

宇宙科学と大規模可視化

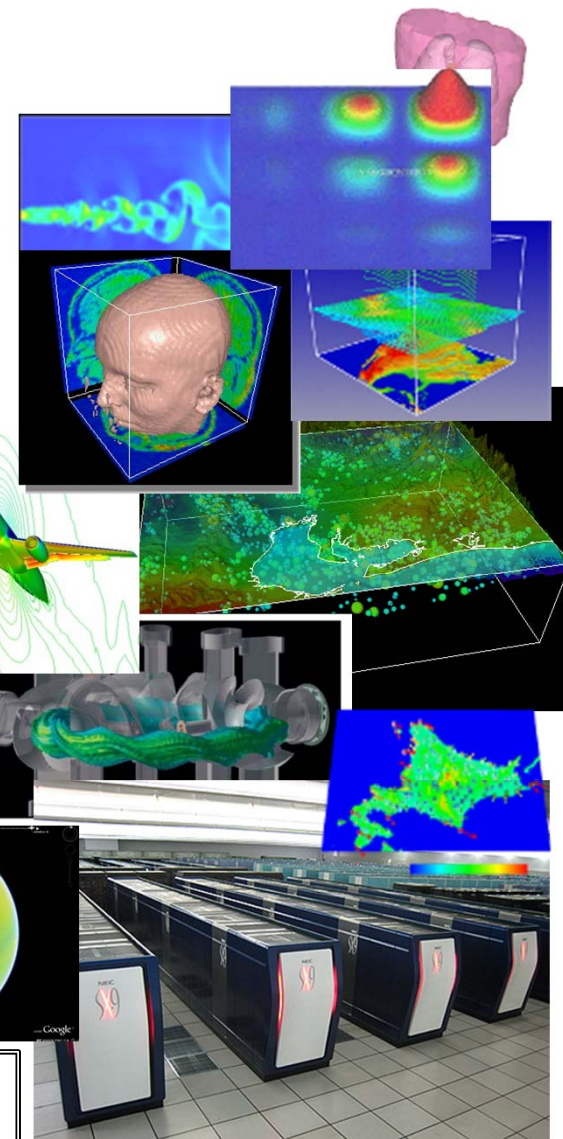
サイエンスクラウドで初めて見えた現象

世界No.1になった「京」コンピュータは一度の計算でPB(ペタバイト)データを計算します。これらのデータを安全に保存する方法はあるのか、大規模データをどのように処理するか、そこから得られる新しい科学は何かなど、NICTサイエンスクラウドが目指すテーマとその成果についてご紹介します。

平成24年12月3日
情報通信研究機構
統合データシステム研究開発室
村田 健史(たけし)



①理論



②観測・実験

科学研究の手法 主要な3つの方法

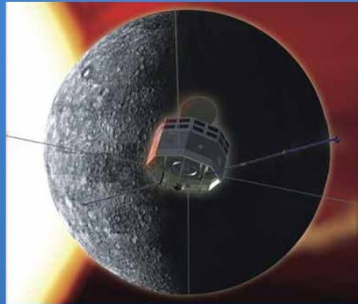
③シミュレーション

宇宙科学のメガサイエンスプロジェクト(1)

衛星観測と数値シミュレーション

科学衛星観測

Scientific Satellite Missions



<http://www.rish.kyoto-u.ac.jp>



<http://www.isas.jaxa.jp>

<http://www.isas.jaxa.jp>

100億円



ペタバイト

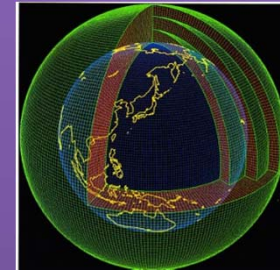
ビッグデータ

数値シミュレーション

Super Computer Projects



<https://aspara.asahi.com/blog/netn-owplus/entry/uSreeI7JOG>



http://www.venus.dti.ne.jp/~inoue-m/on_co2_i/f11_cb100.jpg



<http://pc.watch.impress.co.jp/img/pcw/docs/212/472/nec.jpg>

1,000億円

ペタバイト

ビッグデータ

別々に研究が進められている

メガサイエンスの衝突



<http://cosmos21.exblog.jp/>

第四の科学研究手法 (Jim Gray)

<http://research.microsoft.com/en-us/collaboration/fourthparadigm>

Microsoft
Research

Search Microsoft Research

Videos Projects Publications People Downloads

Home Our Research **Connections** Careers Hub
About Us Research in Action Opportunities Research Accelerators

Tell us what you think.
Take our 6-question site survey.

Microsoft Research Connections > The Fourth Paradigm: Data-Intensive Scientific Discovery

The Fourth Paradigm: Data-Intensive Scientific Discovery

Presenting the first broad look at the rapidly emerging field of data-intensive science

Increasingly, scientific breakthroughs will be powered by advanced computing capabilities that help researchers manipulate and explore massive datasets.

The speed at which any given scientific discipline advances will depend on how well its researchers collaborate with one another, and with technologists, in areas of eScience such as databases, workflow management, visualization, and cloud computing technologies.

In *The Fourth Paradigm: Data-Intensive Scientific Discovery*, the collection of essays expands on the vision of pioneering computer scientist Jim Gray for a new, fourth paradigm of discovery based on data-intensive science and offers insights into how it can be fully realized.

データ指向型科学

Download

- [Full text, low resolution](#) (6 MB)
- [Full text, high resolution](#) (93 MB)
- [By chapter and essay](#)

Purchase from Amazon.com

- [Paperback](#)
- [Kindle version](#)

In the News

- [Sailing on an Ocean of 0s and 1s](#) (*Science Magazine*)
- [A Deluge of Data Shapes a New Era in Computing](#) (*New York Times*)

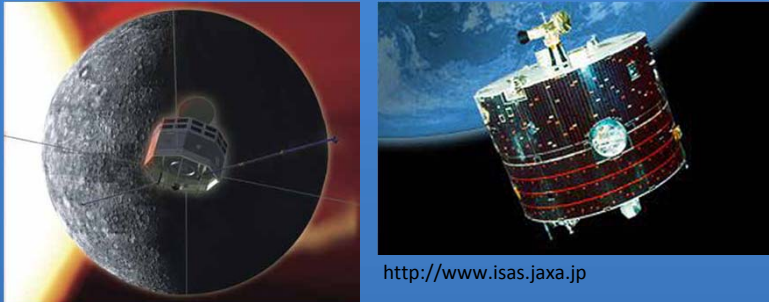


宇宙科学のメガサイエンスプロジェクト(2)

インフォマティクスの役割

科学衛星観測

Scientific Satellite Missions




<http://www.rish.kyoto-u.ac.jp>

<http://www.isas.jaxa.jp>

<http://www.isas.jaxa.jp>

100億円



数値シミュレーション

Super Computer Projects



<https://aspara.asahi.com/blog/netn-owplus/entry/uSreeI7JOG>

http://www.venus.dti.ne.jp/~inoue-m/on_co2_i/f11_cb100.jpg



<http://pc.watch.impress.co.jp/img/pcw/docs/212/472/nec.jpg>

1,000億円

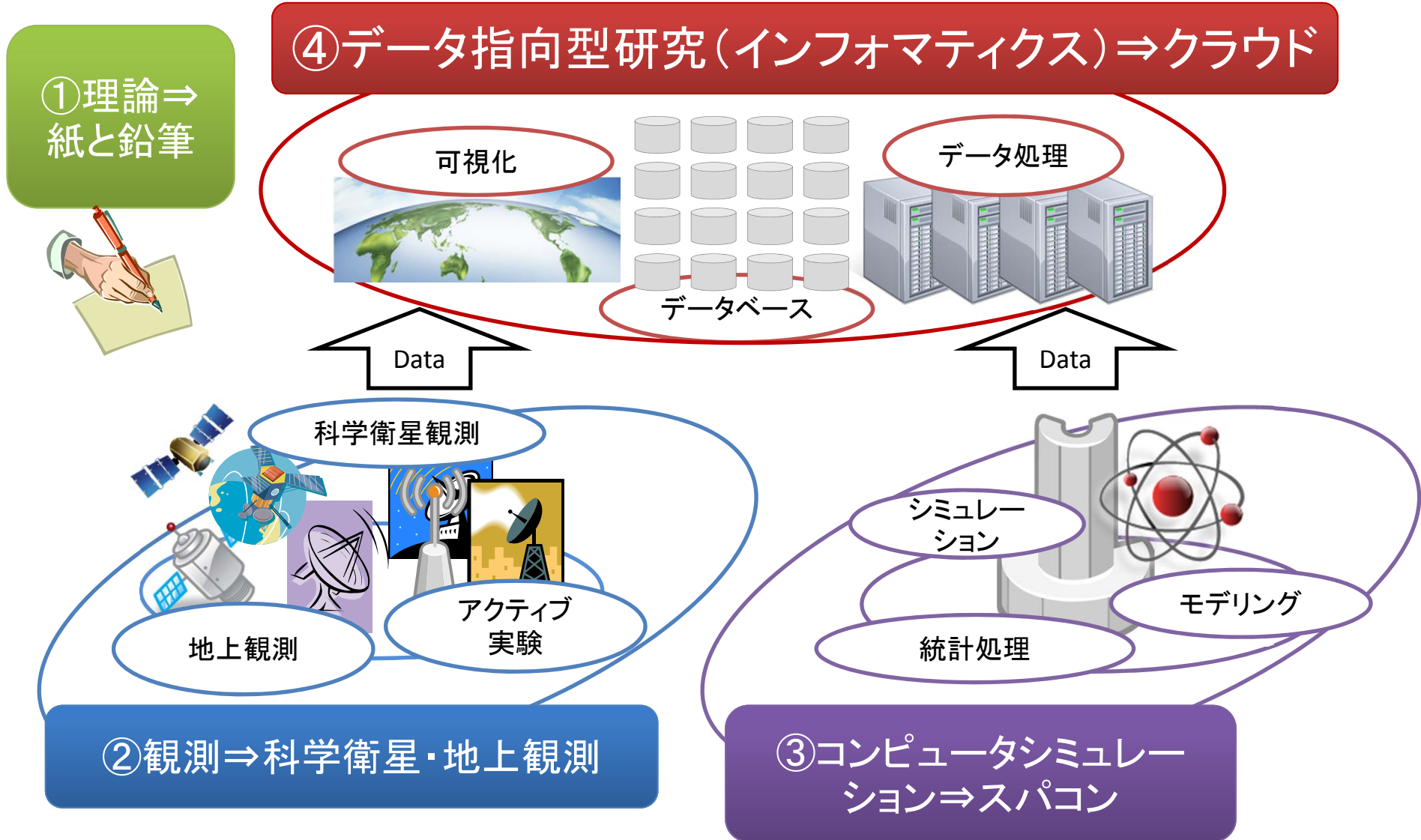
インフォマティクス

データ指向型科学

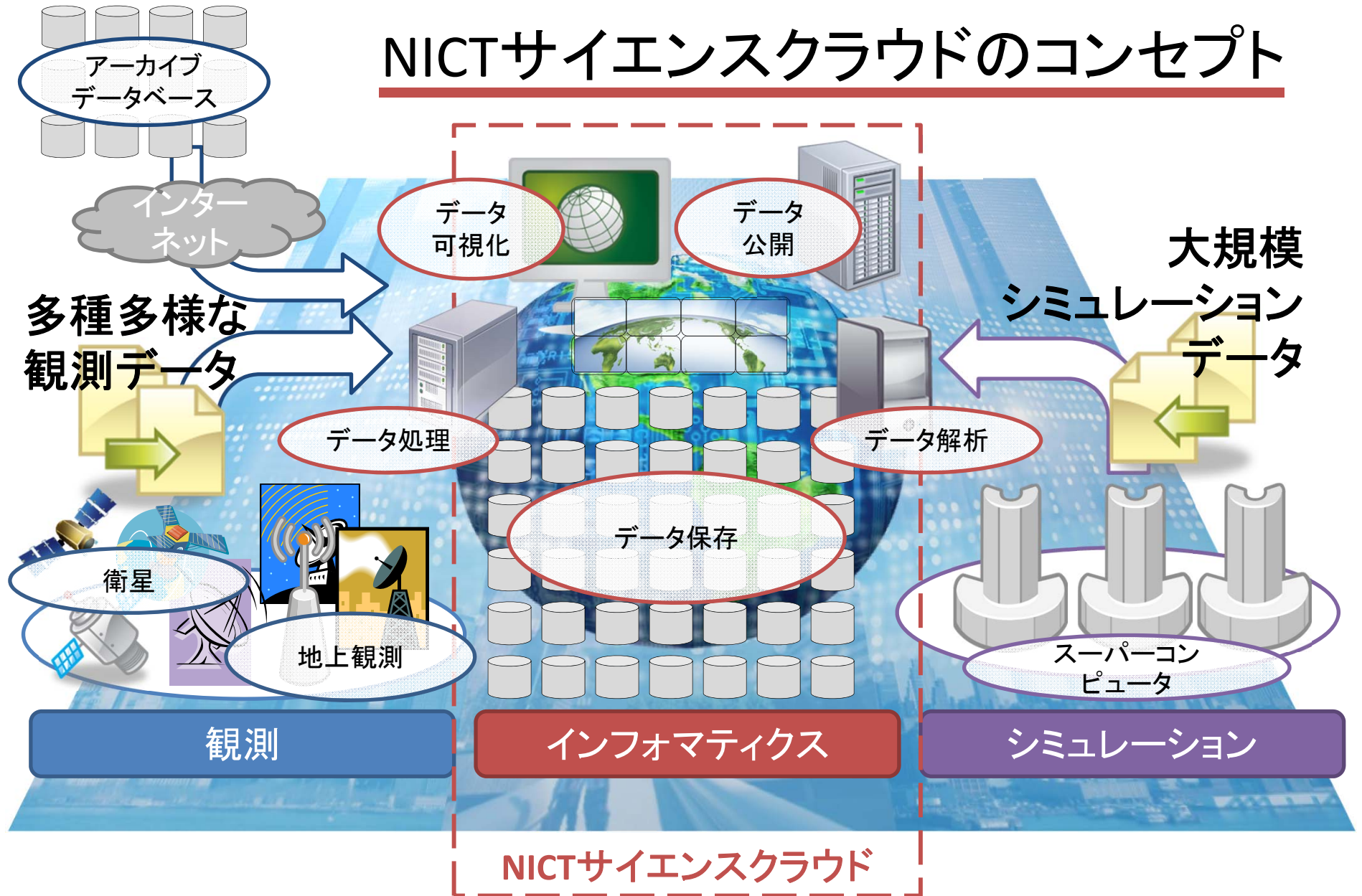
ペタバイト

ビッグデータ

科学を行う方法(宇宙科学の場合)



NICTサイエンスクラウドのコンセプト



あらゆるデータをクラウド上に！

世界の情報ストレージ総容量/増加する科学データ

世界のストレージ動向

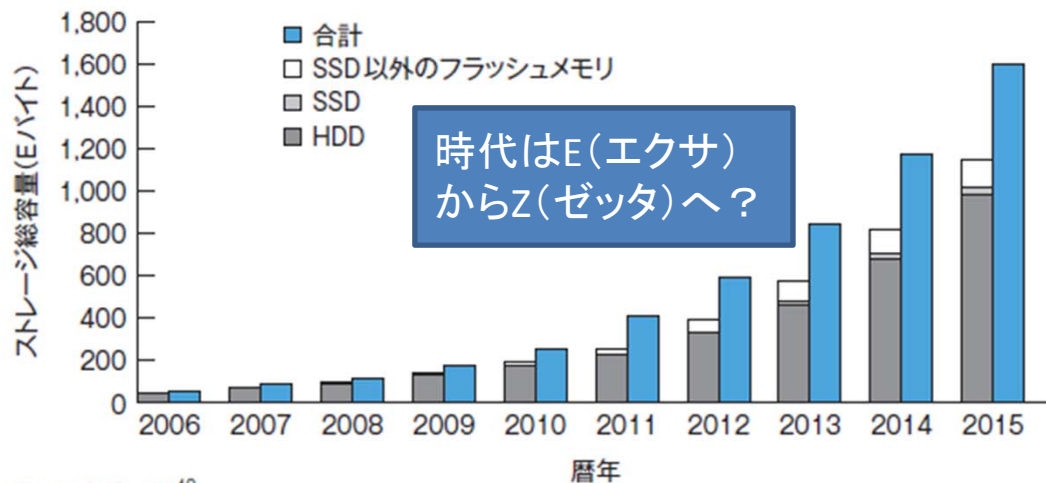
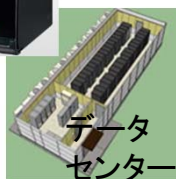
MB(メガバイト)

GB(ギガバイト)

TB(テラバイト)

PB(ペタバイト)

EB(エクサバイト)

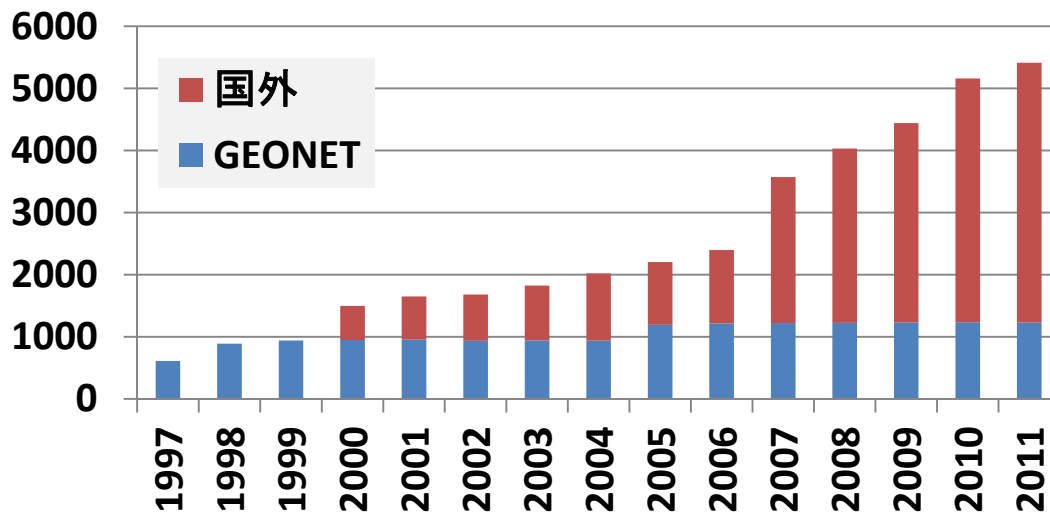
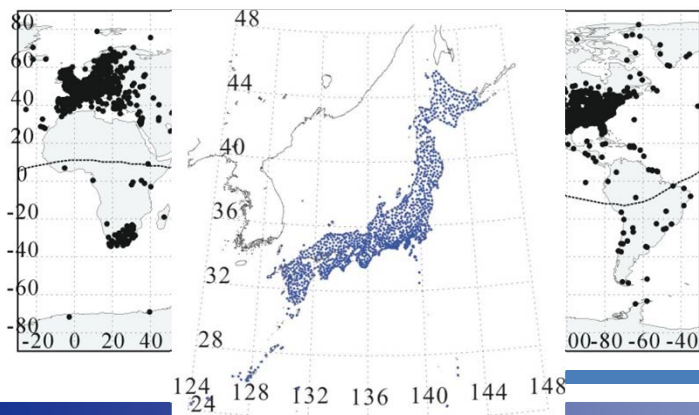


E: エкса, 10^{18}

http://www.toshiba.co.jp/tech/review/2011/08/66_08pdf/b02.pdf

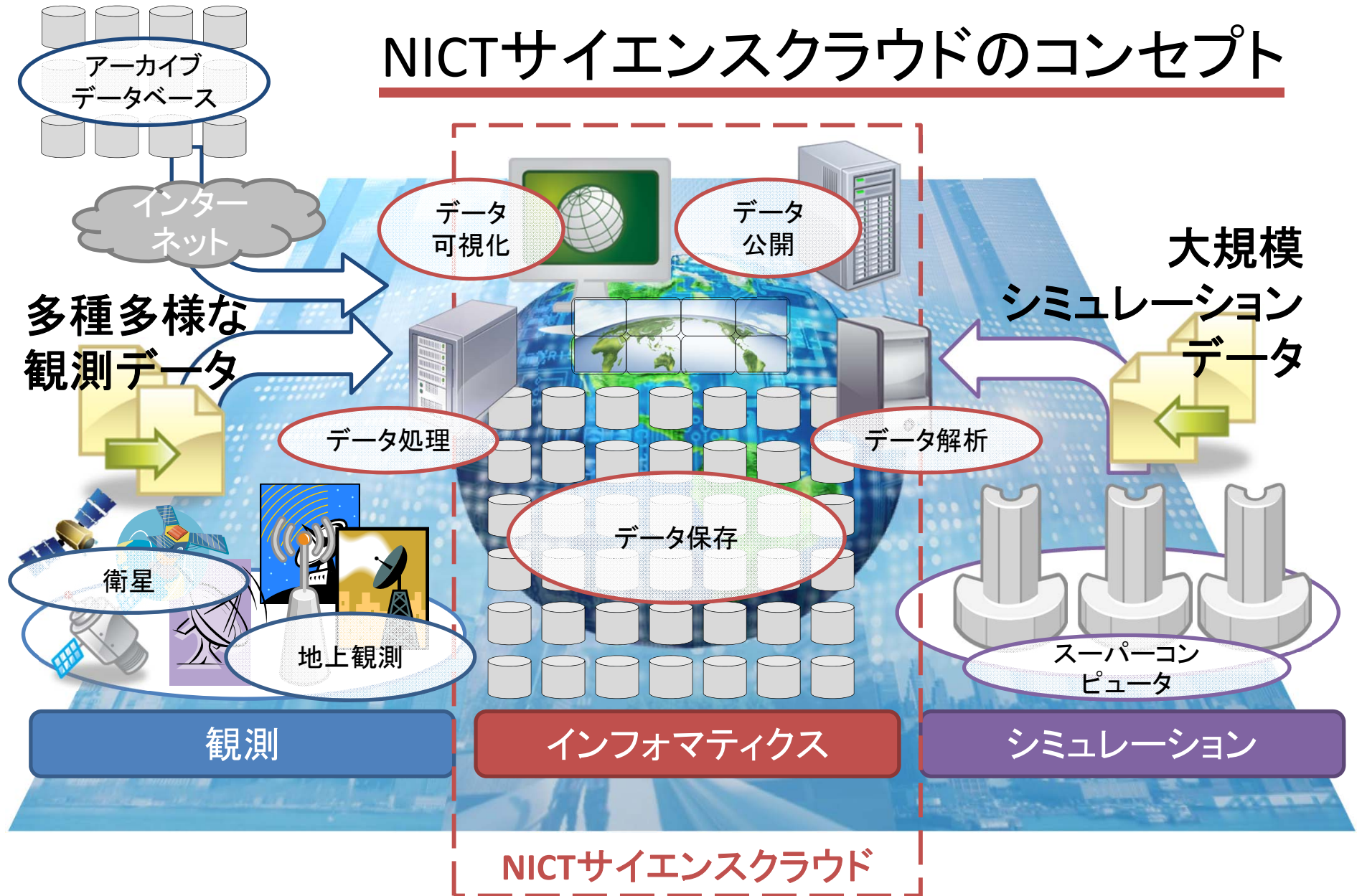
科学研究分野の例

GPS公開データ



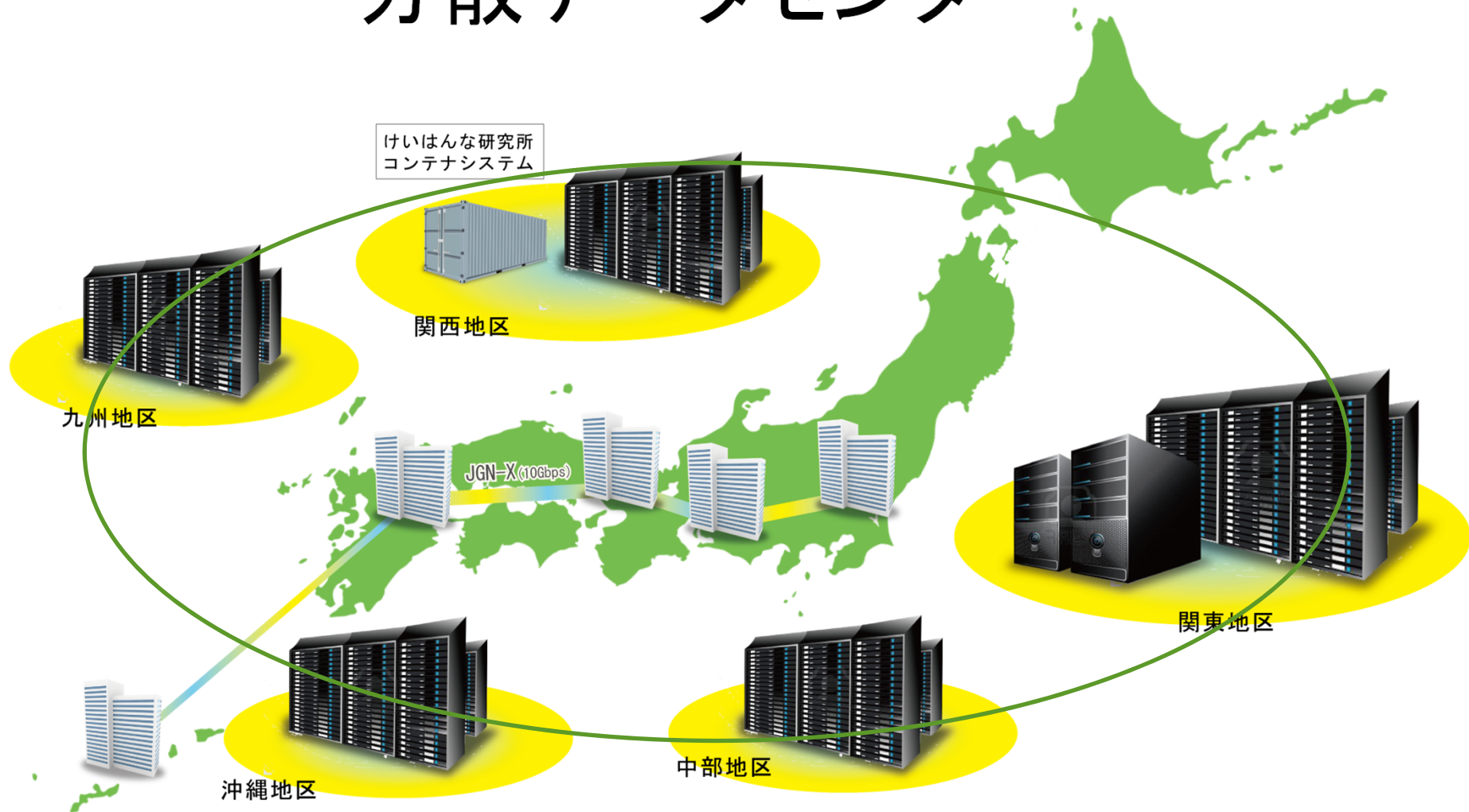
津川(NICT)らによる

NICTサイエンスクラウドのコンセプト

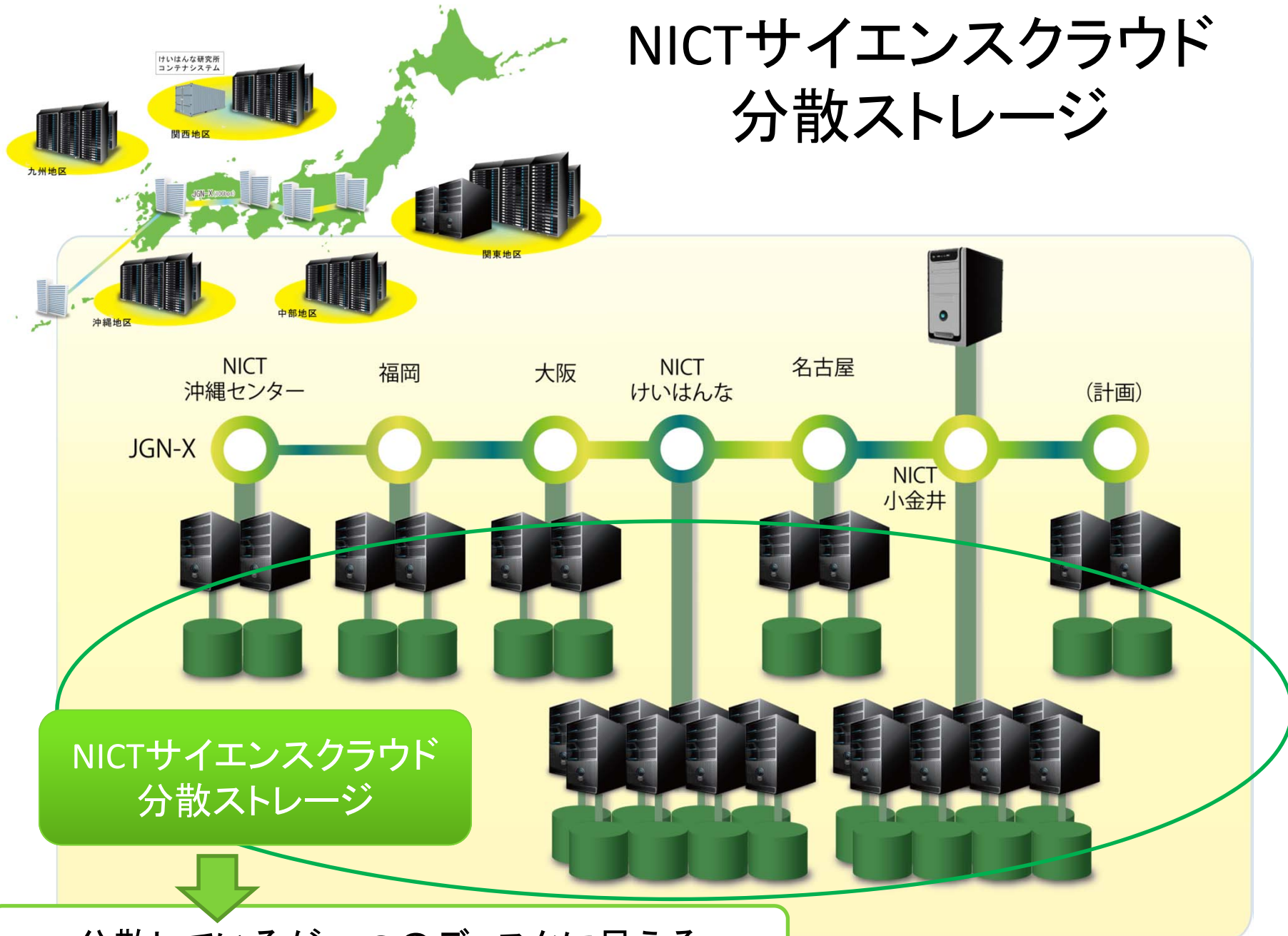


あらゆるデータをクラウド上に！

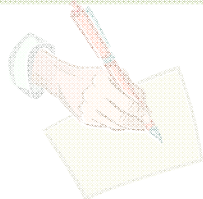
NICTサイエンスクラウド 分散データセンター



NICTサイエンスクラウド 分散ストレージ



①理論→紙と鉛筆



④データ指向型研究(インフォマティクス)⇒クラウド

データ指向型科学を実現するには...

どうやってデータを集めるか？

どうやってデータを保存するか？

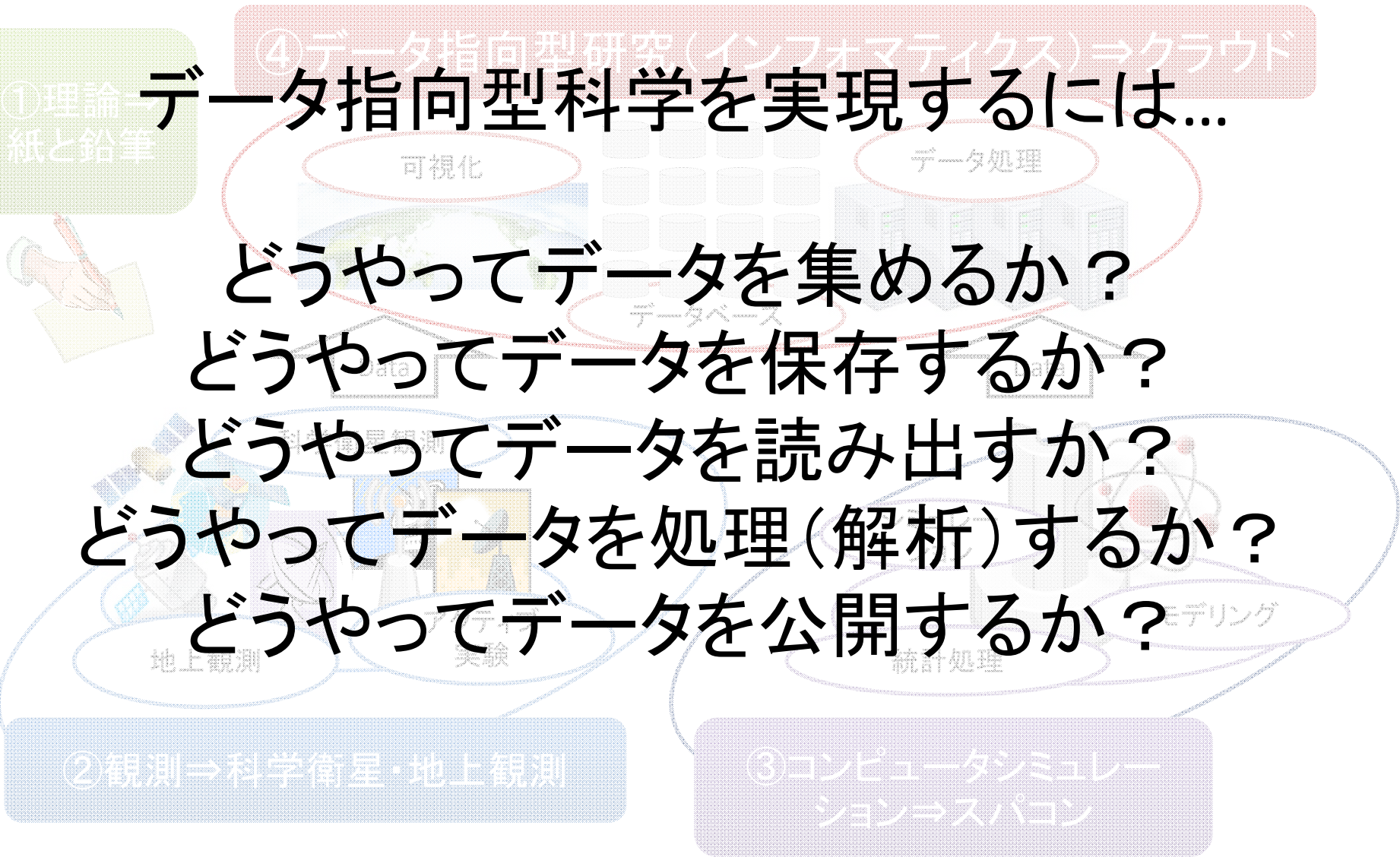
どうやってデータを読み出すか？

どうやってデータを処理(解析)するか？

どうやってデータを公開するか？

②観測⇒科学衛星・地上観測

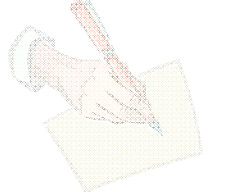
③コンピュータシミュレーション⇒スパコン



①理論→紙と鉛筆

データ指向型科学を実現するには...

④データ指向型研究(インフォマティクス)⇒クラウド



どうやってデータを集めるか？

どうやってデータを保存するか？

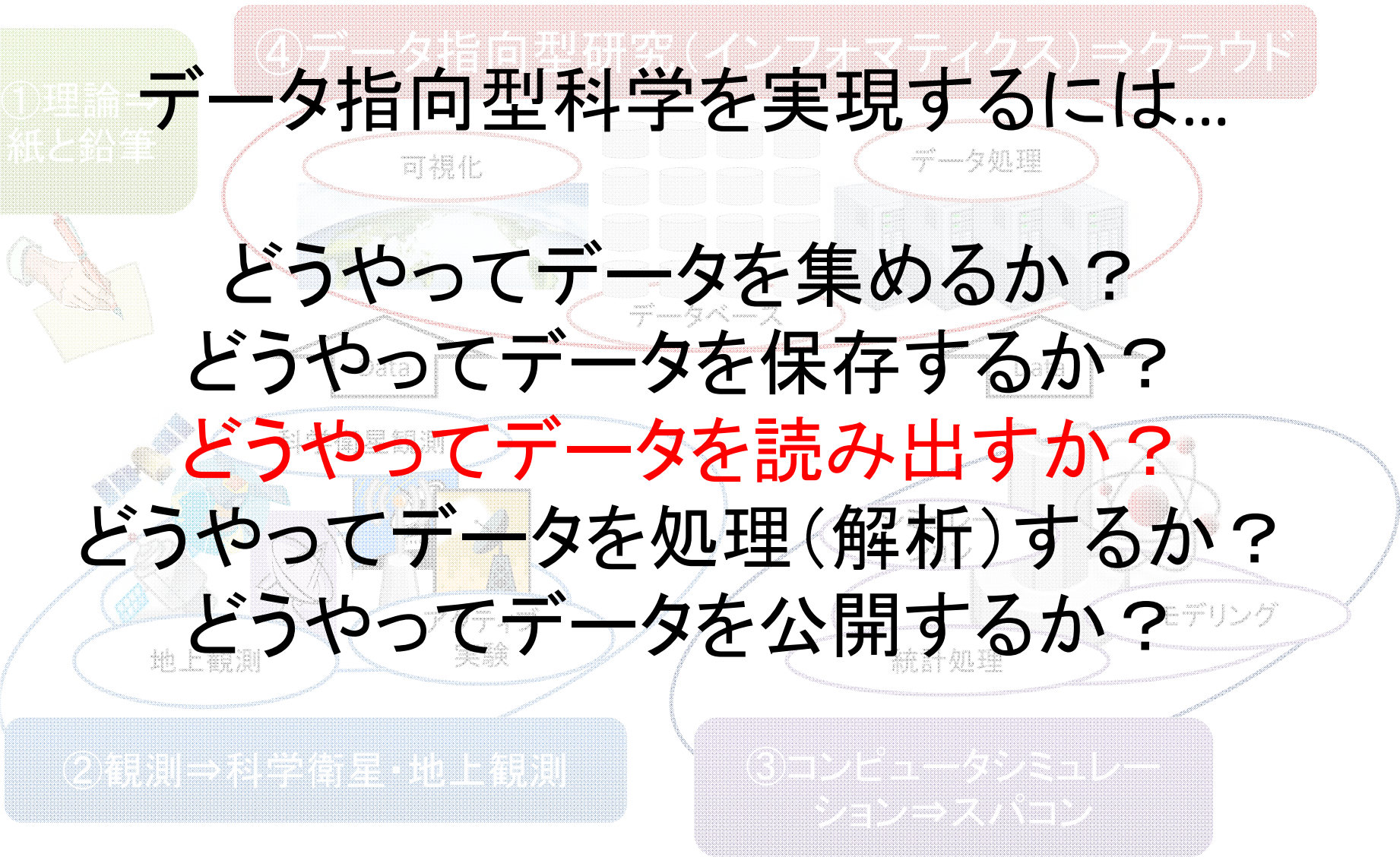
どうやってデータを読み出すか？

どうやってデータを処理(解析)するか？

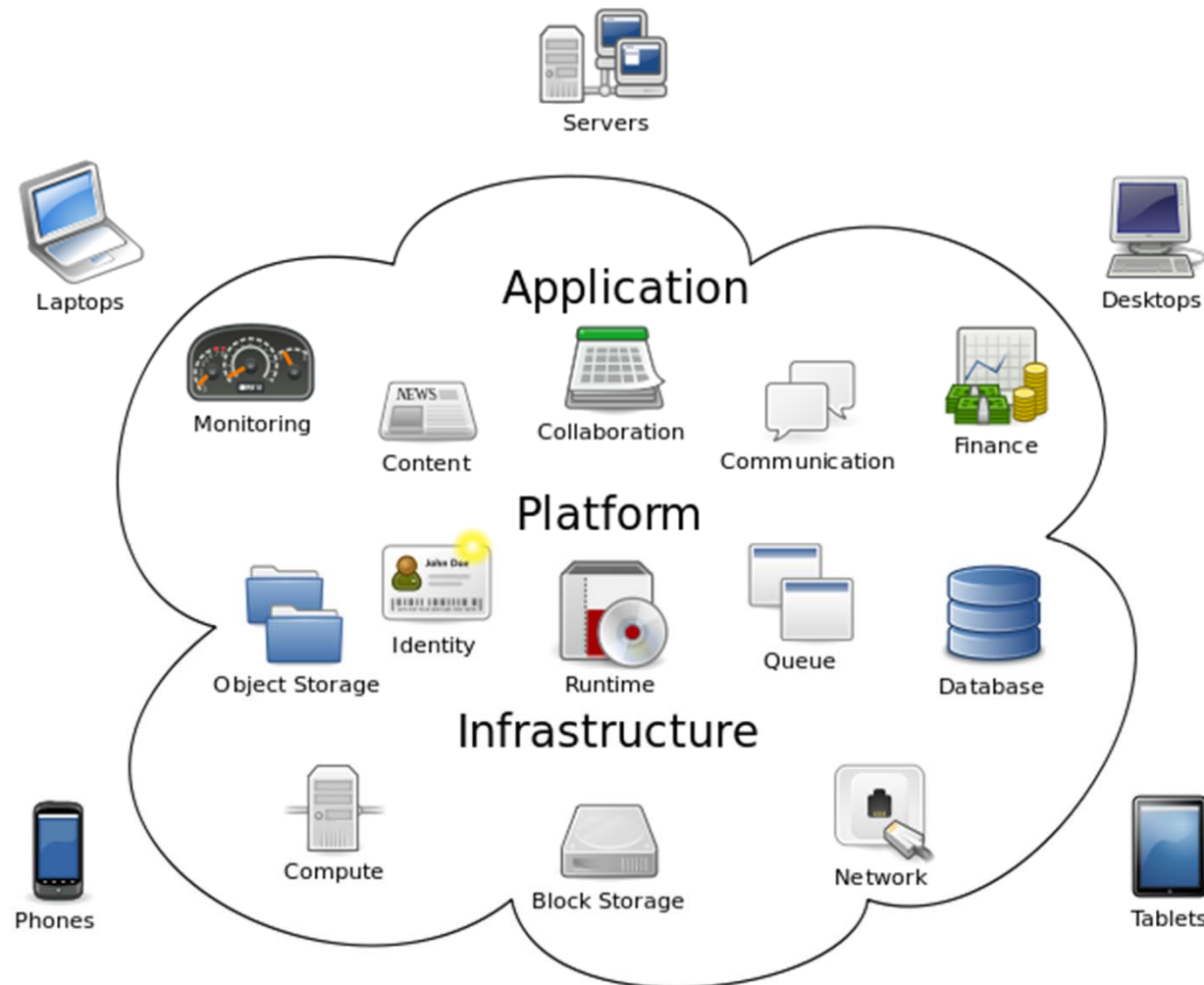
どうやってデータを公開するか？

②観測⇒科学衛星・地上観測

③コンピュータシミュレーション⇒スパコン



クラウドとは？（Wikipediaより）

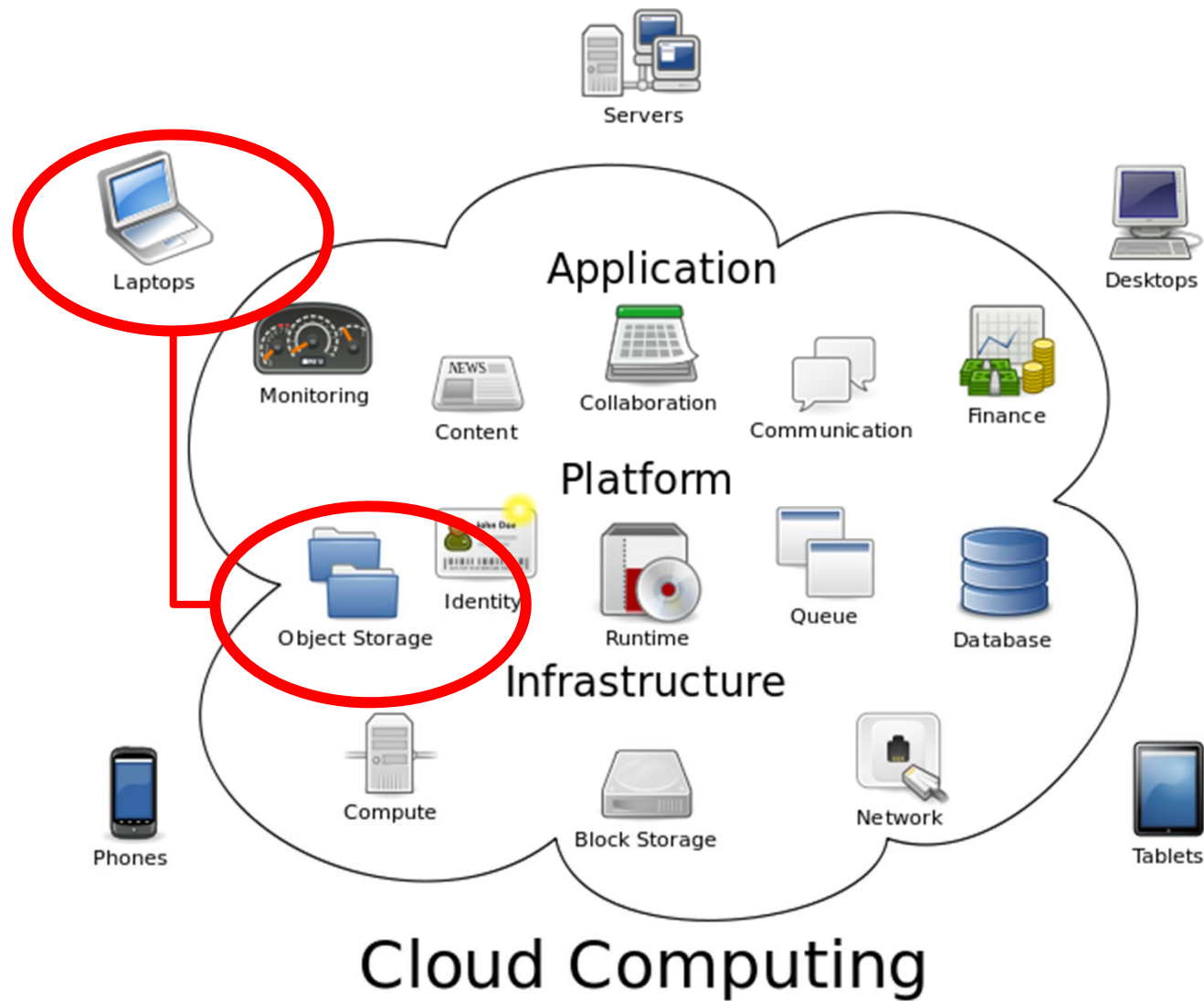


クラウドコンピューティング(英: cloud computing)とは、ネットワーク、特にインターネットをベースとしたコンピュータの利用形態である。ユーザーはコンピュータ処理をネットワーク経由で、サービスとして利用する。

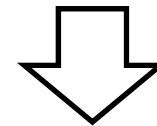
Cloud Computing

http://ja.wikipedia.org/wiki/%E3%83%95%E3%82%A1%E3%82%A4%E3%83%AB:Cloud_computing.svg

NICTサイエンスクラウド利用例



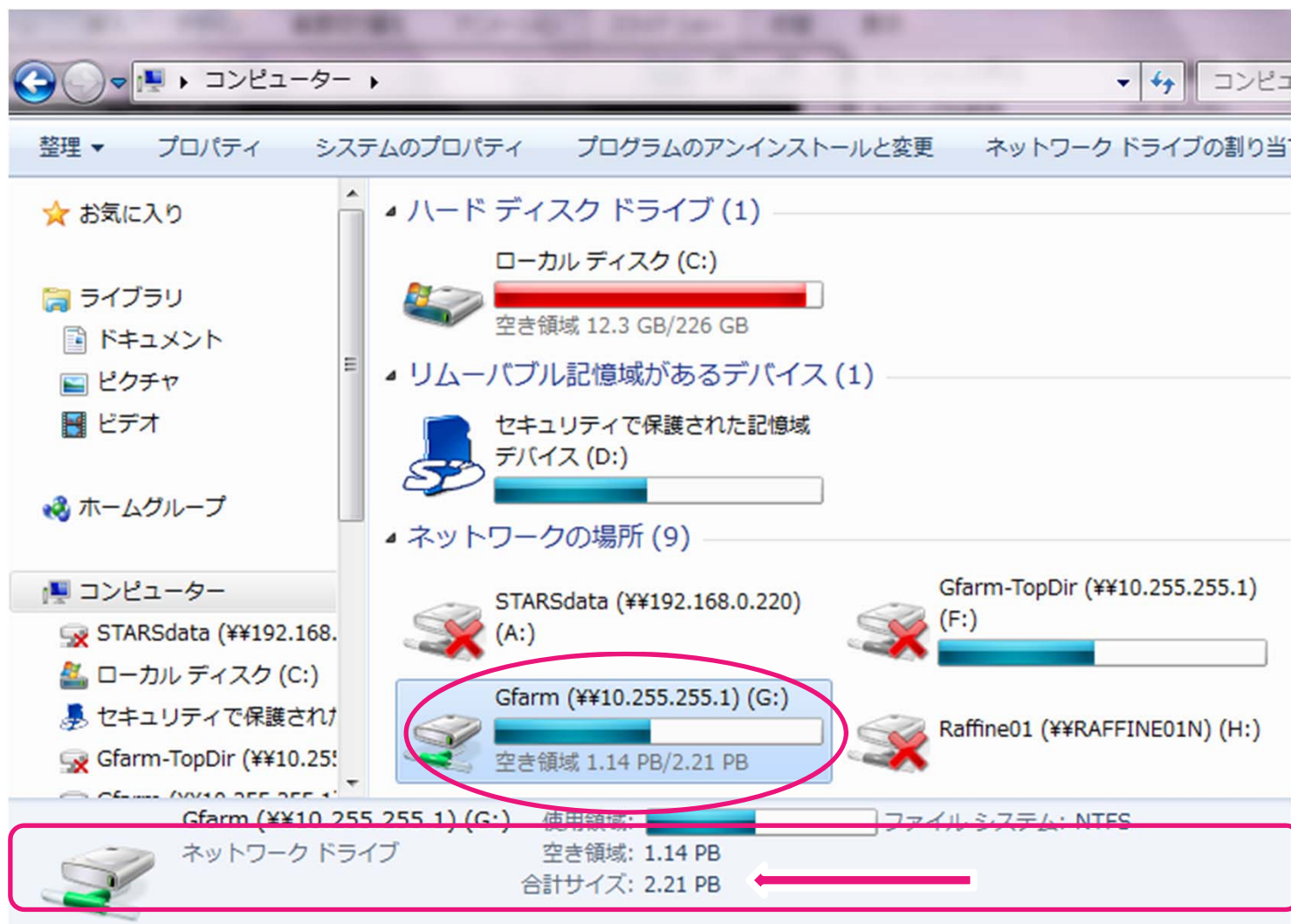
Laptopパソコンからクラウドストレージ(NICTサイエンスクラウド)を利用

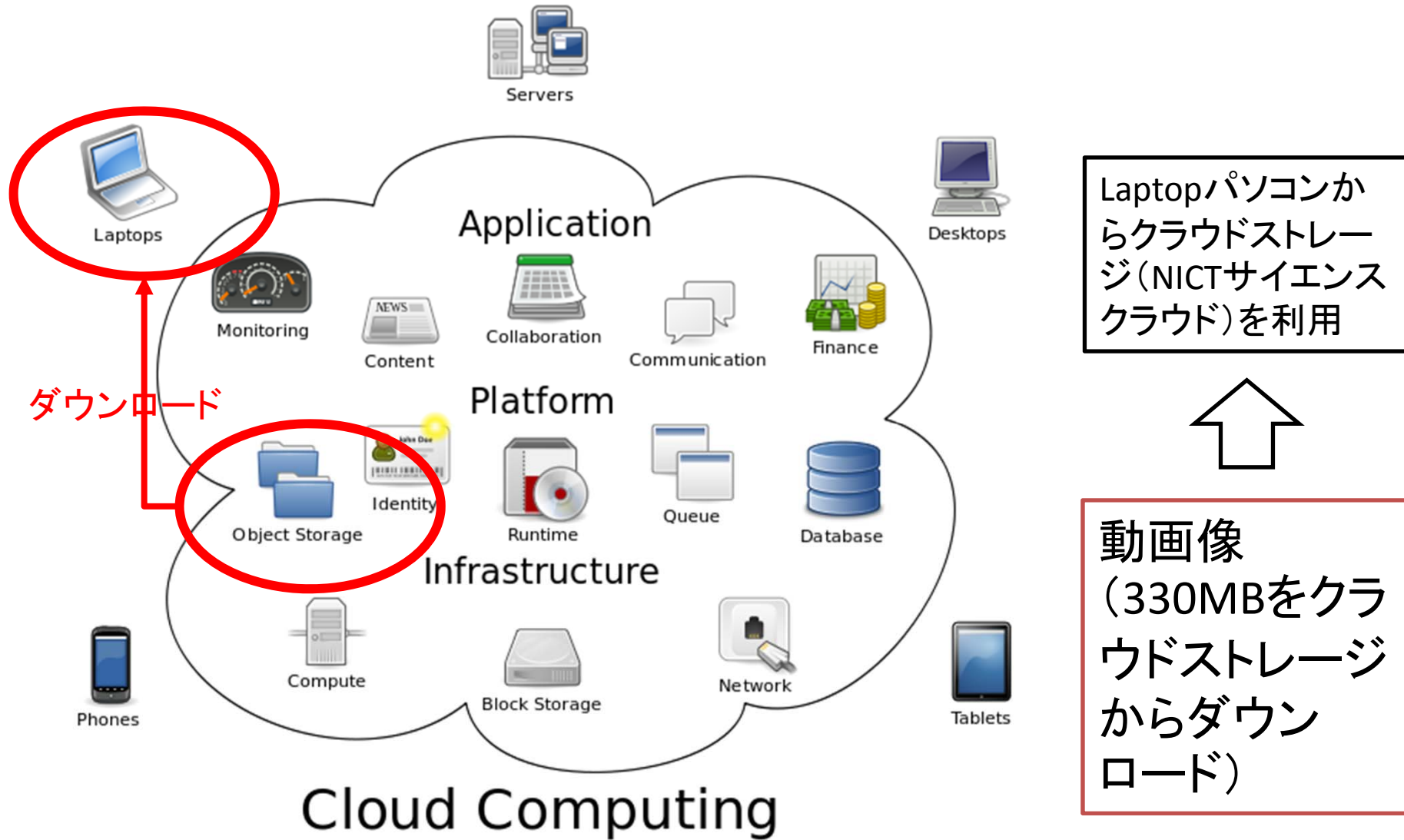


ストレージサイズが問題ではない！

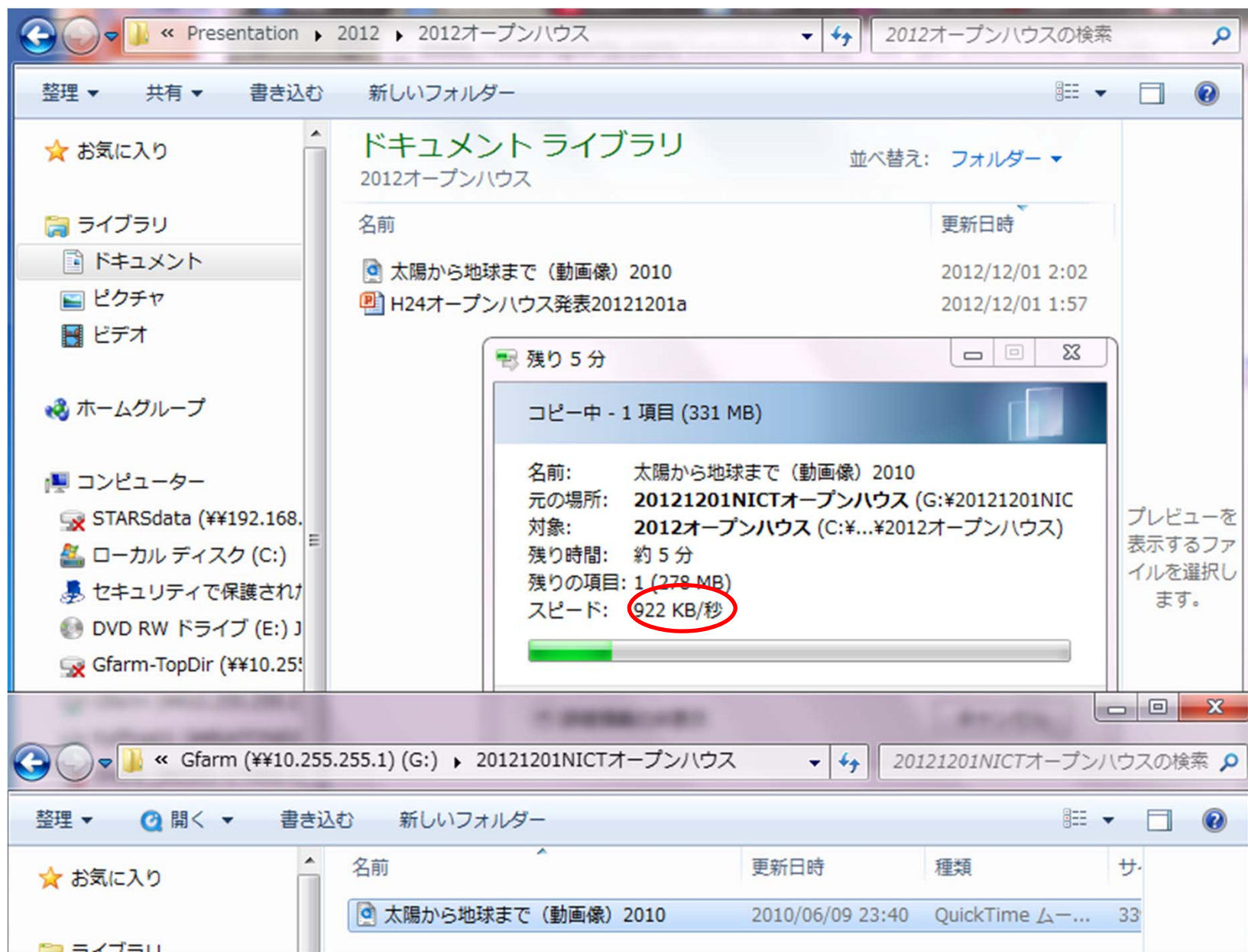
http://ja.wikipedia.org/wiki/%E3%83%95%E3%82%A1%E3%82%A4%E3%83%AB:Cloud_computing.svg

私のラップトップPC

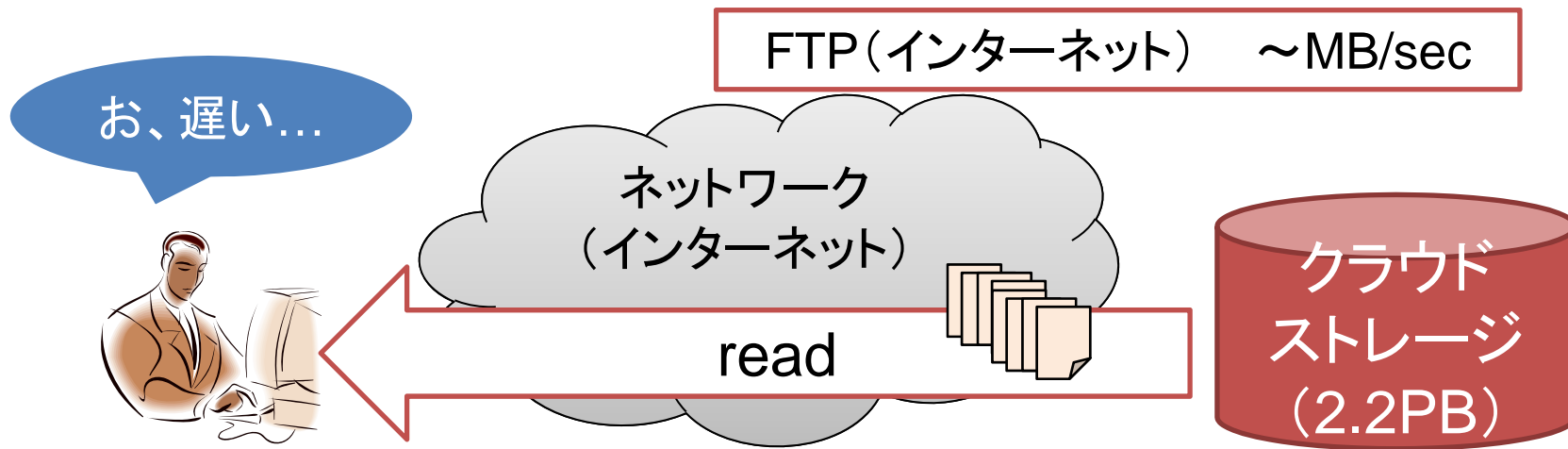




http://ja.wikipedia.org/wiki/%E3%83%95%E3%82%A1%E3%82%A4%E3%83%AB:Cloud_computing.svg



クラウドストレージを科学研究で使うと...



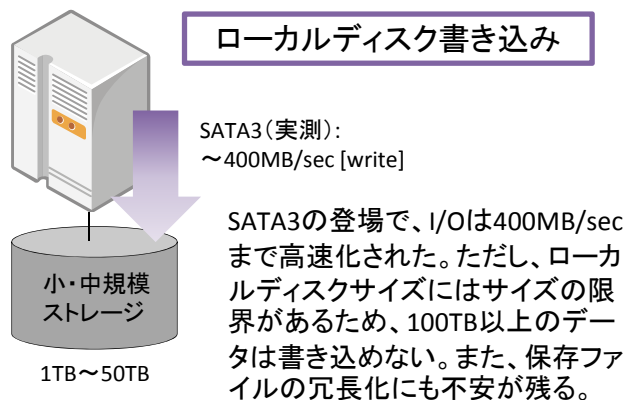
一般的なインターネットが1MB/secと仮定として...

↓
全てのデータファイルを
ダウンロード(2.2PB)

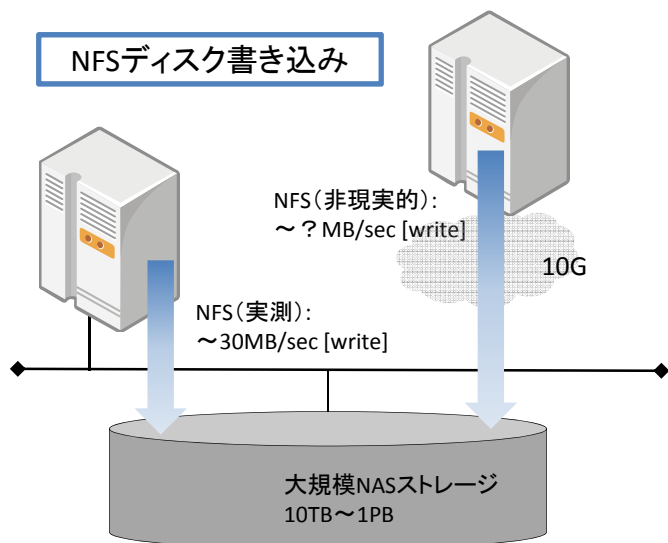
2,200,000,000秒 = 70年 !

データ書き込み方法

ローカルディスク書き込み



NFSディスク書き込み



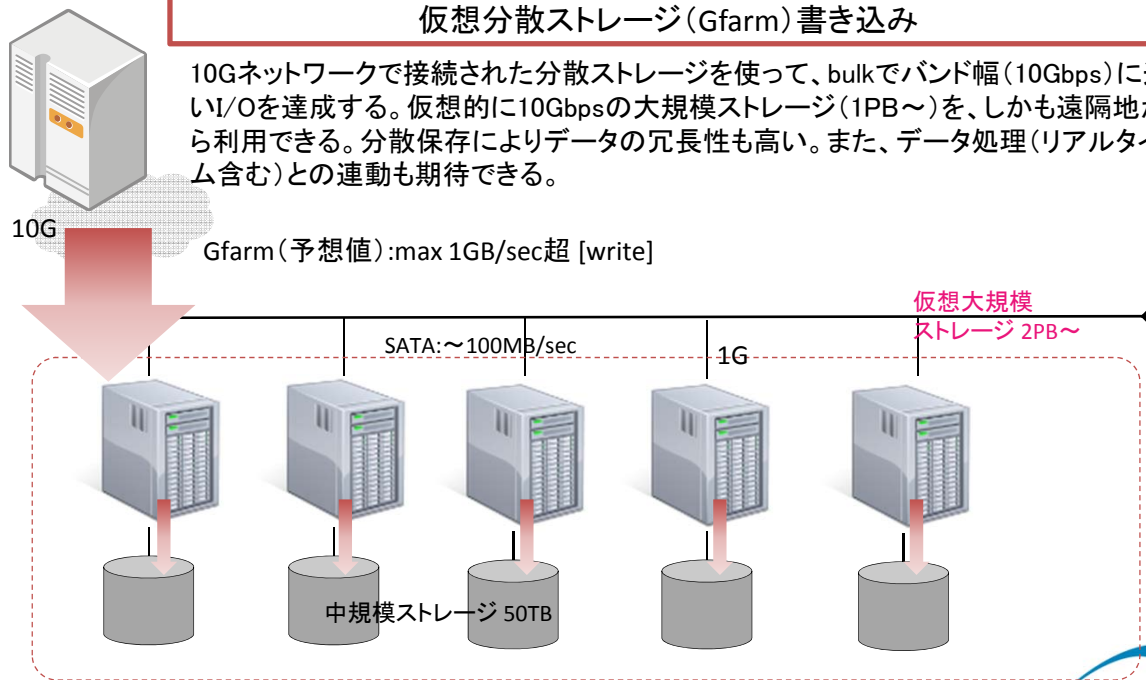
NFSの高速化は容易ではなく、高速ネットワークを使っても50MB/secを超えるのは難しい。また、長距離ネットワーク越えのディスクマウントは現実的ではない。

データ書き込み時間の比較(100TB)

- ローカルディスク書き込み時間
 - 100TBのローカルディスクは(原則)存在しない
 - 250000秒=約3日
- NFS(NAS)ディスク書き込み時間
 - 約3000000秒=約38日
- 仮想分散ディスク書き込み
 - バンド幅(10G)の活用を仮定
 - 10000秒=3時間弱

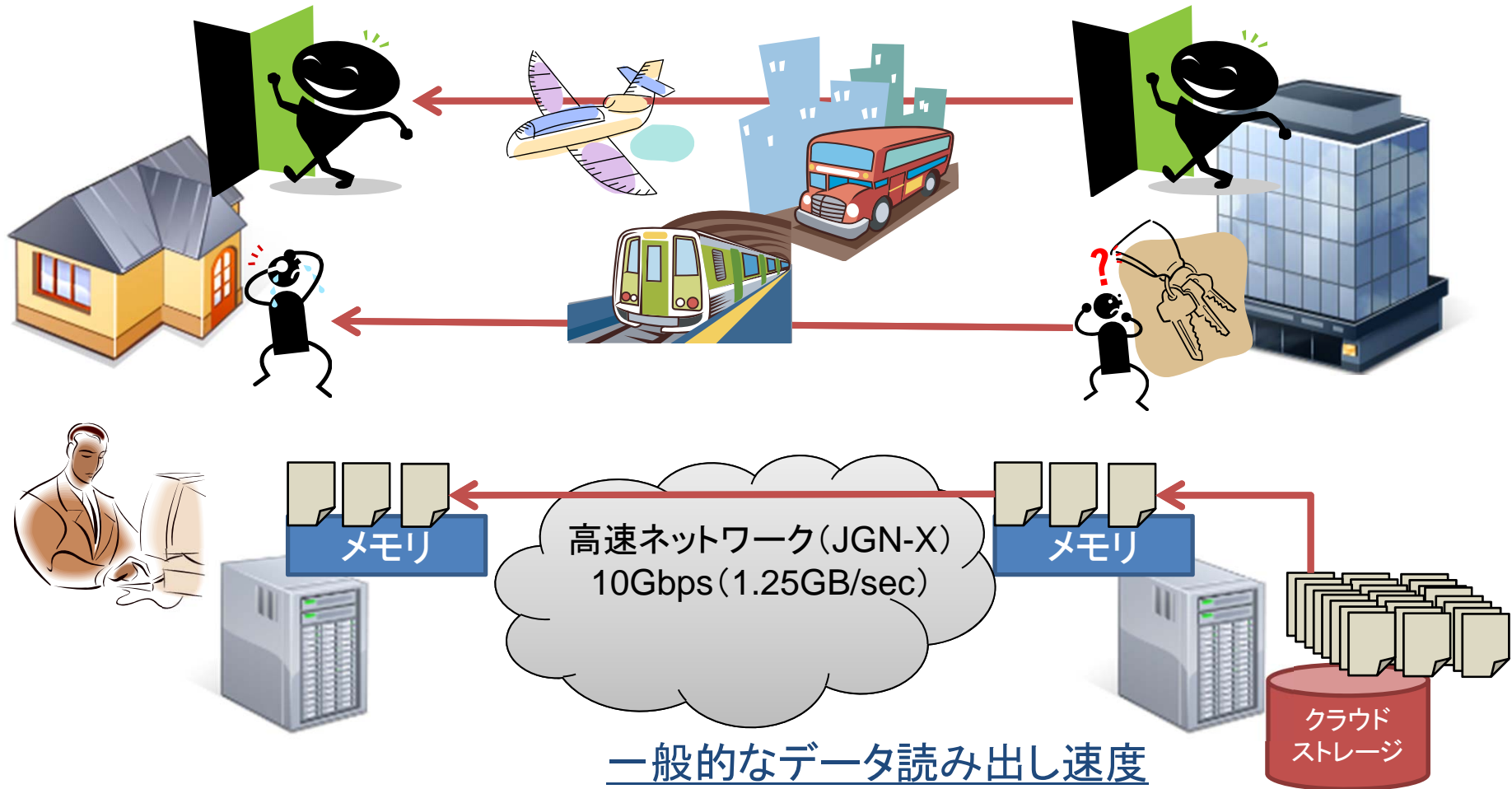
仮想分散ストレージ(Gfarm)書き込み

10Gネットワークで接続された分散ストレージを使って、bulkでバンド幅(10Gbps)に近いI/Oを達成する。仮想的に10Gbpsの大規模ストレージ(1PB~)を、しかも遠隔地から利用できる。分散保存によりデータの冗長性も高い。また、データ処理(リアルタイム含む)との連動も期待できる。



クラウドストレージの高速化

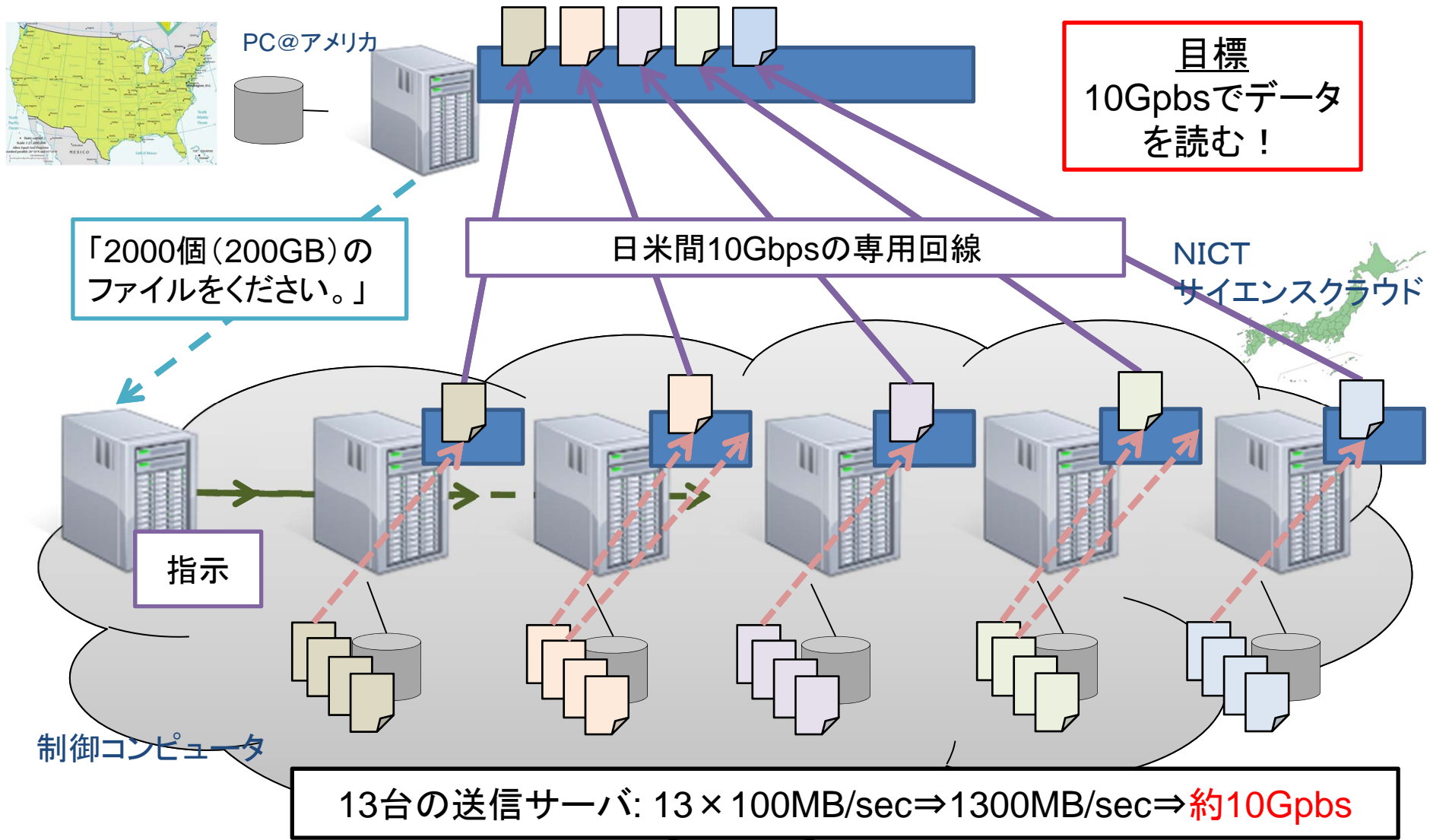
遅いのはネットワークか？



一般的なデータ読み出し速度

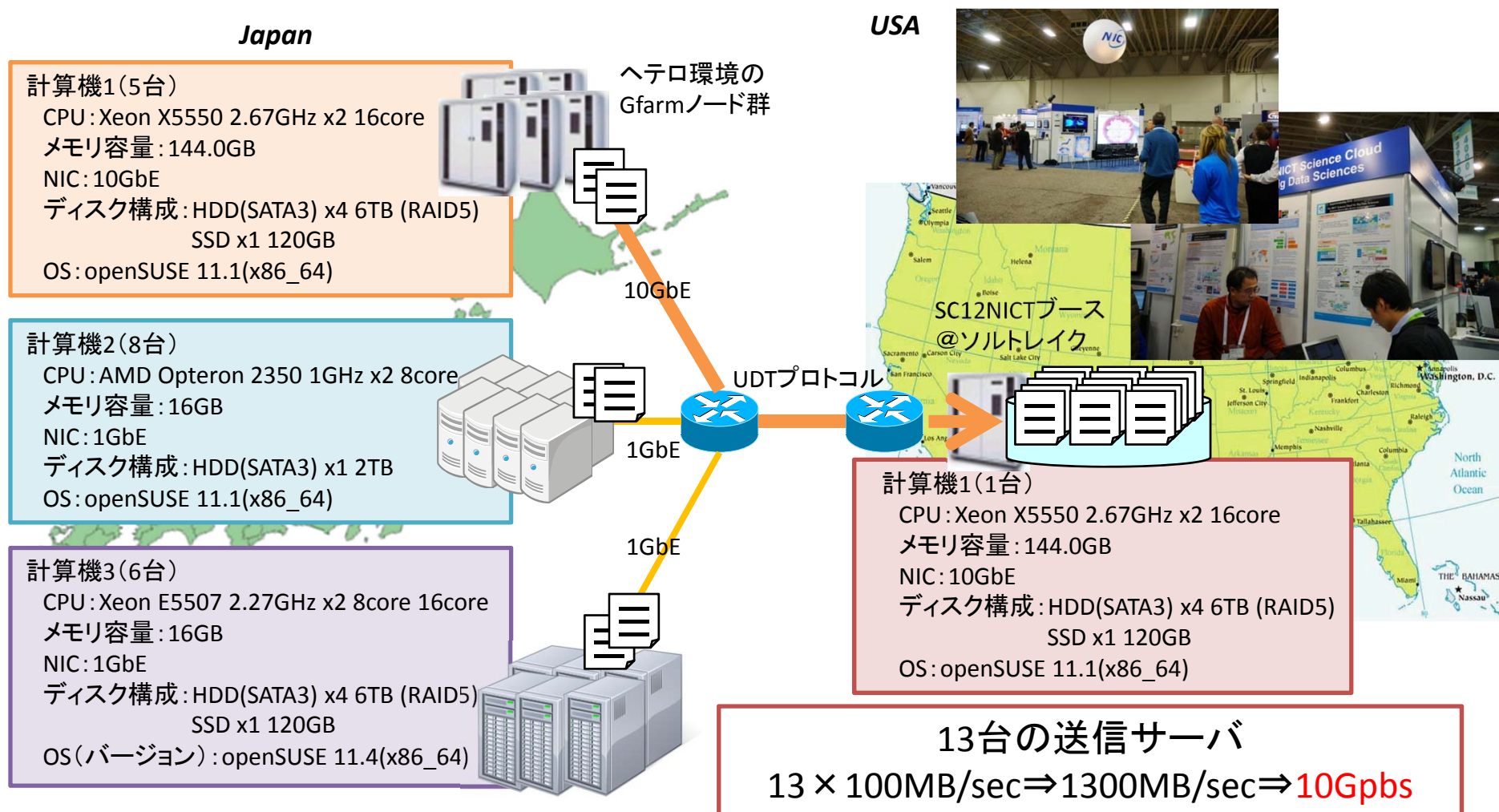
- NAS 400Mbps (50MB/sec)
- SSD disk 1.8Gbps (225MB/sec)
- (SATA3/RAID5 3.8Gbps)
- (RAM disk 3.6Gbps)

日米間高速ストレージ実験①(2012年11月)

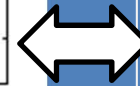
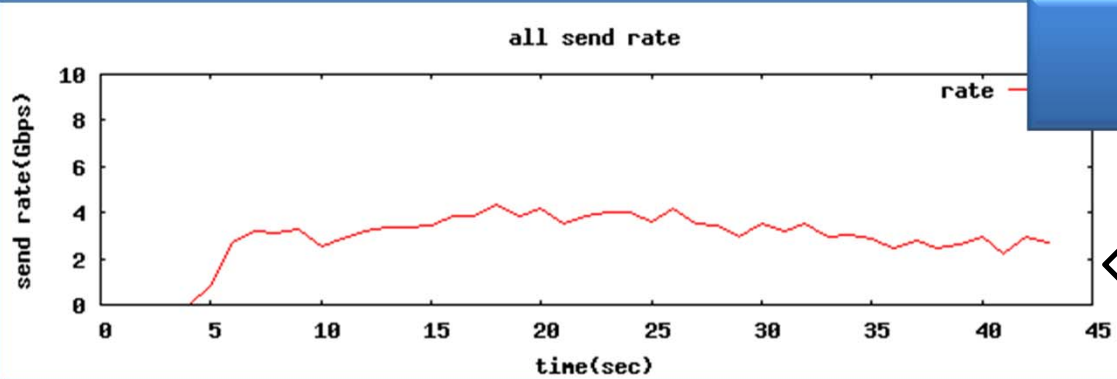


日米間高速ストレージ実験②(2012年11月)

- 遠隔地からクラウドストレージを自分のストレージのように仮想利用
- NICTサイエンスクラウド上の13台のファイルノード上に配置した125MBの複数ファイルを1台のクライアント計算機@米国から読み出し

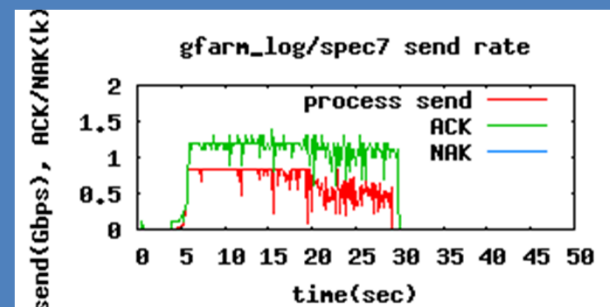
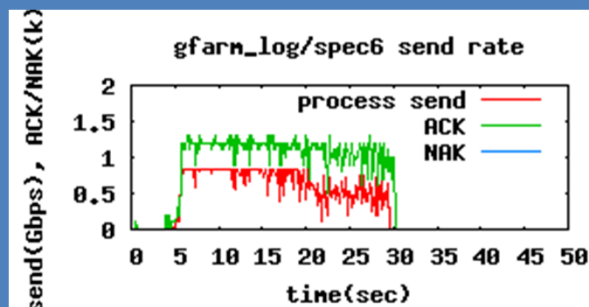
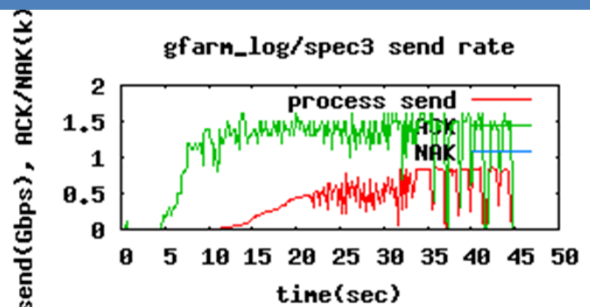
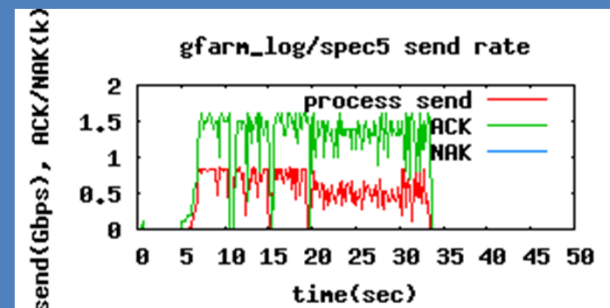
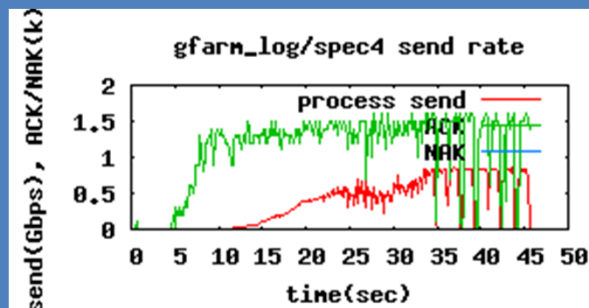
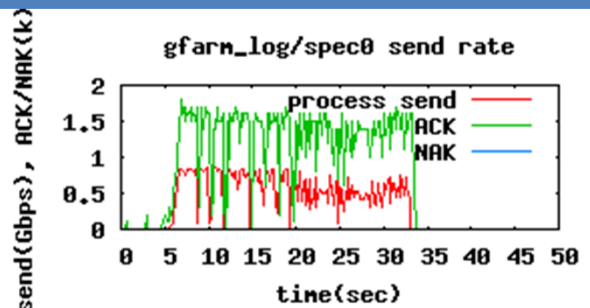
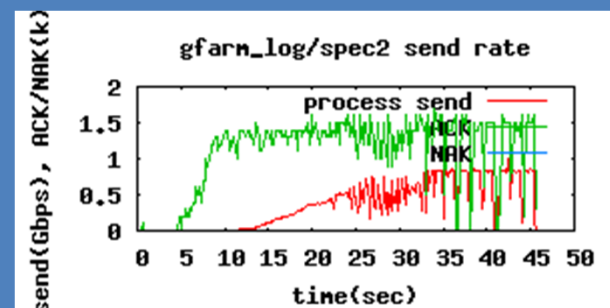
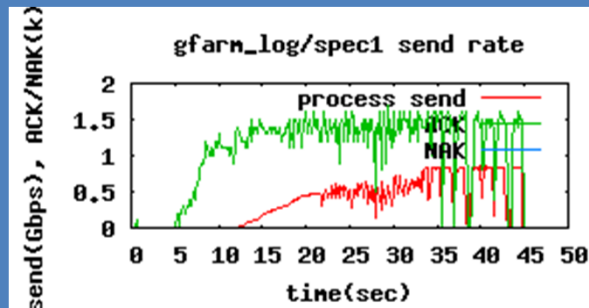


8セッションの場合



- Local Disk (Read)
- SSD disk -> 1.8Gbps
 - SATA3/RAID5 -> 3.8Gbps
 - RAM disk -> 3.6Gbps

- 125MB/150 files
- 8 sessions
- read

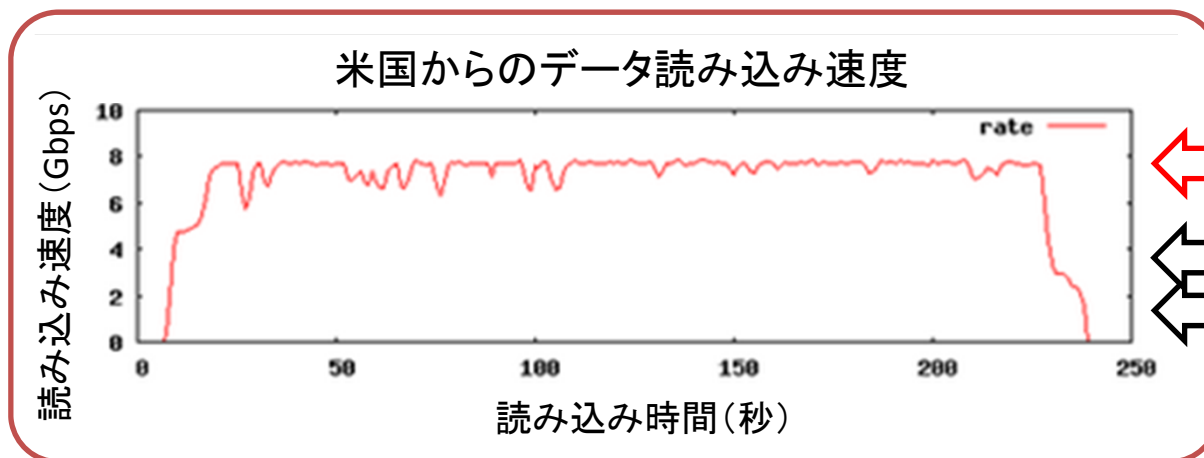


6.95Gbps (平均速度)

※参考: 単体ディスクI/O(Read)性能

SSD: 約1.76Gbps

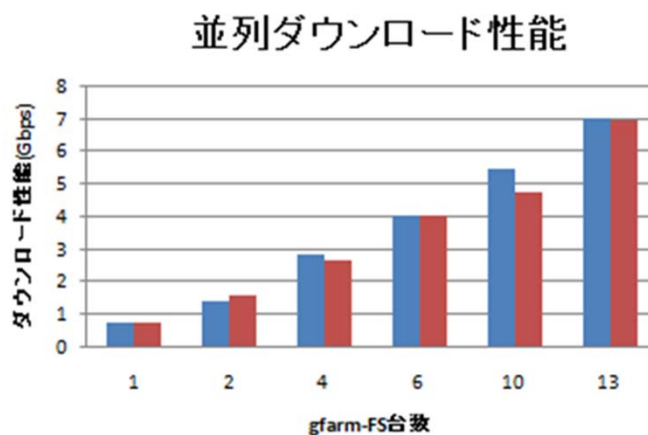
HDD(SATA3・RAID5): 約3.78Gbps



← 約7Gbps

← HDD(SATA3,RAID5): 約3.78Gbps
 ← SSD: 約1.76Gbps

208GB (DVD約44枚分)のデータファイルを日米間で240秒(4分)で伝送した。

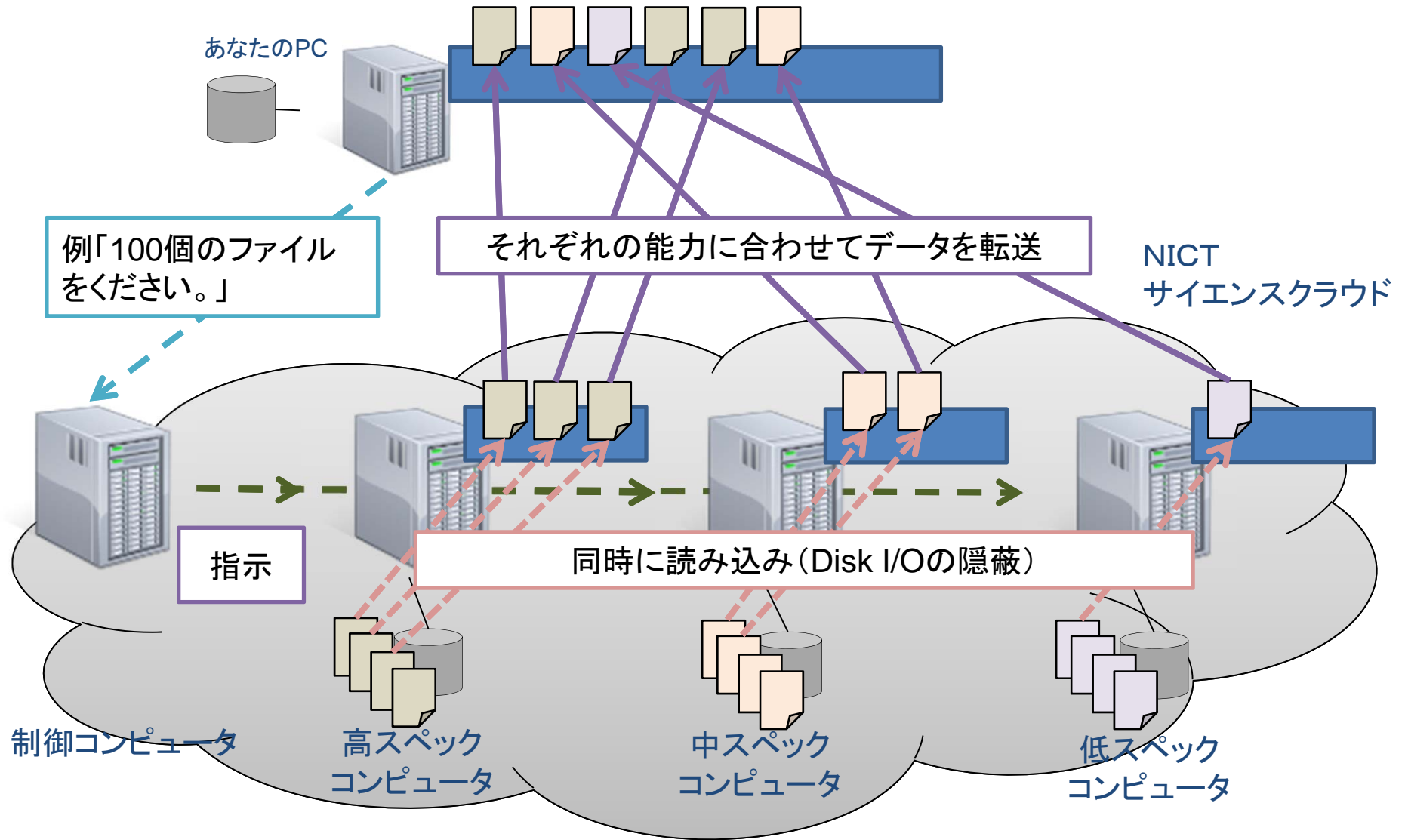


← HDD(SATA3,RAID5): 約3.78Gbps

← SSD: 約1.76Gbps

今後の課題:クラウドストレージの高速化

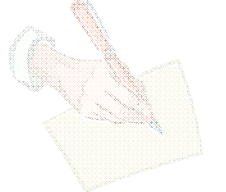
並列データ伝送技術



①理論→紙と鉛筆

データ指向型科学を実現するには...

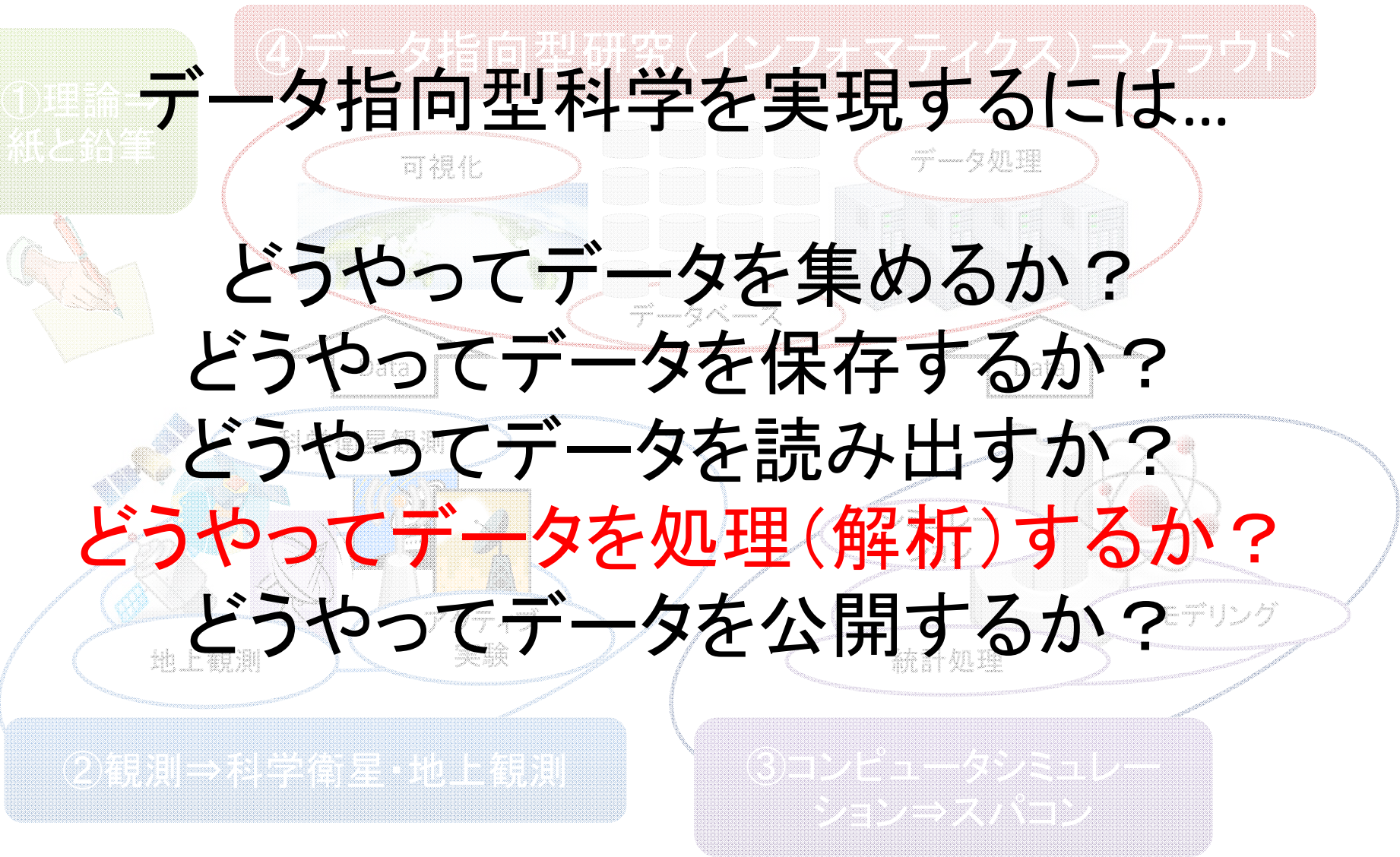
④データ指向型研究(インフォマティクス)⇒クラウド



- どうやってデータを集めるか？
- どうやってデータを保存するか？
- どうやってデータを読み出すか？
- どうやってデータを処理(解析)するか？
- どうやってデータを公開するか？

②観測⇒科学衛星・地上観測

③コンピュータシミュレーション⇒スパコン

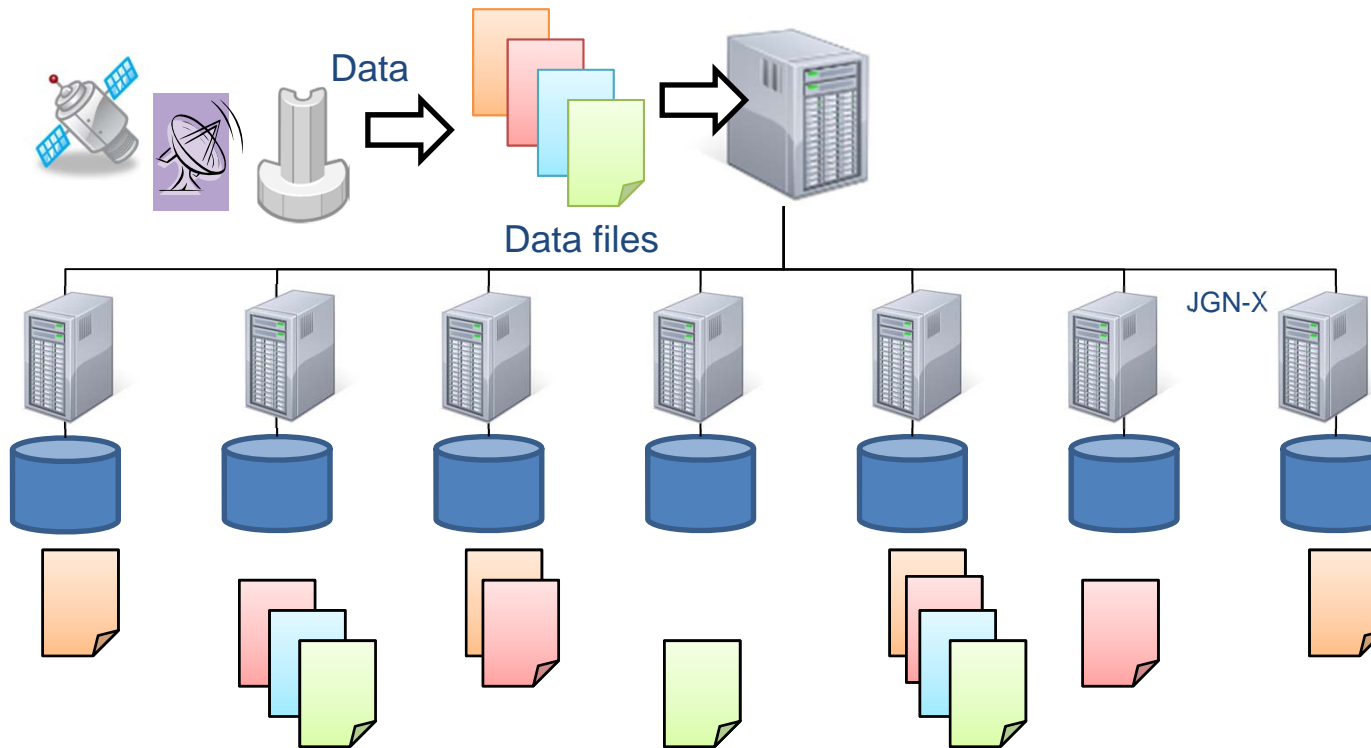


A Virtual Remote Disk for Remote Visualization of Computer Simulations

*Ken. T. Murata, K. Yamamoto, H. Watanabe,
T. Kurosawa, E. Kimura, S. Watari, M. Ishii,
KDDI team, JGN-X team, and H. Tatebe*

*National Institute of Information and Communications Technology
Applied Electromagnetic Research Institute
Space Weather and Environment Informatics Laboratory
4-2-1, Nukui-Kitamachi, Koganei, Tokyo 184-8795, Japan
TEL: +81-42-327-7931 FAX: +81-42-327-6978
E-mail: SciCloud-office@ml.nict.go.jp*

どうやってデータを処理(解析)するか? ①



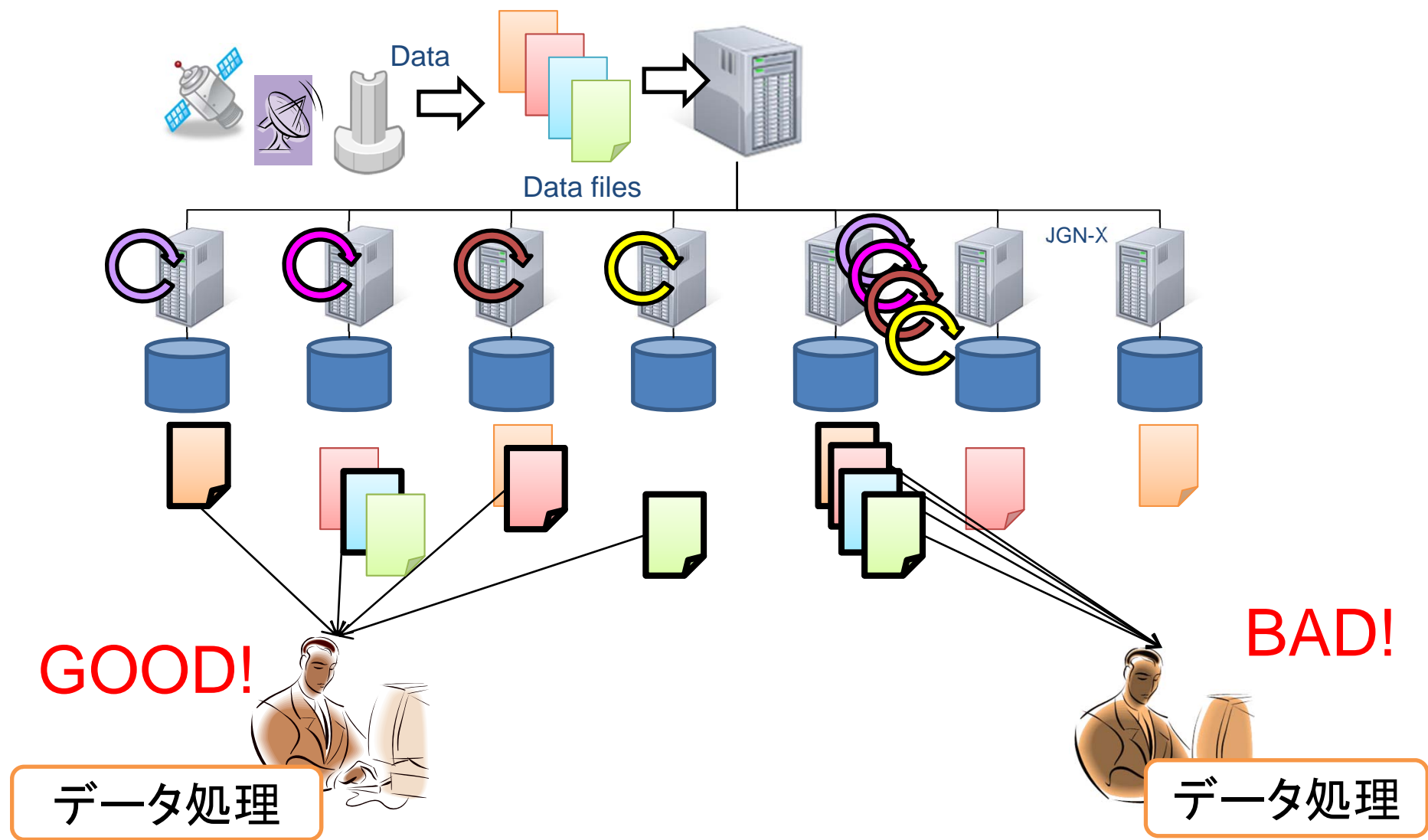
4つのファイルはそれぞれ3か所に保存される。

バックアップレス

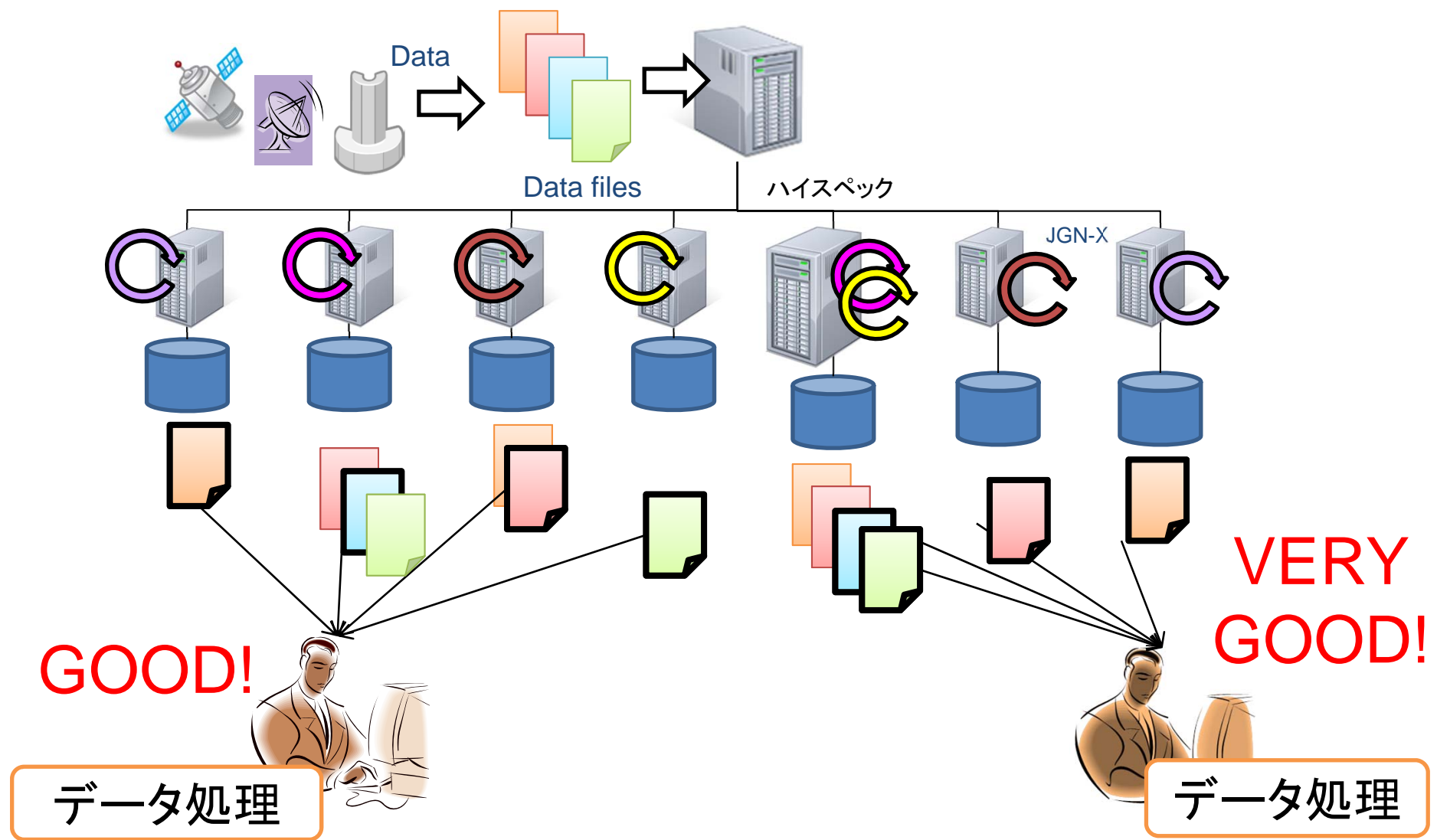
サーバは、①ファイル管理サーバ(ファイルシステムノード)機能と②ファイル処理サーバ(クライアントノード)を兼ねる。

ディスクI/Oの隠蔽

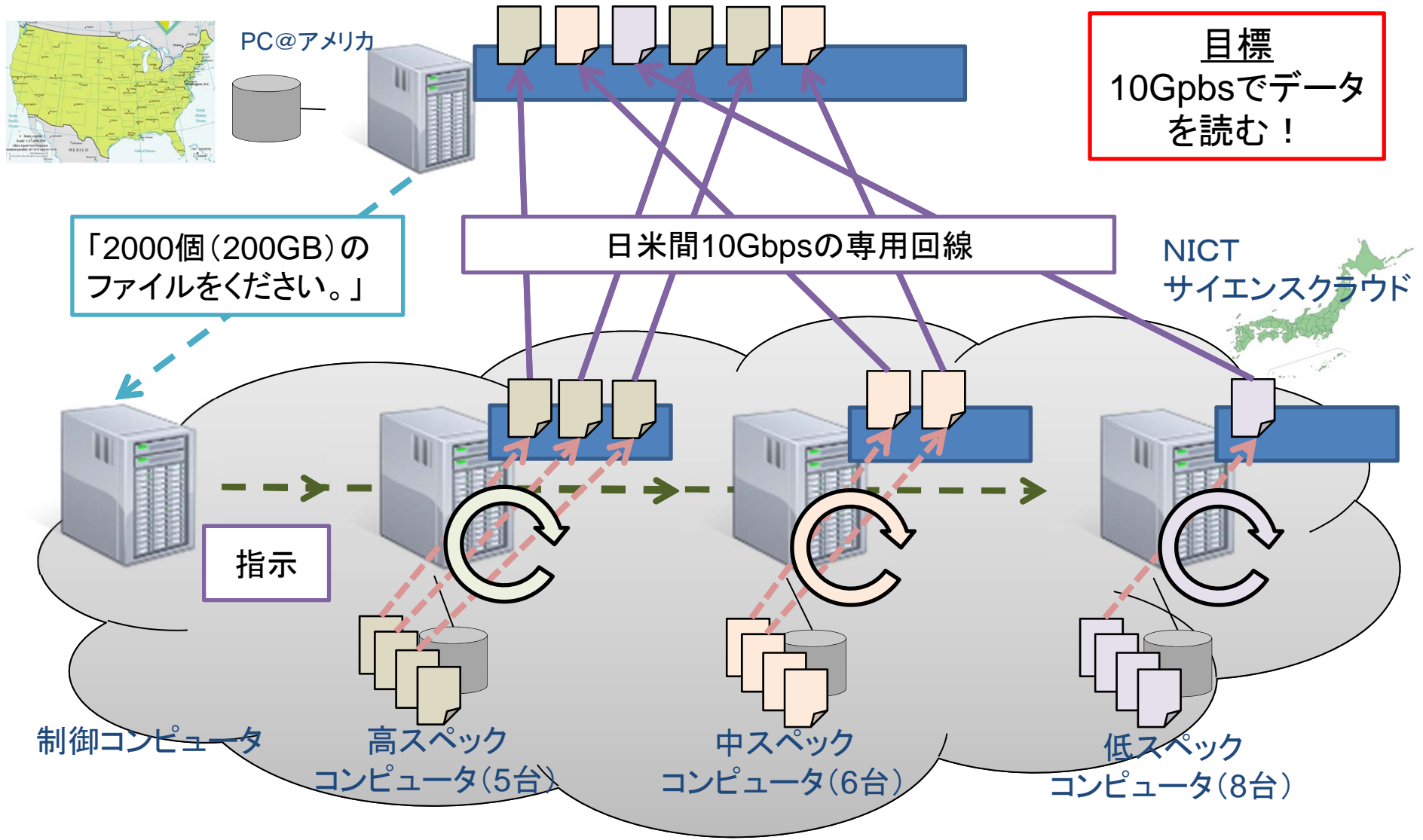
どうやってデータを処理(解析)するか? ②



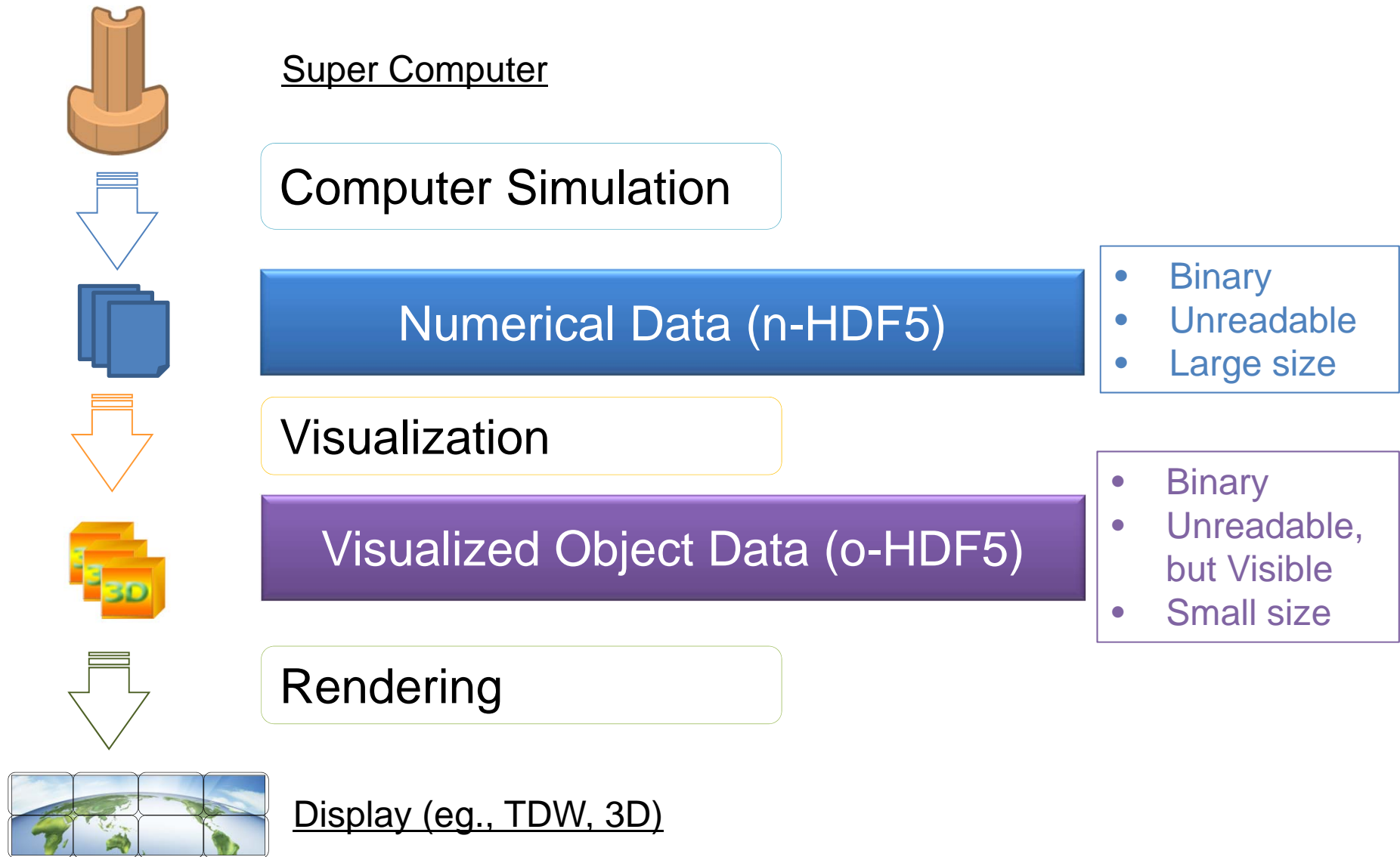
どうやってデータを処理(解析)するか？③



日米間高速データ処理実験①(2012年11月)

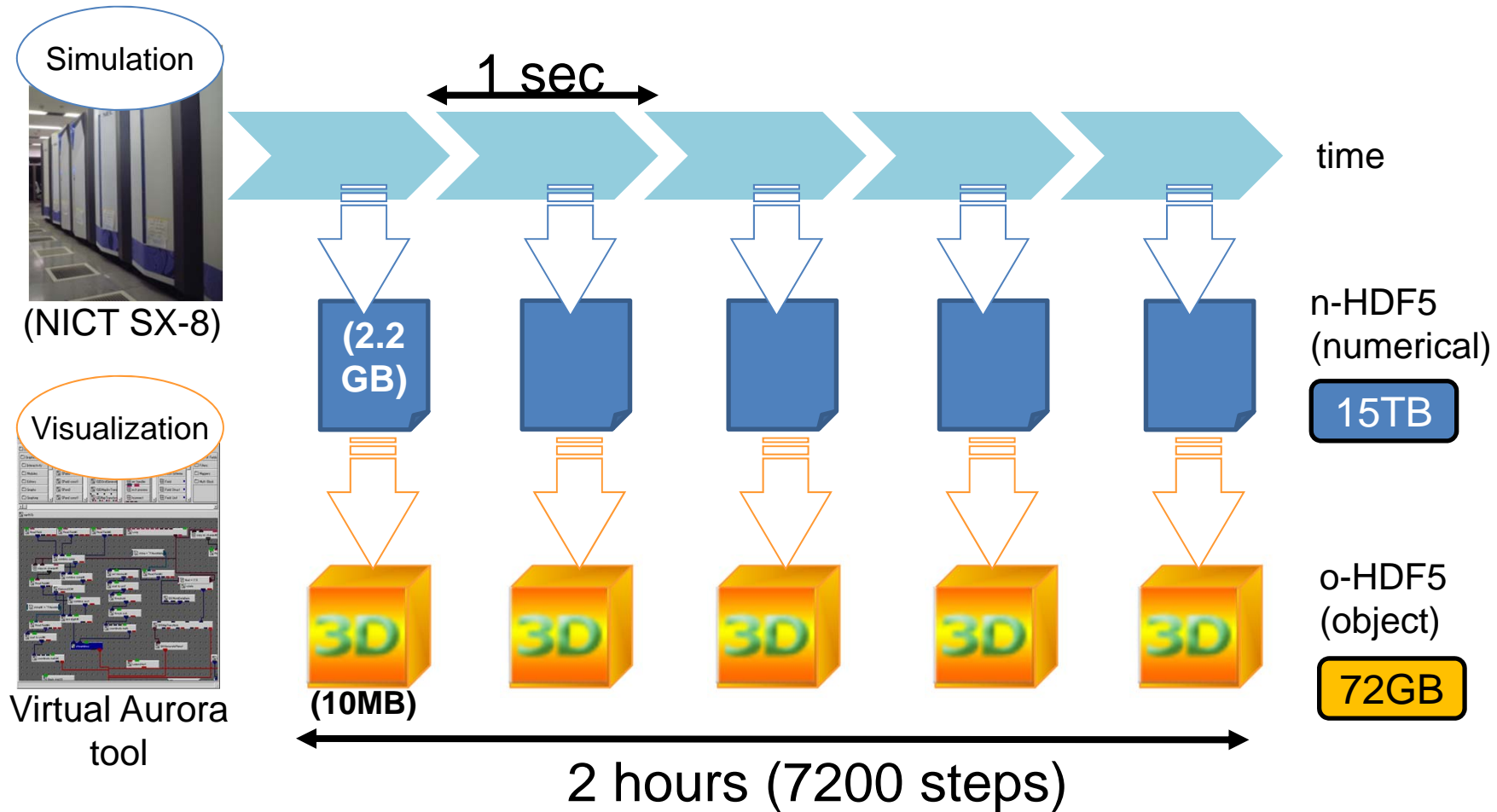


Typical Workflow of Visualization (post processing)



A simulation data size (example)



• Spatial and time resolution
450(x) × 300(y) × 300(z) –uniform grid (dx=0.2Re, dt=0.5sec)



Typical Data Visualization System for Super Computer

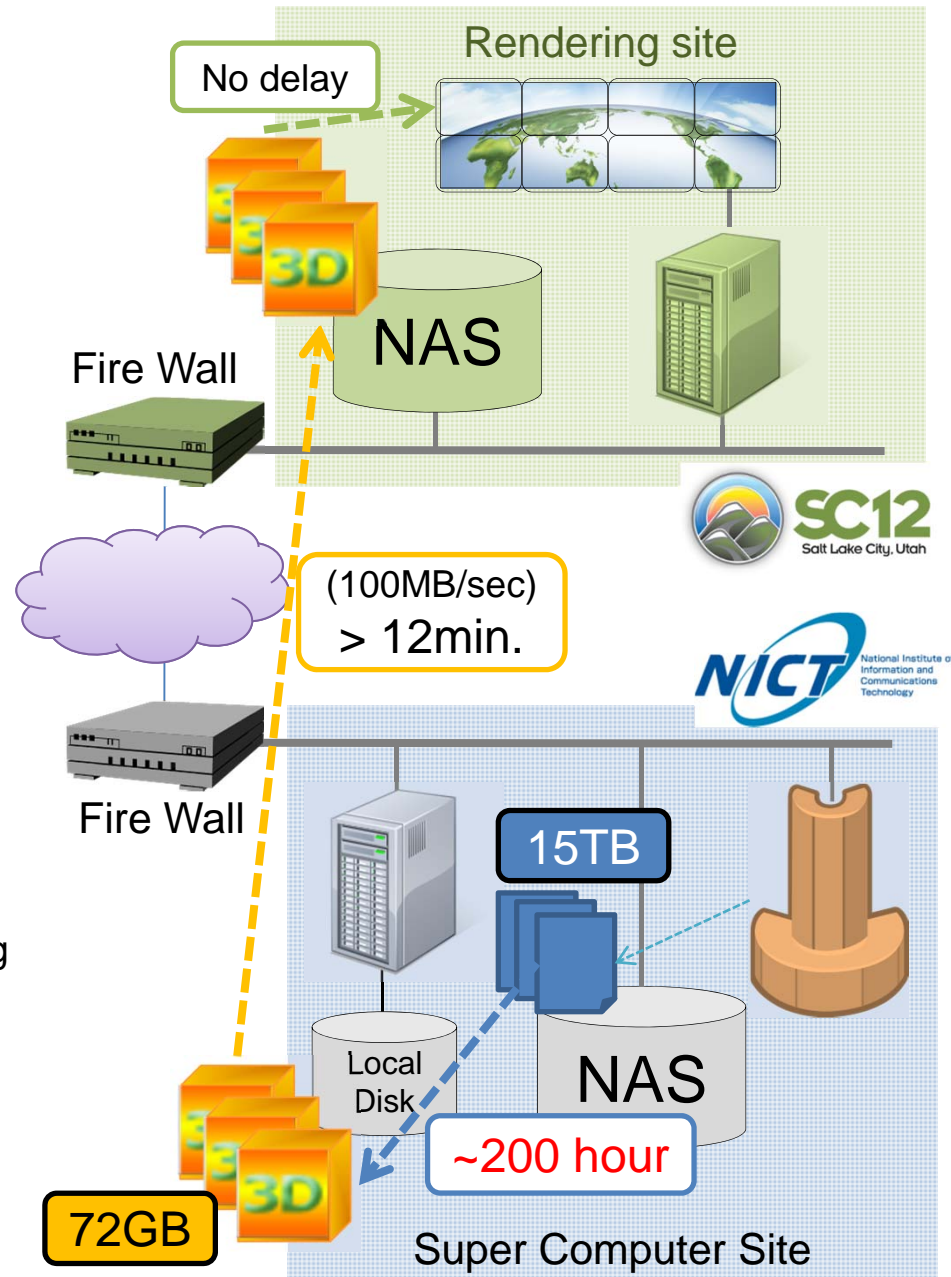
- **Rendering Site**

- The server has two functions: receiving the transferred data (in object HDF5 format) and rendering on the display.

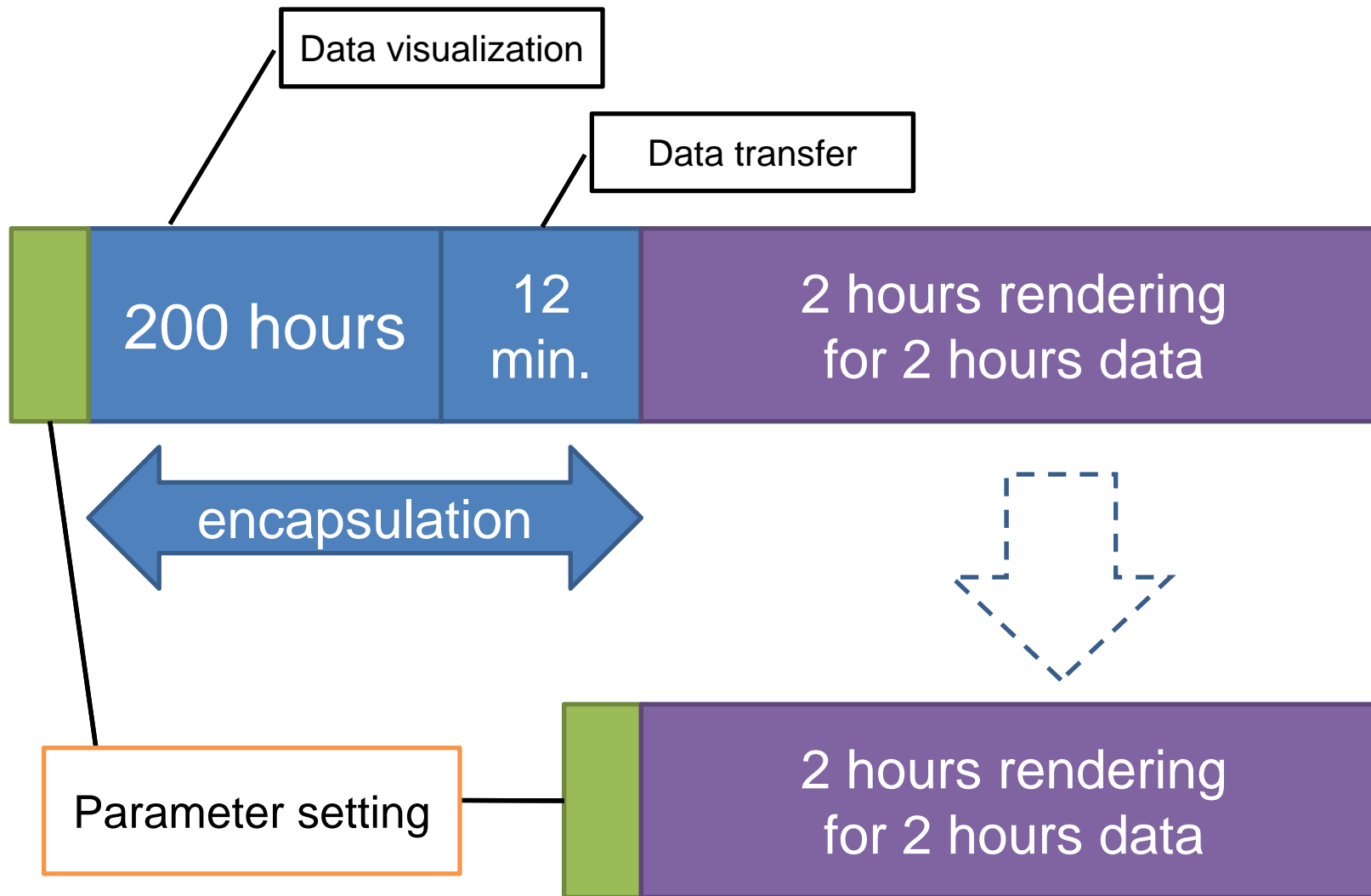
	Numerical data (n-HDF5) 2.2GB/file, file # 7200 (total 15TB)
	Visualized data (o-HDF5) 10MB/file, file # 7200 (total 72GB)

- **Super Computer Site**

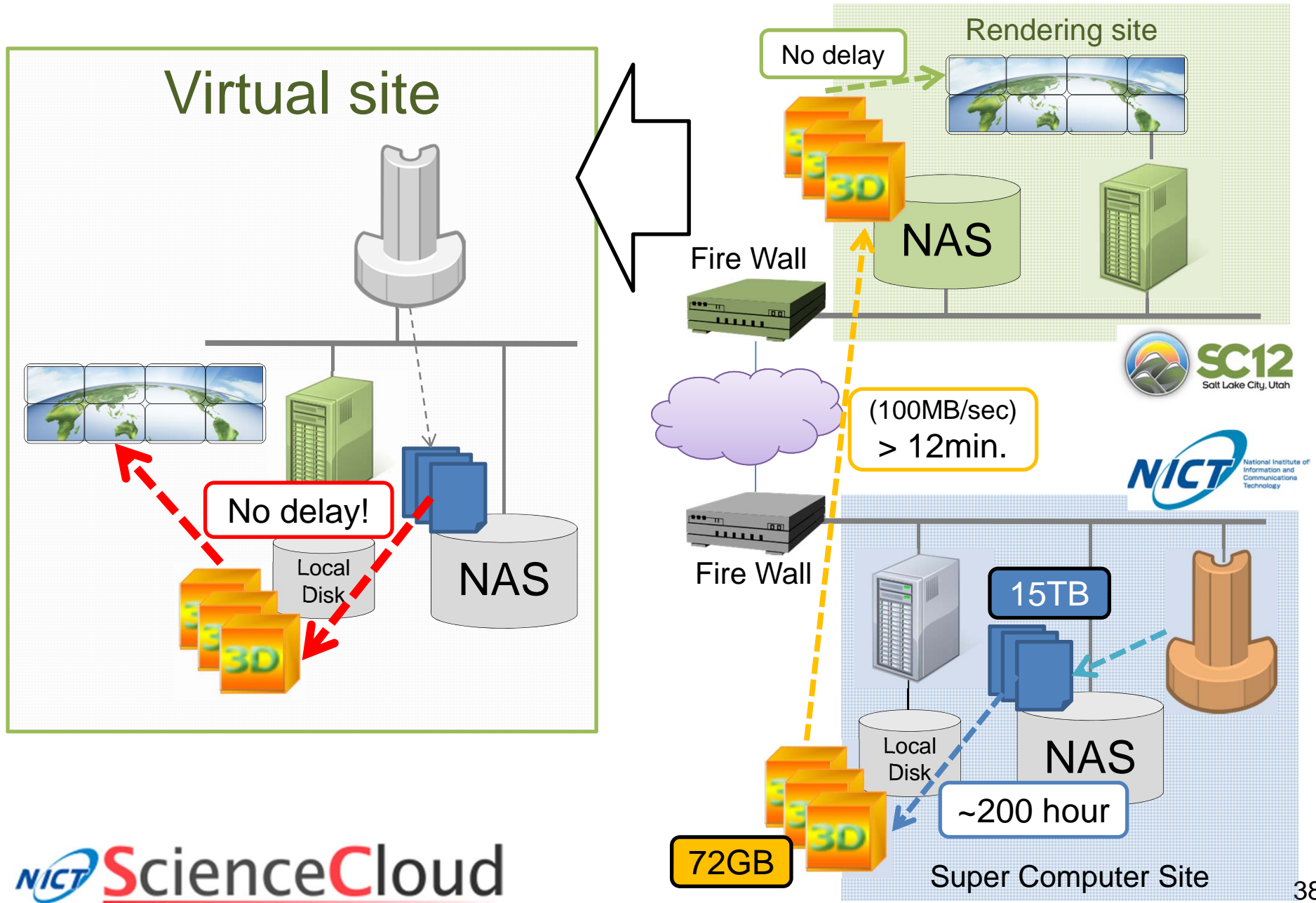
- Numerical simulation data files are storage on storage(s) in SC site. The storage is either local disk or NAS.
- The server has two functions: processing data files on the storage and transfer the data to the rendering site.



Examination: data size and transfer time In case of a 2 hour global MHD simulation

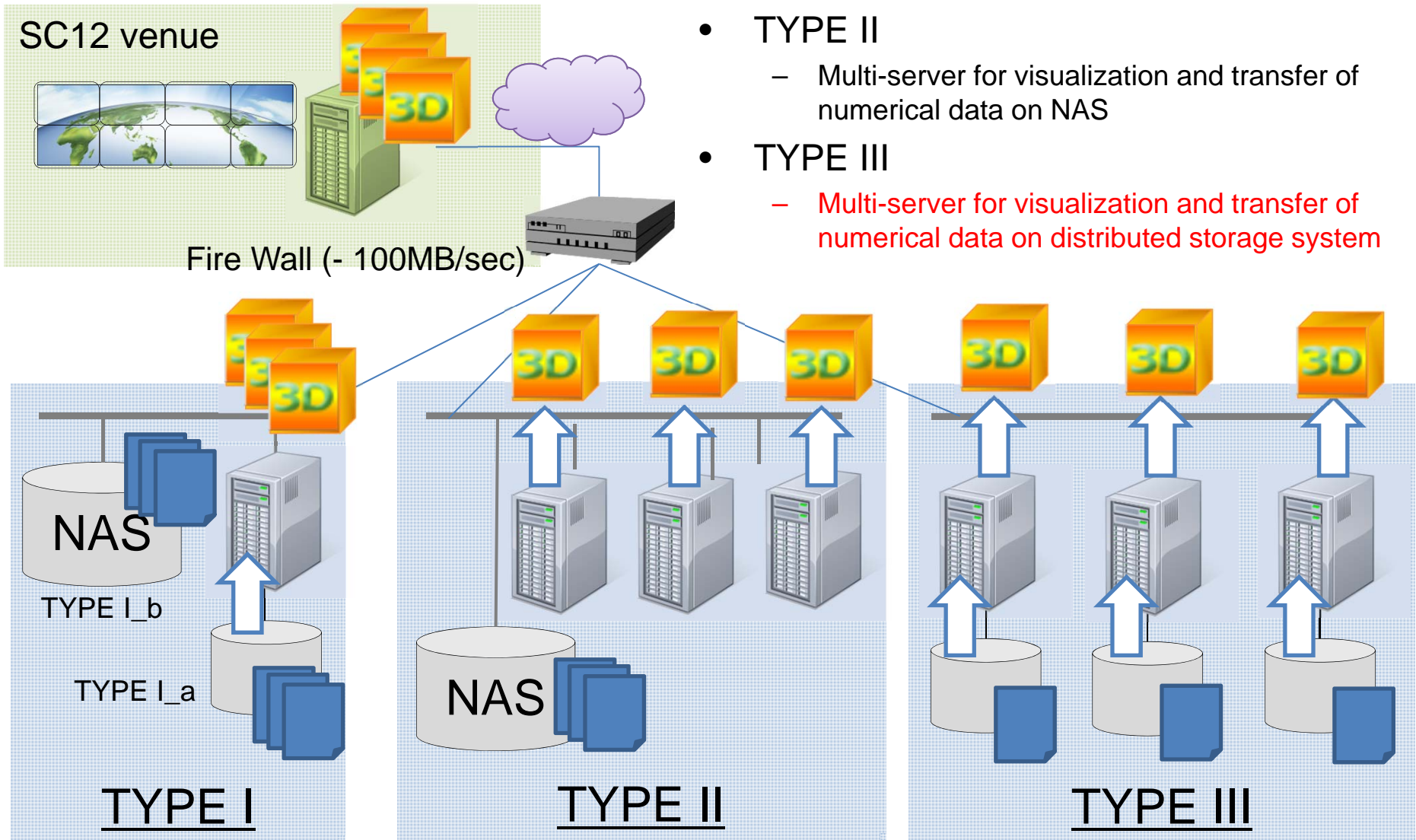


Typical Data Visualization System for Super Computer

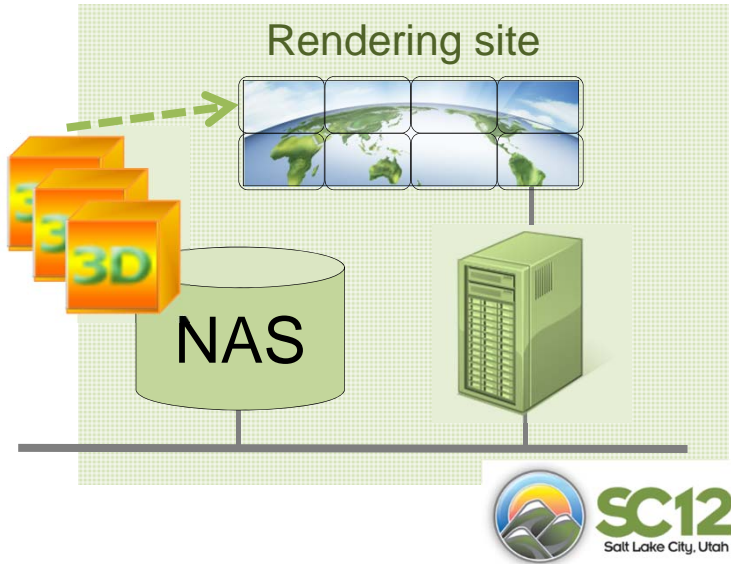


Systems for real-time visualization

- TYPE I
 - Single-server for visualization and transfer of numerical data on NAS/local disk
- TYPE II
 - Multi-server for visualization and transfer of numerical data on NAS
- TYPE III
 - Multi-server for visualization and transfer of numerical data on distributed storage system

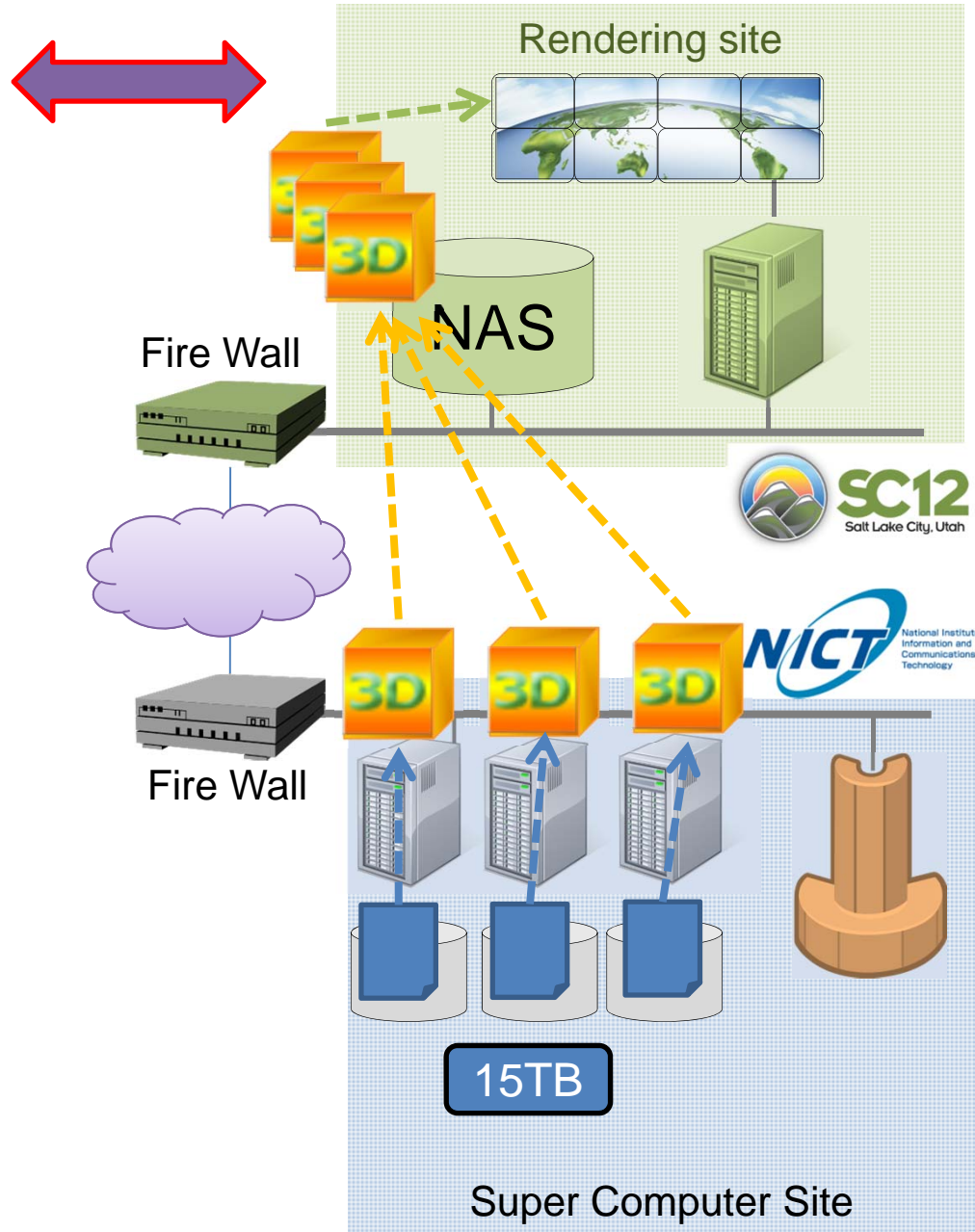


Competition!

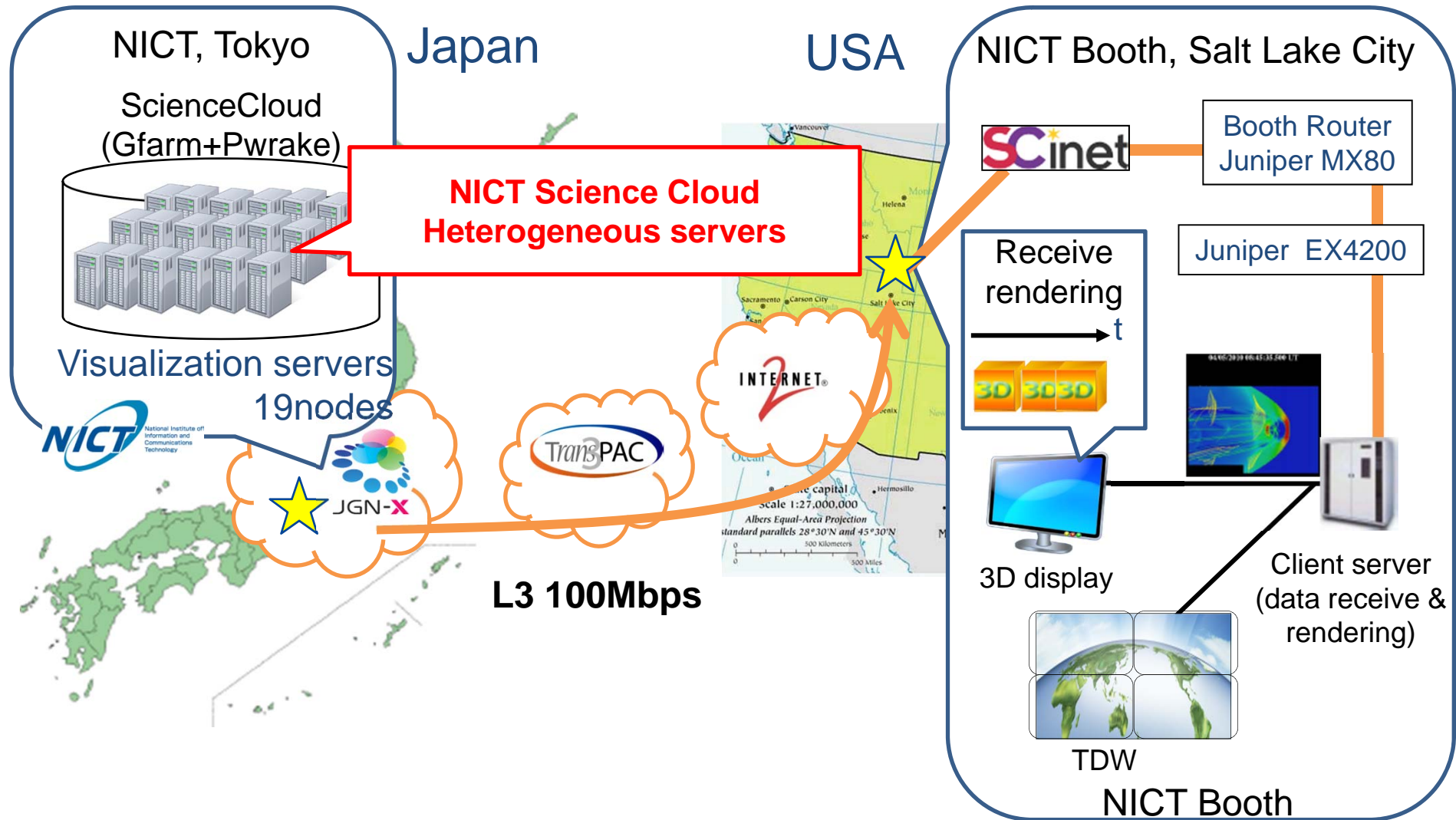


200hours+12min.

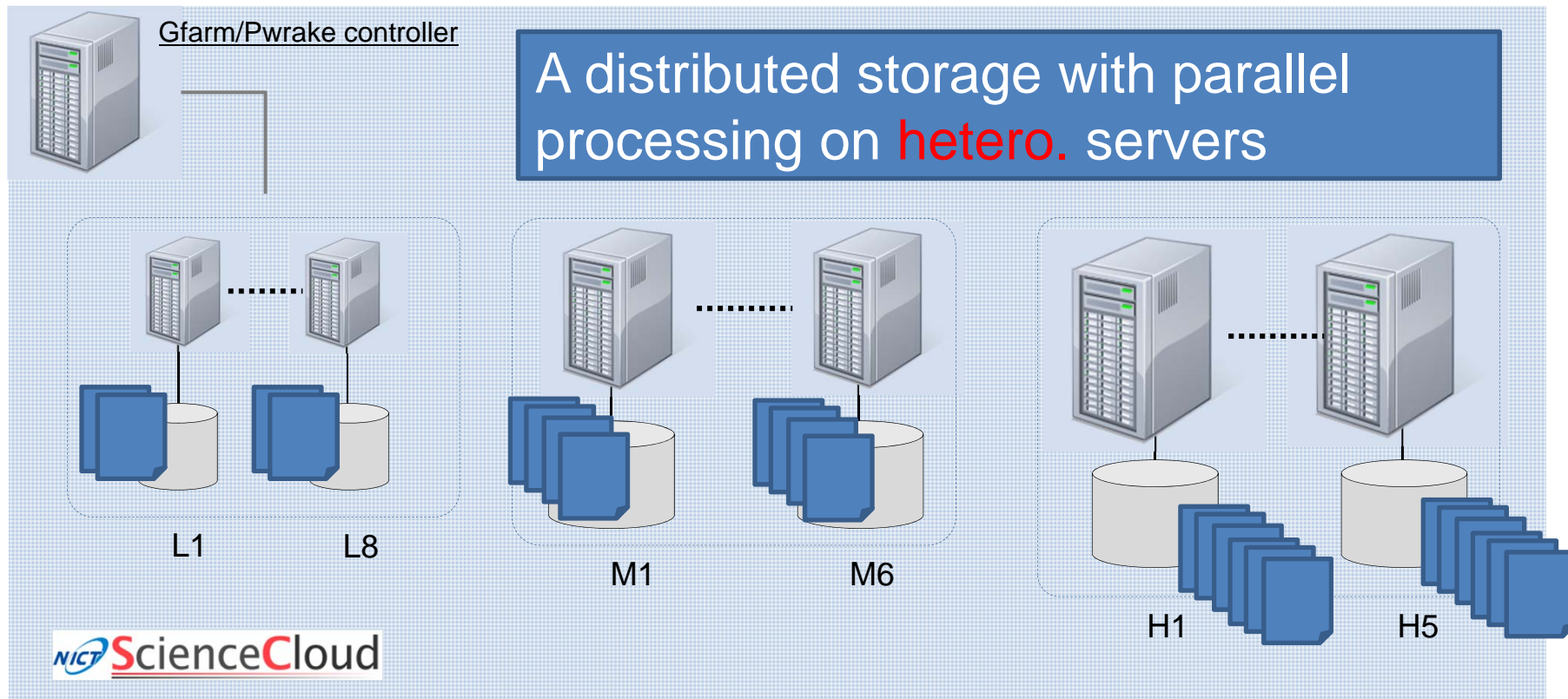
Advantage!



Real-time visualization over international network



Heterogeneous environment at processing site



Low Spec machine (#8)

- CPU #: 8 cores
- CPU: AMD Opteron Processor 2350
- Main memory: 16 GB
- OS: openSUSE 11.1 (x86_64)
- HDD: SATA 2 x1
Read: 127 MB/sec Write: 96MB/sec
- NIC: 1GbE

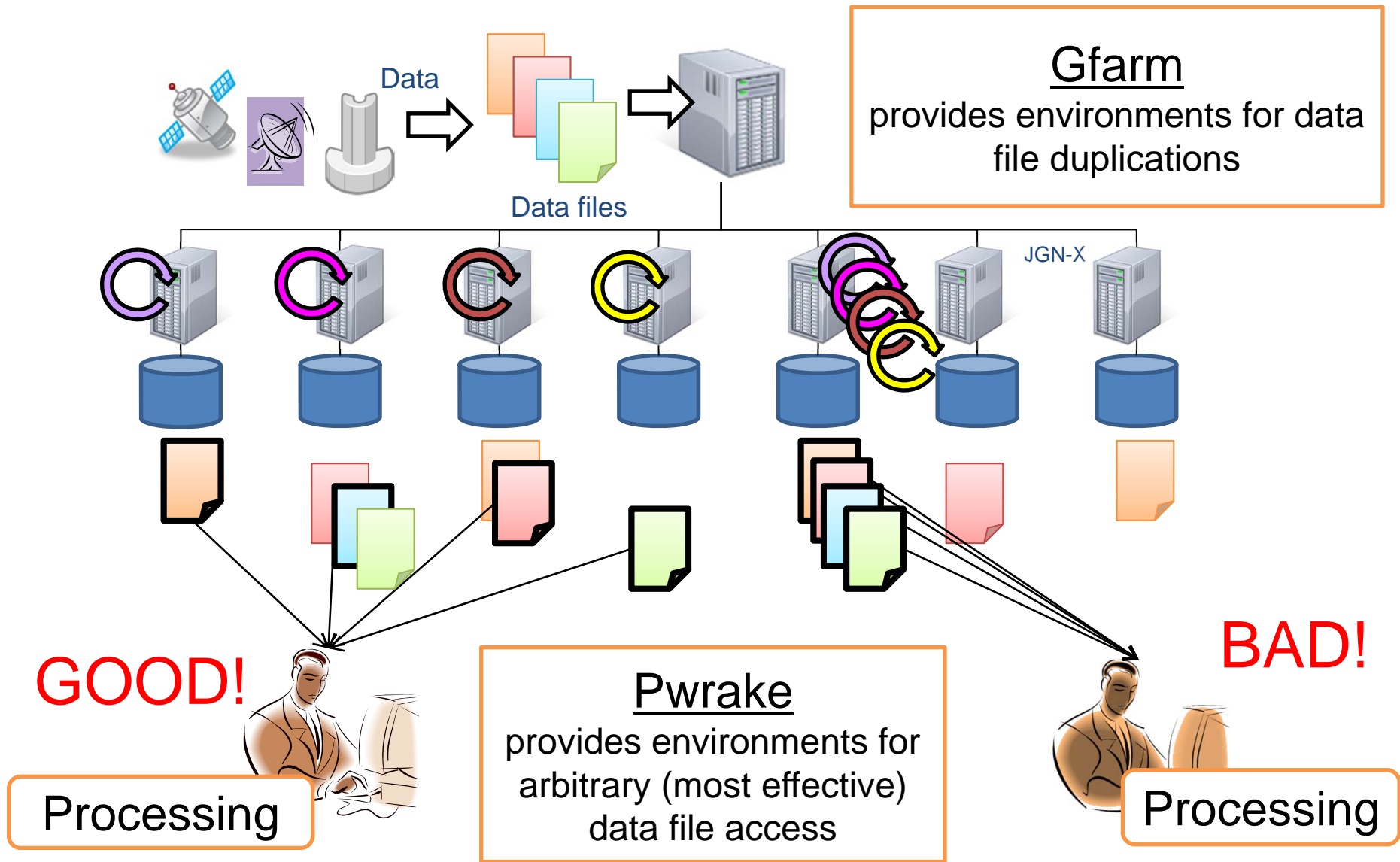
Middle Spec machine(#6)

- CPU #: 8 cores
- CPU: Intel Xeon E5507@2.27GHz
- Main memory: 16 GB
- OS: openSUSE 11.4 (x86_64)
- HDD: SATA 3 x4 (RAID5)
Read: 371 MB/sec Write: 137MB/sec
- NIC: 1GbE

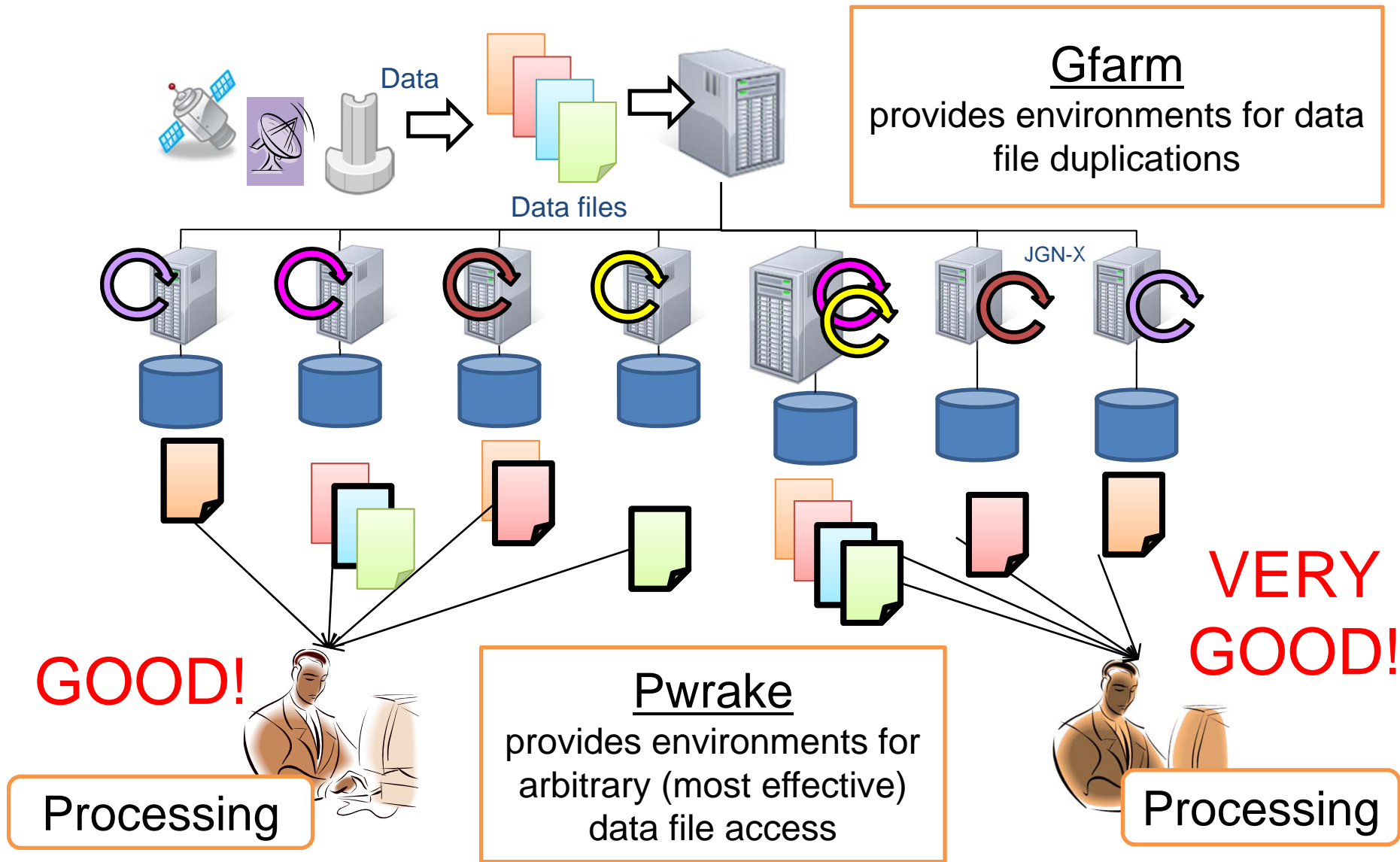
High Spec machine(#5)

- CPU #: 8 cores
- CPU: Intel Xeon X5550@2.67GHz
- Main memory: 144GB
- OS: openSUSE 11.1 (x86_64)
- HDD: SATA 3 x4 (RAID5)
Read: 371 MB/sec Write: 137MB/sec
- NIC: 10GbE

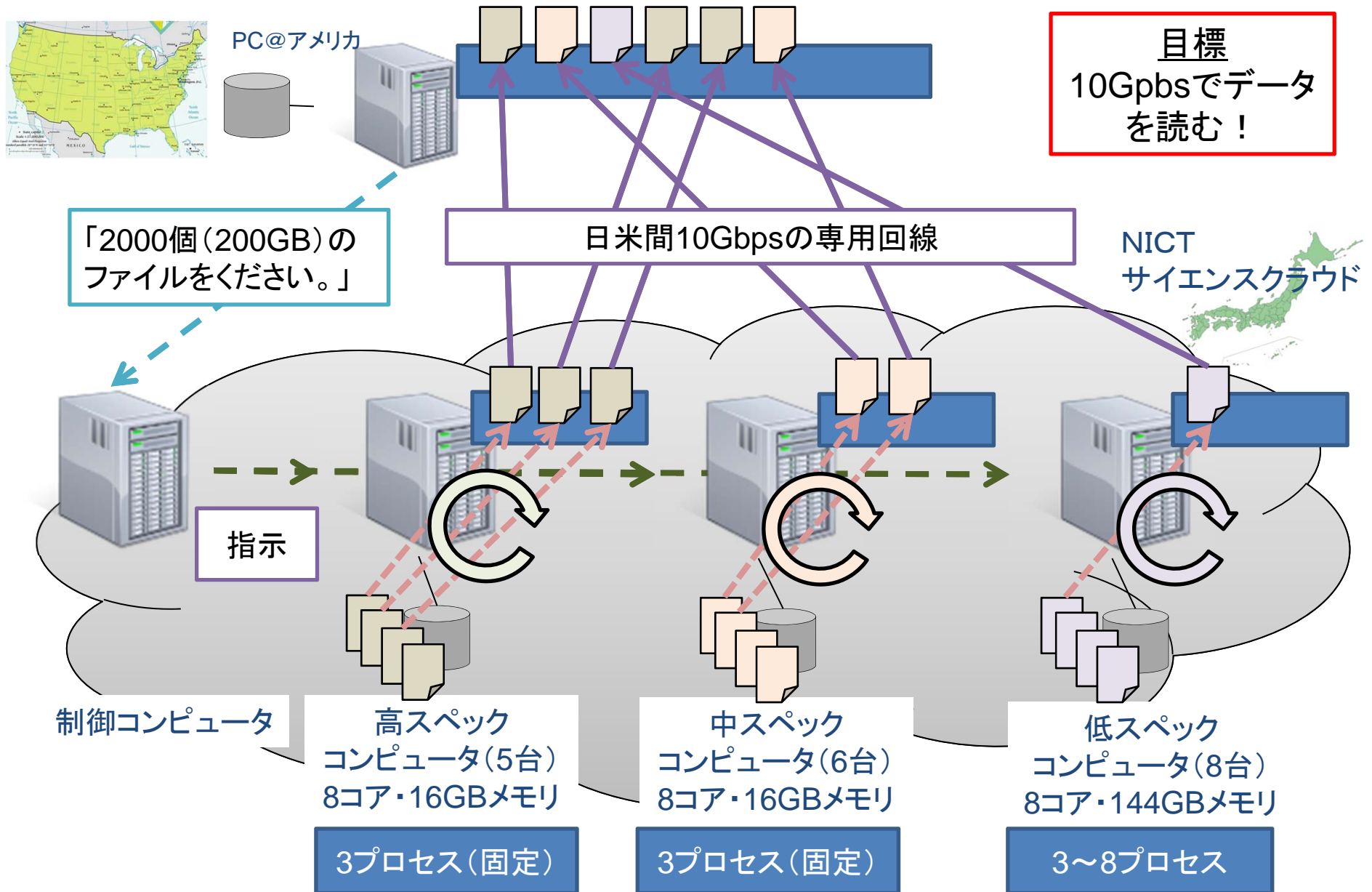
Gfarm (Grid Datafarm)/Pwrake (1)



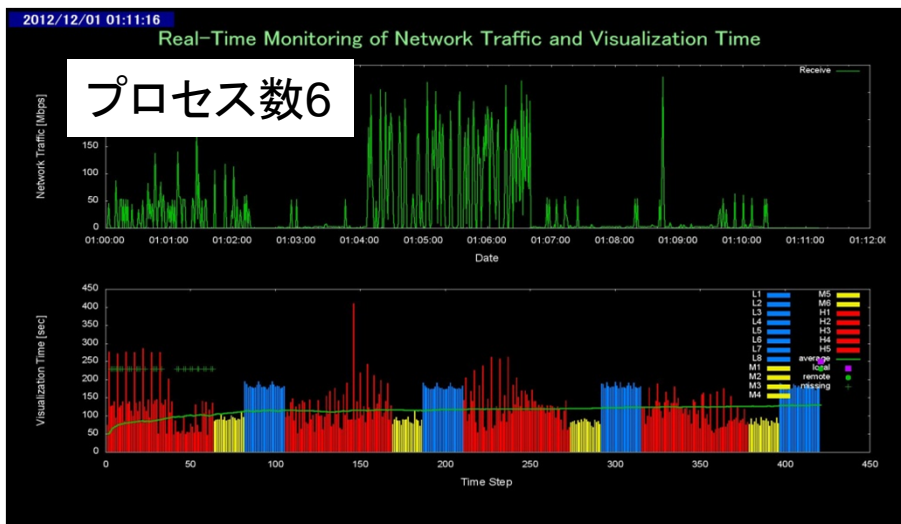
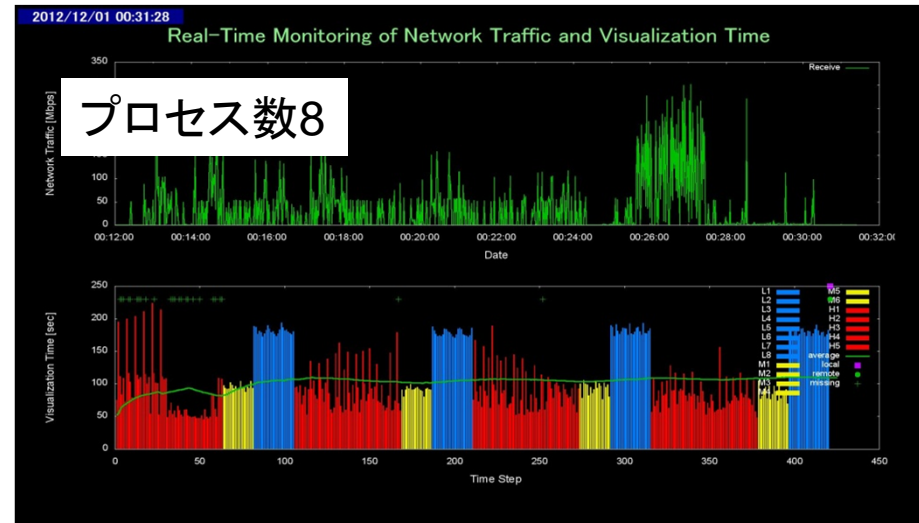
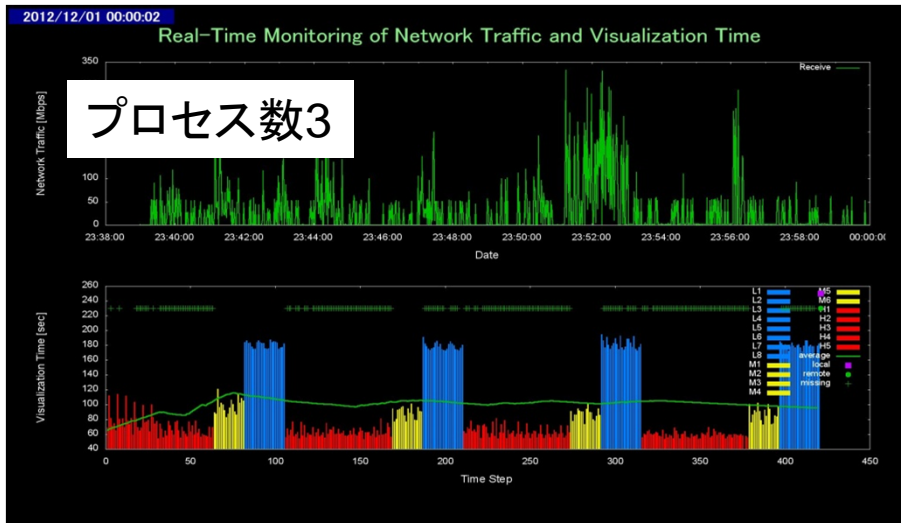
Gfarm (Grid Datafarm)/Pwrake (2)



日米間高速データ処理実験②(2012年11月)



プロセス数3/6/8の場合の比較



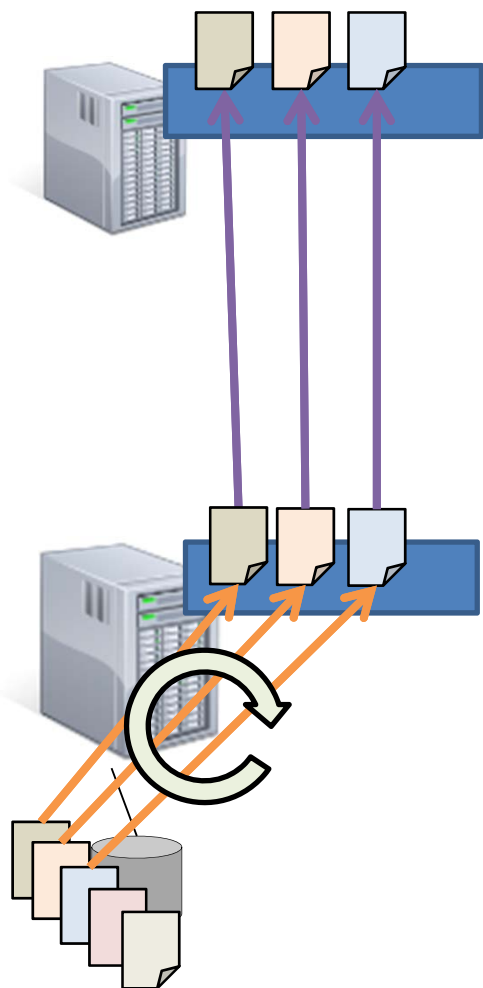
- 時刻ステップごとのデータ処理時間(縦軸は固定していません)
- 青(L)・黄(M)・赤(H)
- プロセス数が増加⇒ディスクアクセス集中が発生



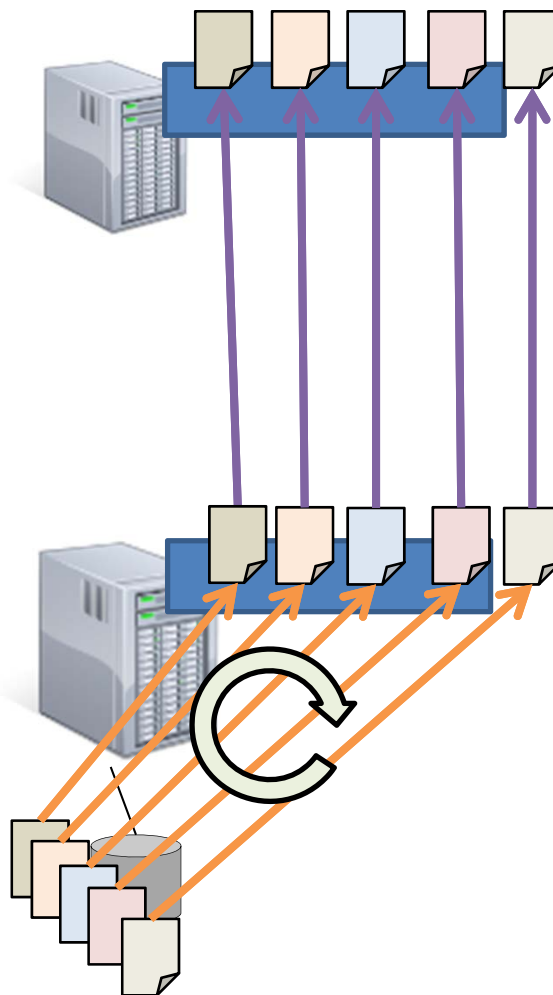
- プロセス数6を超えるとディスクアクセス集中が発生

プロセス数の設定

3プロセス並列



5プロセス並列



データ総量の増加
(メモリ容量超過)

データ通信トラ
フィックの増大

処理の並列化

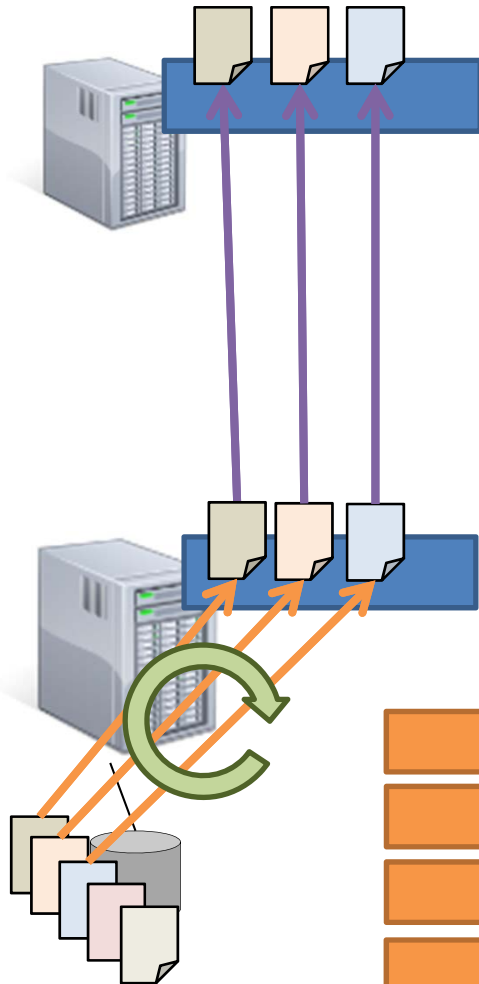
ディスクアクセスの
集中

5プロセスが有利

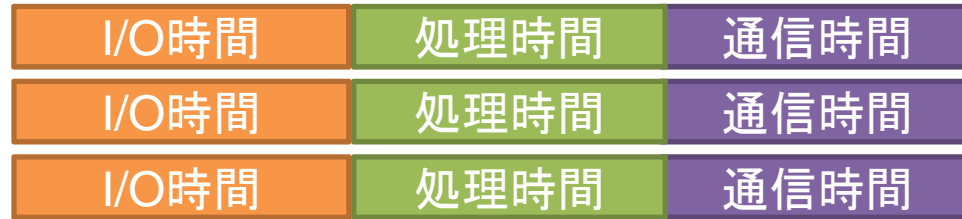
5プロセスが不利

プロセス数の最適化

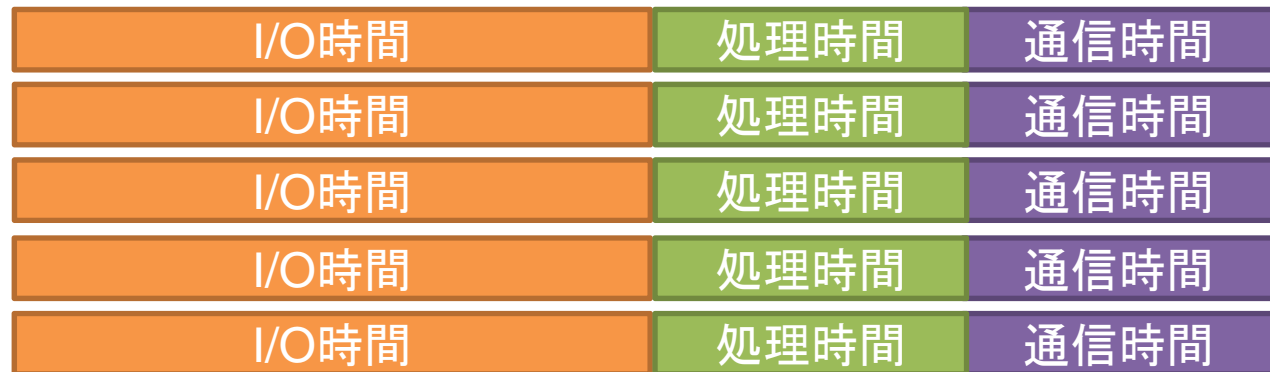
3プロセス並列の場合



3プロセス並列



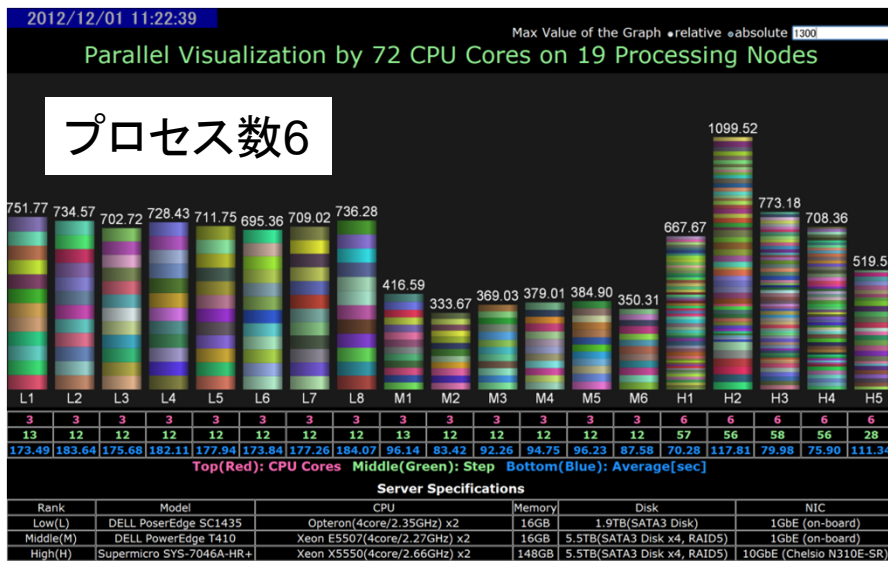
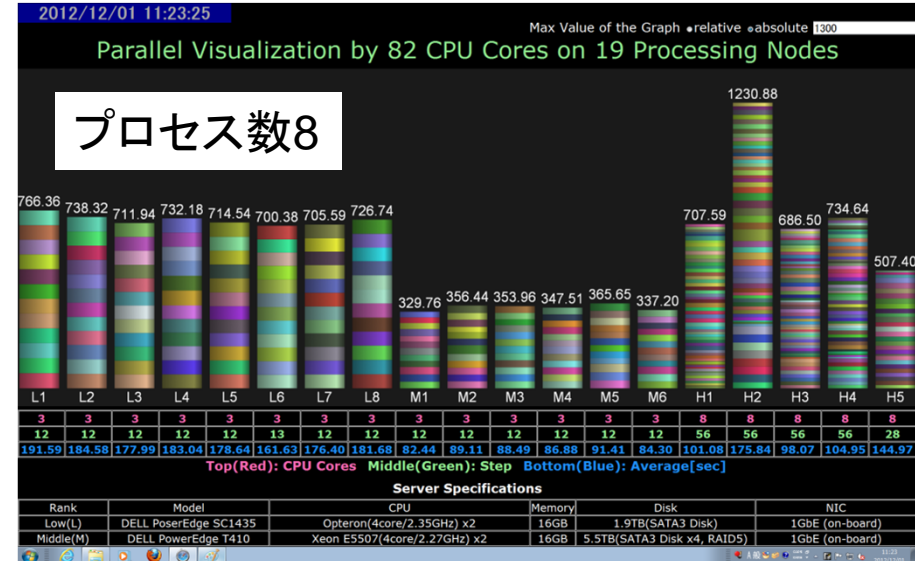
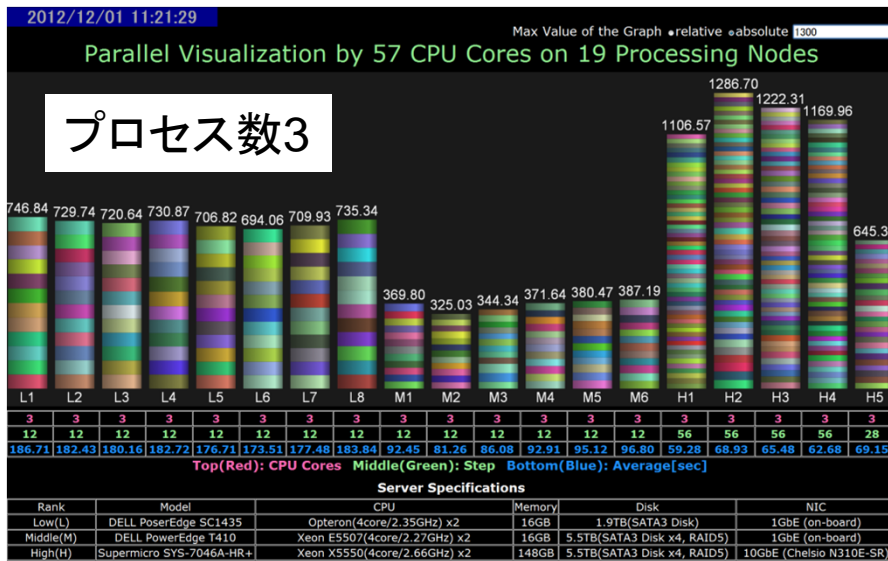
5プロセス並列



処理速度向上

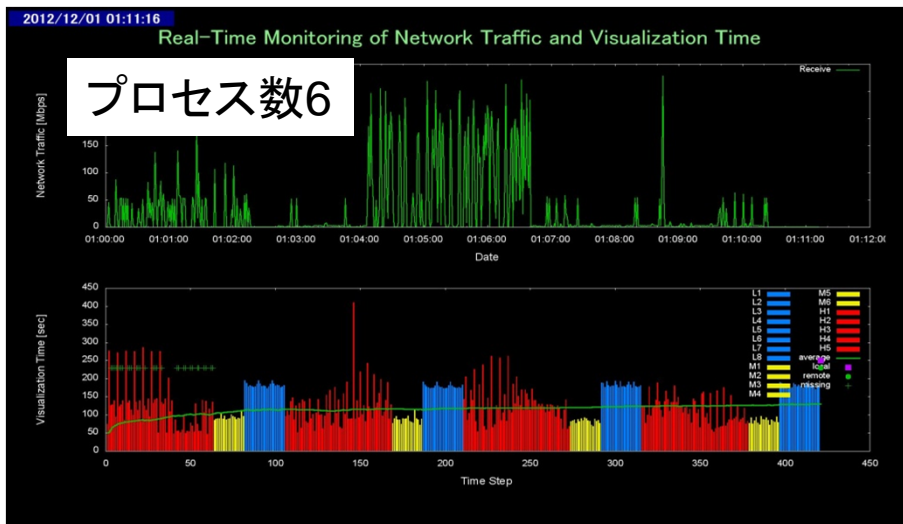
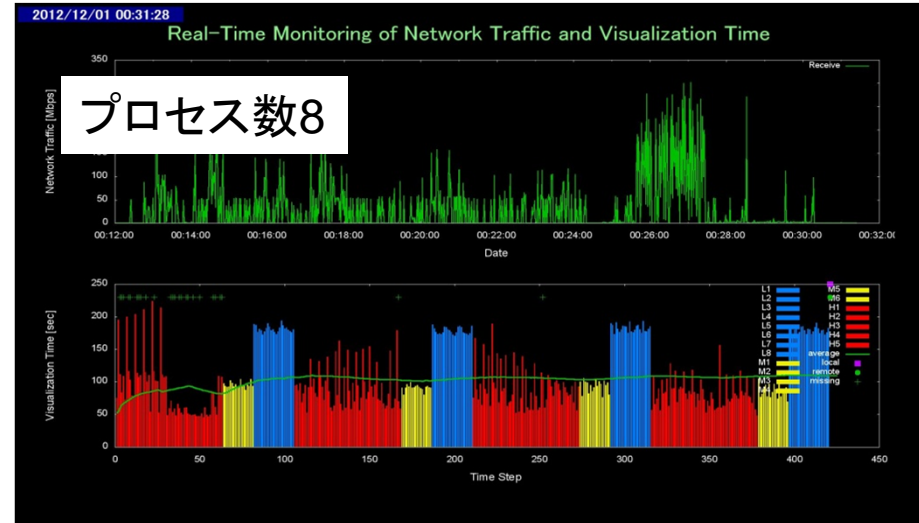
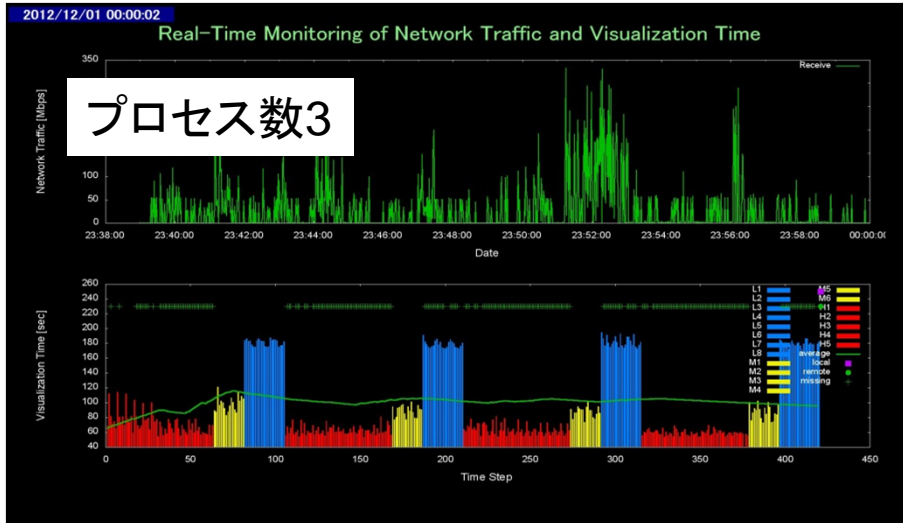
処理速度低下

プロセス数3/6/8の場合の比較: 負荷バランス



- ↓
- プロセス数6を超えるとディスクアクセス集中が発生
 - スケーラビリティとディスク集中のバランスがよいのは6プロセス処理の場合

プロセス数3/6/8の場合の比較:ステップごとのデータ処理時間

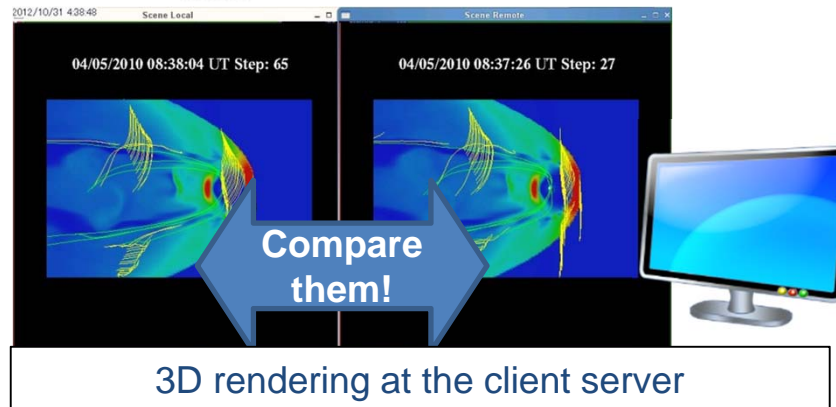
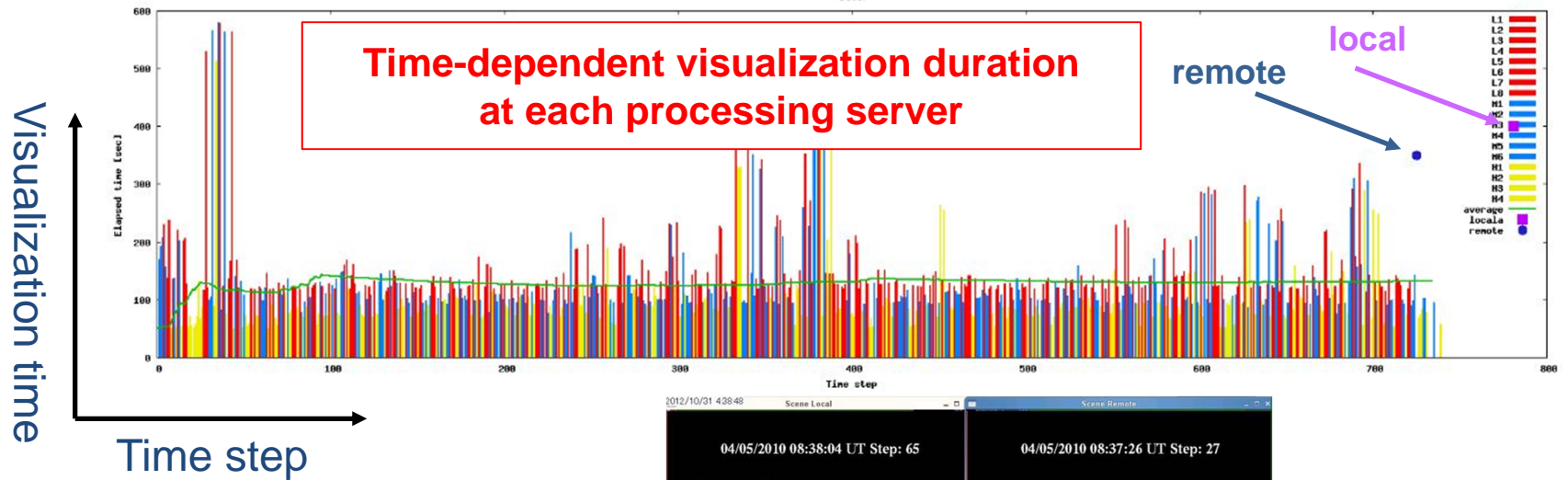
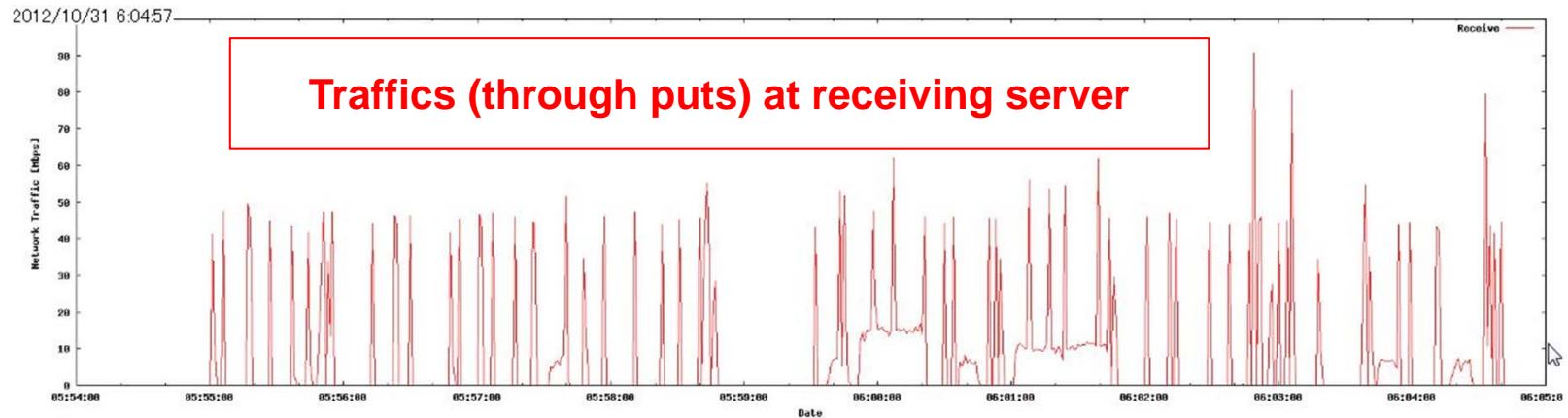


- 時刻ステップごとのデータ処理時間(縦軸は固定していません)
- 青(L)・黄(M)・赤(H)
- プロセス数が増加⇒ディスクアクセス集中が発生

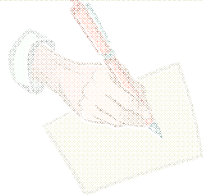


- プロセス数6を超えるとディスクアクセス集中が発生

(参考)FIFO型の場合



①理論→紙と鉛筆



④データ指向型研究(インフォマティクス)⇒クラウド

データ指向型科学を実現するには...

どうやってデータを集めるか？

どうやってデータを保存するか？

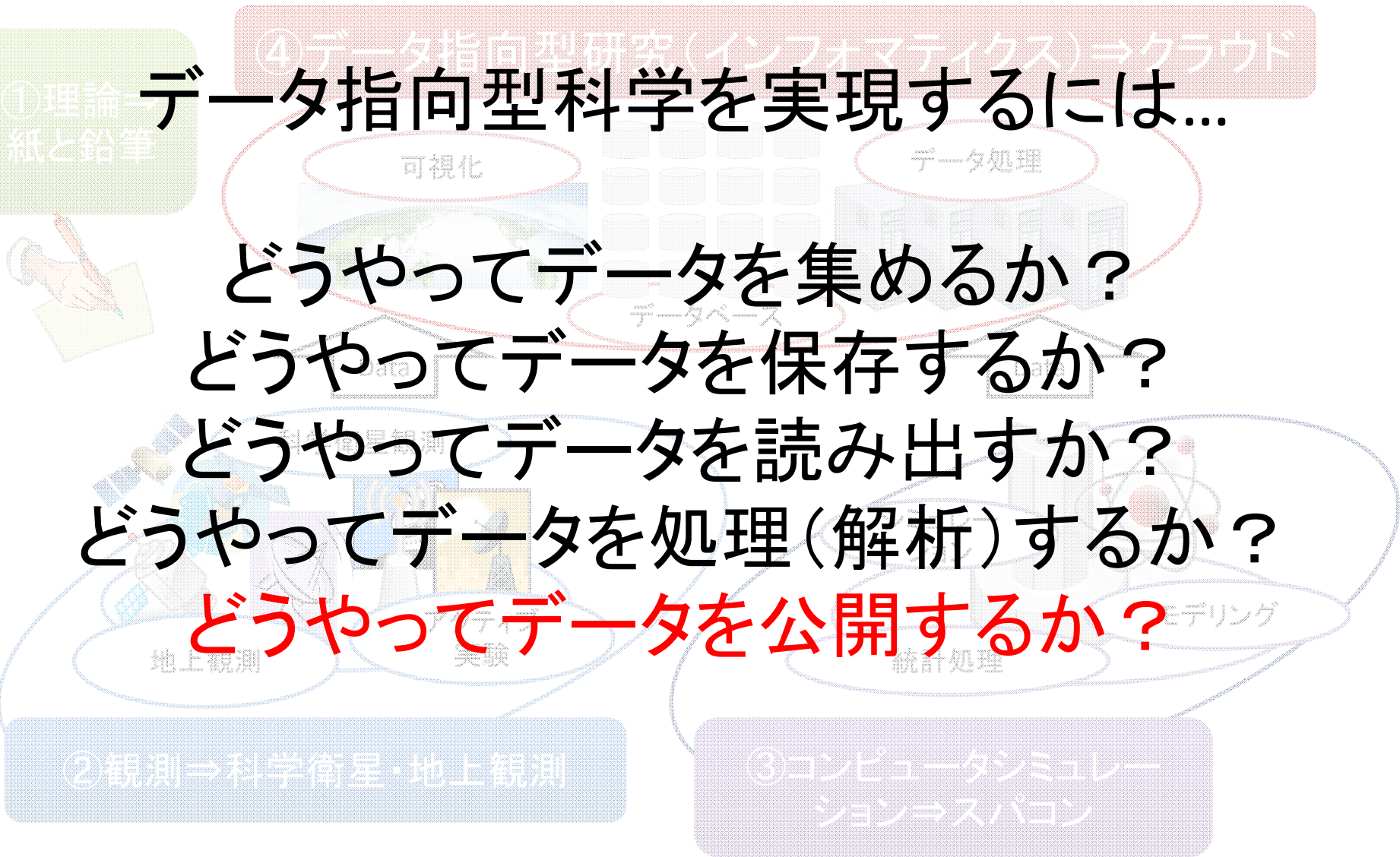
どうやってデータを読み出すか？

どうやってデータを処理(解析)するか？

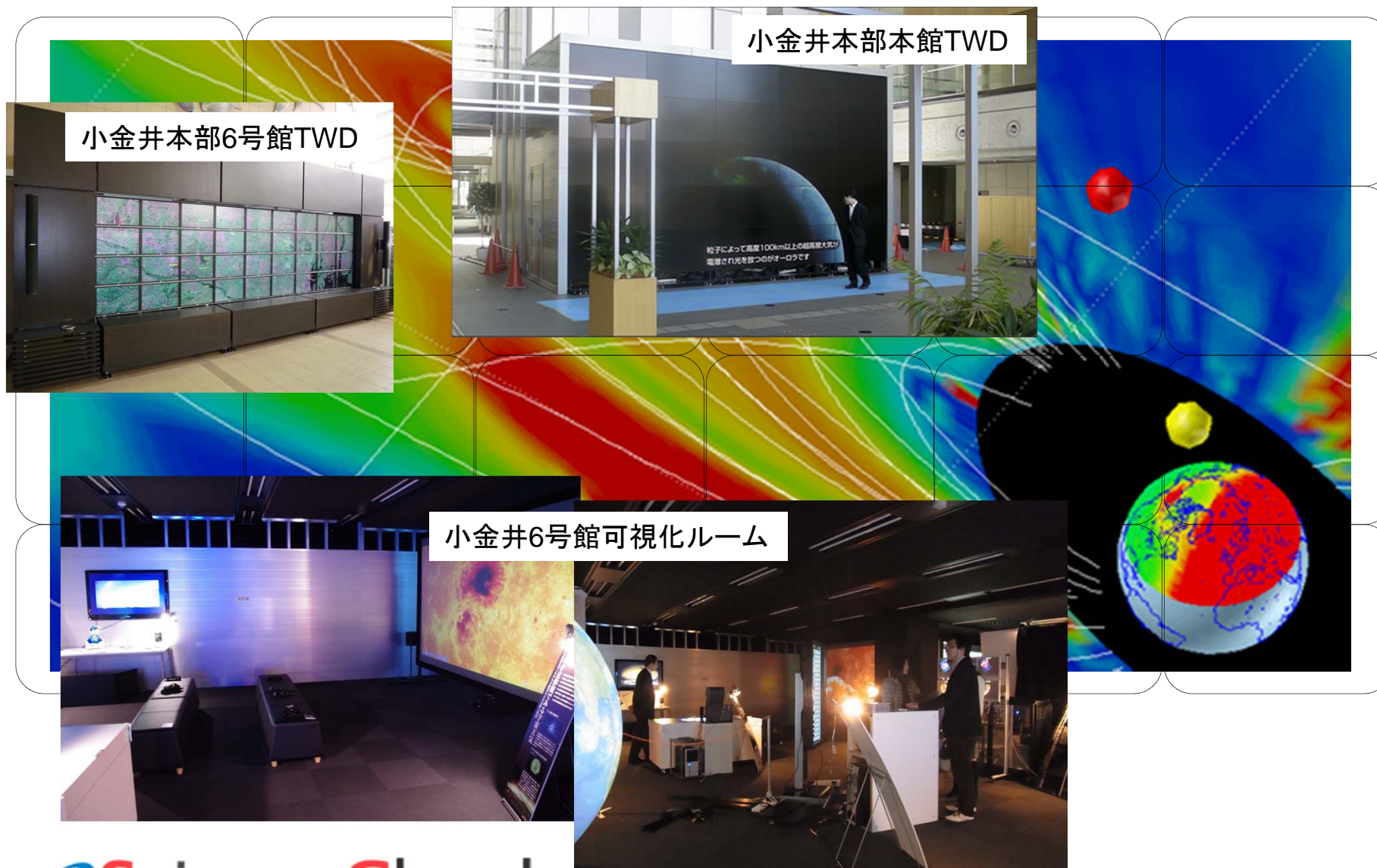
どうやってデータを公開するか？

②観測⇒科学衛星・地上観測

③コンピュータシミュレーション⇒スパコン



どうやってデータを公開するか？



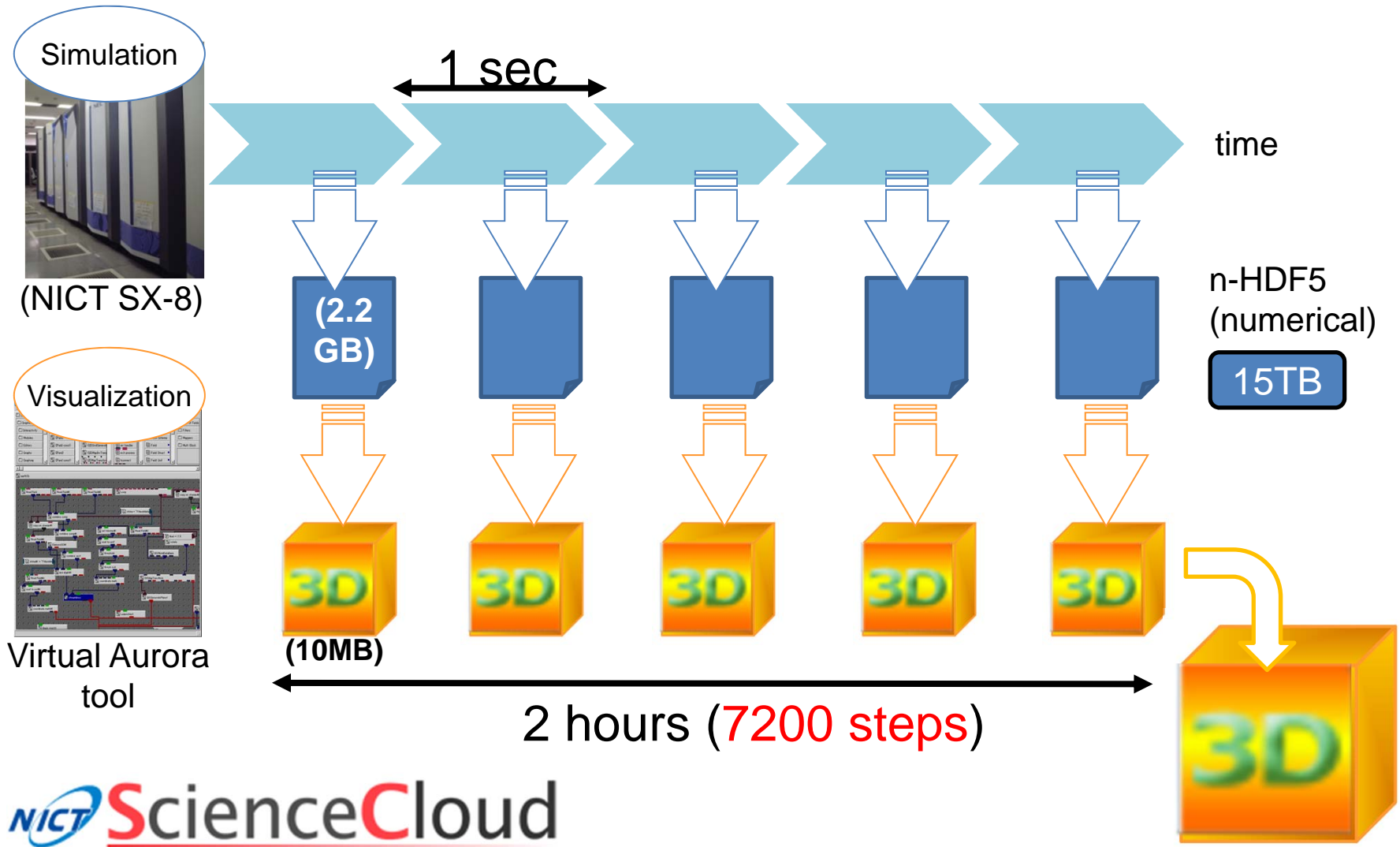
Visualization of global MHD simulation with extremely high time resolution

Y. Kubota and space weather team

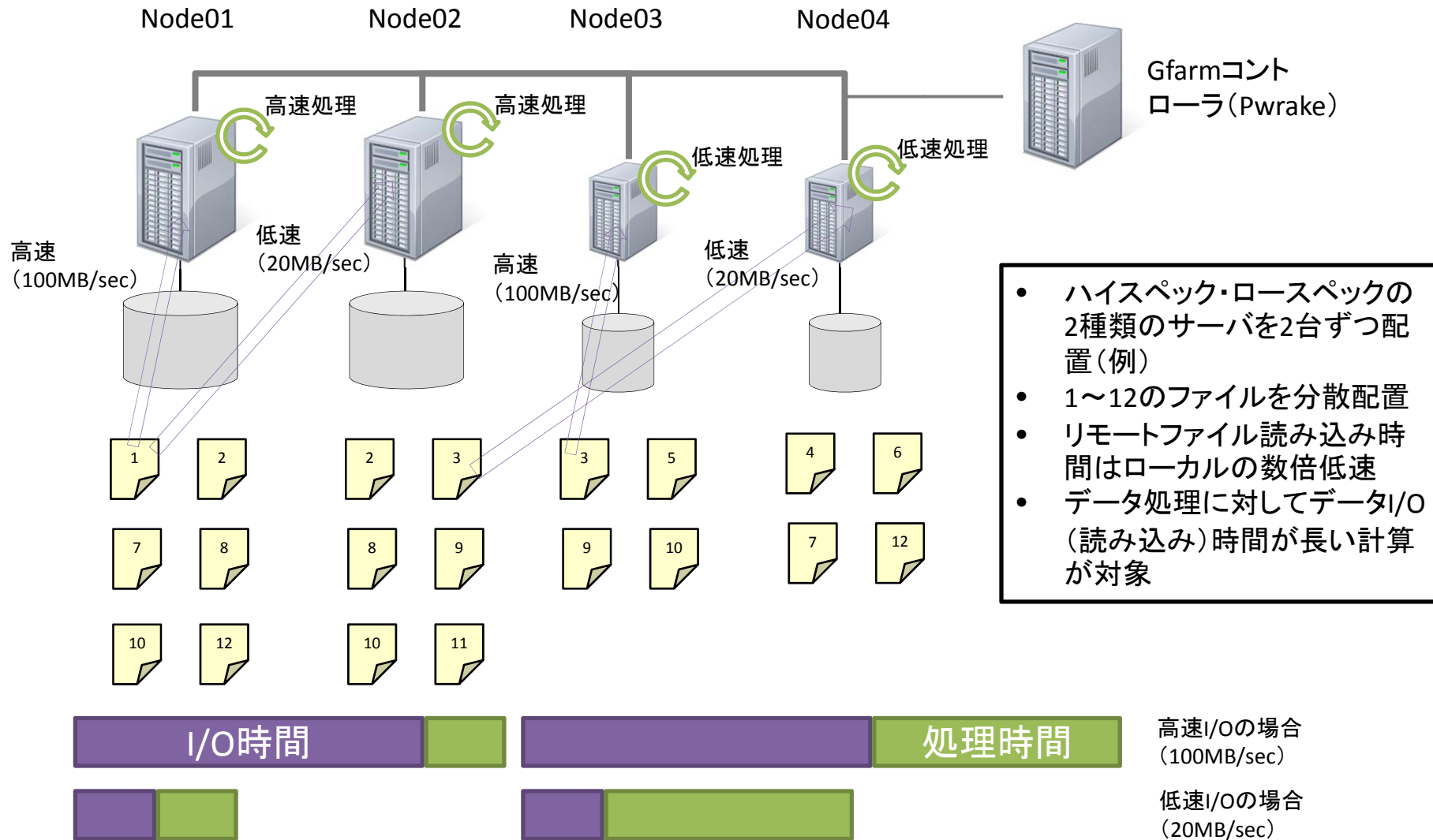
*National Institute of Information and Communications Technology
Applied Electromagnetic Research Institute
Space Weather and Environment Informatics Laboratory
4-2-1, Nukui-Kitamachi, Koganei, Tokyo 184-8795, Japan
TEL: +81-42-327-7931 FAX: +81-42-327-6978
E-mail: SciCloud-office@ml.nict.go.jp*

Simulation data size (typical)

• Spatial and time resolution
450(x) × 300(y) × 300(z) –uniform grid (dx=0.2Re, dt=0.5sec)



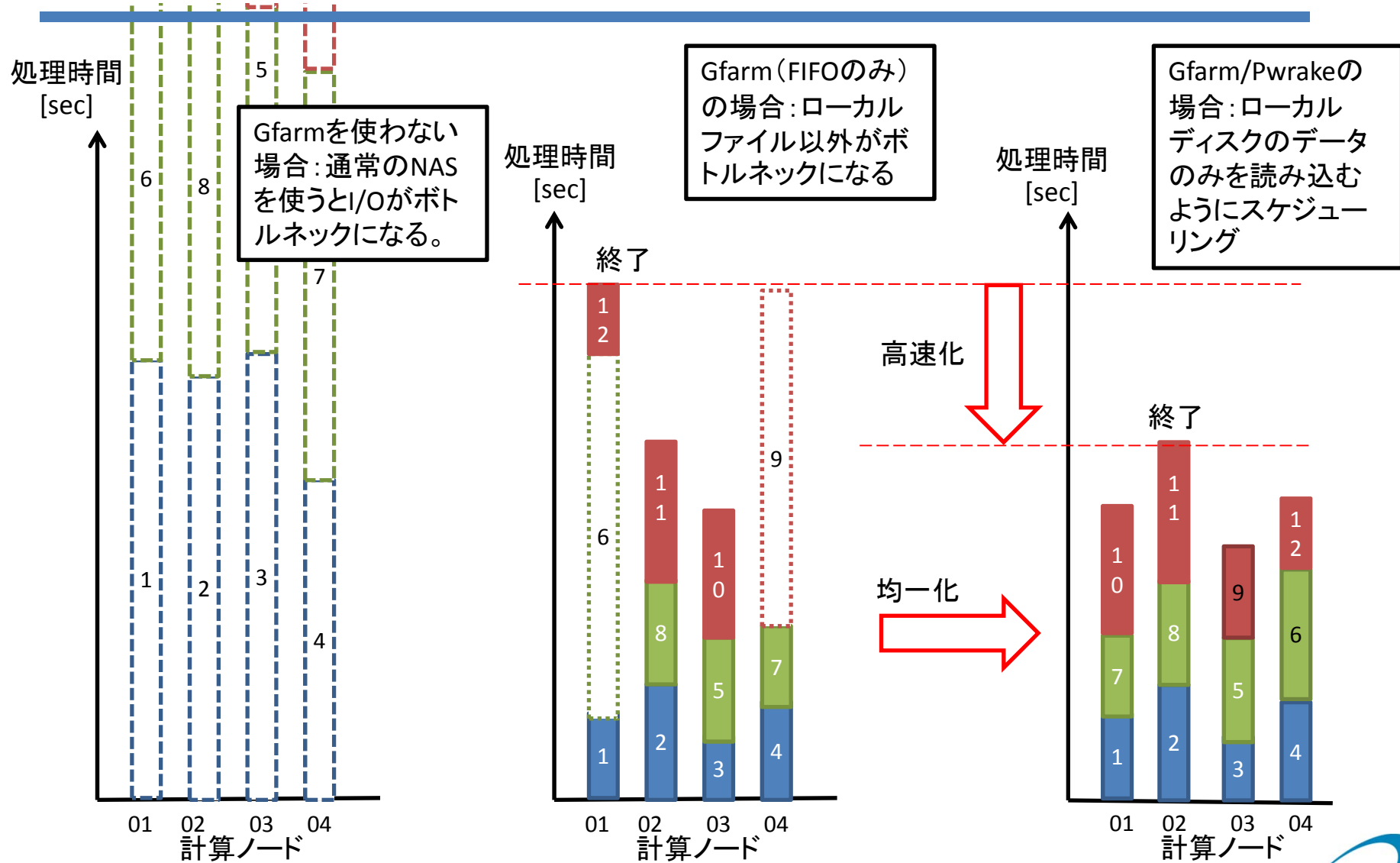
Gfarm/Pwrakeによる並列処理(モデル例)





Gfarm/Pwrakeによるノード処理最適化

FIFO型と非均一分散型の組み合わせ例



Distributed Storage/ Parallel Visualization of 7200 time step data (data size: 15TB)

Traditional Method (36 days)

Data Read Time (I/O Time)
 20MB/sec -> 1.5 million sec. -> 18 days

Data Processing Time
 14400 steps x Tracing time (10 sec.) -> 14 thousands sec. -> 2 days

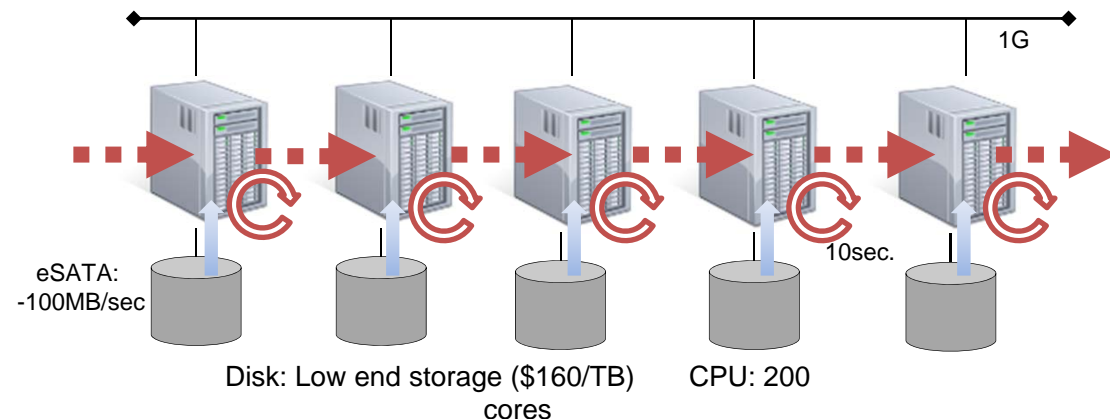
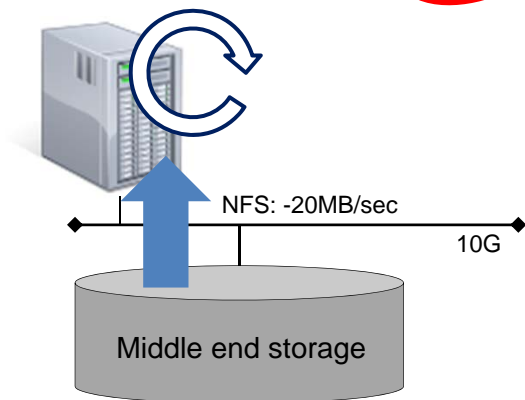
Visualization Time
 14400 steps x 120sec -> 16 days

Parallel Visualization (2 days)

Data Read Time (I/O Time)
 100MB/sec x 200 cores -> 1500 sec.
 -> 30 min.

Data Processing Time
 14400 steps x Chasing time (10 sec.) -> 14 thousands sec. -> 2 days

Visualization Time
 14400 steps x 120sec./200 cores -> 4 hours



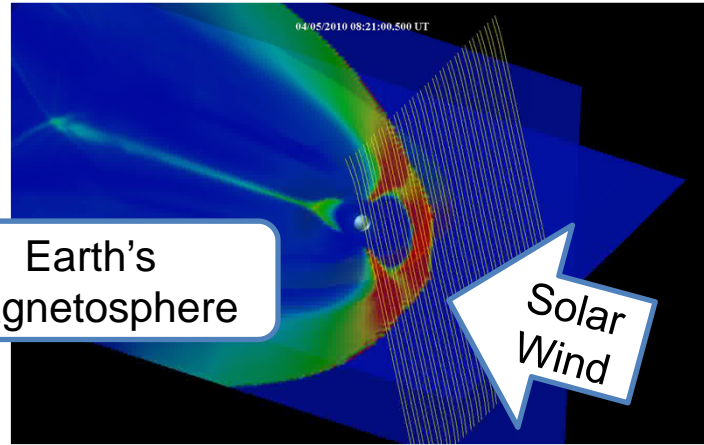
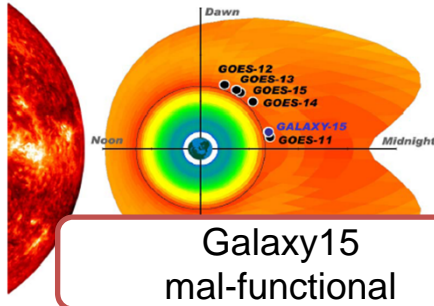


“Crown-milk view” of reconnections at the dayside magnetosphere

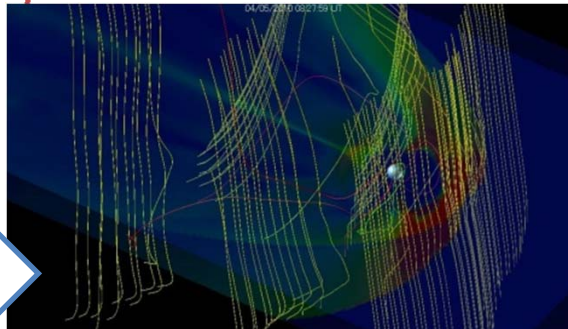
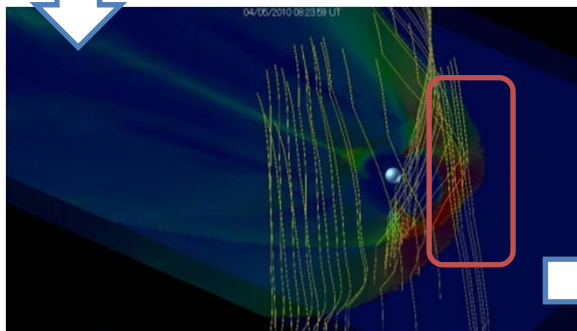
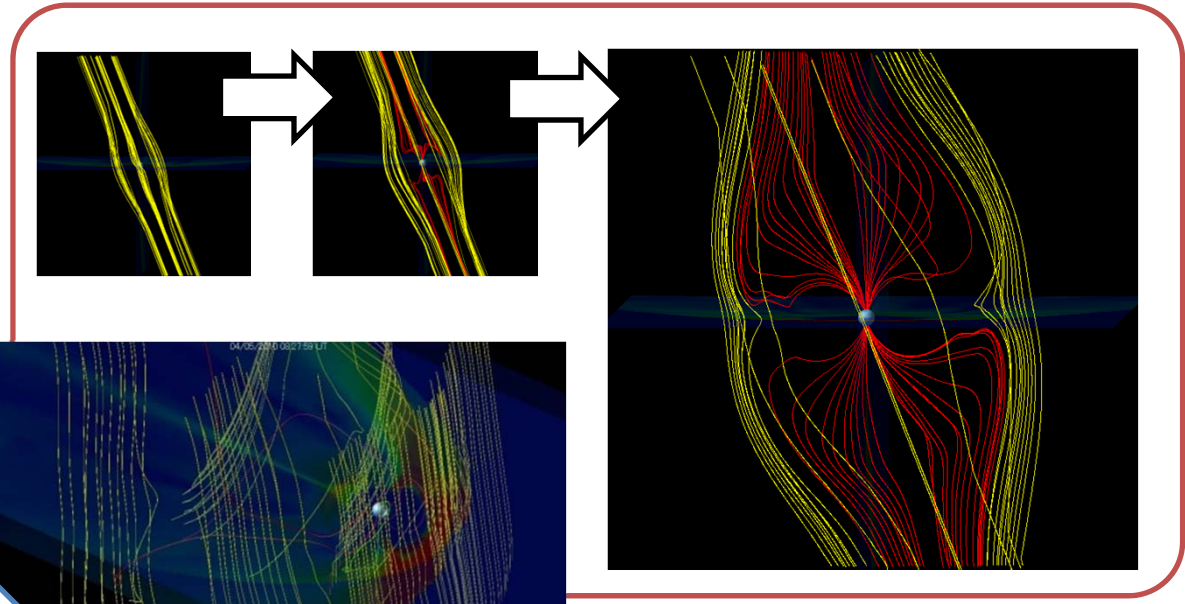
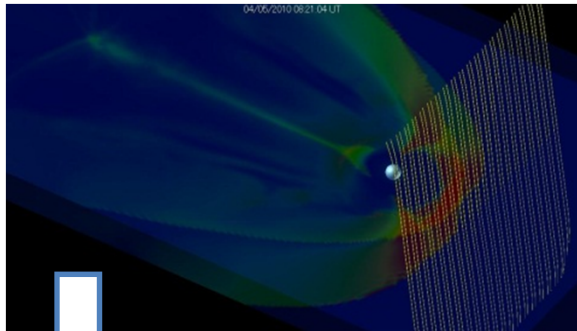
The first successful visualization of solar-wind transfer into magnetosphere

Galaxy 15 (133 W) Anomaly 09:48 UT

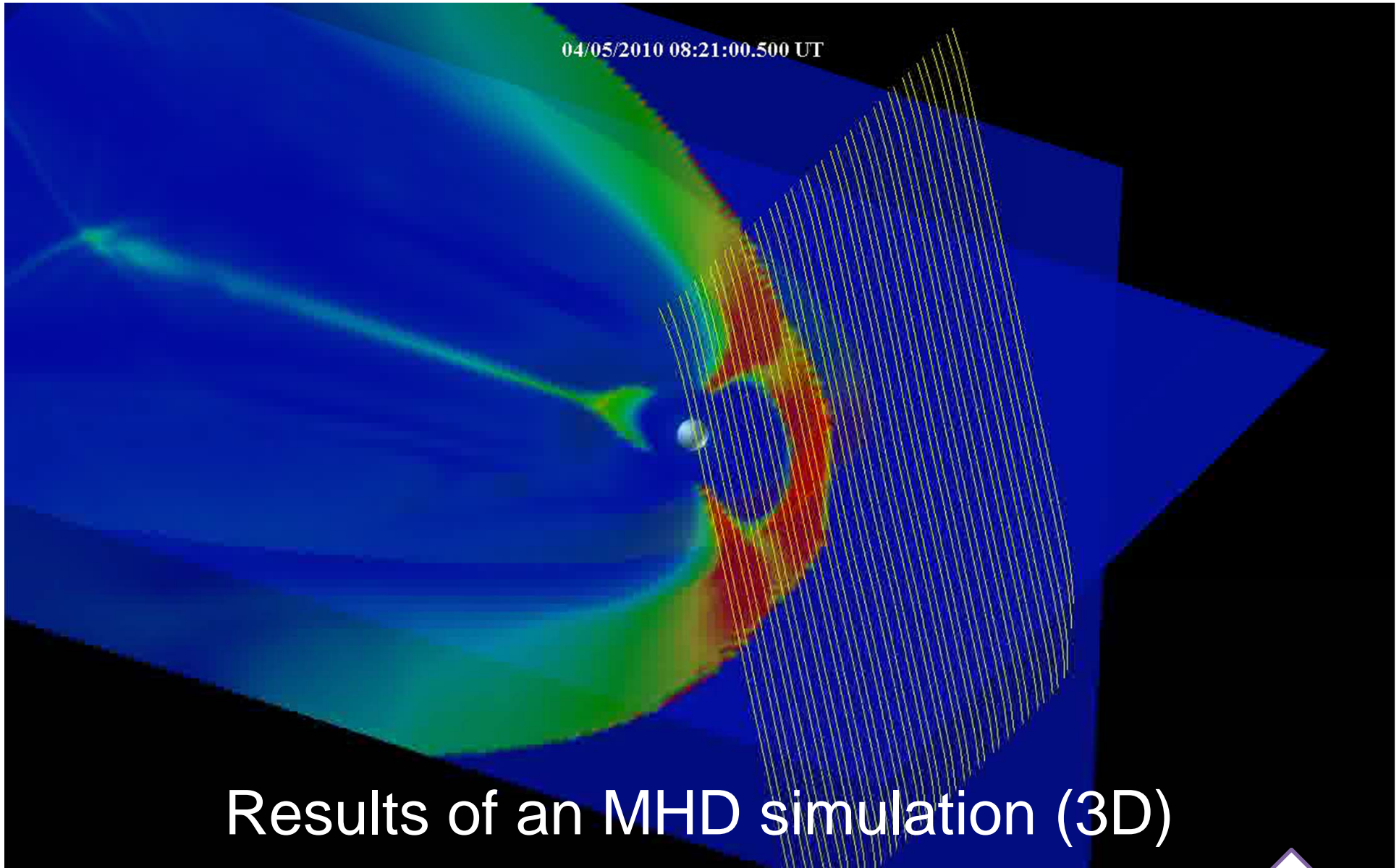
Satellite Locations



Open Question
How does the solar-wind penetrate into Magnetosphere?



04/05/2010 08:21:00.500 UT



Results of an MHD simulation (3D)

Global MHD simulation (1 hour=7200 steps)

 NICT ScienceCloud

64bit 3D player
developed by
NICT

④データ指向型研究(インフォマティクス)⇒クラウド

①理論⇒紙と鉛筆

データ指向型科学を実現するには...

どうやってデータを集めるか？

どうやってデータを保存するか？

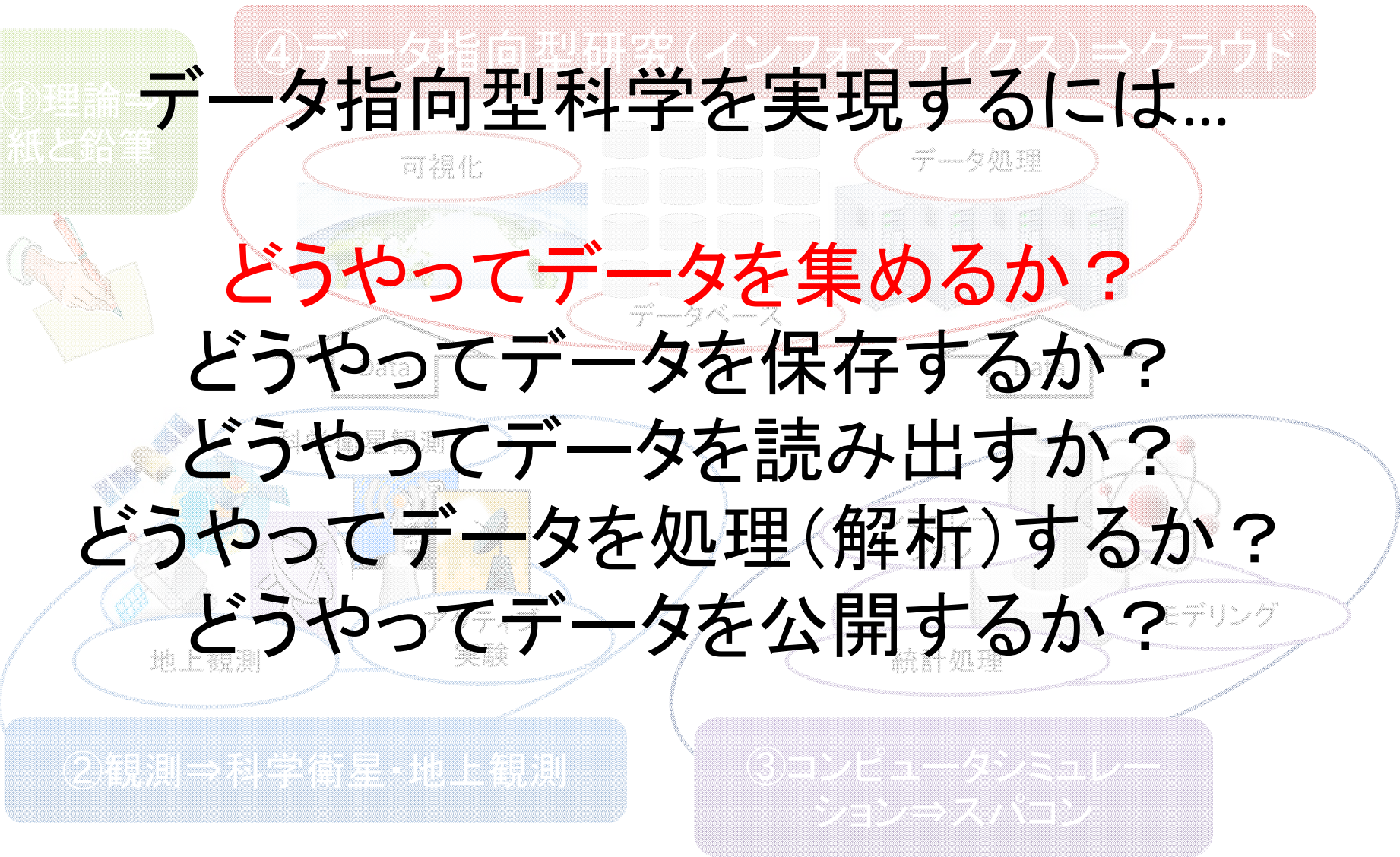
どうやってデータを読み出すか？

どうやってデータを処理(解析)するか？

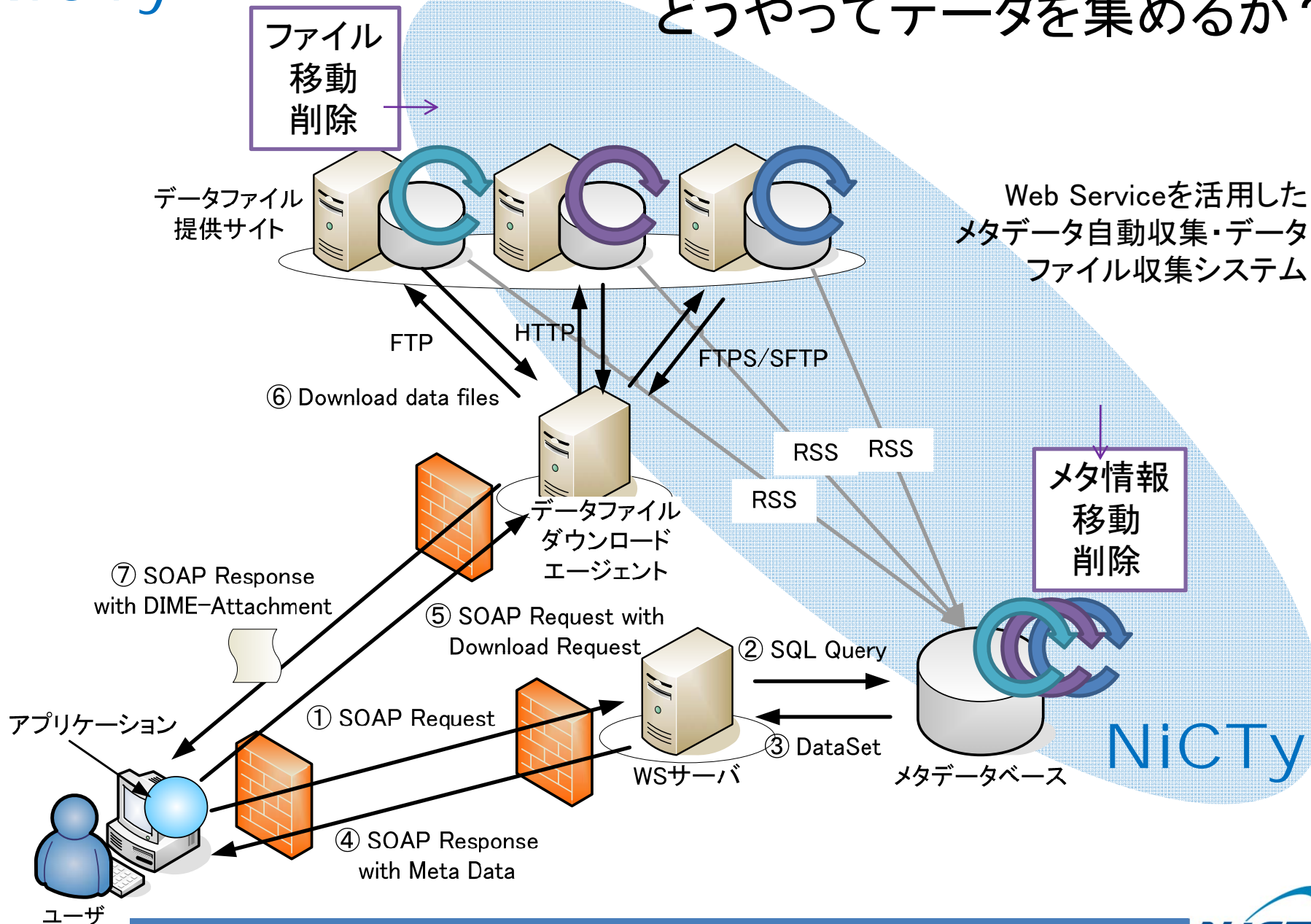
どうやってデータを公開するか？

②観測⇒科学衛星・地上観測

③コンピュータシミュレーション⇒スパコン



どうやってデータを集めるか？



①理論→紙と鉛筆

データ指向型科学を実現するには...

④データ指向型研究(インフォマティクス)⇒クラウド

可視化

データ処理

どうやってデータを集めるか？

どうやってデータを保存するか？

どうやってデータを読み出すか？

どうやってデータを処理(解析)するか？

どうやってデータを公開するか？

②観測⇒科学衛星・地上観測

③コンピュータシミュレーション⇒スパコン

地上観測

実験

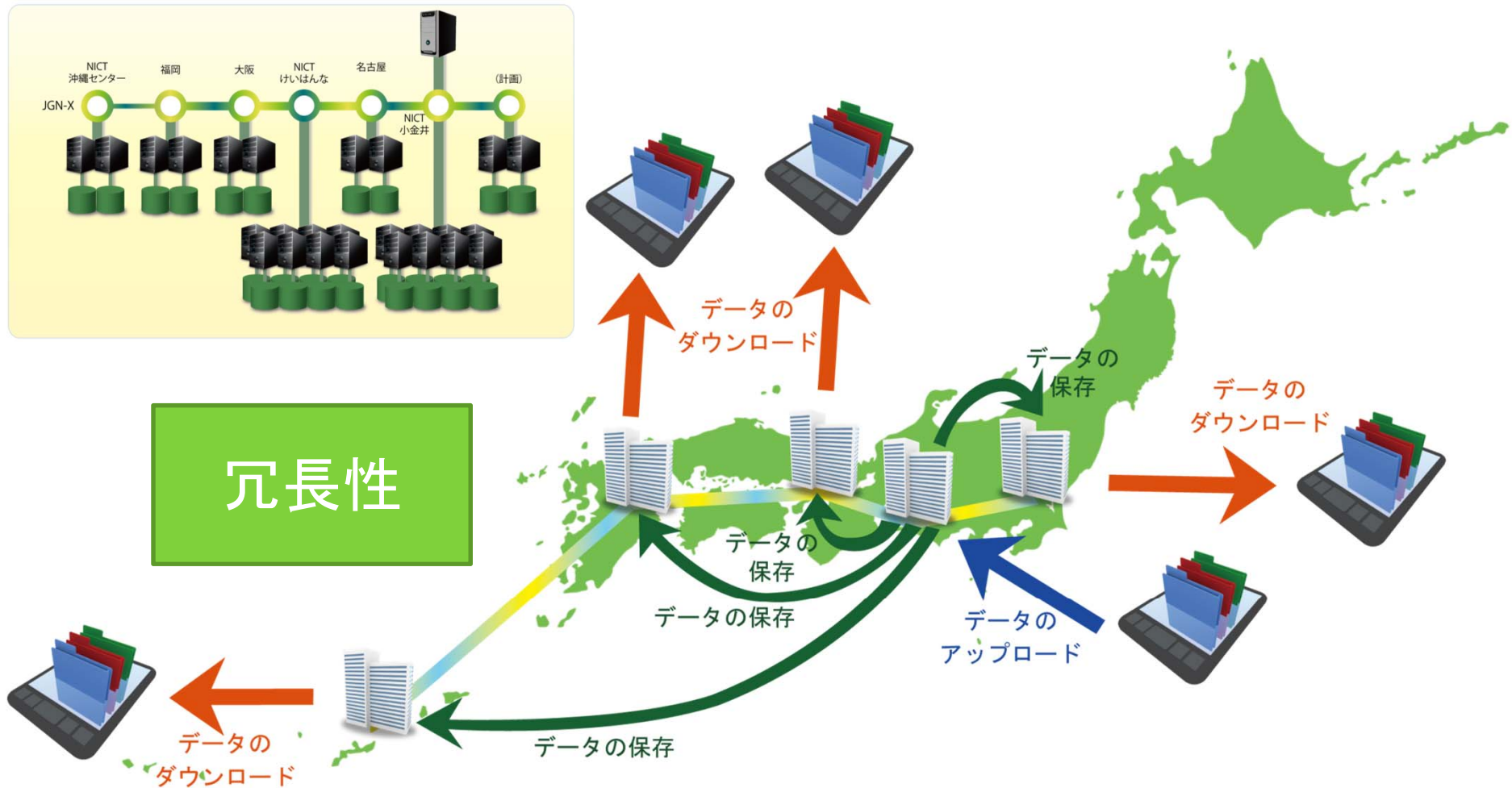
統計処理

モデリング

データベース

どうやってデータを保存するか？

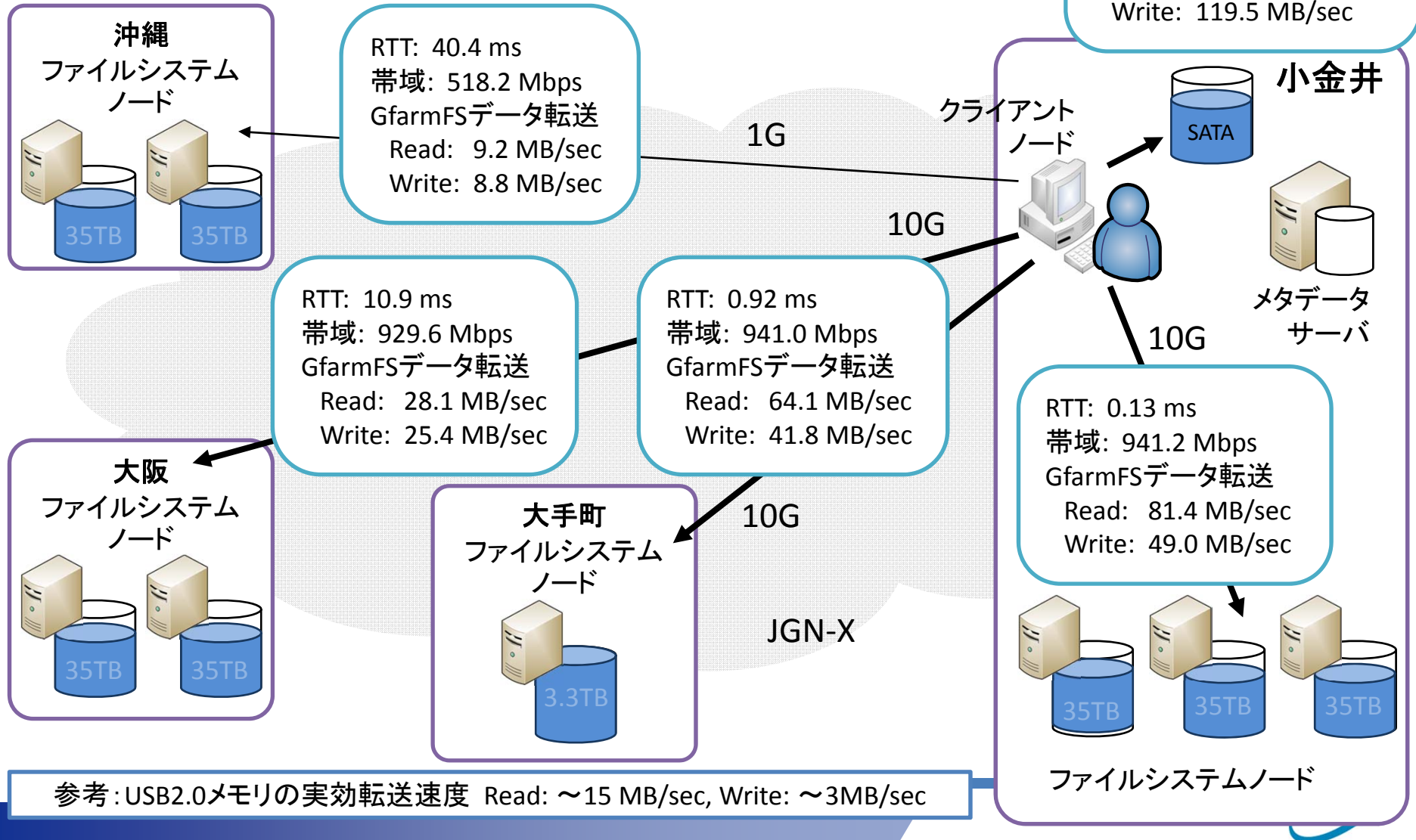
NICTサイエンスクラウド内のデータ保存



ネットワーク試験結果：NICTクラウド内分散ストレージ

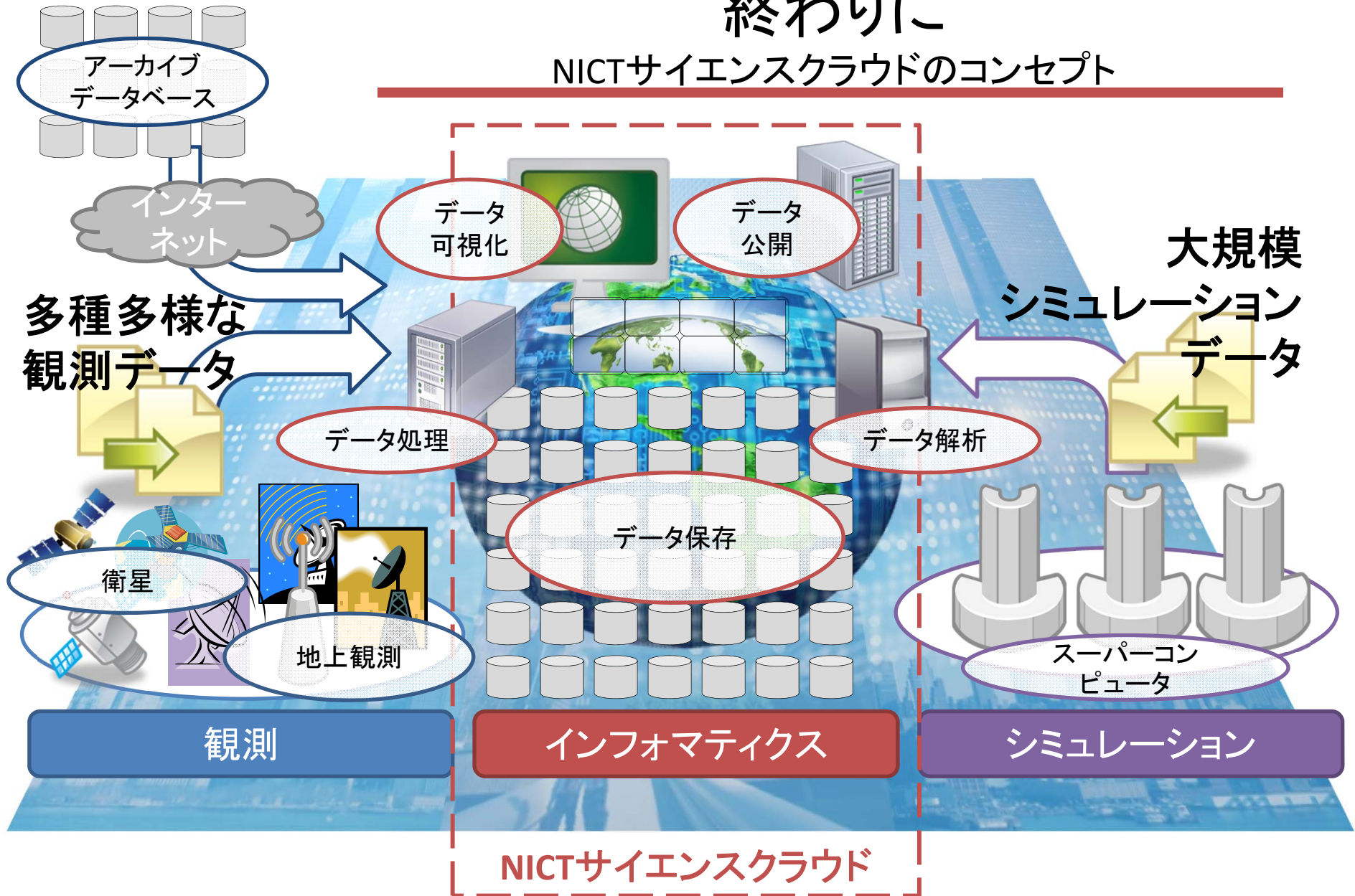
ユーザ@小金井(クライアントノード)が各データサイト(ファイルシステムノード)の10MBのファイルをRead/Write

RTT: 0.0 ms
(SATA理論値:600Mbps)
GfarmFSデータ転送
Read: 358.6 MB/sec
Write: 119.5 MB/sec



終わりに

NICTサイエンスクラウドのコンセプト



クラウドは総合技術！

ご清聴ありがとうございました!