

# Examen Febrero 2022

## Ingeniería de Servidores

Ismael Sallami Moreno

Universidad de Granada  
Doble Grado en Ingeniería Informática y ADE  
<https://elblogdeismael.github.io>

### 2.- (0.75 puntos)

Nuestro computador, que usa la GPU *Nvidia RTX 3090*, tarda 50 segundos en ejecutar un programa monohilo que usa solamente CPU y GPU. Tras vender nuestro otro riñón para cambiar la GPU por la nueva *RTX 3090 Ti* que acaban de anunciar, ahora el programa tarda la mitad del tiempo. Sabiendo que ahora la nueva GPU se usa la mitad del tiempo tras la mejora:

- ¿Cuántas veces es mejor esta GPU respecto a la original? (0,5 puntos)
- ¿Cuál es la ganancia máxima que podríamos conseguir optimizando solo la GPU? (0,25 puntos)

Vamos a resolver el ejercicio utilizando la ley de Amdahl.

- Calculamos cuántas veces es mejor la nueva GPU respecto a la original:

- Tiempo original:  $T_o = 50\text{ s}$
- Tiempo mejorado:  $T_m = 25\text{ s}$
- Tras la mejora, la GPU se usa la mitad del tiempo:  $T_{GPU,m} = T_m/2 = 12,5\text{ s}$

Aplicando la ley de Amdahl:

$$T_o = T_{CPU,o} + T_{GPU,o}$$

$$T_m = T_{CPU,o} + \frac{T_{GPU,o}}{k}$$

Sabemos que tras la mejora:

$$T_{GPU,m} = \frac{T_{GPU,o}}{k} = 12,5\text{ s} \implies T_{GPU,o} = 12,5 \cdot k$$

Sustituyendo:

$$T_o = T_{CPU,o} + 12,5 \cdot k$$

$$T_m = T_{CPU,o} + 12,5$$

De la segunda ecuación:

$$T_{CPU,o} = 25 - 12,5 = 12,5$$

Sustituyendo en la primera:

$$50 = 12,5 + 12,5 \cdot k \implies 12,5 \cdot k = 37,5 \implies k = 3$$

Por tanto, la nueva GPU es **3 veces más rápida** que la original.

b) Calculamos la ganancia máxima teórica optimizando solo la GPU:

La fracción del tiempo dedicada originalmente a la GPU es:

$$f = \frac{T_{GPU,o}}{T_o} = \frac{37,5}{50} = 0,75$$

Según la ley de Amdahl, la ganancia máxima es:

$$S_{\text{máx}} = \frac{1}{1-f} = \frac{1}{0,25} = 4$$

Por tanto, la ganancia máxima teórica mejorando solo la GPU es **4 veces**.

### 3.- (0.75 puntos)

En un servidor con S.O. Linux se tiene instalado el monitor de actividad **sar**. Se sabe que cada activación del monitor implica la ejecución de un total de 1500 instrucciones máquina y almacena un total de 2048 bytes de información en el fichero **/var/log/sa/saDD** del día DD correspondiente. Si el procesador del equipo tiene una velocidad media de ejecución de 75 MIPS:

- ¿Qué valor debe tener el periodo de muestreo (en milisegundos) si se quiere una sobrecarga (overhead) del 1 %? (0,5 puntos)
- Suponiendo ahora que el monitor se activa una vez cada 15 minutos, ¿cuál será el tamaño (en KiB) máximo de cada fichero del directorio **/var/log/sa**? (0,25 puntos)

a) De la definición de MIPS sabemos que:

$$T_{\text{ejec}} = \frac{NI}{MIPS \times 10^6} = \frac{1500}{75 \times 10^6} = 0,00002 \text{ segundos} = 0,02 \text{ milisegundos}$$

De la definición de overhead:

$$1 = \frac{0,00002}{\text{Período de muestreo}} \times 100 \implies \text{Período de muestreo} = 0,00002 \times 100 = 0,002 \text{ s} = 2 \text{ ms}$$

b) Si se activa una vez cada 15 minutos, el número de activaciones por día es:

$$\text{Activaciones por día} = \frac{24 \text{ horas} \times 60 \text{ minutos}}{15 \text{ minutos}} = 96$$

El tamaño total de los ficheros generados en un día es:

$$\text{Tamaño total} = 96 \times 2048 \text{ bytes} = 196608 \text{ bytes}$$

Convertimos a KiB:

$$\text{Tamaño total en KiB} = \frac{196608 \text{ bytes}}{1024} = 192 \text{ KiB}$$

Por tanto, el tamaño máximo de cada fichero del directorio `/var/log/sa` es **192 KiB**.

#### 4.- (2 puntos)

Durante las últimas 24 horas, se ha monitorizado un servidor de base de datos con el fin de obtener un modelo del mismo basado en redes de colas. Como resultado de dicha monitorización, se han obtenido las siguientes medidas (solo hay dos dispositivos en nuestro modelo: CPU y disco duro):

- Se han contabilizado un total de 54000 consultas entrantes al servidor.
- La utilización media del disco duro es del 60 % y la de la CPU es del 70 %.
- Cada consulta al servidor requiere una media de 5 accesos al disco duro.

Nota: indique claramente las definiciones y/o leyes operacionales que ha necesitado utilizar y, en su caso, si se cumplen las condiciones para que pueda usar las leyes operacionales que use.

- a) ¿Está el servidor en equilibrio de flujo? Razone la respuesta. (0,4 puntos)
- b) Calcule la productividad media del disco duro. (0,4 puntos)
- c) Calcule cuánto tiempo, de media, le dedica el disco duro a cada petición de lectura/escritura que le llega. Expresé el resultado en ms. (0,4 puntos)
- d) Calcule la productividad media máxima del servidor. (0,4 puntos)
- e) En las mejores condiciones de carga, calcule cuánto tiempo, de media, tardaría el servidor en responder a una consulta. (0,4 puntos)

- a) **Ley aplicada: definición de equilibrio de flujo y ley de flujo forzado.** El servidor está en equilibrio de flujo si el número de trabajos completados por el servidor coincide con el número de trabajos solicitados, es decir,  $C_0 \approx A_0$ , y también si  $X_0 \approx \lambda_0$ . Además, no hay ningún dispositivo saturado ( $U_i < 1$ ), por lo tanto se pueden aplicar las leyes operacionales, y asumimos que estamos en equilibrio de flujo. Alternativamente, se puede comprobar que  $\lambda_0 < X_0^{\max}$ , como se hará en el apartado (d).

b) **Ley aplicada: Ley del flujo forzado.**

$$X_i = V_i \cdot X_0$$

Donde:

- $X_0 = \lambda_0 = \frac{A_0}{T} = \frac{54000}{24 \times 3600} = 0,625$  consultas/s (Ley de Little)
- $V_{dd} = 5$  (visitas al disco por consulta)

Por tanto:

$$X_{dd} = V_{dd} \cdot X_0 = 5 \cdot 0,625 = \boxed{3,125 \text{ tr/s}}$$

c) **Ley aplicada: Ley de la utilización.** Esta ley afirma que en equilibrio de flujo se cumple:

$$U_i = X_i \cdot S_i \Rightarrow S_i = \frac{U_i}{X_i}$$

Aplicándola al disco duro:

$$S_{dd} = \frac{U_{dd}}{X_{dd}} = \frac{0,6}{3,125} = 0,192 \text{ s} = \boxed{192 \text{ ms}}$$

d) **Ley aplicada: Límite superior de productividad, definida por el cuello de botella.** En equilibrio de flujo, el dispositivo más lento limita el rendimiento, es decir:

$$X_0^{\max} = \frac{1}{D_b} = \frac{1}{D_{\text{CPU}}}$$

Donde:

$$D_{\text{CPU}} = \frac{U_{\text{CPU}}}{X_0} = \frac{0,7}{0,625} = 1,12 \text{ s}$$

Por tanto:

$$X_0^{\max} = \frac{1}{1,12} = \boxed{0,893 \text{ consultas/s}}$$

e) **Ley aplicada: Ley de Little (a nivel de servidor).** En las mejores condiciones de carga, la tasa de llegada es igual a la máxima productividad, por tanto:

$$R_0^{\min} = \frac{N_0}{X_0^{\max}}$$

Como no se proporciona  $N_0$ , se puede aplicar directamente la relación:

$$R_0^{\min} = \sum D_i = D_{\text{CPU}} + D_{\text{DD}} = \frac{0,7}{0,625} + \frac{0,6}{3,125} = 1,12 + 0,192 = \boxed{1,312 \text{ segundos}}$$

(ya que  $D_i = \frac{U_i}{X_0}$  por definición de demanda de servicio).

5.- (1 punto)

Demuestre que, para alta carga, el tiempo medio de respuesta de un servidor modelado mediante una red de colas cerrada interactiva tiende asintóticamente a la recta  $\text{Db} \cdot \text{NT} \cdot \text{Z}$ . Indique también el nombre (si lo tiene) y el significado de NT, Z y Db. Nota: indique claramente las definiciones y/o leyes operacionales que ha necesitado utilizar y, en su caso, si se cumplen las condiciones para que pueda usar las leyes operacionales que use.

Queremos demostrar que, para carga alta ( $N_T \rightarrow \infty$ ), el tiempo medio de respuesta del sistema cerrado tiende a:

$$R_0(N_T) \rightarrow D_b \cdot N_T - Z$$

**Definiciones y significado de los términos:**

- $N_T$ : Número total de trabajos (clientes) en la red cerrada.
- $Z$ : Tiempo de reflexión. Es el tiempo que pasa un cliente fuera del sistema antes de lanzar una nueva petición.
- $D_b$ : Demanda de servicio en el cuello de botella, es decir, el tiempo que el recurso más lento (bottleneck) dedica a cada trabajo.

**Ley operacional utilizada: Ley de Little para redes cerradas:**

$$N_T = X_0 \cdot (R_0 + Z) \Rightarrow R_0 = \frac{N_T}{X_0} - Z$$

**Comportamiento en alta carga:**

En condiciones de alta carga ( $N_T \gg N_T^*$ ), el sistema entra en saturación. El throughput  $X_0$  se estabiliza y no puede superar el máximo definido por el cuello de botella:

$$X_0^{\text{máx}} = \frac{1}{D_b}$$

Sustituyendo en la ecuación de  $R_0$ :

$$R_0(N_T) = \frac{N_T}{X_0^{\text{máx}}} - Z = N_T \cdot D_b - Z$$

**Por tanto, hemos demostrado que:**

$$\lim_{N_T \rightarrow \infty} R_0(N_T) = D_b \cdot N_T - Z$$

Esta recta representa la **asíntota superior** del tiempo medio de respuesta en un sistema cerrado. Es una consecuencia directa del modelo de red cerrada interactiva, del cuello de botella y de la ley de Little. Un *asíntota* es una línea que describe el comportamiento de una función a medida que se aproxima a un valor límite, en este caso, cuando el número de trabajos tiende a infinito.

**6.- (0.75 puntos)**

Responda brevemente a las siguientes cuestiones sobre el benchmark CPU 2017 que ha desarrollado el consorcio SPEC:

- a) ¿Qué componentes del sistema informático evalúa? (0,25 puntos)
- b) Indique cómo se calcula el índice CPU2017IntegerSpeed\_peak (tanto de palabra como poniendo la fórmula). El método de cálculo empleado, ¿satisface todas las exigencias de un buen índice de prestaciones? Razone la respuesta. (0,5 puntos)

- a) El benchmark CPU 2017 desarrollado por el consorcio SPEC evalúa el rendimiento del subsistema procesador y memoria principal. Incluye dos tipos de pruebas: enteros (**Integer**) y en coma flotante (**Floating Point**), cada una en dos modalidades: **speed** (tiempo de ejecución de tareas individuales) y **rate** (rendimiento con múltiples tareas concurrentes).
- b) El índice **CPU2017IntegerSpeed\_peak** se calcula como la media geométrica de los ratios de rendimiento obtenidos en cada uno de los benchmarks individuales, usando las mejores optimizaciones posibles permitidas por SPEC:

$$\text{Índice} = \left( \prod_{i=1}^n \frac{\text{Tiempo de referencia}_i}{\text{Tiempo de ejecución}_i} \right)^{1/n}$$

Este método **satisface parcialmente** los criterios de un buen índice de prestaciones:

- Cumple con la necesidad de **agregación representativa**, al usar media geométrica.
- Es **reproducible** (los resultados pueden ser verificados por terceros bajo las mismas condiciones) y está estandarizado por SPEC (especifica reglas, condiciones de ejecución y conjuntos de pruebas uniformes definidos por SPEC).
- Sin embargo, no refleja bien escenarios de uso mixtos o cargas específicas de usuario, por lo que **no garantiza representatividad general en todos los contextos**.

#### 7.- (0.75 puntos)

Cuestiones (0,25 puntos cada una).

- a) ¿A qué nos referimos por “LANE” cuando hablamos de la interfaz PCIe?
- b) ¿Para qué sirve el **system panel** en una placa base?
- c) Indique las principales secciones de un pliego de prescripciones técnicas para licitar un contrato relacionado con una instalación de servidores junto con una frase explicativa del tipo de información que debe contener cada una de dichas secciones.
- a) En la interfaz PCIe, un **LANE** es una vía de comunicación que consta de dos pares diferenciales: uno para transmisión (TX) y otro para recepción (RX). Cada LANE permite transmitir y recibir datos simultáneamente. Las interfaces PCIe pueden agrupar múltiples LANEs (por ejemplo, x1, x4, x8, x16) para aumentar el ancho de banda.
- b) El **system panel** en una placa base es un conjunto de pines que permite conectar los botones (encendido, reinicio) y los LEDs (actividad del disco, encendido) del chasis del ordenador a la placa base. Facilita el control físico del sistema.
- c) Un pliego de prescripciones técnicas para licitar un contrato relacionado con la instalación de servidores debe contener las siguientes secciones:

- **Objeto del contrato:** Describe el propósito general de la instalación (por ejemplo, despliegue de un nuevo CPD (Centro de Procesamiento de Datos)).
- **Requisitos técnicos:** Especifica las características mínimas exigidas para el hardware, software y entorno físico.
- **Criterios de aceptación:** Define las condiciones bajo las cuales la instalación se considerará satisfactoria (tests, certificaciones, validación).
- **Condiciones de entrega y puesta en marcha:** Detalla los plazos, el lugar de instalación, y responsabilidades en la fase de despliegue.