

# Parallelization

Scientific Computing Winter 2016/2017

Part IV

With material by W. Gropp (<http://wgropp.cs.illinois.edu>) and J. Burkardt (<https://people.sc.fsu.edu/~jburkardt>)

Jürgen Fuhrmann

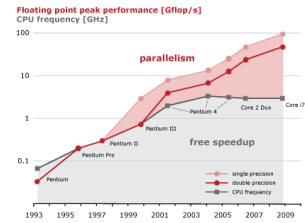
[juergen.fuhrmann@wias-berlin.de](mailto:juergen.fuhrmann@wias-berlin.de)

made wit pandoc



## Why parallelization ?

- ▶ Computers became faster and faster without that...



[Source: spiralgen.com]

- ▶ But: clock rate of processors limited due to physical limits
- ▶ ⇒ parallelization is the main road to increase the amount of data processed
- ▶ Parallel systems nowadays ubiquitous: even laptops and smartphones have multicore processors
- ▶ Amount of accessible memory per processor is limited ⇒ systems with large memory can be created based on parallel processors

## TOP 500 2016 rank 1-6

Based on linpack benchmark: solution of dense linear system. Typical desktop computer:  $R_{max} \approx 100 \dots 1000 \text{ GFlop/s}$

Rank	Site	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power (kW)
1	National Supercomputing Center in Wuxi, China	Sunway TaihuLight - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway NRCFC	10,449,600	93,014.6	125,435.9	15,371
2	National Super Computer Center in Guangzhou, China	Tianhe-2 (MilkyWay-2) - TH-1WB-FEP Cluster, Intel Xeon E5-2692 12C 2.200GHz, TH Express-2, Intel Xeon Phi 3151P, NUDT	3,120,000	33,862.7	54,922.4	17,808
3	ODE/NSA/Oak Ridge National Laboratory, United States	Titan - Cray XK7, Opteron 6274 16C 2.200GHz, Cray Gemini interconnect, NVIDIA K80x, Cray Inc.	540,640	17,590.0	27,112.5	8,209
4	ODE/NSA/LLNL, United States	Sequoia - BlueGene/Q, Power BGC 16C 1.60 GHz, Custom, IBM	1,572,864	17,173.2	20,127.7	7,890
5	RIKEN Advanced Institute for Computational Science (AICS), Japan	K computer, SPARC64 VIIbx 2.0GHz, Tofu interconnect, Fujitsu	705,024	10,510.0	11,280.4	12,640
6	ODE/SC/Argonne National Laboratory, United States	Mira - BlueGene/Q, Power BGC 16C 1.60GHz, Custom, IBM	786,432	8,586.6	10,066.3	3,945

[Source:www.top500.org ]

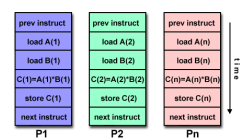
## TOP 500 2016 rank 7-13

Rank	Site	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power (kW)
7	ODE/NSA/LLNL/NSL, United States	Trinity - Cray XC40, Xeon E5-2698v3 14C 2.20GHz, Aries interconnect, Cray Inc.	301,056	8,100.9	11,078.9	4,233
8	Swiss National Supercomputing Centre (SCS), Switzerland	Piz Daint - Cray XC30, Xeon E5-2670 8C 2.600GHz, Aries interconnect, NVIDIA K80x, Cray Inc.	115,984	4,271.0	7,788.9	1,754
9	HLR5 - Höchstleistungsrechenzentrum Stuttgart, Germany	Hazel Hen - Cray XC40, Xeon E5-2680v3 12C 2.5GHz, Aries interconnect, Cray Inc.	165,088	5,642.2	7,403.5	3,415
10	King Abdulaziz University of Science and Technology, Saudi Arabia	Shakeen II - Cray XC40, Xeon E5-2698v3 16C 2.30GHz, Aries interconnect, Cray Inc.	196,608	5,537.0	7,235.2	2,834
11	Total Exploration Production, France	Pangea - SGI ICE X, Xeon Xeon E5-2670 ES-2680v3 12C 2.50GHz, Infiniband FDR, HPE/SGI	220,800	5,283.1	6,712.3	4,150
12	Texas Advanced Computing Center/Rhinc of Texas, United States	Stampede - PowerEdge C8220, Xeon E5-2688 8C 2.700GHz, Infiniband FDR, Intel Xeon Phi SE10P, Dell	462,442	5,168.1	8,520.1	4,510
13	Forschungszentrum Juelich (FZJ), Germany	JUQUEEN - BlueGene/Q, Power BGC 16C 1.600GHz, Custom Interconnect, IBM	458,752	5,008.9	5,872.0	2,301

[Source:www.top500.org ]

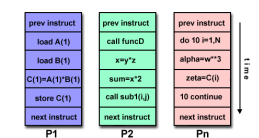
## Parallel paradigms

**SIMD**  
Single Instruction Multiple Data



[Source: computing.llnl.gov/tutorials]

**MIMD**  
Multiple Instruction Multiple Data

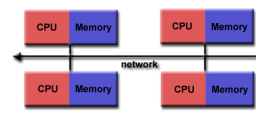


[Source: computing.llnl.gov/tutorials]

- ▶ "classical" vector systems: Cray, Convex ...
- ▶ Graphics processing units (GPU)

- ▶ Shared memory systems
  - ▶ IBM Power, Intel Xeon, AMD Opteron ...
  - ▶ Smartphones ...
  - ▶ Xeon Phi
- ▶ Distributed memory systems
  - ▶ interconnected CPUs

## MIMD Hardware: Distributed memory



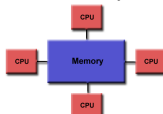
[Source: computing.llnl.gov/tutorials]

```
MPI_Send(buf, count, type, dest, tag, comm)
MPI_Recv(buf, count, type, src, tag, comm, stat)
```

- ▶ "Linux Cluster"
- ▶ "Commodity Hardware"
- ▶ Memory scales with number of CPUs interconnected
- ▶ High latency for communication
- ▶ Mostly programmed using MPI (Message passing interface)
- ▶ Explicit programming of communications: gather data, pack, send, receive, unpack, scatter

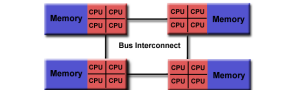
## MIMD Hardware: Shared Memory

Symmetric Multiprocessing (SMP)/Uniform memory access (UMA)



[Source: computing.llnl.gov/tutorials]

Nonuniform Memory Access (NUMA)



[Source: computing.llnl.gov/tutorials]

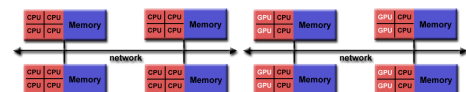
- ▶ Similar processors
- ▶ Similar memory access times

- ▶ Possibly varying memory access latencies
- ▶ Combination of SMP systems
- ▶ ccNUMA: Cache coherent NUMA

- ▶ Shared memory: one (virtual) address space for all processors involved
- ▶ Communication hidden behind memory access
- ▶ Not easy to scale large numbers of CPUs
- ▶ MPI works on these systems as well

## Hybrid distributed/shared memory

- ▶ Combination of shared and distributed memory approach
- ▶ Top 500 computers



[Source: computing.llnl.gov/tutorials]

- ▶ Shared memory nodes can be mixed CPU-GPU
- ▶ Need to master both kinds of programming paradigms

## Shared memory programming: pthreads

- ▶ Thread: lightweight process which can run parallel to others
- ▶ pthreads (POSIX threads): widely distributed
- ▶ cumbersome tuning + synchronization
- ▶ basic structure for more high level interfaces

```
#include <pthread.h>

void *PrintHello(void *threadid)
{ long tid = (long)threadid;
  printf("Hello World! It's me, thread %ld!\n", tid);
  pthread_exit(NULL);
}

int main (int argc, char *argv[])
{ pthread_t threads[NUM_THREADS];
  int rc; long t;

  for(t=0; t<NUM_THREADS; t++){
    printf("In main: creating thread %ld\n", t);
    rc = pthread_create(&threads[t], NULL, PrintHello, (void *)t);
    if (rc) {printf("ERROR: return code from pthread_create() is %d\n", rc); exit(-1);}
    pthread_exit(NULL);
  }
}
```

Source: [computing.lln.gov/tutorials](http://computing.lln.gov/tutorials)

- ▶ compile and link with

```
gcc -pthread -o pthreads pthreads.c
```

9 / 45

## Shared memory programming: C++11 threads

- ▶ Threads introduced into C++ standard with C++11
- ▶ Quite late... many codes already use other approaches
- ▶ But interesting for new applications

```
#include <iostream>
#include <thread>

void call_from_thread(int tid) {
  std::cout << "Launched by thread " << tid << std::endl;
}

int main() {
  std::thread t[num_threads];
  for (int i = 0; i < num_threads; ++i) {
    t[i] = std::thread(call_from_thread, i);
  }
  std::cout << "Launched from the main\n";
  //Join the threads with the main thread
  for (int i = 0; i < num_threads; ++i) {
    t[i].join();
  }
  return 0;
}
```

Source: <https://solarianprogrammer.com/2011/12/16/cpp11-thread-tutorial/>

- ▶ compile and link with

```
g++ -std=c++11 -pthread cpp11threads.cxx -o cpp11threads
```

10 / 45

## Thread programming: mutexes and locking

- ▶ If threads work with common data (write to the same memory address, use the same output channel) access must be synchronized
- ▶ Mutexes allow to define regions in a program which are accessed by all threads in a sequential manner.

```
#include <iostream>
#include <thread>
#include <mutex>
std::mutex mtx;

void call_from_thread(int tid) {
  mtx.lock();
  std::cout << "Launched by thread " << tid << std::endl;
  mtx.unlock();
}

int main() {
  std::thread t[num_threads];
  for (int i = 0; i < num_threads; ++i) {
    t[i] = std::thread(call_from_thread, i);
  }
  std::cout << "Launched from the main\n";
  //Join the threads with the main thread
  for (int i = 0; i < num_threads; ++i) {
    t[i].join();
  }
  return 0;
}
```

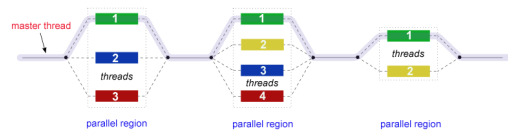
- ▶ **Barrier**: all threads use the same mutex for the same region
- ▶ **Deadlock**: two threads block each other by locking two different locks and waiting for each other to finish

11 / 45

## Shared memory programming: OpenMP

- ▶ Mostly based on pthreads
- ▶ Available in C++, C, Fortran for all common compilers
- ▶ Compiler directives (pragmas) describe *parallel regions*

```
... sequential code ...
#pragma omp parallel
{
  ... parallel code ...
}
(implicit barrier)
... sequential code ...
```



[Source: [computing.lln.gov/tutorials](http://computing.lln.gov/tutorials)]

12 / 45

## Shared memory programming: OpenMP II

```
#include <iostream>
#include <cstdlib>

void call_from_thread(int tid) {
  std::cout << "Launched by thread " << tid << std::endl;
}

int main (int argc, char *argv[])
{
  int num_threads=1;
  if (argc>1) num_threads=atoi(argv[1]);

#pragma omp parallel for
  for (int i = 0; i < num_threads; ++i)
  {
    call_from_thread(i);
  }
  return 0;
}
```

- ▶ compile and link with

```
g++ -fopenmp -o cppomp cppomp.cxx
```

13 / 45

## Example: $u = au + v$ und $s = u \cdot v$

```
double u[n],v[n];
#pragma omp parallel for
for(int i=0; i<n; i++)
  u[i]=a*v[i];

//implicit barrier
double s=0.0;
#pragma omp parallel for reduction(+:s)
for(int i=0; i<n; i++)
  s+=u[i]*v[i];
```

- ▶ Code can be parallelized by introducing compiler directives
- ▶ Compiler directives are ignored if not in parallel mode
- ▶ Write conflict with  $+$   $s$ : several threads may access the same variable
- ▶ In standard situations, reduction variables can be used to avoid conflicts

14 / 45

## Do it yourself reduction

```
#include <omp.h>
int maxthreads=omp_get_max_threads();
double s0(maxthreads);
double u[n],v[n];
for (int ithread=0;ithread<maxthreads; ithread++)
  s0[ithread]=0.0;

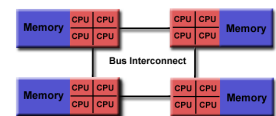
#pragma omp parallel for
for(int i=0; i<n; i++)
{
  int ithread=omp_get_thread_num();
  s0[ithread]+=u[i]*v[i];
}

double s=0.0;
for (int ithread=0;ithread<maxthreads; ithread++)
  s+=s0[ithread];
```

15 / 45

## OpenMP: further aspects

```
double u[n],v[n];
#pragma omp parallel for
for(int i=0; i<n; i++)
  u[i]=a*u[i];
```



[Quelle: [computing.lln.gov/tutorials](http://computing.lln.gov/tutorials)]

- ▶ Distribution of indices with thread is implicit and can be influenced by scheduling directives
- ▶ Number of threads can be set via OMP\_NUM\_THREADS environment variable or call to `omp_set_num_threads()`
- ▶ First Touch Principle (NUMA): first thread which "touches" data triggers the allocation of memory with the processor where the thread is running on

16 / 45

## Parallelization of PDE solution

$$\Delta u = f \text{ in } \Omega, \quad u|_{\partial\Omega} = 0$$

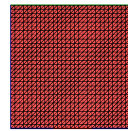
$$\Rightarrow u = \int_{\Omega} f(y)G(x,y)dy.$$

- ▶ Solution in  $x \in \Omega$  is influenced by values of  $f$  in all points in  $\Omega$
- ▶  $\Rightarrow$  global coupling: any solution algorithm needs global communication

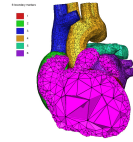
17 / 45

## Structured and unstructured grids

Structured grid



Unstructured grid



[Quelle: tetgen.org]

- ▶ Easy next neighbor access via index calculation
- ▶ Efficient implementation on SIMD/GPU
- ▶ Strong limitations on geometry
- ▶ General geometries
- ▶ Irregular, index vector based access to next neighbors
- ▶ Hardly feasible for SIMD/GPU

18 / 45

## Stiffness matrix assembly for Laplace operator for P1 FEM

$$a_{ij} = a(\phi_i, \phi_j) = \int_{\Omega} \nabla \phi_i \nabla \phi_j \, dx$$

$$= \int_{\Omega} \sum_{K \in \mathcal{T}_h} \nabla \phi_i|_K \nabla \phi_j|_K \, dx$$

Assembly loop:  
Set  $a_{ij} = 0$ .  
For each  $K \in \mathcal{T}_h$ :  
For each  $m, n = 0 \dots d$ :

$$s_{mn} = \int_K \nabla \lambda_m \nabla \lambda_n \, dx$$

$$a_{j_{\text{dof}}(K,m)j_{\text{dof}}(K,n)} = a_{j_{\text{dof}}(K,m)j_{\text{dof}}(K,n)} + s_{mn}$$

19 / 45

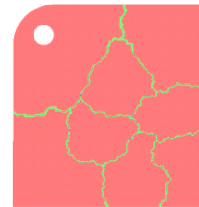
## Mesh partitioning

Partition set of cells in  $\mathcal{T}_h$ , and color the graph of the partitions.

Result:  $\mathcal{C}$ : set of colors,  $\mathcal{P}_c$ : set of partitions of given color. Then:

$$\mathcal{T}_h = \bigcup_{c \in \mathcal{C}} \bigcup_{p \in \mathcal{P}_c} p$$

- ▶ Sample algorithm:
  - ▶ Subdivision of grid cells into equally sized subsets by METIS (Karypis/Kumar)  $\rightarrow$  Partitions of color 1
  - ▶ Create separators along boundaries  $\rightarrow$  Partitions of color 2
  - ▶ "triple points"  $\rightarrow$  Partitions of color 3



- ▶ No interference between assembly loops for partitions of the same color
- ▶ Immediate parallelization without critical regions

20 / 45

## Parallel stiffness matrix assembly for Laplace operator for P1 FEM

Set  $a_{ij} = 0$ .

For each color  $c \in \mathcal{C}$

**#pragma omp parallel for**

For each  $p \in \mathcal{P}_c$ :

For each  $K \in p$ :

For each  $m, n = 0 \dots d$ :

$$s_{mn} = \int_K \nabla \lambda_m \nabla \lambda_n \, dx$$

$$a_{j_{\text{dof}}(K,m)j_{\text{dof}}(K,n)} += s_{mn}$$

- ▶ Similar structure for Voronoi finite volumes, nonlinear operator evaluation, Jacobi matrix assembly

21 / 45

## Linear system solution

- ▶ Sparse matrices
- ▶ Direct solvers are hard to parallelize though many efforts are undertaken
- ▶ Iterative methods easier to parallelize
  - ▶ partitioning of vectors + coloring inherited from cell partitioning
  - ▶ keep loop structure (first touch principle)
  - ▶ parallelize
    - ▶ vector algebra
    - ▶ scalar products
    - ▶ matrix vector products
    - ▶ preconditioners

22 / 45

## MPI - Message passing interface

- ▶ library, can be used from C,C++, Fortran, python
- ▶ de facto standard for programming on distributed memory system (since  $\approx$  1995)
- ▶ highly portable
- ▶ support by hardware vendors: optimized communication speed
- ▶ based on sending/receiving messages over network
  - ▶ instead, shared memory can be used as well
- ▶ very elementary programming model, need to hand-craft communications

23 / 45

## How to install

- ▶ OpenMP/C++11 threads come along with compiler
- ▶ MPI needs to be installed in addition
- ▶ Can run on multiple systems
- ▶ openmpi available for Linux/Mac (homebrew)/ Windows (cygwin)
  - ▶ <https://www.open-mpi.org/faq/?category=mpi-apps>
  - ▶ Compiler wrapper mpic++ - wrapper around (configurable) system compiler - proper flags + libraries to be linked
  - ▶ Process launcher mpirun
- ▶ launcher starts a number of processes which execute statements independently, occasionally waiting for each other

24 / 45

## Threads vs processes

- ▶ Threads are easier to create than processes since they don't require a separate address space.
- ▶ Multithreading requires careful programming since threads share data structures that should only be modified by one thread at a time. Unlike threads, processes don't share the same address space.
- ▶ Threads are considered lightweight because they use far less resources than processes.
- ▶ Processes are independent of each other. Threads, since they share the same address space are interdependent, so caution must be taken so that different threads don't step on each other. This is really another way of stating #2 above.
- ▶ A process can consist of multiple threads.
- ▶ MPI is based on processes, C++11 threads and OpenMP are based on threads.

25 / 45

## MPI Hello world

```
// Initialize MPI.
MPI_Init ( &argc, &argv );

// Get the number of processes.
MPI_Comm_size ( MPI_COMM_WORLD, &nproc );

// Create index vector for processes.
std::vector<unsigned long> idx(nproc+1);

// Determine the rank (number) of this process.
MPI_Comm_rank ( MPI_COMM_WORLD, &iprocc );

if ( iprocc == 0 ) cout << "The number of processes available is " << nproc << "\n";
cout << "Hello from proc " << iprocc << endl;

MPI_Finalize ( );
```

- ▶ Compile with `mpic++ mpi-hello.cpp -o mpi-hello`
- ▶ All MPI programs begin with `MPI_Init()` and end with `MPI_Finalize()`
- ▶ the *communicator* `MPI_COMM_WORLD` designates all processes in the current process group, there may be other process groups etc.
- ▶ The whole program is started *N* times as system process, not as thread: `mpirun -n N mpi-hello`

26 / 45

## MPI hostfile

```
host1 slots=1
host2 slots=2
...
```

- ▶ Distribute code execution over several hosts
- ▶ Need ssh public key access and common file system access for proper execution

27 / 45

## MPI Send

`MPI_Send (start, count, datatype, dest, tag, comm)`

- ▶ The message buffer is described by (start, count, datatype)
- ▶ The target process is specified by dest, which is the rank of the target process in the communicator specified by comm
- ▶ When this function returns, the data has been delivered to the system and the buffer can be reused. The message may not have been received by the target process.
- ▶ The tag codes some type of message

28 / 45

## MPI Receive

`MPI_Recv(start, count, datatype, source, tag, comm, status)`

- ▶ Waits until a matching (on source and tag) message is received from the system, and the buffer can be used.
- ▶ source is rank in communicator specified by comm, or `MPI_ANY_SOURCE`
- ▶ status contains further information
- ▶ Receiving fewer than count occurrences of datatype is OK, but receiving more is an error.

29 / 45

## MPI Broadcast

`MPI_Bcast(start, count, datatype, root, comm)`

- ▶ Broadcasts a message from the process with rank "root" to all other processes of the communicator
- ▶ Root sends, all others receive.

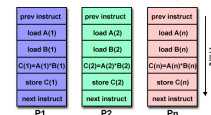
30 / 45

## Differences with OpenMP

- ▶ Programmer has to care about all aspects of communication and data distribution, even in simple situations
- ▶ In simple situations (regularly structured data) OpenMP provides reasonable defaults. For MPI these are not available
- ▶ For PDE solvers (FEM/FVM assembly) on unstructured meshes, in both cases we have to care about data distribution
- ▶ We need explicit handling of data at interfaces

31 / 45

## SIMD Hardware: Graphics Processing Units ( GPU )



[Source: computing.llnl.gov/tutorials]

- ▶ Principle useful for highly structured data
- ▶ Example: textures, triangles for 3D graphics rendering
- ▶ During the 90's, *Graphics Processing Units* (GPUs) started to contain special purpose SIMD hardware for graphics rendering
- ▶ 3D Graphic APIs (DirectX, OpenGL) became transparent to programmers: rendering could be influenced by "shaders" which essentially are programs which are compiled on the host and run on the GPU

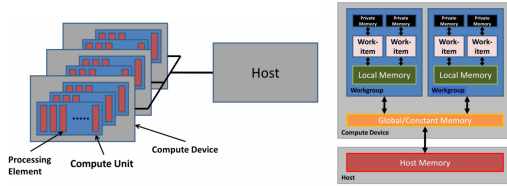


[Source:HardwareZone.com.ph]

32 / 45

## General Purpose Graphics Processing Units (GPGPU)

- ▶ Graphics companies like NVIDIA saw an opportunity to market GPUs for computational purposes
- ▶ Emerging APIs which allow to describe general purpose computing tasks for GPUs: CUDA (Nvidia specific), OpenCL (ATI/AMD designed, general purpose), OpenACC(future ?)
- ▶ GPGPUs are *accelerator cards* added to a computer with own memory and many vector processing pipelines (Nvidia Tesla K40: 12GB + 2880 units)
- ▶ CPU-GPU connection generally via mainbord bus



[Source: amd-dev.wpengine.netdna-cdn.com]

33 / 45

## GPU Programming paradigm

- ▶ CPU:
  - ▶ sets up data
  - ▶ triggers compilation of "kernels": the heavy duty loops to be executed on GPU
  - ▶ sends compiled kernels ("shaders") to GPU
  - ▶ sends data to GPU, initializes computation
  - ▶ receives data back from GPU
- ▶ GPU:
  - ▶ receive data from host CPU
  - ▶ just run the heavy duty loops in local memory
  - ▶ send data back to host CPU
- ▶ CUDA and OpenCL allow explicit management of these steps
- ▶ High efficiency only with good match between data structure and layout of GPU memory (2D rectangular grid)

34 / 45

## Example: OpenCL: computational kernel

```
__kernel void square(
    __global float* input, __global float* output)
{
    size_t i = get_global_id(0);
    output[i] = input[i] * input[i];
}
```

Declare functions with `__kernel` attribute

Defines an entry point or exported method in a program object

Use address space and usage qualifiers for memory

Address spaces and data usage must be specified for all memory objects

Built-in methods provide access to index within compute domain

Use `get_global_id` for unique work-item id, `get_group_id` for work-group, etc

[Source: <http://sa10.idav.ucdavis.edu/docs/sa10-dg-openc1-overview.pdf>]

35 / 45

## OpenCL: Resource build up, kernel creation

```
// Fill our data set with random float values
int count = 1024 * 1024;
for(i = 0; i < count; i++)
    data[i] = rand() / (float)RAND_MAX;

// Connect to a compute device, create a context and a command queue
cl_device_id device;
clGetDeviceIDs(CL_DEVICE_TYPE_GPU, 1, &device, NULL);
cl_context context = clCreateContext(0, 1, &device, NULL, NULL, NULL);
cl_command_queue queue = clCreateCommandQueue(context, device, 0, NULL);

// Create and build a program from our OpenCL-C source code
cl_program program = clCreateProgramWithSource(context, 1, (const char **) &src,
                                              NULL, NULL);
clBuildProgram(program, 0, NULL, NULL, NULL, NULL);

// Create a kernel from our program
cl_kernel kernel = clCreateKernel(program, "square", NULL);
```

[Source: <http://sa10.idav.ucdavis.edu/docs/sa10-dg-openc1-overview.pdf>]

36 / 45

## OpenCL: Data copy to GPU

```
// Allocate input and output buffers, and fill the input with data
cl_mem input = clCreateBuffer(context, CL_MEM_READ_ONLY, sizeof(float) * count,
                             NULL, NULL);

// Create an output memory buffer for our results
cl_mem output = clCreateBuffer(context, CL_MEM_WRITE_ONLY, sizeof(float) * count,
                              NULL, NULL);

// Copy our host buffer of random values to the input device buffer
clEnqueueWriteBuffer(queue, input, CL_TRUE, 0, sizeof(float) * count, data, 0,
                   NULL, NULL);

// Get the maximum number of work items supported for this kernel on this device
size_t global = count; size_t local = 0;
clGetKernelWorkGroupInfo(kernel, device, CL_KERNEL_WORK_GROUP_SIZE, sizeof(int),
                          &local, NULL);
```

[Source: <http://sa10.idav.ucdavis.edu/docs/sa10-dg-openc1-overview.pdf>]

37 / 45

## OpenCL: Kernel execution, result retrieval from GPU

```
// Set the arguments to our kernel, and enqueue it for execution
clSetKernelArg(kernel, 0, sizeof(cl_mem), &input);
clSetKernelArg(kernel, 1, sizeof(cl_mem), &output);
clSetKernelArg(kernel, 2, sizeof(unsigned int), &count);
clEnqueueNDRangeKernel(queue, kernel, 1, NULL, &global, &local, 0, NULL, NULL);

// Force the command queue to get processed, wait until all commands are complete
clFinish(queue);

// Read back the results
clEnqueueReadBuffer(queue, output, CL_TRUE, 0, sizeof(float) * count, results, 0,
                  NULL, NULL);

// Validate our results
int correct = 0;
for(i = 0; i < count; i++)
    correct += (results[i] == data[i] * data[i]) ? 1 : 0;

// Print a brief summary detailing the results
printf("Computed %d/%d correct values!\n", correct, count);
```

[Source: <http://sa10.idav.ucdavis.edu/docs/sa10-dg-openc1-overview.pdf>]

38 / 45

## OpenCL Summary

- ▶ Need good programming experience and system management skills in order to set up tool chains with properly matching versions, vendor libraries etc.
  - ▶ (I was not able to get this running on my laptop in finite time...)
- ▶ Very cumbersome programming, at least as explicit as MPI
- ▶ Data structure restrictions limit class of tasks which can run efficiently on GPUs.

39 / 45

## OpenACC (Open Accelerators)

- ▶ Idea similar to OpenMP: use compiler directives
- ▶ Future merge with OpenMP intended
- ▶ Intended for different accelerator types (GPU, Xeon Phi ...)
- ▶ GCC, Clang implementations on the way (but not yet in the usual repositories)

40 / 45

## OpenACC Sample program

```
#define N 200000000
#define v1 1024
int main(void) {
    double pi = 0.0f;
    long long i;

    #pragma acc parallel vector_length(v1)
    #pragma acc loop reduction(+:pi)
    for (i=0; i<N; i++) {
        double t= (double)((i+0.5)/N);
        pi +=4.0/(1.0+t*t);
    }

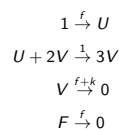
    printf("pi=%11.10f\n",pi/N);
    return 0;
}
```

- ▶ compile with `gcc-5 openacc.c -fopenacc -foffload=nvptx-none -foffload="-O3" -O3 -o openacc-gpu`
- ▶ ... but to do this one has to compile gcc with a special configuration. . .

41 / 45

## Other ways to program GPU

- ▶ WebGL: directly use capabilities of graphics hardware via html, Javascript in the browser
- ▶ Example: Gray-Scott model for Reaction-Diffusion: two chemical species.
  - ▶  $U$  is created with rate  $f$  and decays with rate  $f$
  - ▶  $U$  reacts with  $V$  to more  $V$
  - ▶  $V$  decays with rate  $f+k$ .
  - ▶  $U, V$  move by diffusion



- ▶ Stable states:
  - ▶ No  $V$
  - ▶ "Much of  $V$ ", then it feeds on  $U$  and re-creates itself
- ▶ Reaction-Diffusion equation from mass action law:

$$\begin{aligned} \partial_t u - D_u \Delta u + uv^2 - f(1-u) &= 0 \\ \partial_t v - D_v \Delta v - uv^2 + (f+k)v &= 0 \end{aligned}$$

42 / 45

## Discretization

- ▶ ... GPUs are fast so we choose the explicit Euler method:

$$\begin{aligned} \frac{1}{\tau}(u_{n+1} - u_n) - D_u \Delta u_n + u_n v_n^2 - f(1-u_n) &= 0 \\ \frac{1}{\tau}(v_{n+1} - v_n) - D_v \Delta v_n - u_n v_n^2 + (f+k)v_n &= 0 \end{aligned}$$

- ▶ Finite volume discretization on grid of size  $h$

43 / 45

## The shader

```
<script type="x-webgl/x-fragment-shader" id="timestep-shader">
precision mediump float;
uniform sampler2D u_image;
uniform vec2 u_size;
const float F = 0.05, K = 0.062, D_a = 0.2, D_b = 0.1;
const float TIMESTEP = 1.0;
void main() {
    vec2 p = gl_FragCoord.xy,
        n = p + vec2(0.0, 1.0),
        e = p + vec2(1.0, 0.0),
        s = p + vec2(0.0, -1.0),
        w = p + vec2(-1.0, 0.0);

    vec2 val = texture2D(u_image, p / u_size).xy,
        laplacian = texture2D(u_image, n / u_size).xy
        + texture2D(u_image, e / u_size).xy
        + texture2D(u_image, s / u_size).xy
        + texture2D(u_image, w / u_size).xy
        - 4.0 * val;

    vec2 delta = vec2(D_a * laplacian.x - val.x*val.y*val.y + F * (1.0-val.x),
        D_b * laplacian.y + val.x*val.y*val.y - (K+F) * val.y);

    gl_FragColor = vec4(val + delta * TIMESTEP, 0, 0);
}
</script>
```

- ▶ Embedded as script into html page

44 / 45

## Why does this work so well here ?

- ▶ Data structure fits very well to topology of GPU
  - ▶ rectangular grid
  - ▶ 2 unknowns to be stored in x,y components of vec2
- ▶ GPU speed allows to "break" time step limitation of explicit Euler
- ▶ Data stay within the graphics card: once we loaded the initial value, all computations, and rendering use data which are in the memory of the graphics card.
- ▶ Depending on the application, choose the best way to proceed
- ▶ e.g. deep learning (especially training speed)

45 / 45