

File-backed pages correspond to segment of archivo sobre un disco. Si no contiene archivos que hayan sido escritos recientemente y no hayan regresado al almacenamiento persistente, estos se recuperarian de manera sencilla.

Como regla general, pedir anonymous pages es mas caro que pedir file-backed pages.

Anonymous pages, estos mantienen los datos generados en run-time y usados por un proceso, pedir una pagina anonima, requiere reescribir el contenido swap.

File-backed pages puede realizar operaciones de lectura y escritura desde el almacenamiento persistente a corto o largo plazo, y son almacenadas secuencialmente.

Anonymous pages tienden a disponerse a lo largo del area swap, por lo tanto para un disco que necesite lotar para leer datos, buscar anonymous pages resultaria demasiado caro ya que estos estan dispersas aleatoriamente.

LLEGADA DE SSDs. con esto el swap vuelve a ser interesante debido a las velocidades de acceso aleatorio, aleandadas por estos dispositivos. Aqui el swap vuelve a ser interesante para manejar el overflow y ademas como una extension de memoria para optimizar el balance entre las paginas cache y anonymous hasta una carga moderada.

Reconsidering Swapping

¿Cuando usar swap?

Con esto se denota elswappiness que va de 0 a 100 donde 0 no reclama swap en absoluto, si realmente no es necesario y 100 que reclama de forma igualitaria file-backed pages y anonymous pages. Actualmente Johanness incrementa a 200 el swappiness debido a las altas velocidades que existen actualmente.

Costo de Rotaciones: en cada LRU (last-recent used) el MM (Memory manager) hace lo posible por no reclamar paginas que estan en uso activo, esto lo hace pasando ocasionalmente por la lista y limpiando los bits de referencia de cada pagina, asi, las paginas que vuelven a ser usadas se les volvera a asignar un bit y las que no se les asignen estaran en elbaso. Por lo tanto estas paginas son las primeras en ser reclamadas las paginas que fueron referenciadas de nuevo se añadiran a la cabeza de la lista y pasaran un tiempo antes de que se reclamen de nuevo.

Johannes cambia el sistema de decision para saber si se elegira entre una anonymous page o file-backed pages con lo cual se introduce el concepto de costo por reclamar una pagina. la cual se elige de acuerdo al que tenga el menor costo.

Swapping es cada vez mas atractivo, asi como los dispositivos de almacenamiento, en especial los SSDs. Anteriormente mover una pagina, o remover una pagina del dispositivo de almacenamiento era una tarea sumamente lenta, pero ahora la memoria persistente ha alcanzado niveles de velocidad cercanos a los de acceso a memoria directa en particular para el compute en la nube el overcommitting memory seria insignificante, permitiendo un mejor manejo del sistema entero.

Existe aún otro problema con el radix-tree el cual nos lleva a otro global-lock el cual hace que volver a pedir por el CPU. Para solucionar el "adversus-space" es dividido cada GBMS de swap asignando un reloj a cada página. Esto reduce la contención por un reloj individual.

Aún así hay sobrecarga adquiriendo el bloqueo y puede haber contención en cache-line cuando se esta accediendo al bloque en cualquier otro CPU cluster. Para minimizar el costo se agrega una nueva interfaz para asignar o liberar paginas en conjuntos, cuando un CPU tiene un conjunto de paginas swap puede usarlas sin tomar el bloqueo local del cluster. Paginas liberadas son acumuladas en un cache separado y son regresadas en later.

En 2013 el código de swap fue renovado para un mejor desempeño en SSDs ya que existían dos problemas.

No hay demora de búsqueda en SSDs

Wear-leveling mejora espurciendo los datos al rededor del disco, a diferencia de los discos rotacionales que debían mantener cerca de el para una lectura rapida.

En kernels actuales el dispositivo swap es representado por "swap_info_struct" y dentro de este se encuentra "swap_map" donde se apunta al byte del arreglo donde contiene el contador de referencia para cada pagina almacenada en swap.

En el mismo actualización de 2013 cuando el sistema detectaba que estaba trabajando con un SSDs, este se dividía en clusters.

Un apuntador "per-cpu-cluster" apunta a diferentes clusters de cada CPU. con este, cada CPU puede alojar paginas de swap en su cluster, con el resultado de que estas alojaciones son espaciadas al rededor del dispositivo. (Aunque en la actualidad no se ha alcanzado el potencial esperado).

El problema es que algunas veces se tiene que lidiar con el bloqueo ya que las CPU no tiene acceso exclusivo a cualquier cluster. así que deben adquirir el "lock" de spinlock en el "swap_info_struct" antes de cualquier cambio hecho. En sistemas tipicos solo hay un dispositivo swap, por lo que si hay demasiado swapping el spinlock es muy llamado. Y la contención del spinlock no es el curulo a la escalabilidad.

Making swapping Social

La contención no es necesaria cada cluster es independiente y puede ser adquirido sin tocar a los demás así que no es necesario un bloque global. Para esto se agrega un nuevo bloqueo a una entrada al array de "cluster_info" un solo bit de bloqueo es usado para minimizar el consumo de memoria agregada. Así, cualquier cluster puede asignar paginas (o paginas libres) del cluster sin pelear con las demás.