

# Une architecture de référence pour l'analytique à hautes performances dans le domaine de la santé et des sciences de la vie

*Accélérer les soins personnalisés et d'autres charges de travail biomédicales en utilisant une infrastructure à hautes performances et rentable pour l'analytique des mégadonnées*

---

## Vue d'ensemble

### Défi

Les organisations de santé et de sciences de la vie du monde entier doivent gérer, utiliser, stocker, partager et analyser des mégadonnées dans les limites de leur budget informatique.

### Solution

L'architecture de référence d'IBM pour les soins de santé et les sciences de la vie définit une plateforme permettant d'offrir les meilleurs niveaux de performances pour les charges de travail de mégadonnées tout en diminuant le coût de possession des TI.

---

Les progrès des techniques de profilage moléculaire à haut débit et des systèmes informatiques à hautes performances ont mené à une nouvelle ère de soins personnalisés où le traitement et la prévention des maladies peuvent être adaptés aux profils moléculaires uniques, aux caractéristiques comportementales et aux risques environnementaux des patients individuels. Pour découvrir les thérapies personnalisées convenant aux différentes cohortes de patients – et pour offrir de tels plans de traitement dans un environnement clinique – il faut des plateformes informatiques pouvant analyser les caractéristiques des patients et prédire les résultats cliniques rapidement et avec exactitude. Pour de nombreuses organisations cliniques et de recherche biomédicale, les projets à grande échelle dans les soins personnalisés sont décourageants pour les raisons suivantes :

**Mégadonnées :** Les données biomédicales nécessaires pour soutenir les initiatives de soins personnalisés sont nombreuses, variées et – souvent – non structurées; de plus, la quantité de données recueillies dans le cadre de ces projets augmente généralement de façon exponentielle, et les données doivent être archivées pendant de longues périodes – parfois des décennies. Les sources de données courantes pour les applications de soins de santé personnalisés incluent des séquences entières de génome, les images biomédicales provenant d'instruments cliniques et de laboratoires, les systèmes de dossiers médicaux électroniques, les capteurs physiologiques ou les dispositifs portables, et les collections organisées de textes scientifiques et cliniques. Il arrive souvent que les organisations n'aient pas la capacité de stockage pour suivre le rythme de l'expansion des données.

**Données en silos :** Les soins de santé personnalisés demandent l'agrégation d'informations qui donnent une vue complète des caractéristiques biologiques et comportementales ainsi que des risques environnementaux de chaque patient. Cependant, les données des patients sont généralement disséminées dans des silos de stockage hétérogènes à l'intérieur des systèmes de santé et des organismes de recherche biomédicale. Avant que des ensembles de données disparates puissent être analysés, ils doivent être intégrés dans une base de données commune, ce qui est un processus manuel souvent long et méticuleux.



---

## Composants du système

### Intergiciels de gestion des charges de travail et des données

- IBM Spectrum Scale
- IBM Spectrum LSF
- IBM Spectrum Conductor with Spark

### Matériel de traitement et de stockage

- IBM Power Systems
- IBM Elastic Storage Server

### Réseau

- IBM Aspera

### Environnements infonuagiques

- IBM SoftLayer
  - IBM Cloud Object Storage
- 

**Charges de travail à données et calculs intensifs :** Les charges de travail analytiques peuvent être très intensives en données et en calculs. Des exemples courants incluent les pipelines d'analyse à E-S intensive, qui transforment les données de séquençage brutes de la prochaine génération en fichiers génomiques variants, les techniques d'apprentissage profond permettant de découvrir des tendances dans des ensembles de données biomédicales complexes, et l'exploration de données à grande échelle dans des documents cliniques et scientifiques. L'exécution de ces charges de travail peut prendre des heures, et même des jours, sur les plateformes informatiques existantes.

**Applications et cadres de référence en évolution :** Les initiatives de soins de santé personnalisés doivent souvent prendre en charge à n'importe quel moment des centaines d'applications, y compris celles qui sont liées à l'informatique médicale, à la génomique, à l'analyse d'images et à l'apprentissage profond. Ces applications sont souvent construites sur des cadres et des bases de données qui évoluent constamment, y compris Spark, Hadoop, TensorFlow, Caffe, Docker, MongoDB, HBase et diverses bases de données graphes. Les organismes de recherche biomédicale ont souvent des difficultés à soutenir des versions multiples de ces cadres d'applications et bases de données, qui prolifèrent et évoluent fréquemment – parfois deux fois ou plus par an.

**Collaboration :** Le partage des données entre institutions – et, souvent, au-delà des limites géographiques – est une nécessité croissante pour l'étude des maladies rares et des mécanismes des maladies complexes. De plus en plus de groupes internationaux rassemblant des universitaires, des entreprises privées, des organisations à but non lucratif et des organismes publics apparaissent pour partager des données biomédicales et des analyses connexes. Mais il manque souvent aux partenaires les solutions qui leur permettraient de partager leurs ensembles de données rapidement et de façon rentable sans compromettre l'information protégée sur la santé et les droits de propriété intellectuelle.

Un grand nombre d'organisations du monde entier trouvent difficile de surmonter ces obstacles, surtout dans les limites de leur budget informatique. Elles doivent accéder aux données de recherche clinique et scientifique et les stocker, les analyser, les partager et les archiver de façon rapide et rentable; mais pour de nombreuses organisations de soins de santé, organismes de recherche biomédicale et compagnies pharmaceutiques, les données sont recueillies en telle quantité qu'ils ne peuvent plus les traiter, les stocker de façon adéquate ni les transmettre assez rapidement par les moyens de communication habituels. Pour bon nombre d'entre eux, les silos de traitement et de stockage prolifèrent entre les groupes cliniques et de recherche, tandis que les analystes recueillent des volumes croissants de données et les utilisent dans des charges de travail analytique complexe. Pour transférer des données sur de longues distances, les organisations ont souvent recours à des entreprises de transport qui font parvenir les données brutes sur disque à des centres informatiques externes à des fins de traitement et de stockage, ce qui ralentit leur analyse.

---

## Capacités clés de la plateforme

- **Évolution** permettant de prendre en charge des volumes de mégadonnées qui augmentent de façon exponentielle.
  - **Flexibilité** pour soutenir les applications analytiques construites sur Spark, Hadoop, Docker et d'autres cadres de référence.
  - **Intégration simplifiée** de données biomédicales entre les silos de stockage.
  - Stockage, gestion et analyse de **données non structurées**.
  - Saisie et stockage de **métadonnées** à des fins de recherche, de répétabilité et de vérifiabilité des données et des flux de travaux.
  - **Collaboration facile** malgré les frontières géographiques.
  - **Sécurité des données** de santé protégées et protection des droits de propriété intellectuelle.
  - **Administration des TI** facile et peu coûteuse.
- 

Dans le but d'aider les informaticiens du secteur à surmonter les défis techniques, le groupe Systemes IBM a créé une architecture de référence pour la santé et les sciences de la vie. Cette architecture, qui est fondée sur l'historique de prestation d'IBM des meilleures pratiques en informatique hautes performances (HPC), permet aux organisations des soins de santé et des sciences de la vie de faire évoluer facilement les ressources de traitement et de stockage à mesure que croît la demande, et de soutenir la vaste gamme de cadres de développement et d'applications nécessaires à l'innovation dans le secteur – le tout sans réinvestir dans la technologie. La description de l'architecture de référence d'IBM est fournie dans la prochaine section.

## Architecture de référence d'IBM : une plateforme informatique diversifiée construite sur une infrastructure commune

Que ce soit des chercheurs universitaires qui publient des articles dans des journaux spécialisés afin d'obtenir du financement, des scientifiques travaillant dans des organisations de R et D qui soumettent des médicaments potentiels à des essais cliniques ou des médecins dans les hôpitaux qui prescrivent des traitements visant à obtenir les meilleurs résultats cliniques pour leurs patients, les principaux intervenants qui participent à des projets de soins de santé personnalisés ont besoin d'une plateforme informatique souple et fiable qui réponde à leurs besoins divers en applications. Par l'entremise du travail auprès de nombreux clients et partenaires du secteur de la santé dans le monde, le groupe Systemes IBM a appris que les organisations qui s'intéressent à la génomique, aux soins de santé personnalisés et à d'autres initiatives liées aux mégadonnées en recherche biomédicale auront besoin de systèmes à hautes performances offrant des capacités de plateforme clés.

L'architecture de référence pour la santé et les sciences de la vie (telle qu'elle est illustrée à la Figure 1) a été conçue par le groupe Systemes IBM pour répondre à cet ensemble de besoins communs. Elle reflète l'évolution actuelle de l'informatique à hautes performances (HPC), où les systèmes informatiques doivent traiter les charges de travail par lots de l'informatique HPC traditionnelle, ainsi que l'analytique de longue durée qui porte sur les mégadonnées. Prendre des composantes de traitement et de stockage exécutant des codes algorithmiques HPC, puis leur attribuer des ressources de manière dynamique pour traiter d'autres types d'analyses est une façon moins coûteuse et plus facile à gérer qu'entretenir deux systèmes distincts. À l'ère des charges de travail de calcul diverses, des besoins d'infrastructure divers et des différentes versions des cadres de référence d'applications, il est important d'éviter de créer des groupes de traitement en silos qui, en général, sous-utilisent les ressources informatiques et se traduisent par une mauvaise maîtrise des coûts des TI. La capacité d'utiliser moins de ressources de traitement en faisant appel à des outils intelligents fondés sur des politiques pour la gestion des charges de travail et des données est un élément essentiel à un environnement de recherche dynamique.

Avec des volumes de données à croissance aussi rapide, les systèmes de stockage doivent être évolutifs du point de vue de la capacité et de la performance. Ils doivent aussi être conçus de façon à ce que les mêmes supports de stockage puissent prendre en charge différentes méthodes d'accès d'une façon fiable, mais souple. Par exemple, un serveur de stockage commun qui peut prendre en charge des méthodes d'accès E-S parallèle doit aussi traiter efficacement l'accès aux données non structurées à partir de plateformes telles que Hadoop. De plus, les environnements de traitement technique qui répondent aux exigences toujours changeantes relatives aux charges de travail doivent être faciles à gérer. L'administration simple et centralisée des systèmes de stockage permet non seulement aux chercheurs d'être plus productifs, mais peut aussi diminuer les coûts informatiques et les coûts connexes dans toute l'entreprise.

Pour faire évoluer les plateformes techniques à hautes performances afin de soutenir les volumes de données croissants et les diverses applications tout en continuant à accélérer les charges de travail et à réduire les coûts au minimum, il faut un cadre de référence souple, mais bien coordonné pour l'accès aux données, le traitement et le stockage.

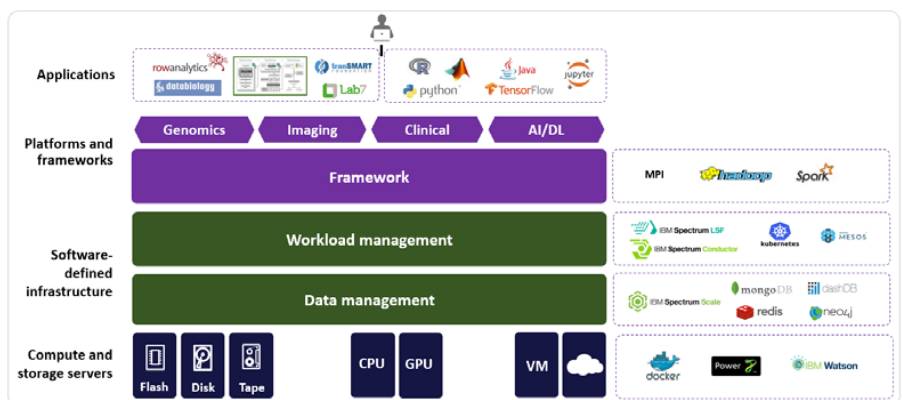


Figure 1. Architecture de référence d'IBM pour la santé et les sciences de la vie

## Éléments de base : Spectrum Computing, Spectrum Scale et Power Systems

L'architecture de référence d'IBM (voir la Figure 1) reflète le travail qui est toujours en cours dans le groupe Systèmes IBM pour intégrer les éléments de la gamme de produits de traitement et de stockage d'IBM afin qu'ils offrent des niveaux élevés de performance pour les mégadonnées, tout en réduisant le coût total de possession des TI. Les deux couches supérieures représentent les applications, les bases de données et les cadres de référence que les chercheurs et les cliniciens utilisent (*boîtes violettes*); la couche inférieure représente les serveurs de traitement et de stockage virtuels ou physiques hautes performances où sont traitées et stockées les données (*boîtes bleues*); les deux couches intermédiaires (*boîtes vertes*) contiennent les logiciels qui permettent à diverses applications biomédicales d'être exécutées rapidement et de façon rentable sur des serveurs de traitement et de stockage partagés.

Pour ce qui est des deux couches intermédiaires, la couche de gestion des charges de travail fait en sorte qu'on peut répartir des milliers de flux de travaux de traitement, en parallèle, dans les serveurs de calcul de l'entreprise d'une manière qui en maximise l'utilisation. La couche de gestion des données permet l'accès aux données avec un faible temps d'attente, la convergence des silos de données hétérogènes, la collecte des métadonnées et la gestion automatisée du cycle de vie de l'information. Elle permet aux organisations de faire évoluer rapidement le traitement et la capacité de stockage à la hausse, ou de diminuer la capacité selon la demande des charges de travail.

### **Gestion des charges de travail : IBM Spectrum Computing**

La couche de gestion des charges de travail attribue de façon dynamique et élastique les tâches de calcul aux divers serveurs de traitement sans intervention de l'utilisateur. Elle consiste en plusieurs ordonnanceurs cohérents de flux des travaux qui sont coordonnés pour placer les tâches de calcul sur des groupes locaux et à distance de manière efficace et rentable. Les ordonnanceurs fondés sur les ressources et les politiques d'IBM comprennent Spectrum LSF (pour les charges de travail HPC par lots), Spectrum Conductor with Spark (pour les charges de travail Spark) et IBM Spectrum Symphony (pour Hadoop MapReduce et les charges de travail en temps quasi réel). Comme on le voit à la Figure 2, ces ordonnanceurs sont étroitement intégrés : si un type de charge de travail n'utilise que quelques ressources dans un groupe, les autres types de charges de travail peuvent utiliser complètement les ressources restantes. La souplesse et l'élasticité de l'utilisation des serveurs entre ces ordonnanceurs éliminent le besoin pour les organisations TI de fournir des groupes dédiés pour chaque type de charge de travail. Quand ils servent un environnement à service partagé, ces ordonnanceurs protègent les clients individuels par l'isolation sécurisée et le contrôle de l'accès fondé sur les rôles. Ils permettent aux charges de travail d'être réparties de façon transparente dans plusieurs environnements physiques et infonuagiques (voir la Figure 3), et ils prennent en charge la répartition des charges de travail déployées dans Docker et d'autres technologies de conteneur.

**IBM Spectrum LSF :** Spectrum LSF est une plateforme de gestion des charges de travail hautement évolutive et sensible aux ressources qui soutient les environnements HPC exigeants, distribués et à mission vitale. Elle a été sélectionnée comme système privilégié de gestion des charges de travail par d'importantes organisations d'analyse de génome pour sa capacité à orchestrer de façon routinière des centaines de milliers de tâches qui sont soumises par lots, ainsi que pour sa capacité à évoluer facilement en fonction de la demande des utilisateurs. Dans le monde entier, les clients travaillent dans des environnements informatiques pris en charge par LSF pour exécuter des centaines de charges de travail génomique, notamment Burrows-Wheeler Aligner (BWA), SAMtools, Picard, GATK, Isaac, CASAVA et d'autres pipelines fréquemment utilisés pour l'analyse génomique.

D'autres fonctions de gestion et de surveillance des charges de travail sont offertes par les logiciels compagnons suivants :

- **IBM Spectrum LSF Process Manager** : Cette application utilisateur simplifie la conception et l'automatisation de flux de travaux de calculs complexes et permet leur partage avec des collaborateurs dans un environnement sécurisé. Les clients d'IBM qui s'intéressent à la recherche en soins de santé personnalisés ont investi dans LSF Process Manager pour simplifier le processus de création de scripts de flux de travaux génomiques qui transforment les données brutes des séquenceurs de la prochaine génération (en format FASTQ) en fichiers variants (par exemple, VCF, SNV, CV) pour l'analyse en aval. Spectrum LSF Process Manager permet aux bio-informaticiens de partager leurs flux de travaux avec certains utilisateurs qui peuvent – ou non – avoir de l'expérience formelle de création de scripts, tout en aidant à maintenir un strict contrôle des versions de ces scripts.
- **IBM Spectrum LSF Application Center** : Cette application fournit aux utilisateurs une interface Web simple et intuitive pour qu'ils accèdent aux flux de travaux créés avec Spectrum LSF Process Manager, et pour qu'ils les soumettent, les gèrent et les surveillent. IBM Spectrum LSF Application Center peut aussi servir de portail d'entreprise à service partagé, où les utilisateurs et les collaborateurs externes peuvent soumettre, exécuter et surveiller les pipelines d'analyses selon des contrôles d'accès à granularité fine fondés sur les rôles.
- **IBM Spectrum LSF RTM** : Ce tableau de bord opérationnel équipe les administrateurs des TI d'outils complets de surveillance et de gestion des charges de travail, ainsi que de production de rapports. Les administrateurs peuvent ainsi faire le suivi des mesures relatives aux charges de travail pour ce qui est de l'utilisation et de l'efficacité des groupes selon une granularité fine – depuis le niveau du groupe jusqu'à celui des utilisateurs spécifiques, des groupes d'utilisateurs ou des types de charges de travail.
- **Spectrum LSF Data Manager** : Cette fonction permet aux administrateurs de créer des politiques qui transfèrent automatiquement les données d'un groupe géré par LSF (sur site ou en nuage) vers un cache résidant près du groupe de calcul d'exécution. De telles politiques peuvent automatiser le déplacement des données à distance, ce qui améliore le débit des données et réduit les durées de traitement.
- **IBM Spectrum Conductor with Spark** : Apache Spark est un cadre d'applications utilisé couramment par des personnes qui font de la recherche computationnelle dans le secteur des soins de santé personnalisés et d'autres domaines des sciences biomédicales. Les plateformes techniques à hautes performances qui prennent en charge les charges de travail HPC par lots traditionnelles doivent aussi prendre en charge un nombre croissant de charges de travail Spark. IBM Spectrum Conductor with Spark est une solution de haut niveau à service partagé pour Apache Spark, qui permet le soutien simultané d'instances multiples de Spark et élimine les silos de ressources qui seraient, sinon, liés à des mises en œuvre Spark distinctes. En offrant le partage évolué de

## IBM Spectrum Scale : principales fonctions

- **Évolutivité extrême** vers des milliards de pétaoctets et des centaines de Gbps.
- **Extension modulaire de la capacité de stockage** avec perturbation minimale du service.
- **Accès aux données avec un faible temps d'attente**, et meilleure performance de sa catégorie pour les charges intensives en E-S.
- **Un seul espace nom mondial** pour les fichiers et les répertoires dans des groupes de stockage hétérogènes.
- **Technologie AFM** (gestion active des fichiers) de mise en cache pour accélérer le partage des données entre les collaborateurs à distance.
- **Gestion du cycle de vie de l'information** (ILM) fondée sur des politiques dans des niveaux de supports de stockage.
- Identification, collecte et recherche des **métadonnées**.
- **Chiffrement des fichiers** selon des politiques.
- **Administration facile** à partir d'un seul tableau de bord.

toutes les charges de travail Spark et des différentes versions de Spark sur une seule plateforme de calcul, Spectrum Conductor permet aux organisations d'augmenter l'utilisation des serveurs de calcul existants et, par conséquent, d'améliorer la maîtrise des coûts des TI. Surtout, il facilite la gestion des charges de travail grâce à une interface conviviale.

- **IBM Spectrum Symphony** : Même s'il n'est pas encore largement utilisé dans les environnements de calculs techniques pour les applications de soins de santé ou la recherche biomédicale (comme l'est IBM Spectrum LSF), IBM Spectrum Symphony est un gestionnaire de grille d'entreprise hautement évolutif qui peut ordonnancer les tâches selon un temps d'attente de l'ordre de quelques millisecondes et, par conséquent, prendre en charge l'analytique en temps quasi réel. IBM Spectrum Symphony Advanced Edition inclut aussi une mise en œuvre de MapReduce compatible avec Apache Hadoop, dont il a été démontré au cours d'un test de performance vérifié qu'elle offre, en moyenne, quatre fois les performances de Hadoop à code source ouvert.

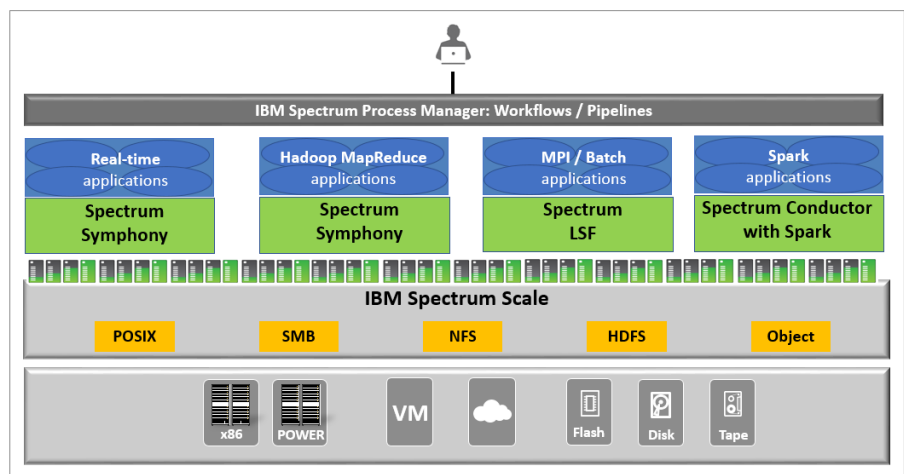


Figure 2. Gestion des charges de travail et des données avec IBM Spectrum Computing et Spectrum Scale

## Gestion des données : IBM Spectrum Scale

Dans l'architecture de référence d'IBM, la gestion des données pour les charges de travail computationnelles est optimisée par IBM Spectrum Scale, connu auparavant sous le nom d'IBM General Parallel File System (IBM GPFS). Spectrum Scale est une solution éprouvée et évolutive, à hautes performances, de gestion de données et de fichiers, qui permet une gestion du stockage de première classe avec une évolutivité extrême. Des centres de génomique, des institutions médicales et des compagnies pharmaceutiques de premier plan investissent déjà dans Spectrum Scale pour stocker, archiver, traiter et gérer une grande quantité de données structurées ou non, y compris des séquences génomiques, des images biomédicales et des dossiers médicaux électroniques. Spectrum Scale convient bien à la gestion des données biomédicales et à l'analytique connexe parce qu'il relève les défis suivants :

- Croissance rapide des volumes de données.
- Charges de travail intensives en E-S, ce qui est souvent nécessaire pour l'analyse des séquences génomiques brutes.
- Groupes hétérogènes de stockage de fichiers et d'objets utilisant différents systèmes d'exploitation, protocoles d'accès au stockage et matériels de stockage.
- Partage des données entre des projets répartis mondialement.
- Identification, collecte et recherche des métadonnées, qui sont nécessaires pour la répétabilité scientifique et clinique, la validation et l'archivage à long terme.
- Besoin de stockage et de suppression sécurisés des données contenant de l'information sensible.

En outre, Spectrum Scale permet la gestion du cycle de vie de l'information (ILM) dans divers niveaux de stockage sur flash, disque et bande, localement ou à distance. Les politiques automatisées permettent aux administrateurs de définir où, quand et sur quels supports les données (ou les métadonnées) seront stockées pour maximiser la performance des charges de travail et réduire au minimum les coûts de stockage globaux. Par exemple, les systèmes flash à faible temps d'attente peuvent offrir le meilleur rapport prix-performance pour les petits volumes de données très utilisés et les données très utilisées pour des charges de travail intensives en E-S; au contraire, les bandes LTFSS offrent le meilleur rapport prix-performance pour les grandes quantités d'ensembles de données moins utilisés qui sont prêts pour un archivage à long terme.

IBM Elastic Storage Server (ESS) est une mise en œuvre de stockage du logiciel IBM Spectrum Scale qui a été intégrée aux serveurs à processeur IBM POWER8 et aux unités de disque IBM. ESS permet de répondre aux demandes des charges de travail de mégadonnées dans la recherche biomédicale tout en offrant tous les avantages de Spectrum Scale. Le traitement multifilière, la vaste bande passante de mémoire et la grande taille de cache offerts par les processeurs IBM POWER8 augmentent le débit de données pour les charges de travail intensives en E-S. De plus, les schémas de protection RAID logiciels dégroupés qui sont installés avec ESS réduisent grandement la durée des reconstructions lors des défaillances de disque.

### **Architecture de nuage hybride construite sur IBM Spectrum Computing et IBM Spectrum Scale**

Les organisations de santé et de sciences de la vie cherchent à diminuer leurs dépenses en capital, à gérer plus facilement leurs TI et à profiter de meilleures fonctions de partage des données et de collaboration externe. Un grand nombre d'entre elles envisagent l'utilisation de ressources infonuagiques à hautes performances pour prendre en charge au moins une partie de leur charge de travail. Les couches de gestion des charges de travail et des données de l'architecture de référence d'IBM permettent aux systèmes IBM de mettre en œuvre des architectures d'entreprise pour les charges de travail de mégadonnées dans des environnements de nuage hybride. L'intégration étroite qui existe entre IBM Spectrum Computing et IBM Spectrum Scale, comme on le voit à la Figure 3, soutient le transfert transparent d'applications pouvant utiliser le nuage entre les serveurs sur site traditionnels et les nuages privés ou publics hors site. Les administrateurs des TI peuvent créer des politiques qui définissent l'utilisation relative des environnements physiques



et virtuels de la façon qui répond le mieux aux besoins techniques, financiers, d'affaires et de réglementation de leur organisation. Pour les cliniciens et les scientifiques biomédicaux qui mettent l'accent sur l'exécution d'applications d'analyse, le transfert de charges de travail dans ces environnements est complètement transparent et crée une expérience utilisateur sans interruption.

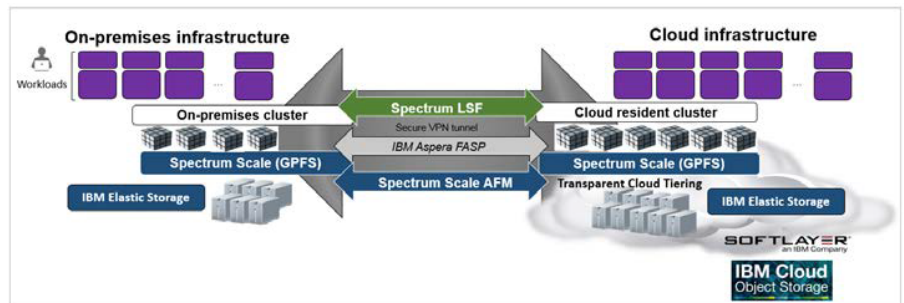


Figure 3. Une solution en nuage hybride

La Figure 3 présente les autres offres d'IBM qui peuvent aider au stockage, au transfert et à la gestion des données biomédicales. Le nuage IBM SoftLayer offre un environnement à hautes performances construit sur les logiciels Spectrum Computing et Spectrum Scale.

IBM Cloud Object Storage, offre maintenant disponible par l'entremise de l'acquisition de Cleversafe, fournit aux organisations une plateforme de stockage axée sur les objets hautement sécurisée et massivement évolutive. Cloud Object Storage peut être déployé sur site ou dans un environnement infonuagique IBM SoftLayer.

La technologie IBM Aspera Fast Adaptive Secure Protocol (IBM FASP) accélère le transfert des grands fichiers et ensembles de données – structurées ou non – sur une infrastructure existante de grand réseau (WAN). Aspera accélère le transfert de données même dans les lieux éloignés où les réseaux sont peu performants, et ce, de manière prévisible, fiable et sûre, quels que soient la taille du fichier, la distance à couvrir et l'état du réseau. IBM Aspera est utilisé actuellement dans le secteur des soins de santé et des sciences de la vie pour transférer des séquences brutes de génome humain (environ 200 Go par génome) et d'autres gros fichiers au sein d'une communauté dispersée géographiquement.

### Traitement accéléré avec les systèmes à processeur IBM POWER8

Le groupe Systèmes IBM est déterminé à offrir des performances supérieures pour les charges de travail computationnelles les plus difficiles. La stratégie d'IBM pour réaliser cet objectif est fondée sur le fait que les systèmes HPC de la prochaine génération devront soutenir des charges de travail à données intensives en plus des charges de travail à calculs intensifs. À mesure que convergent les besoins de l'informatique à hautes performances traditionnelle et ceux plus nouveaux de l'analytique des mégadonnées, le débit du système dépend non seulement des améliorations de la performance des E-S au niveau de l'unité centrale de traitement (CPU), mais aussi de la réduction au minimum des transferts de données au sein de l'architecture. Il

dépend aussi du resserrement de l'intégration du CPU avec les accélérateurs matériels, tels que les unités de traitement graphique (GPU) et les FPGA (Field Programmable Gate Arrays), qui servent souvent à accélérer considérablement les applications des utilisateurs.

**OpenPOWER Foundation :** En 2013, IBM a commencé à proposer des licences de code source ouvert pour les produits liés à l'architecture IBM Power afin de créer une solution pour remplacer les solutions propriétaires et encourager l'innovation en informatique. IBM a lancé la création de l'OpenPOWER Foundation, communauté technique regroupant plus de 130 entreprises et universités collaborant au développement de solutions fondées sur les processeurs IBM POWER qui répondent à des besoins d'affaires spécifiques. Les innovations nées des collaborations OpenPOWER comprennent des systèmes personnalisés pour l'accélération des charges de travail utilisant le GPU, les FPGA et les E-S évoluées.

**Processeurs OpenPOWER :** Les systèmes OpenPOWER sont conçus pour assumer les charges de travail intensives en calculs et en données, comme l'apprentissage machine et les réseaux profonds de neurones. Le processeur IBM POWER8 présente le traitement multifilière simultané (SMT), qui permet jusqu'à huit fils matériels à partir d'un seul cœur physique. Il fait appel au sous-système de mémoire le plus évolué qui soit pour obtenir des performances de pointe – en utilisant un grand nombre de caches de mémoire sur puce et hors puce. Cette conception réduit le temps d'attente de la mémoire et génère des bandes passantes très élevées pour la mémoire et les E-S du système.

**Progrès récents :** Les fonctions suivantes ont été conçues pour augmenter les performances des systèmes POWER8 :

- **Coherent Accelerator Processor Interface (CAPI)**, qui permet aux FPGA et aux autres accélérateurs connectés à un logement PCIe d'accéder directement au bus du processeur avec un faible temps d'attente.
- **NVLink**, interconnecteur de processeur à haute vitesse qui peut être utilisé pour connecter les GPU à des CPU ou à des GPU (Figure 4). NVLink, qui a été développé en partenariat avec NVIDIA, élimine le goulot d'étranglement du PCI, ce qui permet des performances au moins cinq fois supérieures à celles du PCIe. NVLink optimise l'intégration logique de GPU multiples ainsi que de cœurs de CPU et de GPU. Grâce à ces connexions, chaque GPU a un accès direct à la fois à la mémoire du processeur hôte et à celle du GPU frère. Ce modèle de calculs GPU diminue considérablement la complexité de programmation associée aux calculs GPU de base et aux calculs multi-GPU.

**Applications :** Les serveurs fondés sur le POWER8 prennent en charge les distributions standards de Linux, ce qui facilite le portage de codes existants vers la plateforme. Les applications privilégiées du domaine de la santé et des sciences de la vie – comme GATK, BWA, SAMtools, BLAST, MuTect2 et tranSMART – ont déjà été optimisées sur IBM PowerLinux. Les bio-informaticiens des Systemes IBM et les partenaires commerciaux

IBM de la santé et des sciences de la vie effectuent activement le portage et l'optimisation des performances d'autres codes de recherche biomédicale sur POWER. Plus de 100 applications à code source ouvert sont optimisées sur le POWER8.

Dans les domaines de recherche pertinents pour les soins de santé personnalisés, comme la génomique et la bio-informatique, l'adoption de charges de travail optimisées GPU a été lente; cependant, ces charges de travail gagnent du terrain. Un nombre croissant d'organisations menant des recherches biomédicales appliquent des techniques d'apprentissage profond pour découvrir des tendances prédictives dans de très grands ensembles de données complexes, souvent non structurées, comme les images biomédicales et les signaux physiologiques, qui varient avec le temps. Pour soutenir de telles charges de travail, IBM et NVIDIA collaborent sur PowerAI, qui est un nouvel ensemble d'outils d'apprentissage profond. PowerAI est une plateforme facile à déployer qui offre des cadres d'apprentissage profond répandus – notamment Caffe, Torch, et Theano – au sein de l'architecture IBM Power. IBM a optimisé toutes les distributions de ces logiciels d'apprentissage profond de façon à ce qu'elles profitent de la bande passante élevée offerte par le processeur IBM POWER8 et de l'interconnexion NVIDIA NVLink. L'ensemble d'outils utilise aussi les bibliothèques NVIDIA GPUDL, y compris cuDNN, cuBLAS et NCCL, pour offrir l'accélération multi-GPU sur les serveurs IBM.

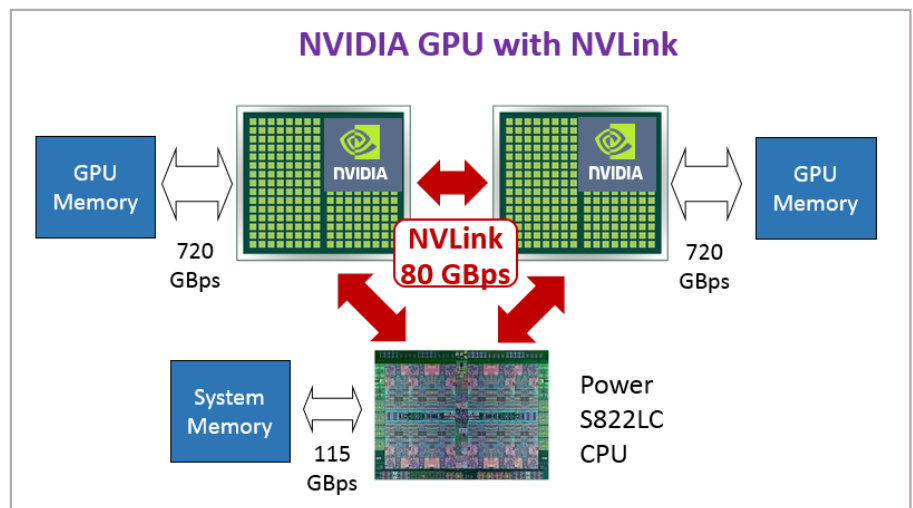


Figure 4. Le processeur POWER8 avec NVIDIA NVLink

## Un ensemble de partenaires commerciaux dans le domaine de la santé et des sciences de la vie

Les experts du secteur de la santé et des sciences de la vie, tant à l'intérieur qu'à l'extérieur d'IBM, contribuent à l'optimisation des algorithmes computationnels, des applications pour utilisateurs et des référentiels de données que les cliniciens et les chercheurs peuvent exécuter sur l'infrastructure des Systèmes IBM (dans la Figure 1, reportez-vous aux boîtes en violet de l'architecture de référence d'IBM). Pour répondre entièrement aux besoins des clients ayant des projets de soins personnalisés, le groupe

Systemes IBM a déjà établi des partenariats d'affaires et de développement avec des leaders du domaine, dont certains sont décrits dans cette section.

**IBM Watson :** La division IBM Watson se consacre au développement de systèmes cognitifs qui peuvent découvrir des tendances répétables et prédictives dans les données les plus complexes, par exemple, le langage naturel, les sons environnementaux et les images. Ces tendances peuvent être utilisées pour avoir un effet positif sur la vie des gens. La sous-division Santé Watson IBM est responsable du développement d'applications de soins de santé et de recherche biomédicale qui appliquent les fonctions de base de Watson en matière de traitement du langage naturel et d'apprentissage machine à l'analyse des sources précieuses de données personnelles comprenant les dossiers médicaux électroniques et les images biomédicales.

**Databiology :** Databiology for Enterprise (DBE) est une plateforme évolutive de gestion de l'information des technologies de biologie moléculaire qui permet la gestion centralisée et sécurisée de ce type de données et leur analyse dans le cadre de divers projets répartis mondialement. DBE offre des outils de requête, d'extraction et de visualisation des données qui simplifient l'expérience utilisateur. Pour soutenir les exigences de conformité relatives à la provenance complète des données et la science reproductible, DBE fournit le suivi automatisé des données et la capacité de traiter à nouveau toute charge de travail analytique en utilisant les données, les logiciels et les paramètres de la charge de travail d'origine. Databiology a intégré DBE aux IBM Power Systems, à Spectrum LSF et à Spectrum Scale pour améliorer la gestion des charges de travail, des ressources et du cycle de vie des données, sur site, hors site et dans des environnements hybrides. DBE est également intégré au logiciel IBM Aspera pour assurer le transfert sécurisé à grande vitesse d'ensembles de données génomiques autour du globe, et avec les environnements de traitement IBM POWER. Reportez-vous à [www.databiology.com](http://www.databiology.com)

**Lab7 :** Enterprise Science Platform (ESP) de Lab7 constitue pour les environnements de laboratoire de petits à intermédiaires une plateforme de gestion des données souple, évolutive, centralisée et définie par l'utilisateur permettant de faire le suivi d'échantillons de laboratoire, de saisir la provenance des données, de traiter des données, de produire des rapports et de gérer les flux de travaux de façon continue. Lab7 a optimisé ESP sur les IBM Power Systems, et gère et met à jour continuellement le service Web BioBuilds (voir à [www.biobuilds.org](http://www.biobuilds.org)), qui offre le déploiement clé en main d'outils bio-informatiques et le soutien de la communauté sur IBM POWER. Les outils pris en charge sur IBM POWER comprennent plusieurs versions de Bowtie, BWA, Picard, Short Oligonucleotide Analysis Package (SOAP), Isaac, PLINK, Bioconductor et de nombreux autres produits. Reportez-vous à [www.lab7.io](http://www.lab7.io)

**RowAnalytics :** Synomics Studio de RowAnalytics constitue une solution de remplacement ultrarapide et hautement évolutive des études traditionnelles d'association à l'échelle du génome. Il permet d'analyser des ensembles de données à grande échelle contenant des données génomiques et cliniques pour trouver des associations entre les combinaisons génomiques variantes de haut niveau et les résultats cliniques. Synomics a recours à des algorithmes massivement parallélisés répartis dans plusieurs unités de calcul GPU, et

demande de grandes quantités de puissance de traitement très dense. On l'a vu offrir des performances supérieures sur une architecture de calculs fondée sur l'architecture de référence d'IBM pour la santé et les sciences de la vie, qui comprend les IBM Power Systems, IBM Spectrum LSF et IBM Spectrum Scale. Reportez-vous à [www.rowanalytics.com](http://www.rowanalytics.com)

**tranSMART :** TranSMART est un système d'entrepôt de données et de gestion des connaissances à code source ouvert qui permet d'intégrer, d'analyser et de partager des données cliniques, génomiques et d'expression des gènes pour de vastes populations de clients. Fondé à l'origine sur le modèle de données pour Informatics for Integrating Biology and the Bedside (i2b2) parrainé par le National Institutes of Health (NIH) et reconnu au niveau international (voir [www.i2b2.org](http://www.i2b2.org), tranSMART a été largement utilisé dans l'industrie pharmaceutique et dans le cadre d'initiatives de recherches translationnelles publiques et privées à grande échelle. Depuis 2015, le groupe Systèmes IBM collabore avec le Data Science Institute de l'Imperial College de Londres pour étudier l'application de nouvelles architectures de traitement et de stockage pouvant améliorer les performances et l'évolutivité de la plateforme tranSMART. Jusqu'à maintenant, l'application de tranSMART à l'IBM POWER8 avec Elastic Storage Server a permis d'obtenir des améliorations de performances des systèmes Intel pendant l'ingestion et l'analyse des données. Reportez-vous à [www.transmartfoundation.org](http://www.transmartfoundation.org)

**WhamTech :** SmartData Fabric (SDF) de WhamTech est une couche de gestion de données distribuées qui s'intègre aux infrastructures informatiques existantes et peut prendre en charge la virtualisation, l'intégration et la fédération de données hautement sécurisées, ainsi que l'analytique pour des silos de données hétérogènes communs aux organisations de santé. Quels que soient le type, le format, la qualité ou la structure des données de la source – qu'elles comprennent des mégadonnées, des bases de données NoSQL, des bases de données relationnelles, des fichiers, des documents de bureau, des courriels ou des appareils IoT (Internet des objets) –, les utilisateurs de SDF peuvent accéder aux données de sources indexées sans copier ni transférer les données d'un endroit à un autre. Les recherches sur les données peuvent être effectuées par rapport aux données indexées, quel que soit leur emplacement, à l'aide du traitement parallèle à hautes performances avec la prise en charge de l'évolutivité extrême. WhamTech participe à plusieurs projets de soins de santé aux États-Unis, au Royaume-Uni et en Australie, et optimise actuellement son logiciel pour les systèmes IBM POWER8 et l'architecture de référence d'IBM pour la santé et les sciences de la vie. Reportez-vous à [www.whamtech.com](http://www.whamtech.com)

## Résumé

L'architecture IBM de référence pour la santé et les sciences de la vie est constituée de composantes d'infrastructure clés de la gamme IBM de produits de traitement et de stockage à haute performance. Elle soutient un groupe croissant de partenaires importants de l'industrie. Elle définit une plateforme très souple et rentable permettant de gérer, de stocker, de partager, d'intégrer et d'analyser les mégadonnées dans le cadre de budgets informatiques limités. Les organisations TI peuvent utiliser cette architecture comme guide général leur permettant de surmonter les défis de la gestion des données et des goulots d'étranglement de traitement qu'il faut souvent relever pour les initiatives de soins de santé personnalisés et d'autres charges biomédicales de traitement à données intensives.

## Obtenez d'autres informations

Pour en savoir plus sur les offres de traitement et de stockage à hautes performances qui constituent l'architecture de référence d'IBM pour la santé et les sciences de la vie, contactez votre représentant ou votre partenaire commercial IBM, ou visitez les sites Web suivants :

- IBM Spectrum Computing  
<http://www.ibm.com/systems/spectrum-computing/>
- IBM Spectrum Scale  
<http://www.ibm.com/systems/storage/spectrum/scale/>
- IBM Power Systems  
<http://www.ibm.com/systems/power/>
- IBM SoftLayer  
<http://www.softlayer.com/>
- IBM Cloud Object Storage  
<https://www.ibm.com/cloud-computing/products/storage/object-storage/cloud/>
- IBM Aspera  
<https://www.ibm.com/software/info/aspera/>

## À propos des auteurs

**Jane Vu, M.D., Ph. D.**, est l'architecte sectorielle mondiale pour la santé et les sciences de la vie du groupe Systèmes IBM. Ayant plus de 25 ans d'expérience en médecine clinique, recherche biomédicale, analytique avancée, ingénierie des systèmes, TI d'entreprise et services-conseils de gestion, M<sup>me</sup> Vu travaille étroitement avec les spécialistes techniques, les partenaires de l'industrie et les clients internationaux d'IBM afin d'offrir les meilleures solutions de traitement et de stockage à hautes performances pour les applications de soins de santé et de recherche biomédicale. Elle se spécialise dans le soutien d'organisations qui mènent des programmes de soins de santé personnalisés et de médecine translationnelle.

**Kathy Tzeng, Ph. D.** groupe Systèmes IBM. Depuis qu'elle s'est jointe à IBM en 2001, M<sup>me</sup> Tzeng travaille avec des équipes techniques d'IBM, de partenaires de l'industrie et de communautés de code source ouvert pour créer des applications de sciences de la vie et en améliorer les performances sur les solutions IBM. Elle a obtenu des brevets et publié des livres rouges IBM, des articles techniques, des chapitres de livre et des articles de journaux spécialisés revus par les pairs.

**Janis Landry-Lane** est la directrice exécutive mondiale des ventes pour les soins de santé et les sciences de la vie dans le groupe Systèmes IBM. Possédant plus de 20 ans d'expérience en informatique hautes performances, elle a collaboré avec de nombreux clients en recherche et développement. M<sup>me</sup> Landry-Lane travaille de près avec une équipe multidisciplinaire de partenaires d'IBM et externes afin de répondre aux exigences techniques des clients en matière de systèmes informatiques à hautes performances. Son équipe fournit des solutions évolutives et extensibles qui sont conçues pour croître dans des environnements sur site et s'intégrer dans des environnements infonuagiques au besoin.



---

© Copyright IBM Corporation, 2017  
© Copyright IBM Canada Ltée, 2017  
IBM Systems  
3039 Cornwallis Road  
RTP, NC 27709

Produit au Canada

IBM, le logo IBM et [ibm.com](http://ibm.com) sont des marques déposées d'International Business Machines Corporation, enregistrées dans de nombreux pays. Si ces marques et d'autres marques d'IBM sont suivies du symbole <sup>MD</sup> ou <sup>MC</sup> à leur première occurrence dans un document, cela signifie qu'il s'agit d'une marque déposée ou de common law aux États-Unis, qui appartenait à IBM au moment où l'information a été publiée. Ces marques peuvent aussi être déposées ou être des marques de common law dans d'autres pays. La liste à jour des marques d'IBM est disponible sur le Web sous «Copyright and trademark information», à [ibm.com/legal/copytrade.shtml](http://ibm.com/legal/copytrade.shtml).

Tous les autres noms de société, de produit ou de service peuvent être des marques de commerce ou des marques de service appartenant à leurs détenteurs respectifs.

La présente publication peut faire référence à des produits ou à des services IBM non annoncés dans votre pays. Cela ne signifie pas qu'IBM ait l'intention de les y annoncer.



Veillez recycler.