# Advanced OpenACC

Steve Abbott <sabbott@nvidia.com>, November 17, 2017

NVIDIA.

# AGENDA
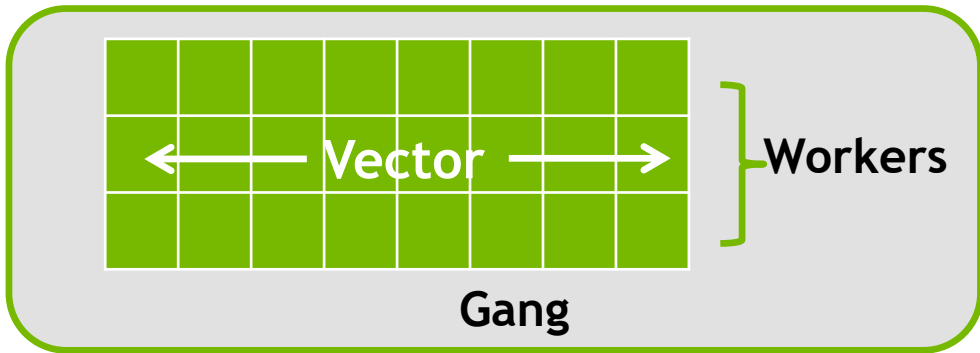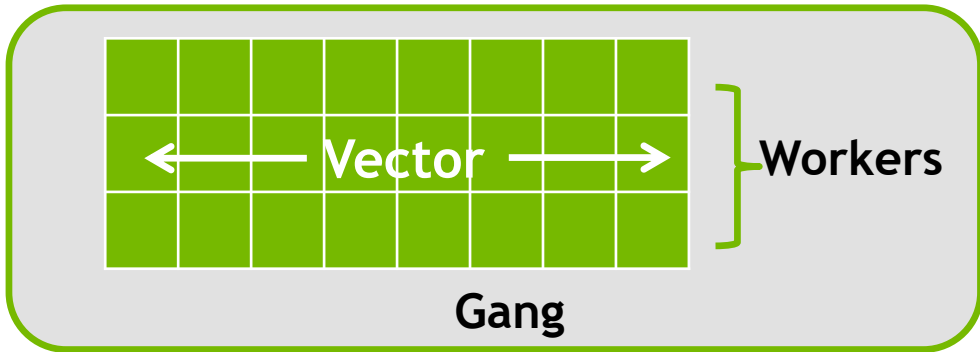
Expressive Parallelism

Pipelining

Routines

# The loop Directive

The `loop` directive gives the compiler additional information about the *next* loop in the source code through several clauses.

- **independent** – all iterations of the loop are independent

- **collapse(N)** – turn the next N loops into one, flattened loop

- **tile(N[,M,…])** - break the next 1 or more loops into *tiles* based on the provided dimensions.

NVIDIA.

# OpenACC: 3 Levels of Parallelism



- *Vector* threads work in lockstep (SIMD/SIMT parallelism)
- *Workers* compute a vector
- *Gangs* have 1 or more workers and share resources (such as cache, the streaming multiprocessor, etc.)
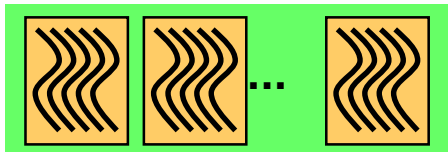- Multiple gangs work independently of each other

# Execution Model

**Software**                **Hardware**

Thread                      Scalar
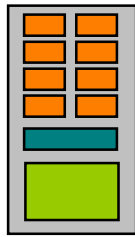                            Processor

Thread                      Multiprocessor
Block

Grid                        Device

Threads are executed by scalar processors

Thread blocks are executed on multiprocessors

Thread blocks do not migrate

Several concurrent thread blocks can reside on one
multiprocessor - limited by multiprocessor
resources (shared memory and register file)

A kernel is launched as a grid of thread blocks

# gang, worker, vector Clauses

*gang*, *worker*, and *vector* can be added to a loop clause

A parallel region can only specify one of each gang, worker, vector

Control the size using the following clauses on the parallel region

num_gangs(n), num_workers(n), vector_length(n)

```
#pragma acc parallel loop gang
 for (int i = 0; i < n; ++i)
  #pragma acc loop vector
   for (int j = 0; j < n; ++j)
    ...
```

```
#pragma acc parallel vector_length(32)
#pragma acc loop gang worker
for (int i = 0; i < n; ++i)
   #pragma acc loop vector
   for (int j = 0; j < n; ++j)
     ...
```

NVIDIA.

# gang, worker, vector Clauses

*gang*, *worker*, and *vector* can be added to a kernels loop clause too

Since different loops in a kernels region may be parallelized differently, fine-tuning is done as a parameter to the gang, worker, and vector clauses.

```
#pragma acc kernels loop gang
 for (int i = 0; i < n; ++i)
   #pragma acc loop vector
     for (int j = 0; j < n; ++j)
       ...
```

```
#pragma acc kernels loop gang worker
for (int i = 0; i < n; ++i)
   #pragma acc loop vector(32)
    for (int j = 0; j < n; ++j)
      ...
```

⬨ **nVIDIA.**

# The collapse Clause

**collapse(n):** Takes the next *n* tightly-nested loops, folds them into one, and applies the OpenACC directives to the new loop.

```
#pragma acc parallel loop \
  collapse(2)
for(int i=0; i<N; i++)
  for(int j=0; j<M; j++)
    ...
```
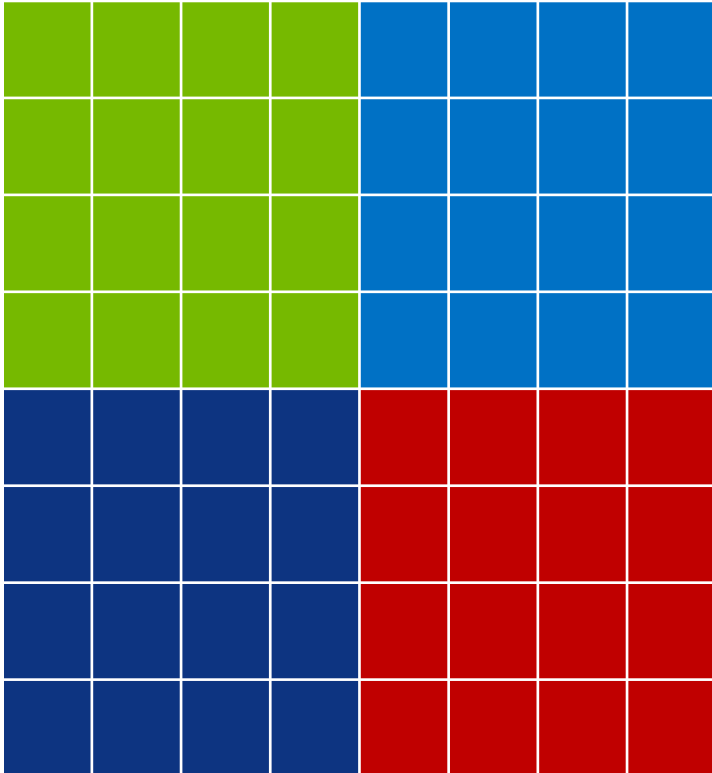
```
#pragma acc parallel loop
for(int ij=0; ij<N*M; ij++)
    ...
```

Why?

- Collapse outer loops to enable creating more gangs.
- Collapse inner loops to enable longer vector lengths.
- Collapse all loops, when possible, to do both.

NVIDIA.

# The tile clause

Operate on smaller blocks of the operation to exploit data locality

```
#pragma acc loop tile(4,4)
  for(i = 1; i <= ROWS; i++) {
    for(j = 1; j <= COLUMNS; j++) {
      Temp[i][j] = 0.25 *
        (Temp_last[i+1][j] +
        Temp_last[i-1][j] +
        Temp_last[i][j+1] +
        Temp_last[i][j-1]);
    }
  }
```

NVIDIA.

# Stride-1 Memory Accesses

```
for(i=0; i<N; i++)
  for(j=0; j<M; j++)
  {
    A[i][j][1] = 1.0f;
    A[i][j][2] = 0.0f;
  }
}
```

```
for(i=0; i<N; i++)
  for(j=0; j<M; j++)
  {
    A[1][i][j] = 1.0f;
    A[2][i][j] = 0.0f;
  }
}
```

The fastest dimension is length 2 and fastest loop strides by 2.

Now the inner loop is the fastest dimension through memory.

# Stride-1 Memory Accesses

```
for(i=0; i<N; i++)
  for(j=0; j<M; j++)
  {
    A[i][j].a = 1.0f;
    A[i][j].b = 0.0f;
  }
}
```

```
for(i=0; i<N; i++)
  for(j=0; j<M; j++)
  {
    Aa[i][j] = 1.0f;
    Ab[i][j] = 0.0f;
  }
}
```

If all threads access the "a" element, they will be accesses every-other memory element.

Now all threads are access contiguous elements of Aa and Ab.

NVIDIA.

# Optimize Loop Performance

```
#pragma acc data copy(A) create(Anew)
while ( err > tol && iter < iter_max ) {
  err=0.0;

#pragma acc parallel loop device_type(nvidia) tile(32,4)
  for( int j = 1; j < n-1; j++) {
    for(int i = 1; i < m-1; i++) {

      Anew[j][i] = 0.25 * (A[j][i+1] + A[j][i-1] +
                           A[j-1][i] + A[j+1][i]);

      err = max(err, abs(Anew[j][i] - A[j][i]));
    }
  }
#pragma acc parallel loop device_type(nvidia) tile(32,4)
  for( int j = 1; j < n-1; j++) {
    for( int i = 1; i < m-1; i++ ) {
      A[j][i] = Anew[j][i];
    }
  }
}
  iter++;
}
```
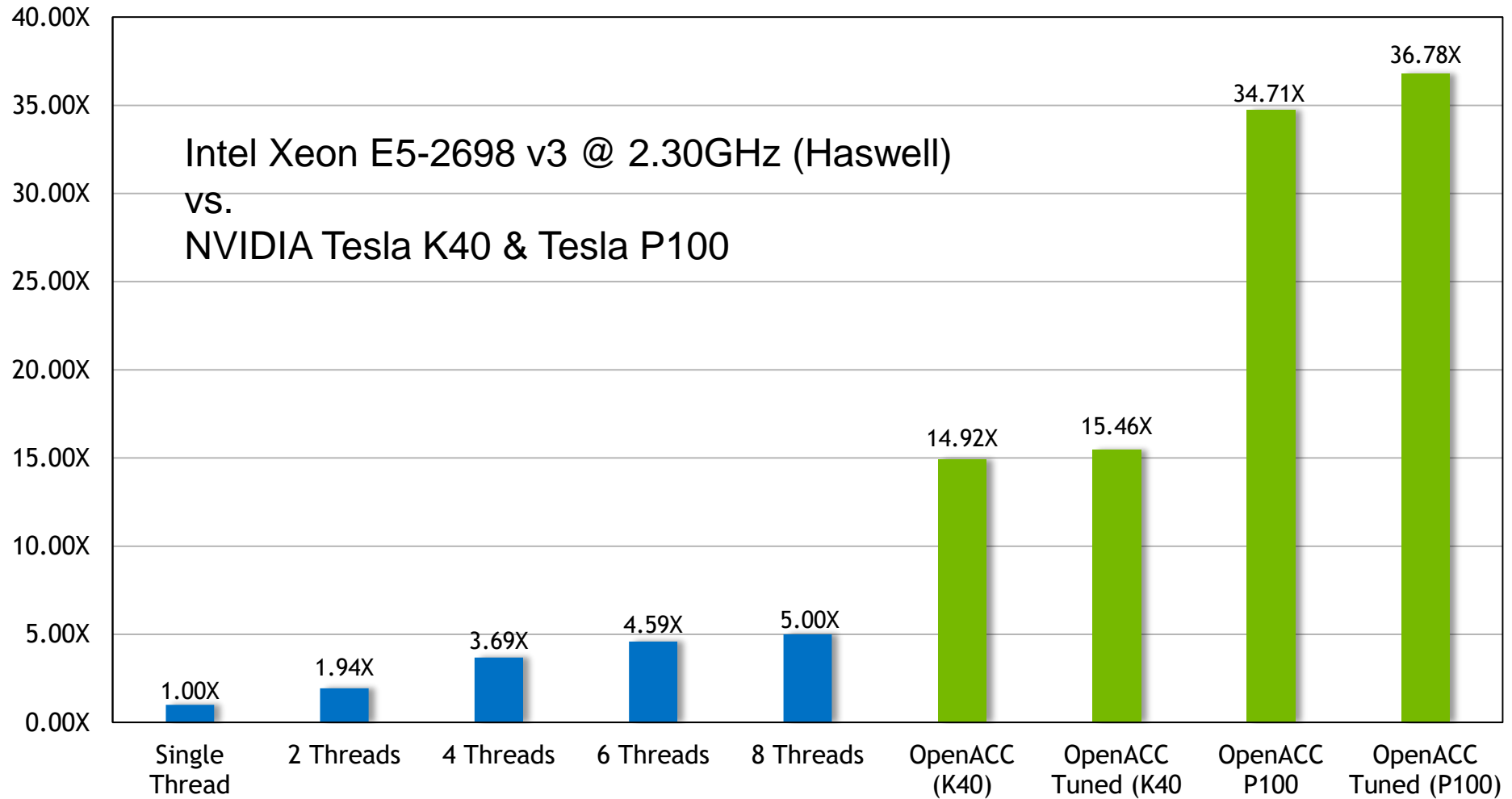
"Tile" the next two loops into 32x4 blocks, but only on NVIDIA GPUs.

# Speed-Up (Higher is Better)

Intel Xeon E5-2698 v3 @ 2.30GHz (Haswell)
vs.
NVIDIA Tesla K40 & Tesla P100

| Category | Speed-Up |
|---|---|
| Single Thread | 1.00X |
| 2 Threads | 1.94X |
| 4 Threads | 3.69X |
| 6 Threads | 4.59X |
| 8 Threads | 5.00X |
| OpenACC (K40) | 14.92X |
| OpenACC Tuned (K40) | 15.46X |
| OpenACC P100 | 34.71X |
| OpenACC Tuned (P100) | 36.78X |

*Compiler: PGI 16.10*

# Asynchronous Programming

Programming multiple operations without immediate synchronization

Real World Examples:

- Cooking a Meal: Boiling potatoes while preparing other parts of the dish.

- Three students working on a project on George Washington, one researches his early life, another his military career, and the third his presidency.

- Automobile assembly line: each station adds a different part to the car until it is finally assembled.
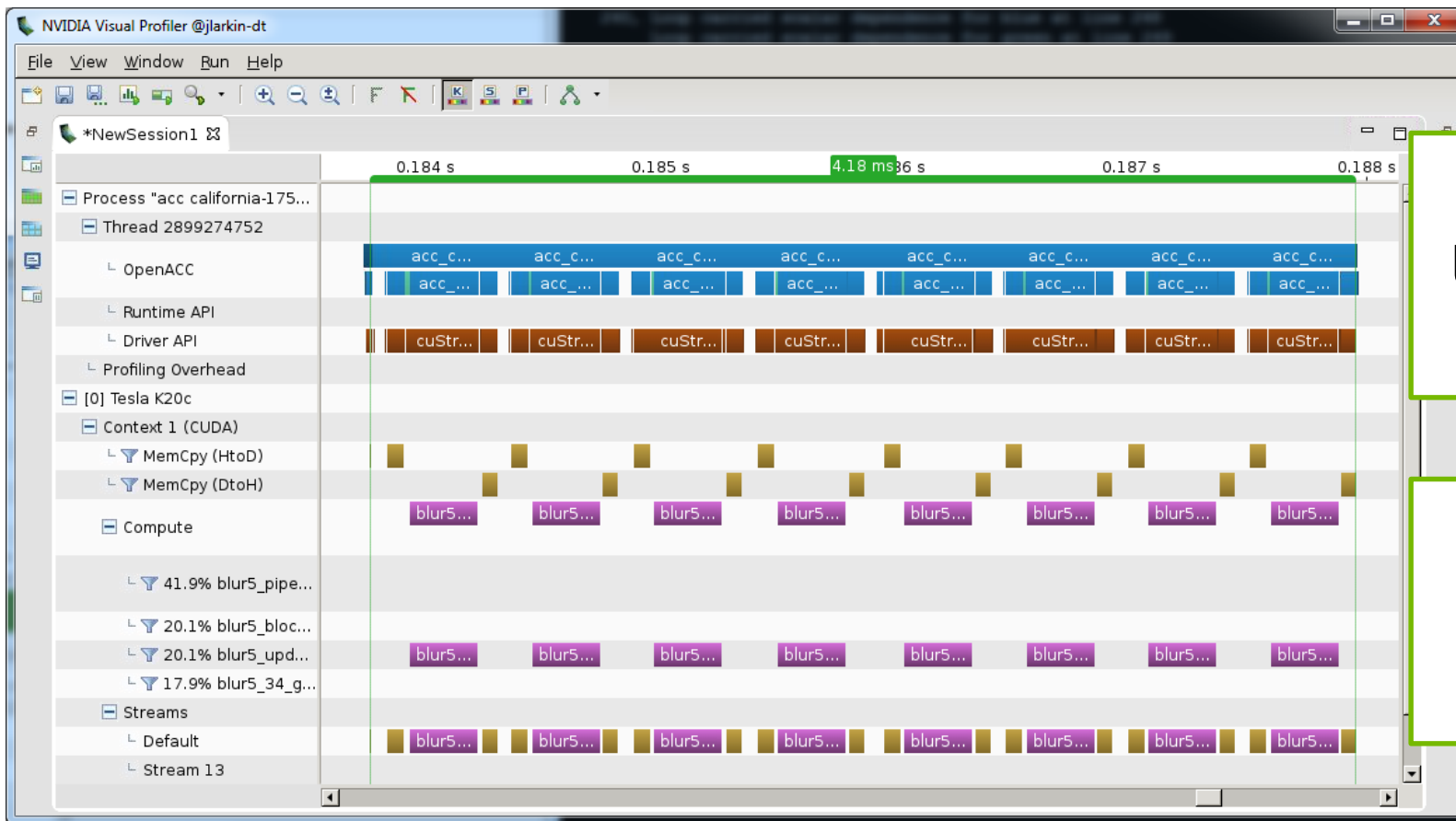
NVIDIA.

# Asynchronous OpenACC

So far, all OpenACC directives have been synchronous with the host

- Host waits for the parallel loop to complete

- Host waits for data updates to complete

Most OpenACC directives can be made asynchronous

- Host issues multiple parallel loops to the device before waiting

- Host performs part of the calculation while the device is busy

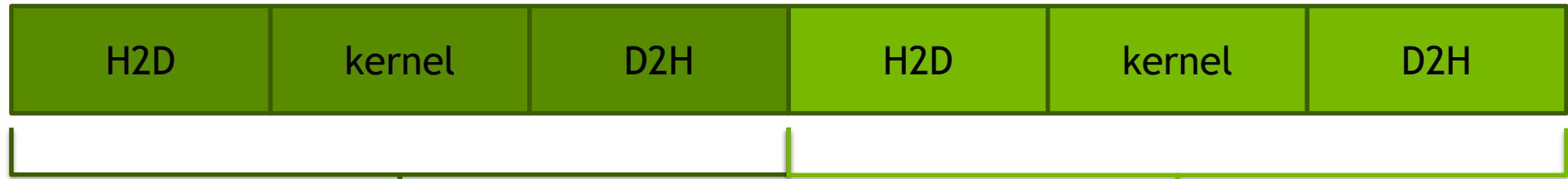- Data transfers can happen before the data is needed
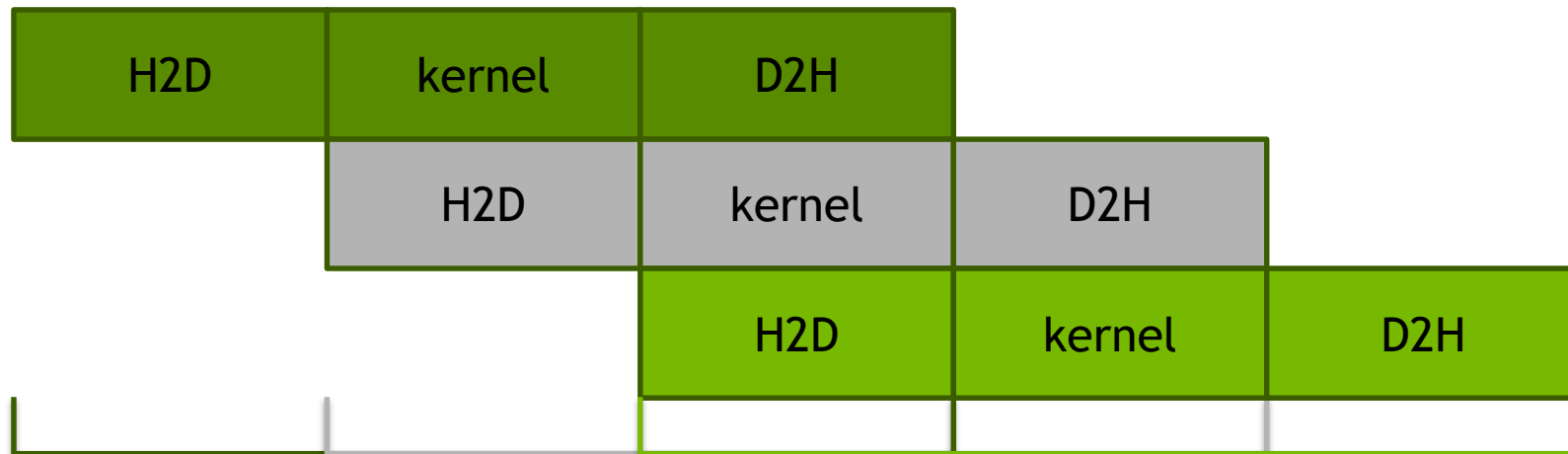
# GPU Timeline Blocked Updates



Compute and Updates happen in blocks.

The last step is to overlap compute and copy.

# Pipelining Data Transfers

| H2D | kernel | D2H | H2D | kernel | D2H |
|-----|--------|-----|-----|--------|-----|

Two Independent Operations Serialized

| H2D | kernel | D2H | | |
|-----|--------|-----|---|---|
| | H2D | kernel | D2H | |
| | | H2D | kernel | D2H |

NOTE: In real applications, your boxes will not be so evenly sized.

Overlapping Copying and Computation

NVIDIA.

# OpenACC async and wait

async(n): launches work asynchronously in *queue n*

wait(n): blocks host until all operations in *queue n* have completed

Work queues operate in-order, serving as a way to express dependencies.

Work queues of different numbers *may* (*or may not*) run concurrently.

```
#pragma acc parallel loop async(1)
...
#pragma acc parallel loop async(1)
for(int i=0; i<N; i++)
  ...
#pragma acc wait(1)
for(int i=0; i<N; i++)
```

If *n* is not specified, *async* will go into a default queue and *wait* will wait all previously queued work.
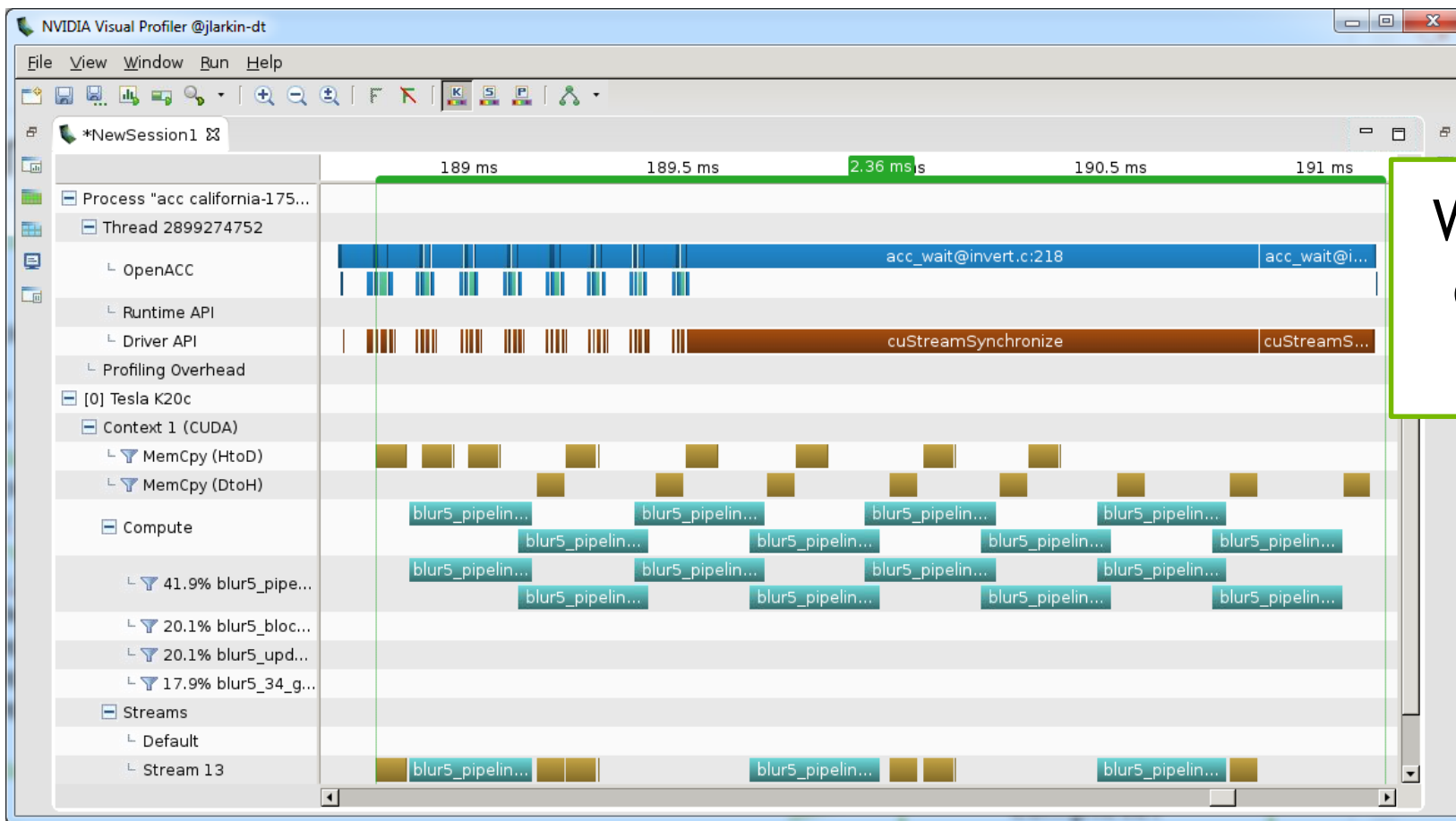
NVIDIA.

# Pipelined Code

```
#pragma acc data create(imgData[w*h*ch],out[w*h*ch])
               copyin(filter)
{
for ( long blocky = 0; blocky < nblocks; blocky++)
{
  long starty = MAX(0,blocky * blocksize - filtersize/2);
  long endy   = MIN(h,starty + blocksize + filtersize/2);
#pragma acc update device(imgData[starty*step:(endy-starty)*step]) async(block%3+1)
  starty = blocky * blocksize;
  endy = starty + blocksize;
#pragma acc parallel loop collapse(2) gang vector async(block%3+1)
  for (y=starty; y<endy; y++) for ( x=0; x<w; x++ ) {
    <filter code ommitted>
    out[y * step + x * ch]      = 255 - (scale * blue);
    out[y * step + x * ch + 1 ] = 255 - (scale * green);
    out[y * step + x * ch + 2 ] = 255 - (scale * red);
  }
#pragma acc update self(out[starty*step:blocksize*step]) async(block%3+1)
}
#pragma acc wait
}
```

Cycle between 3 async queues by blocks.

Wait for all blocks to complete.

# GPU Timeline Pipelined



We're now able to overlap compute and copy.

# OpenACC Routine Directive

Specifies that the compiler should generate a device copy of the function/subroutine and what type of parallelism the routine contains.

**Clauses:**

**gang/worker/vector/seq**
Specifies the level of parallelism contained in the routine.

**bind**
Specifies an optional name for the routine, also supplied at call-site

**no_host**
The routine will only be used on the device

**device_type**
Specialize this routine for a particular device type.

You *must* declare *one* level of parallelism on the routine directive.

NVIDIA.

# Routine Directive: C/C++

```
// foo.h
#pragma acc routine seq
double foo(int i);


// Used in main()
#pragma acc parallel loop
for(int i=0;i<N;i++) {
  array[i] = foo(i);
}
```

- At function source:
  - Function needs to be built for the GPU.
  - It will be called by each thread (sequentially)
- At call the compiler needs to know:
  - Function will be available on the GPU
  - It is a sequential routine

NVIDIA.

# OPENACC Resources

Guides ● Talks ● Tutorials ● Videos ● Books ● Spec ● Code Samples ● Teaching Materials ● Events ● Success Stories ● Courses ● Slack ● Stack Overflow
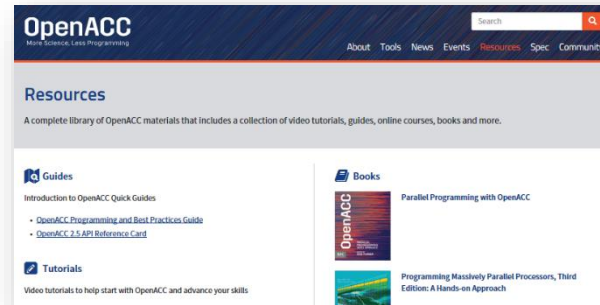
## FREE Compilers



https://www.openacc.org/community#slack

## Resources
https://www.openacc.org/resources



## Success Stories
https://www.openacc.org/success-stories



## Compilers and Tools
https://www.openacc.org/tools



## Events
https://www.openacc.org/events