

[snip]

>>>>> We need to decide whether we want to do per-container memory
>>>>> limitation via these data structures, or whether we do it via
>>>>> a physical scan of some software zone, possibly based on Mel's
>>>>> patches.
>>>>> why not do simple page accounting (as done currently
>>>>> in Linux) and use that for the limits, without
>>>>> keeping the reference from container to page?
>>>> As I've already answered in my previous letter simple
>>>> limiting w/o per-container reclamation and per-container
>>>> oom killer isn't a good memory management. It doesn't allow
>>>> to handle resource shortage gracefully.
>>> per container OOM killer does not require any container
>>> page reference, you know `_what_` tasks belong to the
>>> container, and you know their `_badness_` from the normal
>>> OOM calculations, so doing them for a container is really
>>> straight forward without having any page 'tagging'
>> That's true. If you look at the patches you'll
>> find out that no code in oom killer uses page 'tag'.
>
> so what do we keep the context -> page reference
> then at all?

We need this for

1. keeping page's owner to uncharge to IT when page goes away. Or do you propose to uncharge it to current (i.e. ANY) container like you do all across Vserver accounting which screws up accounting with pages sharing?
2. managing LRU lists for good reclamation. See Balbir's patches for details.
3. possible future uses - correct sharing accounting, dirty pages accounting, etc

>>> for the reclamation part, please elaborate how that will
>>> differ in a (shared memory) guest from what the kernel
>>> currently does ...
>> This is all described in the code and in the
>> discussions we had before.
>
> must have missed some of them, please can you
> point me to the relevant threads ...

lkml.org archives and google will help you :)

> TIA,
> Herbert

Containers mailing list

Containers@lists.osdl.org

<https://lists.osdl.org/mailman/listinfo/containers>
