

Instalación de R en un clúster desplegado en la nube

Pablo Baños López

Facultad de Informática - Universidad de Murcia

Septiembre de 2009

Índice

Motivación

Entorno de trabajo

- Entornos de cómputo

- Librerías científicas

- La plataforma cloud de Rightscale

Herramientas utilizadas

- Sistema de ficheros distribuido

- Servicio de centralización de usuarios

- Sistema gestor de colas

- Librerías de intercambio de mensajes

- R y Rmpi

Evaluación de rendimiento

- Motivación y procedimiento

- Resultados

Conclusiones y vías futuras

Motivación

Proyecto Fin de Carrera en la empresa Luss Innovation S.L.

- ▶ Instalación de software científico (R, Rmpi)
- ▶ Entorno final: cluster de supercomputación
- ▶ Entorno de trabajo: el cloud

Entornos de cómputo

- ▶ Clúster
- ▶ Cloud
- ▶ Supercomputador

Librerías científicas

- ▶ Software científico: Entorno estadístico R y Rmpi
- ▶ Interfaces de intercambio de mensajes MPI: Open MPI y LAM/MPI
- ▶ Compiladores: GNU e Intel

La plataforma cloud de Rightscale

Rightscale nos ofrece un servicio cloud de infraestructura de cómputo, basado en los servicios Amazon Web Service, donde podemos:

- ▶ Levantar instancias eligiendo una imagen con un sistema operativo
- ▶ Automatizar tareas en el arranque y durante la operación
- ▶ Crear plantillas para la configuración de instancias similares
- ▶ Utilizar un dispositivo de almacenamiento permanente pagando en función de lo que se usa.

Herramientas utilizadas

- ▶ Sistema de ficheros distribuido
- ▶ Servicio de centralización de usuarios
- ▶ Sistema gestor de colas
- ▶ Librerías de intercambio de mensajes
- ▶ R y Rmpi

Sistema de ficheros distribuido

Network File System (NFS)

- ▶ Servidor NFS en el nodo maestro
- ▶ Clientes NFS en los nodos esclavo
- ▶ Compartición del directorio /opt

Servicio de centralización de usuarios

Network Information Service (NIS)

- ▶ Servidor NIS en el nodo maestro
- ▶ Clientes NIS en los nodos esclavo

Sistema gestor de colas

Sun Grid Engine (SGE)

- ▶ Se instala desde el maestro
- ▶ Se especifican los nodos y su rol
- ▶ Se instala en el directorio compartido /opt
- ▶ Usuario sgeadmin como administrador del sistema

Librerías de intercambio de mensajes

- ▶ Instalación manual de Open MPI
- ▶ Instalación de LAM/MPI mediante paquete RPM

R y Rmpi

- ▶ Creación de paquetes RPM para la instalación de R
- ▶ Instalación de Rmpi

Motivación

En el entorno desarrollado en el cloud:

- ▶ Determinar con que compilador se obtienen mejores resultados
- ▶ Determinar con que librería MPI se obtienen mejores resultados

Así preparamos el procedimiento para el entorno final.

Procedimiento

Se han realizado 5 ejecuciones para cada una de las configuraciones del entorno, que se han realizado variando:

- ▶ El número de nodos esclavo de 1 a 16 en potencias de 2
- ▶ El compilador utilizado en la instalación de R y Rmpi
- ▶ La implementación MPI utilizada (Open MPI o LAM/MPI)

Comparación según los compiladores

Con ambos compiladores la aplicación ha escalado correctamente, comportandose mejor con los compiladores GNU.

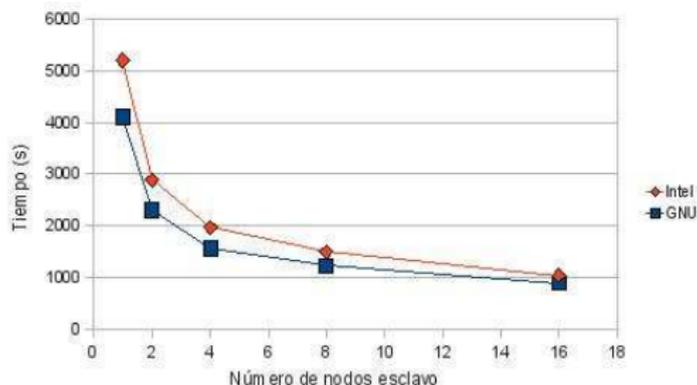


Figura: Tiempos de ejecución utilizando LAM/MPI

Comparación según los compiladores

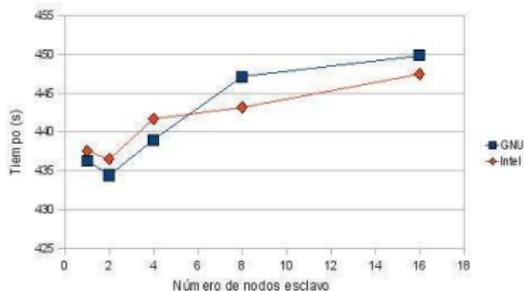


Figura: Tiempo de comunicación inicial utilizando LAM/MPI

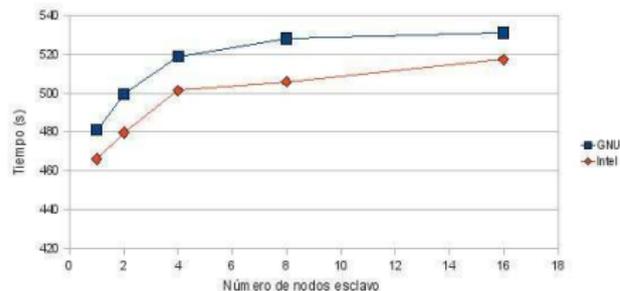


Figura: Tiempo de comunicación inicial utilizando Open MPI

Comparación según la implementación de MPI

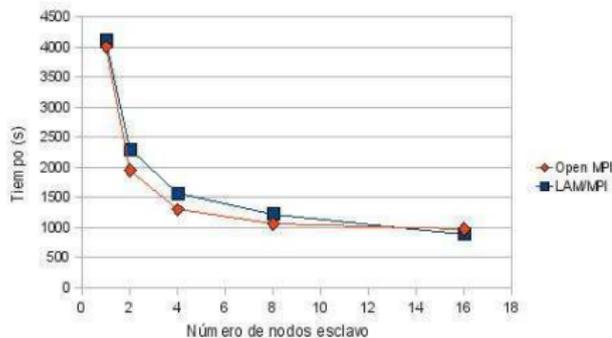


Figura: Tiempos de ejecución utilizando compiladores GNU

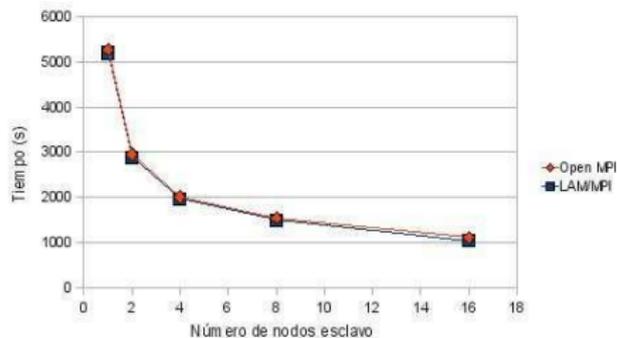


Figura: Tiempos de ejecución utilizando compiladores Intel

Conclusiones

- ▶ Con los compiladores GNU hemos obtenido mejores resultados
- ▶ LAM/MPI, para un número suficiente de nodos, se comporta mejor que Open MPI
- ▶ El entorno construido es adecuado para estudiar el rendimiento de aplicaciones
- ▶ Los resultados obtenidos no son extrapolables a otras arquitecturas

Vías futuras

- ▶ Realizar el estudio usando otros tipos de instancias en Rightscale
- ▶ Realizar el estudio variando el tamaño del problema
- ▶ Configurar el cluster utilizando otras herramientas
- ▶ Instalación y evaluación de rendimientos en el entorno objetivo