

# Tricks, Tips, and Timings: The Data Movement Strategies You Need to Know

David Appelhans

GPU Technology Conference  
March 26, 2018



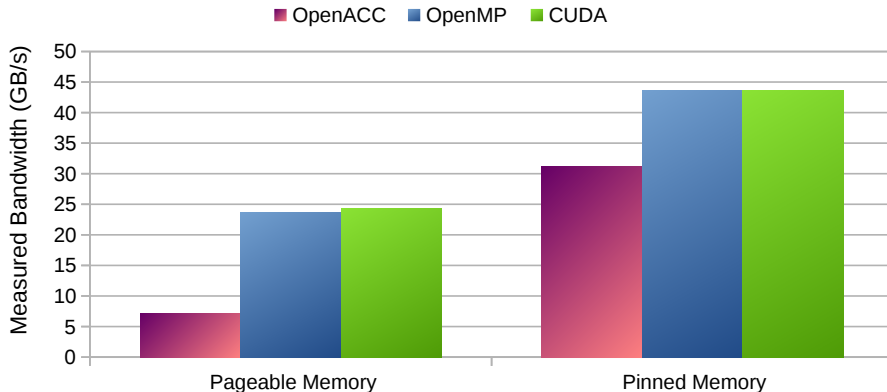
# INTRODUCTION

- My role: readying applications for SUMMIT and SIERRA supercomputers(past 3 years).
- Talk is a summary of data movement techniques, especially when working with NVLINK:
  - Importance of pinned memory. (Interoperability, CUDA+OpenMP+OpenACC)
  - Zero-copy tricks. (Interoperability, CUDA+OpenMP)
  - Dealing with nested data structures. (Efficiency, CUDA)
- All code examples are available on my public Github page.  
<https://github.com/dappelha/gpu-tips/nvtx>

# MOTIVATION: WHY YOU SHOULD PIN YOUR MEMORY

## Pageable vs Pinned HtoD Bandwidth Impact

Dual socket P9 + 6 Volta GPUs



Hint: make sure your task starts in the appropriate socket: `taskset -c 0 ./test`

# PINNED MEMORY OPTION 1:

Use CUDA Fortran<sup>1</sup> pinned attribute to pin at allocation time,

---

```
1 real(kind=8), pinned, allocatable :: p_A(:)
2 allocate ( p_A(N) )
3 !$omp target data map(alloc:p_A)
4 do i=1,samples
5     !$omp target update to(p_A)
6     ...
7 enddo
```

---

Can also check success of pinning:

---

```
1 logical :: pstat
2 allocate ( p_A(N), pinned=pstat)
3 if (.not. pstat) print *, "ERROR: p_A was not pinned"
```

---

---

<sup>1</sup>PGI and XLF compilers both support CUDA Fortran, so the pinned attribute can easily be combined with directives.

## PINNED MEMORY OPTION 2:

Pin already allocated memory,<sup>2</sup>

---

```
1 use, intrinsic :: iso_c_binding
2 use cudafor
3 real, pointer, contiguous :: phi (:,:)
4 allocate ( phi(dim1, dim2) ) ! phi can also be pointer passed from C++
5 istat = cudaHostRegister(C_LOC(phi(1,1)), sizeof(phi), cudaHostRegisterMapped)
6
7 !$acc enter data create (phi)
8 do i=1,samples
9     !$acc update self (phi)
10    ...
11 enddo
```

---

Warning: act of pinning memory is very slow. Memory should only be pinned if it is going to be used for data transfers.

---

<sup>2</sup>This technique is especially useful if the memory was allocated outside the developers control, for example in a C++ calling routine.

# OPENACC INTEROPERABILITY WARNING 1

You must use the flag `-ta=tesla:pinned` in order for OpenACC to benefit from pinned memory.

- ❶ *Compiling* with the flag `-ta=tesla:pinned` forces **all** memory to be pinned memory. This is a big hammer approach.
- ❷ *Linking* the final executable with `-ta=tesla:pinned` causes the OpenACC runtime to check if an array is already pinned. This gives fine grain user control.

## OPENACC INTEROPERABILITY WARNING 2

The OpenACC runtime uses a memory pool on the device to save from repeated allocation/deallocation of device memory. Can cause trouble when mixing CUDA with OpenACC.

---

```
1 integer :: N = 8*gigabyte
2 real(kind=8), allocatable :: A(:)
3 real(kind=8), device, allocatable :: d_A(:)
4 allocate ( A(N) )
5 !$acc enter data create (A)
6 !$acc exit data delete (A) ! <---not truly free'd unless PGI_ACC_MEM_MANAGE=0
7 allocate ( d_A(N) ) ! <----- can then run out of device memory
```

---

To disable this optimization, set the environment flag **PGI\_ACC\_MEM\_MANAGE=0** and the runtime will free the data at the exit data.

# USES OF ZERO COPY

Zero copy refers to accessing host resident pinned memory directly from a GPU without having to copy the data to the device beforehand (i.e. there are zero device copies).

- Quick overlap of data movement and kernel compute (unified/managed memory is better for this purpose)
- Large arrays where only small percent of data is accessed in random pattern.
- All data is accessed, but read/write pattern is strided/not coalesced.
- Efficiently populating components of a structure, avoiding the overhead of many copy API calls by using GPU threads to fetch data directly.



# CUDA ZERO COPY SETUP

To set up zero copy of a basic array in Fortran, use a CUDA API to get a device pointer that points to the pinned host array, and then associate a fortran array with that C device pointer, specifying the Fortran array attributes.

---

```
1 use iso_c_binding ! provides c_f_pointer and C_LOC
2 ! zero copy pointers for psib
3 type(C_DEVPTR)      :: d_psib_p
4 real(adqt), device, allocatable :: pinned_psib (:,:,)
5
6 ! sets up zero copy of psib on device.
7 istat = cudaHostGetDevicePointer(d_psib_p, C_LOC(psib(1,1,1)), 0)
8 ! Translate that C pointer to the fortran array with given dimensions
9 call c_f_pointer(d_psib_p, pinned_psib, [QuadSet%Groups, Size%nbelem, QuadSet%NumAngles])
```

---

## OPENMP ZERO COPY EXAMPLE

Only requires CUDA pinned array and OpenMP `is_device_ptr` clause.

---

```
1 real(kind=8), pinned, allocatable :: A(:,:) ,At(:,:)
2 allocate ( A(nx,ny), At(ny,nx) )
3
4 ! Transpose in the typical way:
5 !$omp target enter data map(alloc:A,At)
6 call transpose(A,At,nx,ny)
7 !$omp target update from(At)
8 !$omp target exit data map(delete:At)
9
10 ! Ensure device has finished for accurate benchmarking
11 ierr = cudaDeviceSynchronize()
12
13 ! Transpose using zero copy for At.
14 ! At is no longer mapped—is_device_ptr(At) will
15 ! allow addressing host pinned memory (zero copy)
16 call transpose_zero_copy(A,At,nx,ny)
```

---

continued on next slide

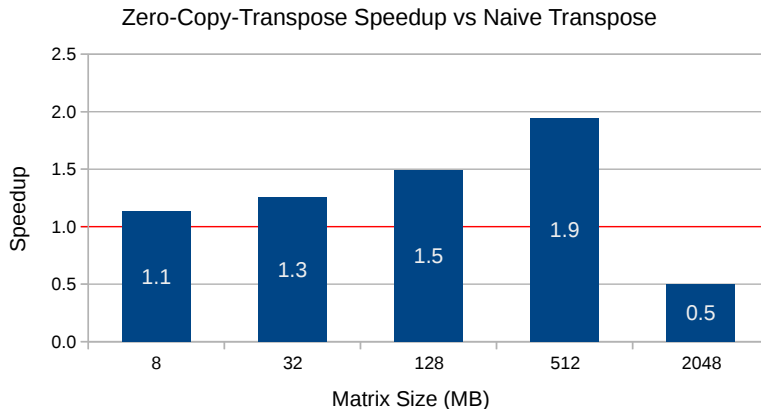
## OPENMP ZERO COPY EXAMPLE CONTINUED

---

```
1  subroutine transpose_zero_copy(A,At,nx,ny)
2    ! example of strided writes to an array that lives on the host
3    implicit none
4    real(kind=8), intent(in) :: A(:, :)
5    real(kind=8), intent(out) :: At(:, :)
6    integer, intent(in) :: nx, ny
7    integer :: i, j
8    !$omp target teams distribute parallel do is_device_ptr(At)
9    do j=1,ny
10       do i=1,nx
11          At(j,i) = A(i,j)
12       enddo
13    enddo
14    return
15 end subroutine transpose_zero_copy
```

---

# OPENMP ZERO COPY TRANSPOSE



**Figure :** Power9 + V100 results of doing a traditional matrix transpose and then copying back from GPU vs doing the transpose directly into pinned host memory.

# NESTED DATA STRUCTURES

## Motivation:

---

```
1  subroutine my_kernel_to_port (...)  
2      ...  
3      element(id)%val(n) = element(id)%x(n)*element(id)%y(n)  
4      ...  
5  end subroutine
```

---

Production codes often have dynamic structures with dynamic components.

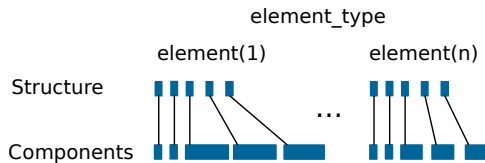
- Flattening data structures is messy (index arrays required for unstructured data) and invasive.
- Would like to keep nested references in compute kernel for portability.
- Often only parts of the data structure need to be used on the GPU.

# NESTED DATA STRUCTURES

Two Topics:

- **How do you make them referencable on the device?**
- How do you efficiently move data into them?

## Often only parts of the data structure need to be used on the GPU



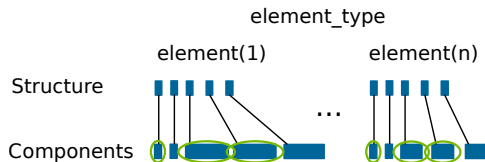

---

```

type, public :: element_type
  integer                                :: Nnodes
  real(kind=8)                           :: volume
  real(kind=8), allocatable, pinned :: x(:)
  real(kind=8), allocatable, pinned :: y(:)
  real(kind=8), allocatable, pinned :: val(:)
  real(kind=8), allocatable              :: old(:)
end type element_type
  
```

---

Often only parts of the data structure need to be used on the GPU




---

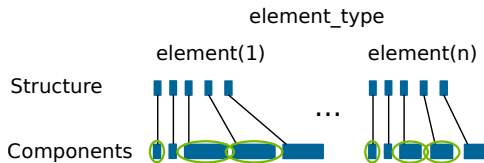
```

type, public :: element_type
  integer                               :: Nnodes
  real(kind=8)                          :: volume
  real(kind=8), allocatable, pinned :: x(:)
  real(kind=8), allocatable, pinned :: y(:)
  real(kind=8), allocatable, pinned :: val(:)
  real(kind=8), allocatable              :: old(:)
end type element_type
  
```

---



Often only parts of the data structure need to be used on the GPU




---

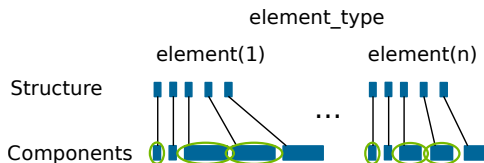
```

type, public :: element_type
  integer                                :: Nnodes
  real(kind=8)                          :: volume
  real(kind=8), allocatable, pinned :: x(:)
  real(kind=8), allocatable, pinned :: y(:)
  real(kind=8), allocatable, pinned :: val(:)
  real(kind=8), allocatable             :: old(:)
end type element_type
  
```

---

Can create a skinny version of the data structure with components that are device variables.

## Often only parts of the data structure need to be used on the GPU




---

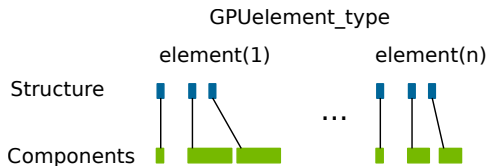
```

type, public :: element_type
  integer                                :: Nnodes
  real (kind=8)                          :: volume
  real (kind=8), allocatable, pinned :: x (:)
  real (kind=8), allocatable, pinned :: y (:)
  real (kind=8), allocatable, pinned :: val (:)
  real (kind=8), allocatable             :: old (:)
end type element_type

```

---

Can create a skinny version of the data structure with components that are device variables.




---

```

type, public :: GPUelement_type
  integer                                :: Nnodes
  real (kind=8), device, allocatable :: x (:)
  real (kind=8), device, allocatable :: y (:)
  real (kind=8), device, allocatable :: val (:)
end type GPUelement_type

```

---

Legend: ■ Host ■ Device ■ Managed

This gives a way to loop through the structure and allocate device components while on the host:

---

```
! can allocate on the host
do id=1,Nelements
  Nnodes = element(id)%Nnodes
  allocate (element(id)% x(Nnodes) )
enddo
```

---

but still cannot use the structure in a device kernel.

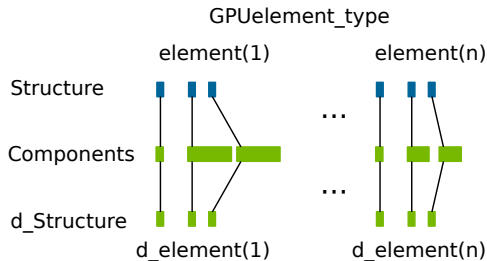
---

```
attributes (global) subroutine cuda_kernel (...)
  ...
  x = element(id)%x ! <- invalid reference of element
  ...
end subroutine cuda_kernel
```

---

Two ways to address this.

# OPTION1: DEVICE STRUCTURE THAT POINTS TO THE SAME DEVICE COMPONENT MEMORY:




---

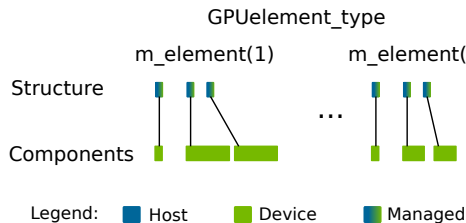
```

1  ! CPU valid version for use on the host:
2  type(GPUelement_type), pinned, allocatable :: GPUelement(:)
3  ! Device structure that will point to same device components:
4  type(GPUelement_type), device, allocatable :: d_GPUelement(:)
5  allocate ( GPUelement(Nelements) )
6  allocate ( d_GPUelement(Nelements) )
7  ! GPUelement%components(:) can be allocated in host code
8  ! copy scalars and addresses of host struct to d struct :
9  cudaMemcpy(d_GPUelement, GPUelement, size(GPUelement))

```

---

# BEST: ALLOCATE THE STRUCTURE AS MANAGED MEMORY:




---

```

1 ! host and device valid structure (managed memory):
2 type(GPUelement_type), managed, allocatable :: m_GPUelement(:)
3
4 allocate ( m_GPUelement(Nelements) )
5 ! m_GPUelement%components(:) can be allocated in host code
6 ! they still live on the device
7 ...
8 ! can prefetch structure to device to avoid pagefault:
9 cudaMemPrefetchAsync(m_GPUelement,size(m_GPUelement),device=0,
stream=0)

```

---

# NESTED DATA STRUCTURES

Two Topics:

- How do you make them referencable on the device?
- **How do you efficiently move data into them?**

# EFFICIENTLY POPULATING NESTED STRUCTURES

A naive implementation for getting data structures populated on the GPU usually looks like this:

---

```
1  ! Host data structure has been created and populated.
2  ! GPU data structure has also been allocated.
3
4  ! still need to populate the values from the host version of the data structure :
5  do id=1, Nelements
6      GPUelement(id)%Nnodes = element(id)%Nnodes ! implicit cudaMemcpy
7      GPUelement(id)%x = element(id)%x           ! implicit cudaMemcpy
8      GPUelement(id)%y = element(id)%y           ! implicit cudaMemcpy
9      GPUelement(id)%val = element(id)%val        ! implicit cudaMemcpy
10 enddo
```

---

**This becomes very slow when Nelements is large.**

# EFFICIENTLY POPULATING NESTED STRUCTURES

There are two ways to fix the naive approach,

- ❶ BETTER: push from host with `cudaMemcpyAsync` calls instead of the numerous blocking calls done above,
- ❷ BEST: pull the data from the GPU by populating the arrays from within a device kernel using zero copy.



# PUSH WITH LOOP OF ASYNC CALLS

Issue copy calls to default stream which is asynchronous to the CPU:

---

```
1 do id=1, Nelements
2   GPUelement(id)%Nnodes = element(id)%Nnodes ! cpu to cpu copy
3   istat =cudaMemcpyAsync(GPUelement(id)%x, element(id)%x, size(element(id)%x), stream=0)
4   istat =cudaMemcpyAsync(GPUelement(id)%y, element(id)%y, size(element(id)%y), stream=0)
5   istat =cudaMemcpyAsync(GPUelement(id)%val, element(id)%val, size(element(id)%val), stream=0)
6 enddo
```

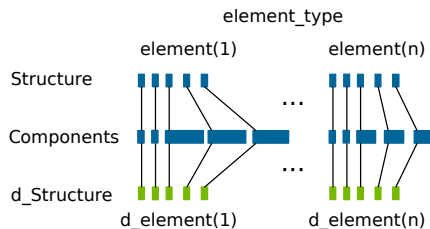
---

Can do similar in OpenMP 4, with update nowait.

Host threading over id loop made no difference in performance.

# PULLING FROM THE GPU

Set up a structure to reference pinned host components from the device (zero copy of structure components):




---

```

1  ! to use element on the device, we have to make a device valid copy called d_element:
2  istat = cudaMemcpyAsync(d_element, element, size(element), 0)
3
4  ! Now we can use these in a CUDA kernel to zero copy
5  ! the component data from d_element into d_GPUElement
6  call  set_elements_kernel <<<blocks,threads>>>(d_GPUElement,d_element,Nelements)

```

---

Launch a kernel to have GPU threads pull data from structures on the host into structures on the device:

---

```

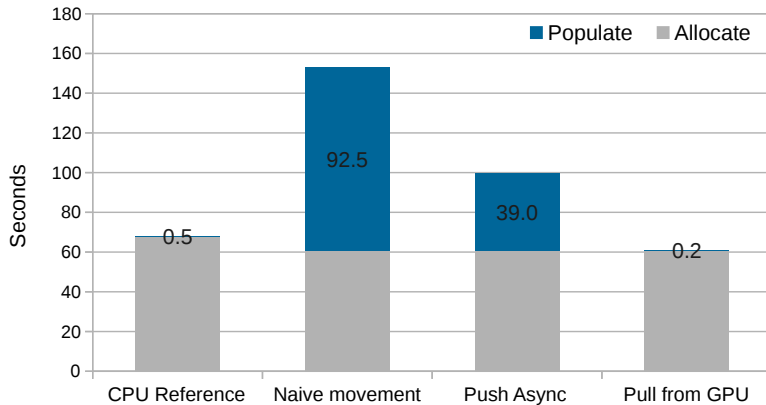
1  attributes (global) subroutine set_elements_kernel (GPUelement,element, Nelements)
2      implicit none
3      ! kernel that uses zero copy to populate the GPUelement structure.
4      type(GPUelement_type), device, intent (inout) :: d_GPUelement(:) ! members are device
5      type(element_type), device, intent (in) :: p_element(:) ! members are pinned host
6      integer , value, intent (in) :: Nelements
7      integer :: id, Nnodes, node
8
9      do id=blockIdx%x,Nelements, gridDim%x
10         Nnodes = p_element(id)%Nnodes
11         d_GPUelement(id)%Nnodes = Nnodes
12         do node = threadIdx%x, Nnodes, blockDim%x
13             d_GPUelement(id)%x(node) = p_element(id)%x(node)
14             d_GPUelement(id)%y(node) = p_element(id)%y(node)
15             d_GPUelement(id)%val(node) = p_element(id)%val(node)
16         enddo
17     enddo
18
19 end subroutine set_elements_kernel

```

---

## Populating Nested Data Structures

POWER9 + V100  
1M elements, 4 nodes/element



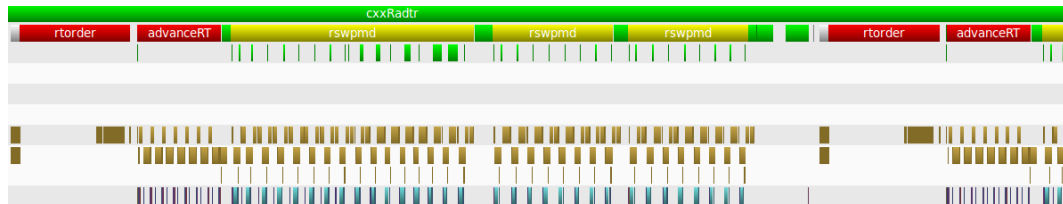
**462x speedup when pulling from the GPU!**

# CLOSING REMARKS

nvtx marker module available for Fortran

- Easily mark regions of host code for viewing in the Nvidia Visual Profiler.
- Works with CUDA, OpenMP, and OpenACC.
- Newly supports non-nested marked regions.
- Very helpful to understand flow of your application.

Available at <https://github.com/dappelha/gpu-tips/nvtx>



# CONCLUSIONS

- Pinning memory is important, even when using directives
- Managed memory makes nested data structure book keeping easier.
- Still important to efficiently populate data structures.
- Zero-copy populating from the device is the fastest method (462x).

Various example codes are available at

<https://github.com/dappelha/gpu-tips>

Questions?

David Appelhans - [dappelh@us.ibm.com](mailto:dappelh@us.ibm.com)