



---

## LIBRO BLANCO

# Big Data, mejor en bare metal

Convierta el rendimiento de big data en una prioridad

## RESUMEN EJECUTIVO

Las empresas actuales crean y capturan cantidades de datos sin precedentes, procedentes de múltiples fuentes y en formatos estructurados y no estructurados. Almacenar, procesar y extraer valor de este “big data” no es tarea fácil. Los profesionales de TI suelen aprovisionar servidores de cloud público para escalar el almacenamiento y potencia de procesamiento necesarios para este continuo flujo de datos, pero los recursos virtualizados no pueden ofrecer el rendimiento y consistencia de los servidores bare metal equivalentes.

IBM Cloud comprobó el rendimiento y consistencia de cargas de trabajo de big data en servidores virtuales y servidores bare metal dedicados para comparar la idoneidad de estas plataformas para aplicaciones que almacenan y procesan muy grandes cantidades de datos. Con estos resultados, los profesionales de TI pueden tomar mejores decisiones a la hora de elegir recursos cloud para cargas de trabajo con uso intensivo de almacenamiento y procesador.

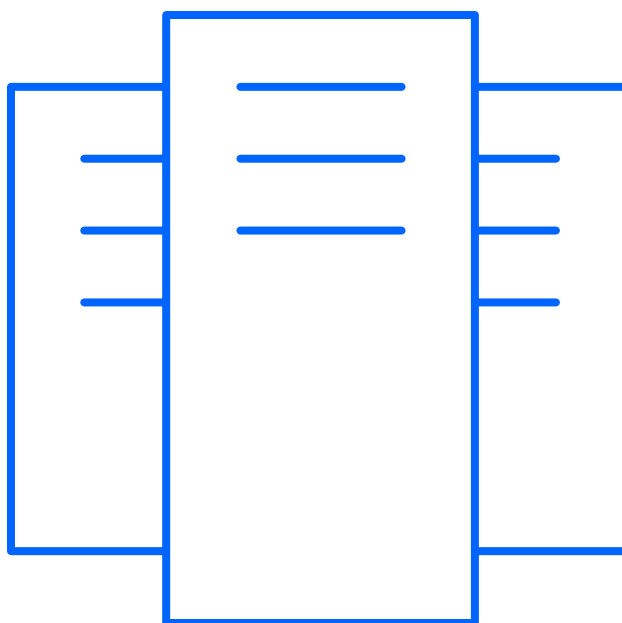


# ¿Qué es el Big Data?

A medida que evolucionan las tecnologías de almacenamiento y la capacidad se convierte en un recurso más accesible, las empresas encuentran nuevas formas de capturar y procesar más información. Estos datos proporcionan información con gran valor potencial de negocio. El reto radica en organizar y analizar los datos para crear nuevas estrategias de negocio y tomar decisiones en la organización.

Hasta hace poco, las herramientas más habituales para organizar y analizar datos eran los sistemas de gestión de bases de datos relacionales (RDBMS) mediante lenguaje de consulta estructurada (SQL). Las soluciones SQL utilizan conjuntos de datos estructurados, generalmente almacenados y manipulados en un mismo servidor. Al aumentar el tamaño del conjunto de datos hasta el techo de capacidad del servidor existente, la solución se amplía verticalmente pasando a un servidor mayor con más potencia de procesamiento, más almacenamiento y RAM. Este escalado puede precisar mucho tiempo y suponer un importante aumento de costes.

En una situación en la que los datos llegan muy rápidamente y desde muy distintas fuentes e innumerables esquemas, los administradores de bases de datos han de maximizar la eficiencia y escalabilidad de sus soluciones. En consecuencia, muchos han comenzado a utilizar bases de datos NoSQL (No Solo SQL), que utilizan conjuntos de datos no relacionales y no estructurados. Esta arquitectura de “big data” permite almacenar datos en múltiples sistemas, de modo que las aplicaciones NoSQL pueden escalar horizontalmente añadiendo gradualmente sistemas comerciales (systems commodity), logrando un aumento de la capacidad bajo demanda y mayor eficacia de costes.



## Estas arquitecturas de big data pueden interpretar muy grandes volúmenes de datos, pero para ello los datos requieren una importante infraestructura:

- Almacenamiento para admitir este volumen de datos
- RAM para mover y cargar los datos a medida que se precise
- Potencia de procesamiento conforme con el nivel de rendimiento requerido por la solución
- Una red capaz de conectar almacenes de datos distribuidos con baja latencia para mejorar el rendimiento

Para abordar estos requisitos, muchas empresas utilizan recursos de computación cloud como infraestructura subyacente para escalar horizontalmente sus entornos de big data. Las piezas de construcción más habituales en estos entornos son los servidores de cloud público virtualizados y los servidores bare metal.

# Las **cuatro V** del Big Data

**Volumen:** Piense a escala de petabytes. Desde historial web hasta registros públicos y documentos privados internos, las empresas lo almacenan todo.

**Variedad:** Grandes volúmenes de datos estructurados y no estructurados, como correo electrónico, redes sociales, vídeo, imágenes, datos atmosféricos, blogs y un largo etcétera.

**Velocidad:** Los datos se generan constantemente con consultas en tiempo real para obtener información significativa bajo demanda.

**Valor:** Información útil significativa derivada de big data que va más allá de los resultados de las consultas e informes tradicionales. Esta información puede ser transformada en analítica predictiva para revelar tendencias y patrones.

# Servidores bare metal frente a servidores virtuales

Considere los servidores bare metal y los servidores virtuales como dos herramientas de la misma caja. No es que una sea inherentemente mejor que la otra; cada una tiene sus puntos fuertes y débiles.

Los servidores bare metal proporcionan a los clientes acceso directo y exclusivo a los recursos de hardware brutos de un servidor. Los servidores virtuales son instancias cloud independientes aprovisionadas por un hipervisor en un nodo de hardware que puede ser público (compartido) o privado.

## Servidores bare metal - potencia bruta

Para cargas de trabajo con uso intensivo de procesador y E/S de disco, son ideales los servidores bare metal (también conocidos como servidores dedicados). Estos servidores son de tenencia única, y están totalmente dedicados a un solo cliente. Esto significa que no hay vecinos ruidosos que obstaculicen el rendimiento.

Asimismo, como los servidores bare metal no se ejecutan sobre un hipervisor, las cargas de trabajo no pagan el “impuesto de supervisor”, la ligera degradación del rendimiento que provoca el hipervisor al actuar como intermediario entre el sistema operativo y el hardware.

Sin hipervisor para abstraer el hardware, el aprovisionamiento y configuración de los servidores bare metal suelen precisar más tiempo que los servidores virtuales. Cuando se precisa escalar una infraestructura rápidamente, suelen evitarse los sistemas bare metal. Para hacer frente a este punto débil, IBM Cloud diseñó sistemas automatizados de despliegue y control de servidores bare metal para completar una selección de configuraciones en línea en solo 20-30 minutos y servidores totalmente personalizados (con su selección de procesador, cores, RAM, almacenamiento, puertos, etc.) en 2-4 horas.

## Servidores virtuales – flexibilidad y escalabilidad

Las aplicaciones y cargas de trabajo con importantes variaciones de tamaño o que precisan mantenerse ágiles en un mercado en constante cambio son ideales para los servidores virtuales. Los servidores virtuales se aprovisionan sobre un hipervisor en un entorno de cloud público con uno o varios tenedores. Los recursos de los servidores virtuales pueden desplegarse en cuestión de cinco minutos con plazos en meses u horas, lo que permite escalar horizontalmente añadiendo muy rápidamente servidores adicionales.

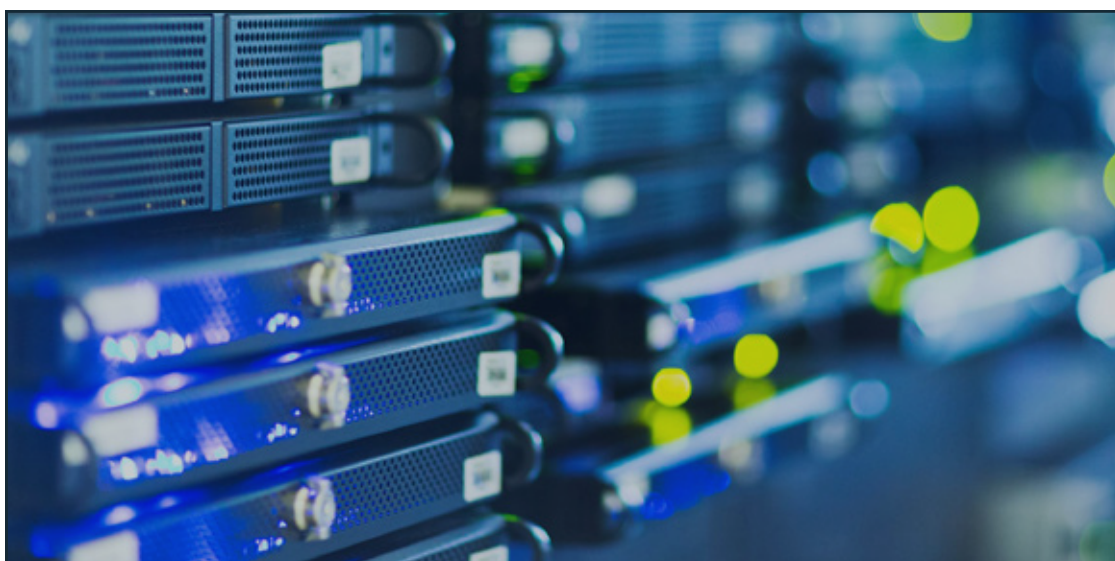
## Servidores bare metal y virtuales en conjunción

IBM Cloud aprovisiona servidores bare metal y virtuales en un mismo entorno de cloud unificado para proporcionar al cliente capacidad de elección y control de los recursos que impulsarán sus muy distintas cargas de trabajo.

**Las necesidades de la empresa cambian. Nuestros productos están diseñados para que usted pueda centrarse en sus necesidades actuales sin preocuparse de cómo serán estas necesidades al cabo de unos días, semanas o meses.**

**La plataforma e infraestructura IBM Cloud se diseñaron y construyeron buscando la máxima escalabilidad:**

- Añada servidores bare metal y virtuales bajo demanda
- Disminuya la escala cuando lo necesite para reducir costes
- Pedidos para plazos en horas o meses a la medida de cada proyecto
- Sin contratos a largo plazo



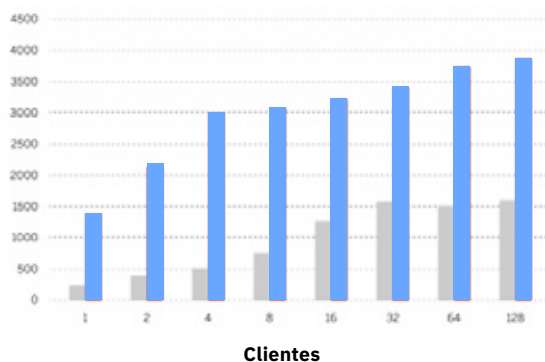
# Big Data - Rendimiento

Para determinar si las aplicaciones de big data son más adecuadas para servidores bare metal o virtuales, preparamos una serie de pruebas para medir el rendimiento y consistencia de ambas plataformas. Para medir el rendimiento con precisión, un ingeniero de IBM Cloud configuró entornos de prueba equivalentes bare metal y virtuales para consultar y actualizar un conjunto de datos de MongoDB utilizando su herramienta de comparativa (información detallada en el Apéndice A).

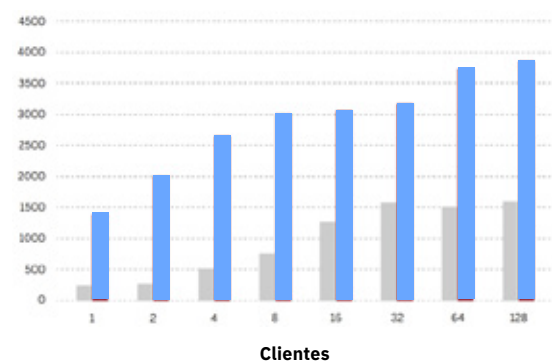
La herramienta de comparativa registró las operaciones de lectura y escritura por segundo en cada clúster, basándose en el número de clientes concurrentes participantes. Los resultados de la prueba fueron muy reveladores. Todos los entornos bare metal superaron en rendimiento al servidor virtual equivalente en términos de promedios de lectura y escritura.

## Servidores virtuales frente a servidores bare metal

**Promedio de operaciones de lectura por segundo por cada cliente concurrente**



**Promedio de operaciones de escritura por segundo por cada cliente concurrente**



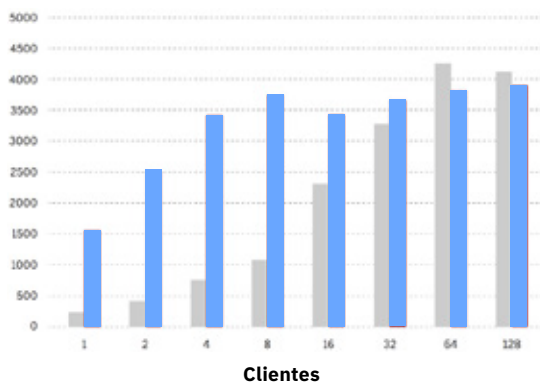
■ Servidores virtuales ■ Servidores bare metal

Como el entorno bare metal pudo utilizar directamente los recursos hardware del servidor y no tuvo que competir por los recursos con otros usuarios, los servidores bare metal ofrecieron un rendimiento hasta seis veces mejor que los servidores virtuales equivalentes.

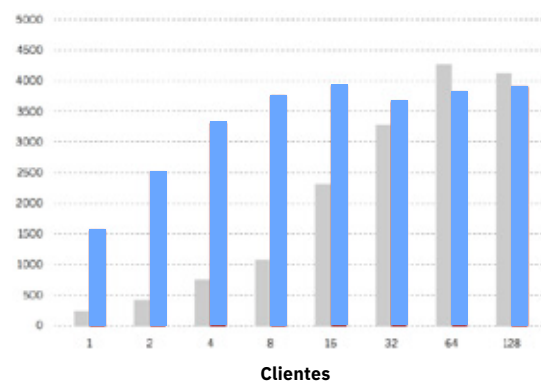
Al medir el promedio de operaciones de lectura y escritura por segundo, la herramienta de comparativa también registró el rendimiento máximo de cada entorno, y estos resultados también son dignos de mención (por distinto motivo):

## Servidores virtuales frente a servidores bare metal

**Promedio de operaciones de lectura por segundo por cada cliente concurrente**



**Promedio de operaciones de escritura por segundo por cada cliente concurrente**



■ Servidores virtuales ■ Servidores bare metal

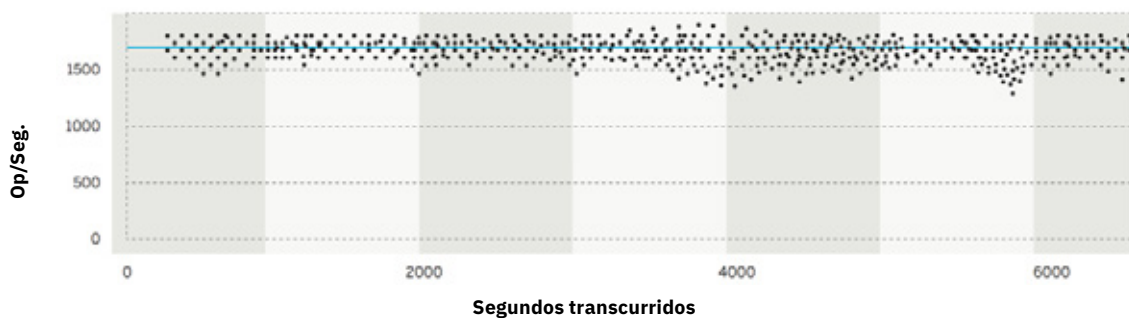
Los resultados máximos para operaciones de lectura y escritura por segundo en los entornos bare metal fueron muy cercanos a los promedios de operaciones de lectura y escritura registrados para ese entorno, pero los resultados máximos fueron enormemente diferentes de los resultados medios en el entorno de servidores virtuales. En dos de los escenarios, los servidores virtuales lograron un máximo superior a los bare metal. Si los tomamos en contexto con el promedio de operaciones por segundo registrado por el entorno de servidores virtuales, los resultados destacan el otro indicador clave del rendimiento para cargas de trabajo de big data: **la homogeneidad.**

# Big Data - Homogeneidad

El rendimiento solo es significativo si es homogéneo. En nuestra prueba de rendimiento, el entorno de servidores virtuales tal vez registrase 4500 operaciones de lectura por segundo como máximo, pero como promedio, este entorno proporcionaba 1500 operaciones de lectura por segundo. Si el rendimiento de un entorno varía de forma tan significativa de un segundo al siguiente, es extremadamente difícil crear un entorno capaz de manejar una carga de trabajo en constante aumento. Para comparar la homogeneidad de los resultados en servidores bare metal en comparación con servidores virtuales, un ingeniero de IBM Cloud configuró dos clústeres Riak de cinco nodos y simuló despliegues de cargas utilizando Basho Bench (información detallada en el Apéndice B). Esta prueba observó y representó las operaciones por segundo a lo largo de un periodo de dos horas:

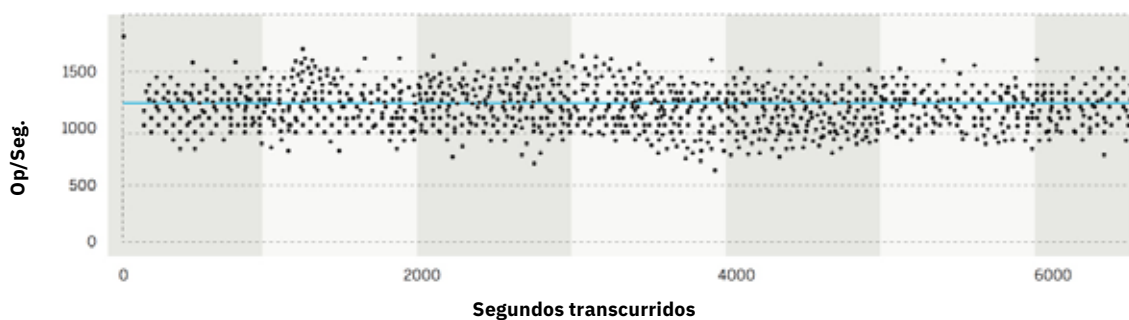
## Servidores bare metal:

Operaciones por segundo bajo carga (2 horas)



## Servidores virtuales:

Operaciones por segundo bajo carga (2 horas)





---

El entorno bare metal promedia más operaciones por segundo en toda la prueba, pero lo más indicativo es que los resultados están mucho más concentrados en torno al promedio. Cuando el rendimiento varía significativamente de un segundo al siguiente en el entorno de servidores virtuales, la planificación de la capacidad se convierte en un problema. ¿Qué estadística debería utilizarse para decidir si escalar un entorno hacia arriba o hacia abajo? Si construye su entorno para admitir los peores resultados registrados, cuando el rendimiento sea bueno estará sobreaprovisionando recursos. Optar por una capacidad básica en base a los mejores recursos probablemente provoque un bajo rendimiento del entorno. Y basar la capacidad en los resultados medios es como elegir entre estas dos alternativas tirando una moneda al aire.

---

Las empresas se basan en **resultados homogéneos** para predecir tendencias, asignar presupuestos y tomar decisiones importantes. Planificar los entornos de infraestructura de cloud no debería ser distinto.

---

# Big Data necesita bare metal

Es posible que sea difícil resistirse a las promesas de despliegue rápido y sencillo. Algunas aplicaciones son más indicadas para ejecutarse en servidores virtuales en un entorno de cloud público, pero el big data no es una de ellas.

Es importante advertir que:

- Las dos características más importantes de un entorno de cloud ejecutando cargas de trabajo con elevado nivel de E/S como big data son **rendimiento** y **homogeneidad**.
- Los servidores bare metal pueden configurarse y optimizarse para **ofrecer un rendimiento incomparable** al servir y **procesar muy elevados volúmenes de datos**.
- Los servidores virtuales con cargas de trabajo con elevado nivel de E/S pueden verse **afectados negativamente por el uso de recursos por parte de otros clientes** cuando múltiples usuarios comparten el mismo nodo host de servidores virtuales.
- Los recursos de servidores bare metal **son locales y no compartidos**, por lo que las cargas de trabajo ofrecen un rendimiento mucho más homogéneo que en entornos de servidores virtuales compartidos y/o en red.
- Los servidores virtuales pueden aprovisionarse rápidamente y escalarse horizontalmente mucho más rápidamente que los servidores bare metal, pero las cargas de trabajo que no precisan ráfagas se beneficiarán del rendimiento y homogeneidad de los servidores bare metal.

---

# ¿Por qué IBM Cloud es un proveedor ideal para cargas de trabajo de big data?

**Tecnología sin rival:** IBM Cloud le proporciona la infraestructura de cloud de más elevado rendimiento del mercado. Tanto si su big data se extiende a nivel global o local, nuestros centros de datos en todo el mundo y nuestros servidores bare metal y virtuales de primera clase están preparados para cualquier tarea.

**Red sin fisuras:** Nuestra red de alta velocidad integra redes públicas, privadas y de gestión interna para ofrecer la máxima velocidad, algo esencial al analizar y transferir big data.

**Gestión y automatización totales:** Hemos desarrollado un tipo distinto de solución de cloud: una plataforma integral automatizada. Cada servidor, dispositivo de almacenamiento y servicio de gestión y seguridad puede controlarse mediante un mismo sistema de gestión, con acceso desde nuestra API, portal del cliente e incluso aplicaciones móviles.

**Ejecute su big data en servidores bare metal. Nuestros expertos de IBM Cloud le ayudarán a crear una infraestructura de cloud de alto rendimiento perfectamente adaptada a sus necesidades de big data.**

Explore los servidores IBM Cloud bare metal y virtuales en <http://ibm.co/bare-metal> y obtenga información detallada sobre soluciones a medida para big data y mejores prácticas para aplicaciones concretas para Riak, Hadoop y MongoDB en <http://ibm.co/big-data>.

¿Desea hacer alguna otra pregunta?

Pregunte a un experto: <http://ibm.co/contact-us> o bien llámenos al: **214-442-0600**.

---

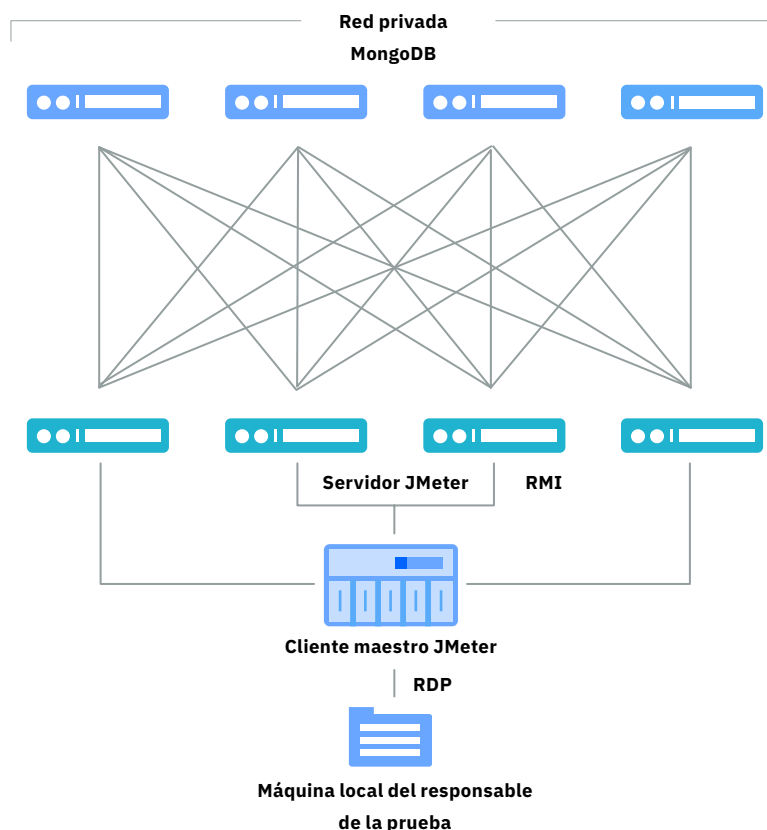
# Anexo A

## Metodología de pruebas de rendimiento de Big Data — MongoDB

Se precargaron conjuntos de datos de documentos de 512 kb en instancias únicas de MongoDB en cada servidor. Los conjuntos de datos se crearon con tamaños variables en comparación con la memoria disponible, con conjuntos de datos mayores (2x) y menores que la memoria disponible. La prueba también comprobó que el conjunto de datos se alterase durante la ejecución de prueba con frecuencia suficiente para evitar que las consultas pusieran todos los datos en memoria caché.

Una vez creados los conjuntos de datos, se utilizaron instancias de servidores JMeter con 4 cores y 16 GB de RAM para ejecutar 'benchrun' desde la herramienta de comparativa de MongoDB. El diagrama inferior ilustra cómo configuramos el entorno de pruebas.

Estos servidores Jmeter funcionan como clientes generando tráfico en las instancias de MongoDB. Cada cliente generó peticiones aleatorias de consultas y actualizaciones con un ratio de seis consultas por cada actualización (las peticiones de actualización de la prueba pretendían asegurar que los datos no se asignasen completamente a memoria caché y nunca ejerciesen lecturas desde disco). Estas pruebas se diseñaron para crear una carga extrema en los servidores desde un número en aumento exponencial de clientes hasta saturar los recursos del sistema, y registramos el rendimiento resultante de la aplicación MongoDB.



### Configuración de las pruebas

- Conjunto de datos (32 GB de documentos de 0,5 mb)
- 200 iteraciones de operaciones consulta-actualización 6:1
- Las conexiones concurrentes de clientes se aumentaron exponencialmente de 1 a 128
- La duración de la prueba se prolongó durante 48 horas

## Apéndice A (continuación)

Metodología de pruebas de rendimiento de Big Data — MongoDB

### Servidores bare metal frente a servidores virtuales

	<b>Nodo de servidores Bare Metal</b>	<b>Nodo de servidores virtuales</b>
Core	CPUs Intel 5670 dual 6 cores	26 unidades de computación virtuales
Sistema operativo	CENTOS de 64 bits	CENTOS de 64 bits
RAM	36 GB de RAM	30 GB de RAM
RAID	2 SSD de 64 GB RAID1 (Journal Mount)	2 x 64 GB de almacenamiento en red RAID1 (Journal Mount)
SAS	4 SSD de 400 GB RAID10 (Data Mount)	4 SSD de 300 GB RAID10 (Data Mount)
Red	Red 1 GB   Conectado	Red 1 GB

# Apéndice B

## Metodología de pruebas de rendimiento de Big Data — Riak

Se crearon clústeres de cinco nodos con Riak 1.3.1 en servidores bare metal y en servidor virtual en cloud público. Para optimizar el rendimiento de Riak, se hicieron ajustes a nivel de SO en cada servidor (ejecutando CentOS de 64 bits):

**Noatime**

**Nodiratime**

**barrier=0**

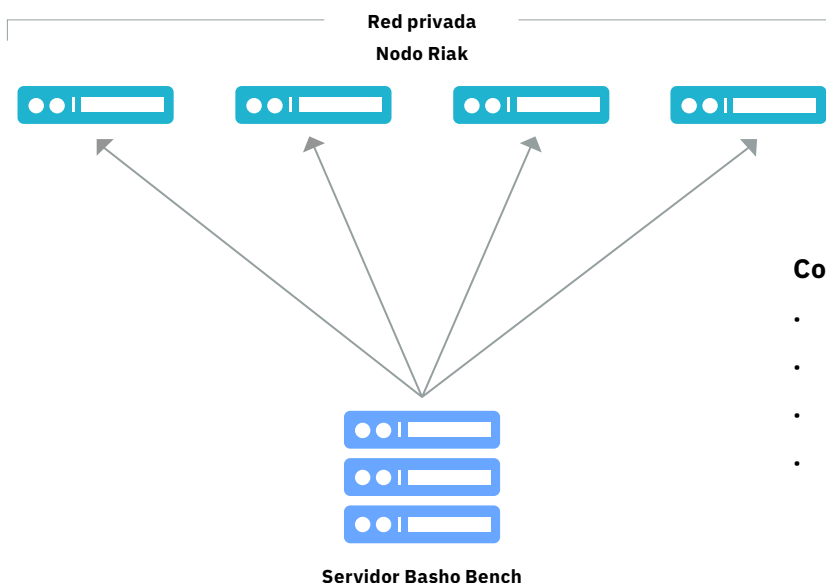
**data=writeback**

**ulimit -n 65536**

Los ajustes comunes de Noatime y Nodiratime eliminan la necesidad de escrituras durante las lecturas para contribuir al rendimiento y desgaste del disco. Los ajustes de barrier y writeback son algo menos comunes y posiblemente no sean los que se ajustarían normalmente. Aunque estos ajustes presentan un muy pequeño riesgo de pérdida de datos en caso de fallo de un disco, recuerde que la solución Riak se despliega en anillos de cinco nodos con datos disponibles redundantemente en múltiples nodos del anillo.

Teniendo esto en cuenta y considerando que cada nodo también se despliega con un array de almacenamiento RAID10, la disminución del riesgo de pérdida de datos en caso de fallo de un único disco en toda la solución no tendría impacto en el conjunto de datos completo (ya que existen numerosas copias redundantes de estos datos). Dado el menor riesgo, los aumentos en el rendimiento de estos dos ajustes justifican el uso.

Con todos los nodos ajustados y configurados en clústeres, configuramos la herramienta de pruebas de Basho (Basho Bench) para simular remotamente carga en los despliegues. Basho Bench permite crear un plan de pruebas configurable para un clúster Riak configurando una serie de trabajadores para utilizar un tipo de controlador y para generar carga. Viene empaquetado como aplicación Erlang con un archivo de configuración de ejemplo que es posible modificar para crear los datos específicos para concurrencia, tamaño de conjuntos de datos y duración de las pruebas. Los resultados pueden verse como datos CSV y existe un paquete opcional de gráficos para generar gráficos. Un gráfico simplificado de nuestro entorno de pruebas es así:



### Configuración de las pruebas

- Conjunto de datos: 400 GB
- Operaciones consulta-actualización 10:1
- 8 conexiones de clientes concurrentes
- Duración de la prueba: 2 horas

## Apéndice B (continuación)

Metodología de pruebas de rendimiento de Big Data — Riak

### Riak - Prueba de homogeneidad

Clúster de 5 nodos bare metal frente a clúster de 5 nodos de servidores virtuales

	<b>Nodo de servidores Bare Metal</b>	<b>Nodo de servidores virtuales</b>
Core	CPUs Intel 5670 dual 6 cores	26 unidades de computación virtuales
Sistema operativo	CENTOS de 64 bits	CentOS de 64 bits
RAM	36 GB de RAM	30 GB de RAM
RAID	4 x SAS 15K 300 GB   RAID10	4 x 300 GB de almacenamiento en red
SAS	Red 1 GB   Conectado	Red 1 GB



IBM España S.A.  
Sta. Hortensia 26-28  
28002 Madrid  
España

El sitio web de IBM está disponible en [ibm.com/es](http://ibm.com/es)

IBM, el logotipo de IBM, [ibm.com](http://ibm.com) y SPSS son marcas comerciales de International Business Machines Corp., registradas en numerosas jurisdicciones de todo el mundo. Otros nombres de productos y servicios pueden ser marcas registradas de IBM u otras compañías. Bajo el epígrafe “Información sobre Copyright y marcas comerciales” puede consultar la lista actualizada de las marcas comerciales de IBM en la página web [www.ibm.com/legal/copytrade.shtml](http://www.ibm.com/legal/copytrade.shtml)

Este documento está actualizado en la fecha de publicación original y puede ser modificado por IBM en cualquier momento. No todas las ofertas están disponibles en todos los países en los que opera IBM.

Los ejemplos de cliente que se citan se presentan solo a título ilustrativo. Los resultados de rendimiento reales pueden variar según las configuraciones y condiciones de operación específicas. Es responsabilidad del usuario evaluar y verificar la operación de cualquier otro producto o programa con los productos y programas IBM. LA INFORMACIÓN DE ESTE DOCUMENTO SE PROPORCIONA “TAL CUAL”, SIN NINGUNA GARANTÍA, NI EXPLÍCITA NI IMPLÍCITA, INCLUYENDO LAS GARANTÍAS DE COMERCIALIZACIÓN, IDONEIDAD PARA UN FIN DETERMINADO Y NO INCUMPLIMIENTO. Los productos IBM están garantizados de acuerdo con los términos y condiciones de los acuerdos en virtud de los cuales se proporcionen.

© Copyright IBM Corporation 2018



Por favor, recicle