**.
8. On the **Review** page, review and confirm the settings. Check the box acknowledging that the template creates AWS Identity and Access Management (IAM) resources.
9. Choose **Create stack** to deploy the stack.

You can view the status of the stack in the AWS CloudFormation Console in the **Status** column. You should receive a CREATE_COMPLETE status in approximately 30 minutes.

Step 2. (Optional) Test the waiting room

If you deployed the getting-started stack, the following steps help you test the waiting room's functionality. To complete testing, you need AWS keys with permissions to call the IAM secured APIs in the core stack.

Generate AWS keys to call the IAM secured APIs

1. [Create](#) or use an IAM user in the AWS account where the `aws-virtual-waiting-room-getting-started` template CloudFormation template was deployed.
2. Grant the [IAM user programmatic access](#). When creating a new set of access keys for the IAM user, download the key file when presented. You need the IAM user's Access Key ID and Secret Access Key to test the waiting room.
3. [Add the IAM user to the ProtectedAPIGroup IAM user group](#) created by the template.

Open up the sample waiting room's control panel

1. Sign in to the [AWS CloudFormation console](#) and select the solution's getting-started stack.
2. Choose the **Outputs** tab.
3. Under the **Key** column, locate **ControlPanelURL**, and select the corresponding value.
4. Open the control panel in a new tab or browser window.
5. In the control panel, expand the **Configuration** section.
6. Enter the Access key ID and Secret Access Key you retrieved in [Generate AWS keys to call the IAM secured APIs \(p. 23\)](#). The endpoints and event ID are filled in from the URL parameters.
7. Choose **Use**. The button activates after you have supplied the credentials.

Test the sample waiting room

1. In the [AWS CloudFormation console](#), select the solution's getting-started stack.
2. Choose the **Outputs** tab.
3. Under the **Key** column, locate **WaitingRoomURL**, and select the corresponding value.
4. Open the waiting room, then choose **Reserve** to enter the waiting room.
5. Navigate back to the browser tab that has the control panel.

6. Under **Increment Serving Counter**, choose **Change**. This allows 100 users to move on from the waiting room to the target site.
7. Navigate back to the waiting room and choose **Check out now!** You will now get redirected to the target site.
8. Choose **Purchase now** to finish your transaction at the target site.

Deploying separate stacks

The core stack is the only required stack to get the main functionality of the waiting room. All other stacks are optional. Launch the authorizers stack if you don't already have a way to validate waiting room-issued tokens or protect any APIs you might already have. Launch the OpenID stack if you require an OpenID identity provider for waiting room integration with authorization interfaces. The sample inlet strategy stack provides a couple of examples on how and when to allow more users into the site you are trying to protect.

1. Launch the core stack

Time to deploy: Approximately 20 minutes

This automated AWS CloudFormation template deploys Virtual Waiting Room on AWS in the AWS Cloud. You must complete the [Prerequisites \(p. 21\)](#) before launching the stack.

Note

You are responsible for the cost of the AWS services used while running this solution. For more details, visit the [Cost \(p. 2\)](#) section in this guide, and refer to the pricing webpage for each AWS service used in this solution.

1. Sign in to the AWS Management Console and select the button to launch the `aws-virtual-waiting-room-on-aws.template` AWS CloudFormation template.



Alternatively, you can [download the template](#) as a starting point for your own implementation.

2. The template launches in the US East (N. Virginia) Region by default. To launch the solution in a different AWS Region, use the Region selector in the console navigation bar.
3. On the **Create stack** page, verify that the correct template URL is in the **Amazon S3 URL** text box and choose **Next**.
4. On the **Specify stack details** page, assign a name to your solution stack. For information about naming character limitations, refer to [IAM and STS Limits](#) in the *AWS Identity and Access Management User Guide*.
5. Under **Parameters**, review the parameters for this solution template and modify them as necessary. This solution uses the following default values.

Parameter	Default	Description
Event ID	Sample	Unique ID for this instance of the Waiting Room, GUID format suggested.
Validity Period	3600	Token validity period in seconds.

Parameter	Default	Description
Enable Events Generation	false	If set to true, metrics related to the waiting room are written to its event bus every minute.
Redis Port	1785	The port number to use for connecting to ElastiCache for Redis server. Recommended not to use the default ElastiCache for Redis port of 6379.
EnableQueuePositionExpiry	true	If set to false, queue position expiry period is not applied.
QueuePositionExpiryPeriod	900	It is the time interval in seconds beyond which a queue position is ineligible to generate a token.
IncrSvcOnQueuePositionExpiry	false	If set to true, the serving counter is automatically advanced based on expired queue positions that did not successfully generate tokens.

6. Choose **Next**.
7. On the **Configure stack options** page, choose **Next**.
8. On the **Review** page, review and confirm the settings. Check the box acknowledging that the template creates AWS Identity and Access Management (IAM) resources.
9. Choose **Create stack** to deploy the stack.

You can view the status of the stack in the AWS CloudFormation Console in the **Status** column. You should receive a CREATE_COMPLETE status in approximately 20 minutes.

2. (Optional) Launch the Authorizers stack

Time to deploy: Approximately five minutes

1. Sign in to the AWS Management Console and select the button to launch the `aws-virtual-waiting-room-on-aws-authorizers.template` AWS CloudFormation template.



Alternatively, you can [download the template](#) as a starting point for your own implementation.

2. The template launches in the US East (N. Virginia) Region by default. To launch the solution in a different AWS Region, use the Region selector in the console navigation bar.
3. On the **Create stack** page, verify that the correct template URL is in the **Amazon S3 URL** text box and choose **Next**.
4. On the **Specify stack details** page, assign a name to your solution stack. For information about naming character limitations, refer to [IAM and STS Limits](#) in the *AWS Identity and Access Management User Guide*.

- Under **Parameters**, review the parameters for this solution template and modify them as necessary. This solution uses the following default values.

Parameter	Default	Description
Public API Endpoint	<Requires input>	Public endpoint for the virtual waiting room APIs.
Waiting Room Event ID	Sample	Event ID of waiting room.
Issuer URI	<Requires input>	Issuer URI of the public keys and tokens.

- Choose **Next**.
- On the **Configure stack options** page, choose **Next**.
- On the **Review** page, review and confirm the settings. Check the box acknowledging that the template creates AWS Identity and Access Management (IAM) resources.
- Choose **Create stack** to deploy the stack.

You can view the status of the stack in the AWS CloudFormation Console in the **Status** column. You should receive a CREATE_COMPLETE status in approximately five minutes.

3. (Optional) Launch the OpenID stack

Time to deploy: Approximately five minutes

- Sign in to the AWS Management Console and select the button to launch the `aws-virtual-waiting-room-on-aws-openid` template AWS CloudFormation template.



Alternatively, you can [download the template](#) as a starting point for your own implementation.

- The template launches in the US East (N. Virginia) Region by default. To launch the solution in a different AWS Region, use the Region selector in the console navigation bar.
- On the **Create stack** page, verify that the correct template URL is in the **Amazon S3 URL** text box and choose **Next**.
- On the **Specify stack details** page, assign a name to your solution stack. For information about naming character limitations, refer to [IAM and STS Limits](#) in the *AWS Identity and Access Management User Guide*.
- Under **Parameters**, review the parameters for this solution template and modify them as necessary. This solution uses the following default values.

Parameter	Default	Description
Public API Endpoint	<Requires input>	Public endpoint URL for the virtual waiting room APIs.
Private API Endpoint	<Requires input>	Private endpoint URL for the virtual waiting room APIs.

Parameter	Default	Description
API Region	<Requires input>	AWS region name for the public and private waiting room APIs.
Event ID	Sample	Event ID of waiting room.

6. Choose **Next**.
7. On the **Configure stack options** page, choose **Next**.
8. On the **Review** page, review and confirm the settings. Check the box acknowledging that the template creates AWS Identity and Access Management (IAM) resources.
9. Choose **Create stack** to deploy the stack.

You can view the status of the stack in the AWS CloudFormation Console in the **Status** column. You should receive a CREATE_COMPLETE status in approximately five minutes.

4. (Optional) Launch the sample inlet strategy stack

Time to deploy: Approximately two minutes

1. Sign in to the AWS Management Console and select the button to launch the `aws-virtual-waiting-room-sample-inlet-strategy` template AWS CloudFormation template.



Alternatively, you can [download the template](#) as a starting point for your own implementation.

2. The template launches in the US East (N. Virginia) Region by default. To launch the solution in a different AWS Region, use the Region selector in the console navigation bar.
3. On the **Create stack** page, verify that the correct template URL is in the **Amazon S3 URL** text box and choose **Next**.
4. On the **Specify stack details** page, assign a name to your solution stack. For information about naming character limitations, refer to [IAM and STS Limits](#) in the *AWS Identity and Access Management User Guide*.
5. Under **Parameters**, review the parameters for this solution template and modify them as necessary. This solution uses the following default values.

Parameter	Default	Description
Event ID	Sample	Event ID of waiting room.
Private Core API Endpoint	<Requires input>	Private endpoint URL for the virtual waiting room APIs.
Core API Region	<Requires input>	AWS Region where the core API is installed.
Inlet Strategy	Periodic	Inlet strategy to be deployed. Periodic increments the

Parameter	Default	Description
		serving number every minute. MaxSize increments serving number based on the maximum number of transactions that the downstream target site can handle at a given time.
Increment By	<Requires input>	How much the serving counter should be incremented every minute. Required if selecting periodic inlet strategy.
Start Time	<Requires input>	Timestamp on when to start incrementing the serving number (epoch time in seconds). Required if selecting periodic inlet strategy.
End Time	<Requires input>	Timestamp on when to stop incrementing the serving number (epoch time in seconds). If left 0, serving number is incremented indefinitely. Required if selecting periodic inlet strategy.
CloudWatch Alarm Name	<Requires input>	Optional CloudWatch alarm name to be associated with periodic inlet strategy. If provided and in alarming state, serving number is not incremented. Applicable to periodic inlet strategy only.
Max Size	<Requires input>	The maximum number of transactions that the downstream target site can process at a time (MaxSize Strategy).

6. Choose **Next**.
7. On the **Configure stack options** page, choose **Next**.
8. On the **Review** page, review and confirm the settings. Check the box acknowledging that the template creates AWS Identity and Access Management (IAM) resources.
9. Choose **Create stack** to deploy the stack.

You can view the status of the stack in the AWS CloudFormation Console in the **Status** column. You should receive a CREATE_COMPLETE status in approximately two minutes.

5. (Optional) Launch the sample waiting room stack

Time to deploy: Approximately five minutes

1. Sign in to the AWS Management Console and select the button to launch the `aws-virtual-waiting-room-sample.template` AWS CloudFormation template.



Alternatively, you can [download the template](#) as a starting point for your own implementation.

2. The template launches in the US East (N. Virginia) Region by default. To launch the solution in a different AWS Region, use the Region selector in the console navigation bar.
3. On the **Create stack** page, verify that the correct template URL is in the **Amazon S3 URL** text box and choose **Next**.
4. On the **Specify stack details** page, assign a name to your solution stack. For information about naming character limitations, refer to [IAM and STS Limits](#) in the *AWS Identity and Access Management User Guide*.
5. Under **Parameters**, review the parameters for this solution template and modify them as necessary. This solution uses the following default values.

Parameter	Default	Description
API Gateway Region	<i><Requires input></i>	AWS Region name of the API Gateway.
Authorizer ARN	<i><Requires input></i>	ARN of the API Gateway Lambda authorizer.
Event ID	Sample	Event ID of the waiting room.
Private API Endpoint	<i><Requires input></i>	Private endpoint URL for the virtual waiting room APIs.
Public API Endpoint	<i><Requires input></i>	Public endpoint URL for the virtual waiting room APIs.

6. Choose **Next**.
7. On the **Configure stack options** page, choose **Next**.
8. On the **Review** page, review and confirm the settings. Check the box acknowledging that the template creates AWS Identity and Access Management (IAM) resources.
9. Choose **Create stack** to deploy the stack.

You can view the status of the stack in the AWS CloudFormation Console in the **Status** column. You should receive a `CREATE_COMPLETE` status in approximately five minutes.

Updating the stack from a previous version

If you are updating the stack from a previous version of the Virtual Waiting Room on AWS solution to version 1.1, then a full reset is required once the update is completed. This ensures all the functionalities of the newer version become available.

Performance data

Virtual Waiting Room on AWS has been load tested with a tool called [Locust](#). The simulated event sizes ranged from 10,000 to 100,000 clients. The load testing environment consisted of the following configuration:

- Locust 2.x with customizations for AWS Cloud deployments
- Four AWS Regions (us-west-1, us-west-2, us-east-1, us-east-2)
- 10 c5.4xlarge Amazon EC2 hosts per Region (40 total)
- 32 Locust processes per host
- Simulated users were evenly spread among the 1,280 processes

The end-to-end API test steps for each user process:

1. Call `assign_queue_num` and receive a request ID.
2. Loop `queue_num` with the request ID until it returns the user's queue position (short time).
3. Loop `serving_num` until the returned value is \geq user's queue position (long time).
4. Infrequently call `waiting_room_size` to retrieve the number of waiting users.
5. Call `generate_token` and receive a JWT for use in the target site.

Findings

There is no practical upper limit to the number of clients that can be processed through the waiting room.

The rate at which users enter the waiting room impacts Lambda function concurrent run quotas for the Region in which it is deployed.

The load test was not able to exceed the default API Gateway request limits of 10,000 requests per second with the caching policies used with CloudFront.

The `get_queue_num` Lambda function has a near 1:1 invocation rate to the rate of incoming users to the waiting room. This Lambda function may be throttled during high rates of incoming users due to concurrency limits or burst limits. Throttling caused by a large number of `get_queue_num` Lambda function invocations can impact other Lambda functions as a side-effect. The overall system continues operating if the client software can respond appropriately to this type of temporary scaling error with retry/back-off logic.

The CloudFront distribution configured by the core stack in a default quota configuration can handle a waiting room holding 250,000 users with each user polling the `serving_num` API at least every second.

Troubleshooting

4xx response status from APIs

- This can be caused from an incorrect Event ID or Request ID or both. This occurs in the CloudWatch Logs for the related Lambda function.
- Private APIs are IAM authenticated and the client needs AWS keys that have rights to invoke the private APIs. This occurs in the CloudWatch Logs for API Gateway.

5xx response status from APIs

- Response from throttled Lambda or API Gateway, check `<LambdaFunctionName>ThrottlesAlarm` CloudWatch alarm.
- Misconfiguration on back-end, check `<LambdaFunctionName>ErrorsAlarm` CloudWatch alarm and CloudWatch Logs for details.

5XXErrorPublic/PrivateApiAlarm

- This alarm state is ALARM when the API returns a 5XX status to the caller within a 60-second period.
- This alarm returns to OK when no 5xx status is returned for 60 seconds.
- This alarm can be started by a Lambda function or Lambda runtime returning an error to API Gateway.

4XXErrorPublic/PrivateApiAlarm

- This alarm state is ALARM when the API return a 4XX status to the caller within a 60 second period.
- This alarm returns to OK when to 4XX status is returned for 60 seconds.
- This alarm can be initiated by an incorrect API URL.

`<LambdaFunctionName>ThrottlesAlarm`

- This alarm state is ALARM when the named Lambda encounters a concurrent run limit within a 60-second period.
- This alarm returns to OK if no throttles are encountered for 60 seconds.
- You may need to increase the concurrency limit for your account's Region.
- You may be encountering the burst limit for Lambda, which requires some retry logic on your client.

`<LambdaFunctionName>ErrorsAlarm`

- This alarm state is ALARM when the named Lambda encounters a runtime run error within a 60-second period.
- This alarm returns to OK if no errors are encountered for 60 seconds.
- This can be caused by a misconfiguration on the backend.
- This can be caused by a bug in the Lambda's code.

Additional resources

AWS services	
• AWS CloudFormation	• Amazon DynamoDB
• Amazon Simple Storage Service	• Amazon API Gateway
• AWS Lambda	• AWS Secrets Manager
• Amazon CloudFront	• Amazon Simple Queue Service
• Amazon EventBridge	• Amazon CloudWatch
• Amazon ElastiCache for Redis	• Amazon Comprehend
• Amazon Virtual Private Cloud	• AWS Identity and Access Management

Uninstall the solution

You can uninstall the Virtual Waiting Room on AWS solution from the AWS Management Console or by using the AWS Command Line Interface. You must manually delete the S3 buckets used to store logs by various resources created by this solution. AWS Solutions Implementations do not automatically delete these S3 buckets so you still have the ability to review logs event after the solution has been deleted.

If you have manually added an IAM user to the ProtectedAPIGroup IAM user group created by the solution, [remove the IAM user from the IAM user group](#) before uninstalling the solution. Otherwise, the IAM user group and the associated IAM policy fails to delete.

For each of the stacks deployed, follow the instructions below.

Using the AWS Management Console

1. Sign in to the [AWS CloudFormation console](#).
2. On the **Stacks** page, select this solution's installation stack.
3. Choose **Delete**.

Using AWS Command Line Interface

Determine whether the AWS Command Line Interface (AWS CLI) is available in your environment. For installation instructions, refer to [What Is the AWS Command Line Interface?](#) in the *AWS CLI User Guide*. After confirming that the AWS CLI is available, run the following command.

```
$ aws cloudformation delete-stack --stack-name <installation-stack-name>
```

Deleting the Amazon S3 buckets

This solution is configured to retain the solution-created Amazon S3 bucket (for deploying in an opt-in Region) if you decide to delete the AWS CloudFormation stack to prevent accidental data loss. After uninstalling the solution, you can manually delete this S3 bucket if you do not need to retain the data. Follow these steps to delete the Amazon S3 bucket.

1. Sign in to the [Amazon S3 console](#).
2. Choose **Buckets** from the left navigation pane.
3. Locate the <stack-name> S3 buckets.
4. Select the S3 bucket and choose **Delete**.

To delete the S3 bucket using AWS CLI, run the following command:

```
$ aws s3 rb s3://<bucket-name> --force
```

Source code

Visit our [GitHub repository](#) to download the source files for this solution and to share your customizations with others.

Revisions

Date	Change
November 2021	Initial release
September 2022	Version 1.1: Automatic serving counter increment based on expired queue positions. Relocate some Redis usage to DynamoDB. Public API endpoint to get remaining queue position expiry time. For more information, refer to the CHANGELOG.md file in the GitHub repository
April 2023	Version 1.1.1: Mitigated impact caused by new default settings for S3 Object Ownership (ACLs disabled) for all new S3 buckets. For more information, refer to the CHANGELOG.md file in the GitHub repository

Contributors

- Jim Thario
- Thyag Ramachandran
- Joan Morgan
- Justin Pirtle
- Allen Moheimani

Notices

Customers are responsible for making their own independent assessment of the information in this document. This document: (a) is for informational purposes only, (b) represents AWS current product offerings and practices, which are subject to change without notice, and (c) does not create any commitments or assurances from AWS and its affiliates, suppliers or licensors. AWS products or services are provided “as is” without warranties, representations, or conditions of any kind, whether express or implied. AWS responsibilities and liabilities to its customers are controlled by AWS agreements, and this document is not part of, nor does it modify, any agreement between AWS and its customers.

Virtual Waiting Room on AWS is licensed under the terms of the of the Apache License Version 2.0 available at [The Apache Software Foundation](https://www.apache.org/licenses/LICENSE-2.0).

AWS glossary

For the latest AWS terminology, see the [AWS glossary](#) in the *AWS Glossary Reference*.