



Diseño para la alta disponibilidad y la resiliencia en las aplicaciones de Amazon EKS

# AWS Guía prescriptiva



# AWS Guía prescriptiva: Diseño para la alta disponibilidad y la resiliencia en las aplicaciones de Amazon EKS

Copyright © 2026 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

Las marcas comerciales y la imagen comercial de Amazon no se pueden utilizar en relación con ningún producto o servicio que no sea de Amazon, de ninguna manera que pueda causar confusión entre los clientes y que menosprecie o desacredite a Amazon. Todas las demás marcas registradas que no son propiedad de Amazon son propiedad de sus respectivos propietarios, que pueden o no estar afiliados, conectados o patrocinados por Amazon.

# Table of Contents

Introducción .....	1
Diseño de alta disponibilidad y resiliencia .....	2
Distribuya las cargas de trabajo .....	3
Utilice las restricciones de dispersión de la topología de los módulos .....	3
Afinidad y antiafinidad entre los pods .....	7
Presupuesto para la disrupción de cápsulas .....	9
Sondas y controles de estado .....	10
Sonda de inicio .....	10
Sonda de vitalidad .....	11
Sonda de preparación .....	11
Comprobaciones de estado del balanceador de carga y recursos de Ingress .....	11
Ganchos para el ciclo de vida .....	12
Comprenda el desalojo de los pods durante las interrupciones zonales .....	14
Implementación del cambio zonal de Amazon EKS para mejorar la resiliencia .....	15
Comprender el mecanismo de cambio zonal .....	15
Métodos de activación por turnos zonales .....	16
Requisitos previos para un cambio zonal efectivo .....	17
Recomendaciones para la resiliencia a las disrupciones zonales .....	17
Finalización y recuperación de los turnos .....	18
Conclusión .....	19
Recursos .....	20
Historial de documentos .....	21
Glosario .....	22
# .....	22
A .....	23
B .....	26
C .....	28
D .....	31
E .....	36
F .....	38
G .....	40
H .....	41
I .....	43
L .....	45

M .....	46
O .....	51
P .....	54
Q .....	57
R .....	57
S .....	60
T .....	64
U .....	66
V .....	66
W .....	67
Z .....	68
	lxix

# Diseño para una alta disponibilidad y resiliencia en las aplicaciones de Amazon EKS

Haofei Feng, Frank Fan y Rus Kalakutskiy, Amazon Web Services ()AWS

Octubre de 2025 ([historia del documento](#))

Garantizar la alta disponibilidad (HA) y la resiliencia en el diseño de las aplicaciones es fundamental para lograr un objetivo de punto de recuperación (RPO) y un objetivo de tiempo de recuperación (RTO) próximos a cero. A medida que las organizaciones migran y modernizan cada vez más sus aplicaciones a los entornos de Kubernetes, la demanda de soluciones sólidas y escalables sigue aumentando. Amazon Elastic Kubernetes Service (Amazon EKS) le ayuda a gestionar de forma eficiente las aplicaciones en contenedores a escala.

Esta guía profundiza en un conjunto de recomendaciones y prácticas recomendadas ampliamente reconocidas para diseñar y administrar aplicaciones de microservicios de Amazon EKS. Basados en una amplia experiencia y en implementaciones reales, estos conocimientos ofrecen una valiosa orientación para arquitectos y desarrolladores. Implemente estas recomendaciones para lograr un alto rendimiento, confiabilidad y escalabilidad de sus aplicaciones basadas en Kubernetes y lograr operaciones sólidas.

# Consideraciones de diseño de alta disponibilidad y resiliencia

El modelo de responsabilidad compartida se vuelve más complejo con Kubernetes. Amazon Web Services (AWS) gestiona la disponibilidad y la resiliencia del plano de control de Amazon EKS. Su organización administra el plano de datos, lo que puede afectar significativamente al rendimiento y la disponibilidad de sus aplicaciones de microservicios.

Al diseñar una aplicación flexible y de alta disponibilidad en Amazon EKS, tenga en cuenta los siguientes componentes:

- La aplicación de microservicios: sus cápsulas y contenedores
- El plano de datos de la carga de trabajo: controlador de ingreso, pod, componentes del sistema, como la [interfaz de red de contenedores \(CNI\) de Amazon Virtual Private Cloud \(Amazon VPC\)](#), [sidecars](#) de malla de servicios y kube-proxy
- La capa de administración de la carga de trabajo: controladores, controladores de admisión, motores de políticas de red y almacenamiento persistente de datos para estos componentes
- El plano de control de Kubernetes
- Infraestructura: nodos, red y dispositivos de red

En cuanto a las tres primeras consideraciones, que se refieren a los componentes que se ejecutan en un clúster de Kubernetes, en esta guía se tratan los siguientes temas:

- [Distribuir las cargas de trabajo entre los nodos y las zonas de disponibilidad](#)
- [Proteja las cargas de trabajo críticas con una PDB](#)
- [Configuración de sondas y controles de estado](#)
- [Configuración de los ganchos del ciclo de vida](#)
- [Cómo entender el desalojo de cápsulas durante las interrupciones zonales](#)

## Distribuya las cargas de trabajo entre los nodos y las zonas de disponibilidad

La distribución de una carga de trabajo entre [los dominios de error](#), como las zonas de disponibilidad y los nodos, mejora la disponibilidad de los componentes y reduce las probabilidades de que se produzcan errores en las aplicaciones escalables horizontalmente. En las siguientes secciones, se presentan formas de distribuir las cargas de trabajo entre los nodos y las zonas de disponibilidad.

### Utilice las restricciones de dispersión de la topología de los módulos

[Las restricciones de dispersión de la topología de los pods de Kubernetes](#) indican al programador de Kubernetes que distribuya los pods gestionados por diferentes dominios de error ReplicaSet o StatefulSet entre ellos (zonas de disponibilidad, nodos y tipos de hardware). Cuando utilizas las restricciones de dispersión de la topología de los pods, puedes hacer lo siguiente:

- Distribuya o concentre los pods en diferentes dominios de error en función de los requisitos de la aplicación. Por ejemplo, puede distribuir los pods para aumentar la resiliencia y concentrarlos para mejorar el rendimiento de la red.
- Combine diferentes condiciones, como la distribución entre zonas de disponibilidad y la distribución entre nodos.
- Especifique la acción preferida si no se pueden cumplir las condiciones:
  - Úselo `whenUnsatisfiable: DoNotSchedule` con una combinación de `maxSkew` y `minDomains` para crear requisitos estrictos para el programador.
  - Úselo `whenUnsatisfiable: ScheduleAnyway` para reducir `maxSkew`.

Si una zona de error deja de estar disponible, los pods de esa zona dejan de estar en buen estado. Kubernetes reprograma los pods y, si es posible, respeta la restricción de dispersión.

En el siguiente código se muestra un ejemplo del uso de las restricciones de distribución de la topología de los pods en las zonas de disponibilidad o en los nodos:

```
...
spec:
  selector:
    matchLabels:
      app: <your-app-label>
```

```
replicas: 3
template:
  metadata:
    labels: <your-app-label>
  spec:
    serviceAccountName: <ServiceAccountName>
...
  topologySpreadConstraints:
    - labelSelector:
        matchLabels:
          app: <your-app-label>
      maxSkew: 1
      topologyKey: topology.kubernetes.io/zone # <---spread those pods evenly over all availability zones
      whenUnsatisfiable: ScheduleAnyway
    - labelSelector:
        matchLabels:
          app: <your-app-label>
      maxSkew: 1
      topologyKey: kubernetes.io/hostname # <---spread those pods evenly over all nodes
      whenUnsatisfiable: ScheduleAnyway
```

## Restricciones de dispersión de la topología predeterminadas en todo el clúster

De forma predeterminada, Kubernetes proporciona un [conjunto de restricciones de distribución de la topología para distribuir los pods entre los](#) nodos y las zonas de disponibilidad:

```
defaultConstraints:
  - maxSkew: 3
    topologyKey: "kubernetes.io/hostname"
    whenUnsatisfiable: ScheduleAnyway
  - maxSkew: 5
    topologyKey: "topology.kubernetes.io/zone"
    whenUnsatisfiable: ScheduleAnyway
```

### Note

Las aplicaciones que necesitan diferentes tipos de restricciones topológicas pueden anular la política a nivel de clúster.

Las restricciones predeterminadas establecen un valor máximo maxSkew, lo que no resulta útil para las implementaciones que tienen un número reducido de pods. A partir de ahora, [no se KubeSchedulerConfiguration puede cambiar](#) en Amazon EKS. Si necesita aplicar otros conjuntos de restricciones de dispersión topológica, considere la posibilidad de utilizar un controlador de admisión mutante, como se indica en la sección siguiente. También puede controlar las restricciones de dispersión de la topología predeterminadas si ejecuta un programador alternativo. Sin embargo, la administración de programadores personalizados añade complejidad y puede tener implicaciones en la resiliencia del clúster y en la alta disponibilidad. Por estas razones, no recomendamos usar un programador alternativo únicamente para las restricciones de dispersión de la topología.

## La política de Gatekeeper sobre las restricciones de dispersión de la topología

[Otra opción para aplicar las restricciones de dispersión topológica es utilizar una política del proyecto Gatekeeper.](#) Las políticas de Gatekeeper se definen a nivel de aplicación.

Los siguientes ejemplos de código muestran el uso de una Gatekeeper OPA política para la implementación. Puede modificar la política según sus necesidades. Por ejemplo, aplique la política solo a las implementaciones que tengan la etiqueta HA=true o escriba una política similar utilizando un controlador de políticas diferente.

En este primer ejemplo ConstraintTemplate se muestra el uso de `conk8stopologyspreadrequired_template.yaml`:

```
apiVersion: templates.gatekeeper.sh/v1
kind: ConstraintTemplate
metadata:
  name: k8stopologyspreadrequired
spec:
  crd:
    spec:
      names:
        kind: K8sTopologySpreadRequired
      validation:
        openAPIV3Schema:
          type: object
          properties:
            message:
              type: string
  targets:
    - target: admission.k8s.gatekeeper.sh
```

```
rego: |
  package k8stopologyspreadrequired

  get_message(parameters, _default) =3D msg {
    not parameters.message
    msg :=_default
  }

  get_message(parameters, _default) =3D msg {
    msg := parameters.message
  }

  violation[{"msg": msg}] {
    input.review.kind.kind ="Deployment"
    not input.review.object.spec.template.spec.topologySpreadConstraint
    def_msg :"Pod Topology Spread Constraints are required for Deployments"
    msg :get_message(input.parameters, def_msg)
  }
}
```

El siguiente código muestra el manifiesto `k8stopologyspreadrequired_constraint.yaml` de `constraints` YAML:

```
apiVersion: constraints.gatekeeper.sh/v1beta1
kind: K8sTopologySpreadRequired
metadata:
  name: require-topologyspread-for-deployments
spec:
  match:
    kinds:
      - apiGroups: ["apps"]
        kinds: ["Deployment"]
  namespaces: ## Without these two lines will apply to the whole cluster
    - "example"
```

## Cuándo usar las restricciones de dispersión topológica

Considere la posibilidad de utilizar restricciones de dispersión topológica en los siguientes escenarios:

- Cualquier aplicación escalable horizontalmente (por ejemplo, servicios web sin estado)

- Aplicaciones con réplicas activas-activas o activas-pasivas (por ejemplo, bases de datos o cachés NoSQL)
- Aplicaciones con réplicas en espera (por ejemplo, controladores)

Entre los componentes del sistema que se pueden utilizar en el escenario de escalabilidad horizontal, por ejemplo, se incluyen los siguientes:

- [Autoscaler y Karpenter en clúster](#) (con y) `replicaCount > 1 leader-elect = true`
- [AWS Controlador Load Balancer](#)
- [CoreDNS](#)

## Afinidad y antiafinidad entre los pods

En algunos casos, es conveniente asegurarse de que no se ejecute más de un pod de un tipo específico en un nodo. Por ejemplo, para evitar programar varios pods con un uso intensivo de la red en el mismo nodo, puedes usar la regla de antiafinidad con la etiqueta `o. Ingress Network-heavy`. Cuando la utilices `anti-affinity`, también puedes usar una combinación de las siguientes opciones:

- Contaminaciones en los nodos optimizados para la red
- Las tolerancias correspondientes en los módulos con un uso intensivo de la red
- Afinidad de nodos o selector de nodos para garantizar que los pods con un uso intensivo de la red utilicen instancias optimizadas para la red

Como ejemplo, se utilizan pods con un uso intensivo de la red. Es posible que tengas requisitos diferentes, como la GPU, la memoria o el almacenamiento local. Para ver otros ejemplos de uso y opciones de configuración, consulta la documentación de [Kubernetes](#).

## Reequilibre los pods

En esta sección, se analizan dos enfoques para reequilibrar los pods en un clúster de Kubernetes. El primero usa el Descheduler para Kubernetes. El Descheduler ayuda a mantener la distribución de los pods al aplicar estrategias para eliminar los pods que infrinjan las restricciones de distribución topológica o las reglas de antiafinidad. El segundo enfoque utiliza la función de consolidación y empaquetado de contenedores de Karpenter. La consolidación evalúa y optimiza continuamente el

uso de los recursos al consolidar las cargas de trabajo en un menor número de nodos agrupados de manera más eficiente.

Le recomendamos que utilice Descheduler si no utiliza Karpenter. Si usa Karpenter y Cluster Autoscaler juntos, puede usar Descheduler con Cluster Autoscaler para grupos de nodos.

#### Descheduler para nodos sin grupo

No hay garantía de que se cumplan las restricciones topológicas cuando se eliminan los pods. Por ejemplo, reducir la escala de una implementación podría provocar un desequilibrio en la distribución de los módulos. Sin embargo, dado que Kubernetes utiliza las restricciones de dispersión de la topología de los pods solo en la fase de programación, los pods quedan desequilibrados en todo el dominio del error.

Para mantener una distribución de pods equilibrada en estos escenarios, puedes usar Descheduler for Kubernetes. Descheduler es una herramienta útil para múltiples propósitos, como hacer cumplir la edad máxima o el tiempo de vida (TTL) de los pods o mejorar el uso de la infraestructura. En el contexto de la resiliencia y la alta disponibilidad (HA), considere las siguientes estrategias de Descheduler:

- [RemovePodsViolatingTopologySpreadConstraint](#)
- [RemovePodsViolatingInterPodAntiAffinity](#)
- [RemoveDuplicates](#)

#### Función de consolidación y empaquetado de contenedores de Karpenter

En el caso de las cargas de trabajo que utilizan Karpenter, puede utilizar la funcionalidad de consolidación y empaquetado de contenedores para optimizar la utilización de los recursos y reducir los costes en los clústeres de Kubernetes. Karpenter evalúa continuamente la ubicación de los módulos y la utilización de los nodos, e intenta consolidar las cargas de trabajo en un menor número de nodos agrupados de forma más eficiente siempre que sea posible. Este proceso implica analizar los requisitos de recursos, tener en cuenta las limitaciones, como las reglas de afinidad de los módulos, y el posible traslado de los módulos entre nodos para mejorar la eficiencia general del clúster. En el siguiente código se proporciona un ejemplo:

```
apiVersion: karpenter.sh/v1beta1
kind: NodePool
metadata:
  name: default
```

```
spec:  
  disruption:  
    consolidationPolicy: WhenUnderutilized  
    expireAfter: 720h
```

Para `consolidationPolicy`, puede usar `WhenUnderutilized` o `WhenEmpty`:

- Si `consolidationPolicy` se establece en `WhenUnderutilized`, Karpenter considera todos los nodos para su consolidación. Cuando Karpenter descubre un nodo que está vacío o infráutilizado, intenta quitarlo o reemplazarlo para reducir los costes.
- Cuando `consolidationPolicy` se establece en `WhenEmpty`, Karpenter solo considera para la consolidación los nodos que no contienen módulos de carga de trabajo.

Las decisiones de consolidación de Karpenter no se basan únicamente en los porcentajes de utilización de la CPU o la memoria que pueden observarse en las herramientas de supervisión. En su lugar, Karpenter utiliza un algoritmo más complejo basado en las solicitudes de recursos de los módulos y en las posibles optimizaciones de costes. Para obtener más información, consulte la documentación de [Karpenter](#).

## Proteja las cargas de trabajo críticas con una PDB

Un presupuesto de interrupción del módulo (PDB) es una característica esencial para mantener la alta disponibilidad de las aplicaciones en un clúster. La PDB especifica un tamaño objetivo, que es la disponibilidad mínima para un tipo concreto de módulo. Esto significa que se debe ejecutar un número mínimo de réplicas de un tipo de pod concreto en un momento dado. Si el número de réplicas en ejecución es inferior al tamaño objetivo, Kubernetes evita que se produzcan más interrupciones en las réplicas restantes hasta que se alcance el tamaño objetivo. PDBs ayudan a garantizar que las cargas de trabajo no se vean afectadas por estos eventos y puedan seguir funcionando sin interrupciones. Cuando se produce una interrupción, Kubernetes intenta desalojar correctamente los pods de los nodos afectados y, al mismo tiempo, mantener el número de réplicas especificado en la PDB.

Puedes usar una PDB para declarar la cantidad y el número de réplicas. `minAvailable` `maxUnavailable` Por ejemplo, si quieras que estén disponibles al menos tres copias de tu aplicación, crea una PDB similar al ejemplo siguiente:

```
apiVersion: policy/v1beta1  
kind: PodDisruptionBudget
```

```
metadata:  
  name: my-svc-pdb  
spec:  
  minAvailable: 3  
  selector:  
    matchLabels:  
      app: my-svc
```

La configuración PDBs correcta de las aplicaciones ayuda a minimizar las interrupciones durante los eventos planificados o imprevistos. Puede usar la regla de antiafinidad para programar los módulos de una implementación en diferentes nodos y evitar demoras en la PDB durante las actualizaciones de los nodos.

## Configure las sondas y las comprobaciones del estado del equilibrador de carga

Kubernetes ofrece varias formas de realizar comprobaciones del estado de las aplicaciones, además de las comprobaciones del estado del balanceador de cargas. Puedes ejecutar las siguientes sondas integradas en Kubernetes junto con la comprobación del estado del balanceador de carga como un comando en el contexto del pod o como una sonda a kubelet o a la dirección IP del host. HTTP/TCP

Las sondas de actividad y de preparación deben ser diferentes e independientes (o, al menos, tener valores de tiempo de espera diferentes). Si una aplicación tiene un problema temporal, la sonda de disponibilidad marcará el módulo como no preparado hasta que se resuelva el problema. Si la configuración de la sonda de operatividad no es correcta, la sonda de operatividad podría cerrar el módulo.

### Sonda de inicio

Utilice sondas de inicio para proteger las aplicaciones que tienen ciclos de inicialización largos. Hasta que la sonda de inicio se realice correctamente, las demás sondas se desactivarán.

Puedes definir un tiempo máximo que Kubernetes debe esperar para que se inicie la aplicación. Si, transcurrido el tiempo máximo configurado, el pod sigue fallando, las sondas de inicio, la aplicación se cierra y se crea un nuevo pod.

Utilice sondeos de inicio cuando la hora de inicio de una aplicación sea impredecible. Si sabe que la aplicación necesita 10 segundos para iniciarse, utilice en su lugar una sonda de actividad o una sonda de disponibilidad. `initialDelaySeconds`

## Sonda de vitalidad

Utilice las sondas de dinámica para detectar problemas en las aplicaciones o si el proceso se está ejecutando sin problemas. Una sonda de actividad puede detectar situaciones de bloqueo en las que el proceso continúa ejecutándose pero la aplicación deja de responder. Cuando utilice una sonda de actividad, haga lo siguiente:

- Se utiliza `initialDelaySeconds` para retrasar la primera sonda.
- No establezca la misma especificación para las sondas de vitalidad y preparación.
- No configure una sonda de actividad para que dependa de un factor externo al pod (por ejemplo, una base de datos).
- Configura la sonda de vitalidad para una sonda específica. `terminationGracePeriodSeconds` Para obtener más información, consulte [Documentación Kubernetes](#) en la documentación de Kubernetes.

## Sonda de preparación

Utilice una sonda de preparación para detectar lo siguiente:

- Si la aplicación está lista para aceptar tráfico
- Disponibilidad parcial, en la que la aplicación puede no estar disponible temporalmente, pero se espera que vuelva a estar en buen estado una vez finalizada una determinada operación

Los sondeos de disponibilidad ayudan a garantizar que la configuración y las dependencias de la aplicación se ejecuten sin problemas ni errores, de modo que la aplicación pueda atender el tráfico. Sin embargo, una sonda de disponibilidad mal configurada puede provocar una interrupción en lugar de evitarla. Las sondas de disponibilidad que dependen de factores externos, como la conectividad a una base de datos, pueden provocar que todos los módulos no pasen la sonda. Estos errores pueden provocar una interrupción y provocar una falla en cascada desde un servicio de back-end a otros servicios que utilizaron los pods defectuosos.

## Comprobaciones de estado del balanceador de carga y recursos de Ingress

Application Load Balancer y ingress Kubernetes ofrecen funciones de comprobación de estado. Para las comprobaciones de estado de Application Load Balancer, especifique los puertos y la ruta de destino.

### Note

En el caso de Kubernetesingress, habrá una latencia de anulación del registro. El valor predeterminado de Application Load Balancer es de 300 segundos. Considere configurar el estado del recurso de entrada o del equilibrador de carga con los mismos valores que utilizó para la sonda de preparación.

NGINX también proporciona un chequeo de estado. Para obtener más información, consulte la documentación de [NGINX](#).

Las pasarelas de entrada y salida de Istio no tienen un mecanismo de comprobación de estado comparable al de la comprobación de estado HTTP de NGINX. Sin embargo, puede lograr una funcionalidad similar utilizando el disyuntor [Istio](#) o la detección de valores atípicos.

#### DestinationRule

Para obtener más información, consulte [Disponibilidad y ciclo de vida del pod](#) en la Guía de prácticas recomendadas de Amazon EKS.

## Configure los enlaces del ciclo de vida de

Durante un cierre correcto de un contenedor, la aplicación debería responder a una SIGTERM señal iniciando el cierre para que los clientes no sufran ningún tiempo de inactividad. La aplicación debe ejecutar procedimientos de limpieza como los siguientes:

- Guardar datos
- Cerrar los descriptores de los archivos
- Cerrar conexiones a bases de datos
- Completar correctamente las solicitudes durante el vuelo
- Salir puntualmente para cumplir con la solicitud de cierre del módulo

Establece un período de gracia que sea lo suficientemente largo como para que finalice la limpieza. Para saber cómo responder a la SIGTERM señal, consulte la documentación del lenguaje de programación que utilice para la aplicación.

[Los ganchos de ciclo](#) de vida de los contenedores permiten a los contenedores estar al tanto de los eventos de su ciclo de vida de administración. Los contenedores pueden ejecutar el código

implementado en un controlador cuando se ejecuta el enlace del ciclo de vida correspondiente. Los enlaces del ciclo de vida de los contenedores ofrecen una solución alternativa a la naturaleza asíncrona de Kubernetes y la nube. Este enfoque puede evitar la pérdida de las conexiones que se reenvían al pod de destino antes que el recurso de entrada y que `iptables` se actualizan para no enviar tráfico nuevo al pod.

El ciclo de vida del contenedor y `EndpointSlice` forman parte de un ciclo de vida diferente. `Endpoint APIs` Es importante organizarlos APIs. Sin embargo, cuando se cierra un pod, la API de Kubernetes notifica simultáneamente tanto al kubelet (para el ciclo de vida del contenedor) como al controlador. `EndpointSlice` Para obtener más información, incluido un diagrama, consulte [Gestionar correctamente las solicitudes de los clientes](#) en la Guía de prácticas recomendadas de Amazon EKS.

Cuando se kubelet envía SIGTERM al pod, el `EndpointSlice` controlador termina el `EndpointSlice` objeto. Esta terminación notifica a los servidores de la API de Kubernetes que notifiquen la actualización `kube-proxy` de cada nodo. `iptables` Aunque estas acciones se producen al mismo tiempo, no hay dependencias ni secuencias entre ellas. Existe una alta probabilidad de que el contenedor reciba la SIGKILL señal mucho antes de que `kube-proxy` en cada nodo se actualicen `iptables` las reglas locales. En ese caso, los posibles escenarios incluyen los siguientes:

- Si su solicitud descarta de forma inmediata y sin rodeos las solicitudes y conexiones durante el vuelo al recibirlasSIGTERM, los clientes verán errores. 500
- Si su solicitud garantiza que todas las solicitudes y conexiones durante el vuelo se procesen por completo al recibirlasSIGTERM, durante el período de gracia, las solicitudes de nuevos clientes seguirán enviándose al contenedor de solicitudes, ya que es posible que `iptables` las reglas aún no se hayan actualizado. Hasta que el procedimiento de limpieza cierre el socket del servidor del contenedor, esas nuevas solicitudes generarán nuevas conexiones. Cuando finalice el período de gracia, las nuevas conexiones que se establecieron después del envío SIGTERM se eliminarán incondicionalmente.

Para abordar los escenarios anteriores, puedes implementar la integración en la aplicación o el enlace del PreStop ciclo de vida. Para obtener más información, incluido un diagrama, consulte [Cierre correctamente las aplicaciones](#) en la Guía de prácticas recomendadas de Amazon EKS.

### Note

Independientemente de si la aplicación se cierra correctamente o como resultado del bloqueo, los preStop contenedores de la aplicación terminan cerrándose al final del período de gracia. SIGKILL

Usa el preStop gancho con un sleep comando para retrasar el envío. SIGTERM Esto ayudará a seguir aceptando las nuevas conexiones mientras el objeto de entrada las dirige al pod. Pruebe el valor temporal del sleep comando para asegurarse de que se tiene en cuenta cualquier latencia de Kubernetes y otras dependencias de la aplicación, como se muestra en el siguiente ejemplo:

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: nginx
spec:
  containers:
    - name: nginx
      lifecycle:
        # This "sleep" preStop hook delays the Pod shutdown until
        # after the Ingress Controller removes the matching Endpoint or EndpointSlice
        preStop:
          exec:
            command:
              - /bin/sleep
              - "20"
            # This period should be turned to Ingress/Service Mesh update latency
```

Para obtener más información, consulte [Container Hooks](#) en la documentación de Kubernetes y [Cierre correctamente las aplicaciones](#) en la Guía de prácticas recomendadas de Amazon EKS.

## Comprenda el desalojo de los pods durante las interrupciones zonales

Cuando se produce una interrupción total de la zona de disponibilidad (es decir, cuando todos los nodos de esa zona de disponibilidad pierden la conectividad con el plano de control de Kubernetes), el [controlador del ciclo de vida de los nodos de Kubernetes detecta la situación](#) y expulsa los pods de la zona afectada. Los pods de los nodos inalcanzables se marcan como nodos en buen

estado de las zonas de disponibilidad disponibles Terminating y se programan nuevos pods para los nodos en buen estado. Durante este período, los nodos afectados muestran un NotReady estado, el planificador impide que se coloquen nuevos módulos en esos nodos y el EndpointSlice controlador elimina del enrutamiento del servicio los puntos finales que están asociados a la zona de disponibilidad dañada hasta que se restablezca la conectividad.

En los casos en los que se producen fallos parciales en los nodos de una zona (en los que solo queda inalcanzable un subconjunto de nodos), el controlador del ciclo de vida de los nodos aplica un comportamiento de desalojo diferente. Si la interrupción persiste más allá del período de tolerancia configurado (de forma predeterminada, cinco minutos), los módulos de los nodos desconectados se marcan como Terminating y se programan nuevos módulos en los nodos en buen estado de las zonas de disponibilidad disponibles.

## Implementación del cambio zonal de Amazon EKS para mejorar la resiliencia

El [cambio zonal de Amazon EKS](#), que se integra con Amazon Application Recovery Controller (ARC), proporciona un mecanismo para gestionar el tráfico de forma proactiva durante las alteraciones de la zona de disponibilidad. Esta capacidad permite redirigir temporalmente el tráfico de red desde una zona de disponibilidad en mal estado hacia zonas en buen estado dentro de la misma, Región de AWS a fin de minimizar las interrupciones del servicio.

### Comprender el mecanismo de cambio zonal

El cambio zonal de Amazon EKS aborda el tráfico este-oeste (comunicación entre módulos dentro del clúster). Cuando el cambio zonal se configura con balanceadores de carga de aplicaciones o balanceadores de carga de red, también es compatible con el enrutamiento del tráfico de entrada. El mecanismo funciona mediante la coordinación de varios componentes de Kubernetes y del plano de AWS control para redirigir el tráfico de forma segura sin interrumpir las cargas de trabajo en ejecución. Durante un cambio zonal activo, Amazon EKS realiza automáticamente las siguientes acciones coordinadas:

- Acordonamiento de nodos: todos los nodos de la zona de disponibilidad afectada están acordonados. Esto evita que el programador de Kubernetes coloque nuevos módulos en los nodos mientras mantiene las cargas de trabajo existentes.
- Suspensión del reequilibrio de la zona de disponibilidad: en el caso de los grupos de nodos gestionados, las operaciones de reequilibrio de la zona de disponibilidad se suspenden y los

grupos de Auto Scaling se actualizan para lanzar nuevos nodos del plano de datos exclusivamente en zonas de disponibilidad en buen estado. Esto garantiza que no se aprovisione nueva capacidad en la zona afectada.

- Eliminación de los puntos finales: el EndpointSlice controlador elimina los puntos finales del módulo situados en la zona de disponibilidad dañada de todos los puntos de conexión relevantes. EndpointSlices Esto garantiza que los mecanismos de descubrimiento de servicios y equilibrio de carga dirijan el tráfico únicamente a los pods que se ejecutan en zonas de disponibilidad en buen estado.
- Preservación de la carga de trabajo: Amazon EKS se abstiene de cerrar los nodos o desalojar los pods de la zona de disponibilidad afectada. Mantiene toda la capacidad en la zona afectada para que, cuando el cambio zonal caduque o se cancele, el tráfico pueda regresar de forma segura sin necesidad de operaciones de escalado adicionales.

## Métodos de activación por turnos zonales

Puede elegir entre dos enfoques para iniciar los cambios zonales, según su modelo operativo:

- El [cambio zonal manual](#) proporciona un control controlado por el operador cuando se detectan problemas específicos en las zonas de disponibilidad mediante la supervisión, las alertas o los informes de los clientes. Este método requiere una acción explícita a través de la consola ARC, AWS Command Line Interface (AWS CLI) o un cambio APIs zonal, en el que los operadores especifican la zona de disponibilidad afectada y definen una hora de caducidad para el turno. Los cambios manuales son adecuados cuando los equipos cuentan con capacidades específicas de monitoreo y de guardia y prefieren mantener un control directo sobre las decisiones de gestión del tráfico.
- El [cambio automático zonal](#) autoriza AWS a iniciar automáticamente los turnos cuando ARC detecta posibles fallos o deficiencias en las zonas de disponibilidad en función de la telemetría interna y las señales de estado en varias zonas, Servicios de AWS incluidas las métricas de red, Amazon Elastic Compute Cloud (Amazon EC2) y Elastic Load Balancing. AWS finaliza automáticamente un cambio automático cuando los indicadores muestran que el problema se ha resuelto. Si desea obtener la máxima disponibilidad con una intervención manual mínima, le recomendamos este enfoque, ya que permite responder en menos de un minuto a las deficiencias detectadas en la zona de disponibilidad.

## Requisitos previos para un cambio zonal efectivo

Para que el cambio zonal proteja correctamente las aplicaciones durante las interrupciones de las zonas de disponibilidad, debe diseñar sus clústeres para que sean resilientes en zonas de disponibilidad múltiples antes de activar la función de cambio zonal:

- Distribución de nodos en zonas de disponibilidad múltiples: aprovisione los nodos de trabajo en al menos tres zonas de disponibilidad para garantizar una redundancia suficiente cuando una zona deje de estar disponible.
- Planificación de la capacidad: aprovisione previamente suficiente capacidad de cómputo en las zonas de disponibilidad en buen estado para dar cabida a toda la carga de trabajo cuando una zona de disponibilidad quede fuera del servicio, ya que si se escalan las operaciones durante una interrupción activa, es posible que la capacidad no sea suficiente.
- Distribución de pods y preescalado: [Implemente múltiples réplicas de cada aplicación en todas las zonas de disponibilidad y escale previamente los componentes críticos del sistema, como CoreDNS, en cada zona.](#) Esto ayuda a garantizar que quede suficiente capacidad después de que una zona se haya desplazado.

## Recomendaciones para la resiliencia a las disruptpciones zonales

- Habilite el cambio zonal al crear el clúster: en el caso de los nuevos clústeres de EKS, habilite la integración del cambio zonal con ARC durante el aprovisionamiento inicial a través de la consola Amazon EKS o de herramientas de infraestructura como código (iAC) AWS CLI, como AWS CloudFormation [Los clústeres del modo automático de EKS](#) que se crean con una configuración rápida tienen activado el cambio zonal de forma predeterminada.
- Seleccione el método de activación adecuado: elija el cambio automático zonal para los entornos de producción que requieren la máxima disponibilidad con una respuesta automática, especialmente para las aplicaciones orientadas al cliente, en las que los minutos de inactividad durante una zona de disponibilidad reducida pueden tener un impacto empresarial significativo. Utilice el cambio zonal manual en entornos en los que los equipos de operaciones prefieran dar su aprobación explícita antes de los cambios de tráfico o en los que las aplicaciones aún estén en curso de prueba y validación.
- Pruebe la resiliencia antes del despliegue de producción: valide el comportamiento del clúster en caso de pérdida en una única zona de disponibilidad iniciando manualmente los turnos zonales de prueba o habilitando la práctica de cambio automático zonal para comprobar que las

aplicaciones mantienen la disponibilidad, el rendimiento sigue siendo aceptable y la capacidad es suficiente cuando funcionan con un número reducido de zonas de disponibilidad. Recomendamos encarecidamente realizar estas pruebas para poder identificar las brechas de configuración antes de que se produzcan alteraciones reales en la zona de disponibilidad.

- Coordine con la configuración del balanceador de carga: en el caso de las aplicaciones que reciben tráfico externo, habilite el cambio zonal ARC en los balanceadores de carga de aplicaciones y en los balanceadores de carga de red asociados para garantizar que tanto el tráfico de entrada como el tráfico del clúster de este a oeste se desplacen al mismo tiempo cuando la zona de disponibilidad se deteriora. Esta coordinación evita situaciones en las que las solicitudes externas lleguen a los módulos en buen estado, pero estos no puedan comunicarse con las dependencias de la zona apartada.
- Supervise las operaciones de turno: después de activar el cambio zonal, configure la supervisión y las alertas de los eventos de turno, incluidas las activaciones de los turnos automáticos, los inicios de los turnos manuales y los vencimientos de los turnos, a fin de mantener la visibilidad operativa de las acciones de gestión del tráfico y su impacto en el comportamiento de las aplicaciones.

## Finalización y recuperación de los turnos

Cuando un turno zonal caduca en función de la duración configurada o se cancela manualmente una vez que se resuelve el deterioro de la zona de disponibilidad, el EndpointSlice controlador actualiza automáticamente todo EndpointSlices para volver a incorporar los puntos finales a la zona de disponibilidad restaurada. El tráfico regresa gradualmente a la zona afectada anteriormente a medida que los clientes actualizan la información de los puntos finales y establecen nuevas conexiones. Esto permite utilizar toda la capacidad del clúster sin necesidad de intervención manual ni reprogramar los módulos.

# Conclusión

Cuando diseñe su arquitectura para una alta disponibilidad y resiliencia de las aplicaciones, tenga en cuenta los siguientes componentes:

- La aplicación de microservicios (sus cápsulas y contenedores)
- El plano de datos de la carga de trabajo (controlador de ingreso, pod, componentes del sistema como el CNI de [Amazon VPC](#), sidecars de malla de servicios y kube-proxy)
- La capa de administración de la carga de trabajo (controladores, controladores de admisión, motores de políticas de red y almacenamiento persistente de datos para estos componentes)
- El plano de control de Kubernetes
- Infraestructura (nodos, red y dispositivos de red)

Para abordar esas consideraciones sobre los componentes, utilice las siguientes estrategias clave:

- Para garantizar una alta disponibilidad y una tolerancia a los errores, distribuya las cargas de trabajo entre los nodos y las zonas de disponibilidad.
- Para proteger las cargas de trabajo críticas, mantenga la estabilidad de las aplicaciones durante las interrupciones utilizando los presupuestos de interrupción de los módulos (). PDBs
- Para garantizar que los pods estén funcionando y gestionando el tráfico correctamente, configure las sondas de inicio, las sondas de actividad, las sondas de disponibilidad y las comprobaciones del estado del balanceador de carga.
- Para gestionar las transiciones de estado de los contenedores de forma eficiente, configura los enlaces del ciclo de vida de los contenedores.
- Para controlar el proceso de desalojo durante las fallas o el mantenimiento de los nodos, configure la hora de desalojo del pod.

Al implementar estas prácticas, puede mejorar considerablemente la confiabilidad y la resiliencia de las aplicaciones que se ejecutan en Amazon EKS, lo que garantiza un rendimiento sólido y una alta disponibilidad.

## Recursos

- Restricciones de dispersión de la topología del [pod de Kubernetes \(documentación de Kubernetes\)](#)
- [FAQsKarpenter](#) (documentación de Karpenter)
- [Descheduler para Kubernetes \(repositorio GitHub\)](#)
- [Disponibilidad y ciclo de vida de los pods](#): Guía de prácticas recomendadas de Amazon EKS
- [Cierre las aplicaciones sin problemas](#): Guía de mejores prácticas de Amazon EKS
- [\[EKS\] \[solicitud\]: Capacidad de configuración pod-eviction-timeout y soluciones alternativas](#) (repositorio de la hoja de ruta de contenedores)

## Historial de documentos

En la siguiente tabla, se describen cambios significativos de esta guía. Si quiere recibir notificaciones de futuras actualizaciones, puede suscribirse a las [notificaciones RSS](#).

Cambio	Descripción	Fecha
<a href="#"><u>Actualización</u></a>	Se revisó la sección sobre los <a href="#"><u>desalojos de cápsulas durante interrupciones zonales</u></a> .	29 de octubre de 2025
<a href="#"><u>Actualización</u></a>	Se revisó la sección <a href="#"><u>Utilizar restricciones de dispersión de la topología de cápsulas</u></a> .	27 de enero de 2025
<a href="#"><u>Publicación inicial</u></a>	—	23 de octubre de 2024

# AWS Glosario de orientación prescriptiva

Los siguientes son términos de uso común en las estrategias, guías y patrones proporcionados por la Guía AWS prescriptiva. Para sugerir entradas, utilice el enlace Enviar comentarios al final del glosario.

## Números

### Las 7 R

Siete estrategias de migración comunes para trasladar aplicaciones a la nube. Estas estrategias se basan en las 5 R que Gartner identificó en 2011 y consisten en lo siguiente:

- Refactorizar/rediseñar: traslade una aplicación y modifique su arquitectura mediante el máximo aprovechamiento de las características nativas en la nube para mejorar la agilidad, el rendimiento y la escalabilidad. Por lo general, esto implica trasladar el sistema operativo y la base de datos. Ejemplo: migre su base de datos Oracle local a la edición compatible con PostgreSQL de Amazon Aurora.
- Redefinir la plataforma (transportar y redefinir): traslade una aplicación a la nube e introduzca algún nivel de optimización para aprovechar las capacidades de la nube. Ejemplo: migre su base de datos Oracle local a Amazon Relational Database Service (Amazon RDS) para Oracle en el Nube de AWS
- Recomprar (readquirir): cambie a un producto diferente, lo cual se suele llevar a cabo al pasar de una licencia tradicional a un modelo SaaS. Ejemplo: migre su sistema de gestión de relaciones con los clientes (CRM) a Salesforce.com.
- Volver a alojar (migrar mediante lift-and-shift): traslade una aplicación a la nube sin realizar cambios para aprovechar las capacidades de la nube. Ejemplo: migre su base de datos Oracle local a Oracle en una EC2 instancia del Nube de AWS
- Reubicar: (migrar el hipervisor mediante lift and shift): traslade la infraestructura a la nube sin comprar equipo nuevo, reescribir aplicaciones o modificar las operaciones actuales. Los servidores se migran de una plataforma local a un servicio en la nube para la misma plataforma. Ejemplo: migrar una Microsoft Hyper-V aplicación a AWS.
- Retener (revisitar): conserve las aplicaciones en el entorno de origen. Estas pueden incluir las aplicaciones que requieren una refactorización importante, que desee posponer para más adelante, y las aplicaciones heredadas que desee retener, ya que no hay ninguna justificación empresarial para migrarlas.

- Retirar: retire o elimine las aplicaciones que ya no sean necesarias en un entorno de origen.

## A

### ABAC

Consulte control de [acceso basado en atributos](#).

### servicios abstractos

Consulte [servicios gestionados](#).

### ACID

Consulte [atomicidad, consistencia, aislamiento y durabilidad](#).

### migración activa-activa

Método de migración de bases de datos en el que las bases de datos de origen y destino se mantienen sincronizadas (mediante una herramienta de replicación bidireccional o mediante operaciones de escritura doble) y ambas bases de datos gestionan las transacciones de las aplicaciones conectadas durante la migración. Este método permite la migración en lotes pequeños y controlados, en lugar de requerir una transición única. Es más flexible, pero requiere más trabajo que la migración [activa-pasiva](#).

### migración activa-pasiva

Método de migración de bases de datos en el que las bases de datos de origen y destino se mantienen sincronizadas, pero solo la base de datos de origen gestiona las transacciones de las aplicaciones conectadas mientras los datos se replican en la base de datos de destino. La base de datos de destino no acepta ninguna transacción durante la migración.

### función de agregación

Función SQL que opera en un grupo de filas y calcula un único valor de retorno para el grupo. Entre los ejemplos de funciones agregadas se incluyen SUM y MAX.

### IA

Véase [inteligencia artificial](#).

### AIOps

Consulte las [operaciones de inteligencia artificial](#).

## anonimización

El proceso de eliminar permanentemente la información personal de un conjunto de datos. La anonimización puede ayudar a proteger la privacidad personal. Los datos anonimizados ya no se consideran datos personales.

## antipatrones

Una solución que se utiliza con frecuencia para un problema recurrente en el que la solución es contraproducente, ineficaz o menos eficaz que una alternativa.

## control de aplicaciones

Un enfoque de seguridad que permite el uso únicamente de aplicaciones aprobadas para ayudar a proteger un sistema contra el malware.

## cartera de aplicaciones

Recopilación de información detallada sobre cada aplicación que utiliza una organización, incluido el costo de creación y mantenimiento de la aplicación y su valor empresarial. Esta información es clave para [el proceso de detección y análisis de la cartera](#) y ayuda a identificar y priorizar las aplicaciones que se van a migrar, modernizar y optimizar.

## inteligencia artificial (IA)

El campo de la informática que se dedica al uso de tecnologías informáticas para realizar funciones cognitivas que suelen estar asociadas a los seres humanos, como el aprendizaje, la resolución de problemas y el reconocimiento de patrones. Para más información, consulte [¿Qué es la inteligencia artificial?](#)

## operaciones de inteligencia artificial (AIOps)

El proceso de utilizar técnicas de machine learning para resolver problemas operativos, reducir los incidentes operativos y la intervención humana, y mejorar la calidad del servicio. Para obtener más información sobre cómo AIOps se utiliza en la estrategia de AWS migración, consulte la [guía de integración de operaciones](#).

## cifrado asimétrico

Algoritmo de cifrado que utiliza un par de claves, una clave pública para el cifrado y una clave privada para el descifrado. Puede compartir la clave pública porque no se utiliza para el descifrado, pero el acceso a la clave privada debe estar sumamente restringido.

## atomicidad, consistencia, aislamiento, durabilidad (ACID)

Conjunto de propiedades de software que garantizan la validez de los datos y la fiabilidad operativa de una base de datos, incluso en caso de errores, cortes de energía u otros problemas.

control de acceso basado en atributos (ABAC)

La práctica de crear permisos detallados basados en los atributos del usuario, como el departamento, el puesto de trabajo y el nombre del equipo. Para obtener más información, consulte [ABAC AWS en la](#) documentación AWS Identity and Access Management (IAM).

## origen de datos fidedigno

Ubicación en la que se almacena la versión principal de los datos, que se considera la fuente de información más fiable. Puede copiar los datos del origen de datos autorizado a otras ubicaciones con el fin de procesarlos o modificarlos, por ejemplo, anonimizarlos, redactarlos o seudonimizarlos.

## Zona de disponibilidad

Una ubicación distinta dentro de una Región de AWS que está aislada de los fallos en otras zonas de disponibilidad y que proporciona una conectividad de red económica y de baja latencia a otras zonas de disponibilidad de la misma región.

## AWS Marco de adopción de la nube (AWS CAF)

Un marco de directrices y mejores prácticas AWS para ayudar a las organizaciones a desarrollar un plan eficiente y eficaz para migrar con éxito a la nube. AWS CAF organiza la orientación en seis áreas de enfoque denominadas perspectivas: negocios, personas, gobierno, plataforma, seguridad y operaciones. Las perspectivas empresariales, humanas y de gobernanza se centran en las habilidades y los procesos empresariales; las perspectivas de plataforma, seguridad y operaciones se centran en las habilidades y los procesos técnicos. Por ejemplo, la perspectiva humana se dirige a las partes interesadas que se ocupan de los Recursos Humanos (RR. HH.), las funciones del personal y la administración de las personas. Desde esta perspectiva, AWS CAF proporciona orientación para el desarrollo, la formación y la comunicación de las personas a fin de preparar a la organización para una adopción exitosa de la nube. Para obtener más información, consulte la [Página web de AWS CAF](#) y el [Documento técnico de AWS CAF](#).

## AWS Marco de calificación de la carga de trabajo (AWS WQF)

Herramienta que evalúa las cargas de trabajo de migración de bases de datos, recomienda estrategias de migración y proporciona estimaciones de trabajo. AWS WQF se incluye con AWS

Schema Conversion Tool (.AWS SCT) Analiza los esquemas de bases de datos y los objetos de código, el código de las aplicaciones, las dependencias y las características de rendimiento y proporciona informes de evaluación.

## B

Un bot malo

Un [bot](#) destinado a interrumpir o causar daño a personas u organizaciones.

BCP

Consulte la [planificación de la continuidad del negocio](#).

gráfico de comportamiento

Una vista unificada e interactiva del comportamiento de los recursos y de las interacciones a lo largo del tiempo. Puede utilizar un gráfico de comportamiento con Amazon Detective para examinar los intentos de inicio de sesión fallidos, las llamadas sospechosas a la API y acciones similares. Para obtener más información, consulte [Datos en un gráfico de comportamiento](#) en la documentación de Detective.

sistema big-endian

Un sistema que almacena primero el byte más significativo. Véase también [endianness](#).

clasificación binaria

Un proceso que predice un resultado binario (una de las dos clases posibles). Por ejemplo, es posible que su modelo de ML necesite predecir problemas como “¿Este correo electrónico es spam o no es spam?” o “¿Este producto es un libro o un automóvil?”.

filtro de floración

Estructura de datos probabilística y eficiente en términos de memoria que se utiliza para comprobar si un elemento es miembro de un conjunto.

implementación azul/verde

Una estrategia de despliegue en la que se crean dos entornos separados pero idénticos. La versión actual de la aplicación se ejecuta en un entorno (azul) y la nueva versión de la aplicación en el otro entorno (verde). Esta estrategia le ayuda a revertirla rápidamente con un impacto mínimo.

**bot**

Aplicación de software que ejecuta tareas automatizadas a través de Internet y simula la actividad o interacción humana. Algunos bots son útiles o beneficiosos, como los rastreadores web que indexan información en Internet. Algunos otros bots, conocidos como bots malos, tienen como objetivo interrumpir o causar daños a personas u organizaciones.

**botnet**

Redes de [bots](#) que están infectadas por [malware](#) y que están bajo el control de una sola parte, conocida como pastor u operador de bots. Las botnets son el mecanismo más conocido para escalar los bots y su impacto.

**branch**

Área contenida de un repositorio de código. La primera rama que se crea en un repositorio es la rama principal. Puede crear una rama nueva a partir de una rama existente y, a continuación, desarrollar características o corregir errores en la rama nueva. Una rama que se genera para crear una característica se denomina comúnmente rama de característica. Cuando la característica se encuentra lista para su lanzamiento, se vuelve a combinar la rama de característica con la rama principal. Para obtener más información, consulte [Acerca de las sucursales](#) (GitHub documentación).

**acceso con cristales rotos**

En circunstancias excepcionales y mediante un proceso aprobado, un usuario puede acceder rápidamente a un sitio para el Cuenta de AWS que normalmente no tiene permisos de acceso. Para obtener más información, consulte el indicador [Implemente procedimientos de rotura de cristales en la guía Well-Architected AWS](#).

**estrategia de implementación sobre infraestructura existente**

La infraestructura existente en su entorno. Al adoptar una estrategia de implementación sobre infraestructura existente para una arquitectura de sistemas, se diseña la arquitectura en función de las limitaciones de los sistemas y la infraestructura actuales. Si está ampliando la infraestructura existente, puede combinar las estrategias de implementación sobre infraestructuras existentes y de [implementación desde cero](#).

**caché de búfer**

El área de memoria donde se almacenan los datos a los que se accede con más frecuencia.

## capacidad empresarial

Lo que hace una empresa para generar valor (por ejemplo, ventas, servicio al cliente o marketing). Las arquitecturas de microservicios y las decisiones de desarrollo pueden estar impulsadas por las capacidades empresariales. Para obtener más información, consulte la sección [Organizado en torno a las capacidades empresariales](#) del documento técnico [Ejecutar microservicios en contenedores en AWS](#).

## planificación de la continuidad del negocio (BCP)

Plan que aborda el posible impacto de un evento disruptivo, como una migración a gran escala en las operaciones y permite a la empresa reanudar las operaciones rápidamente.

## C

### CAF

[Consulte el marco AWS de adopción de la nube.](#)

### despliegue canario

El lanzamiento lento e incremental de una versión para los usuarios finales. Cuando se tiene confianza, se despliega la nueva versión y se reemplaza la versión actual en su totalidad.

### CCoE

[Consulte Cloud Center of Excellence.](#)

### CDC

[Consulte la captura de datos de cambios.](#)

### captura de datos de cambio (CDC)

Proceso de seguimiento de los cambios en un origen de datos, como una tabla de base de datos, y registro de los metadatos relacionados con el cambio. Puede utilizar los CDC para diversos fines, como auditar o replicar los cambios en un sistema de destino para mantener la sincronización.

### ingeniería del caos

Introducir intencionalmente fallos o eventos disruptivos para poner a prueba la resiliencia de un sistema. Puedes usar [AWS Fault Injection Service \(AWS FIS\)](#) para realizar experimentos que estresen tus AWS cargas de trabajo y evalúen su respuesta.

## CI/CD

Consulte la [integración continua y la entrega continua](#).

## clasificación

Un proceso de categorización que permite generar predicciones. Los modelos de ML para problemas de clasificación predicen un valor discreto. Los valores discretos siempre son distintos entre sí. Por ejemplo, es posible que un modelo necesite evaluar si hay o no un automóvil en una imagen.

## cifrado del cliente

Cifrado de datos localmente, antes de que el objetivo los Servicio de AWS reciba.

## Centro de excelencia en la nube (CCoE)

Equipo multidisciplinario que impulsa los esfuerzos de adopción de la nube en toda la organización, incluido el desarrollo de las prácticas recomendadas en la nube, la movilización de recursos, el establecimiento de plazos de migración y la dirección de la organización durante las transformaciones a gran escala. Para obtener más información, consulte las [publicaciones de CCoE](#) en el blog de estrategia Nube de AWS empresarial.

## computación en la nube

La tecnología en la nube que se utiliza normalmente para la administración de dispositivos de IoT y el almacenamiento de datos de forma remota. La computación en la nube suele estar conectada a la tecnología de [computación perimetral](#).

## modelo operativo en la nube

En una organización de TI, el modelo operativo que se utiliza para crear, madurar y optimizar uno o más entornos de nube. Para obtener más información, consulte [Creación de su modelo operativo de nube](#).

## etapas de adopción de la nube

Las cuatro fases por las que suelen pasar las organizaciones cuando migran a Nube de AWS:

- Proyecto: ejecución de algunos proyectos relacionados con la nube con fines de prueba de concepto y aprendizaje
- Fundamento: realizar inversiones fundamentales para escalar su adopción de la nube (p. ej., crear una landing zone, definir una CCoE, establecer un modelo de operaciones)

- Migración: migración de aplicaciones individuales
- Reinención: optimización de productos y servicios e innovación en la nube

Stephen Orban definió estas etapas en la entrada del blog [The Journey Toward Cloud-First & the Stages of Adoption en el](#) blog Nube de AWS Enterprise Strategy. Para obtener información sobre su relación con la estrategia de AWS migración, consulte la guía de [preparación para la migración.](#)

## CMDB

Consulte la [base de datos de administración de la configuración.](#)  
repositorio de código

Una ubicación donde el código fuente y otros activos, como documentación, muestras y scripts, se almacenan y actualizan mediante procesos de control de versiones. Los repositorios en la nube más comunes incluyen GitHub o Bitbucket Cloud. Cada versión del código se denomina rama. En una estructura de microservicios, cada repositorio se encuentra dedicado a una única funcionalidad. Una sola canalización de CI/CD puede utilizar varios repositorios.

## caché en frío

Una caché de búfer que está vacía no está bien poblada o contiene datos obsoletos o irrelevantes. Esto afecta al rendimiento, ya que la instancia de la base de datos debe leer desde la memoria principal o el disco, lo que es más lento que leer desde la memoria caché del búfer.

## datos fríos

Datos a los que se accede con poca frecuencia y que suelen ser históricos. Al consultar este tipo de datos, normalmente se aceptan consultas lentas. Trasladar estos datos a niveles o clases de almacenamiento de menor rendimiento y menos costosos puede reducir los costos.

## visión artificial (CV)

Campo de la [IA](#) que utiliza el aprendizaje automático para analizar y extraer información de formatos visuales, como imágenes y videos digitales. Por ejemplo, Amazon SageMaker AI proporciona algoritmos de procesamiento de imágenes para CV.

## desviación de configuración

En el caso de una carga de trabajo, un cambio de configuración con respecto al estado esperado. Puede provocar que la carga de trabajo deje de cumplir las normas y, por lo general, es gradual e involuntario.

## base de datos de administración de configuración (CMDB)

Repositorio que almacena y administra información sobre una base de datos y su entorno de TI, incluidos los componentes de hardware y software y sus configuraciones. Por lo general, los datos de una CMDB se utilizan en la etapa de detección y análisis de la cartera de productos durante la migración.

## paquete de conformidad

Conjunto de AWS Config reglas y medidas correctivas que puede reunir para personalizar sus comprobaciones de conformidad y seguridad. Puede implementar un paquete de conformidad como una entidad única en una región Cuenta de AWS y, o en una organización, mediante una plantilla YAML. Para obtener más información, consulta los [paquetes de conformidad](#) en la documentación. AWS Config

## integración y entrega continuas (CI/CD)

El proceso de automatización de las etapas de origen, compilación, prueba, puesta en escena y producción del proceso de publicación del software. CI/CD se describe comúnmente como una canalización. CI/CD puede ayudarlo a automatizar los procesos, mejorar la productividad, mejorar la calidad del código y entregar más rápido. Para obtener más información, consulte [Beneficios de la entrega continua](#). CD también puede significar implementación continua. Para obtener más información, consulte [Entrega continua frente a implementación continua](#).

## CV

Vea la [visión artificial](#).

## D

### datos en reposo

Datos que están estacionarios en la red, como los datos que se encuentran almacenados.

### clasificación de datos

Un proceso para identificar y clasificar los datos de su red en función de su importancia y sensibilidad. Es un componente fundamental de cualquier estrategia de administración de riesgos de ciberseguridad porque lo ayuda a determinar los controles de protección y retención adecuados para los datos. La clasificación de datos es un componente del pilar de seguridad

del AWS Well-Architected Framework. Para obtener más información, consulte [Clasificación de datos](#).

#### desviación de datos

Una variación significativa entre los datos de producción y los datos que se utilizaron para entrenar un modelo de machine learning, o un cambio significativo en los datos de entrada a lo largo del tiempo. La desviación de los datos puede reducir la calidad, la precisión y la imparcialidad generales de las predicciones de los modelos de machine learning.

#### datos en tránsito

Datos que se mueven de forma activa por la red, por ejemplo, entre los recursos de la red.

#### malla de datos

Un marco arquitectónico que proporciona una propiedad de datos distribuida y descentralizada con administración y gobierno centralizados.

#### minimización de datos

El principio de recopilar y procesar solo los datos estrictamente necesarios. Practicar la minimización de los datos Nube de AWS puede reducir los riesgos de privacidad, los costos y la huella de carbono de la analítica.

#### perímetro de datos

Un conjunto de barreras preventivas en su AWS entorno que ayudan a garantizar que solo las identidades confiables accedan a los recursos confiables desde las redes esperadas. Para obtener más información, consulte [Crear un perímetro de datos sobre](#). AWS

#### preprocesamiento de datos

Transformar los datos sin procesar en un formato que su modelo de ML pueda analizar fácilmente. El preprocesamiento de datos puede implicar eliminar determinadas columnas o filas y corregir los valores faltantes, incoherentes o duplicados.

#### procedencia de los datos

El proceso de rastrear el origen y el historial de los datos a lo largo de su ciclo de vida, por ejemplo, la forma en que se generaron, transmitieron y almacenaron los datos.

#### titular de los datos

Persona cuyos datos se recopilan y procesan.

## almacenamiento de datos

Un sistema de administración de datos que respalte la inteligencia empresarial, como el análisis.

Los almacenes de datos suelen contener grandes cantidades de datos históricos y, por lo general, se utilizan para consultas y análisis.

## lenguaje de definición de datos (DDL)

Instrucciones o comandos para crear o modificar la estructura de tablas y objetos de una base de datos.

## lenguaje de manipulación de datos (DML)

Instrucciones o comandos para modificar (insertar, actualizar y eliminar) la información de una base de datos.

## DDL

Consulte el [lenguaje de definición de bases](#) de datos.

## conjunto profundo

Combinar varios modelos de aprendizaje profundo para la predicción. Puede utilizar conjuntos profundos para obtener una predicción más precisa o para estimar la incertidumbre de las predicciones.

## aprendizaje profundo

Un subcampo del ML que utiliza múltiples capas de redes neuronales artificiales para identificar el mapeo entre los datos de entrada y las variables objetivo de interés.

## defense-in-depth

Un enfoque de seguridad de la información en el que se distribuyen cuidadosamente una serie de mecanismos y controles de seguridad en una red informática para proteger la confidencialidad, la integridad y la disponibilidad de la red y de los datos que contiene. Al adoptar esta estrategia AWS, se añaden varios controles en diferentes capas de la AWS Organizations estructura para ayudar a proteger los recursos. Por ejemplo, un defense-in-depth enfoque podría combinar la autenticación multifactorial, la segmentación de la red y el cifrado.

## administrador delegado

En AWS Organizations, un servicio compatible puede registrar una cuenta de AWS miembro para administrar las cuentas de la organización y gestionar los permisos de ese servicio. Esta

cuenta se denomina administrador delegado para ese servicio. Para obtener más información y una lista de servicios compatibles, consulte [Servicios que funcionan con AWS Organizations](#) en la documentación de AWS Organizations .

## Implementación

El proceso de hacer que una aplicación, características nuevas o correcciones de código se encuentren disponibles en el entorno de destino. La implementación abarca implementar cambios en una base de código y, a continuación, crear y ejecutar esa base en los entornos de la aplicación.

entorno de desarrollo

Consulte [entorno](#).

control de detección

Un control de seguridad que se ha diseñado para detectar, registrar y alertar después de que se produzca un evento. Estos controles son una segunda línea de defensa, ya que lo advierten sobre los eventos de seguridad que han eludido los controles preventivos establecidos. Para obtener más información, consulte [Controles de detección](#) en Implementación de controles de seguridad en AWS.

asignación de flujos de valor para el desarrollo (DVSM)

Proceso que se utiliza para identificar y priorizar las restricciones que afectan negativamente a la velocidad y la calidad en el ciclo de vida del desarrollo de software. DVSM amplía el proceso de asignación del flujo de valor diseñado originalmente para las prácticas de fabricación ajustada. Se centra en los pasos y los equipos necesarios para crear y transferir valor a través del proceso de desarrollo de software.

gemelo digital

Representación virtual de un sistema del mundo real, como un edificio, una fábrica, un equipo industrial o una línea de producción. Los gemelos digitales son compatibles con el mantenimiento predictivo, la supervisión remota y la optimización de la producción.

tabla de dimensiones

En un [esquema en estrella](#), tabla más pequeña que contiene los atributos de datos sobre los datos cuantitativos de una tabla de hechos. Los atributos de la tabla de dimensiones suelen ser campos de texto o números discretos que se comportan como texto. Estos atributos se utilizan habitualmente para restringir consultas, filtrar y etiquetar conjuntos de resultados.

## desastre

Un evento que impide que una carga de trabajo o un sistema cumplan sus objetivos empresariales en su ubicación principal de implementación. Estos eventos pueden ser desastres naturales, fallos técnicos o el resultado de acciones humanas, como una configuración incorrecta involuntaria o un ataque de malware.

## recuperación de desastres (DR)

La estrategia y el proceso que se utilizan para minimizar el tiempo de inactividad y la pérdida de datos ocasionados por un [desastre](#). Para obtener más información, consulte [Recuperación ante desastres de cargas de trabajo en AWS: Recovery in the Cloud in the AWS Well-Architected Framework](#).

## DML

Consulte el lenguaje de manipulación de [bases de datos](#).

## diseño basado en el dominio

Un enfoque para desarrollar un sistema de software complejo mediante la conexión de sus componentes a dominios en evolución, o a los objetivos empresariales principales, a los que sirve cada componente. Este concepto lo introdujo Eric Evans en su libro, Diseño impulsado por el dominio: abordando la complejidad en el corazón del software (Boston: Addison-Wesley Professional, 2003). Para obtener información sobre cómo utilizar el diseño basado en dominios con el patrón de higos estranguladores, consulte [Modernización gradual de los servicios web antiguos de Microsoft ASP.NET \(ASMX\) mediante contenedores y Amazon API Gateway](#).

## DR

Consulte [recuperación ante desastres](#).

## detección de deriva

Seguimiento de las desviaciones con respecto a una configuración de referencia. Por ejemplo, puedes usarlo AWS CloudFormation para [detectar desviaciones en los recursos del sistema](#) o puedes usarlo AWS Control Tower para [detectar cambios en tu landing zone](#) que puedan afectar al cumplimiento de los requisitos de gobierno.

## DVSM

Consulte [el mapeo del flujo de valor del desarrollo](#).

# E

## EDA

Consulte el [análisis exploratorio de datos](#).

## EDI

Véase [intercambio electrónico de datos](#).

## computación en la periferia

La tecnología que aumenta la potencia de cálculo de los dispositivos inteligentes en la periferia de una red de IoT. En comparación con [la computación en nube, la computación perimetral](#) puede reducir la latencia de la comunicación y mejorar el tiempo de respuesta.

## intercambio electrónico de datos (EDI)

El intercambio automatizado de documentos comerciales entre organizaciones. Para obtener más información, consulte [Qué es el intercambio electrónico de datos](#).

## cifrado

Proceso informático que transforma datos de texto plano, legibles por humanos, en texto cifrado.

## clave de cifrado

Cadena criptográfica de bits aleatorios que se genera mediante un algoritmo de cifrado. Las claves pueden variar en longitud y cada una se ha diseñado para ser impredecible y única.

## endianidad

El orden en el que se almacenan los bytes en la memoria del ordenador. Los sistemas big-endianos almacenan primero el byte más significativo. Los sistemas Little-Endian almacenan primero el byte menos significativo.

## punto de conexión

[Consulte el punto final del servicio](#).

## servicio de punto de conexión

Servicio que puede alojar en una nube privada virtual (VPC) para compartir con otros usuarios. Puede crear un servicio de punto final AWS PrivateLink y conceder permisos a otros directores

Cuentas de AWS o a AWS Identity and Access Management (IAM). Estas cuentas o entidades principales pueden conectarse a su servicio de punto de conexión de forma privada mediante la creación de puntos de conexión de VPC de interfaz. Para obtener más información, consulte [Creación de un servicio de punto de conexión](#) en la documentación de Amazon Virtual Private Cloud (Amazon VPC).

## planificación de recursos empresariales (ERP)

Un sistema que automatiza y gestiona los procesos empresariales clave (como la contabilidad, el [MES](#) y la gestión de proyectos) de una empresa.

## cifrado de sobre

El proceso de cifrar una clave de cifrado con otra clave de cifrado. Para obtener más información, consulte el [cifrado de sobres](#) en la documentación de AWS Key Management Service (AWS KMS).

## entorno

Una instancia de una aplicación en ejecución. Los siguientes son los tipos de entornos más comunes en la computación en la nube:

- entorno de desarrollo: instancia de una aplicación en ejecución que solo se encuentra disponible para el equipo principal responsable del mantenimiento de la aplicación. Los entornos de desarrollo se utilizan para probar los cambios antes de promocionarlos a los entornos superiores. Este tipo de entorno a veces se denomina entorno de prueba.
- entornos inferiores: todos los entornos de desarrollo de una aplicación, como los que se utilizan para las compilaciones y pruebas iniciales.
- entorno de producción: instancia de una aplicación en ejecución a la que pueden acceder los usuarios finales. En un CI/CD proceso, el entorno de producción es el último entorno de implementación.
- entornos superiores: todos los entornos a los que pueden acceder usuarios que no sean del equipo de desarrollo principal. Esto puede incluir un entorno de producción, entornos de preproducción y entornos para las pruebas de aceptación por parte de los usuarios.

## epopeya

En las metodologías ágiles, son categorías funcionales que ayudan a organizar y priorizar el trabajo. Las epopeyas brindan una descripción detallada de los requisitos y las tareas de implementación. Por ejemplo, las epopeyas AWS de seguridad de CAF incluyen la gestión de identidades y accesos, los controles de detección, la seguridad de la infraestructura, la protección

de datos y la respuesta a incidentes. Para obtener más información sobre las epopeyas en la estrategia de migración de AWS , consulte la [Guía de implementación del programa](#).

## PERP

Consulte [planificación de recursos empresariales](#).

análisis de datos de tipo exploratorio (EDA)

El proceso de analizar un conjunto de datos para comprender sus características principales. Se recopilan o agregan datos y, a continuación, se realizan las investigaciones iniciales para encontrar patrones, detectar anomalías y comprobar las suposiciones. El EDA se realiza mediante el cálculo de estadísticas resumidas y la creación de visualizaciones de datos.

## F

tabla de datos

La tabla central de un [esquema en forma de estrella](#). Almacena datos cuantitativos sobre las operaciones comerciales. Normalmente, una tabla de hechos contiene dos tipos de columnas: las que contienen medidas y las que contienen una clave externa para una tabla de dimensiones.

fallan rápidamente

Una filosofía que utiliza pruebas frecuentes e incrementales para reducir el ciclo de vida del desarrollo. Es una parte fundamental de un enfoque ágil.

límite de aislamiento de fallas

En el Nube de AWS, un límite, como una zona de disponibilidad Región de AWS, un plano de control o un plano de datos, que limita el efecto de una falla y ayuda a mejorar la resiliencia de las cargas de trabajo. Para obtener más información, consulte [Límites de AWS aislamiento](#) de errores.

rama de característica

Consulte la [sucursal](#).

características

Los datos de entrada que se utilizan para hacer una predicción. Por ejemplo, en un contexto de fabricación, las características pueden ser imágenes que se capturan periódicamente desde la línea de fabricación.

## importancia de las características

La importancia que tiene una característica para las predicciones de un modelo. Por lo general, esto se expresa como una puntuación numérica que se puede calcular mediante diversas técnicas, como las explicaciones aditivas de Shapley (SHAP) y los gradientes integrados. Para obtener más información, consulte [Interpretabilidad del modelo de aprendizaje automático con AWS](#).

## transformación de funciones

Optimizar los datos para el proceso de ML, lo que incluye enriquecer los datos con fuentes adicionales, escalar los valores o extraer varios conjuntos de información de un solo campo de datos. Esto permite que el modelo de ML se beneficie de los datos. Por ejemplo, si divide la fecha del “27 de mayo de 2021 00:15:37” en “jueves”, “mayo”, “2021” y “15”, puede ayudar al algoritmo de aprendizaje a aprender patrones matizados asociados a los diferentes componentes de los datos.

## indicaciones de unos pocos pasos

Proporcionar a un [LLM](#) un pequeño número de ejemplos que demuestren la tarea y el resultado deseado antes de pedirle que realice una tarea similar. Esta técnica es una aplicación del aprendizaje contextual, en el que los modelos aprenden a partir de ejemplos (planos) integrados en las instrucciones. Las indicaciones con pocas tomas pueden ser eficaces para tareas que requieren un formato, un razonamiento o un conocimiento del dominio específicos. [Consulte también el apartado de mensajes sin intervención.](#)

## FGAC

Consulte el control [de acceso detallado](#).

## control de acceso preciso (FGAC)

El uso de varias condiciones que tienen por objetivo permitir o denegar una solicitud de acceso.

## migración relámpago

Método de migración de bases de datos que utiliza la replicación continua de datos mediante la [captura de datos modificados](#) para migrar los datos en el menor tiempo posible, en lugar de utilizar un enfoque gradual. El objetivo es reducir al mínimo el tiempo de inactividad.

## FM

Consulte el [modelo básico](#).

## modelo de base (FM)

Una gran red neuronal de aprendizaje profundo que se ha estado entrenando con conjuntos de datos masivos de datos generalizados y sin etiquetar. FMs son capaces de realizar una amplia variedad de tareas generales, como comprender el lenguaje, generar texto e imágenes y conversar en lenguaje natural. Para obtener más información, consulte [Qué son los modelos básicos](#).

# G

## IA generativa

Un subconjunto de modelos de [IA](#) que se han entrenado con grandes cantidades de datos y que pueden utilizar un simple mensaje de texto para crear contenido y artefactos nuevos, como imágenes, vídeos, texto y audio. Para obtener más información, consulte [Qué es la IA generativa](#).

## bloqueo geográfico

Consulta [las restricciones geográficas](#).

## restricciones geográficas (bloqueo geográfico)

En Amazon CloudFront, una opción para impedir que los usuarios de países específicos accedan a las distribuciones de contenido. Puede utilizar una lista de permitidos o bloqueados para especificar los países aprobados y prohibidos. Para obtener más información, consulta [la sección Restringir la distribución geográfica del contenido](#) en la CloudFront documentación.

## Flujo de trabajo de Gitflow

Un enfoque en el que los entornos inferiores y superiores utilizan diferentes ramas en un repositorio de código fuente. El flujo de trabajo de Gitflow se considera heredado, y el [flujo de trabajo basado en enlaces troncales](#) es el enfoque moderno preferido.

## imagen dorada

Instantánea de un sistema o software que se utiliza como plantilla para implementar nuevas instancias de ese sistema o software. Por ejemplo, en la fabricación, una imagen dorada se puede utilizar para aprovisionar software en varios dispositivos y ayuda a mejorar la velocidad, la escalabilidad y la productividad de las operaciones de fabricación de dispositivos.

## estrategia de implementación desde cero

La ausencia de infraestructura existente en un entorno nuevo. Al adoptar una estrategia de implementación desde cero para una arquitectura de sistemas, puede seleccionar todas las tecnologías nuevas sin que estas deban ser compatibles con una infraestructura existente, lo que también se conoce como [implementación sobre infraestructura existente](#). Si está ampliando la infraestructura existente, puede combinar las estrategias de implementación sobre infraestructuras existentes y de implementación desde cero.

## barrera de protección

Una regla de alto nivel que ayuda a regular los recursos, las políticas y el cumplimiento en todas las unidades organizativas (OUs). Las barreras de protección preventivas aplican políticas para garantizar la alineación con los estándares de conformidad. Se implementan mediante políticas de control de servicios y límites de permisos de IAM. Las barreras de protección de detección detectan las vulneraciones de las políticas y los problemas de conformidad, y generan alertas para su corrección. Se implementan mediante Amazon AWS Config AWS Security Hub CSPM GuardDuty AWS Trusted Advisor, Amazon Inspector y AWS Lambda cheques personalizados.

# H

## HA

Consulte la [alta disponibilidad](#).

## migración heterogénea de bases de datos

Migración de la base de datos de origen a una base de datos de destino que utilice un motor de base de datos diferente (por ejemplo, de Oracle a Amazon Aurora). La migración heterogénea suele ser parte de un esfuerzo de rediseño de la arquitectura y convertir el esquema puede ser una tarea compleja. [AWS ofrece AWS SCT](#), lo cual ayuda con las conversiones de esquemas.

## alta disponibilidad (HA)

La capacidad de una carga de trabajo para funcionar de forma continua, sin intervención, en caso de desafíos o desastres. Los sistemas de alta disponibilidad están diseñados para realizar una conmutación por error automática, ofrecer un rendimiento de alta calidad de forma constante y gestionar diferentes cargas y fallos con un impacto mínimo en el rendimiento.

## modernización histórica

Un enfoque utilizado para modernizar y actualizar los sistemas de tecnología operativa (TO) a fin de satisfacer mejor las necesidades de la industria manufacturera. Un histórico es un tipo de base de datos que se utiliza para recopilar y almacenar datos de diversas fuentes en una fábrica.

## datos retenidos

Parte de los datos históricos etiquetados que se ocultan de un conjunto de datos que se utiliza para entrenar un modelo de aprendizaje [automático](#). Puede utilizar los datos de reserva para evaluar el rendimiento del modelo comparando las predicciones del modelo con los datos de reserva.

## migración homogénea de bases de datos

Migración de la base de datos de origen a una base de datos de destino que comparte el mismo motor de base de datos (por ejemplo, Microsoft SQL Server a Amazon RDS para SQL Server). La migración homogénea suele formar parte de un esfuerzo para volver a alojar o redefinir la plataforma. Puede utilizar las utilidades de bases de datos nativas para migrar el esquema.

## datos recientes

Datos a los que se accede con frecuencia, como datos en tiempo real o datos traslacionales recientes. Por lo general, estos datos requieren un nivel o una clase de almacenamiento de alto rendimiento para proporcionar respuestas rápidas a las consultas.

## hotfix

Una solución urgente para un problema crítico en un entorno de producción. Debido a su urgencia, las revisiones suelen realizarse fuera del flujo de trabajo habitual de las versiones. DevOps

## periodo de hiperatención

Periodo, inmediatamente después de la transición, durante el cual un equipo de migración administra y monitorea las aplicaciones migradas en la nube para solucionar cualquier problema. Por lo general, este periodo dura de 1 a 4 días. Al final del periodo de hiperatención, el equipo de migración suele transferir la responsabilidad de las aplicaciones al equipo de operaciones en la nube.

|

**IaC**

Vea [la infraestructura como código.](#)

políticas basadas en identidades

Política asociada a uno o más directores de IAM que define sus permisos en el Nube de AWS entorno.

aplicación inactiva

Aplicación que utiliza un promedio de CPU y memoria de entre 5 y 20 por ciento durante un periodo de 90 días. En un proyecto de migración, es habitual retirar estas aplicaciones o mantenerlas en las instalaciones.

**IIoT**

Consulte [Internet de las cosas industrial.](#)

infraestructura inmutable

Un modelo que implementa una nueva infraestructura para las cargas de trabajo de producción en lugar de actualizar, parchear o modificar la infraestructura existente. [Las infraestructuras inmutables son intrínsecamente más consistentes, fiables y predecibles que las infraestructuras mutables.](#) Para obtener más información, consulte las prácticas recomendadas para [implementar con una infraestructura inmutable](#) en Well-Architected Framework AWS .

VPC entrante (de entrada)

En una arquitectura de AWS cuentas múltiples, una VPC que acepta, inspecciona y enruta las conexiones de red desde fuera de una aplicación. La [arquitectura AWS de referencia de seguridad](#) recomienda configurar la cuenta de red con entradas, salidas e inspección VPCs para proteger la interfaz bidireccional entre la aplicación y el resto de Internet.

migración gradual

Estrategia de transición en la que se migra la aplicación en partes pequeñas en lugar de realizar una transición única y completa. Por ejemplo, puede trasladar inicialmente solo unos pocos microservicios o usuarios al nuevo sistema. Tras comprobar que todo funciona correctamente, puede trasladar microservicios o usuarios adicionales de forma gradual hasta que pueda retirar su sistema heredado. Esta estrategia reduce los riesgos asociados a las grandes migraciones.

## Industria 4.0

Un término que [Klaus Schwab](#) introdujo en 2016 para referirse a la modernización de los procesos de fabricación mediante avances en la conectividad, los datos en tiempo real, la automatización, el análisis y la inteligencia artificial/aprendizaje automático.

### infraestructura

Todos los recursos y activos que se encuentran en el entorno de una aplicación.

### infraestructura como código (IaC)

Proceso de aprovisionamiento y administración de la infraestructura de una aplicación mediante un conjunto de archivos de configuración. La IaC se ha diseñado para ayudarlo a centralizar la administración de la infraestructura, estandarizar los recursos y escalar con rapidez a fin de que los entornos nuevos sean repetibles, fiables y consistentes.

### Internet de las cosas industrial (T) IIo

El uso de sensores y dispositivos conectados a Internet en los sectores industriales, como el productivo, el eléctrico, el automotriz, el sanitario, el de las ciencias de la vida y el de la agricultura. Para obtener más información, consulte [Creación de una estrategia de transformación digital de la Internet de las cosas \(IIoT\) industrial](#).

### VPC de inspección

En una arquitectura de AWS cuentas múltiples, una VPC centralizada que gestiona las inspecciones del tráfico de red VPCs entre Internet y las redes locales (en una misma o Regiones de AWS diferente). La [arquitectura AWS de referencia de seguridad](#) recomienda configurar su cuenta de red con entrada, salida e inspección VPCs para proteger la interfaz bidireccional entre la aplicación e Internet en general.

### Internet de las cosas (IoT)

Red de objetos físicos conectados con sensores o procesadores integrados que se comunican con otros dispositivos y sistemas a través de Internet o de una red de comunicación local. Para obtener más información, consulte [¿Qué es IoT?](#).

### interpretabilidad

Característica de un modelo de machine learning que describe el grado en que un ser humano puede entender cómo las predicciones del modelo dependen de sus entradas. Para obtener más información, consulte Interpretabilidad del [modelo de aprendizaje automático](#) con AWS

## IoT

Consulte [Internet de las cosas.](#)

## biblioteca de información de TI (ITIL)

Conjunto de prácticas recomendadas para ofrecer servicios de TI y alinearlos con los requisitos empresariales. La ITIL proporciona la base para la ITSM.

## administración de servicios de TI (ITSM)

Actividades asociadas con el diseño, la implementación, la administración y el soporte de los servicios de TI para una organización. Para obtener información sobre la integración de las operaciones en la nube con las herramientas de ITSM, consulte la [Guía de integración de operaciones.](#)

## ITIL

Consulte la [biblioteca de información de TI.](#)

## ITSM

Consulte [Administración de servicios de TI.](#)

## L

## control de acceso basado en etiquetas (LBAC)

Una implementación del control de acceso obligatorio (MAC) en la que a los usuarios y a los propios datos se les asigna explícitamente un valor de etiqueta de seguridad. La intersección entre la etiqueta de seguridad del usuario y la etiqueta de seguridad de los datos determina qué filas y columnas puede ver el usuario.

## zona de aterrizaje

Una landing zone es un AWS entorno multicuenta bien diseñado, escalable y seguro. Este es un punto de partida desde el cual las empresas pueden lanzar e implementar rápidamente cargas de trabajo y aplicaciones con confianza en su entorno de seguridad e infraestructura. Para obtener más información sobre las zonas de aterrizaje, consulte [Configuración de un entorno de AWS seguro y escalable con varias cuentas.](#)

## modelo de lenguaje grande (LLM)

Un modelo de [IA](#) de aprendizaje profundo que se entrena previamente con una gran cantidad de datos. Un LLM puede realizar múltiples tareas, como responder preguntas, resumir documentos, traducir textos a otros idiomas y completar oraciones. [Para obtener más información, consulte Qué son. LLMs](#)

## migración grande

Migración de 300 servidores o más.

## LBAC

Consulte control de [acceso basado en etiquetas](#).

## privilegio mínimo

La práctica recomendada de seguridad que consiste en conceder los permisos mínimos necesarios para realizar una tarea. Para obtener más información, consulte [Aplicar permisos de privilegio mínimo](#) en la documentación de IAM.

## migrar mediante lift-and-shift

Ver [7 Rs.](#)

## sistema little-endian

Un sistema que almacena primero el byte menos significativo. Véase también [endianness](#).

## LLM

Véase un modelo de lenguaje [amplio](#).

## entornos inferiores

Véase [entorno](#).

## M

## machine learning (ML)

Un tipo de inteligencia artificial que utiliza algoritmos y técnicas para el reconocimiento y el aprendizaje de patrones. El ML analiza y aprende de los datos registrados, como los datos del

Internet de las cosas (IoT), para generar un modelo estadístico basado en patrones. Para más información, consulte [Machine learning](#).

rama principal

Ver [sucursal](#).

malware

Software diseñado para comprometer la seguridad o la privacidad de la computadora. El malware puede interrumpir los sistemas informáticos, filtrar información confidencial u obtener acceso no autorizado. Algunos ejemplos de malware son los virus, los gusanos, el ransomware, los troyanos, el spyware y los keyloggers.

servicios gestionados

Servicios de AWS para los que AWS opera la capa de infraestructura, el sistema operativo y las plataformas, y usted accede a los puntos finales para almacenar y recuperar datos. Amazon Simple Storage Service (Amazon S3) y Amazon DynamoDB son ejemplos de servicios gestionados. También se conocen como servicios abstractos.

sistema de ejecución de fabricación (MES)

Un sistema de software para rastrear, monitorear, documentar y controlar los procesos de producción que convierten las materias primas en productos terminados en el taller.

MAP

Consulte [Migration Acceleration Program](#).

mecanismo

Un proceso completo en el que se crea una herramienta, se impulsa su adopción y, a continuación, se inspeccionan los resultados para realizar ajustes. Un mecanismo es un ciclo que se refuerza y mejora a sí mismo a medida que funciona. Para obtener más información, consulte [Creación de mecanismos](#) en el AWS Well-Architected Framework.

cuenta de miembro

Todas las Cuentas de AWS demás cuentas, excepto la de administración, que forman parte de una organización. AWS Organizations Una cuenta no puede pertenecer a más de una organización a la vez.

MES

Consulte el [sistema de ejecución de la fabricación](#).

## Transporte telemétrico de Message Queue Queue (MQTT)

[Un protocolo de comunicación ligero machine-to-machine \(M2M\), basado en el patrón de publicación/suscripción, para dispositivos de IoT con recursos limitados.](#)

### microservicio

Un servicio pequeño e independiente que se comunica a través de una red bien definida APIs y que, por lo general, es propiedad de equipos pequeños e independientes. Por ejemplo, un sistema de seguros puede incluir microservicios que se adapten a las capacidades empresariales, como las de ventas o marketing, o a subdominios, como las de compras, reclamaciones o análisis. Los beneficios de los microservicios incluyen la agilidad, la escalabilidad flexible, la facilidad de implementación, el código reutilizable y la resiliencia. Para obtener más información, consulte [Integrar microservicios mediante AWS servicios sin servidor.](#)

### arquitectura de microservicios

Un enfoque para crear una aplicación con componentes independientes que ejecutan cada proceso de la aplicación como un microservicio. Estos microservicios se comunican a través de una interfaz bien definida mediante un uso ligero. APIs Cada microservicio de esta arquitectura se puede actualizar, implementar y escalar para satisfacer la demanda de funciones específicas de una aplicación. Para obtener más información, consulte [Implementación de microservicios en AWS](#).

### Programa de aceleración de la migración (MAP)

Un AWS programa que proporciona soporte de consultoría, formación y servicios para ayudar a las organizaciones a crear una base operativa sólida para migrar a la nube y para ayudar a compensar el costo inicial de las migraciones. El MAP incluye una metodología de migración para ejecutar las migraciones antiguas de forma metódica y un conjunto de herramientas para automatizar y acelerar los escenarios de migración más comunes.

### migración a escala

Proceso de transferencia de la mayoría de la cartera de aplicaciones a la nube en oleadas, con más aplicaciones desplazadas a un ritmo más rápido en cada oleada. En esta fase, se utilizan las prácticas recomendadas y las lecciones aprendidas en las fases anteriores para implementar una fábrica de migración de equipos, herramientas y procesos con el fin de agilizar la migración de las cargas de trabajo mediante la automatización y la entrega ágil. Esta es la tercera fase de la [estrategia de migración de AWS](#).

## fábrica de migración

Equipos multifuncionales que agilizan la migración de las cargas de trabajo mediante enfoques automatizados y ágiles. Los equipos de las fábricas de migración suelen incluir a analistas y propietarios de operaciones, empresas, ingenieros de migración, desarrolladores y DevOps profesionales que trabajan a pasos agigantados. Entre el 20 y el 50 por ciento de la cartera de aplicaciones empresariales se compone de patrones repetidos que pueden optimizarse mediante un enfoque de fábrica. Para obtener más información, consulte la [discusión sobre las fábricas de migración](#) y la [Guía de fábricas de migración a la nube](#) en este contenido.

## metadatos de migración

Información sobre la aplicación y el servidor que se necesita para completar la migración. Cada patrón de migración requiere un conjunto diferente de metadatos de migración. Algunos ejemplos de metadatos de migración son la subred de destino, el grupo de seguridad y AWS la cuenta.

## patrón de migración

Tarea de migración repetible que detalla la estrategia de migración, el destino de la migración y la aplicación o el servicio de migración utilizados. Ejemplo: realoje la migración a Amazon EC2 con AWS Application Migration Service.

## Migration Portfolio Assessment (MPA)

Una herramienta en línea que proporciona información para validar el modelo de negocio para migrar a Nube de AWS. La MPA ofrece una evaluación detallada de la cartera (adecuación del tamaño de los servidores, precios, comparaciones del costo total de propiedad, análisis de los costos de migración), así como una planificación de la migración (análisis y recopilación de datos de aplicaciones, agrupación de aplicaciones, priorización de la migración y planificación de oleadas). La [herramienta MPA](#) (requiere iniciar sesión) está disponible de forma gratuita para todos los AWS consultores y consultores asociados de APN.

## Evaluación de la preparación para la migración (MRA)

Proceso que consiste en obtener información sobre el estado de preparación de una organización para la nube, identificar sus puntos fuertes y débiles y elaborar un plan de acción para cerrar las brechas identificadas mediante el AWS CAF. Para obtener más información, consulte la [Guía de preparación para la migración](#). La MRA es la primera fase de la [estrategia de migración de AWS](#).

## estrategia de migración

El enfoque utilizado para migrar una carga de trabajo a Nube de AWS Para obtener más información, consulte la entrada de las [7 R](#) de este glosario y consulte [Movilice a su organización para acelerar las migraciones a gran escala](#).

## ML

[Consulte el aprendizaje automático.](#)

## modernización

Transformar una aplicación obsoleta (antigua o monolítica) y su infraestructura en un sistema ágil, elástico y de alta disponibilidad en la nube para reducir los gastos, aumentar la eficiencia y aprovechar las innovaciones. Para obtener más información, consulte [Estrategia para modernizar las aplicaciones en el Nube de AWS](#).

## evaluación de la preparación para la modernización

Evaluación que ayuda a determinar la preparación para la modernización de las aplicaciones de una organización; identifica los beneficios, los riesgos y las dependencias; y determina qué tan bien la organización puede soportar el estado futuro de esas aplicaciones. El resultado de la evaluación es un esquema de la arquitectura objetivo, una hoja de ruta que detalla las fases de desarrollo y los hitos del proceso de modernización y un plan de acción para abordar las brechas identificadas. Para obtener más información, consulte [Evaluación de la preparación para la modernización de las aplicaciones en el Nube de AWS](#).

## aplicaciones monolíticas (monolitos)

Aplicaciones que se ejecutan como un único servicio con procesos estrechamente acoplados. Las aplicaciones monolíticas presentan varios inconvenientes. Si una característica de la aplicación experimenta un aumento en la demanda, se debe escalar toda la arquitectura. Agregar o mejorar las características de una aplicación monolítica también se vuelve más complejo a medida que crece la base de código. Para solucionar problemas con la aplicación, puede utilizar una arquitectura de microservicios. Para obtener más información, consulte [Descomposición de monolitos en microservicios](#).

## MAPA

Consulte [la evaluación de la cartera de migración](#).

## MQTT

Consulte [Message Queue Queue Telemetría y Transporte](#).

## clasificación multiclasé

Un proceso que ayuda a generar predicciones para varias clases (predice uno de más de dos resultados). Por ejemplo, un modelo de ML podría preguntar “¿Este producto es un libro, un automóvil o un teléfono?” o “¿Qué categoría de productos es más interesante para este cliente?”.

## infraestructura mutable

Un modelo que actualiza y modifica la infraestructura existente para las cargas de trabajo de producción. Para mejorar la coherencia, la fiabilidad y la previsibilidad, el AWS Well-Architected Framework recomienda el uso [de una infraestructura inmutable](#) como práctica recomendada.

# O

## OAC

[Consulte el control de acceso de origen.](#)

## OAI

[Consulte la identidad de acceso de origen.](#)

## OCM

[Consulte gestión del cambio organizacional.](#)

## migración fuera de línea

Método de migración en el que la carga de trabajo de origen se elimina durante el proceso de migración. Este método implica un tiempo de inactividad prolongado y, por lo general, se utiliza para cargas de trabajo pequeñas y no críticas.

## OI

[Consulte integración de operaciones.](#)

## OLA

Véase el [acuerdo a nivel operativo.](#)

## migración en línea

Método de migración en el que la carga de trabajo de origen se copia al sistema de destino sin que se desconecte. Las aplicaciones que están conectadas a la carga de trabajo pueden seguir

funcionando durante la migración. Este método implica un tiempo de inactividad nulo o mínimo y, por lo general, se utiliza para cargas de trabajo de producción críticas.

## OPC-UA

Consulte [Open Process Communications: arquitectura unificada](#).

### Comunicaciones de proceso abierto: arquitectura unificada (OPC-UA)

Un protocolo de comunicación machine-to-machine (M2M) para la automatización industrial. El OPC-UA proporciona un estándar de interoperabilidad con esquemas de cifrado, autenticación y autorización de datos.

## acuerdo de nivel operativo (OLA)

Acuerdo que aclara lo que los grupos de TI operativos se comprometen a ofrecerse entre sí, para respaldar un acuerdo de nivel de servicio (SLA).

## revisión de la preparación operativa (ORR)

Una lista de preguntas y las mejores prácticas asociadas que le ayudan a comprender, evaluar, prevenir o reducir el alcance de los incidentes y posibles fallos. Para obtener más información, consulte [Operational Readiness Reviews \(ORR\)](#) en AWS Well-Architected Framework.

## tecnología operativa (OT)

Sistemas de hardware y software que funcionan con el entorno físico para controlar las operaciones, los equipos y la infraestructura industriales. En la industria manufacturera, la integración de los sistemas de TO y tecnología de la información (TI) es un enfoque clave para las transformaciones de [la industria 4.0](#).

## integración de operaciones (OI)

Proceso de modernización de las operaciones en la nube, que implica la planificación de la preparación, la automatización y la integración. Para obtener más información, consulte la [Guía de integración de las operaciones](#).

## registro de seguimiento organizativo

Un registro creado por el AWS CloudTrail que se registran todos los eventos para todos Cuentas de AWS los miembros de una organización AWS Organizations. Este registro de seguimiento se crea en cada Cuenta de AWS que forma parte de la organización y realiza un seguimiento de la actividad en cada cuenta. Para obtener más información, consulte [Crear un registro para una organización](#) en la CloudTrail documentación.

## administración del cambio organizacional (OCM)

Marco para administrar las transformaciones empresariales importantes y disruptivas desde la perspectiva de las personas, la cultura y el liderazgo. La OCM ayuda a las empresas a prepararse para nuevos sistemas y estrategias y a realizar la transición a ellos, al acelerar la adopción de cambios, abordar los problemas de transición e impulsar cambios culturales y organizacionales. En la estrategia de AWS migración, este marco se denomina aceleración del personal, debido a la velocidad de cambio que requieren los proyectos de adopción de la nube. Para obtener más información, consulte la [Guía de OCM](#).

## control de acceso de origen (OAC)

En CloudFront, una opción mejorada para restringir el acceso y proteger el contenido del Amazon Simple Storage Service (Amazon S3). El OAC admite todos los buckets de S3 Regiones de AWS, el cifrado del lado del servidor AWS KMS (SSE-KMS) y las solicitudes dinámicas PUT y DELETE dirigidas al bucket de S3.

## identidad de acceso de origen (OAI)

En CloudFront, una opción para restringir el acceso y proteger el contenido de Amazon S3. Cuando utiliza OAI, CloudFront crea un principal con el que Amazon S3 puede autenticarse. Los directores autenticados solo pueden acceder al contenido de un bucket de S3 a través de una distribución específica. CloudFront Consulte también el [OAC](#), que proporciona un control de acceso más detallado y mejorado.

## ORR

Consulte la revisión de [la preparación operativa](#).

## OT

Consulte la [tecnología operativa](#).

## VPC saliente (de salida)

En una arquitectura de AWS cuentas múltiples, una VPC que gestiona las conexiones de red que se inician desde una aplicación. La [arquitectura AWS de referencia de seguridad](#) recomienda configurar la cuenta de red con entradas, salidas e inspección VPCs para proteger la interfaz bidireccional entre la aplicación e Internet en general.

# P

## límite de permisos

Una política de administración de IAM que se adjunta a las entidades principales de IAM para establecer los permisos máximos que puede tener el usuario o el rol. Para obtener más información, consulte [Límites de permisos](#) en la documentación de IAM.

## información de identificación personal (PII)

Información que, vista directamente o combinada con otros datos relacionados, puede utilizarse para deducir de manera razonable la identidad de una persona. Algunos ejemplos de información de identificación personal son los nombres, las direcciones y la información de contacto.

## PII

Consulte la [información de identificación personal](#).

## manual de estrategias

Conjunto de pasos predefinidos que capturan el trabajo asociado a las migraciones, como la entrega de las funciones de operaciones principales en la nube. Un manual puede adoptar la forma de scripts, manuales de procedimientos automatizados o resúmenes de los procesos o pasos necesarios para operar un entorno modernizado.

## PLC

Consulte [controlador lógico programable](#).

## PLM

Consulte la [gestión del ciclo de vida del producto](#).

## policy

Un objeto que puede definir los permisos (consulte la [política basada en la identidad](#)), especifique las condiciones de acceso (consulte la [política basada en los recursos](#)) o defina los permisos máximos para todas las cuentas de una organización AWS Organizations (consulte la política de control de [servicios](#)).

## persistencia políglota

Elegir de forma independiente la tecnología de almacenamiento de datos de un microservicio en función de los patrones de acceso a los datos y otros requisitos. Si sus microservicios tienen la misma tecnología de almacenamiento de datos, pueden enfrentarse a desafíos de

implementación o experimentar un rendimiento deficiente. Los microservicios se implementan más fácilmente y logran un mejor rendimiento y escalabilidad si utilizan el almacén de datos que mejor se adapte a sus necesidades. Para obtener más información, consulte [Habilitación de la persistencia de datos en los microservicios](#).

#### evaluación de cartera

Proceso de detección, análisis y priorización de la cartera de aplicaciones para planificar la migración. Para obtener más información, consulte la [Evaluación de la preparación para la migración](#).

#### predicate

Una condición de consulta que devuelve `true` o `false`, normalmente, se encuentra en una cláusula `WHERE`.

#### pulsar un predicado

Técnica de optimización de consultas de bases de datos que filtra los datos de la consulta antes de transferirlos. Esto reduce la cantidad de datos que se deben recuperar y procesar de la base de datos relacional y mejora el rendimiento de las consultas.

#### control preventivo

Un control de seguridad diseñado para evitar que ocurra un evento. Estos controles son la primera línea de defensa para evitar el acceso no autorizado o los cambios no deseados en la red. Para obtener más información, consulte [Controles preventivos](#) en Implementación de controles de seguridad en AWS.

#### entidad principal

Una entidad AWS que puede realizar acciones y acceder a los recursos. Esta entidad suele ser un usuario raíz para un Cuenta de AWS rol de IAM o un usuario. Para obtener más información, consulte Entidad principal en [Términos y conceptos de roles](#) en la documentación de IAM.

#### privacidad desde el diseño

Un enfoque de ingeniería de sistemas que tiene en cuenta la privacidad durante todo el proceso de desarrollo.

#### zonas alojadas privadas

Un contenedor que contiene información sobre cómo desea que Amazon Route 53 responda a las consultas de DNS de un dominio y sus subdominios dentro de uno o más VPCs. Para obtener más información, consulte [Uso de zonas alojadas privadas](#) en la documentación de Route 53.

## control proactivo

Un [control de seguridad](#) diseñado para evitar el despliegue de recursos no conformes. Estos controles escanean los recursos antes de aprovisionarlos. Si el recurso no cumple con el control, significa que no está aprovisionado. Para obtener más información, consulte la [guía de referencia de controles](#) en la AWS Control Tower documentación y consulte [Controles proactivos](#) en Implementación de controles de seguridad en AWS.

## gestión del ciclo de vida del producto (PLM)

La gestión de los datos y los procesos de un producto a lo largo de todo su ciclo de vida, desde el diseño, el desarrollo y el lanzamiento, pasando por el crecimiento y la madurez, hasta el rechazo y la retirada.

## entorno de producción

Consulte [el entorno](#).

## controlador lógico programable (PLC)

En la fabricación, una computadora adaptable y altamente confiable que monitorea las máquinas y automatiza los procesos de fabricación.

## encadenamiento rápido

Utilizar la salida de una solicitud de [LLM](#) como entrada para la siguiente solicitud para generar mejores respuestas. Esta técnica se utiliza para dividir una tarea compleja en subtareas o para refinar o ampliar de forma iterativa una respuesta preliminar. Ayuda a mejorar la precisión y la relevancia de las respuestas de un modelo y permite obtener resultados más detallados y personalizados.

## seudonimización

El proceso de reemplazar los identificadores personales de un conjunto de datos por valores de marcadores de posición. La seudonimización puede ayudar a proteger la privacidad personal. Los datos seudonimizados siguen considerándose datos personales.

## publish/subscribe (pub/sub)

Un patrón que permite las comunicaciones asíncronas entre microservicios para mejorar la escalabilidad y la capacidad de respuesta. Por ejemplo, en un [MES](#) basado en microservicios, un microservicio puede publicar mensajes de eventos en un canal al que se puedan suscribir otros microservicios. El sistema puede añadir nuevos microservicios sin cambiar el servicio de publicación.

## Q

plan de consulta

Serie de pasos, como instrucciones, que se utilizan para acceder a los datos de un sistema de base de datos relacional SQL.

regresión del plan de consulta

El optimizador de servicios de la base de datos elige un plan menos óptimo que antes de un cambio determinado en el entorno de la base de datos. Los cambios en estadísticas, restricciones, configuración del entorno, enlaces de parámetros de consultas y actualizaciones del motor de base de datos PostgreSQL pueden provocar una regresión del plan.

## R

Matriz RACI

Véase [responsable, responsable, consultado, informado \(RACI\)](#).

RAG

Consulte [Recuperación y generación aumentada](#).

ransomware

Software malicioso que se ha diseñado para bloquear el acceso a un sistema informático o a los datos hasta que se efectúe un pago.

Matriz RASCI

Véase [responsable, responsable, consultado, informado \(RACI\)](#).

RCAC

Consulte control de [acceso por filas y columnas](#).

réplica de lectura

Una copia de una base de datos que se utiliza con fines de solo lectura. Puede enrutar las consultas a la réplica de lectura para reducir la carga en la base de datos principal.

rediseñar

Ver [7 Rs.](#)

## objetivo de punto de recuperación (RPO)

La cantidad de tiempo máximo aceptable desde el último punto de recuperación de datos. Esto determina qué se considera una pérdida de datos aceptable entre el último punto de recuperación y la interrupción del servicio.

## objetivo de tiempo de recuperación (RTO)

La demora máxima aceptable entre la interrupción del servicio y el restablecimiento del servicio.

## refactorizar

Ver [7 Rs.](#)

## Region

Una colección de AWS recursos en un área geográfica. Cada uno Región de AWS está aislado e independiente de los demás para proporcionar tolerancia a las fallas, estabilidad y resiliencia. Para obtener más información, consulte [Regiones de AWS Especificar qué cuenta puede usar.](#)

## regresión

Una técnica de ML que predice un valor numérico. Por ejemplo, para resolver el problema de “¿A qué precio se venderá esta casa?”, un modelo de ML podría utilizar un modelo de regresión lineal para predecir el precio de venta de una vivienda en función de datos conocidos sobre ella (por ejemplo, los metros cuadrados).

## volver a alojar

Consulte [7 Rs.](#)

## versión

En un proceso de implementación, el acto de promover cambios en un entorno de producción.

## trasladarse

Ver [7 Rs.](#)

## redefinir la plataforma

Ver [7 Rs.](#)

## recompra

Ver [7 Rs.](#)

## resiliencia

La capacidad de una aplicación para resistir las interrupciones o recuperarse de ellas. [La alta disponibilidad y la recuperación ante desastres](#) son consideraciones comunes a la hora de planificar la resiliencia en el. Nube de AWS Para obtener más información, consulte [Nube de AWS Resiliencia](#).

## política basada en recursos

Una política asociada a un recurso, como un bucket de Amazon S3, un punto de conexión o una clave de cifrado. Este tipo de política especifica a qué entidades principales se les permite el acceso, las acciones compatibles y cualquier otra condición que deba cumplirse.

## matriz responsable, confiable, consultada e informada (RACI)

Una matriz que define las funciones y responsabilidades de todas las partes involucradas en las actividades de migración y las operaciones de la nube. El nombre de la matriz se deriva de los tipos de responsabilidad definidos en la matriz: responsable (R), contable (A), consultado (C) e informado (I). El tipo de soporte (S) es opcional. Si incluye el soporte, la matriz se denomina matriz RASCI y, si la excluye, se denomina matriz RACI.

## control receptivo

Un control de seguridad que se ha diseñado para corregir los eventos adversos o las desviaciones con respecto a su base de seguridad. Para obtener más información, consulte [Controles receptivos](#) en Implementación de controles de seguridad en AWS.

## retain

Consulte [7 Rs.](#)

## jubilarse

Ver [7 Rs.](#)

## Generación aumentada de recuperación (RAG)

Tecnología de [inteligencia artificial generativa](#) en la que un máster [hace referencia](#) a una fuente de datos autorizada que se encuentra fuera de sus fuentes de datos de formación antes de generar una respuesta. Por ejemplo, un modelo RAG podría realizar una búsqueda semántica en la base de conocimientos o en los datos personalizados de una organización. Para obtener más información, consulte [Qué es el RAG](#).

## rotación

Proceso de actualizar periódicamente un [secreto](#) para dificultar el acceso de un atacante a las credenciales.

## control de acceso por filas y columnas (RCAC)

El uso de expresiones SQL básicas y flexibles que tienen reglas de acceso definidas. El RCAC consta de permisos de fila y máscaras de columnas.

## RPO

Consulte el [objetivo del punto de recuperación](#).

## RTO

Consulte el [objetivo de tiempo de recuperación](#).

## manual de procedimientos

Conjunto de procedimientos manuales o automatizados necesarios para realizar una tarea específica. Por lo general, se diseñan para agilizar las operaciones o los procedimientos repetitivos con altas tasas de error.

# S

## SAML 2.0

Un estándar abierto que utilizan muchos proveedores de identidad (IdPs). Esta función permite el inicio de sesión único (SSO) federado, de modo que los usuarios pueden iniciar sesión Consola de administración de AWS o llamar a las operaciones de la AWS API sin tener que crear un usuario en IAM para todos los miembros de la organización. Para obtener más información sobre la federación basada en SAML 2.0, consulte [Acerca de la federación basada en SAML 2.0](#) en la documentación de IAM.

## SCADA

Consulte el [control de supervisión y la adquisición de datos](#).

## SCP

Consulte la [política de control de servicios](#).

## secreta

Información confidencial o restringida, como una contraseña o credenciales de usuario, que almacene de forma cifrada. AWS Secrets Manager Se compone del valor secreto y sus metadatos. El valor secreto puede ser binario, una sola cadena o varias cadenas. Para obtener más información, consulta [¿Qué hay en un secreto de Secrets Manager?](#) en la documentación de Secrets Manager.

## seguridad desde el diseño

Un enfoque de ingeniería de sistemas que tiene en cuenta la seguridad durante todo el proceso de desarrollo.

## control de seguridad

Barrera de protección técnica o administrativa que impide, detecta o reduce la capacidad de un agente de amenazas para aprovechar una vulnerabilidad de seguridad. Existen cuatro tipos principales de controles de seguridad: [preventivos, de detección](#), con [capacidad de respuesta](#) y [proactivos](#).

## refuerzo de la seguridad

Proceso de reducir la superficie expuesta a ataques para hacerla más resistente a los ataques. Esto puede incluir acciones, como la eliminación de los recursos que ya no se necesitan, la implementación de prácticas recomendadas de seguridad consistente en conceder privilegios mínimos o la desactivación de características innecesarias en los archivos de configuración.

## sistema de información sobre seguridad y administración de eventos (SIEM)

Herramientas y servicios que combinan sistemas de administración de información sobre seguridad (SIM) y de administración de eventos de seguridad (SEM). Un sistema de SIEM recopila, monitorea y analiza los datos de servidores, redes, dispositivos y otras fuentes para detectar amenazas y brechas de seguridad y generar alertas.

## automatización de la respuesta de seguridad

Una acción predefinida y programada que está diseñada para responder automáticamente a un evento de seguridad o remediarlo. Estas automatizaciones sirven como controles de seguridad [detectables](#) o [adaptables](#) que le ayudan a implementar las mejores prácticas AWS de seguridad. Algunos ejemplos de acciones de respuesta automatizadas incluyen la modificación de un grupo de seguridad de VPC, la aplicación de parches a una EC2 instancia de Amazon o la rotación de credenciales.

## cifrado del servidor

Cifrado de los datos en su destino, por parte de quien Servicio de AWS los recibe.  
política de control de servicio (SCP)

Política que proporciona un control centralizado de los permisos de todas las cuentas de una organización en AWS Organizations. SCPs defina barreras o establezca límites a las acciones que un administrador puede delegar en usuarios o roles. Puede utilizarlas SCPs como listas de permitidos o rechazados para especificar qué servicios o acciones están permitidos o prohibidos. Para obtener más información, consulte [las políticas de control de servicios](#) en la AWS Organizations documentación.

## punto de enlace de servicio

La URL del punto de entrada de un Servicio de AWS. Para conectarse mediante programación a un servicio de destino, puede utilizar un punto de conexión. Para obtener más información, consulte [Puntos de conexión de Servicio de AWS](#) en Referencia general de AWS.

## acuerdo de nivel de servicio (SLA)

Acuerdo que aclara lo que un equipo de TI se compromete a ofrecer a los clientes, como el tiempo de actividad y el rendimiento del servicio.

## indicador de nivel de servicio (SLI)

Medición de un aspecto del rendimiento de un servicio, como la tasa de errores, la disponibilidad o el rendimiento.

## objetivo de nivel de servicio (SLO)

Una métrica objetivo que representa el estado de un servicio, medido mediante un indicador de nivel de servicio.

## modelo de responsabilidad compartida

Un modelo que describe la responsabilidad que comparten con respecto a la seguridad y AWS el cumplimiento de la nube. AWS es responsable de la seguridad de la nube, mientras que usted es responsable de la seguridad en la nube. Para obtener más información, consulte el [Modelo de responsabilidad compartida](#).

## SIEM

Consulte [la información de seguridad y el sistema de gestión de eventos](#).

## punto único de fallo (SPOF)

Una falla en un único componente crítico de una aplicación que puede interrumpir el sistema.

### SLA

Consulte el acuerdo [de nivel de servicio](#).

### SLI

Consulte el indicador de [nivel de servicio](#).

### SLO

Consulte el objetivo de nivel de [servicio](#).

### split-and-seed modelo

Un patrón para escalar y acelerar los proyectos de modernización. A medida que se definen las nuevas funciones y los lanzamientos de los productos, el equipo principal se divide para crear nuevos equipos de productos. Esto ayuda a ampliar las capacidades y los servicios de su organización, mejora la productividad de los desarrolladores y apoya la innovación rápida. Para obtener más información, consulte [Enfoque gradual para modernizar las aplicaciones en el Nube de AWS](#).

### SPOT

Consulte el [punto único de falla](#).

### esquema en forma de estrella

Estructura organizativa de una base de datos que utiliza una tabla de datos grande para almacenar datos transaccionales o medidos y una o más tablas dimensionales más pequeñas para almacenar los atributos de los datos. Esta estructura está diseñada para usarse en un [almacén de datos](#) o con fines de inteligencia empresarial.

### patrón de higo estrangulador

Un enfoque para modernizar los sistemas monolíticos mediante la reescritura y el reemplazo gradual de las funciones del sistema hasta que se pueda desmantelar el sistema heredado. Este patrón utiliza la analogía de una higuera que crece hasta convertirse en un árbol estable y, finalmente, se apodera y reemplaza a su host. El patrón fue [presentado por Martin Fowler](#) como una forma de gestionar el riesgo al reescribir sistemas monolíticos. Para ver un ejemplo con la aplicación de este patrón, consulte [Modernización gradual de los servicios web antiguos de Microsoft ASP.NET \(ASMX\) mediante contenedores y Amazon API Gateway](#).

subred

Un intervalo de direcciones IP en la VPC. Una subred debe residir en una sola zona de disponibilidad.

supervisión, control y adquisición de datos (SCADA)

En la industria manufacturera, un sistema que utiliza hardware y software para monitorear los activos físicos y las operaciones de producción.

cifrado simétrico

Un algoritmo de cifrado que utiliza la misma clave para cifrar y descifrar los datos.

pruebas sintéticas

Probar un sistema de manera que simule las interacciones de los usuarios para detectar posibles problemas o monitorear el rendimiento. Puede usar [Amazon CloudWatch Synthetics](#) para crear estas pruebas.

indicador del sistema

Una técnica para proporcionar contexto, instrucciones o pautas a un [LLM](#) para dirigir su comportamiento. Las indicaciones del sistema ayudan a establecer el contexto y las reglas para las interacciones con los usuarios.

## T

etiquetas

Pares clave-valor que actúan como metadatos para organizar los recursos. AWS Las etiquetas pueden ayudarle a administrar, identificar, organizar, buscar y filtrar recursos. Para obtener más información, consulte [Etiquetado de los recursos de AWS](#).

variable de destino

El valor que intenta predecir en el ML supervisado. Esto también se conoce como variable de resultado. Por ejemplo, en un entorno de fabricación, la variable objetivo podría ser un defecto del producto.

lista de tareas

Herramienta que se utiliza para hacer un seguimiento del progreso mediante un manual de procedimientos. La lista de tareas contiene una descripción general del manual de

procedimientos y una lista de las tareas generales que deben completarse. Para cada tarea general, se incluye la cantidad estimada de tiempo necesario, el propietario y el progreso. entorno de prueba

[Consulte entorno.](#)

entrenamiento

Proporcionar datos de los que pueda aprender su modelo de ML. Los datos de entrenamiento deben contener la respuesta correcta. El algoritmo de aprendizaje encuentra patrones en los datos de entrenamiento que asignan los atributos de los datos de entrada al destino (la respuesta que desea predecir). Genera un modelo de ML que captura estos patrones. Luego, el modelo de ML se puede utilizar para obtener predicciones sobre datos nuevos para los que no se conoce el destino.

puerta de enlace de tránsito

Un centro de tránsito de red que puede usar para interconectar sus VPCs redes con las locales. Para obtener más información, consulte [Qué es una pasarela de tránsito](#) en la AWS Transit Gateway documentación.

flujo de trabajo basado en enlaces troncales

Un enfoque en el que los desarrolladores crean y prueban características de forma local en una rama de característica y, a continuación, combinan esos cambios en la rama principal. Luego, la rama principal se adapta a los entornos de desarrollo, preproducción y producción, de forma secuencial.

acceso de confianza

Otorgar permisos a un servicio que especifique para realizar tareas en su organización AWS Organizations y en sus cuentas en su nombre. El servicio de confianza crea un rol vinculado al servicio en cada cuenta, cuando ese rol es necesario, para realizar las tareas de administración por usted. Para obtener más información, consulte [AWS Organizations Utilización con otros AWS servicios](#) en la AWS Organizations documentación.

ajuste

Cambiar aspectos de su proceso de formación a fin de mejorar la precisión del modelo de ML. Por ejemplo, puede entrenar el modelo de ML al generar un conjunto de etiquetas, incorporar etiquetas y, luego, repetir estos pasos varias veces con diferentes ajustes para optimizar el modelo.

## equipo de dos pizzas

Un DevOps equipo pequeño al que puedes alimentar con dos pizzas. Un equipo formado por dos integrantes garantiza la mejor oportunidad posible de colaboración en el desarrollo de software.

## U

### incertidumbre

Un concepto que hace referencia a información imprecisa, incompleta o desconocida que puede socavar la fiabilidad de los modelos predictivos de ML. Hay dos tipos de incertidumbre: la incertidumbre epistémica se debe a datos limitados e incompletos, mientras que la incertidumbre aleatoria se debe al ruido y la aleatoriedad inherentes a los datos. Para más información, consulte la guía [Cuantificación de la incertidumbre en los sistemas de aprendizaje profundo](#).

### tareas indiferenciadas

También conocido como tareas arduas, es el trabajo que es necesario para crear y operar una aplicación, pero que no proporciona un valor directo al usuario final ni proporciona una ventaja competitiva. Algunos ejemplos de tareas indiferenciadas son la adquisición, el mantenimiento y la planificación de la capacidad.

### entornos superiores

Ver [entorno](#).

## V

### succión

Una operación de mantenimiento de bases de datos que implica limpiar después de las actualizaciones incrementales para recuperar espacio de almacenamiento y mejorar el rendimiento.

### control de versión

Procesos y herramientas que realizan un seguimiento de los cambios, como los cambios en el código fuente de un repositorio.

## Emparejamiento de VPC

Una conexión entre dos VPCs que le permite enrutar el tráfico mediante direcciones IP privadas. Para obtener más información, consulte [¿Qué es una interconexión de VPC?](#) en la documentación de Amazon VPC.

## vulnerabilidad

Defecto de software o hardware que pone en peligro la seguridad del sistema.

## W

### caché caliente

Un búfer caché que contiene datos actuales y relevantes a los que se accede con frecuencia. La instancia de base de datos puede leer desde la caché del búfer, lo que es más rápido que leer desde la memoria principal o el disco.

### datos templados

Datos a los que el acceso es infrecuente. Al consultar este tipo de datos, normalmente se aceptan consultas moderadamente lentas.

### función de ventana

Función SQL que realiza un cálculo en un grupo de filas que se relacionan de alguna manera con el registro actual. Las funciones de ventana son útiles para procesar tareas, como calcular una media móvil o acceder al valor de las filas en función de la posición relativa de la fila actual.

### carga de trabajo

Conjunto de recursos y código que ofrece valor comercial, como una aplicación orientada al cliente o un proceso de backend.

### flujo de trabajo

Grupos funcionales de un proyecto de migración que son responsables de un conjunto específico de tareas. Cada flujo de trabajo es independiente, pero respalda a los demás flujos de trabajo del proyecto. Por ejemplo, el flujo de trabajo de la cartera es responsable de priorizar las aplicaciones, planificar las oleadas y recopilar los metadatos de migración. El flujo de trabajo de la cartera entrega estos recursos al flujo de trabajo de migración, que luego migra los servidores y las aplicaciones.

## GUSANO

Mira, [escribe una vez, lee muchas.](#)

## WQF

Consulte el [marco AWS de calificación de la carga](#) de trabajo.

escribe una vez, lee muchas (WORM)

Un modelo de almacenamiento que escribe los datos una sola vez y evita que se eliminen o modifiquen. Los usuarios autorizados pueden leer los datos tantas veces como sea necesario, pero no pueden cambiarlos. Esta infraestructura de almacenamiento de datos se considera [inmutable](#).

## Z

ataque de día cero

Un ataque, normalmente de malware, que aprovecha una vulnerabilidad de [día cero](#).

vulnerabilidad de día cero

Un defecto o una vulnerabilidad sin mitigación en un sistema de producción. Los agentes de amenazas pueden usar este tipo de vulnerabilidad para atacar el sistema. Los desarrolladores suelen darse cuenta de la vulnerabilidad a raíz del ataque.

aviso de tiro cero

Proporcionar a un [LLM](#) instrucciones para realizar una tarea, pero sin ejemplos (imágenes) que puedan ayudar a guiarla. El LLM debe utilizar sus conocimientos previamente entrenados para realizar la tarea. La eficacia de las indicaciones cero depende de la complejidad de la tarea y de la calidad de las indicaciones. [Consulte también las indicaciones de pocos pasos.](#)

aplicación zombi

Aplicación que utiliza un promedio de CPU y memoria menor al 5 por ciento. En un proyecto de migración, es habitual retirar estas aplicaciones.

Las traducciones son generadas a través de traducción automática. En caso de conflicto entre la traducción y la versión original de inglés, prevalecerá la versión en inglés.