



LE GUIDE DÉFINITIF DE

la VDI sur une infrastructure hyperconvergée

NUTANIX™

Brian Suhr [auteur] a plus de vingt ans d'expérience en IT dans la conception, la mise en œuvre et l'administration des infrastructures d'entreprise. Il a fourni une expertise en architecture et en ingénierie dans divers projets de virtualisation, de datacenter et de cloud tout en travaillant avec des équipes techniques hautement qualifiées dans des environnements internationaux. En tant que contributeur des blogs DataCenterZombie et VirtualizeTips, Brian se concentre sur la création de contenu axé sur la virtualisation, l'automatisation, l'infrastructure et la promotion de produits et services qui bénéficient à la communauté technologique. Suivez Brian sur Twitter : [@bsuhr](#)

Sachin Chheda [expert] est le directeur des solutions et du marketing vertical chez Nutanix. Évoluant dans le secteur des technologies de l'information depuis longtemps, il a occupé des postes en ingénierie, gestion et marketing dans les entreprises les plus innovantes du secteur - développant et commercialisant des produits et solutions destinés aux entreprises les plus prestigieuses et visionnaires du monde. Suivez Sachin sur Twitter : [@StorSC](#)

Publication officielle version « 1.0 »

Copyright 2017 Nutanix, Inc. Tous droits réservés. Ce produit est protégé par les lois américaines et internationales sur le droit d'auteur et la propriété intellectuelle. Nutanix est une marque commerciale de Nutanix, Inc. aux États-Unis et/ou dans d'autres juridictions. Toutes les autres marques et noms mentionnés dans ce document peuvent être des marques de commerce de leurs sociétés respectives.

| | |
|----------------------|---|
| Auteur | 2 |
| À propos de ce guide | 4 |
| Introduction : | 5 |

PRINCIPES ARCHITECTURAUX

| | |
|-----------------------|----|
| Point d'entrée | 6 |
| Évolutivité | 7 |
| Performances | 8 |
| Capacité | 9 |
| Surveillance | 10 |
| Blocs de construction | 12 |

ALTERNATIVES D'INFRASTRUCTURE

| | |
|----------------------------------|----|
| Build Your Own | 14 |
| Infrastructure convergée | 16 |
| Infrastructure hyperconvergente | 18 |
| Exigences en matière de stockage | 20 |

TYPES DE STOCKAGE

| | |
|--------------------------------------|----|
| Architectures multi-niveaux héritées | 22 |
| 100 % flash | 23 |
| Flash hybride | 23 |

| | |
|---|----|
| Dimensionnement de la couche de calcul | 24 |
| Conception de cluster de virtualisation | 28 |
| Premiers pas | 30 |

| | |
|---------------------|----|
| À PROPOS DE NUTANIX | 32 |
|---------------------|----|

À PROPOS DE CE GUIDE

Ce livre est axé sur la conception d'infrastructure pour la VDI et les environnements EUC (End-User Computing). Le contenu de ce guide a été tiré du chapitre consacré à l'infrastructure du livre « Architecting and Designing End-User Computing Solutions », à paraître prochainement.

INTRODUCTION :

Après avoir choisi la bonne stratégie et le bon fournisseur de logiciel de services et d'applications EUC, les choix d'infrastructure constituent la prochaine décision cruciale pour les projets de virtualisation d'applications et de postes de travail.

L'infrastructure de calcul et de stockage est la base sur laquelle il est possible de construire des services. Comme pour les services électriques et l'eau courante, nous comptons sur eux, et nous nous attendons à ce qu'ils fonctionnent simplement en tournant un robinet ou en actionnant un interrupteur.

Sans une infrastructure stable, hautement disponible et performante au cœur de la conception, le département IT sera confronté à un certain nombre de défis supplémentaires lors des phases de déploiement et d'exploitation de votre projet EUC. Cela renforce le constat que l'infrastructure est cruciale, mais ne justifie pas de dépenser une grande partie des ressources du département IT. Les architectes et les ingénieurs doivent se concentrer sur la livraison de services et d'applications EUC, plutôt que de gérer la maintenance de l'infrastructure.

Plusieurs facteurs importants doivent être pris en compte dans le processus de conception des infrastructures EUC. En combinant ces facteurs avec les exigences de l'organisation, vous pourrez évaluer plus efficacement les alternatives d'architecture. Dans les projets EUC, les facteurs à prendre en compte lors de l'évaluation des alternatives architecturales et des options de fournisseurs sont les suivants :

- Point d'entrée
- Évolutivité
- Performances
- Surveillance
- Capacité

POINT D'ENTRÉE

Le point d'entrée ou de départ de l'infrastructure est souvent une décision cruciale pour la réussite d'un projet. Il s'agit de la quantité d'infrastructure et des coûts requis pour démarrer le déploiement de la virtualisation d'applications/de postes de travail, et la distribution des implémentations, en fonction de la taille de ces points de départ.

Si le projet est prévu pour atteindre 10 000 utilisateurs une fois la phase initiale de 5 000 utilisateurs pleinement déployée, l'entreprise sera probablement moins susceptible d'être choquée par les coûts initiaux. Le raisonnement est le suivant : selon le type d'infrastructure sélectionné, le coût par utilisateur peut ne pas avoir de sens avant que vous ayez déployé quelques milliers d'utilisateurs.

Le revers de la médaille, c'est que si une organisation veut servir 10 000 utilisateurs, mais n'a l'intention que de commencer avec 500 utilisateurs (qui augmenteront par la suite à un rythme régulier au cours du projet), elle aura tendance à examiner attentivement le coût initial de déploiement de l'infrastructure de cette taille plutôt que de franchir une première étape plus importante. À cette taille, le coût par utilisateur peut rester stable à mesure que l'environnement évolue ou peut sembler biaisé à première vue, en raison de l'augmentation des dépenses initiales d'infrastructure.

Bien que le coût par utilisateur puisse être considéré comme vague et presque non pertinent pour déterminer les coûts d'infrastructure, il sera abordé lorsque vous tenterez de vendre le projet au reste de l'entreprise ou de justifier votre choix d'infrastructure à la direction. Si vous choisissez une alternative dont le coût par utilisateur initial est plus élevé, vous devez être prêt à expliquer en détails les raisons de votre choix. Évaluez les solutions qui, selon vous, conviendraient mieux à votre environnement. Sinon, soyez prêt à définir la décision sur la manière dont les coûts vont se jouer. Un exemple de ces deux scénarios est présenté dans la Figure 1.

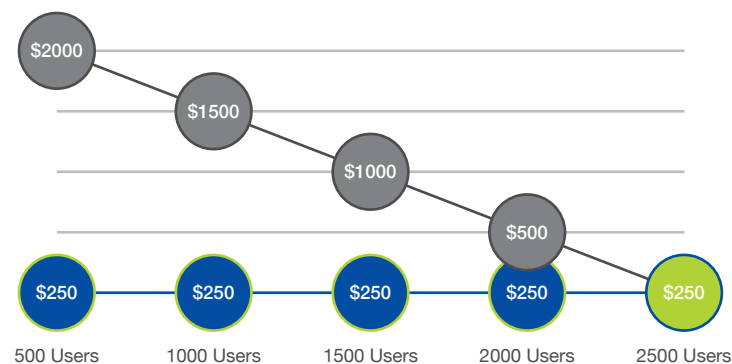


Figure 1 :
Points d'entrée par poste de travail

ÉVOLUTIVITÉ

L'évolutivité de l'architecture est un facteur important dans l'évaluation de la viabilité d'un projet. Un architecte devra comprendre les options de dimensionnement initial des différentes alternatives ; ce qui nous ramène au sujet du point d'entrée. L'alternative pourra-t-elle facilement permettre un démarrage à petite échelle si besoin, ou l'organisation devra-t-elle acheter plus de capacité d'infrastructure que nécessaire pour répondre à la taille initiale d'un projet, sans être capable de l'utiliser pleinement jusqu'à ce que le projet prenne de l'ampleur ?

Outre la petite échelle à laquelle l'alternative peut commencer, il est tout aussi important de considérer jusqu'à quelle échelle l'alternative peut s'étendre. Si vous souhaitez commencer à 500 utilisateurs et pouvoir par la suite étendre à 10 000 utilisateurs, à quoi ressemblera l'alternative aux deux extrémités de ce spectre ? L'organisation sera-t-elle satisfaite des points les plus élevés, plus bas, ou des deux ?

L'argument de l'évolutivité n'est pas seulement valable pour le stockage, mais aussi pour la puissance de calcul, la mise en réseau et à éventuellement d'autres niveaux de la conception. Si des ajustements sont apportés à la configuration de la couche de calcul pour obtenir une densité plus faible de machines virtuelles par serveur hôte, comment cela pourrait-il affecter les différents choix de conception lors de la mise à l'échelle ? Par exemple, si la conception initiale de l'hôte commence avec 128 Go de mémoire par hôte et que l'objectif final est de 256 Go ou plus, il faudra vous assurer d'utiliser des modules DIMM de taille appropriée pour assurer l'évolutivité de la configuration à l'avenir. Si de mauvais choix sont faits au départ pour réduire les coûts, cela affectera la densité, ou coûtera plus cher à long terme en raison des modules DIMM non réutilisables.

L'architecte doit se concentrer sur la manière dont la solution pourra commencer à petite échelle, tout en étant capable d'évoluer au maximum. Mais il n'est pas question d'ignorer tous les points intermédiaires, car selon la manière dont le déploiement est dimensionné, il peut y avoir de nombreux points de mise à l'échelle entre le début et la fin du projet. Idéalement, recherchez une solution qui permettra à la conception de s'adapter facilement aux buckets d'utilisateurs identifiés dans le projet, sans dépasser le calendrier et les capacités de déploiement. La taille idéale de mise à l'échelle d'un bucket pour un projet peut être atteinte par ajouts progressifs de 100-200 utilisateurs. Mais si l'alternative architecturale choisie propose une mise à l'échelle plus vaste, vous devez comprendre que cela affectera les coûts et le déploiement.

PERFORMANCES

Les performances EUC mesurées par l'expérience de l'utilisateur final doivent toujours être soigneusement prises en compte. L'architecture choisie doit pouvoir répondre aux exigences à toutes les phases du projet. Selon les alternatives, le défi peut être difficile à relever : si une solution est dimensionnée à la baisse pour répondre au minimum aux exigences de l'utilisateur initial, cela peut entraîner une dégradation des performances si elle n'est pas évolutive de manière linéaire. Les architectes ne veulent pas faire de compromis sur l'architecture pour répondre à un petit point de départ susceptible d'affecter l'ensemble des options de performances maximales d'une solution. Passer plus de temps au départ pour prendre la bonne décision vous aidera à éviter des problèmes à l'avenir.

La conception d'une solution EUC comporte généralement de nombreuses exigences de performance différentes. Choisissez une alternative d'architecture suffisamment flexible pour répondre à toutes les exigences de performance en une seule option. Que la conception fournisse plusieurs types de services EUC ou se concentre uniquement sur la virtualisation d'applications et de postes de travail, les multiples besoins de performance doivent être pris en compte. Comprendre comment chaque solution alternative peut répondre ou non aux exigences de performance individuelle influencera grandement le processus d'évaluation et de conception.

CAPACITÉ

Le discours sur la capacité est assez similaire à celui sur la performance. Les projets EUC doivent répondre à plusieurs exigences en matière de capacité. La solution nécessitera l'exécution de serveurs VM, de postes de travail VM, d'applications, de profils utilisateurs et de données utilisateurs pour ce type d'architecture. Chaque niveau du projet peut avoir des besoins en capacité très différents : certains utilisent de grandes quantités de données qui sont généralement faciles à dédupliquer. D'autres, comme les profils et les données utilisateurs, consistent en de plus petites quantités de données compressibles par utilisateur, mais qui, multipliées par des milliers d'utilisateurs, s'avèrent au final très importantes.

Au cours des dernières années, l'achat d'une trop grande ou d'une trop petite capacité pour atteindre les niveaux de performance requis a été un gros problème. Il s'agit alors d'étudier soigneusement les alternatives architecturales pendant la phase de conception pour voir comment elles seront en mesure de fournir la capacité requise, tout en vous assurant que les exigences minimales de performance sont respectées. La solution alternative ne devrait pas fournir plus de deux à trois fois la capacité requise pour répondre aux exigences de performance de stockage ou ajouter une performance supplémentaire importante pour répondre aux exigences de capacité. La solution idéale est celle qui permet une flexibilité suffisante pour faire évoluer les performances et la capacité dans des proportions similaires, de sorte qu'aucune des deux ne soit trop éloignée de l'autre.

Par le passé, ce sujet a suscité beaucoup de discussions et de problèmes. De nombreuses organisations se sont heurtées à des difficultés de planification des performances et de la capacité en raison d'une mise à l'échelle de la capacité plus rapide que celle de

la performance. Ce n'est pas parce que la solution dispose de 5 To d'espace libre qu'elle peut évoluer en ajoutant 500 utilisateurs supplémentaires. Ce scénario peut entraîner une dégradation significative des performances. Les administrateurs et dirigeants informatiques qui n'ont pas une bonne compréhension de la façon dont une solution se met à l'échelle peuvent tomber dans ce piège.

SURVEILLANCE

La surveillance est très importante et souvent négligée. Lorsqu'il s'agit de surveiller l'infrastructure dans un environnement EUC, les administrateurs se concentrent généralement sur la performance. Ils doivent être capables de comprendre ce qui relève d'un fonctionnement normal et quand il y a un problème.

L'utilisation des systèmes de surveillance devrait être simple, tout en fournissant un grand nombre d'informations détaillées. Malheureusement, ce n'est pas le cas pour de nombreux fournisseurs, et vous devez donc examiner attentivement l'expérience de surveillance pour chaque alternative.

Une autre exigence est la capacité d'assurer une surveillance des performances au niveau de la machine virtuelle. Malheureusement, la majorité des fournisseurs d'infrastructure ne sont toujours pas en mesure d'offrir ce niveau de visibilité dans l'environnement de virtualisation. La capacité d'examiner rapidement la couche de stockage et de déterminer si son problème de performance se situe à l'échelle globale ou s'il est isolé au niveau d'un hôte, d'un groupe de VM, ou d'une seule VM est désormais incontournable.

Grâce à la gestion des performances de stockage au niveau de la VM, vous pouvez utiliser une approche similaire pour gérer les performances du CPU et de la mémoire d'une machine virtuelle au niveau de l'hôte. Les administrateurs doivent savoir si une machine virtuelle utilise temporairement des performances supplémentaires ou si elle consomme régulièrement des performances de stockage supérieures à celle des utilisateurs classiques. Cela vous aidera à comprendre quand il y a un pic et quand d'autres recherches sont nécessaires pour identifier le problème.

Un bloc de construction est un ensemble prédéfini d'infrastructures qui correspond à une quantité spécifique de ressources ou un nombre déterminé d'utilisateurs. Cette approche est l'un des meilleurs moyens d'aborder la conception d'infrastructures avec l'EUC.

En utilisant cette approche, il est possible de développer une architecture qui offre un modèle évolutif de capacité et de performance, avec des coûts prévisibles. Lorsque vous déterminez la taille de vos blocs de construction, vous choisissez les incréments dont vous avez besoin pour mettre à l'échelle le nombre d'utilisateurs et comment l'infrastructure sélectionnée pourra s'adapter à vos choix. Par exemple, vous voudrez peut-être échelonner le nombre d'utilisateurs par incréments de 50 à 100, sans que l'infrastructure de votre choix ne soit capable de gérer de si petits incréments. Cela peut forcer le projet à évoluer par incréments plus importants de 500 ou 1 000 utilisateurs à la fois. Si l'infrastructure que vous avez choisie est modulable par grands blocs, vous pouvez choisir de vous y adapter ou tout simplement accepter que les coûts d'infrastructure ne seront pas dimensionnés de la même manière que les blocs de déploiement utilisateur. Cela signifie simplement que l'entreprise achèterait l'infrastructure par blocs de 1 000 utilisateurs et ne l'implanterait que par groupes de 50 à 100 utilisateurs.

Ce critère rend les coûts des postes de travail virtuels ou des sessions utilisateurs plus élevés lors de l'achat du gros bloc pour déployer un plus petit nombre d'utilisateurs, mais cela se rééquilibre si l'organisation implémente tous les utilisateurs prévus.

Les architectures en blocs de construction sont utiles dans tout projet de conception. Mais les déploiements EUC ont toujours des blocs d'utilisateurs communs et des cas d'utilisation présentant des caractéristiques similaires et déployés en groupes. Pour continuer avec l'exemple d'un bloc de 100 utilisateurs, la compréhension des besoins en ressources de ces 100 utilisateurs permet de s'assurer que le bloc d'infrastructure est capable de fournir tout ce dont ces utilisateurs ont besoin.

Si chaque utilisateur a besoin de 15 IOPS en régime permanent, 30 Go de capacité de stockage, 2 Go de mémoire et 200 MHz de CPU, l'architecte sait alors que les blocs de construction doivent fournir 1500 IOPS, 3 To de capacité, 200 Go de mémoire et 20 GHz de CPU. L'architecte peut concevoir les blocs de construction de manière à contenir des ressources supplémentaires, mais aucune d'entre elles ne peut être inférieure à ces valeurs. En outre, nous ne voulons pas gaspiller les ressources en incluant inutilement trop de surplus qui ne peuvent être utilisés dans chaque bloc.

Avec cette approche et cette granularité dans la conception, il est désormais possible de dimensionner l'environnement en plus petits groupes de 50 à 100 utilisateurs. Cela permet une approche lente et régulière et fournit des valeurs prévisibles que les organisations peuvent planifier pour le déploiement, les performances, la capacité et les coûts. Si les organisations souhaitent s'étendre plus largement et plus rapidement, elles peuvent simplement insérer plusieurs blocs de construction à la fois.

Enfin, l'approche par blocs de construction s'est révélée particulièrement avantageuse, car la plupart des clients souhaitent commencer par des déploiements plus petits et évoluer progressivement par la suite. Le modèle « commencez petit et payez à mesure que vous grandissez », leur permet d'investir moins de capital au départ et d'acquérir de l'expérience à mesure que le déploiement progresse. La section suivante couvre les différents types d'architecture d'infrastructures disponibles aujourd'hui, et la façon dont chacun d'entre eux prend en charge ou non l'approche par blocs de construction.

Il existe actuellement trois alternatives principales d'architecture pour la virtualisation d'applications/de postes de travail ou, plus largement, pour les solutions EUC. Les alternatives sont : la construction autonome de votre propre infrastructure, appelée Build Your Own (BYO), l'infrastructure convergée (CI) et l'infrastructure hyperconvergée (HCI).

BUILD YOUR OWN

L'infrastructure BYO implique littéralement ce que son nom signifie : l'architecte ou les équipes choisissent de manière indépendante les produits qu'ils préfèrent ou considèrent comme les meilleurs sur le marché. Cette alternative se traduit par une augmentation significative de la période initiale de planification et de recherche, car l'équipe doit évaluer chaque produit séparément et déterminer comment ils peuvent ou non fonctionner ensemble.

Cette alternative permet également de sélectionner et de suivre une architecture de référence publiée par un fournisseur pour le type de solution en cours de construction. Ces architectures de référence sont généralement publiées par un seul fournisseur et se concentrent sur leur produit. Ces architectures de référence à faire soi-même (DIY) peuvent vous faire gagner du temps et réduire certains risques, mais elles ne sont pas toujours conformes à vos exigences de conception, à vos cas d'utilisation et à votre environnement.

Au minimum, une alternative BYO pour un projet basé sur l'EUC contiendra des ressources de calcul et de stockage. Il est possible que vous puissiez utiliser la connectivité réseau existante, donc ce composant peut ne pas être inclus dans cette alternative. La Figure 2 illustre un exemple simple des pièces constitutives d'une alternative BYO. Avec la flexibilité de la mise à l'échelle, les coûts sont assez prévisibles, exception faite du stockage. Selon la taille maximale du projet et le type de stockage choisi, il peut être nécessaire de disposer de plusieurs baies de stockage ou appliances. Au fur et à mesure que vous augmentez votre espace de stockage et que vous devez ajouter une nouvelle baie ou une nouvelle appliance, les coûts augmenteront substantiellement.

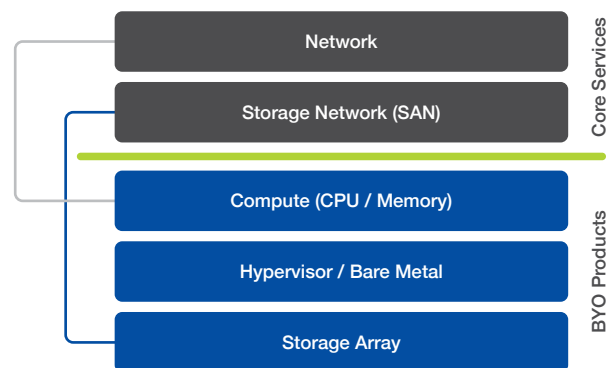


Figure 2 :
Construction autonome de votre infrastructure (Build your own - BYO)

Chaque fois que vous assemblez plusieurs produits du même fournisseur ou de plusieurs fournisseurs sans expérience préalable, le risque que vous prenez est plus grand. Il y aura un niveau d'incertitude quant aux performances et à la fiabilité de la solution jusqu'à ce que l'infrastructure réelle soit achetée et déployée de la manière décrite dans l'architecture.

Si vous êtes en mesure d'accepter les imprévus et des risques accrus, l'option BYO maximise la flexibilité puisqu'elle permet de choisir n'importe quel fournisseur et produit capable de fonctionner ensemble. De fait, vous pouvez continuer à travailler avec des fournisseurs avec lesquels vous êtes à l'aise, et changer de fournisseurs dans d'autres domaines.

L'alternative BYO est capable de mettre à l'échelle les ressources de calcul et de stockage indépendamment. Les seules limites à la méthode de mise à l'échelle ou de taille maximale seraient une contrainte sur le choix du produit individuel. Étant donné que les produits sont achetés séparément, il n'y a pas de minimum ni de quantité définie pour les produits. Cela permet de faire preuve de souplesse en essayant de tenir compte de l'approche par blocs de construction mentionnée ci-dessus.

INFRASTRUCTURE CONVERGÉE

L'infrastructure convergée (CI) est une architecture introduite sur le marché vers 2010. Les offres d'infrastructure convergée comprennent généralement les mêmes produits que ceux qui pourraient être sélectionnés dans le cadre de l'alternative BYO, et les regroupent dans une solution complète. Cela signifie qu'un fournisseur de CI inclura le calcul, le stockage et la mise en réseau dans son offre. En règle générale, la plupart des offres de CI contiennent des produits de plusieurs fournisseurs et les incluent dans une solution unique. Un fournisseur peut également proposer toutes les couches d'une offre de CI à partir de sa propre gamme de produits. La Figure 3 illustre un exemple simple d'une alternative d'infrastructure convergée.

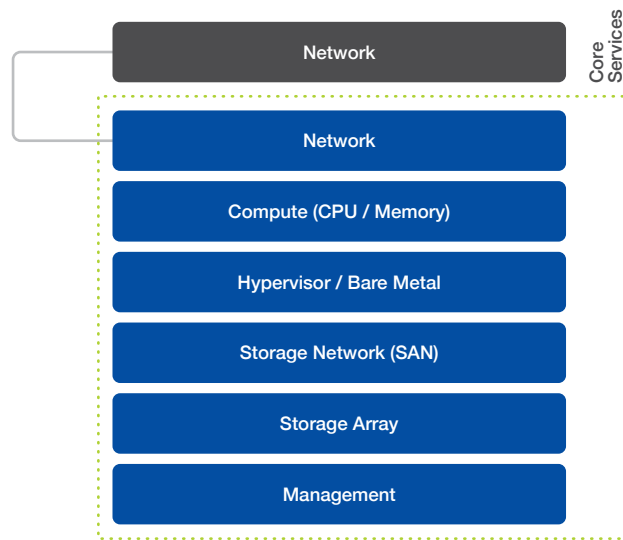


Figure 3 :
Infrastructure convergée

Une offre d'infrastructure convergée vous permettra d'acheter des produits que vous connaissez déjà, regroupés dans une solution unique. Cela peut également s'appliquer à une architecture de référence qui peut être achetée en tant que produit. En fonction du produit CI évalué, le produit peut, ou non, offrir une convergence supplémentaire par rapport aux produits achetés séparément dans une alternative BYO.

En règle générale, la plupart des fournisseurs et des produits CI offrent la possibilité d'acheter toutes les pièces d'infrastructure dans une seule unité de gestion des stocks. Le fournisseur de CI doit être en mesure d'offrir un support unifié pour l'ensemble de la solution de CI, ce qui signifie qu'il peut prendre en charge tous les produits de la solution. C'est un avantage supplémentaire, car il permet aux clients d'éliminer le besoin de travailler avec plusieurs fournisseurs en cas de problème.

La plupart des offres de CI proposent un nombre limité de produits dans leur solution. Cela permet au fournisseur de pré-tester et de valider toutes les pièces afin d'assurer qu'elles fonctionnent ensemble correctement, éliminant ainsi une grande partie du risque de l'option BYO.

Même plusieurs années après la mise sur le marché des produits CI, les fournisseurs ont peu fait pour simplifier leur gestion. Avec des offres de CI incluant les mêmes produits que les alternatives BYO, il est généralement possible de gérer les deux alternatives de manière similaire et dispersée. Cette alternative peut faire converger l'achat et/ou certains des produits, mais ne parvient généralement pas à faire converger la gestion opérationnelle quotidienne de la solution.

Un produit d'infrastructure convergée doit pouvoir faire évoluer les ressources internes indépendamment les unes des autres. Cela signifie que vous ne pouvez ajouter que de la puissance de calcul supplémentaire, même par incréments minimaux. Le stockage est une autre ressource évolutive dans une offre de CI, et il dépend fortement du type de solution de stockage choisi dans le cadre de l'offre de CI. Un produit d'infrastructure convergée aura une taille maximale, ce qui signifie qu'il aura une limite sur le nombre de serveurs supportés et une limite de stockage en fonction de la baie de stockage incluse.

Les limites de mise à l'échelle d'une offre de CI sont généralement assez grandes, mais elles atteindront leur maximum un jour si les ressources au sein du produit de CI sont mises à l'échelle. Pour continuer à faire évoluer la conception à ce stade, il faudra acheter un produit CI supplémentaire. Cela entraînera d'importants pics de coûts d'infrastructure à différents moments du processus de dimensionnement en fonction de la taille maximale de votre conception.

INFRASTRUCTURE HYPERCONVERGÉE

L'architecture hyperconvergée a été introduite sur le marché environ un an après la CI. Les véritables architectures hyperconvergées sont obtenues en faisant converger les ressources de calcul, de stockage et la couche de gestion en un seul produit. Il est possible de déployer une solution hyperconvergée dans une méthode d'architecture BYO ou de référence, mais pour être vraiment hyperconvergé, le produit doit inclure une appliance matérielle.

En incluant une appliance matérielle dans le produit, le fournisseur peut désormais inclure la gestion de l'infrastructure avec les autres ressources convergées dans le produit. La Figure 4 illustre un exemple simple d'une alternative d'infrastructure hyperconvergée.

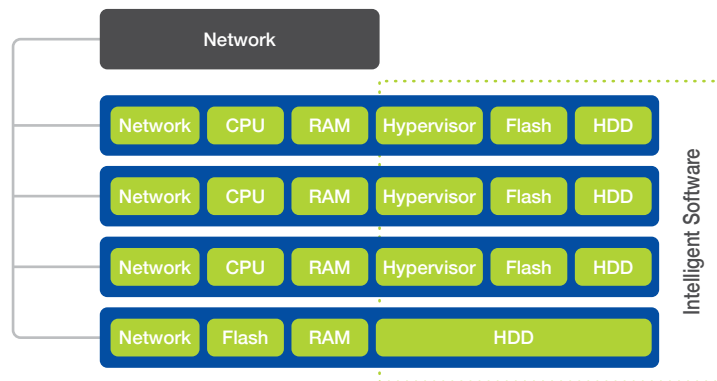


Figure 4 :
Infrastructure hyperconvergée

Un produit véritablement hyperconvergé offre un certain nombre d'avantages que d'autres architectures de référence sont incapables d'offrir :

Installation simple - Les principaux produits HCI doivent installer des nœuds en quelques minutes ou quelques heures tout au plus, et non pas en quelques jours ou semaines, grâce à un processus hautement automatisé.

Évolutivité facile - Le produit doit pouvoir être facilement mis à l'échelle. L'ajout de nouveaux nœuds à l'environnement doit se faire facilement et rapidement via l'interface de gestion.

Gestion moderne - Une interface de gestion moderne doit se concentrer sur la machine virtuelle (VM) en tant que point de gestion. Un administrateur doit être en mesure de comprendre les performances des machines virtuelles, la quantité de ressources consommée par chaque machine virtuelle, si des événements ou des erreurs spécifiques se produisent, et fournir la possibilité d'extraire facilement des rapports pour chaque VM.

Extensibilité - Vous devez être capable d'intégrer facilement l'infrastructure avec d'autres parties de la solution et de la contrôler via logiciel. Pour cela, le produit HCI doit proposer une API et éventuellement une autre méthode, telle que les commandes PowerShell. Avec une API, vous serez en mesure d'automatiser la communication et le contrôle entre les produits afin de réduire davantage les efforts et d'augmenter la précision de l'environnement.

Les performances ont été intentionnellement exclues de la liste des avantages de la HCI car tout le monde s'attend à ce qu'une solution hybride moderne ou basée sur Flash fonctionne correctement. La HCI consiste à créer une couche d'infrastructure simple et efficace. Elle permet aux équipes d'arrêter de perdre du temps à peaufiner l'infrastructure, et d'apporter une valeur ajoutée à l'entreprise aux niveaux de l'automatisation ou des applications.

Il existe de nombreux besoins différents en ressources de stockage pour tout projet EUC. Il faut prendre en compte les VM basées sur serveur, les données utilisateurs et l'infrastructure de postes de travail virtuels (VDI, ou machines virtuelles utilisateur). Les besoins de stockage associés seront les plus exigeants dans l'environnement, et ce sont aussi ceux qui causent le plus d'échecs de projets ou d'expériences en général.

Pour cette raison, la partie de cet eBook consacrée au stockage sera axée sur les besoins du service VDI de la solution choisie. Les besoins de chaque poste de travail virtuel peuvent souvent sembler petits et insignifiants, mais lorsqu'ils sont rassemblés en grands groupes pendant la mise à l'échelle du stockage, les exigences de performances peuvent facilement submerger un stockage qui n'a pas été correctement conçu pour répondre à ces besoins.

Si chaque poste de travail virtuel présente une moyenne de 15 IOPS avec une latence raisonnable et que 2 000 utilisateurs simultanés sont prévus, cela représente 30 000 IOPS, un nombre assez élevé qui pourrait surcharger la baie de stockage moyenne. Mais vous ne pouvez pas simplement concevoir la solution de stockage pour répondre à la moyenne des E/S de l'environnement : la conception doit prendre en compte les pics, y compris les démarrages des postes de travail et les connexions utilisateurs.

Une charge de travail de poste de travail virtuel est très différente des autres types de charges de travail exécutées dans un datacenter d'entreprise moyen, car les postes de travail virtuels peuvent souvent provoquer des pics d'E/S. Par exemple, l'ouverture d'une application comme Outlook pour la première fois peut générer jusqu'à 1 000 IOPS pour cette session utilisateur. C'est bien au-delà de la moyenne de 15 IOPS évoquée précédemment. La Figure 5 présente un exemple des différents impacts des IOP sur les applications.

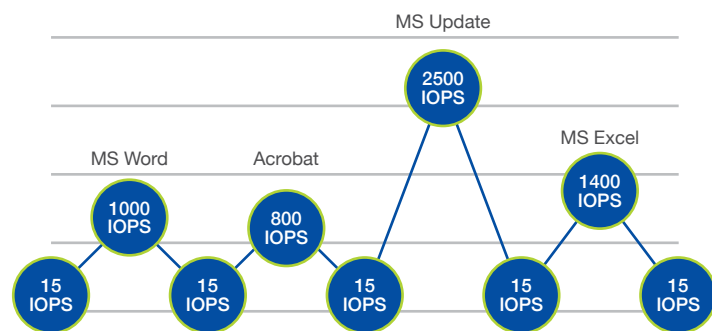


Figure 5 :
IOPS VDI

D'autres éléments opérationnels et de déploiement, tels que les correctifs et les mises à jour d'environnement, peuvent également créer des pics d'activité d'IOPS et affecter les performances s'ils ne sont pas calculés et planifiés avec précision. Par exemple, la mise en œuvre de 50 postes de travail virtuels supplémentaires est une action qui peut créer un pic significatif dans les E/S. Pour ces raisons, il est nécessaire de prendre en compte les opérations de maintenance dans l'architecture de stockage pour les pics d'IOPS.

Il existe plusieurs manières de concevoir des solutions VDI avec des clones complets ou une image partagée, et chacune peut avoir des effets différents sur les exigences de stockage en termes de capacité et de performances. Étant donné que les clones complets consomment plus de capacité et de stockage, la déduplication sera importante. Les clones complets doivent également être corrigés indépendamment, ce qui augmentera les E/S au cours de ces opérations.

L'approche par image partagée proposée par Citrix avec MCS ou PVS et VMware avec des clones liés présente différents défis d'E/S. Par nature, ces approches par image partagée nécessitent moins de capacité de stockage car l'image principale est partagée et que chaque poste de travail virtuel nécessite moins d'espace pour ses données uniques. L'image partagée a des exigences de performance différentes de celles de la VM typique. Elle est désormais utilisée par des centaines ou des milliers de postes de travail virtuels et doit pouvoir générer de grandes quantités d'IOPS pour gérer des situations telles que les boot storm. Si l'image partagée est un goulot d'étranglement, tous les postes de travail virtuels l'utilisant seront affectés et l'expérience utilisateur en pâtira.

En tenant compte de ces considérations pour les pics et différents types d'architectures d'applications et de postes de travail, il est nécessaire de sélectionner et de concevoir une solution de stockage capable de répondre aux besoins accrus au démarrage et à la connexion, ainsi qu'aux exigences de l'environnement à l'état stationnaire. Pour comprendre les exigences de stockage de votre projet, il convient de procéder à une évaluation du poste de travail sur l'environnement physique du PC existant. Cette évaluation permettra de recueillir des détails sur les performances et la capacité réelles de la base d'utilisateurs afin que vous puissiez les appliquer à vos calculs de conception.

Une dernière considération sur les exigences de stockage liées à la virtualisation d'applications et de postes de travail est qu'en plus d'être très imprévisible au niveau des E/S, les charges de travail des postes de travail sont également très lourdes à écrire. Contrairement à de nombreuses charges de travail de serveurs qui lisent principalement des données et les fournissent aux utilisateurs, les postes de travail prennent généralement plus de temps à écrire sur le disque. Les opérations d'écriture sont plus intensives pour les baies de stockage que les lectures. Une charge de travail typique du serveur peut comprendre 80 % de lectures et 20 % d'écritures, tandis que ce rapport peut être inversé pour la charge de travail d'un poste de travail virtuel à l'état stationnaire. Lorsque vous évaluez vos choix de stockage, veillez à porter une attention particulière à la façon dont la solution de stockage met en mémoire tampon et valide les écritures, au lieu de vous concentrer uniquement sur des promesses de stockage qui font « un excellent travail » dans la mise en cache des blocs de lecture couramment lus pour gérer les boot storm.

TYPES DE STOCKAGE

Il existe différents types de stockage. Les principales alternatives disponibles à ce jour sont les baies de stockage multi-niveaux, les baies Flash hybrides et les baies 100 % Flash. Chaque solution adopte une approche différente pour assurer performance et capacité aux charges de travail. Au sein de chaque alternative, les fournisseurs adoptent différentes approches dans la construction de leurs offres, c'est pourquoi nous avons listé ci-dessous une brève explication de chacune d'entre elles.

Architectures multi-niveaux héritées - Il s'agit des baies d'entreprise héritées utilisées pour les charges de travail basées sur serveur depuis 10 à 20 ans. Ce sont généralement des architectures basées sur deux contrôleurs qui ont été modifiées au cours de la dernière décennie pour permettre d'inclure plusieurs niveaux de performance et de capacité disque dans l'architecture. Différents niveaux de disques sont fournis pour essayer de répondre aux besoins de capacité et de performance des charges de travail dispersées. Deux options sont alors disponibles : vous pouvez concevoir de manière à favoriser la performance en créant des pools dédiés de disques hautement performants pour une charge de travail, mais cela peut être très coûteux et contraignant. L'autre option consiste à essayer de tirer parti de la hiérarchisation qui a été ajoutée à cette architecture pour demander à la baie de promouvoir ou de rétrograder des blocs de données en fonction de la demande. Or, le problème de cette hiérarchisation automatique est que la prise de ces décisions pour les charges de travail VDI est souvent trop longue.

100 % Flash - Les baies de stockage 100% Flash sont entièrement composées de stockage Flash. Il existe de nombreux types de Flash pouvant être utilisés dans ces baies de stockage. Les baies modernes 100 % Flash ont été conçues pour tirer parti des caractéristiques du stockage Flash, ce qui signifie que le système d'exploitation et le système de fichiers ont été conçus pour Flash. Certains produits ont adopté une conception de baie héritée et ont simplement remplacé les disques rotatifs par des disques 100 % Flash. Bien qu'il s'agisse d'une option beaucoup plus performante que la précédente, le produit final n'a pas été conçu à cet effet.

Les baies de stockage 100% Flash sont très rapides, avec un niveau de performance unique dans le produit. Pour assurer que la baie peut également fournir la capacité requise pour la conception à un prix abordable, vous devez rechercher des baies qui offrent déduplication et compression. Bien que presque toutes les baies 100 % Flash modernes soient plus faciles à gérer que les baies héritées, elles n'offrent pas toujours la même facilité de gestion globale et de gestion par machine virtuelle que la plupart des solutions Flash hybrides.

Flash hybride - Les baies de stockage hybrides sont des architectures modernes conçues pour utiliser efficacement une combinaison de lecteurs Flash et de disques rotatifs. Les fournisseurs ont adopté différentes approches architecturales pour utiliser la capacité et les performances dans leurs baies, mais les résultats finaux sont similaires. Ils sont tous capables d'offrir des performances impressionnantes à partir d'une quantité réduite de Flash, tout en fournissant une grande capacité par le stockage des données sur de grands disques rotatifs utilisés dans la baie. Les alternatives idéales d'architecture de stockage hybride utilisent l'intelligence intégrée pour hiérarchiser automatiquement les données sur des lecteurs Flash et disques en fonction de la demande, éliminant le besoin de réglage manuel et les pièges potentiels de performance.

Les architectures les mieux adaptées à un projet de VDI moderne sont les architectures de stockage hybrides et 100 % Flash. Elles peuvent fournir les performances requises pour les environnements VDI et offrent généralement aussi les expériences de gestion modernes décrites précédemment. Les charges de travail VDI sont, par nature, très imprévisibles. Si votre solution de stockage met trop de temps pour prendre des décisions de stockage ou promouvoir des blocs à un certain niveau de mise en cache, les performances se dégraderont considérablement et l'expérience utilisateur en pâtira.

Il existe différentes approches pour déterminer le dimensionnement de la couche de calcul du projet. La première est l'approche « scale up » qui utilise quelques grands hôtes pour fournir des ressources, tandis que l'approche « scale out » utilise plusieurs petits hôtes pour fournir des ressources. La méthode idéale se situe quelque part entre les deux approches, et utilise deux socket hôtes et les rend aussi denses que possible sans violer les ratios de consolidation définis dans le cadre du projet. Le but de cet eBook est d'aider à dimensionner les ressources informatiques pour la charge de travail VDI.

Il y a trois calculs principaux à prendre en compte lors du dimensionnement des ressources de calcul à l'étape de la conception : la quantité de mémoire physique dans chaque hôte, la fréquence d'horloge du CPU, le nombre de cœurs du CPU et leur ratio de CPU pour chacun. Tout d'abord, il ne faut jamais surcharger la mémoire dans un projet VDI. Le non-respect de cette règle n'aura que très peu d'utilité et n'entraînera que des problèmes de performance dans l'environnement.

Le calcul de la fréquence d'horloge du CPU dépend fortement des informations recueillies précédemment lors de l'évaluation du poste de travail. Les rapports d'évaluation indiqueront la quantité de CPU utilisée en moyenne et aux heures de pointe par session utilisateur. On utilisera ces détails ainsi que les informations sur la mémoire contenus dans l'évaluation pour effectuer les calculs.

Quelques autres conseils sur les clusters d'hôtes et la virtualisation : ne dépassez jamais 80 % d'utilisation d'un hôte et dimensionnez toujours votre cluster à N+1. La règle des 80 % d'utilisation de l'hôte n'est pas seulement une question de déploiements de virtualisation d'applications et de postes de travail, c'est une recommandation qui s'applique à toute charge de travail s'exécutant sur un hyperviseur. Si vous utilisez vos hôtes au-delà de la barre des 80 %, vous aurez très peu d'espace disponible pour les pics d'activité et vos ressources ne seront peut-être pas suffisantes en cas de défaillance d'un hôte, selon la taille de votre cluster. Le deuxième élément à prendre en compte dans le dimensionnement N+1 du cluster est de s'assurer qu'il dispose de suffisamment de ressources pour compenser la défaillance d'un seul hôte, afin de garantir que toutes les machines virtuelles puissent continuer à fonctionner, et que celles qui échouent redémarrent sans problème. La défaillance d'un seul hôte est le niveau de résilience le plus courant ; seul un nombre limité de clients ont besoin d'un niveau N+2 pour répondre à des exigences SLA plus strictes.

Le dernier élément relatif au dimensionnement des capacités de traitement est le ratio de CPU, soit le nombre de CPU virtuels par CPU physique (vCPU/pCPU). Ce ratio est crucial car s'il est trop élevé, il est fort probable qu'un problème de planification du CPU survienne, et affecte considérablement les performances et l'expérience utilisateur. Lorsqu'un tel problème survient sur les hôtes vSphere, la durée de disponibilité du CPU augmente, ce qui indique que le planificateur a des difficultés à programmer tous les vCPU sur des pCPU. Le vCPU devra alors attendre, même s'il est prêt. Le ratio de CPU est très varié pour les différents types de charges de travail virtualisées sur les clusters VMware. En règle générale, les charges de travail de serveur et de database ont un ratio beaucoup plus faible, tandis que les charges de travail VDI peuvent avoir un ratio plus élevé.

L'utilisation des vCPU n'est pas un calcul linéaire, dans le sens où il est possible de créer un hôte qui présente un ratio de consolidation plus élevé si toutes les machines virtuelles ont un seul vCPU. Lorsque plusieurs VM ont deux vCPU ou plus, cela affecte les calculs : vous ne pouvez pas simplement diviser par deux pour représenter deux fois plus de vCPU. La Figure 6 représente un intervalle qui s'est avéré efficace pour des déploiements réels chez les clients. Les fournisseurs qui effectuent des tests de synthèse peuvent afficher des ratios plus élevés. Soyez prudent avec ce type de résultats, car ils ne s'appliquent pas toujours aux projets du monde réel.

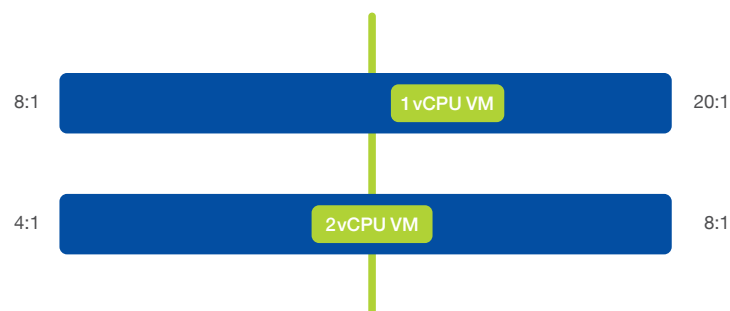


Figure 6 :
La consolidation de la VDI dépend largement de la proportion de vCPU avec lequel vos postes de travail virtuels seront configurés. Le graphique représente un intervalle dont l'expérience a démontré la fiabilité.

L'intervalle pour le fonctionnement normal d'un seul poste de travail virtuel vCPU est de 8:1 à 20:1. Il s'agit d'un intervalle large, et le point où vous pouvez vous situer dans cet intervalle est déterminé par plusieurs choix. Le premier concerne la taille des hôtes, le nombre de machines virtuelles par hôte et le niveau de confort du client par rapport à ce nombre. Par exemple, un hôte double socket avec deux CPU 18 coeurs pourrait prendre en charge plus de 700 VM, à condition que vous disposiez de la quantité de mémoire appropriée et de la fréquence d'horloge disponible. Généralement, avoir autant de machines virtuelles sur un seul hôte effraie la plupart des clients. Il y a donc deux choix possibles dans ce scénario : le premier consiste à choisir une densité plus faible en imposant une limitation artificielle.

Choisir l'extrémité inférieure du ratio équivaldrait à 288 VM sur le même hôte. La deuxième option serait de choisir des CPU avec moins de coeurs, mais de choisir un ratio quelque part au milieu de l'intervalle. Si vous choisissez des CPU de 12 coeurs et utilisez un ratio de 12:1, vous obtiendrez 288 VM. Cette décision est généralement une combinaison de feedback client, de recommandations d'architectes et de prix de l'infrastructure. Le choix de différentes configurations de CPU physiques peut générer des économies de coûts significatives

Les calculs pour un poste de travail virtuel avec deux vCPU sont similaires, sauf que vous devez traiter deux fois plus de vCPU. L'intervalle de fonctionnement dans ce cas est compris entre 4:1 et 8:1. Certains fournisseurs promettent de meilleurs résultats, mais nos recommandations sont basées sur des déploiements réels par les clients. Vous devez utiliser les mêmes points de décision que dans l'exemple précédent, mais avec un intervalle de ratio de CPU différent.

Une autre chose à garder à l'esprit est que si vous sélectionnez un ratio de CPU situé entre ces intervalles, vous serez libre de redimensionner la densité de consolidation si l'environnement continue à fonctionner dans les tolérances. Une chose à noter est qu'il n'existe aucun endroit spécifique pour configurer ces ratios de CPU en tant que paramètre dans les outils disponibles aujourd'hui. Ces attributs doivent être inscrits dans la conception et devenir des données à prendre en compte dans la gestion et la mise à l'échelle de l'environnement. Tout comme la mémoire et la fréquence d'horloge, le ratio de CPU doit être pris en considération dans la décision d'ajouter plus de machines virtuelles à un cluster, et du moment opportun pour ajouter un autre hôte à un cluster afin de fournir plus de ressources.

Il est possible de gérer le ratio CPU par des calculs manuels en collectant des données. Certains administrateurs utilisent un script PowerShell qui collecte les données et présente le ratio comme résultat du script. Avec un script, il peut être exécuté quotidiennement comme un processus planifié pour assurer le respect du ratio et l'absence de danger dans les clusters.

La fréquence de la RAM ou du bus mémoire est également associée au dimensionnement de calcul. La règle lors du dimensionnement de la mémoire est de viser la densité maximale avec la vitesse de bus la plus élevée possible en fonction de votre budget. Le problème souvent rencontré avec la mémoire est que son ralentissement peut entraîner des cycles de CPU inactifs en attente de la fin des opérations de lecture/écriture en RAM.

Il existe plusieurs raisons de créer différents clusters de virtualisation dans un projet EUC. La décision d'avoir des clusters différents vient généralement de charges de travail et de tailles de cluster différentes. Nous ne nous étendrons pas sur ce propos dans cet eBook, mais voici quelques recommandations qui s'appuient sur les sujets abordés ailleurs dans le livre et en ligne.

Tout d'abord, lors de la construction d'un projet VDI impliquant plus de quelques centaines d'utilisateurs, il est essentiel de séparer l'infrastructure de gestion de la virtualisation de la charge de travail VDI. Cela signifie que tous les serveurs de gestion, les brokers VDI, les serveurs de fichiers, les serveurs de gestion d'applications et toutes les autres fonctions qui ne sont pas des postes de travail virtuels doivent s'exécuter sur un cluster différent. La question de savoir si le cluster de gestion doit être uniquement dédié à la conception de l'EUC dépendra de la taille de l'environnement. Si le projet est plus petit, vous pouvez exécuter des VM de gestion dans un cluster de virtualisation de serveurs existant.

Vous pouvez redimensionner ces clusters de postes de travail virtuels pour atteindre une taille comprise entre 16 et 32 hôtes. Cet intervalle permet de créer un plus grand pool de ressources pour les machines virtuelles, et pousse également la plupart des clients à adopter un cluster plus volumineux que leurs tailles habituelles. Les mises à jour récentes de l'hyperviseur permettent de créer des clusters allant jusqu'à 64 hôtes, mais il faudra du temps avant que les architectes et clients se sentent à l'aise avec des clusters si grands. Si l'environnement est suffisamment grand pour que le nombre d'hôtes dépasse ces intervalles, vous pouvez avoir besoin de plus d'un cluster VDI.

Une autre raison de penser le projet pour plusieurs clusters de virtualisation, en dehors de la taille de l'environnement, est la multiplicité des charges de travail. Il existe différentes charges de travail dans les clusters VDI. S'il existe une quantité significative de postes de travail virtuels, 1 vCPU et 2 vCPU, il convient de concevoir un cluster distinct pour chacun. La Figure 7 illustre une approche de conception multi-cluster. Cela permet de gérer le ratio de CPU de manière différente dans chaque cluster, facilitant ainsi la conception. Si vous deviez fusionner les différentes configurations de CPU, il faudrait alors calculer un nouveau ratio fusionné, et cela ne ferait qu'embrouiller les choses.

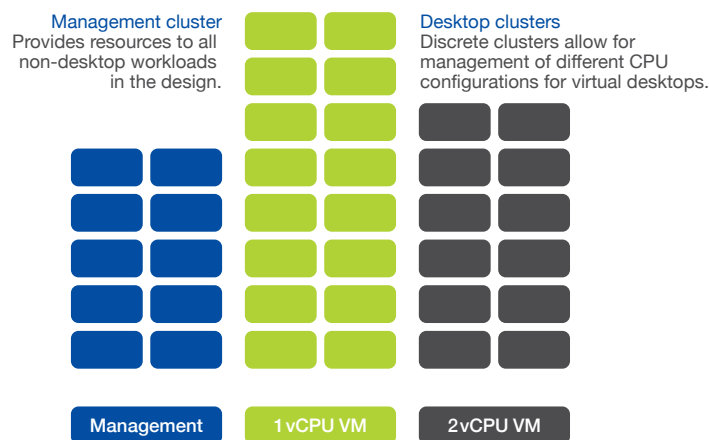


Figure 7 :
Clusters de gestion et de postes de travail

Si vous avez tout lu jusqu'à présent, nous espérons que vous serez au moins intrigué par les possibilités de la virtualisation de postes de travail et d'applications. Si vous et votre organisation êtes prêts à apprendre comment déployer avec succès la virtualisation d'applications et de postes de travail, sachez que Nutanix est là pour vous accompagner. L'infrastructure invisible de Nutanix peut grandement simplifier votre parcours grâce à son architecture web-scale maintes fois récompensée comme étant la meilleure plateforme VDI.

LE MEILLEUR MOMENT POUR SE LANCER

Cela ne vous surprendra probablement pas : Nutanix a beaucoup réfléchi aux meilleurs moyens d'assurer le succès des déploiements de virtualisation de postes de travail et d'applications.

Comprendre votre environnement actuel :

Le processus commence par une compréhension complète de votre environnement d'utilisateur final actuel, notamment :

- Les statistiques sur l'utilisateur final : rassemblez les profils de l'utilisateur final et les facteurs associés tels que les applications utilisées, les périphériques d'accès, la localisation et la connectivité.
- Les services réseau et les indicateurs spécifiques à l'infrastructure : collectez les informations appropriées sur différents services destinés aux utilisateurs finaux, tels que les services de fichiers, l'authentification et le contrôle d'accès, ainsi que l'équilibre pare-feu/chargement. Rassemblez également les mesures de performance, de latence, de débit, etc.
- Faites correspondre chaque élément à ses propriétaires : la responsabilité est un facteur clé de succès.

Dimensionnement du nouvel environnement :

Maintenant que vous disposez de toutes ces informations, vous pouvez dimensionner avec précision votre nouvel environnement. Nutanix Sizer simplifie cette tâche, mais voici quelques directives à garder à l'esprit :

- Prévoyez toujours une haute disponibilité pour les serveurs et postes de travail clés.
- Des infrastructures ou clusters supplémentaires peuvent être requis au regard de ces considérations :
 - Business : SLA, licences, sécurité, budget, politiques
 - Planification de la transition : suivez les meilleures pratiques et directives de Nutanix et de l'industrie en matière de migration P2V et accordez une attention particulière à la création d'une image maître qui sera utilisée pour créer d'autres postes de travail. Si vous migrez un déploiement existant, nous vous recommandons d'utiliser des outils partenaires Nutanix ou des outils natifs lorsque cela est possible.

Naturellement, Nutanix Global Services peut vous aider à franchir toutes ces étapes pour assurer le succès de votre infrastructure. Grâce à notre organisation Global Services, Nutanix offre la seule solution du secteur pour éliminer le risque de mise à l'échelle de votre infrastructure pour les projets de virtualisation de postes de travail.

Dans le cadre du programme VDI Assurance, Nutanix veille à ce que vos postes de travail virtuels disposent toujours des ressources de traitement (CPU virtuel et mémoire) et de stockage (performances et capacité) nécessaires pour répondre aux attentes VDI de l'utilisateur final. Il suffit de déterminer le type et le nombre d'utilisateurs de VDI dans leur environnement et de transférer le risque de dimensionnement des besoins en infrastructure à Nutanix en utilisant VDI Assurance.

Si vous souhaitez en savoir plus sur l'infrastructure invisible pour les applications d'entreprise, contactez-nous à info@nutanix.com, envoyez-nous un DM sur [Twitter @nutanix](https://twitter.com/nutanix) ou une demande à www.nutanix.com/demo pour configurer votre propre briefing et démonstration personnalisés. Nous pourrions voir ensemble comment les solutions validées et certifiées de Nutanix peuvent vous aider à tirer le meilleur parti de vos applications d'entreprise.

Restez en contact avec les experts et les clients de Nutanix sur la communauté en ligne Nutanix Next (next.nutanix.com).

Nutanix fournit une infrastructure invisible pour la dernière génération d'informatique d'entreprise, permettant aux équipes IT de se concentrer sur les applications et les services stratégiques. La plateforme pilotée par logiciel Xtreme Computing fait converger nativement calcul, virtualisation et stockage dans une solution unique pour apporter plus de simplicité au datacenter. En utilisant Nutanix, les clients bénéficient de performances prévisibles, d'une évolutivité linéaire et d'une consommation d'énergie similaire à celle du cloud. En savoir plus sur www.nutanix.fr ou suivez-nous sur Twitter @NutanixFrance.

NUTANIX™

T. +33 (0)1 82 88 15 90
contact-france@nutanix.com | www.nutanix.fr | [Twitter @NutanixFrance](https://twitter.com/NutanixFrance)