

Informe de Materia de Graduación

"Uso de la plataforma Pig sobre Hadoop como alternativa a una RDBMS para el análisis de datos masivos. Prueba de concepto utilizando registros de detalles de llamadas"

Presentado por:

**Romeo Cabrera Arévalo
Fabricio Medina Palacios**

Profesora: Ing. Cristina Abad

Introducción

- ◆ En la actualidad se ha dado una explosión y alta penetración en la telefonía celular.
- ◆ La información de registros (CDRs) generada por el uso de los servicios es del orden de los terabytes al mes.
- ◆ Enfoque tradicional: Almacenar esta información en RDBMSs para su procesamiento.
- ◆ El paradigma “en la nube” para realizar procesamiento paralelo masivo de información surge como una alternativa.

Objetivos

1. Comprobar la escalabilidad y adecuación al uso de la herramienta Pig sobre Hadoop para el procesamiento de cantidades masivas de registros.
2. Comparar la razón costo/rendimiento entre el uso de Pig-Hadoop en un clúster en la nube contra el uso de un RDBMS comercial para procesar registros masivamente.
3. Demostrar la facilidad de crear consultas ad-hoc usando Pig para diversas y cambiantes necesidades de análisis de información.

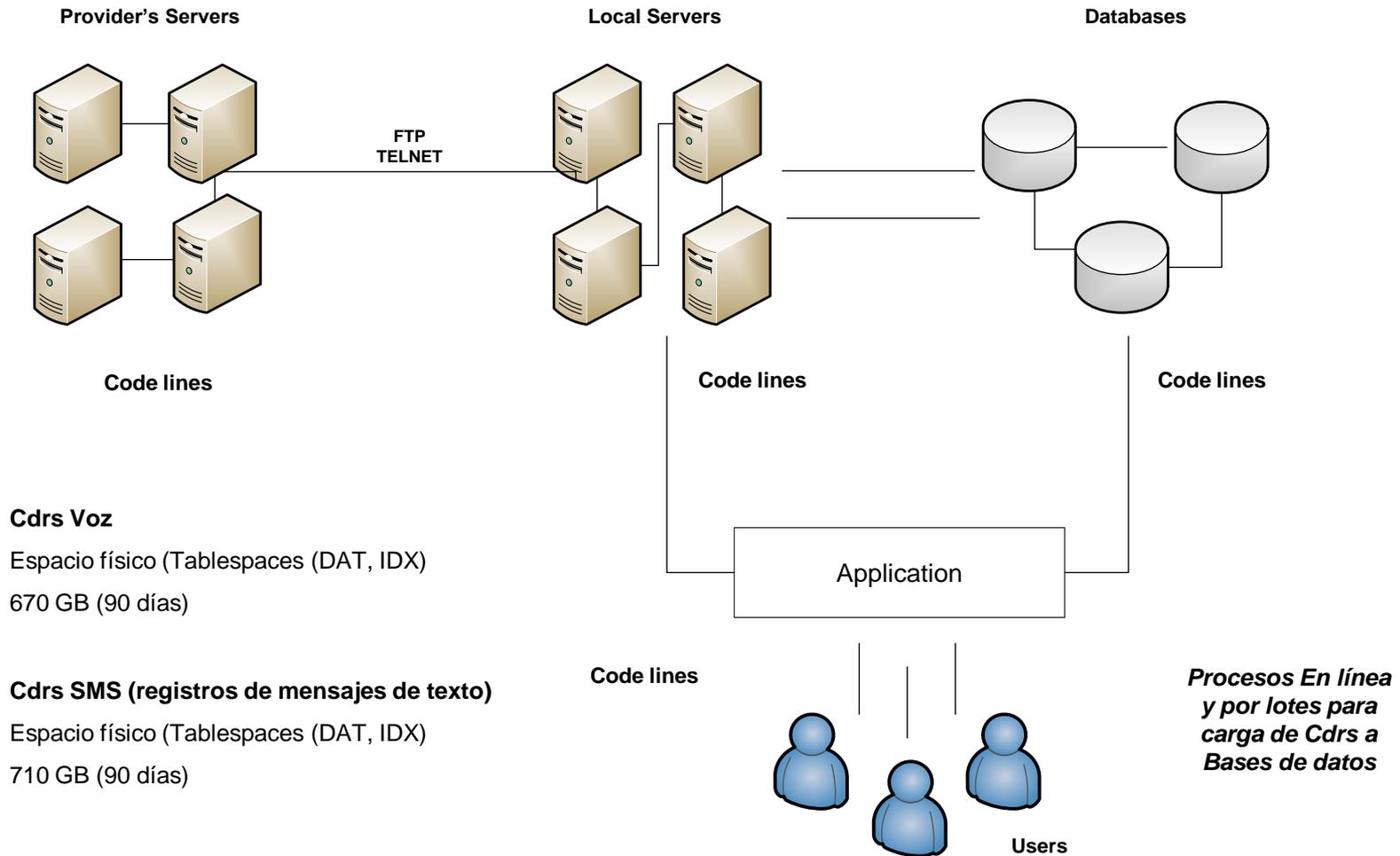
Alcance

1. Instalar una plataforma basada en Pig sobre Hadoop para el procesamiento de registros de detalles de llamadas.
2. Configurar esta plataforma para que sea ejecutada sobre Amazon EC2.
3. Desarrollar scripts en Pig que realicen análisis de esta información.
4. Comparar esta solución frente al uso de una herramienta comercial (Oracle) en una empresa de telefonía celular local.
5. Desarrollar un sencilla interfaz Web para poder ingresar y/o subir scripts en Pig para esta plataforma, y para realizar una visualización de los resultados.

Archivos CDRs

- ◆ Registros informáticos generados por una central telefónica, los cuales contienen detalles de los eventos que han pasado a través de ella.
- ◆ Ejemplo:
 - 10888,59390730123,6-29-2009,86067,6-30-2009,17667,202,304,1222F832640858180069073012300980,1436,0,0,0,319,,108,0,0,0,701,1002,101,-1,0,,0,9,35,9/28/200817:55:0,1/1/1900 0:0:0,9/28/200817:55:0,0,0,0,0,0,,0,0,0,0,0,0,1,2,1,350,16,50,-1,0,59385818006,740010115532771,59397995028,2,1
- ◆ Especifica datos como: Número origen, Número destino, Hora de evento, Celda origen, Duración del evento, Perfil del suscriptor, Tipo de Tarificación del evento, entre otros.

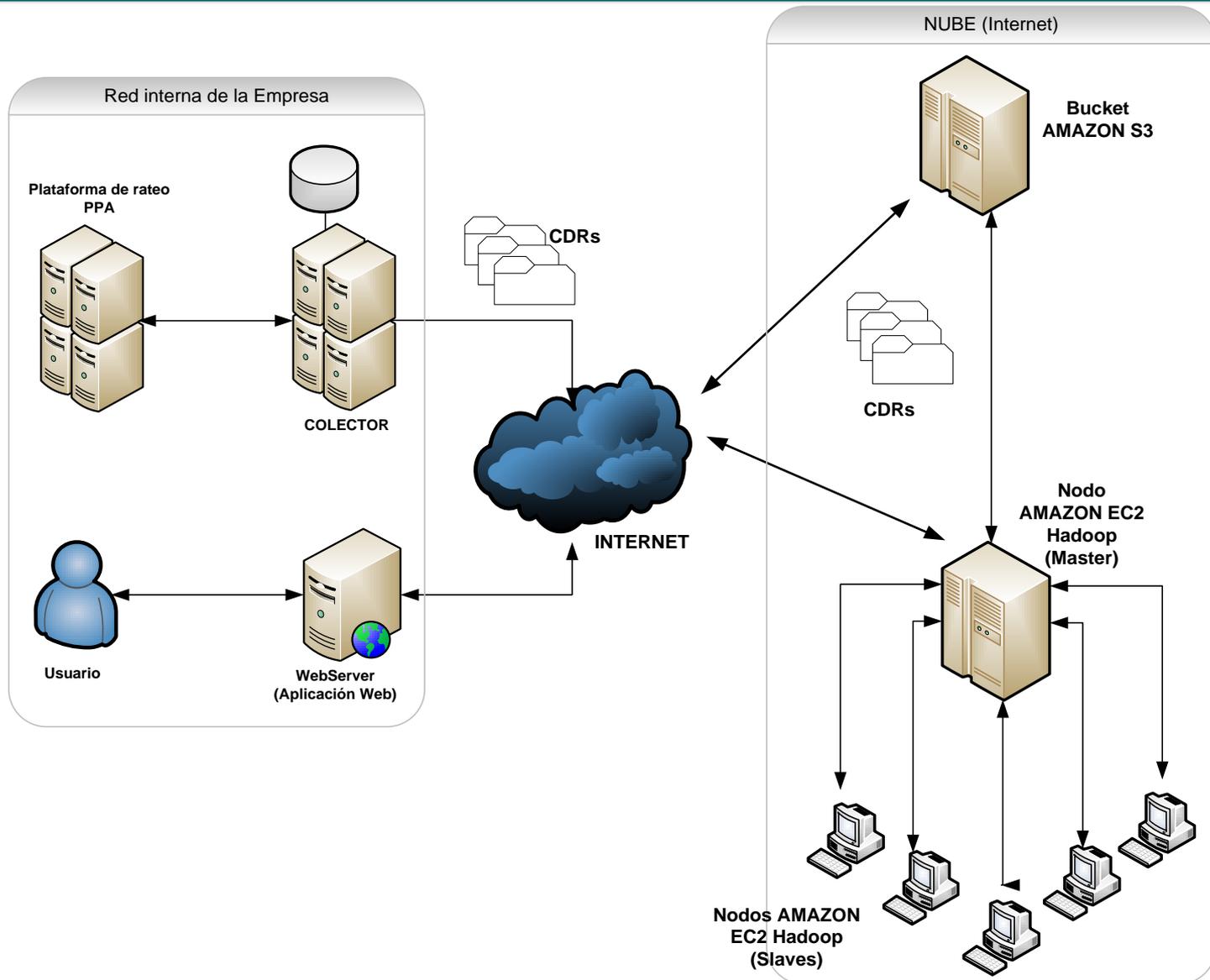
Arquitectura sistema actual



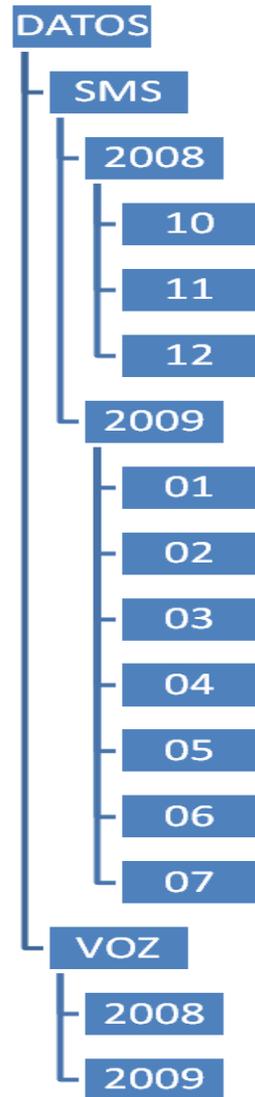
Diseño

- ◆ Hadoop: Implementación Open Source del paradigma MapReduce de computación distribuida.
- ◆ Pig: Capa de software que recibe scripts en un language de flujo de datos y los convierte en trabajos MapReduce.
- ◆ Uso de archivos comprimidos BZ2
- ◆ Uso de servicios Web de Amazon
 - Simple Storage Service (S3)
 - Elastic Cloud Computing (EC2)

Arquitectura de la solución



Esquema de directorios



Script en Pig

```
A = LOAD 'cdrs/data/sms/2009/01'  
  USING PigStorage(',');  
Y = GROUP A BY $1;  
Z = FOREACH Y GENERATE  
  $0,COUNT(A);  
N = ORDER Z BY $1 DESC;  
STORE N INTO 'output/sms' USING  
  PigDump();
```

Pantalla de consultas predeterminadas

The screenshot shows a Mozilla Firefox browser window titled 'Consultator 1.0 - Mozilla Firefox'. The address bar displays 'http://localhost/mysite/consultaPredeterminada.php'. The page content includes a navigation menu with 'Inicio', 'Esquema propuesto', 'Bibliografía', and 'Contáctenos', and a date 'Jueves, 29 de Octubre de 2009'. A sidebar menu on the left has 'Consultas' selected, with sub-items 'Predeterminadas' and 'Ingresar consulta'. The main content area is titled 'Consultas predeterminadas:' and lists three options: '1) Cantidad de llamadas por rango de fechas', '2) Cantidad de sms por rango de fechas', and '3) Conteo de llamadas por duración, por fechas'. The footer contains the text '© Consultator 1.0. Derechos reservados. Diseñado por: Fabricio Medina.' and a 'Terminado' status bar at the bottom left.

Consultator 1.0

Inicio Esquema propuesto Bibliografía Contáctenos Jueves, 29 de Octubre de 2009

Consultas
Predeterminadas
Ingresar consulta

Consultas predeterminadas:

Seleccione consulta predeterminada a utilizar:

- 1) [Cantidad de llamadas por rango de fechas](#)
- 2) [Cantidad de sms por rango de fechas](#)
- 3) [Conteo de llamadas por duración, por fechas](#)

© Consultator 1.0. Derechos reservados. Diseñado por: Fabricio Medina.

Terminado

Ingreso de script por pantalla

Consultator 1.0 - Mozilla Firefox (Build 20100115144158)

Archivo Editar Ver Historial Marcadores Herramientas Ayuda

http://ecua.com/hadoop/mysite/ingresaScript.php

Consultator 1.0

Consultator^{1.0}

[Inicio](#) [Esquema propuesto](#) [Bibliografía](#) [Contáctenos](#) Domingo 28 de Febrero de 2010

Consultas
Predeterminadas
Ingresar consulta

Reportes
Ver Resultados

Ingreso por teclado:

Por favor ingrese script en PIG a procesar:

Descripción:

```
A = LOAD 'cdrs/data/SMS/2009/01' USING PigStorage(',');
Y = GROUP A BY $1;
Z = FOREACH Y GENERATE $0,COUNT(A);
N = ORDER Z BY $1 DESC;
STORE N INTO 'output' USING PigDump();
```

S3Fox

Ingreso de script por archivo

Consultator 1.0 - Mozilla Firefox (Build 20100115144158)

Archivo Editar Ver Historial Marcadores Herramientas Ayuda

http://ecua.com/hadoop/mysite/subeArchivo.php

Consultator 1.0 Consultator 1.0 Consultator ... http://...=REQ13 http://...=REQ13 http://...=REQ13 http://...=REQ13



Consultator^{1.0}

[Inicio](#) [Esquema propuesto](#) [Bibliografía](#) [Contáctenos](#) Domingo 28 de Febrero de 2010

Consultas
Predeterminadas
Ingresar consulta ▶

Reportes
Ver Resultados

Ingreso por archivo:

Seleccione ubicación y nombre de archivo que contiene script en FIG. Condiciones: .txt; 1Mb max.

Descripción:

Archivo:

© Consultator 1.0. Derechos reservados. Diseñado por: Fabricio Medina.

Terminado S3fox

Consulta de requerimientos

Consultator 1.0 - Mozilla Firefox (Build 20100115144158)

Archivo Editar Ver Historial Marcadores Herramientas Ayuda

http://ecua.com/hadoop/mysite/consultaResultados.php

Consultator 1.0 Consultator 1.0 http://e...D=REQ13 http://e...D=REQ13 http://e...D=REQ13 http://e...D=REQ13

Consulta de requerimientos

Buscar requerimiento:

Requerimiento	Detalle Requerimiento	Fecha Requerimiento	Estado
REQ13	Prueba 13	2009-09-14 14:07:35	Finalizado
REQ15	Prueba 15	2009-09-18 00:00:00	Finalizado
REQ10	Prueba 10	2009-09-14 13:44:43	Procesando
REQ14	Prueba 14	2009-09-18 00:00:00	Procesando
REQ16	Prueba 16	2009-09-18 13:48:24	Procesando
REQ17	Prueba 17	2009-09-18 13:55:35	Procesando
REQ26	Cantidad de llamadas por rango de fechas 2009/09/08 - 2009/09/23	2010-01-08 13:38:45	Procesando
REQ28	Prueba nueva de ingreso teclado	2010-01-08 13:56:40	Procesando
REQ29	Prueba nueva de ingreso teclado 2	2010-01-08 14:24:09	Procesando
REQ30	Prueba nueva de ingreso teclado 3	2010-01-08 14:24:49	Procesando
REQ31	Prueba nueva de ingreso por archivo	2010-01-08 14:25:48	Procesando

Total de registros: 11

Terminado

S3 Fox

- Consultas
 - Predeterminadas
 - Ingresar consulta
- Reportes
 - Ver Resultados

Consulta de detalle de requerimiento

Consultator 1.0 - Mozilla Firefox (Build 20100115144158)

Archivo Editar Ver Historial Marcadores Herramientas Ayuda

http://ecua.com/hadoop/mysite/detalleRequerimiento.php?codigo=REQ13

Consultator 1.0 Consultator... Consultator 1.0 http://...REQ13 http://...REQ13 http://...REQ13 Consultator 1.0

Detalle de requerimientos

Requerimiento: REQ13

Detalle Requerimiento: Prueba 13

Consulta:

```
A = LOAD 'cdrs/data/VOZ/2009/01' USING PigStorage(','); X = FOREACH A GENERATE $4/10 AS decimo:int , 1 AS cuenta:int; Y = GROUP X BY decimo; Z = FOREACH Y GENERATE $0,COUNT(X); DUMP Z; STORE Z INTO 'stats' USING PigDump();
```

Fecha Requerimiento: 2009-09-14 14:07:35

Ubicación: /opt/lampp/htdocs/hadoop/mysite/archivos/

Estado: Finalizado

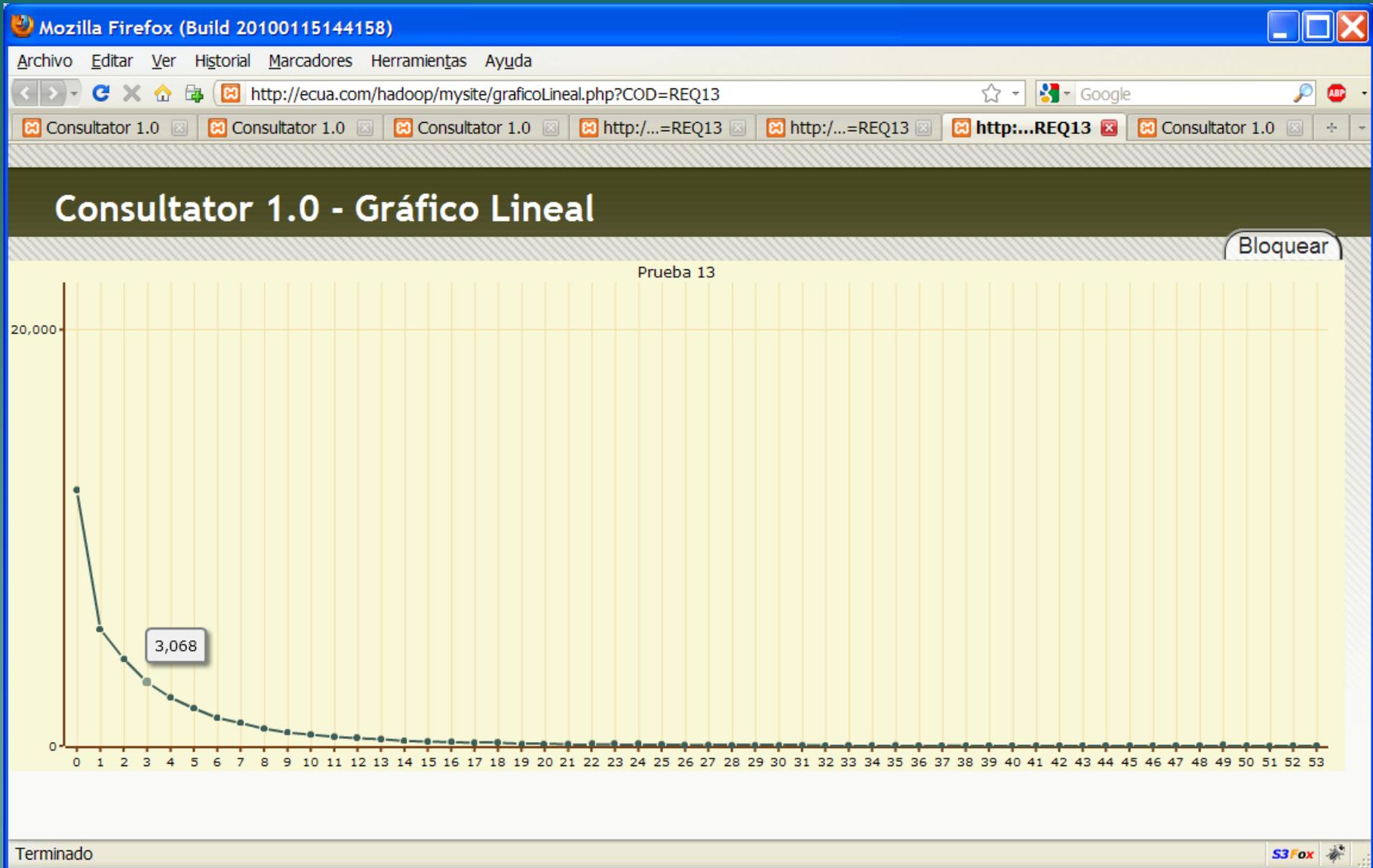
Verificar Resultados

Gráfico Lineal	Gráfico de barras	Gráfico Pie	Tabla de valores
--------------------------------	-----------------------------------	-----------------------------	----------------------------------

Terminado

S3 Fox

Gráfico de resultado

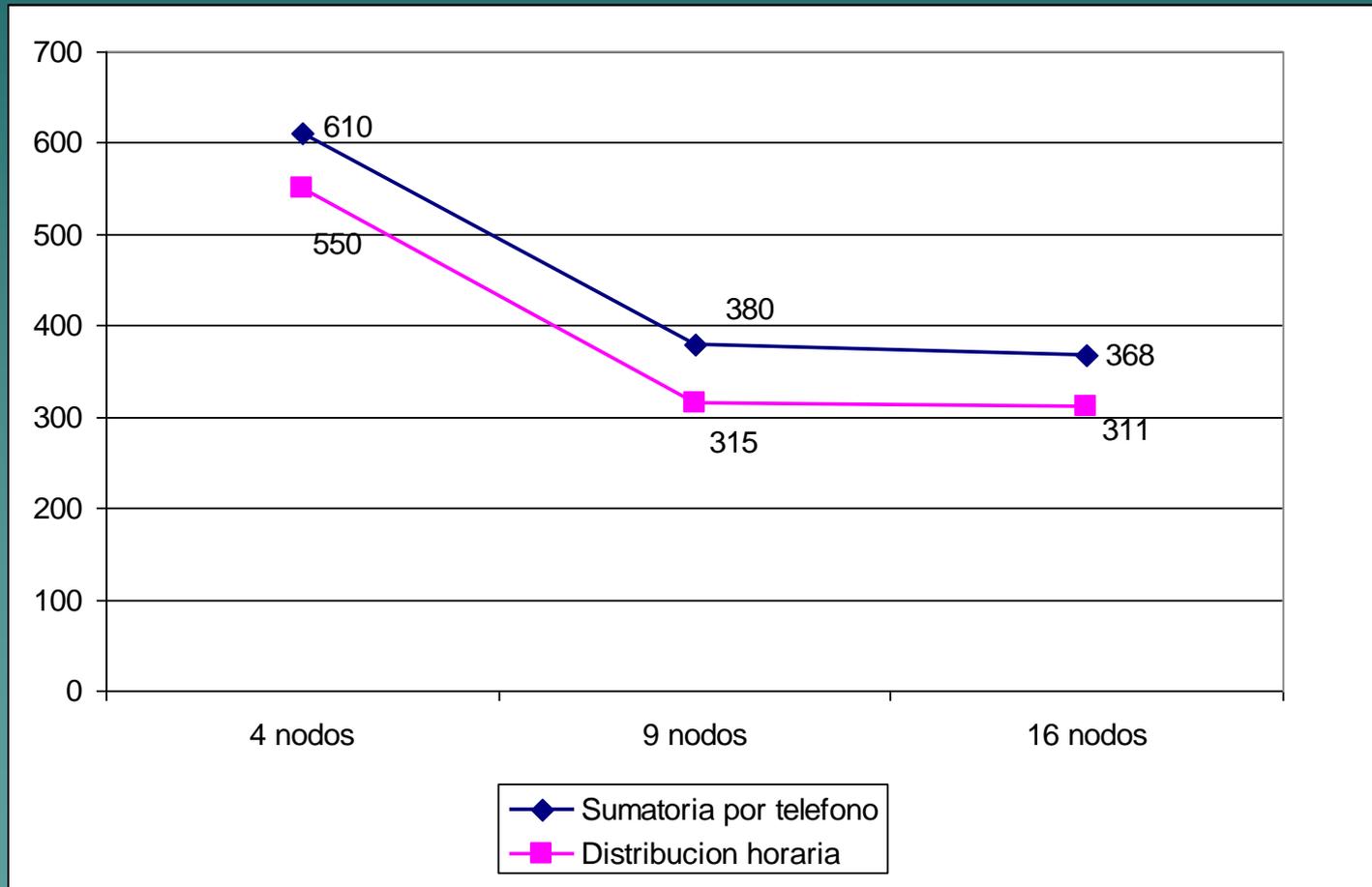


Pruebas

- ◆ Se probó con archivos CDRs de 3,15 GB. (1,54 GB comprimidos)
- ◆ Dos procesos:
 - Sumatoria de SMS.
 - Distribución de envío de mensajes.
- ◆ Nodos Linux “High CPU medium” (1.7 GB de memoria, 5 EC2 computing units, 350 GB de almacenamiento local y plataforma de 32 bits.)
- ◆ 1 EC2 C.U. = Xeon 2007 de 1.0-1.2 GHz

Resultados

Tiempo de ejecución de procesos (en segundos)



Costo monetario

Esquema tradicional

Costo licencia RDBMS por procesador:	\$47.500
No. De procesadores:	24
Costo total de licencia RDBMS para servidor:	\$1'140.000
Costos de soporte anual RDBMS:	\$10,450
Costo anual soporte de hardware	\$30,000
TOTAL	~\$1'120,000

Costo anual de soporte HW: \$250.000

Costo del HW: \$1 Millón

Nuevo esquema

Almacenamiento	
2500 GB /12 meses x \$0,15 /GB	\$4500 anuales
Transferencia	
210 GB / meses x 12 meses x \$0,17 /GB	\$428 anuales
Procesamiento	
Pago único anual:	\$455
Horas al año x instancia:	8760 (uso 24/7)
Número de instancias:	20
Horas utilizadas al año total:	175 200
Costo por nodo por hora, esquema "Reserved Instances":	\$0,06
Total procesamiento:	\$10512
TOTAL GENERAL:	~\$15000

Conclusiones

- ◆ La solución presentada permite el realizar consultas y análisis sobre volúmenes de información que no hubieran sido posibles en un esquema de RDMS convencional.
- ◆ Esta tecnología no implica un reemplazo de una RDBMS tradicional, más bien la complementa.
- ◆ Pig minimiza el tiempo necesario para implementar un requerimiento ad-hoc.
- ◆ El costo monetario de almacenamiento y procesamiento en un clúster en la nube es dos órdenes de magnitud inferiores al de una solución tradicional.
- ◆ El uso de las plataformas Pig, MapReduce, EC2, simplifican el desarrollo de aplicaciones distribuidas.

Recomendaciones

- ◆ Uso de EBS (Elastic Block Store) como alternativa a S3.
- ◆ Actualizar a Hadoop 0.20 (se utiliza 0.18). Mejoras varias en rendimiento y permite referenciar directamente archivos S3.
- ◆ Levantar los nodos en demanda, y en una cantidad óptima para cada script.
- ◆ Uso de SQS (Simple Queue System) o similar para control y priorización de trabajos enviados a procesar.
- ◆ Usos adicionales: Minería de datos para marketing, detección de patrones de fraude, etc.

Preguntas



¡Gracias!

