

Avis d'expert : quels indicateurs pour suivre son infrastructure virtualisée ? (1)

Lors d'un [précédent article](#), je vous ai présenté les quatre changements majeurs dans la manière de dimensionner une infrastructure virtuelle par rapport à son homologue physique. Ces quatre différences se situent au niveau des processeurs, de la mémoire, du stockage et du réseau. Dans la continuité de la conception, il y a la vie de l'infrastructure, c'est-à-dire l'exploitation au quotidien. Est-ce que l'on supervise au quotidien une infrastructure virtuelle comme on le fait en environnement physique ? La mutualisation de nombreuses ressources implique la prise en considération d'autres paramètres que nous allons détailler ici.

Le suivi de la performance des processeurs

Même si la notion de processeur virtuel (vCPU) pour un hyperviseur (VMWare ou Microsoft) est proche de la notion de cœur physique, **un processeur virtuel** est beaucoup moins puissant (indice de puissance) qu'un processeur physique. La charge processeur visible au sein d'un serveur virtuel pourra donc être implicitement plus importante qu'en environnement physique. A cela rien d'inquiétant si le serveur a été dimensionné pour répondre aux pics de charge. En revanche, une conservation des seuils d'alerte de charge n'a souvent plus de sens et il est judicieux de les ajuster.

Tous les hyperviseurs ne sont pas forcément identiques au sein d'une ferme. Les différences peuvent se situer au niveau des processeurs (génération, fréquence, cache) ou d'autres caractéristiques techniques. C'est un aspect à prendre en compte car un serveur virtuel peut être migré à chaud d'un hyperviseur à un autre (vMotion chez VMWare, Live Motion chez Microsoft). Après un déplacement à chaud d'un serveur virtuel, le taux d'utilisation du processeur peut alors varier. Est-ce que **les seuils d'alerte** qui ont été définis prennent en compte ces changements de contexte ?

J'ai expliqué dans l'article précédemment cité le rôle de l'ordonnanceur dans l'accès aux processeurs physiques par les serveurs virtuels. Ces accès ne sont pas immédiats. Plus il y a de serveurs virtuels, plus il y a d'accès et plus le temps d'attente (latence) est important. Le suivi des performances ne se limite pas uniquement à identifier la manière dont la puissance est utilisée, mais aussi à pouvoir détecter quand elle est disponible. Il est possible d'imaginer, dans des cas extrêmes, un serveur virtuel avec une faible charge processeur, mais lent. Cet indicateur de latence processeur, propre à chaque serveur virtuel, est révélateur d'une bonne utilisation de la puissance disponible. Il ne faut pas croire que la puissance augmente en ajoutant des processeurs virtuels au serveur, c'est en réalité plus complexe et c'est souvent l'effet contraire de ce qui est recherché qui se produit. Nous adressons généralement cette problématique par une analyse du nombre total de processeurs virtuels sur chaque hyperviseur.

Le suivi de l'utilisation de la mémoire

Si la mémoire est partagée entre les serveurs virtuels exécutés au sein d'un même hyperviseur, il convient de distinguer mémoire utilisée et non utilisée. La mémoire occupée est généralement allouée statiquement au serveur virtuel concerné alors que la mémoire non utilisée est mutualisée. C'est d'ailleurs pour cette raison qu'il est possible d'exécuter des serveurs virtuels dont le total de mémoire est supérieur à la capacité mémoire d'hyperviseur. Le **sur-provisionnement mémoire** des serveurs virtuels sur l'hyperviseur n'est pas anodin. C'est une sorte de prise de risque, de pari que tous les serveurs ne vont pas utiliser leur mémoire au même moment. Il n'y a pas de considération de type « c'est bien » ou « c'est mal », cela dépend beaucoup du dimensionnement des serveurs virtuels (au juste besoin ou confortable). Cependant, surveiller permet d'éviter que cela ne devienne une source de dégradation des performances.

Une première conséquence importante est que le sur-provisionnement empêche le démarrage de tous les serveurs virtuels en même temps. En effet, un hyperviseur vérifie qu'il peut affecter la totalité de la mémoire d'un serveur virtuel avant de le démarrer. Il est toutefois possible de les démarrer avec une temporisation afin que chacun libère sa mémoire non utilisée, dans un ordre précis.

Le ralentissement d'accès à la mémoire est une seconde conséquence. VMWare a implémenté une méthode de **recupération mémoire** (appelée 'ballooning') auprès des serveurs virtuels. Elle se déclenche lorsque l'hyperviseur doit fournir de la mémoire à un serveur virtuel et que les capacités disponibles sont insuffisantes. L'hyperviseur force la libération de mémoire auprès des serveurs virtuels. Cette mémoire libérée est alors distribuée aux serveurs selon les besoins. Ce mécanisme n'est pas immédiat, il y a un peu de latence entre la demande de mémoire et son allocation effective. C'est la raison pour laquelle les ralentissements sont perceptibles. Cela nuit à la bonne performance des serveurs.

Autre conséquence, le ballooning peut aussi changer totalement le comportement mémoire des serveurs. Les systèmes Linux sont réputés pour utiliser la mémoire disponible et la conserver en cache lorsqu'elle est libérée par les processus (mémoire qui n'apparaît alors pas comme libre). Le taux d'occupation mémoire d'un tel serveur est plutôt stable et élevé. Ce taux va diminuer avec le ballooning et ne sera pas stable (des libérations fréquentes puis des augmentations).

J'ai évoqué, toujours dans l'article précédemment cité, qu'il existe des cas où la mémoire utilisée peut être mutualisée entre les serveurs virtuels. Cela apparaît lorsque le ballooning ne suffit plus. Les performances sont alors grandement dégradées, c'est une situation à proscrire.

Le suite sur la seconde partie (2) de notre avis d'expert, avec le suivi des performances du stockage, le suivi des performances du réseau, et d'autres facteurs induits.

En complément :

Un autre avis d'expert de Sébastien Truttet : [Pour un meilleur dimensionnement des applications et des infrastructures](#)

Voir aussi

[Silicon.fr étend son site dédié à l'emploi IT](#)

[Silicon.fr en direct sur les smartphones et tablettes](#)