

Caso de Estudio: "El Dilema del Servidor Único"

Escenario:

Tu empresa, "DataGenius", ha crecido rápidamente y tu equipo de IA necesita más capacidad de cómputo. Actualmente, toda la infraestructura reside en un único servidor físico muy potente. Este servidor ejecuta tres aplicaciones críticas:

1. **Un servidor de bases de datos:** Contiene todos los datos de entrenamiento y las predicciones históricas.
2. **Un servidor web:** Sirve la API que permite a los clientes enviar datos y recibir predicciones de tu modelo de IA.
3. **Un entorno de desarrollo de IA:** Es un espacio donde tu equipo de científicos de datos experimenta y entrena nuevos modelos.

Tu jefe de tecnología te ha llamado a una reunión urgente. Los científicos de datos se quejan de que sus entrenamientos, que consumen muchos recursos de CPU, están ralentizando la base de datos y haciendo que la API web responda muy lentamente. Un día, una actualización de software en el entorno de desarrollo causó un fallo que derribó todo el servidor, dejando la API y la base de datos inaccesibles durante horas.

El jefe te pregunta: "¿Cómo podemos resolver este caos y garantizar que cada equipo tenga los recursos que necesita sin afectar a los demás? Además, necesitamos una solución que nos proteja de futuros fallos de software."

Preguntas para resolver

4. **La Solución de Virtualización:** Propón una solución usando máquinas virtuales (VMs) e hipervisores. Dibuja un diagrama simple que muestre cómo la arquitectura pasaría de un único servidor a una arquitectura virtualizada.
5. **El Rol del Hipervisor:** Explica cómo el hipervisor resolvería los problemas de interferencia entre los tres equipos. ¿Qué concepto clave (visto en clase) permite que esto funcione?
6. **La Importancia del Aislamiento:** Justifica cómo esta nueva arquitectura protege a la base de datos y al servidor web de un fallo de software en el entorno de desarrollo de IA.
7. **Ventajas Adicionales:** Además de resolver los problemas inmediatos, ¿qué otras dos ventajas clave aportaría la virtualización a la empresa a largo

plazo? Piensa en conceptos como la portabilidad o la eficiencia de recursos.

Actividad: El Caso de la Empresa de IA

Objetivo: Ayudar a una empresa ficticia a elegir las instancias de EC2 correctas para su nuevo proyecto.

Escenario

"InnovA Inc." es una startup que se dedica al análisis de imágenes para detectar defectos de fabricación en una cadena de producción. Tienen tres cargas de trabajo principales:

1. **Servidor web:** Un servidor que recibe las imágenes de las cámaras y sirve la página web de la empresa. La carga de tráfico es muy baja la mayor parte del día, pero tiene picos de uso durante la mañana, cuando los supervisores inician sesión para ver los informes.
 2. **Entrenamiento del modelo:** El equipo de científicos de datos está entrenando un modelo de deep learning para mejorar la precisión de la detección de defectos. Este proceso es muy intensivo y debe completarse lo antes posible.
 3. **Inferencia en tiempo real:** Una vez que el modelo está entrenado, se despliega una API que procesa cada imagen en tiempo real y devuelve el resultado (con o sin defecto). La latencia debe ser lo más baja posible para evitar cuellos de botella en la producción.
-

Preguntas para los equipos

- **Pregunta 1:** ¿Qué familia y qué tipo de instancia (por ejemplo, t3.large) recomendarían para el **servidor web**? ¿Por qué es adecuada para los picos de tráfico? Expliquen el concepto de **burstable** y cómo se aplica aquí.
- **Pregunta 2:** ¿Qué familia de instancia usarían para el **entrenamiento del modelo**? Nombren un ejemplo de instancia de esa familia y expliquen por qué es la mejor opción para esta tarea.
- **Pregunta 3:** Para la **inferencia en tiempo real**, ¿qué familia de instancia específica de AWS usarían? ¿Cuál es la ventaja principal de esta instancia para la inferencia, en comparación con una que usarían para el entrenamiento?

- **Pregunta 4 (Desafío):** Suponiendo que la empresa tiene un gran banco de datos de imágenes que procesar cada día, ¿qué tipo de instancia de la familia **I** o **D** podría ser útil? Justifiquen la elección.
-

Caso de Estudio: "La API de IA de 'MegaMind'"

Escenario:

Tu startup, "MegaMind", ha creado una API de inteligencia artificial para analizar el sentimiento de los comentarios en redes sociales en tiempo real. La API es un éxito, y los clientes no paran de llegar. Sin embargo, te enfrentas a tres desafíos críticos:

1. **Tráfico impredecible:** El uso de la API es muy variable. Hay picos de tráfico enormes durante eventos o campañas de marketing que duran unos minutos, y luego el uso vuelve a ser muy bajo. No puedes permitir que la API se caiga o se vuelva lenta.
 2. **Presupuesto limitado:** Como startup, no puedes permitirte pagar por una capacidad masiva las 24 horas del día. Necesitas una solución que sea lo más rentable posible.
 3. **Punto de entrada fijo:** Tus clientes se conectan a tu API usando una dirección IP que no puede cambiar. Si tu instancia principal falla, la dirección de tu API no puede desaparecer.
-

Preguntas para resolver

- **1. El Problema del Punto Fijo:** ¿Qué servicio de AWS usarías para garantizar que la dirección IP de tu API sea siempre la misma, sin importar qué instancia esté funcionando detrás? ¿Qué ventaja clave ofrece este servicio sobre una IP pública normal?
- **2. La Solución de Escalabilidad:** ¿Qué dos servicios trabajarían juntos para manejar automáticamente los picos de tráfico sin que la API se sature? Explica brevemente el rol de cada uno.
- **3. La Estrategia de Coste:** Para mantener los costes al mínimo, ¿qué modelo de precios de EC2 usarías para las instancias que ejecutan tu API? ¿Por qué esta opción es la más adecuada, a pesar de sus riesgos? ¿Cómo mitigarías esos riesgos?

- **4. El Desafío Final:** Dibuja un diagrama simple que muestre cómo estos tres servicios interactúan entre sí para crear una solución de alta disponibilidad, escalable y rentable. Nombra cada componente.

SOLUCIONES

Caso 1: "El Dilema del Servidor Único" (DataGenius)

1. La Solución de Virtualización

Se propone la implementación de un **Hipervisor de Tipo 1 (Bare Metal)** instalado directamente sobre el hardware del servidor actual. Sobre este hipervisor, se crearán tres **Máquinas Virtuales (VMs)** independientes:

- **VM 1:** Servidor de Bases de Datos.
- **VM 2:** Servidor Web / API.
- **VM 3:** Entorno de Desarrollo de IA.

Diagrama de Arquitectura:

Plaintext

[Hardware del Servidor Físico Potente]

|
[Hipervisor]

/ | \
| | |

[VM 1] [VM 2] [VM 3]

(DB) (Web API) (Entorno IA)

2. El Rol del Hipervisor y Concepto Clave

El hipervisor actúa como un administrador de recursos que distribuye la CPU, memoria y almacenamiento de forma aislada a cada VM.

- **Resolución de interferencias:** Permite asignar cuotas fijas de recursos. Así, aunque el entrenamiento de IA sea intensivo, el hipervisor garantiza que la base de datos y la web siempre tengan su porción de CPU reservada.

- **Concepto clave: Abstracción de Recursos** (o Aislamiento de Recursos), que permite que el software se desvincule del hardware físico y funcione de forma independiente.

3. La Importancia del Aislamiento

El aislamiento garantiza que los fallos de software queden confinados dentro de la VM donde ocurren. Si el entorno de desarrollo sufre un fallo crítico por una actualización, solo la **VM 3** se verá afectada. Las **VM 1 y 2** continuarán operando con normalidad, ya que sus sistemas operativos y dependencias son totalmente independientes, protegiendo así la disponibilidad de la API y la base de datos.

4. Ventajas Adicionales a Largo Plazo

- **Eficiencia de Recursos:** Permite maximizar el uso del servidor físico al ejecutar múltiples cargas de trabajo que de otro modo requerirían hardware separado.
- **Portabilidad:** Las VMs pueden ser respaldadas o movidas a otro servidor físico fácilmente, facilitando la recuperación ante desastres o migraciones futuras.

Caso 2: El Caso de la Empresa de IA (InnovA Inc.)

Pregunta 1: Servidor Web

- **Recomendación:** Familia **T** (ej. **t3.large**).
- **Justificación:** Es ideal porque el tráfico es bajo con picos puntuales. Estas instancias son **"burstable"** (con capacidad de ráfaga): acumulan "créditos de CPU" cuando la carga es baja y los consumen para aumentar el rendimiento por encima del nivel base durante los picos de la mañana.

Pregunta 2: Entrenamiento del Modelo

- **Recomendación:** Familia **P** (Instancias de Computación Acelerada) o **G**.
- **Ejemplo:** p3.2xlarge.
- **Justificación:** El entrenamiento de Deep Learning es extremadamente intensivo en cómputo paralelo; estas instancias cuentan con **GPUs de alto rendimiento** que procesan estos modelos mucho más rápido que una CPU estándar.

Pregunta 3: Inferencia en Tiempo Real

- **Recomendación:** Familia **Inf** (AWS Inferentia).
- **Ventaja principal:** A diferencia de las instancias de entrenamiento (familia P) que están diseñadas para aprendizaje masivo, las instancias Inf están optimizadas específicamente para **baja latencia y alto rendimiento de inferencia** a un coste mucho menor, lo que evita cuellos de botella en la producción.

Pregunta 4: Procesamiento de Grandes Datos (I o D)

- **Recomendación:** Familia **I** (Optimización de Almacenamiento, ej. i3.large).
- **Justificación:** Dado que se procesa un "gran banco de datos de imágenes" diariamente, se requiere un alto rendimiento de **IOPS (operaciones de entrada/salida por segundo)** y baja latencia de lectura/escritura en disco local, algo en lo que la familia I destaca gracias a sus discos NVMe.

Caso 3: "La API de IA de MegaMind"

1. El Problema del Punto Fijo

- **Servicio:** Elastic IP (EIP).
- **Ventaja:** A diferencia de una IP pública normal que cambia si la instancia se detiene o falla, la Elastic IP es una dirección estática que puedes reasignar rápidamente a otra instancia en caso de fallo, manteniendo el punto de entrada invisible para el cliente.

2. La Solución de Escalabilidad

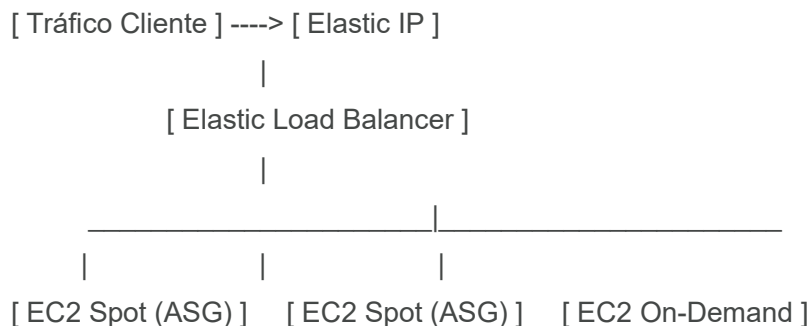
- **Servicios:** **Elastic Load Balancer (ELB)** y **Auto Scaling Group (ASG)**.
- **Roles:** El **ASG** añade o elimina instancias automáticamente según la demanda (picos de tráfico), mientras que el **ELB** distribuye de forma equitativa el tráfico entrante entre todas las instancias activas para evitar saturaciones.

3. La Estrategia de Coste

- **Modelo:** **Instancias Spot**.
- **Justificación:** Ofrecen descuentos de hasta el 90% sobre el precio bajo demanda, lo cual es vital para el presupuesto de una startup.
- **Mitigación de riesgos:** El riesgo es que AWS puede retirar la instancia con poco aviso. Se mitiga configurando el **Auto Scaling** para que combine instancias Spot con algunas **Instancias On-Demand** (bajo demanda) para asegurar una disponibilidad mínima siempre.

4. Diagrama Final de Alta Disponibilidad

Plaintext



(Componentes: Elastic IP como punto fijo, ELB como distribuidor, y ASG gestionando una mezcla de instancias para coste y disponibilidad).