

eksv2

Guía Práctica de Kubernetes para Implementar una Pasarela de Pagos en EKS con Fargate Cumpliendo PCI DSS

En esta guía, te explicaré **los conceptos clave de Kubernetes** que necesitas conocer y los **pasos y comandos** específicos para implementar las configuraciones necesarias en tu pasarela de pagos en **EKS con Fargate**, cumpliendo con **PCI DSS**. Me enfocaré en las configuraciones y herramientas de Kubernetes más relevantes para tu situación.

1. Namespaces: Organización y Aislamiento

- Un **namespace** en Kubernetes es un espacio de nombres que te permite **organizar** y **aislar** recursos dentro de un clúster. Es ideal para separar las cargas de trabajo de diferentes aplicaciones o equipos.

Ejemplo:

- `kube-system`: Namespace para recursos del sistema de Kubernetes.
- `mi-app`: Namespace donde estará tu aplicación de pasarela de pagos.

Crear un namespace para tu aplicación:

```
kubectl create namespace mi-app
```

2. Pods, Deployments, y Services

- Un **pod** es la unidad más pequeña de Kubernetes. Un **deployment** es una forma de administrar y escalar pods.
- Un **service** es un recurso de Kubernetes que expone tus pods a través de una IP estable. Puede ser de tipo `ClusterIP`, `NodePort`, o `LoadBalancer`.

Crear un Deployment básico para tu aplicación:

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: payments-api
  namespace: mi-app
```

```
spec:
  replicas: 3 # Número de pods
  selector:
    matchLabels:
      app: payments-api
  template:
    metadata:
      labels:
        app: payments-api
    spec:
      containers:
        - name: payments-container
          image: your-payment-app-image:latest
          ports:
            - containerPort: 8080
```

Aplicar el deployment:

```
kubectl apply -f deployment.yaml
```

Crear un Service para exponer la aplicación a través de un LoadBalancer:

```
apiVersion: v1
kind: Service
metadata:
  name: payments-service
  namespace: mi-app
spec:
  type: LoadBalancer
  selector:
    app: payments-api
  ports:
    - protocol: TCP
      port: 443 # Puerto en el ALB (HTTPS)
      targetPort: 8080 # Puerto en los pods
```

Aplicar el Service:

```
kubectl apply -f service.yaml
```

3. ResourceQuota y LimitRange : Control de Recursos

ResourceQuota :

- Define los **recursos máximos** que un namespace puede consumir, controlando el total de CPU y memoria utilizados por los pods.

```
apiVersion: v1
kind: ResourceQuota
metadata:
  name: fargate-quota
  namespace: mi-app
spec:
  hard:
    requests.cpu: "4" # 4 CPUs en total
    requests.memory: "8Gi" # 8 GiB de memoria en total
    limits.cpu: "8" # Límite máximo de CPU en el namespace
    limits.memory: "16Gi" # Límite máximo de memoria en el namespace
```

Aplicar el ResourceQuota :

```
kubectl apply -f resource-quota.yaml
```

LimitRange :

- Establece los **límites predeterminados** y las **solicitudes mínimas** de recursos para los pods dentro de un namespace.

```
apiVersion: v1
kind: LimitRange
metadata:
  name: fargate-limits
  namespace: mi-app
spec:
  limits:
    - default:
        cpu: "500m"
        memory: "512Mi"
      defaultRequest:
        cpu: "200m"
        memory: "256Mi"
      type: Container
```

Aplicar el LimitRange :

```
kubectl apply -f limit-range.yaml
```

4. Horizontal Pod Autoscaler (HPA): Escalado Automático

El **HPA** ajusta automáticamente el número de réplicas de pods según la demanda de CPU o memoria. Esto es esencial para que tu aplicación escale dinámicamente según la carga.

Configurar HPA para escalar basado en el uso de CPU:

```
apiVersion: autoscaling/v2beta2
kind: HorizontalPodAutoscaler
metadata:
  name: payments-hpa
  namespace: mi-app
spec:
  scaleTargetRef:
    apiVersion: apps/v1
    kind: Deployment
    name: payments-api
  minReplicas: 2
  maxReplicas: 10 # Máximo número de pods
  metrics:
    - type: Resource
      resource:
        name: cpu
        target:
          type: Utilization
          averageUtilization: 50 # Escala cuando el uso de CPU supera el 50%
```

Aplicar HPA:

```
kubectl apply -f hpa.yaml
```

5. TLS y Seguridad con Application Load Balancer (ALB)

Para cumplir con PCI DSS, debes asegurar que todo el tráfico esté cifrado utilizando **TLS 1.2 o superior**. Kubernetes puede integrar un **ALB** para exponer el tráfico HTTPS a tus pods.

Crear el ALB con HTTPS:

- Usarás anotaciones en el **Service** de tipo `LoadBalancer` para configurar el ALB para soportar **TLS 1.2+** y manejar un certificado SSL.

```
apiVersion: v1
kind: Service
metadata:
  name: payments-service
  namespace: mi-app
  annotations:
    alb.ingress.kubernetes.io/scheme: internet-facing
    alb.ingress.kubernetes.io/target-type: ip # Tipo IP para Fargate
    alb.ingress.kubernetes.io/certificate-arn: arn:aws:acm:region:account-id:certificate/certificate-id # Certificado SSL
spec:
  type: LoadBalancer
  ports:
    - protocol: TCP
      port: 443 # HTTPS
      targetPort: 8080 # Puerto en los pods
  selector:
    app: payments-api
```

Aplicar el Service con ALB:

```
kubectl apply -f alb-service.yaml
```

6. Monitoreo y Auditoría: CloudTrail y VPC Flow Logs

PCI DSS requiere monitoreo continuo. Aunque **CloudTrail** y **VPC Flow Logs** se configuran fuera de Kubernetes, es importante mencionarlos porque se usan para auditar todas las actividades en tu infraestructura.

Habilitar CloudTrail (fuera de Kubernetes):

```
aws cloudtrail create-trail --name MyTrail --s3-bucket-name MyLogsBucket
aws cloudtrail start-logging --name MyTrail
```

Habilitar VPC Flow Logs (fuera de Kubernetes):

```
aws ec2 create-flow-logs --resource-type VPC --resource-id vpc-id --traffic-type ALL --log-group-name FlowLogsGroup --deliver-logs-permission-arn
```

```
arn:aws:iam::account-id:role/FlowLogsIAMRole
```

7. Web Application Firewall (WAF): Protección de la Aplicación

Para proteger tu pasarela de pagos contra amenazas comunes, como inyecciones SQL y ataques XSS, debes implementar **AWS WAF** en el ALB. Esto se hace en la consola de AWS o mediante CloudFormation.

Agregar un Web ACL (fuera de Kubernetes, pero importante para PCI DSS):

```
aws wafv2 create-web-acl --name "WebACL" --scope REGIONAL --default-action
Allow={} --visibility-config
SampledRequestsEnabled=true,CloudWatchMetricsEnabled=true,MetricName=WebACL --
rules ' [{"Name": "SQLInjectionRule", "Priority": 1, "Statement":
{"SqliMatchStatement": {"FieldToMatch": {"UriPath": {}},
"TextTransformations": [{"Priority": 0, "Type": "NONE"}]}}, {"Action":
{"Block": {}}, "VisibilityConfig": {"SampledRequestsEnabled": true,
"CloudWatchMetricsEnabled": true, "MetricName": "SQLInjectionRule"}}]' --
region <your-region>
```

8. Verificación y Monitoreo

Después de aplicar todas las configuraciones, es importante verificar que todo esté funcionando como se espera.

Verificar pods y servicios:

```
kubectl get pods -n mi-app
kubectl get services -n mi-app
```

Verificar el HPA:

```
kubectl get hpa -n mi-app
```

Verificar el tráfico de red:

- Puedes usar **CloudWatch** para revisar los logs de **VPC Flow Logs** y las métricas de **CloudTrail**.

Conclusión

Este conjunto de comandos y configuraciones te guía a través de los conceptos clave de Kubernetes, como namespaces, pods, deployments, y escalado automático, todos aplicados específicamente para una **pasarela de pagos** que necesita cumplir con los estándares de **PCI DSS**.

Los puntos principales que debes implementar incluyen:

1. **Control de recursos con ResourceQuota y LimitRange .**
2. **Cifrado de tráfico con TLS 1.2+ usando ALB.**
3. **Monitoreo y auditoría continuos** con CloudTrail y VPC Flow

Logs.

4. **Escalabilidad automática con HPA** para manejar variaciones en la demanda.

Esta guía te ayudará a implementar una infraestructura segura, escalable y conforme a las normativas.

En **Amazon EKS** con **Fargate**, no tienes control directo sobre los nodos subyacentes como en los clústeres de Kubernetes que usan instancias EC2. Sin embargo, existen configuraciones que puedes aplicar para **controlar los recursos** y establecer **límites** tanto a nivel de los pods individuales como a nivel de las configuraciones del clúster para asegurarte de que no se excedan ciertos límites de uso de Fargate.

1. Fargate Resource Quotas en EKS

Fargate, por sí mismo, no tiene un límite estricto de recursos predefinido a nivel de clúster o espacio de nombres, pero puedes configurar ciertas **cuotas** y **límites** dentro de Kubernetes (a nivel de namespaces o clúster) para limitar el uso de CPU y memoria por los pods que se ejecutan en Fargate.

Opciones para limitar los recursos en Fargate en EKS:

- **ResourceQuotas:** Puedes utilizar los objetos de **Kubernetes ResourceQuota** para establecer límites a nivel de namespace. Esto es particularmente útil para evitar que un namespace consuma más recursos de los que debería.
- **LimitRange:** Puedes usar **LimitRange** para definir los límites por defecto de CPU y memoria que cada pod o contenedor puede solicitar en un namespace.

2. Configuración de ResourceQuota para limitar los recursos de un namespace

Un **ResourceQuota** en Kubernetes te permite establecer límites en la cantidad total de CPU, memoria, y otros recursos que los pods en un **namespace** específico pueden consumir.

Ejemplo de ResourceQuota para limitar los recursos:

```
apiVersion: v1
kind: ResourceQuota
metadata:
  name: fargate-quota
  namespace: mi-app # Namespace en el que deseas aplicar la cuota
spec:
  hard:
    requests.cpu: "4" # El namespace no puede solicitar más de 4 CPU en total
    requests.memory: "8Gi" # El namespace no puede solicitar más de 8Gi de
memoria en total
    limits.cpu: "8" # El namespace no puede exceder 8 CPU como límite
    limits.memory: "16Gi" # El namespace no puede exceder 16Gi de memoria
como límite
```

Explicación:

- **requests.cpu**: La cantidad total de CPU solicitada por los pods en el namespace no puede exceder **4 vCPU**.
- **limits.cpu**: El límite máximo de CPU que los pods pueden usar en total no puede exceder **8 vCPU**.
- **requests.memory** y **limits.memory**: Funciona de manera similar para la memoria. Aquí estamos limitando el uso total de la memoria a 16 GiB.

Esto asegura que los pods que se ejecuten en **Fargate** dentro del namespace `mi-app` no consuman más recursos de los permitidos.

3. Configuración de LimitRange para Pods individuales

Mientras que **ResourceQuota** limita el uso total de recursos dentro de un namespace, **LimitRange** define los límites de recursos para **cada pod** o **contenedor** dentro de un namespace.

Ejemplo de LimitRange :

```
apiVersion: v1
kind: LimitRange
metadata:
  name: fargate-limits
```



```
namespace: mi-app # Namespace donde aplicar los límites
spec:
  limits:
    - default:
        cpu: "500m" # Límite por defecto para CPU por pod
        memory: "512Mi" # Límite por defecto para memoria por pod
      defaultRequest:
        cpu: "200m" # Solicitud por defecto para CPU por pod
        memory: "256Mi" # Solicitud por defecto para memoria por pod
      type: Container
```

Explicación:

- **default:** Establece los límites predeterminados de CPU y memoria que un contenedor puede consumir si no se especifican en la configuración del pod.
- **defaultRequest:** Establece los recursos mínimos que Kubernetes debe asignar a los contenedores si no se especifican valores de `requests` explícitamente.

Esto ayuda a evitar que los pods consuman recursos de manera descontrolada y asegura una distribución más equitativa de los recursos en el namespace.

4. Configuraciones por Defecto en Fargate

Por defecto, **Fargate** no aplica un límite rígido a nivel de clúster en cuanto a cuántos recursos (CPU y memoria) puede utilizar, pero:

- **Cada pod en Fargate** necesita declarar explícitamente sus **requests** y **limits** de recursos. Si no lo haces, Kubernetes podría intentar agotar todos los recursos disponibles.
- **Configuración de HPA (Horizontal Pod Autoscaler):** Si necesitas escalabilidad automática, es recomendable implementar **HPA** (Horizontal Pod Autoscaler), que ajusta automáticamente el número de réplicas de pods según el uso de recursos, manteniendo un control dinámico sobre el uso de recursos.

5. Cotas del Clúster EKS

A nivel del clúster EKS, no hay una forma directa de establecer límites en los recursos que pueden usar los **pods en Fargate**, pero puedes controlar indirectamente esto con los siguientes mecanismos:

- **Quota Service (AWS Service Quotas):** AWS impone ciertos límites a nivel de cuenta para EKS, como el número máximo de clústeres, nodos, y pods por clúster. Estos límites se pueden modificar a través del **AWS Service Quotas**.

Para Fargate, los **pods** no tienen restricciones estrictas a nivel de "nodos", ya que Fargate gestiona la infraestructura subyacente, pero es importante estar al tanto de los límites predeterminados de AWS para la cuenta.

6. Escalabilidad y Autoscaling en Fargate

En lugar de tener un límite rígido en los recursos a nivel de clúster, puedes controlar la **escalabilidad** y el **autoscaling** de tu aplicación con la configuración de **Horizontal Pod Autoscaler (HPA)**. Esto te permitirá ajustar automáticamente el número de réplicas de los pods en función de la demanda, sin preocuparte de gestionar la infraestructura subyacente.

Ejemplo de HPA para escalar pods en Fargate:

```
apiVersion: autoscaling/v2beta2
kind: HorizontalPodAutoscaler
metadata:
  name: payments-hpa
  namespace: mi-app
spec:
  scaleTargetRef:
    apiVersion: apps/v1
    kind: Deployment
    name: payments-api
  minReplicas: 2
  maxReplicas: 10 # Límite superior para el número de réplicas
  metrics:
    - type: Resource
      resource:
        name: cpu
        target:
          type: Utilization
          averageUtilization: 50 # Escalar cuando el uso de CPU supere el 50%
```

Resumen:

- **Fargate no tiene límites de recursos predefinidos por clúster** de forma nativa, pero puedes implementar **ResourceQuotas** y **LimitRanges** a nivel de **namespace** para controlar cuántos recursos se pueden utilizar en total o por pod.
- Usar **Horizontal Pod Autoscaler (HPA)** es la mejor práctica para asegurar que los pods se escalen automáticamente sin saturar los recursos disponibles.
- **LimitRanges** y **ResourceQuotas** permiten que Kubernetes gestione el uso de CPU y memoria dentro de los namespaces de manera eficiente, optimizando costos y manteniendo la estabilidad.

- También puedes ajustar los **Service Quotas** de AWS si es necesario ampliar el uso de recursos en toda tu cuenta.

Esta combinación de herramientas te permite tener un control granular y seguro sobre los recursos que consume tu aplicación en **EKS con Fargate**, ayudando a evitar excesos de recursos y cumpliendo con las mejores prácticas para escalar y controlar costos.

Si tu aplicación no ha salido a producción y aún no tienes métricas reales para determinar el consumo de recursos, el enfoque que debes adoptar debe tener en cuenta tanto la **escalabilidad y optimización de costos** como el **cumplimiento de PCI DSS**. Este enfoque será iterativo, lo que significa que puedes empezar con configuraciones iniciales y luego ajustar según las métricas que obtengas en pruebas y producción. Aquí te presento una aproximación general para manejar esta situación:

1. Estimaciones Iniciales de Recursos Basadas en Pruebas y Carga Simulada

Dado que no tienes datos de producción, una buena práctica es hacer **pruebas de carga** en tu ambiente de preproducción o staging para tener una idea preliminar del uso de recursos.

- **Pruebas de carga simulada:** Usa herramientas como **Apache JMeter**, **Gatling** o **k6** para simular diferentes volúmenes de tráfico y ver cómo responde tu aplicación. Estas herramientas te ayudarán a identificar cuántos recursos (CPU y memoria) necesitas por pod para manejar diferentes niveles de carga.

Pasos para realizar pruebas de carga:

1. **Configura un entorno staging** que sea lo más similar posible a producción.
2. Simula varios escenarios de carga que incluyan tráfico normal, picos de uso y condiciones límite.
3. Monitorea el uso de recursos (CPU, memoria) en los pods y ajusta las configuraciones de **requests** y **limits** según el resultado.

Esto te permitirá obtener una idea preliminar de cuántos recursos necesita tu aplicación sin estar ya en producción.

2. Configuración Inicial de Requests y Limits Basada en Suposiciones

Si no tienes datos reales de producción, puedes usar estimaciones basadas en el tipo de aplicación que estás desarrollando. A continuación, te doy una guía aproximada que puedes

usar como punto de partida, dependiendo de si tu aplicación es de tipo **web**, **procesamiento de pagos**, **API**, etc.

- **Aplicaciones web ligeras (frontend):**
 - **CPU:** 200–300m (millicores)
 - **Memoria:** 256–512Mi
- **Aplicaciones backend/APIs** (como una pasarela de pagos):
 - **CPU:** 300–500m
 - **Memoria:** 512–1024Mi

Ejemplo de configuración inicial de recursos:

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: payments-api
  namespace: mi-app
spec:
  replicas: 3
  selector:
    matchLabels:
      app: payments-api
  template:
    metadata:
      labels:
        app: payments-api
    spec:
      containers:
        - name: payments-container
          image: payments-app-image:latest
          resources:
            requests:
              memory: "512Mi" # Estimación de memoria inicial
              cpu: "300m" # Estimación de CPU inicial
            limits:
              memory: "1024Mi" # Límite de memoria
              cpu: "500m" # Límite de CPU
```

3. Cumplimiento de PCI DSS

El cumplimiento de **PCI DSS** implica que tu infraestructura debe estar diseñada para asegurar los datos de los titulares de tarjetas y cumplir con estrictos requisitos de seguridad. Aquí te dejo las áreas más relevantes para implementar en tu infraestructura basada en **Fargate** y **EKS**:

3.1 Seguridad en la Red

- **Subnets Privadas:** Asegúrate de que las subnets donde se ejecutan tus pods (especialmente los de `kube-system` y `mi-app`) sean **subnets privadas** para que los componentes internos no estén directamente expuestos a internet.

```
FargateProfileMiApp:
  Type: AWS::EKS::FargateProfile
  Properties:
    ClusterName: !Ref EKSCluster
    FargateProfileName: mi-app-profile
    PodExecutionRoleArn: !GetAtt PodExecutionRole.Arn
    Subnets:
      - !Ref PrivateSubnet1
      - !Ref PrivateSubnet2
    Selectors:
      - Namespace: mi-app
```

3.2 Cifrado de Datos en Tránsito

- **TLS 1.2 o superior:** Cumple con los requisitos de PCI DSS configurando el **ALB (Application Load Balancer)** para que utilice **TLS 1.2 o superior** en todas las comunicaciones externas.

```
ALBListenerHTTPS:
  Type: AWS::ElasticLoadBalancingV2::Listener
  Properties:
    LoadBalancerArn: !Ref ALB
    Port: 443
    Protocol: HTTPS
    SslPolicy: ELBSecurityPolicy-TLS-1-2-2017-01 # Forzar TLS 1.2
    Certificates:
      - CertificateArn: !Ref SSLCertificateArn
```

3.3 Cifrado de Datos en Reposo

- **Cifrado con AWS KMS:** Todos los datos que almacenas, ya sea en bases de datos, S3, o cualquier otro servicio, deben estar cifrados utilizando **AWS KMS** u otro mecanismo de cifrado que cumpla con PCI DSS.

3.4 Monitoreo y Registro (Logging)

- **CloudTrail y VPC Flow Logs:** Debes habilitar **AWS CloudTrail** para monitorear todas las operaciones en tu infraestructura y **VPC Flow Logs** para auditar el tráfico de red. Estos registros son críticos para cumplir con los requisitos de auditoría de PCI DSS.

```
CloudTrail:
  Type: AWS::CloudTrail::Trail
  Properties:
    IsLogging: true
    S3BucketName: !Ref TrailLogsBucket
```

3.5 Web Application Firewall (WAF)

Implementa un **AWS WAF** para proteger tu aplicación de amenazas comunes como inyecciones SQL, XSS (Cross-Site Scripting), y ataques DDoS, que son requisitos de PCI DSS.

```
WAFWebACL:
  Type: AWS::WAFv2::WebACL
  Properties:
    DefaultAction:
      Allow: {}
    Scope: REGIONAL
    Rules:
      - Name: SQLInjectionRule
        Priority: 1
        Statement:
          SqliMatchStatement:
            FieldToMatch:
              UriPath: {}
            TextTransformations:
              - Priority: 0
                Type: NONE
        Action:
          Block: {}
    VisibilityConfig:
      SampledRequestsEnabled: true
      CloudWatchMetricsEnabled: true
      MetricName: WebACL
```

4. Configuración de Autoscaling para Optimización

Dado que no tienes métricas, es recomendable utilizar el **Horizontal Pod Autoscaler (HPA)** para que tu aplicación pueda escalar de forma automática según las necesidades de recursos

en tiempo real. Aunque tus límites de recursos iniciales pueden ser conservadores, el HPA ajustará dinámicamente las réplicas según la carga.

Ejemplo de HPA:

```
apiVersion: autoscaling/v2beta2
kind: HorizontalPodAutoscaler
metadata:
  name: payments-hpa
  namespace: mi-app
spec:
  scaleTargetRef:
    apiVersion: apps/v1
    kind: Deployment
    name: payments-api
  minReplicas: 2
  maxReplicas: 10 # Establece un límite inicial en el número de réplicas
  metrics:
    - type: Resource
      resource:
        name: cpu
        target:
          type: Utilization
          averageUtilization: 50 # Escala cuando el uso de CPU supera el 50%
```

5. Pruebas y Validación de Cumplimiento PCI DSS

- **Escaneo de vulnerabilidades:** Antes de lanzar la aplicación a producción, realiza **escaneos de vulnerabilidades** en tu infraestructura para asegurar que cumple con las normativas de **PCI DSS**.
- **Penetration Testing:** Realiza pruebas de penetración para garantizar que tu infraestructura sea segura.

Resumen de los pasos:

1. **Pruebas de carga:** Ejecuta pruebas de carga en un entorno de preproducción para estimar el uso de recursos.
2. **Configuraciones iniciales:** Basándote en las pruebas de carga o el tipo de aplicación, define configuraciones iniciales conservadoras de `requests` y `limits`.
3. **Cumplimiento de PCI DSS:** Asegura el cifrado de datos en tránsito y en reposo, y utiliza CloudTrail, WAF y VPC Flow Logs para monitorear y registrar la actividad.
4. **Escalabilidad:** Implementa **Horizontal Pod Autoscaler (HPA)** para permitir que tu aplicación escale automáticamente según la carga.

5. **Monitoreo continuo:** Una vez en producción, ajusta las configuraciones de recursos según las métricas reales obtenidas.

Este enfoque equilibrado te permitirá tener un entorno seguro, conforme con PCI DSS, y escalable en base a las necesidades de tu aplicación, sin incurrir en costos innecesarios antes de tener métricas reales.

Para conectar una **pasarela de pagos** con el exterior (es decir, permitir que los clientes puedan acceder a los servicios de pago desde internet de forma segura), debes implementar una arquitectura que garantice **seguridad, disponibilidad, escalabilidad**, y que cumpla con **PCI DSS**. Existen varias opciones para exponer la pasarela de pagos, pero lo más adecuado sería utilizar un **Application Load Balancer (ALB)**, que actúe como punto de entrada seguro y balancee el tráfico entre los diferentes servicios backend.

Opción recomendada: Application Load Balancer (ALB)

El **ALB (Application Load Balancer)** es la mejor opción para conectar tu pasarela de pagos al exterior porque ofrece:

1. **Cifrado TLS/SSL:** Puedes configurar el ALB para que acepte solo tráfico HTTPS, cumpliendo así con los requisitos de **PCI DSS** que exigen que todo el tráfico que involucre datos financieros esté cifrado.
2. **Balanceo de carga:** El ALB puede distribuir el tráfico entre los pods de Kubernetes (o contenedores en Fargate), asegurando que el tráfico entrante se maneje de manera equilibrada y eficiente, y que no se sobrecargue un solo pod o nodo.
3. **Compatibilidad con múltiples microservicios:** El ALB puede manejar aplicaciones basadas en microservicios, dirigiendo las solicitudes a diferentes grupos de destino o servicios según la ruta o el contenido.
4. **Escalabilidad:** El ALB funciona en conjunto con **Auto Scaling**, lo que significa que puede manejar aumentos repentinos de tráfico ajustando automáticamente el número de instancias o pods disponibles.
5. **Integración con WAF (Web Application Firewall):** El ALB puede ser protegido con **AWS WAF**, que te permite agregar reglas de seguridad adicionales como protección contra inyecciones SQL, ataques de fuerza bruta, y más. Esto es clave para asegurar que tu aplicación esté protegida contra las amenazas comunes de aplicaciones web.

Arquitectura recomendada usando ALB:

1. **ALB con HTTPS (TLS 1.2 o superior):**
Configura el ALB para aceptar tráfico en el puerto 443 (HTTPS). Usa un certificado de **AWS Certificate Manager (ACM)** para cifrar las conexiones. Solo permite conexiones HTTPS para cumplir con los requisitos de PCI DSS.

2. ALB dirigiendo el tráfico a los pods de Kubernetes (Fargate):

El ALB balancea el tráfico hacia los pods backend que procesan las solicitudes de la pasarela de pagos. Los **Fargate Profiles** deben estar configurados para manejar estos pods.

3. WAF para protección adicional:

Utiliza **AWS WAF** junto con el ALB para bloquear ataques maliciosos y amenazas conocidas, como ataques de inyección de código SQL o denegación de servicio (DDoS).

4. Seguridad con grupos de seguridad:

Configura grupos de seguridad en el ALB para aceptar solo tráfico HTTPS y, a nivel de backend (Fargate), asegúrate de que solo se permita el tráfico interno entre el ALB y los servicios backend. El ALB no debe exponer directamente los pods al tráfico de internet.

Pasos para implementar un ALB con tu pasarela de pagos

1. Configura el ALB con HTTPS

Para configurar el ALB y habilitar HTTPS, utiliza **AWS Certificate Manager (ACM)** para generar un certificado TLS/SSL y asócialo al ALB. Aquí tienes un ejemplo de cómo hacerlo en **CloudFormation**:

```
ALBSecurityGroup:
  Type: AWS::EC2::SecurityGroup
  Properties:
    GroupDescription: "Security group for ALB"
    VpcId: !Ref VPC
    SecurityGroupIngress:
      - IpProtocol: tcp
        FromPort: 443 # HTTPS traffic
        ToPort: 443
        CidrIp: 0.0.0.0/0 # Public access for HTTPS
    SecurityGroupEgress:
      - IpProtocol: -1
        CidrIp: 0.0.0.0/0 # Allow outbound traffic
    Tags:
      - Key: Name
        Value: ALBSecurityGroup

ALB:
  Type: AWS::ElasticLoadBalancingV2::LoadBalancer
  Properties:
    Name: PaymentGatewayALB
    Scheme: internet-facing
```

```
Subnets:
  - !Ref PublicSubnet1
  - !Ref PublicSubnet2
SecurityGroups:
  - !Ref ALBSecurityGroup
Tags:
  - Key: Name
    Value: ALB-PaymentGateway
```

```
ALBListener:
  Type: AWS::ElasticLoadBalancingV2::Listener
  Properties:
    LoadBalancerArn: !Ref ALB
    Port: 443 # HTTPS port
    Protocol: HTTPS
    Certificates:
      - CertificateArn: !Ref SSLCertificateArn # ACM certificate
    DefaultActions:
      - Type: forward
        TargetGroupArn: !Ref PaymentTargetGroup
```

2. Configura el Target Group que dirige el tráfico a tus pods

El **Target Group** define hacia dónde va el tráfico balanceado por el ALB. En este caso, dirigirá el tráfico hacia los pods que gestionan la pasarela de pagos.

```
PaymentTargetGroup:
  Type: AWS::ElasticLoadBalancingV2::TargetGroup
  Properties:
    Name: PaymentTargetGroup
    Port: 8080 # Traffic to backend pods
    Protocol: HTTP
    VpcId: !Ref VPC
    TargetType: ip # Targeting Fargate pods
    HealthCheckProtocol: HTTP
    HealthCheckPath: /health
    HealthCheckIntervalSeconds: 30
    HealthCheckTimeoutSeconds: 5
    HealthyThresholdCount: 2
    UnhealthyThresholdCount: 2
    Tags:
      - Key: Name
        Value: PaymentTargetGroup
```

3. Configura los Pods en Fargate (o nodos EC2)

Asegúrate de que tus **Pods** estén correctamente configurados para procesar el tráfico que viene del ALB. A continuación, un ejemplo de cómo exponer un pod con un **Service** tipo `LoadBalancer` que se conectará al ALB:

```
apiVersion: v1
kind: Service
metadata:
  name: payments-service
  namespace: payments
  annotations:
    alb.ingress.kubernetes.io/scheme: internet-facing
    alb.ingress.kubernetes.io/target-type: ip
spec:
  type: LoadBalancer
  ports:
    - port: 8080
      targetPort: 8080
      protocol: TCP
  selector:
    app: payments-app
```

4. Habilita WAF para proteger el ALB

Configura **AWS WAF** para agregar una capa adicional de protección contra ataques comunes como inyecciones SQL, XSS, y ataques de denegación de servicio. Aquí tienes un ejemplo de regla básica de **WAF** en CloudFormation:

```
WAFWebACL:
  Type: AWS::WAFv2::WebACL
  Properties:
    DefaultAction:
      Allow: {}
    Scope: REGIONAL
    Rules:
      - Name: SQLInjectionRule
        Priority: 1
        Statement:
          SqliMatchStatement:
            FieldToMatch:
              UriPath: {}
            TextTransformations:
              - Priority: 0
```

```
    Type: NONE
  Action:
    Block: {}
  VisibilityConfig:
    SampledRequestsEnabled: true
    CloudWatchMetricsEnabled: true
    MetricName: PaymentWAFACL
```

Resumen de la Arquitectura Recomendada

1. **ALB con HTTPS:** El ALB acepta tráfico HTTPS, usando un certificado de AWS ACM, asegurando que todo el tráfico esté cifrado.
2. **ALB con WAF:** AWS WAF se integra con el ALB para proteger contra ataques comunes y garantizar la seguridad adicional.
3. **Fargate Pods/EC2 con Load Balancer:** El ALB dirige el tráfico a los pods backend (que pueden estar ejecutados en Fargate o en nodos EC2) utilizando un **Target Group**.
4. **Seguridad con Grupos de Seguridad:** Todos los componentes están protegidos por **Security Groups** bien configurados para permitir únicamente el tráfico necesario.

¿Por qué esta arquitectura es la mejor opción?

- **Cumplimiento PCI DSS:** Al utilizar HTTPS, WAF, y seguridad en la red a través de Security Groups, cumples con los requisitos de PCI DSS.
- **Escalabilidad:** El ALB y Fargate son completamente escalables. El ALB ajusta el tráfico entre los pods según la carga, y Kubernetes con Fargate se ajusta dinámicamente para manejar más solicitudes en momentos de alta demanda.
- **Protección contra ataques:** AWS WAF te permite definir reglas específicas para bloquear tráfico malicioso, asegurando la integridad de la pasarela de pagos.

Si necesitas más detalles o ajustes adicionales para optimizar tu arquitectura, no dudes en mencionarlo. ¡Estaré encantado de ayudarte!