# VIP Briefing on GPU Accelerator Technology

Steve Scott, CTO of Tesla
Ian Buck, GM of GPU Computing
Dr. Dirk Pleiter, Juelich

# Long Term Goals for Tesla
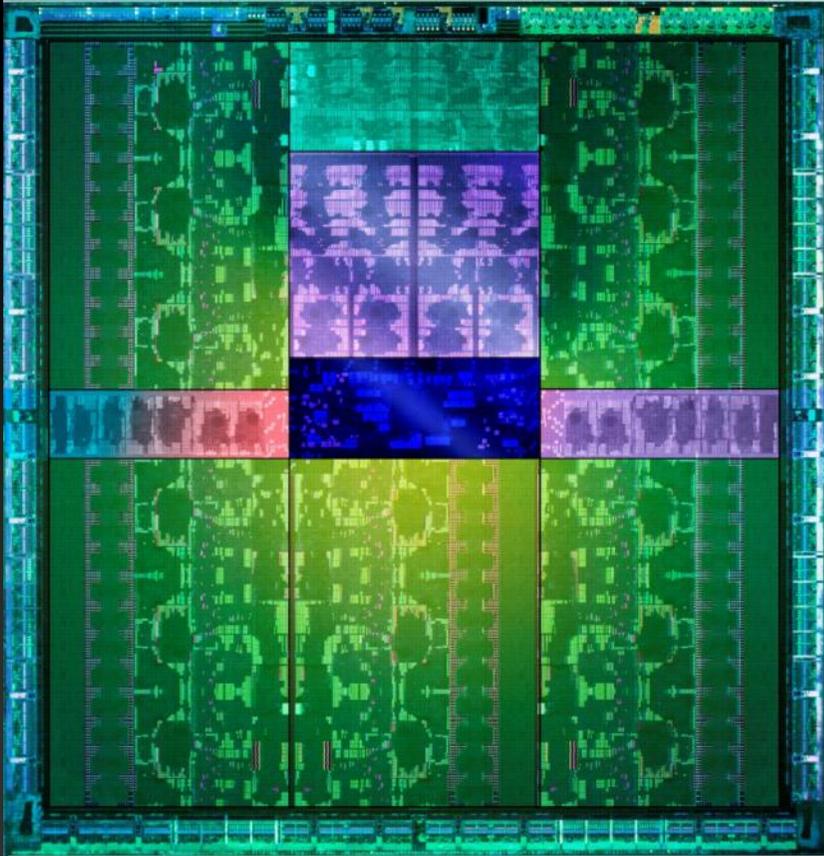


Power
Efficiency

Ease of
Programming
And Portability

Application
Space
Coverage

# KEPLER
## THE WORLD'S FASTEST, MOST EFFICIENT HPC ACCELERATOR

SMX                            *(power efficiency)*

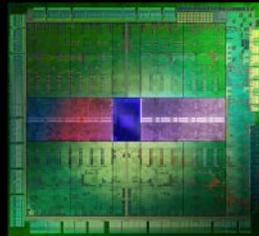Hyper-Q                        *(programmability and application coverage)*

Dynamic Parallelism

# Tesla K10

Dual GK104 GPUs

3x Single Precision
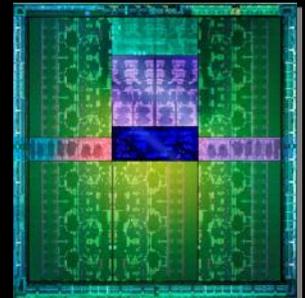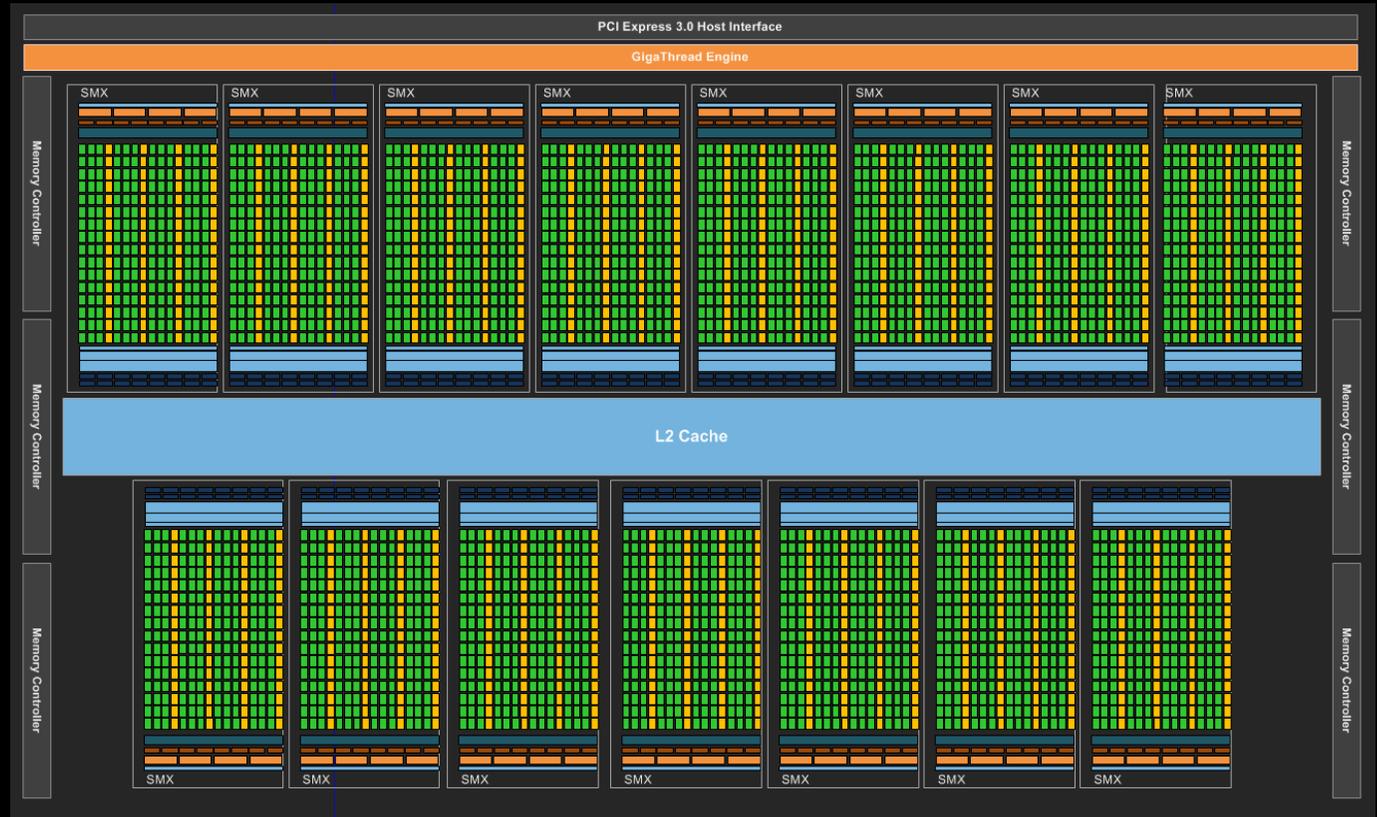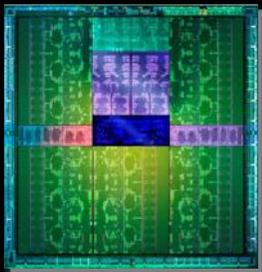Video, Signal, Life Sciences, Seismic

Available Now

# Tesla K20

GK110 GPU

3x Double Precision
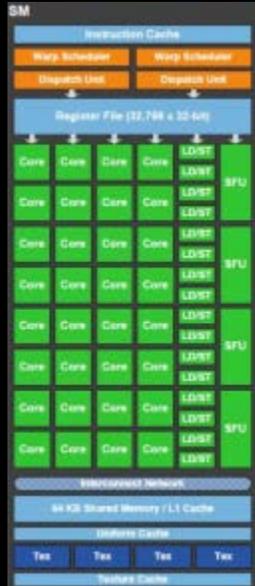Hyper-Q & Dynamic Parallelism
CFD, FEA, Finance, Physics, etc.

Available Q4 2012

4

# Kepler GK110 Block Diagram

- **7.1B Transistors**
- **15 SMX units**
- **> 1 TFLOP FP64**
- **1.5 MB L2 Cache**
- **384-bit GDDR5**
  - **~250 GB/s**
- **PCI Express Gen3**



PCI Express 3.0 Host Interface

GigaThread Engine

SMX SMX SMX SMX SMX SMX SMX SMX

Memory Controller

L2 Cache

SMX SMX SMX SMX SMX SMX SMX

# Kepler GK110 SMX vs Fermi SM
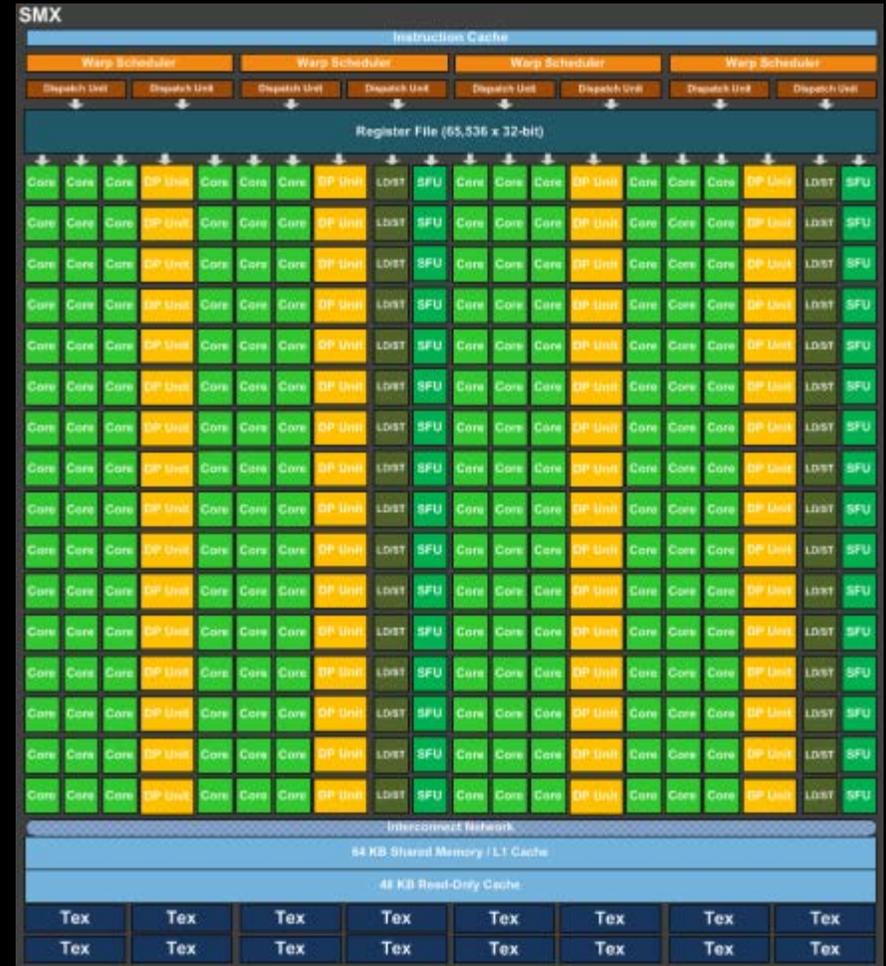


**3x sustained perf/W**

**Ground up redesign for perf/W**
6x the SP FP units
4x the DP FP units
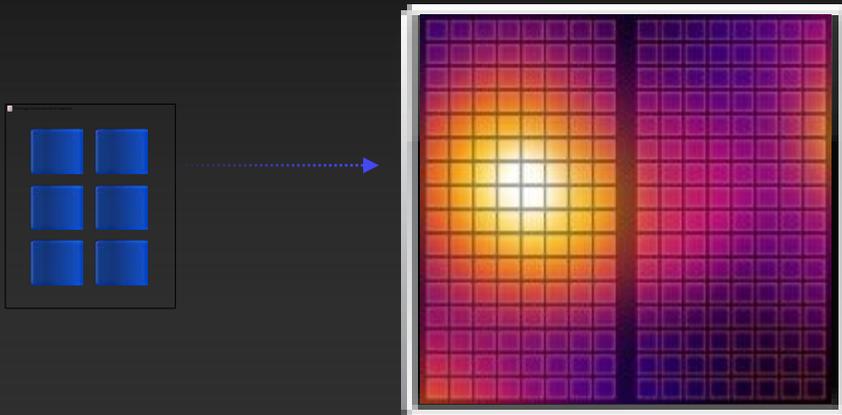Significantly slower FU clocks

# Selected Kepler ISA Enhancements

- **Larger number of registers per thread**
  - **63 in Fermi → 255 in Kepler**
  - **Common performance limited in Fermi due to register spilling**
  - **Significant performance improvement for some codes (e.g.: 5.3x on Quda QCD!)**
- **Atomic operations**
  - **Added int64 to match int32**
  - **Added functional units → 2-10x performance gains**
- **SHFL instruction for data exchange amongst threads of a warp**
  - **Broadcast, shifts, butterflies**
  - **Useful for sorts, reductions, etc.**
- **Loads through texture memory**
  - **Higher bandwidth and flexibility for read-only data  (const__restrict)**
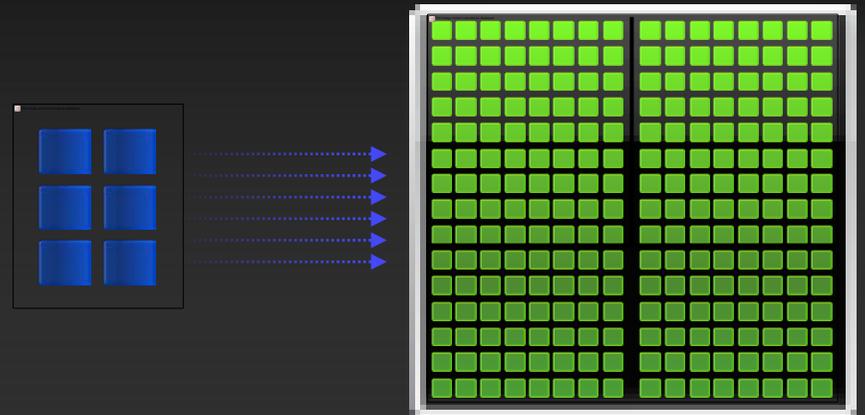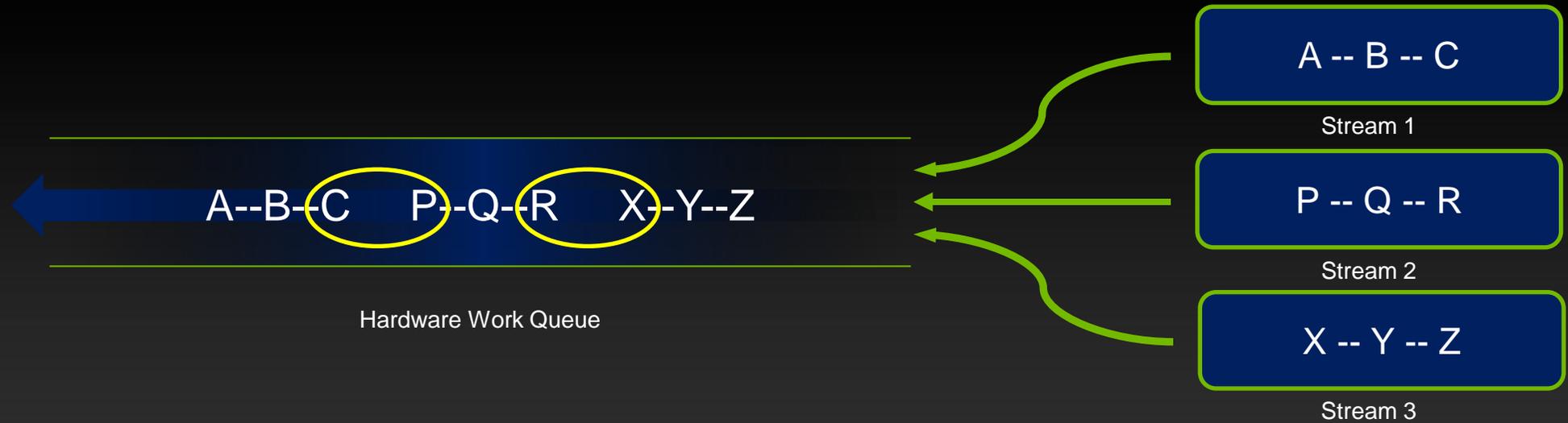
# Hyper-Q



## FERMI
### 1 Work Queue

## KEPLER
### 32 Concurrent Work Queues

# Fermi Concurrency

A -- B -- C
Stream 1

P -- Q -- R
Stream 2

X -- Y -- Z
Stream 3

A--B--C   P--Q--R   X--Y--Z

Hardware Work Queue
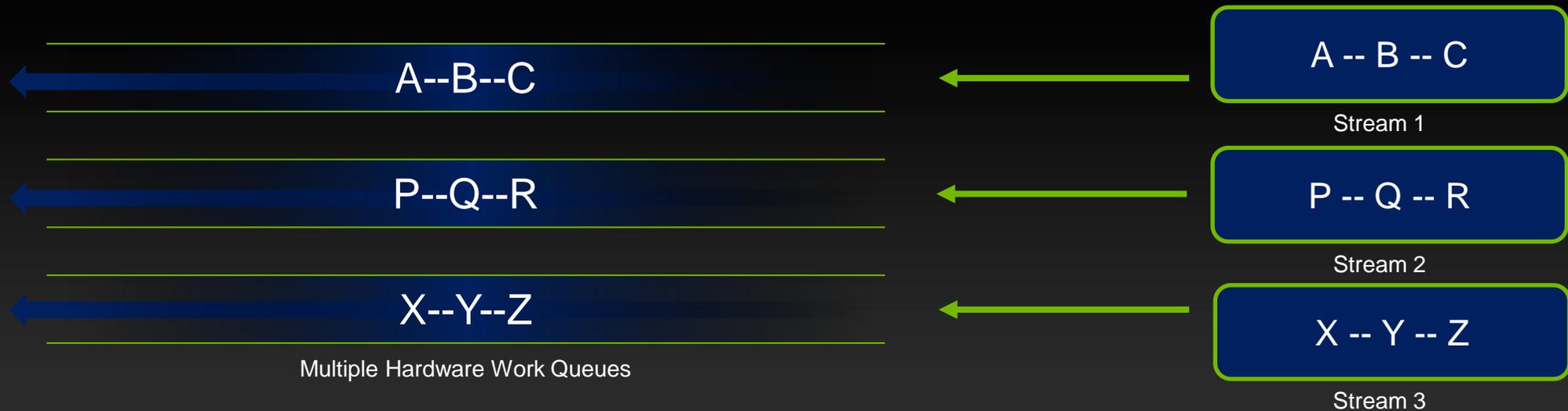
## Fermi allows 16-way concurrency

- Up to 16 grids can run at once
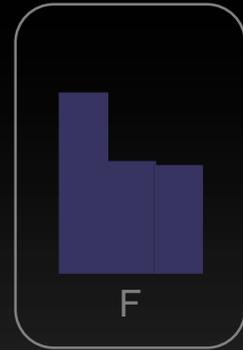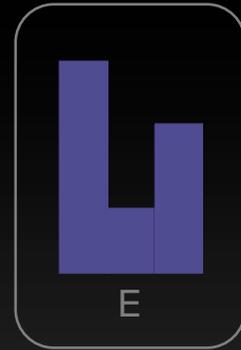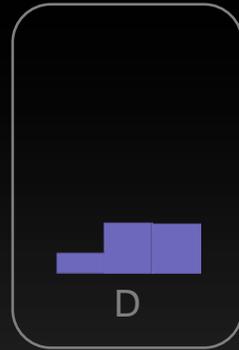- But CUDA streams multiplex into a single queue
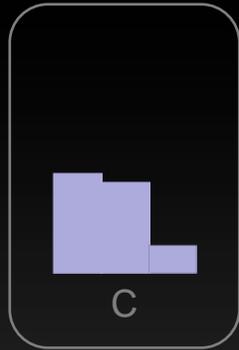- Overlap only at stream edges

# Kepler Improved Concurrency

A--B--C

A -- B -- C

Stream 1

P--Q--R
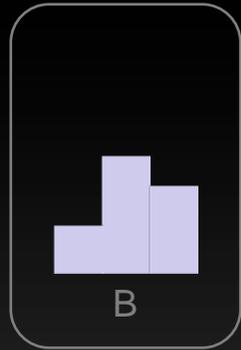
P -- Q -- R

Stream 2

X--Y--Z

X -- Y -- Z

Stream 3

Multiple Hardware Work Queues
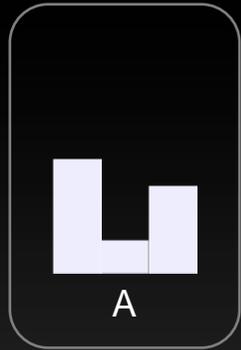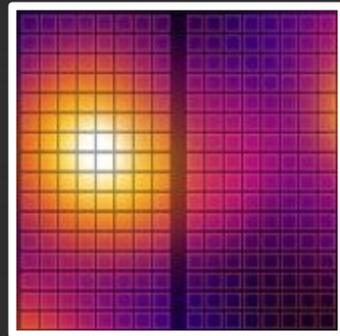
## Kepler allows 32-way concurrency

- One work queue per stream
- Concurrency at full-stream level
- No inter-stream dependencies

# Fermi: Time-Division Multiprocess



CPU Processes

Shared GPU

# Fermi: Time-Division Multiprocess



CPU Processes

Shared GPU

# Fermi: Time-Division Multiprocess



CPU Processes
Shared GPU

# Fermi: Time-Division Multiprocess



CPU Processes
Shared GPU

# Fermi: Time-Division Multiprocess



A    B    C    D    E    F
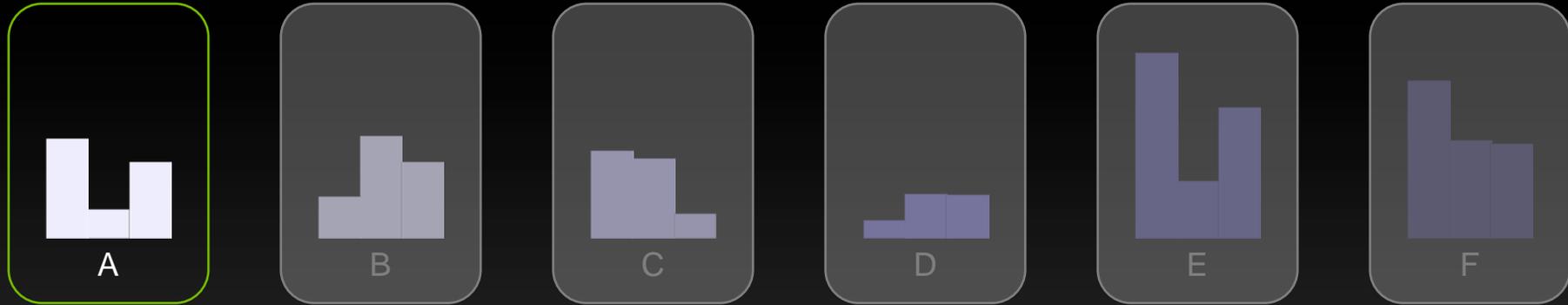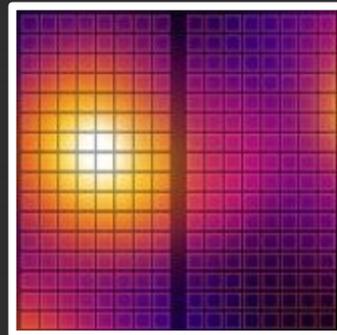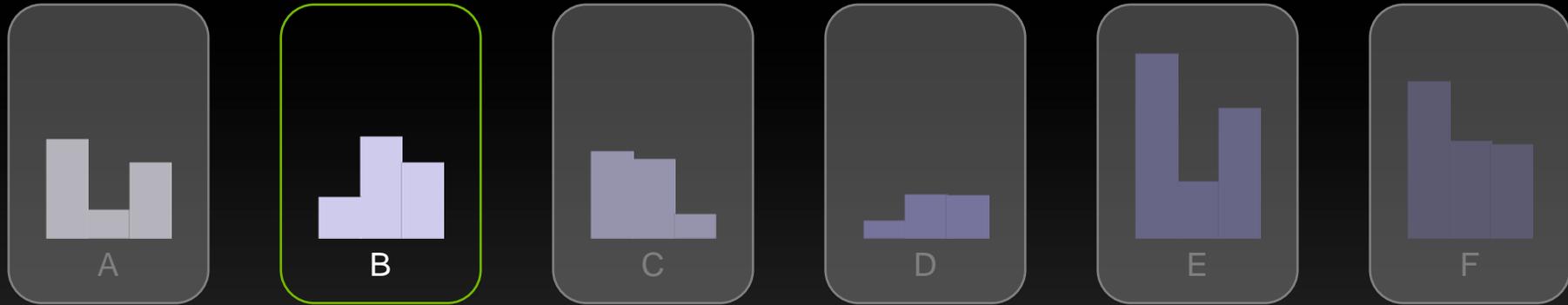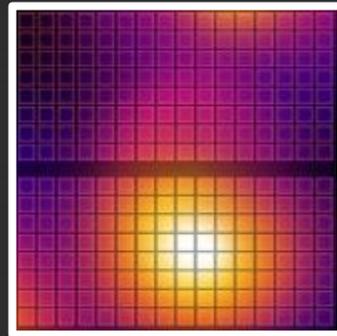
CPU Processes

Shared GPU

# Fermi: Time-Division Multiprocess



CPU Processes

Shared GPU

# Fermi: Time-Division Multiprocess



CPU Processes
Shared GPU

# Kepler Hyper-Q: Simultaneous Multiprocess



CPU Processes

Shared GPU

# With Hyper-Q



Easier threaded parallelism

Multi-rank MPI parallelism

Better strong scaling

# Dynamic Parallelism

The ability for any GPU thread to launch a parallel GPU kernel

- Dynamically
- Simultaneously
- Independently



**Fermi: Only CPU can generate GPU work**

**Kepler: GPU can generate work for itself**

# Dynamic Work Generation

Coarse grid — Higher Performance Lower Accuracy

Fine grid — Higher Accuracy Lower Performance

Dynamic grid — Target performance where accuracy is required

ears:        3898.34
DIES:        280875
DIES/SEC:7310286

# Familiar Syntax and Programming Model

```
int main() {
    float *data;
    setup(data);

    A <<< ... >>> (data);
    B <<< ... >>> (data);
    C <<< ... >>> (data);

    cudaDeviceSynchronize();
    return 0;
}
```

```
__global__ void B(float *data)
{
    do_stuff(data);

    X <<< ... >>> (data);
    Y <<< ... >>> (data);
    Z <<< ... >>> (data);
    cudaDeviceSynchronize();

    do_more_stuff(data);
}
```

# Simpler Code: LU Example

## LU decomposition (Fermi)

```
dgetrf(N, N) {
   for j=1 to N
      for i=1 to 64
         idamax<<<>>>
         memcpy
         dswap<<<>>>
         memcpy
         dscal<<<>>>
         dger<<<>>>
      next i

      memcpy
      dlaswap<<<>>>
      dtrsm<<<>>>
      dgemm<<<>>>
   next j
}
```

```
idamax();
```
```
dswap();
```
```
dscal();
```
```
dger();
```
```
dlaswap();
```
```
dtrsm();
```
```
dgemm();
```

**CPU Code**                          GPU Code

## LU decomposition (Kepler)

```
dgetrf(N, N) {
   dgetrf<<<>>>
```

CPU is Free

```
   synchronize();
}
```

```
dgetrf(N, N) {
   for j=1 to N
      for i=1 to 64
         idamax<<<>>>
         dswap<<<>>>
         dscal<<<>>>
         dger<<<>>>
      next i
      dlaswap<<<>>>
      dtrsm<<<>>>
      dgemm<<<>>>
   next j
}
```

**CPU Code**                          GPU Code

# CUDA By the Numbers:

| | |
|---|---|
| **>375,000,000** | CUDA-Capable GPUs |
| **>1,000,000** | Toolkit Downloads |
| **>120,000** | Active Developers |
| **>500** | Universities Teaching CUDA |

# CUDA 5

## Nsight™ for Linux & Mac

## NVIDIA GPUDirect™

## Library Object Linking

**Preview Release
Now Available**

# NVIDIA® Nsight™ Eclipse Edition



**CUDA-Aware Editor**

- Automated CPU to GPU code refactoring
- Semantic highlighting of CUDA code
- Integrated code samples & docs

**Nsight Debugger**

- Simultaneously debug of CPU and GPU
- Inspect variables across CUDA threads
- Use breakpoints & single-step debugging

**Nsight Profiler**

- Quickly identifies performance issues
- Integrated expert system
- Automated analysis
- Source line correlation

# Available for Linux and Mac OS

# Kepler Enables Full NVIDIA GPUDirect™



Server 1

Server 2

# GPU Computing with LLVM

**Developers want to build front-ends for**
**Java, Python, R, DSLs**

**Target other processors like**
**ARM, FPGA, GPUs, x86**

```
┌─────────────────┐   ┌─────────────────┐
│ CUDA            │   │ New Language    │
│ C, C++, Fortran │   │ Support         │
└────────┬────────┘   └────────┬────────┘
         │                     │
         ▼                     │
┌──────────────────────────────┐
│       LLVM Compiler          │
└──┬──────────┬────────────┬───┘
   ▼          ▼            ▼
┌───────┐ ┌───────┐ ┌──────────────┐
│NVIDIA │ │ x86   │ │New Processor │
│GPUs   │ │ CPUs  │ │Support       │
└───────┘ └───────┘ └──────────────┘
```

LLVM
COMPILER
INFRASTRUCTURE

# OpenACC Directives

**CPU**

**GPU**

```
Program myscience
    ... serial code ...
!$acc kernels
    do k = 1,n1
        do i = 1,n2
            ... parallel code ...
        enddo
    enddo
!$acc end kernels
    ...
End Program myscience
```

OpenACC
Compiler
Hint

Your original
Fortran or C code

Simple Compiler hints

Compiler Parallelizes code

Portability, Productivity,
Performance

# Performance: Leveraging GPU

## Reading DNA nucleotide sequences

*Shanghai JiaoTong University*

4 directives

16x faster

## Designing circuits for quantum computing

*UIST, Macedonia*

1 week

40x faster

## Extracting image features in real-time

*Aselsan*

3 directives

4.1x faster

## HydroC- Galaxy Formation

*PRACE Benchmark Code, CAPS*

1 week

3x faster

## Real-time Derivative Valuation

*Opel Blue, Ltd*

Few hours

70x faster

## Matrix Matrix Multiply

*Independent Research Scientist*

4 directives

6.4x faster

# Enabling ARM Ecosystem: CARMA DevKit
## CUDA on ARM

CUDA GPU     Tegra ARM CPU



Tegra 3 Quad-core ARM A9
Quadro 1000M (96 CUDA cores)
Ubuntu

Gigabit Ethernet
SATA Connector
HDMI, DisplayPort, USB

# The Day Job That Makes It All Possible...

- Leverage volume graphics market to serve HPC
  - HPC needs outstrip HPC market's ability to fund the development
  - Computational graphics and compute are *highly* aligned



Tegra

GeForce

Quadro

# Jülich-NVIDIA Application Lab

19. June 2012 | Dirk Pleiter (JSC)

# Supercomputing at Forschungszentrum Jülich

Role of the **Jülich Supercomputing Centre (JSC)**:

- **Operation** of supercomputers for local, national and European scientists

- **User support** including support of research communities by means of simulation laboratories

- **R&D** on future IT technologies, algorithms, tools, GRID, etc.
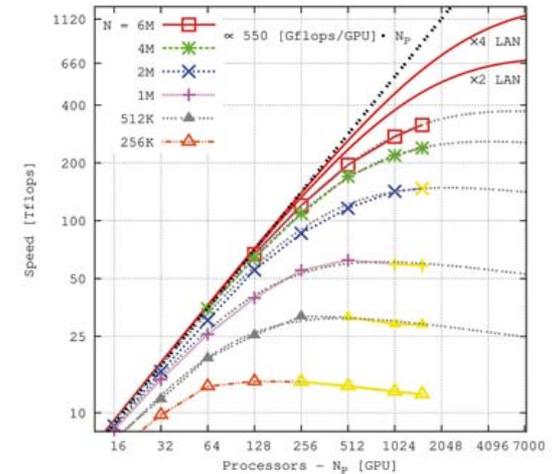
- **Education** and training of users



JUQUEEN

JUROPA

JUPACE

# Our view on GPU computing

- Performance acceleration for a significant set of relevant scientific applications

- **JUDGE** = Jülich Dedicated GPU Environment

  - 206 node IBM iDataPlex cluster

  - Dual-CPU, dual-GPU nodes

  - About 240 TFlops (peak)

  - Partitions dedicated to astrophysics and brain research

- Large potential for **energy efficient computing**

  - JUDGE is #14 on Green500 (Nov. 2011)

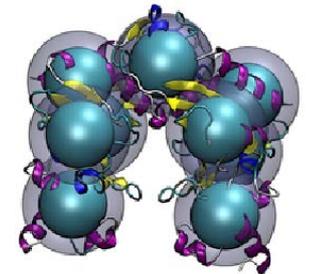  - Need for efficient utilisation of all computing devices

# Jülich-NVIDIA Application Lab
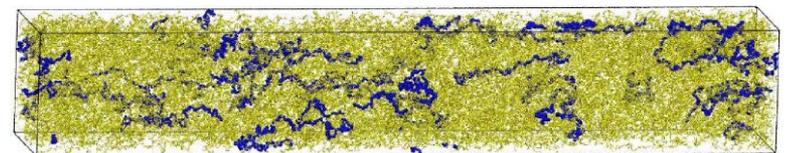


[R. Spurzem et al., 2012]

- Lab hosted at JSC
- **Mission statement**
  - Enable scientific applications for GPU-based architectures
  - Provide support for optimization
  - Investigate performance and scaling
- **Targeted research areas**
  - Astrophysics and astronomy
  - Computational medicine and neuroscience
  - Elementary particle physics
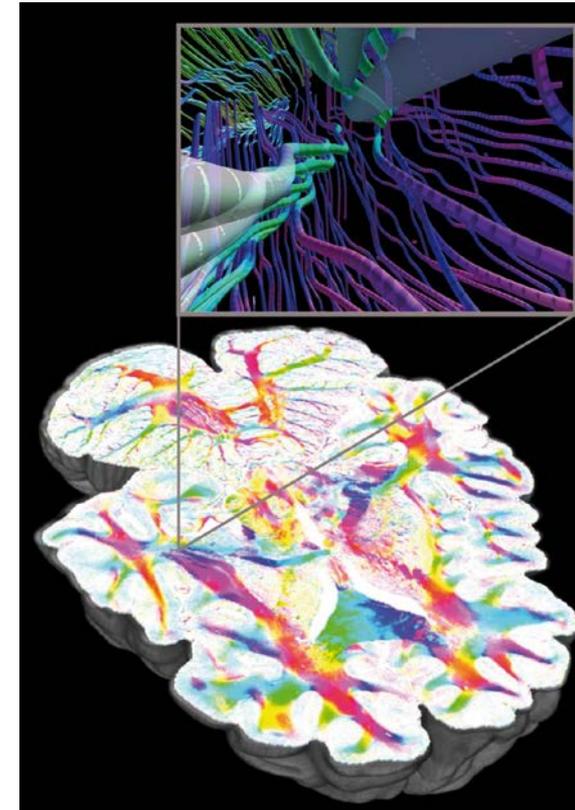  - Material science
  - Protein folding



[O. Zimmerrmann, 2011]



[G. Sutmann et al., 2011]

# Pilot application: JuBrain

- The **Jülich Brain Model** will display selected aspects of the brain's structural organization such as cortical areas and fiber tracts

  – Improve understanding of fiber operation

  – Help treating neurological disease

- Procedure

  – Preparation of brain sections

  – Image processing

  – 3D reconstruction and fiber tractography

- Already today significant speed-up using GPUs



[M. Axer et al., 2012]

# Questions?