



**UNIVERSIDAD
DE BURGOS**

Microsoft HPC

V 1.0

José M. Cámara
(checam@ubu.es)

Introducción

Microsoft HPC (High Performance Computing) es la solución de Microsoft a la computación de alto rendimiento.

Está enfocado principalmente a clusters de consumo.

Se pretende obtener el máximo rendimiento de los recursos existentes para minimizar la inversión.

El sistema es altamente escalable debido a la incorporación de recursos internos y en la nube.

El rendimiento del conjunto se ve afectado por la inversión realizada en los distintos subsistemas: hardware de los nodos e interconexión.

Arquitectura del cluster

Nodo de cabecera
Windows Server + HPC Pack



WAN



Nodos en la nube
Windows Azure

Nodos de computación

Windows 7, 8, (Pro. Ent.), Server + HPC Pack



LAN / SAN

Estaciones de trabajo

Windows 7, 8, (Pro. Ent.), Server + HPC Pack



Arquitectura del cluster II

Tipos de nodos

- Nodo de cabecera: no tiene especiales características hardware (doble tarjeta de red).
- Nodo de computación: dedicado a HPC.
- Estación de trabajo: computador de propósito general que ocasionalmente se une al cluster para apoyo a la computación. Se pueden incorporar manualmente o de acuerdo con un horario preconfigurado.
- Nodo en la nube: computador virtual contratado con Microsoft.

Sistema operativo

- El nodo de cabecera debe está equipado con Windows Server
- El resto de nodos admiten la versión de servidor de Windows pero también sistemas operativos convencionales, habitualmente restringidos a versiones: profesional, enterprise o ultimate.
- Los nodos en la nube se integran bajo Windows Azure.

Arquitectura de red

- La conexión a los recursos en la nube y/o de gestión remota se realiza a través de una conexión WAN convencional. Obviamente un elevado ancho de banda es deseable.
- La conexión a los nodos internos se realiza mediante redes LAN o SAN; la mayoría de las soluciones comerciales son integrables. Los nodos pueden estar interconectados mediante redes privadas, red corporativa o ambas.

Despliegue del cluster

Instalación del servidor

Instalación del SO.
Configuración del dominio activo.
Instalación del paquete de computación.

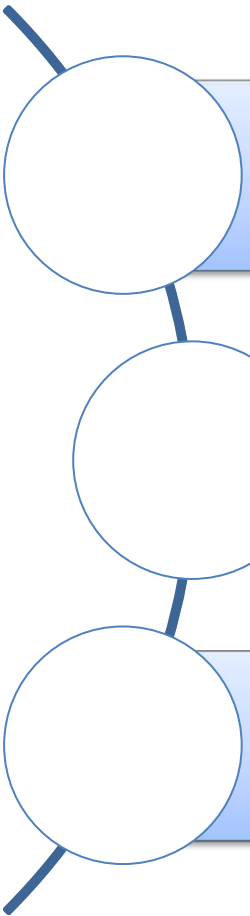
Configuración del nodo de cabecera

Definición de la topología del cluster.
Configuración de las comunicaciones.
Alta de usuarios.

Incorporación de nodos

Instalación del SO.
Incorporación de nodos al dominio.
Instalación del paquete de computación.
Inclusión de nodos en el cluster.

Dominio activo



Para que todas las operaciones a realizar en el cluster se puedan llevar a cabo de manera segura, todos los nodos deben pertenecer a un dominio común o a dominios con relaciones de confianza establecidas.

El dominio se configura como un nuevo bosque y habitualmente el servidor es promovido a controlador del dominio; a no ser que se emplee otro servidor para esta función.

Todos los nodos y usuarios deben pertenecer al dominio para poder ser añadidos al cluster.

Paquete de computación

El paquete de computación (HPC Pack) es proporcionado por Microsoft de forma gratuita.

El mismo paquete permite configurar el nodo de cabecera y las estaciones de trabajo.

Al seleccionar el tipo de nodo a incorporar, se determinan las herramientas necesarias.

Tanto en el nodo de cabecera como en el resto, deja instalado un grupo de herramientas para gestión del cluster.

Topologías

Microsoft utiliza este concepto para definir una serie de alternativas de interconexión. No tiene relación con el grafo de la red.

Red corporativa



Topología 1: los nodos de computación y las estaciones de trabajo se conectan solo a la red privada.

Red privada

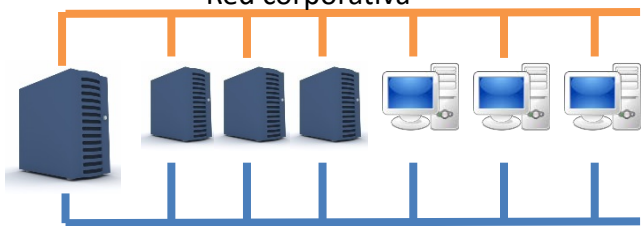
Red corporativa



Topología 1: los nodos de computación se conectan a la red privada. Las estaciones de trabajo se conectan a la red corporativa.

Red privada

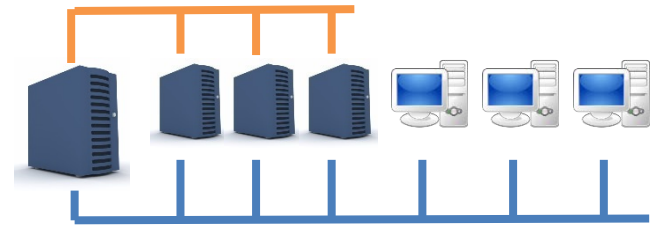
Red corporativa



Topología 2: los nodos de computación se conectan a la red privada y corporativa. Las estaciones de trabajo a ambas o solo a la privada.

Red privada

Red corporativa



Red privada

Red corporativa



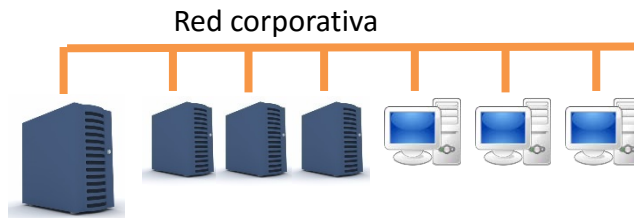
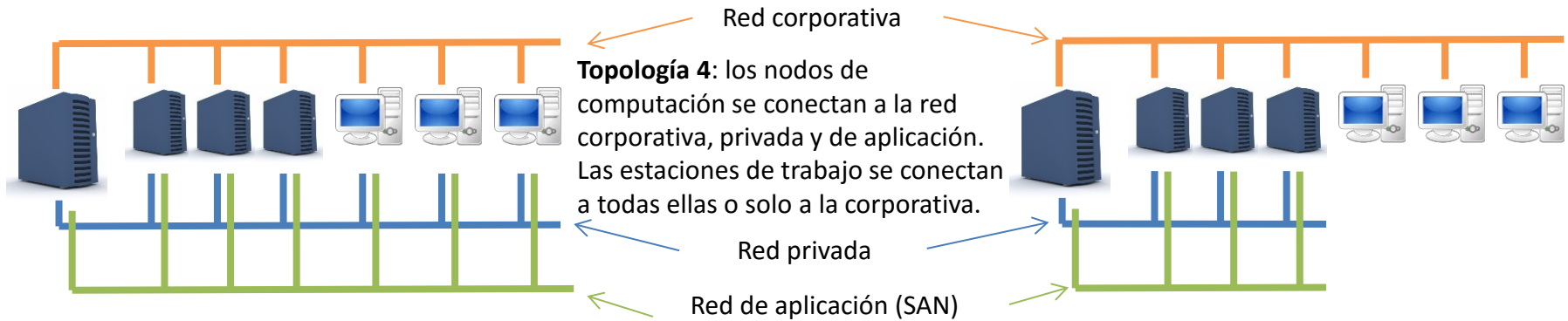
Topología 3: los nodos de computación se conectan a la red privada y a la de aplicación. Las estaciones de trabajo se conectan a ambas o a la red corporativa.

Red privada

Red de aplicación (SAN)



Topologías II



Topología 5: no hay red privada ni de aplicación. Todos los nodos se conectan a la red corporativa.



El usuario no está obligado a ajustarse a la configuración real de la red. Puede existir una conexión a la red corporativa pero no ser interesante que las estaciones de trabajo la utilicen.

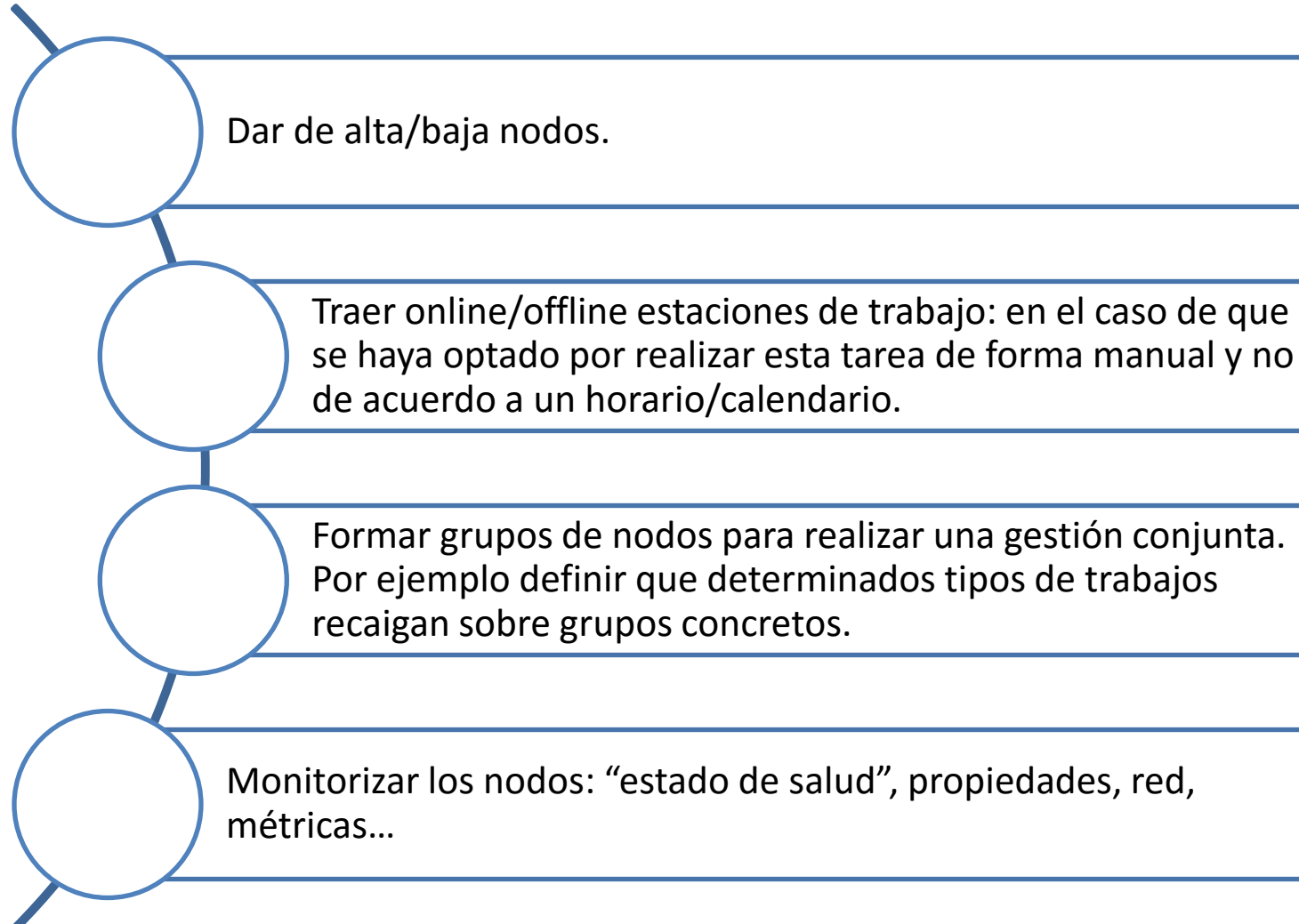
Gestión del cluster

El administrador del sistema debe realizar una serie de tareas de configuración y supervisión del cluster.

El Gestor del Cluster (Cluster Manager) es la herramienta que le permite realizar estas operaciones:

- Gestión de nodos.
- Gestión de usuarios.
- Configuración del planificador.

Gestión de nodos



Dar de alta/baja nodos.

Traer online/offline estaciones de trabajo: en el caso de que se haya optado por realizar esta tarea de forma manual y no de acuerdo a un horario/calendario.

Formar grupos de nodos para realizar una gestión conjunta. Por ejemplo definir que determinados tipos de trabajos recaigan sobre grupos concretos.

Monitorizar los nodos: “estado de salud”, propiedades, red, métricas...

Gestión de usuarios

Dar de alta/baja nuevos usuarios.

Gestionar grupos de usuario.

Asignar roles de usuarios:

- **Usuario:** puede administrar sus propios trabajos.
- **Administrador:** puede administrar trabajos y recursos.
- **Administrador de trabajo:** puede administrar trabajos , pero no recursos.
- **Operador de trabajo:** puede administrar trabajos de una manera restringida (ver, cancelar, finalizar, re-encolar).

Gestión del planificador

El planificador decide qué trabajos lanzar en cada momento y contra qué recursos.

Lo hace de acuerdo a una serie de políticas:

- Encolada.
- Balanceada.

Permite la utilización de “backfilling”: adelantar trabajos sobre uno más prioritario que está en espera de obtener los recursos que necesita, siempre que éste no se vea retrasado por ello.

Políticas de gestión de trabajos

Subpolíticas

Precedencia: trabajos más prioritarios pueden detraer recursos de trabajos menos prioritarios ya en curso.

- **Elegante:** detraer recursos de tareas que terminan.
- **Inmediata:** cancelar todas las tareas en ejecución.
- **A nivel de tarea:** cancelar tareas individualmente.

Asignación dinámica de recursos: los recursos asignados a un trabajo se pueden alterar en tiempo de ejecución.

- **Incremento automático:** priorizar la asignación de nuevos recursos sobre el inicio de trabajos menos prioritarios.
- **Decremento automático:** detraer recursos no usados por trabajos que no tienen tareas pendientes.

Encolada

- Intenta arrancar los trabajos en orden de llegada.
- Optimiza trabajos voluminosos.
- Precedencia elegante por defecto.
- Habilitados incremento y decremento automáticos por defecto.

Políticas

Balanceada

- Inicia los trabajos tan pronto como disponen del mínimo de recursos solicitados.
- Nuevos recursos disponibles son asignados a los trabajos en curso de acuerdo con su prioridad.
- Optimiza tareas cortas e interactivas.
- Precedencia inmediata por defecto.

Gestor de Trabajos



Los usuarios envían trabajos al cluster.

Los trabajos están formados por tareas.

El Gestor de Trabajos (Job Manager) es la herramienta que permite crear y enviar trabajos.

Tipos de tareas

Básica

- Ejecuta una única instancia de una aplicación serie o paralela.

Barrido paramétrico

- Ejecuta múltiples instancias de una aplicación.
- El número exacto lo determina un parámetro que toma valores diferentes.
- El valor del parámetro tiene un significado para cada instancia.

Preparación de nodo

- Un comando o “script” para ser ejecutado en cada nodo.
- Se ejecuta antes que cualquier otra tarea del trabajo.

Liberación de nodo

- Un comando o “script” para ser ejecutado en cada nodo.
- Se ejecuta cuando el nodo es desasignado al trabajo.

Servicio

- Ejecuta un comando o servicio en cada recurso asignado al trabajo.
- Si una instancia del comando finaliza y el recurso sigue asignado, una nueva instancia se inicia.

Tipos de trabajos

Trabajo

- Es el tipo más general.
- Los otros dos pueden ser definidos a partir de él también.
- Estos trabajos pueden incluir múltiples tareas y barridos.

Single-task job

- Pensado para facilitar la configuración de trabajos simples.
- Se utiliza para trabajos de una única tarea.

Parametric sweep job

- Pensado para facilitar la configuración de trabajos simples.
- Se utiliza para trabajos con una única tarea paramétrica.

Propiedades de los trabajos

ID del trabajo

- Identificación numérica del trabajo.
- Asignada por el Job Manager.

Nombre del trabajo

- Identificación textual del trabajo.
- Asignada por el usuario.

Plantilla de trabajo

- Nombre de la plantilla bajo la que se quiere configurar el trabajo.
- Las plantillas establecen valores y restricciones por defecto y son creadas por el administrador.

Prioridad

- Valor numérico que establece la prioridad del trabajo.
- Varía de 0 a 4000, siendo 4000 la prioridad más alta.

Tiempo de ejecución

- Tiempo máximo para completar el trabajo.
- Si se supera el trabajo es cancelado.

Memoria

- Mínima cantidad de memoria necesaria para que un nodo pueda ser asignado.

Licencias

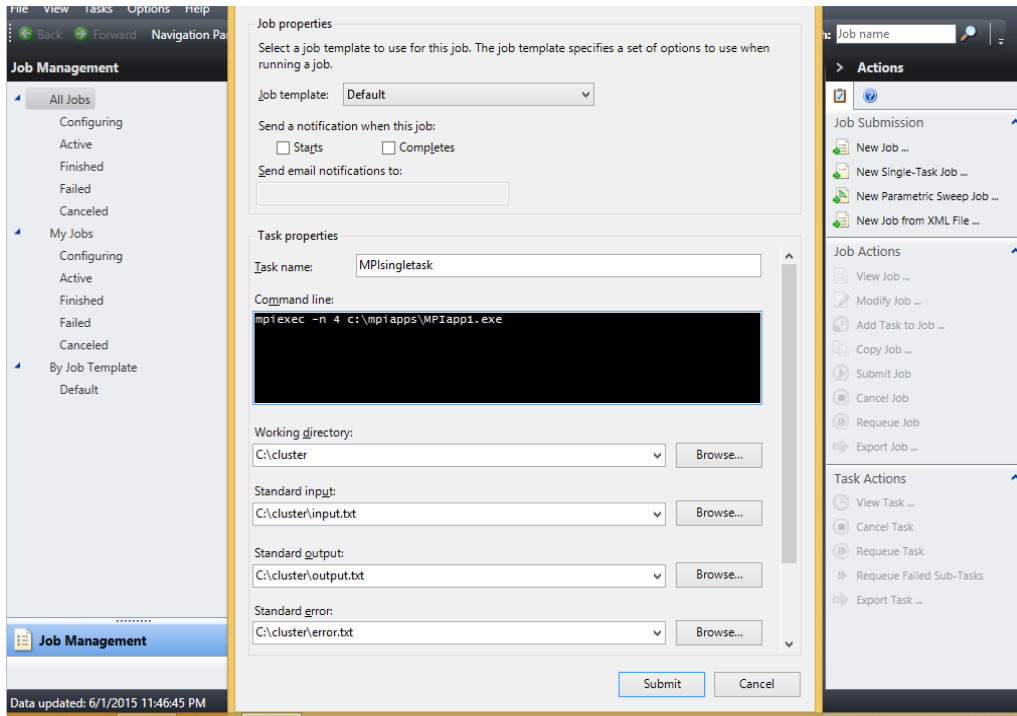
- Licencias solicitadas por el trabajo.

Aparte del “nombre”, no es obligatorio establecer el resto de propiedades.

Hay varias propiedades más. Para una lista completa consultar:

<https://technet.microsoft.com/es-es/library/ff919649>

Envío y monitorización.



Envío de trabajos desde el Job Manager.

Los trabajos pueden ser creados y enviados desde un nodo a través del Job Manager.

Elige el tipo de trabajo y rellena los campos.

Define la línea de comandos y los ficheros de salida.

Los valores por defecto se pueden respetar si no son necesarios otros.

Los trabajos se pueden enviar remotamente. Para ello se debe establecer una conexión de escritorio remoto a uno de los nodos.

Una referencia completa de cómo hacer esto está disponible en:

<https://technet.microsoft.com/en-us/library/gg315415%28v=ws.10%29.aspx>

Planificación encolada.

Scheduling mode:

- Queued - Attempt to assign the maximum amount of requested resources to running jobs.
- Balanced - Start as many jobs as possible with the minimum amount of requested resources for each. If additional resources are available on the cluster, grow jobs based on their priority and the Priority Bias setting.

Pre-emption options

- Graceful pre-emption - To enable higher priority jobs to start sooner, take resources away from lower priority jobs as their tasks complete.
- Immediate pre-emption - To enable higher priority jobs to start sooner, take resources away from lower priority jobs by canceling running jobs.
 - Task level pre-emption - To enable pre-emption of individual tasks instead of whole jobs.
- No pre-emption - Jobs will continue to run until completion, even if higher priority jobs are waiting for resources.

Adjust resources automatically

- Increase resources automatically (grow) - Use available resources to grow higher priority, running jobs to their maximum before starting lower priority jobs.
 - Grow by pre-emption - To help grow higher priority, running jobs, use pre-emption to take resources from lower priority, running jobs.
- Decrease resources automatically (shrink) - Automatically release unused job resources over time when a job holds resources that it cannot use.

[More about policy configuration](#)

Intenta asignar el mayor número de recursos posibles a los trabajos.

Es posible configurar opciones de precedencia.

También se puede decidir cómo se ajustan los recursos asignados a los trabajos.

Principales opciones de planificación encolada.

Planificación balanceada.

Scheduling mode:

- Queued - Attempt to assign the maximum amount of requested resources to running jobs.
- Balanced** - Start as many jobs as possible with the minimum amount of requested resources for each. If additional resources are available on the cluster, grow jobs based on their priority and the Priority Bias setting.

Pre-emption options

- Immediate pre-emption (Recommended)** - To enable additional jobs to start, take resources away from running jobs by canceling running tasks
- Graceful pre-emption (Advanced) - To enable additional jobs to start, take resources away from running jobs as tasks exit
- i** For most cluster workloads, immediate pre-emption in Balanced mode enables more jobs to start in a time period.

Priority bias

Priority Bias controls how additional resources are allocated to running jobs. A higher bias level allocates more resources to higher priority jobs.

Priority Bias level:

- High bias
- Medium bias**
- No bias

Rebalancing interval

The job scheduler rebalances resource allocation at a constant time interval. Jobs can grow and shrink in order to start new jobs, fill available resources, and balance resource allocation according to the Priority Bias level.

Seconds between rebalancing:

[More about policy configuration](#)

Intenta arrancar tantos trabajos como sea posible aunque no pueda asignarles el máximo de recursos.

Es posible configurar opciones de precedencia.

También se puede ajustar el sesgo aplicado a la prioridad a la hora de asignar nuevos recursos disponibles.

El periodo de rebalanceo de recursos es también configurable.

Principales opciones de planificación balanceada.

Referencias

- Cluster network topologies: <https://technet.microsoft.com/en-us/library/gg145543.aspx>
- Node management: <https://technet.microsoft.com/es-es/library/ff919378>
- Managing cluster users: <https://msdn.microsoft.com/en-us/library/ff919335.aspx>
- New cluster user roles:
<http://blogs.technet.com/b/windowshpc/archive/2013/09/04/hpc-pack-2012-sp1-available.aspx>
- Job scheduler: <https://technet.microsoft.com/es-es/library/ff919436>
- Job manager: <https://technet.microsoft.com/es-es/library/ff919691>