

# Operációs Rendszerek MSc

## Szuperszámítógépek

2019/2020/I.

Dr. Vincze Dávid  
Miskolci Egyetem, IIT  
[vincze.david@iit.uni-miskolc.hu](mailto:vincze.david@iit.uni-miskolc.hu)

# Operációs Rendszerek MSc

⇒ Miről lesz szó?

⇒ Szuperszámítógépek

- Symmetric Multiprocessing (SMP)
- Single System Image (SSI)
- Massive(ly) Parallel Processing (MPP)
- Interconnect

⇒ Klaszterek

# Operációs Rendszerek MSc

- ⇒ Szuperszámítógépek
- ⇒ Egy helyre koncentrált hatalmas erőforrás
  - Sok CPU, sok RAM, nagy I/O teljesítmény
  - Speciális igényekhez módosított architektúra
    - IBM zSeries
    - Sun Enterprise
    - stb.
  - Egy backplane-re csatlakoznak a modulok
    - általában passzív, és nem redundáns
      - de fizikailag tönkre tehető... :)
  - Szervízprocesszor felügyel
  - Magas rendelkezésre állás, redundancia

# Operációs Rendszerek MSc

- ⇒ Szuperszámítógépek
  - Több fizikai gépből felépítve
  - Speciális nagy sávszélességű **interconnect**
    - pl. InfiniBand, Omni-Path, Tofu, NUMALink (CrayLink)
  - Ez a **Massively Parallel Processing** (MPP)
    - Manapság többnyire ilyenek a mai szuperszámítógépek
  - GPU (Graphical Processing Unit)
    - pl. NVIDIA Tesla / Volta
  - MIC (Many Integrated Core)
    - pl. Xeon Phi
  - NUMA architektúra
    - Non-Uniform Memory Access

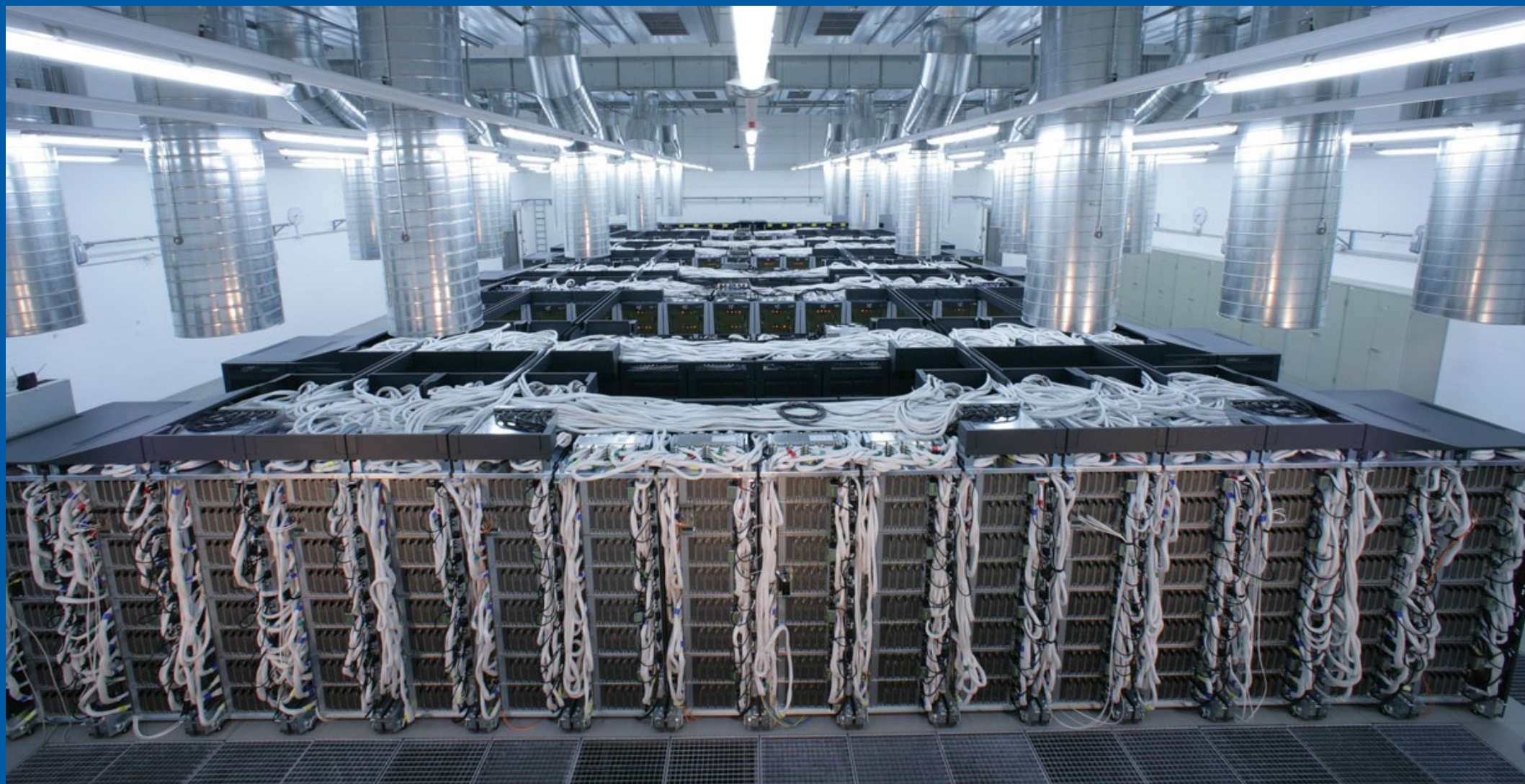


Blue Gene/P

Blue Gene/P

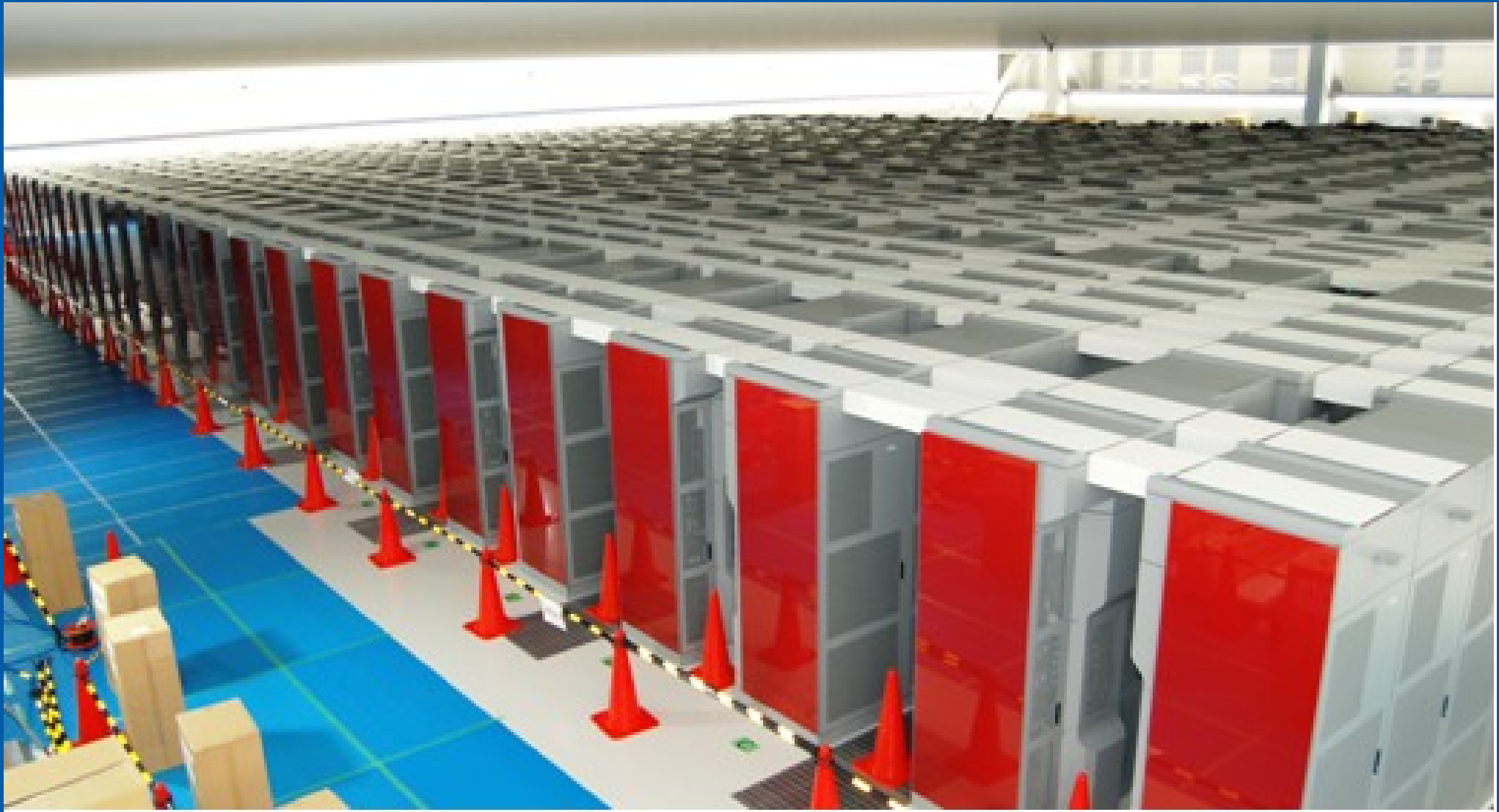
Blue Gene/P

# Operációs Rendszerek MSc



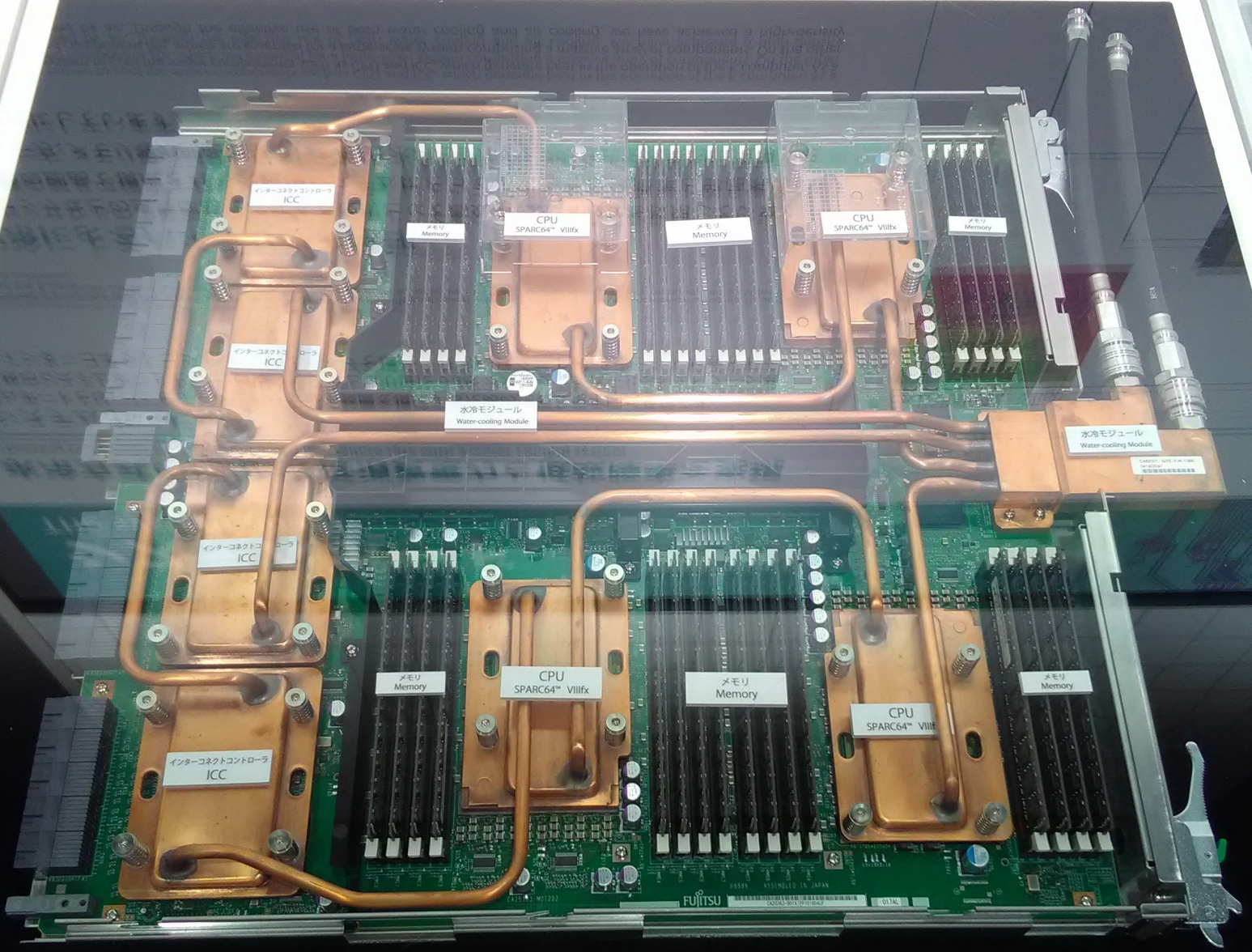


# K Computer



YouTube <https://www.youtube.com/watch?v=UJPslu9OaTc>





システムボード  
System Board



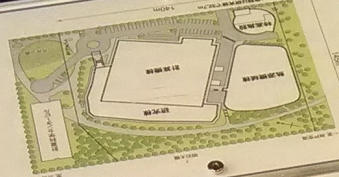


### 世界最高水準の「ハイコベ」施設として

研究棟  
研究棟は、最先端の研究設備を備え、最新の研究成果を生み出すための重要な役割を果たします。

特設設備  
特設設備は、最新の研究成果を生み出すための重要な役割を果たします。

熱環境制御  
熱環境制御は、最新の研究成果を生み出すための重要な役割を果たします。



1 耐震性 (Seismicity)  
震害防止の観点から、耐震設計を徹底して実施しています。

2 耐久性 (Durability)  
建築環境に合わせた素材と工法を採用し、長寿命を実現しています。

3 省エネルギー (Energy Saving)  
最新の省エネルギー設備を導入し、環境負荷を低減しています。

4 自然環境との調和 (Harmony with Nature)  
自然環境との調和を図り、快適な環境を実現しています。

5 安全安心 (Safety and Security)  
最新のセキュリティ設備を導入し、安全安心を実現しています。

6 快適性 (Comfort)  
最新の快適設備を導入し、快適な環境を実現しています。

7 先進性 (Advanced)  
最新の先進設備を導入し、最先端の研究環境を実現しています。

8 持続可能性 (Sustainability)  
最新の持続可能性設備を導入し、持続可能な環境を実現しています。

研究棟の概要  
研究棟は、最先端の研究設備を備え、最新の研究成果を生み出すための重要な役割を果たします。



<https://www.google.com/maps/place/K%C3%B3be,+Hj%C3%B3go+prefekt%C3%B4ra,+Jap%C3%A1n/@34.6527583,135.2229209,279a,35y,292.43h,39.3t/data=!3m1!1e3!4m5!3m4!1s0x60007d812aed89d9:0xc7126106c2f670f4!8m2!3d34.690083!4d135.1955112>

天河

天河二号

TH-2 High Performance Computer System

天河

天河

天河

天河

天河

天河

天河

天河

天河

天河

天河

天河

天河

天河

天河

天河

天河

天河

天河

天河

天河

2014/04/22

每秒 5.49 亿亿次、持续计算速度每秒 3.39 亿亿次、能效比每瓦特 19 亿次...  
值计算速度、持续计算速度以及综合技术水平处于国际领先地位，是我国超级计算技术发展取得的重大  
进展。

The MilkyWay II supercomputer system, developed by the National University of Defense Technology (NUDT), is the outstanding achievement of the National 863 Program. Partly funded by Guangdong province and Guangzhou city, the MilkyWay II will be installed in the Guangzhou Supercomputing Center. Capable of a peak performance of 54.9PFlops, the MilkyWay II achieves a sustained performance of 33.9PFlops with a performance-per-watt of 1.9GFlops/W. The peak performance, the sustained performance and the level of comprehensive technology of the MilkyWay II all are in the lead. The MilkyWay II makes a great advance in the research and development of supercomputing technology in China.

天河二号由 170 个机柜组成，包括 125 个计算机柜、8 个服务机柜、13 个通信机柜和 24 个存储机柜，占地面积 720 平方米。内存总容量 1.4PB，存储总容量 12.4PB，最大运行功耗 17.8 兆瓦。

The MilkyWay II consists of a total sum of 170 cabinets including 125 compute cabinets, 8 service cabinets, 13 communication cabinets and 24 storage cabinets, covering 720 square meters. It has 1.4PB memory, 12.4PB storage capacity, and the power consumption of 17.8MW at its peak.

## 系统特点 System Features

(1) 高性能：天河二号计算速度世界领先，排名世界超级计算机 500 强第一位。  
High Performance: The performance of the MilkyWay II is world leading, which is ranking No.1 on the TOP500 list.

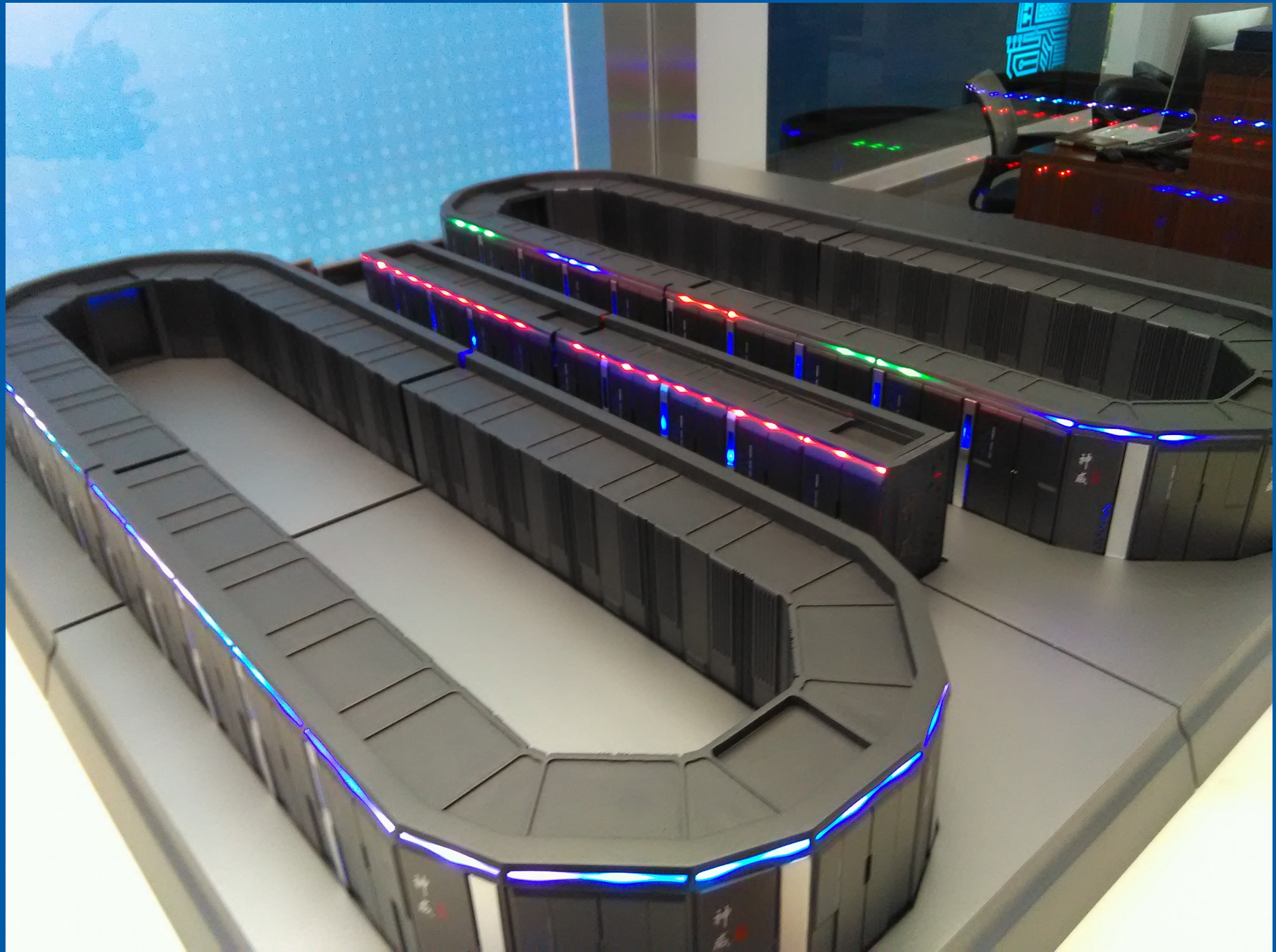
(2) 高能效：天河二号节能降耗水平处于国际先进行列。  
High Energy Efficiency: The energy efficiency of the MilkyWay II is at advanced level in the world.

(3) 应用面广：天河二号采用新型异构多态体系结构，可有效满足科学工程计算、大数据处理、高吞吐率和高安全的信息服务等多类应用需求。  
Massive Applications: A Multipurpose-Heterogeneous Architecture is adopted, which effectively meets the needs of big data processing, high throughput and high security information service as well as scientific and engineering computing.

(4) 易用性好：天河二号采用新型微异构计算...  
低了应用软件编程难度...

2014/04/22

Ren



You  <https://www.youtube.com/watch?v=KEdsrT1mFAU>

You  <https://www.youtube.com/watch?v=OU68MstXsrl>



YouTube <https://www.youtube.com/watch?v=RIKZyF9WIH>  
4



# HP cluster

Két szuperszámítógép a Miskolci Egyetemen  
(SGI UV2000, HP cluster)

26-node HP cluster (312-core, 1,2 TB RAM)

- 24 node
  - 2x 6-core Intel Xeon X5650 2.67GHz
  - 48GB ECC RAM
- 2 node
  - 2x 6-core Intel Xeon X5670 2.93GHz
  - 48GB ECC RAM
  - NVIDIA GPU
- InfiniBand



# SGI UV2000

Két szuperszámítógép a Miskolci Egyetemen  
(SGI UV2000, HP cluster)

SGI UV2000 – SSI

- 3x8x2x8 mag (Intel Xeon E5-4627 v2 @ 3.33 GHz)
- 1.4TB RAM
- 240TB HDD
- NUMALink 6 (4x 6.5Gbyte/s)

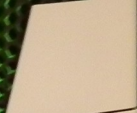
SGI UV 2000

Altix UV 1000

sgi

sgi

INFORMACIONEN  
& FELSŐKONTAKTUS  
KÖZLEMÉNYEK  
2009. ÁPRILIS





sgi

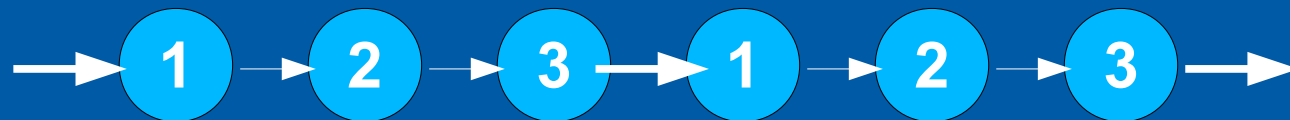
# Alkalmazás

⇒ Mire lehet használni?

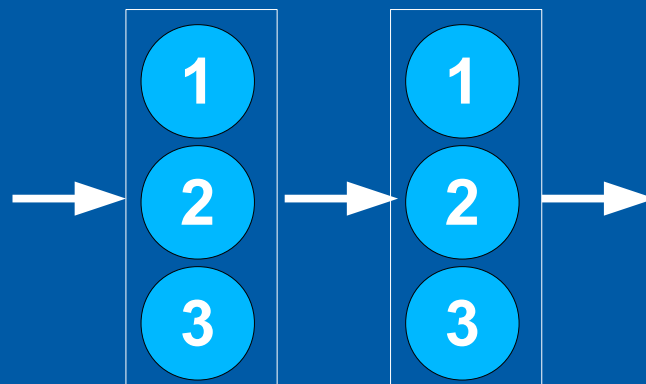
- Folyamat modellje

$$\begin{aligned} q_{i_1 i_2 \dots i_N u}^{k+1} &= q_{i_1 i_2 \dots i_N u}^k + \prod_{n=1}^N \mu_{i_n, n}(s_n) \cdot \mu_u(a) \cdot \Delta \tilde{Q}_{i, u}^{k+1} = \\ &= q_{i_1 i_2 \dots i_N u}^k + \prod_{n=1}^N \mu_{i_n, n}(s_n) \cdot \mu_u(a) \cdot \alpha_{i, u}^k \cdot \left( g_{i, u, j} + \gamma \cdot \max_{v \in U} \tilde{Q}_{j, v}^{k+1} - \tilde{Q}_{i, u}^k \right) \end{aligned}$$

- Szekvenciális probléma



- Párhuzamosítható probléma



# Alkalmazás

## ⇒ Deep Blue vs. Garry Kasparov

- IBM Deep Blue
- 1996/1997
- 11.38 GFLOPS
  - TOP500: 259.
- Sakkra specializált segédprocesszorok
- 1996: 4-2 GK – DB
- 1997: 3½–2½ DB – GK



# Alkalmazás

## ⇒ Filmgyártás / animáció

- Képkockák végleges renderelése
- Képkocka szinten párhuzamosítható
- Maguk a képkockák részei is egymástól függetlenül előállíthatóak
- Monsters University (2013):
  - 24 000 processzormagos sz.sz.





# Alkalmazás

⇒ SETI at home



⇒ Search for Extra-Terrestrial Intelligence  
„földön kívüli intelligencia keresése”

⇒ Elosztott rendszer

- „virtuális sz.sz.”

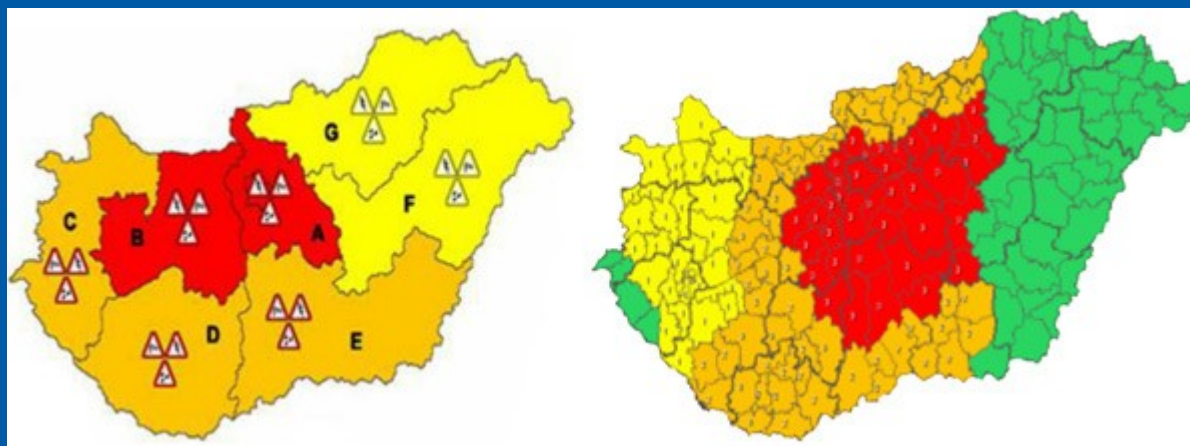
⇒ Bárki csatlakozhat

⇒ Rádióteleszkópok által vett jelek feldolgozása, mintázatok keresése



# Alkalmazás

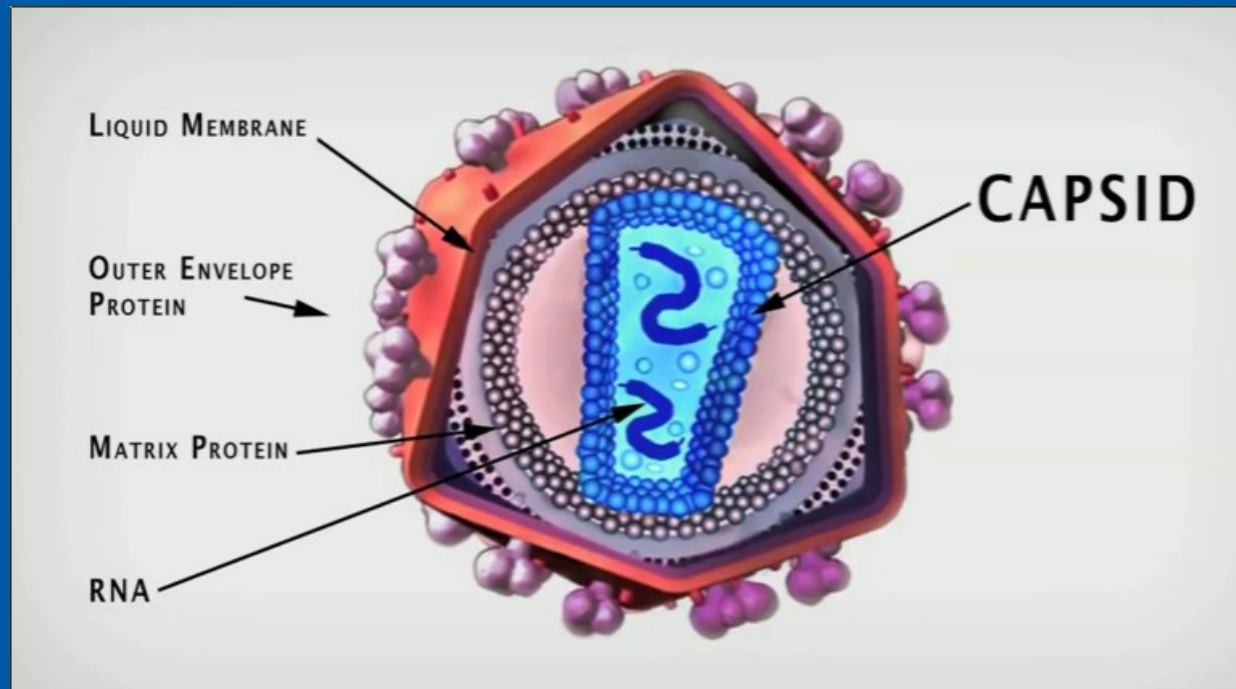
- ⇒ Időjárás előrejelzés
  - Adatok: műholdak, repülés, földi mérőállomások + Modell
- ⇒ OMSZ
- ⇒ 2010-ben új szuperszámítógép:
  - 7 régió helyett 174 kistérség
- ⇒ Riasztások kiadása
- ⇒ 2018.01.10. - új gép – HP Apollo 6000



# Alkalmazás

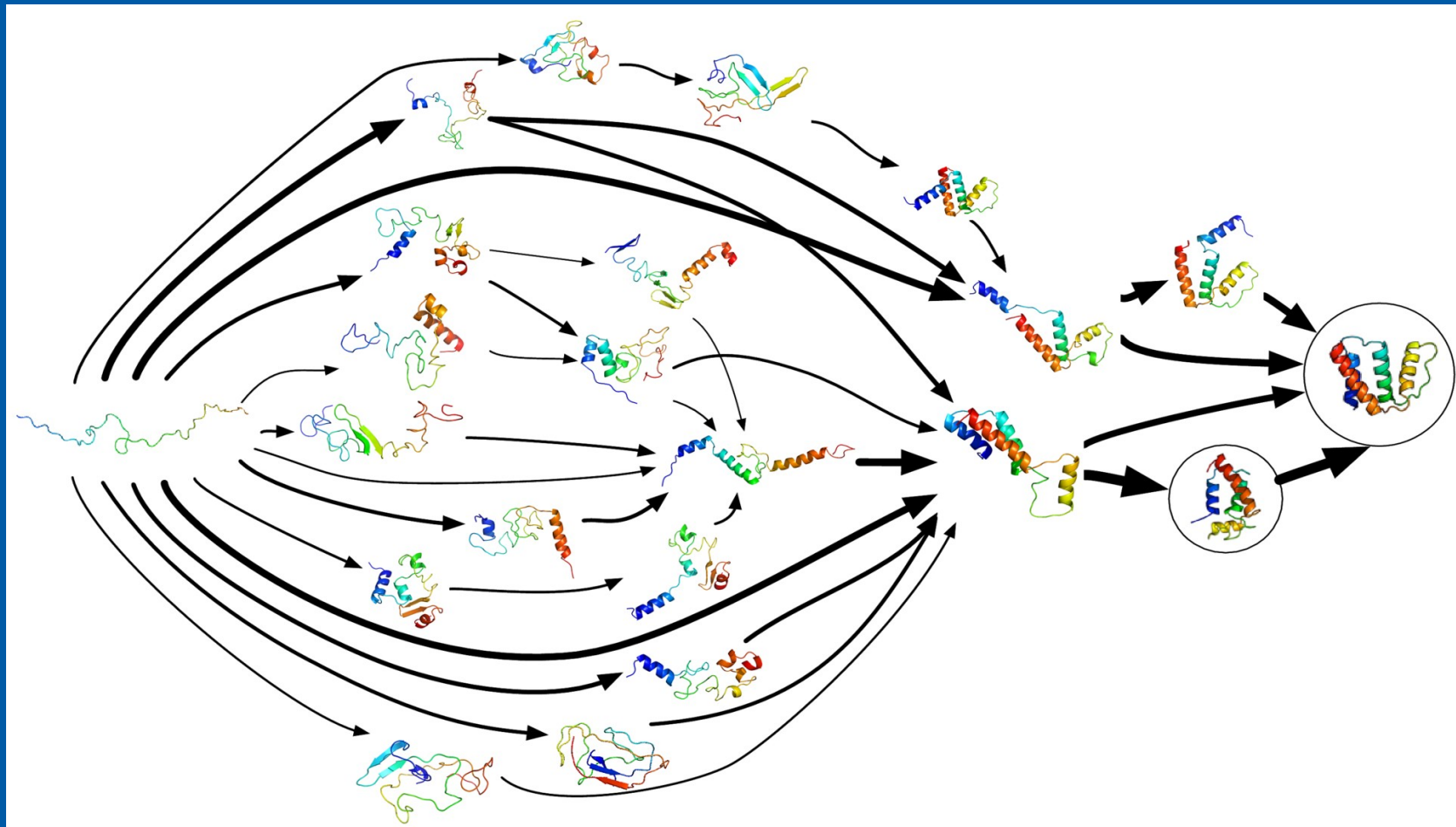
## ⇒ Betegségek kutatása

- Blue Waters sz.sz.
- 13,3 PFLOPS; 1,5 PB mem
- HIV-1 vírus fehérjeburkának feltérképezése
- 1300 különböző alakú protein
- 64 millió atom interakciójának szimulációja



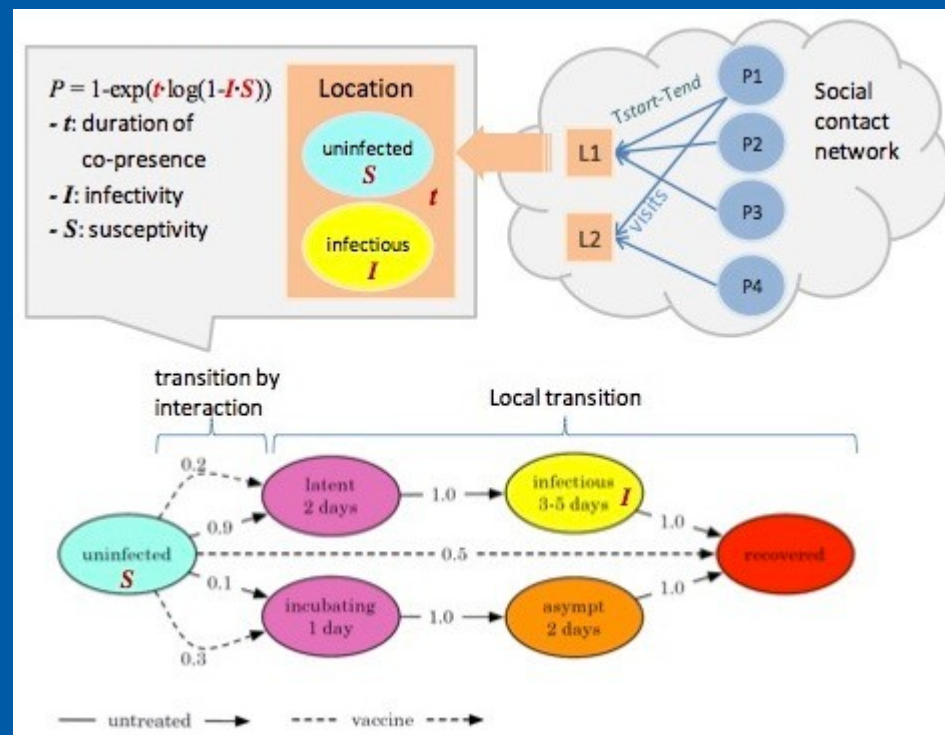
# Alkalmazás

- ➔ Protein folding
  - Folding at Home
  - kb. 40 PFLOPS, elosztott rendszer

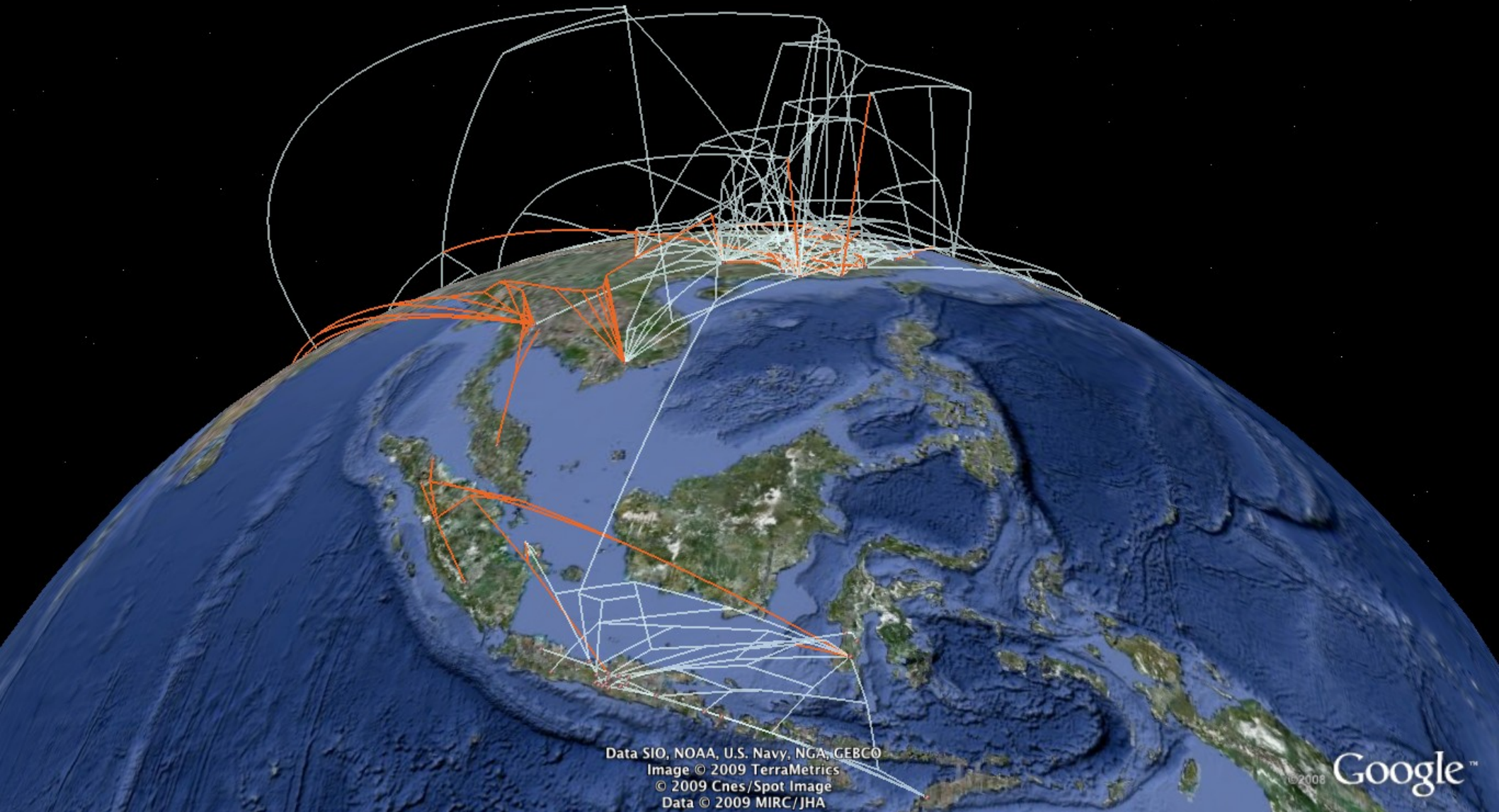


# Alkalmazás

- ⇒ Fertőzések terjedésének modellezése
  - pl. H1N1, H5N1
  - Begyűjtött adatok alapján előrejelzés készítése
  - „EpiSimdemics” programcsomag
    - 100 milliós nagyságrendű szociális kapcsolatrendszer kezelése

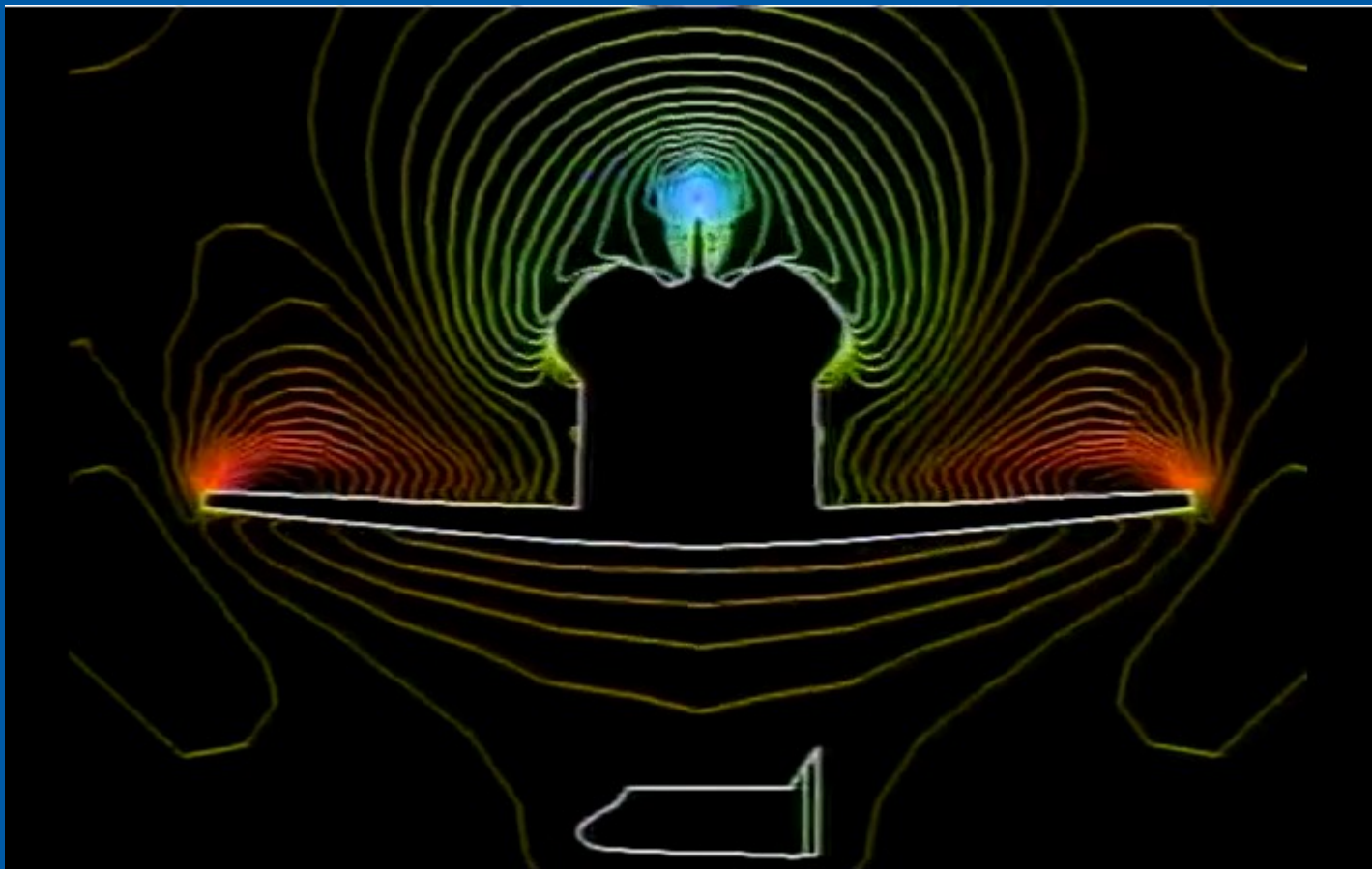


# Alkalmazás



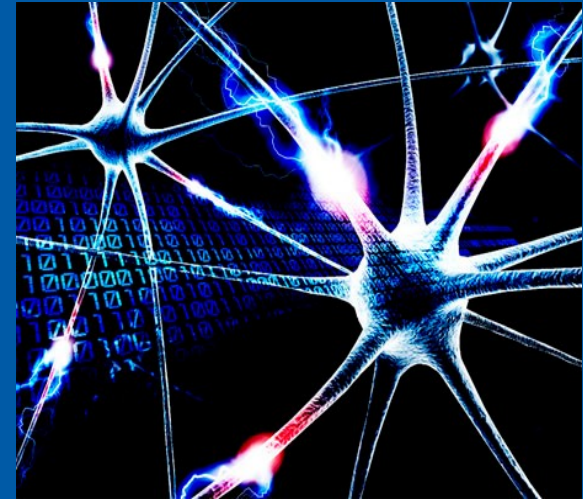
# Alkalmazás

- ⇒ NASA, 1986, CRAY-2, 1.9 GFLOPS
  - Áramlási szimulációk
  - Űrsikló



# Alkalmazás

- ⇒ Agy modellezés
- ⇒ „K Computer”, Japán, Kobe
- ⇒ 10.5 PFLOPS (TOP500: 4.)
  
- ⇒ Szimuláció:
  - 1.73 milliárd virtuális idegsejt
  - 10.4 billió virtuális szinapszis
  - Kb. 1% egy átlagos emberi agynak
  - 1 másodpercnyi aktivitást 40 percnyi számítást vesz igénybe

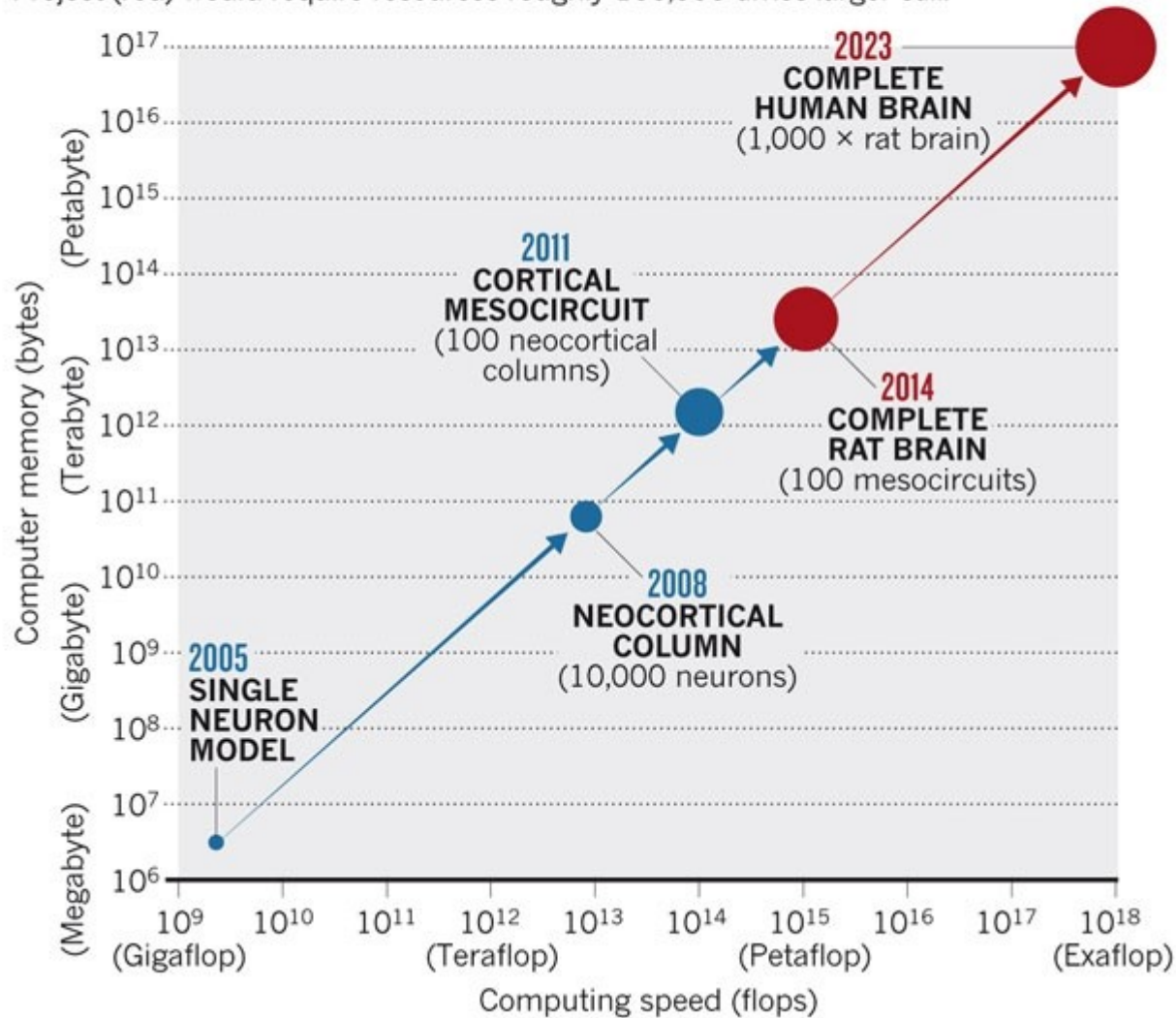




# Alkalmazás

## FAR TO GO

The Blue Brain Project has steadily increased the scale of its cortical simulations through the use of cutting-edge supercomputers and ever-increasing memory resources. But the full-scale simulation called for in the proposed Human Brain Project (red) would require resources roughly 100,000 times larger still.



# Alkalmazás

- ⇒ Kőolaj / földgáz lelőhelyek kutatása
- ⇒ Mérési adatok feldolgozása
  - mesterséges rezgés-keltés (hanghullám)
  - visszaverődések rögzítése (számos ponton)
- ⇒ Szeizmikus modell
- ⇒ Kőzetminták tanulmányozása
  - nagy felbontású scanner
  - áramlás

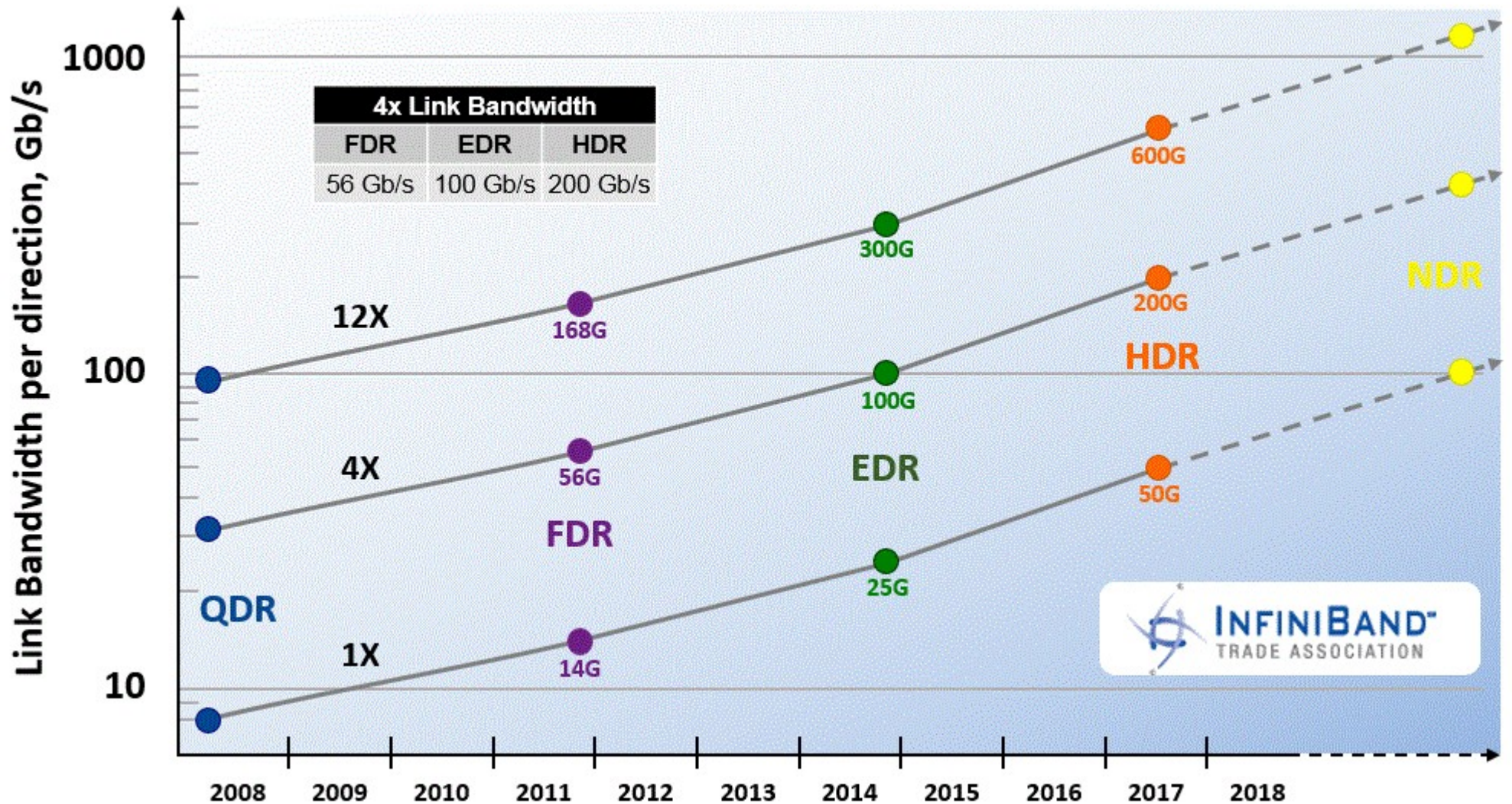


# InterConnect - InfiniBand

## ⇒ Interconnect - InfiniBand

- sok vállalat kooperációja révén
- nagyon elterjedt
- alacsony késleltetés
- alapvetően soros átvitel (de lehet több *lane*)
- nagy sáv szélesség
  - Eredetileg 2.5Gbit/s az alap sebesség (1db lane)
    - 4 pin (8b/10b kódolás)
  - InfiniBand 4X – 10Gbit/sec, 16 pin
  - InfiniBand 12X – 30Gbit/sec, 48 pin
  - Alapsebesség növekszik
  - Elterjedt manapság az 56Gbit/sec és a 100Gbit/s (4x 14Gbit / 25Gbit – 64b/66b kódolás)
  - De piacon van már a 200 (HDR) Gbit/sec
- 64k db eszköz címezhető

# InterConnect - InfiniBand



# InterConnect - InfiniBand

## ⇒ Interconnect - Infiniband

- switched-fabric topológia
- HCA (host channel adapter)
- TCA (target channel adapter) – I/O eszközökhöz
- Switch
  
- Multipath
- Multicast
  
- Médiák:
  - rézkábel (~17m, gyártófüggő)
  - üvegszál (km-ek)
  - PCB (NYÁK)

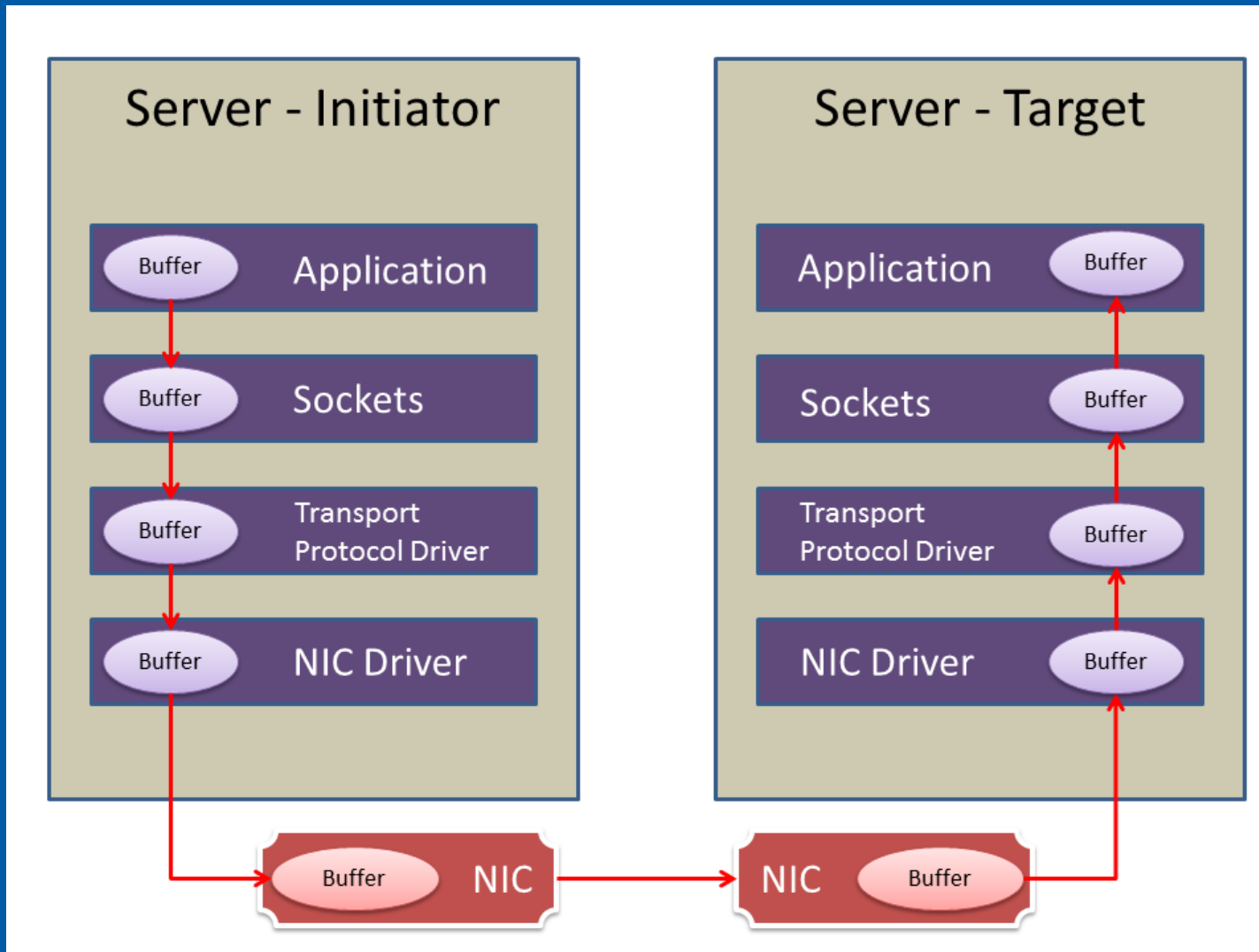
# InterConnect - InfiniBand

## ➔ InfiniBand

- **RDMA** (Remote DMA)
  - Távoli gép memóriájának közvetlen elérése egy processzből! (Driver + OS + VirtMem)
- QoS
- Maximum payload: 4k
- Rétegelt protokoll
  - Érdekesség, hogy IPv6-ot használ a network layer-ben
- Route-olható

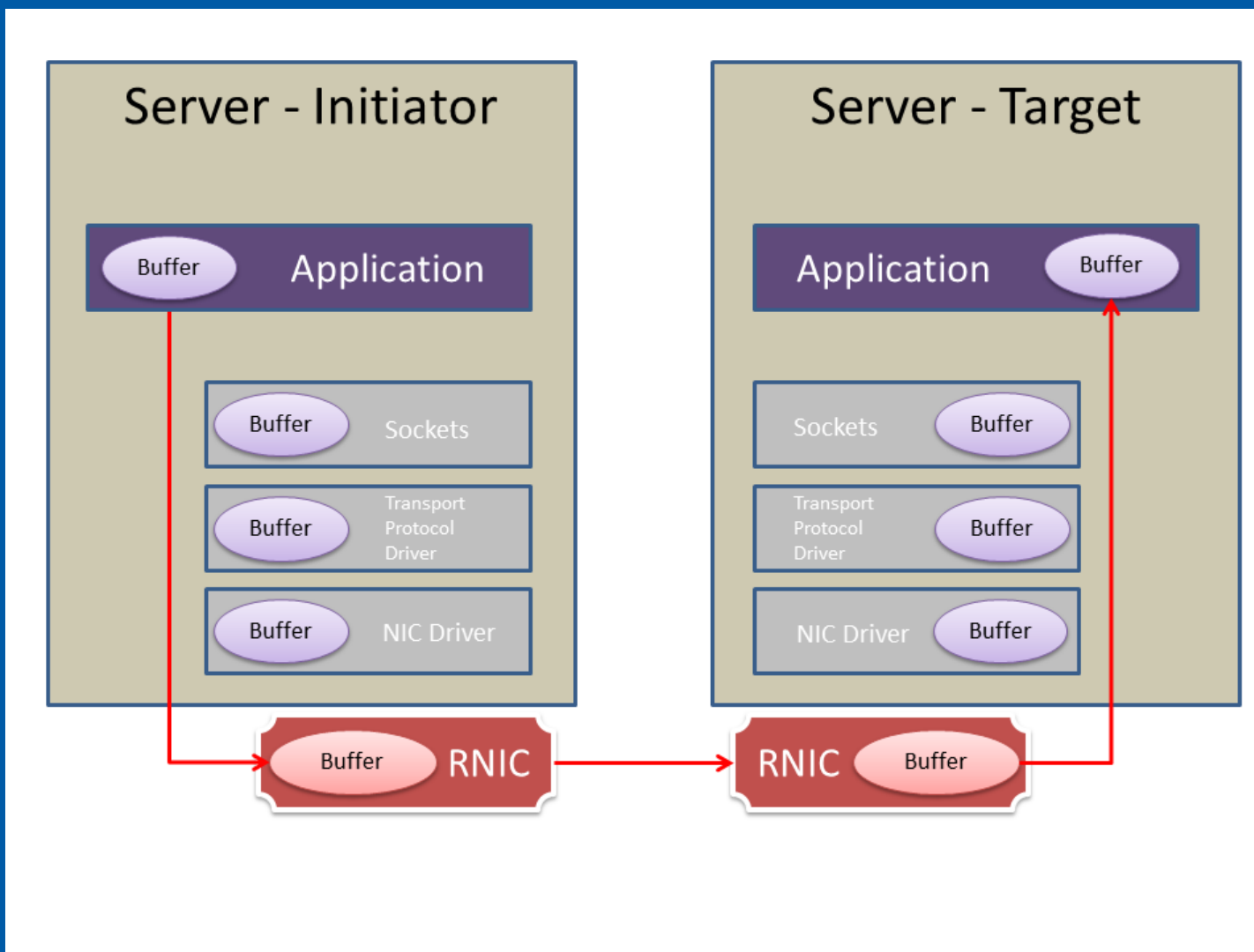
# InfiniBand - RDMA

- ➔ InfiniBand RDMA
  - „Hagyományos” hálózati átvitel

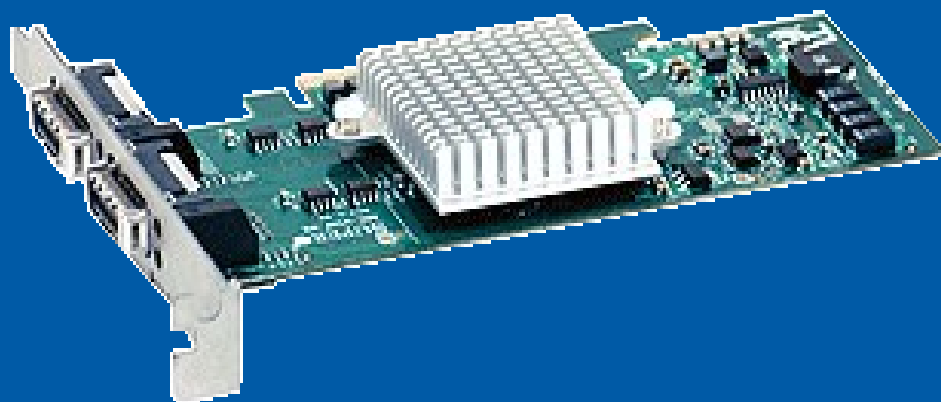
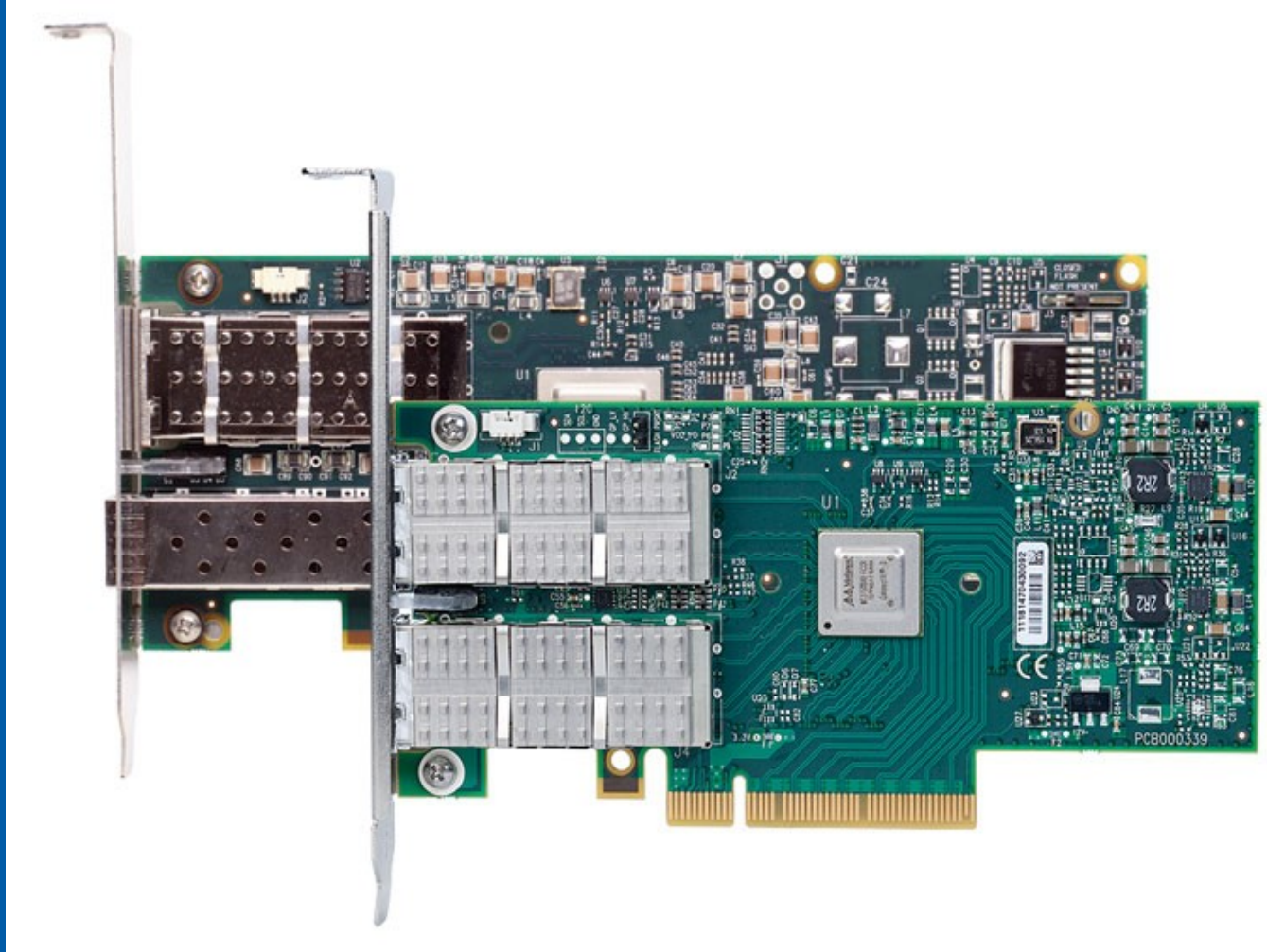


# InfiniBand - RDMA

- ⇒ InfiniBand RDMA
  - RDMA hálózati átvitel







# InfiniBand - Switch



Photo by Dr. Dávid Vincze

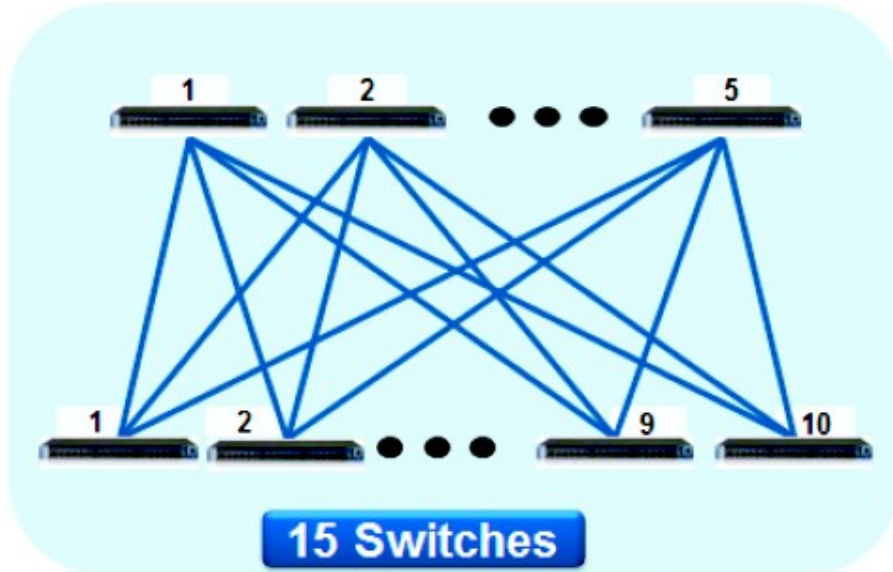
# InterConnect - Omni-Path

## ⇒ Interconnect – Egyebek

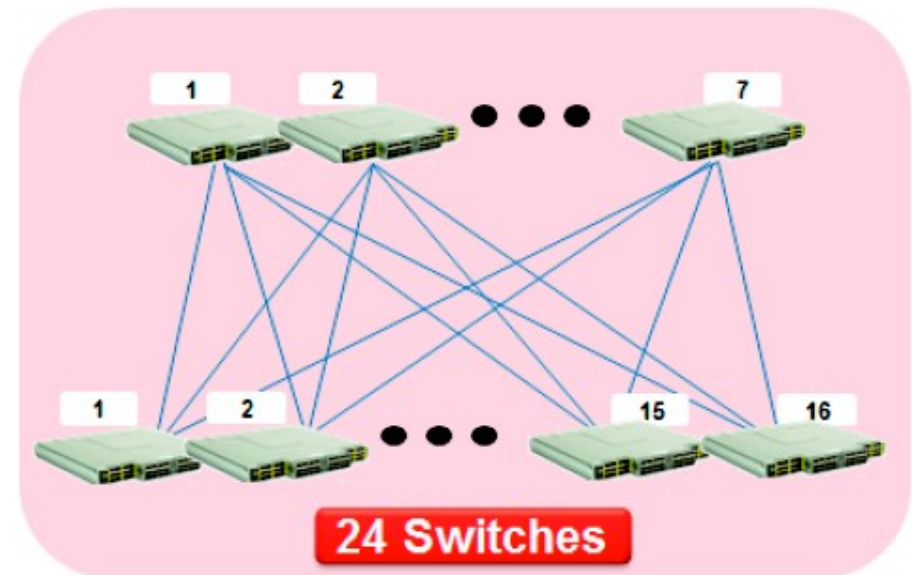
- Omni-Path (Intel)
  - 100Gbit/s, 64b/66b kódolás
  - 24-bit címzés
  - Retransmit-kor nem muszáj az egész csomagot újraküldeni
  - Intel: „HPC optimalizált”, „hatékonyabb”,
    - Kevesebb hely, kevesebb energia, olcsóbb, stb.
  - Mellanox: „not competitive”, nincs off-load, stb.
  - 48-port (OP) vs. 36-port (IB) switch
    - Kevesebb switch-re van szükség (mondja az Intel...)

# OmniPath vs. InfiniBand

400-Node 100G InfiniBand Platform



384-node 100G Omni-Path Platform



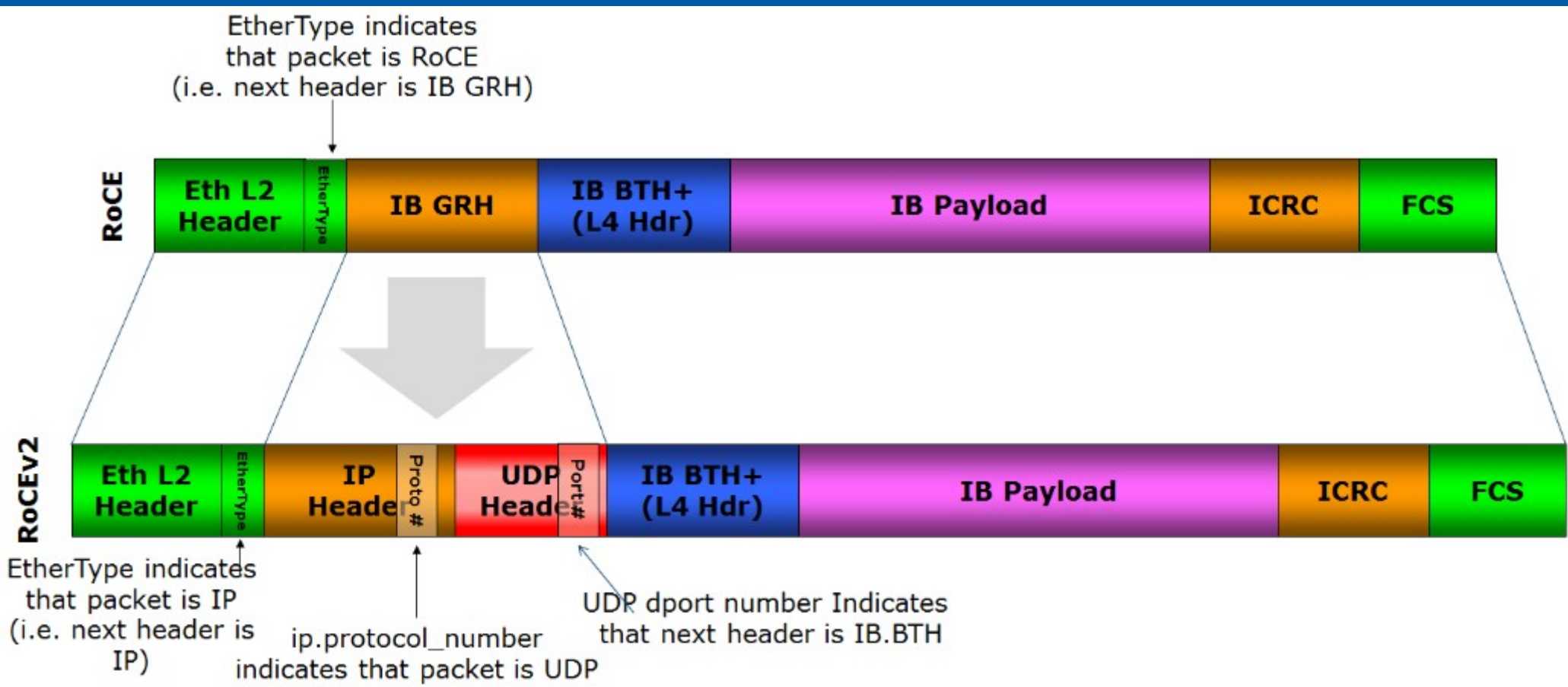
**Figure 2.** HDR100 Requires 1.6X Fewer Switches for 400-Nodes

[http://www.mellanox.com/related-docs/whitepapers/WP\\_Introducing\\_200G\\_HDR\\_InfiniBand\\_Solutions.pdf](http://www.mellanox.com/related-docs/whitepapers/WP_Introducing_200G_HDR_InfiniBand_Solutions.pdf)

# InterConnect - RoCE

## ➔ Interconnect – Egyebek

- RDMA over Converged Ethernet (RoCE)
  - Ethernet felett (pl. 10/25/50/100Gbit/s)
  - RoCEv2
    - IP+UDP
  - IB transport protocol



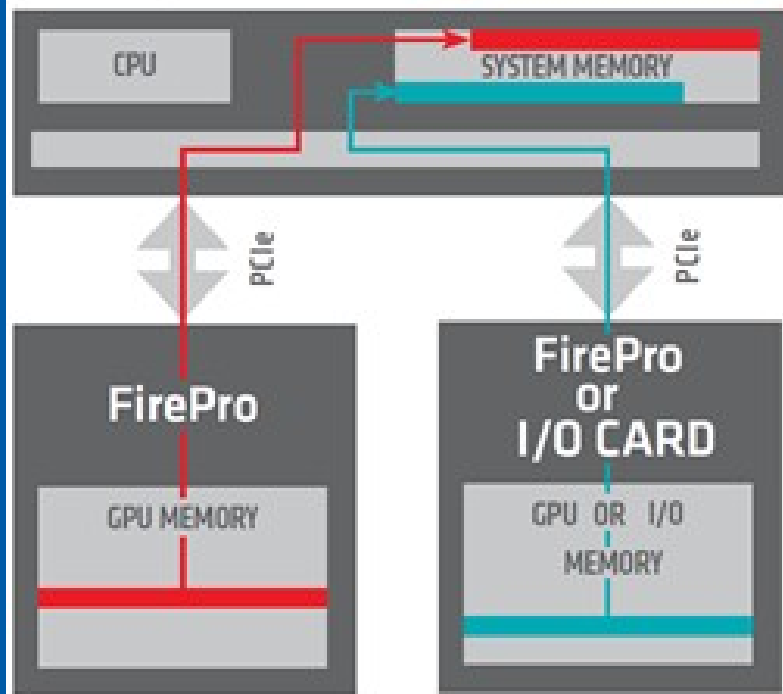
# InterConnect - GPU

## ⇒ Interconnect – Egyebek

- GPUDirect / DirectGMA
  - Nem a központi memóriát lehet RDMA-zni
  - Hanem perifériák memóriáját. pl. GPU
  - Persze maga a komm. mehet IB-n pl.

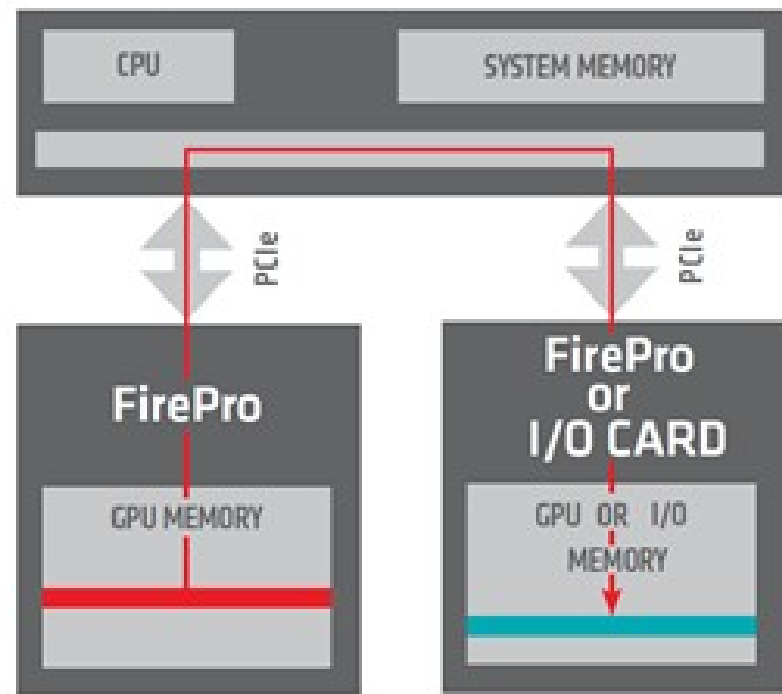
### Traditional two-step copy approach

Total transfer time: 36.4 ms



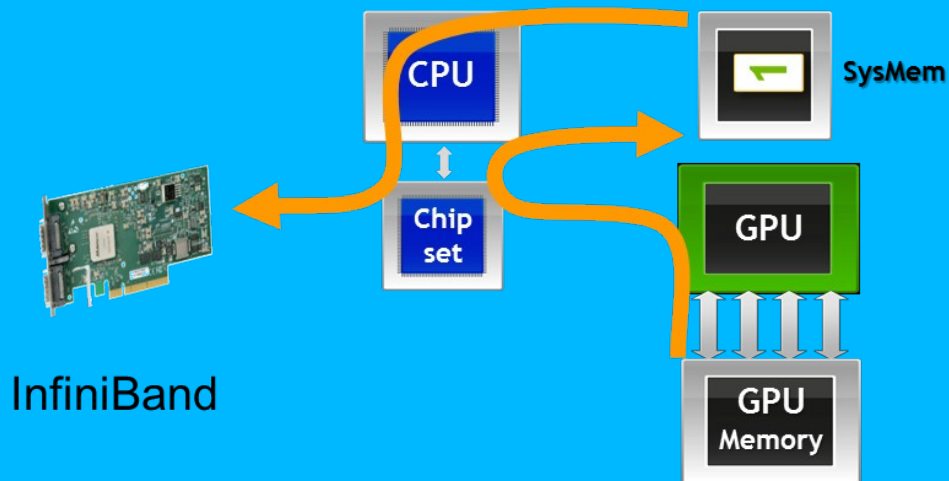
### Approach with AMD FirePro™ and DirectGMA

Total transfer time: 28 ms

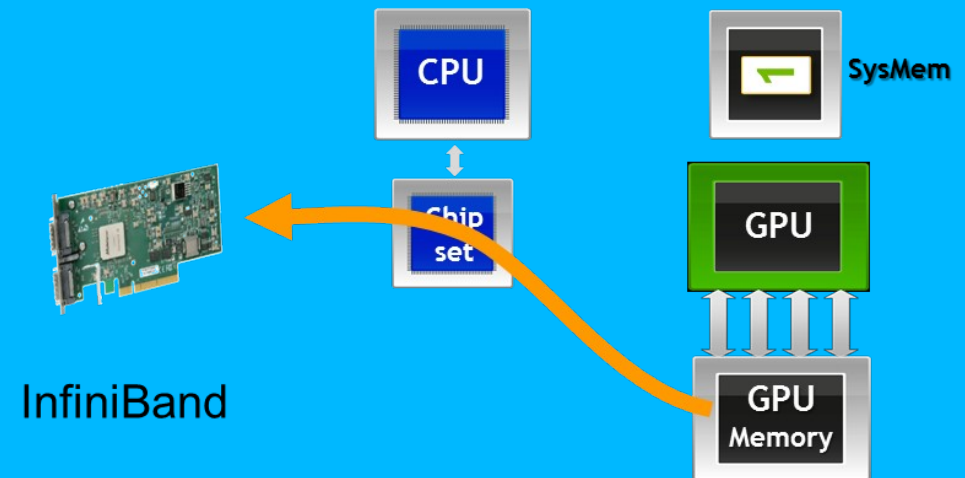


# InterConnect - GPU

*No GPUDirect RDMA*



*GPUDirect RDMA*



# Operációs Rendszerek MSc

## ➔ Klaszterek

- Több (akár heterogén) számítógép „egyesítése” adott feladatra
- HA-LB
  - High Availability – Load Balancing
- HA
  - Valamelyik erőforrás kiesik, maradjon működőképes
  - Kiszolgálók, eszközök többszörözése
  - Adatok folyamatos szinkronizációja
  - Master-Slave, Master-Master
  - Hogyan detektáljuk, hogy valami kiesett?
    - biztos, hogy nem „én” estem ki?
    - mi a helyzet, ha menet közben visszajön?
    - kinek friss az adata?
  - Fencing



# Operációs Rendszerek MSc

## ⇒ Klaszterek

### ● HA

- jellemző, hogy több hálózati címük van
  - Minden gépnek egy saját
  - Egy, amin a szolgáltatás fut
  - Az veszi fel a „szolgáltatás” címet, amelyik aktív

### ● LB

- A terhelést valamilyen szempont szerint szétosztani
- Kiszolgálók, eszközök többszörözése
- Adatok folyamatos szinkronizációja
- Master-Slave, Master-Master
  - read only node-ok
- **Hogyan detektáljuk, hogy valami kiesett?**
  - biztos, hogy nem „én” estem ki?
  - mi a helyzet, ha menet közben visszajön?

# Operációs Rendszerek MSc

## ⇒ Klaszterek

- HA és LB sokszor együtt
- Single Point of Failure (SPoF)
- Elosztás több szempont szerint
  - Alkalmazás
    - Egyik szerver Web, másik SQL
    - DNS Round Robin
    - stb.
  - Hálózati réteg
    - Linux Virtual Server
    - HW Load Balancer
    - stb.
  - OS
  - HW

# Operációs Rendszerek MSc

## ⇒ Single System Image (SSI)

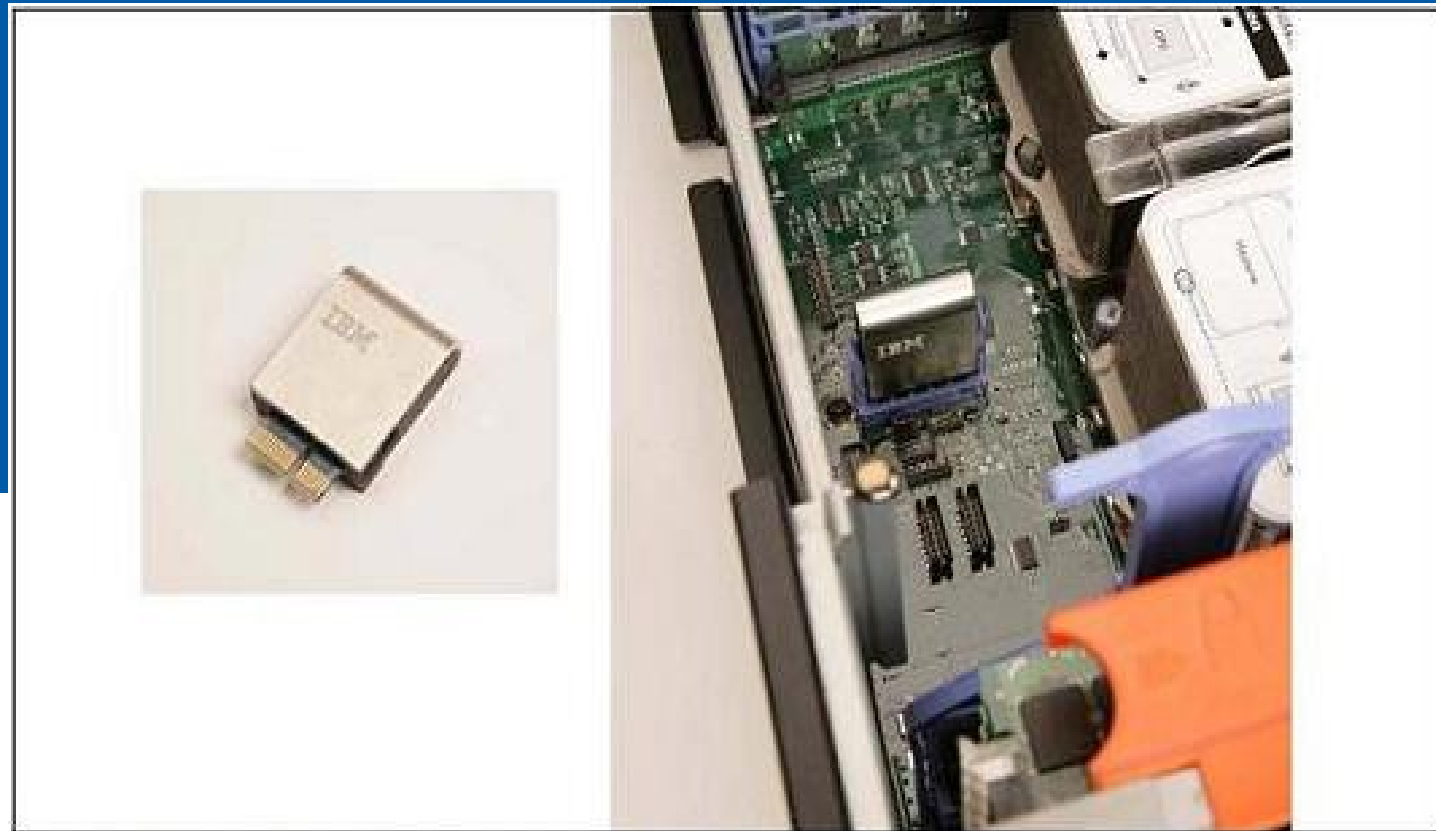
- Erőforrások egyesítése
- Több fizikai eszköz hardveres vagy szoftveres egyesítése egy logikai erőforrássá
- Magyarul: több gép egynek látszódik
- Processzek különböző node-okon
- Memóriaterület elosztva (NUMA!)
- Háttértár elosztva  
(jusson eszünkbe a RAID, konzisztencia, SAN...)
- Checkpoint
- Processz migráció
- Szint: HW, OS, middleware, alkalmazás

# Operációs Rendszerek MSc

## ⇒ HW

- Régen Digital Alpha rendszerek
  - OpenVMS galaxy
- Korábban pl. IBM x3850M2 / x3950M2
- 4 node fűzhető össze
  - saját interfész (ScaleXpander)
- Egy gép 256GB RAM, 4x6-core CPU
- Max 1TB RAM, 4\*4-way, 4\*4\*6 mag
  - Maga a gép is elég drága + ScaleXpander - \$3999
- SGI UV 2000/3000

# Operációs Rendszerek MSc



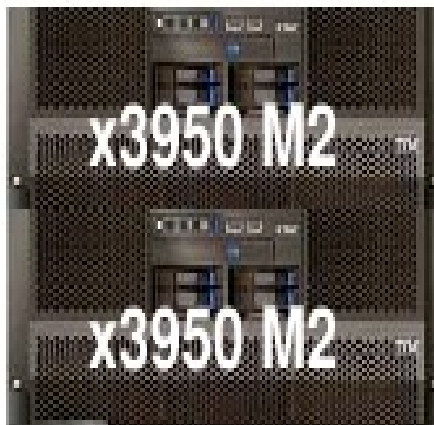
*Figure 2 The ScaleXpander chip installed in the x3850 M2 enables the server to scale*



**One node**  
**2-way or 4-way**  
Up to 256 GB RAM



**Two nodes**  
**4-way or 8-way**  
(Each node is  
2-way or 4-way)  
Up to 512 GB RAM



**Three nodes\***  
**6-way or 12-way**  
(Each node is  
2-way or 4-way)  
Up to 768 GB RAM



**Four nodes\***  
**8-way or 16-way**  
(Each node is  
2-way or 4-way)  
Up to 1 TB RAM



\* 3-node and 4-node configurations are planned to be supported in 2Q08