

# Performance Portability of a Coarray Fortran atmospheric Model

Damian Rouson



Ethan Gutman  
Alessandro Fanfarillo



Brian Friesen



# Overview



**Fortran 2018 in a Nutshell**



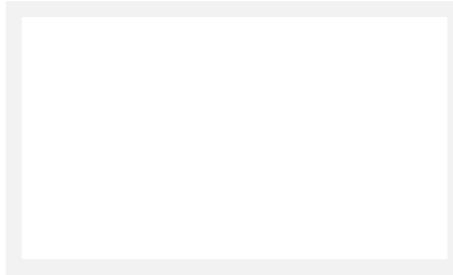
**ICAR & Coarray ICAR**



**Results**



**Conclusions**



# Overview



## Fortran 2018 in a Nutshell

- Motivation: exascale challenges
- Background: SPMD & PGAS in Fortran 2018
- Beyond CAF



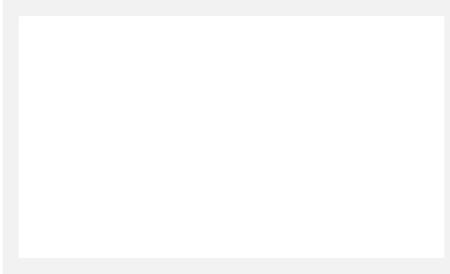
## ICAR & Coarray ICAR



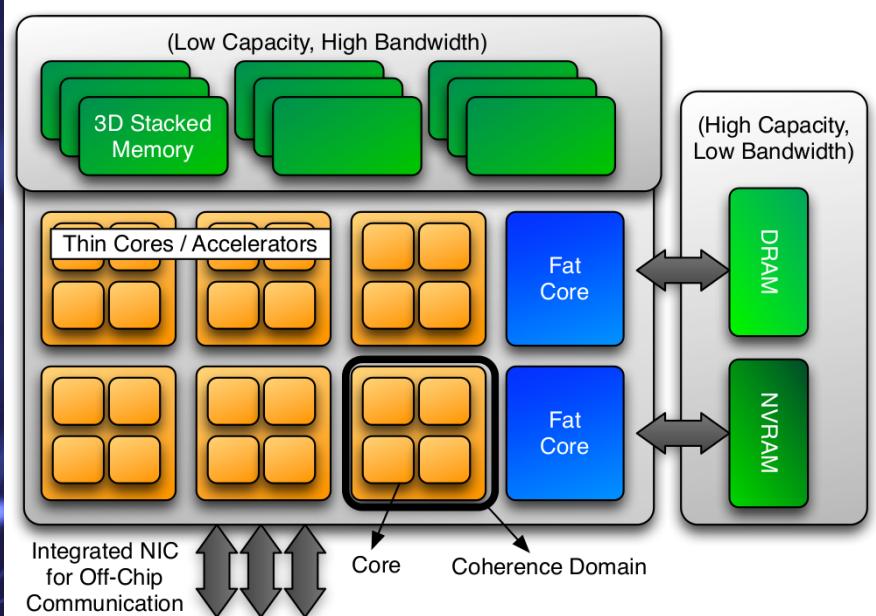
## Results



## Conclusions



# EXASCALE CHALLENGES & Fortran 2018 Response



Source: Ang, J.A., Barrett, R.F., Benner, R.E., Burke, D., Chan, C., Cook, J., Donofrio, D., Hammond, S.D., Hemmert, K.S., Kelly, S.M. and Le, H., 2014, November. Abstract machine models and proxy architectures for exascale computing. In *Hardware-Software Co-Design for High Performance Computing (Co-HPC)*, 2014 (pp. 25-32). IEEE.



## Billion-way concurrency with high levels of on-chip parallelism

- events
- collective subroutines
- richer set of atomic subroutines
- teams



## Higher failure rates

- failed-image detection



## Expensive data movement

- one-sided communication
- teams (locality control)



## Heterogeneous hardware on processor

- events

# Single Program Multiple Data (SPMD)

Images



# Single Program Multiple Data (SPMD)

Images →



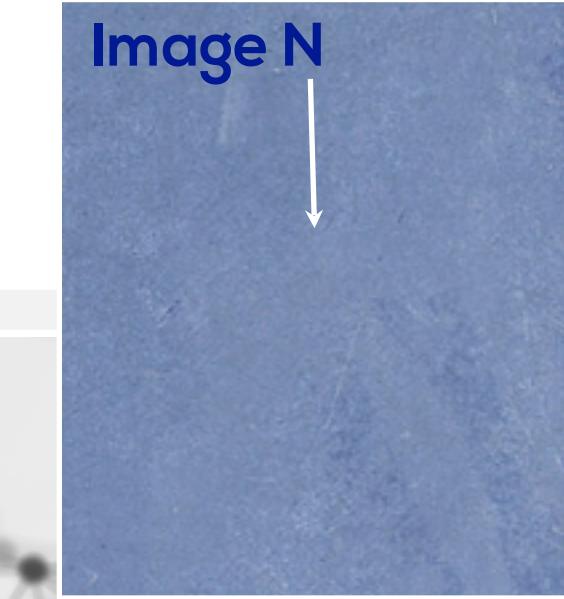
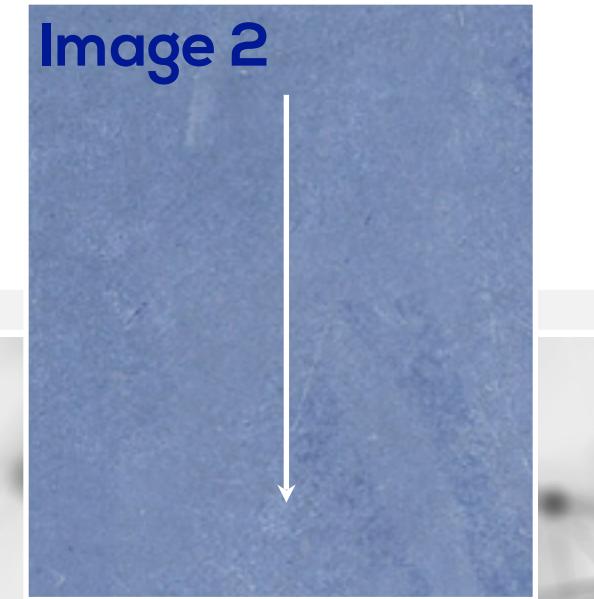
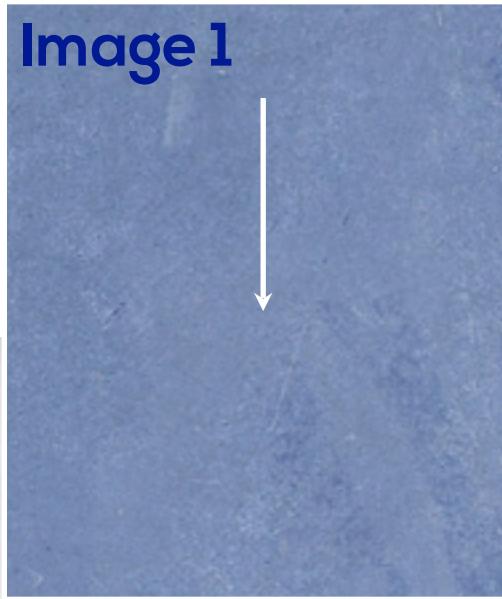
# Single Program Multiple Data (SPMD)

Images → {processes | threads}



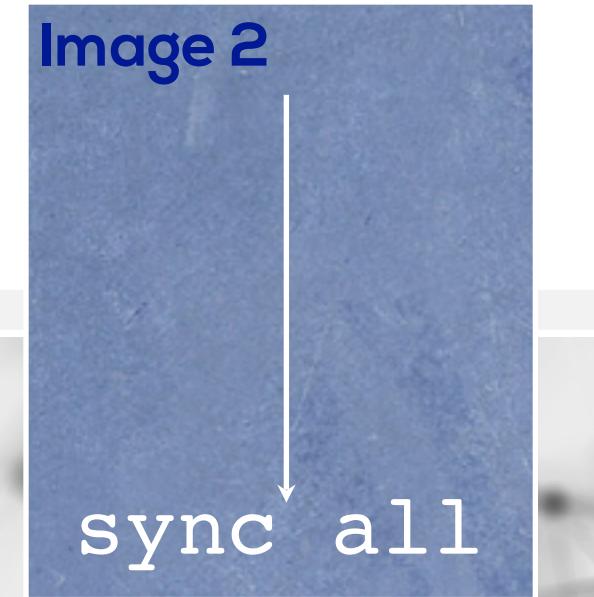
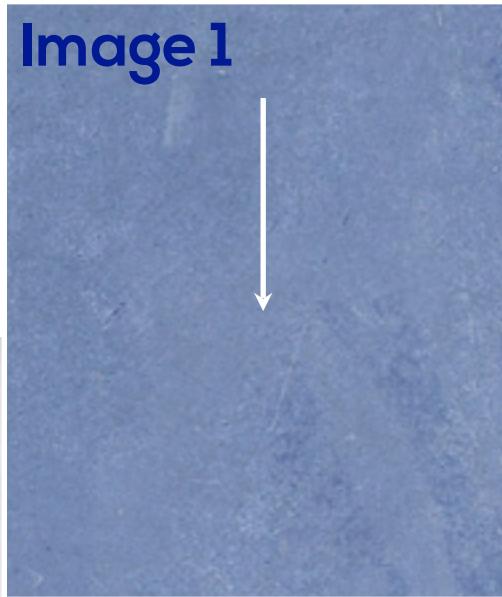
# Single Program Multiple Data (SPMD)

Images → {processes | threads}



# Single Program Multiple Data (SPMD)

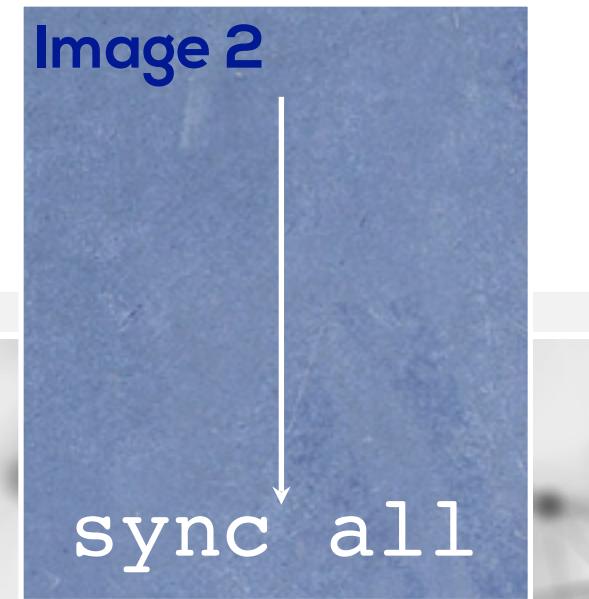
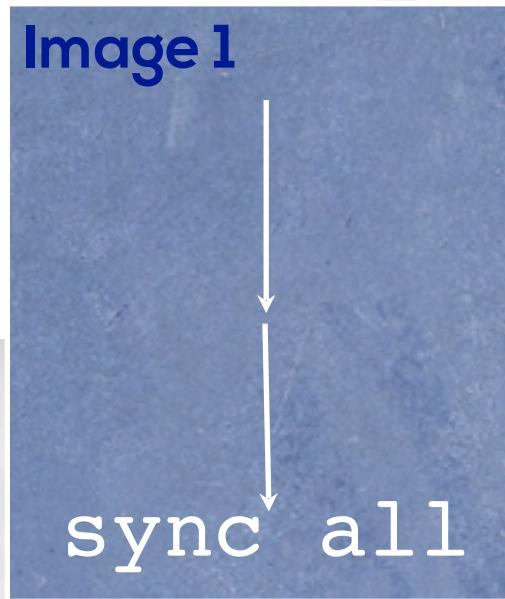
Images → {processes | threads}



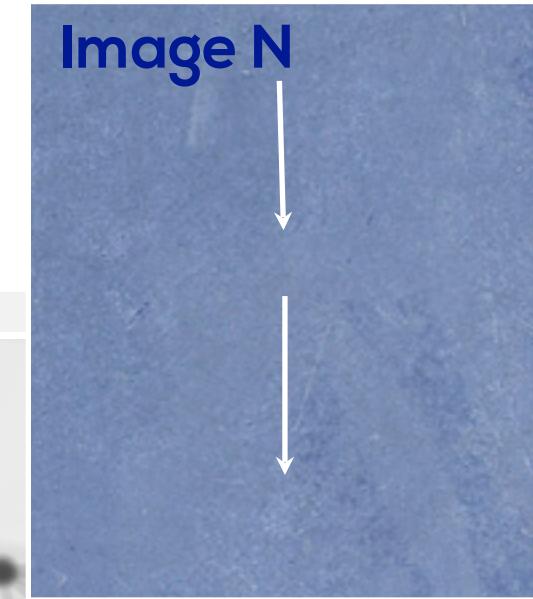
...

# Single Program Multiple Data (SPMD)

Images → {processes | threads}

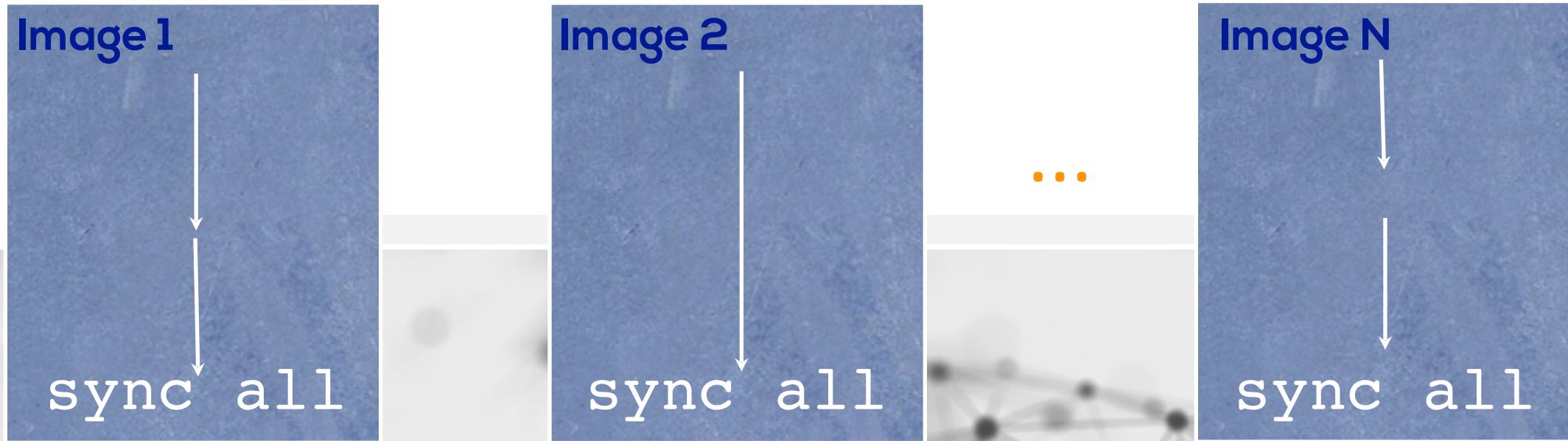


...



# Single Program Multiple Data (SPMD)

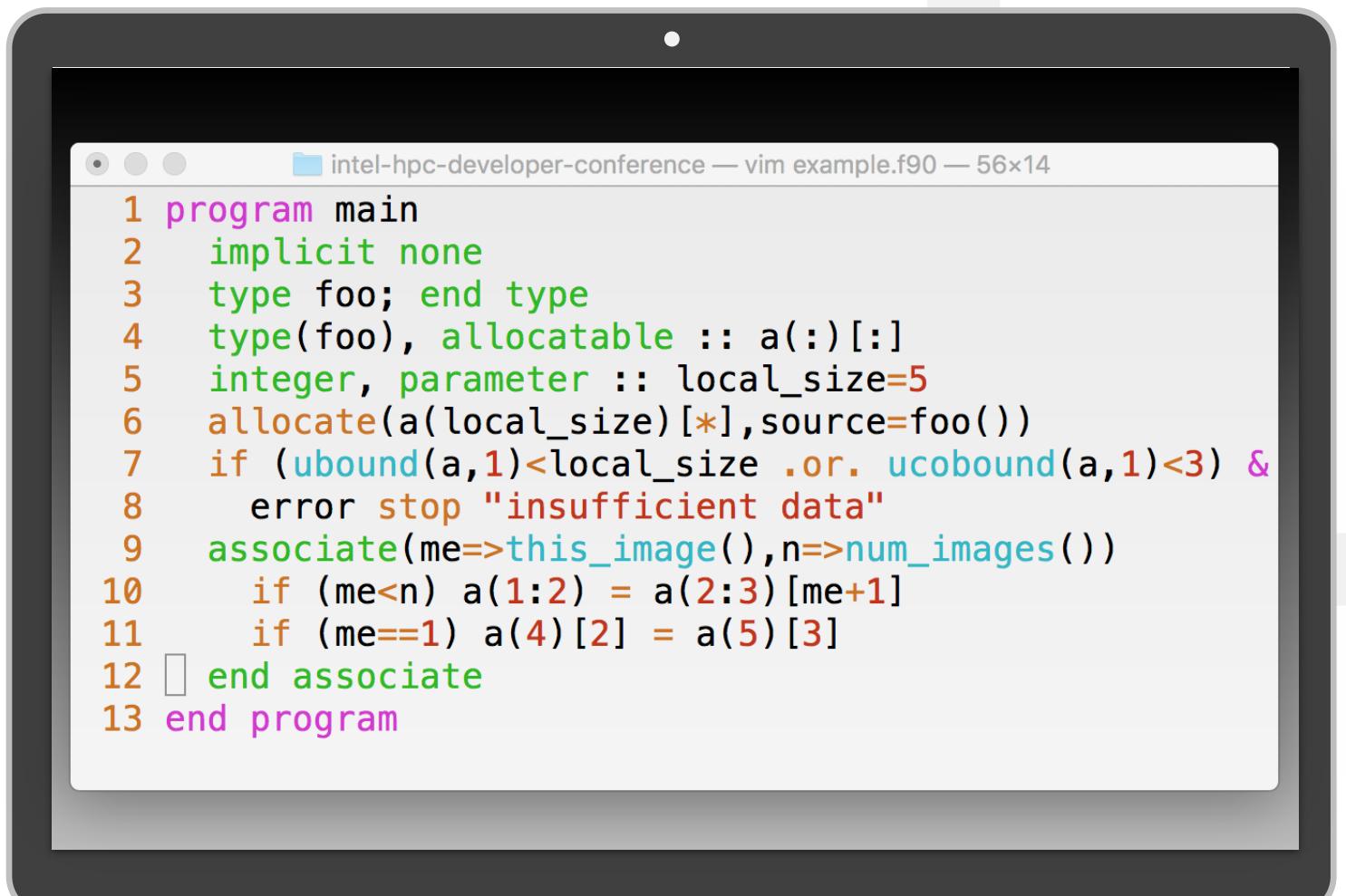
Images → {processes | threads}



Images execute asynchronously up to a programmer-specified synchronization:

sync all  
sync images  
allocate/deallocate

# Partitioned Global Address Space (PGAS)



A screenshot of a laptop screen showing a terminal window titled "intel-hpc-developer-conference — vim example.f90 — 56x14". The window displays the following Fortran code:

```
1 program main
2 implicit none
3 type foo; end type
4 type(foo), allocatable :: a(:)[:]
5 integer, parameter :: local_size=5
6 allocate(a(local_size)[*],source=foo())
7 if (ubound(a,1)<local_size .or. ucobound(a,1)<3) &
8   error stop "insufficient data"
9 associate(me=>this_image(),n=>num_images())
10  if (me<n) a(1:2) = a(2:3)[me+1]
11  if (me==1) a(4)[2] = a(5)[3]
12 end associate
13 end program
```

# Partitioned Global Address Space (PGAS)

Coarrays integrate with other languages features:



Fortran 90 array syntax works on local data.



Communicate objects

A screenshot of a terminal window titled "intel-hpc-developer-conference — vim example.f90 — 56x14". The window displays the following Fortran 90 code:

```
1 program main
2 implicit none
3 type foo; end type
4 type(foo), allocatable :: a(:,:)
5 integer, parameter :: local_size=5
6 allocate(a(local_size)*,source=foo())
7 if (ubound(a,1)<local_size .or. ucobound(a,1)<3) &
8   error stop "insufficient data"
9 associate(me=>this_image(),n=>num_images())
10  if (me<n) a(1:2) = a(2:3)[me+1]
11  if (me==1) a(4)[2] = a(5)[3]
12 end associate
13 end program
```

The code demonstrates the use of coarrays (array objects) and communication objects (like the "associate" statement) within a Fortran 90 program.

# Partitioned Global Address Space (PGAS)

Coarrays integrate with other languages features:

- Fortran 90 array syntax works on local data.
- Communicate objects

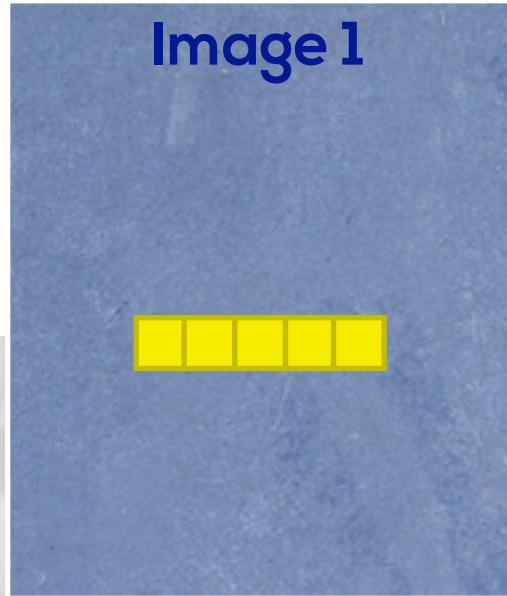
The ability to drop the square brackets harbors two important implications:

- Easily determine where communication occurs.
- Parallelize legacy code with minimal revisions.

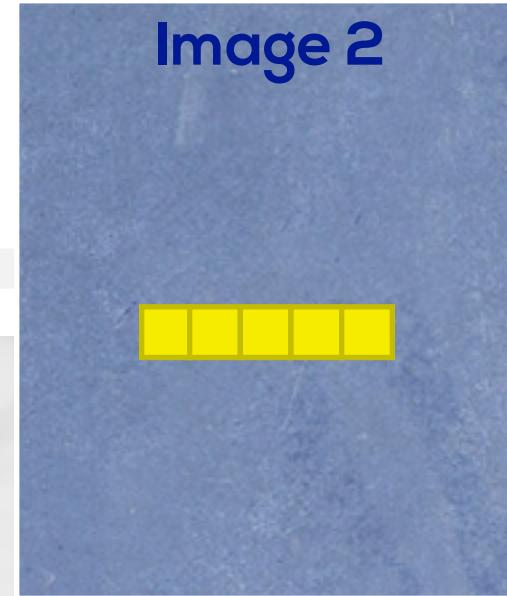
```
intel-hpc-developer-conference — vim example.f90 — 56x14
1 program main
2 implicit none
3 type foo; end type
4 type(foo), allocatable :: a(:,:)
5 integer, parameter :: local_size=5
6 allocate(a(local_size)*[],source=foo())
7 if (ubound(a,1)<local_size .or. ucobound(a,1)<3) &
8   error stop "insufficient data"
9 associate(me=>this_image(),n=>num_images())
10  if (me<n) a(1:2) = a(2:3)[me+1]
11  if (me==1) a(4)[2] = a(5)[3]
12 end associate
13 end program
```

`if (me<n) a(1:2) = a(2:3)[me+1]`

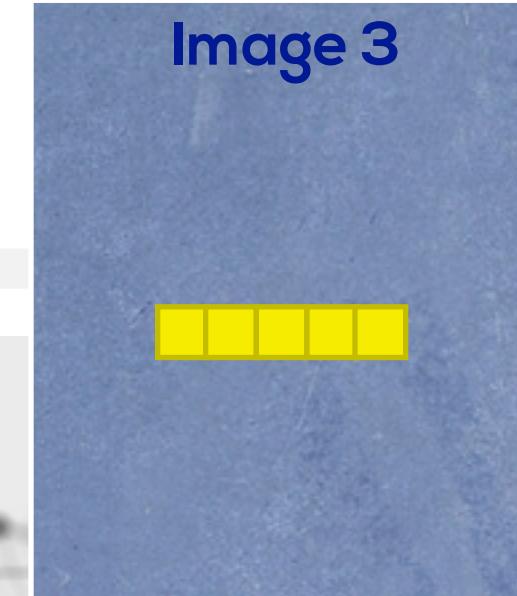
**Image 1**



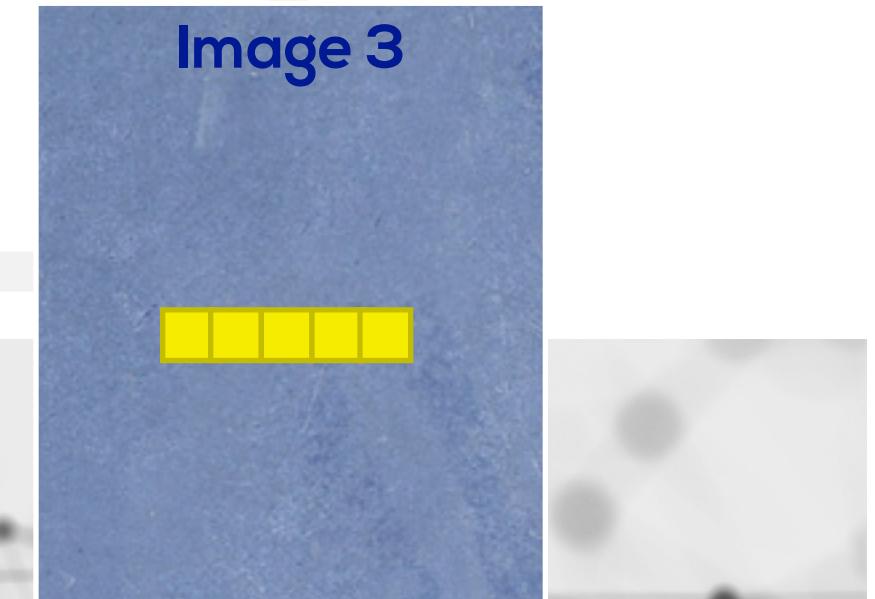
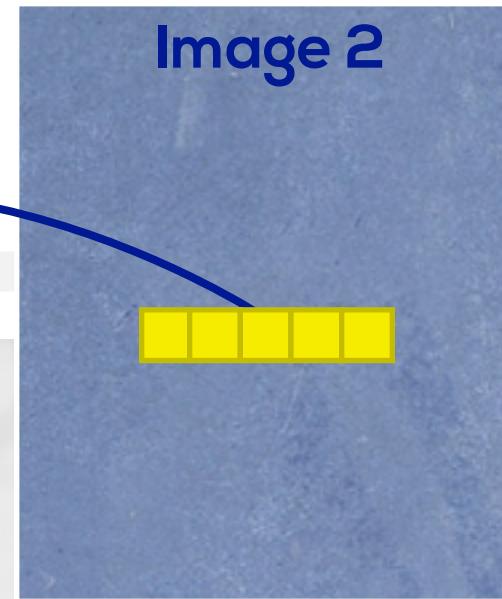
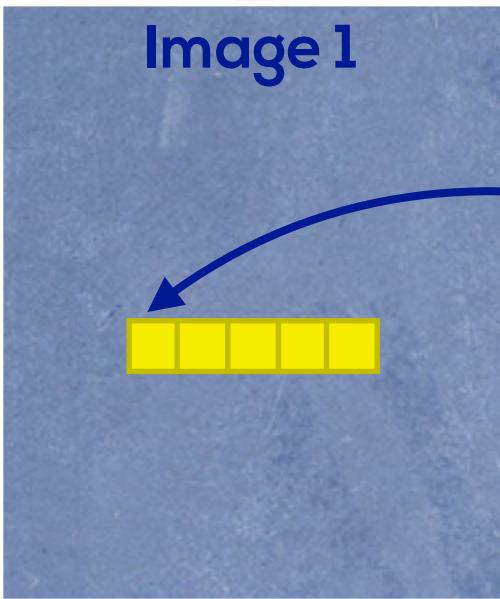
**Image 2**



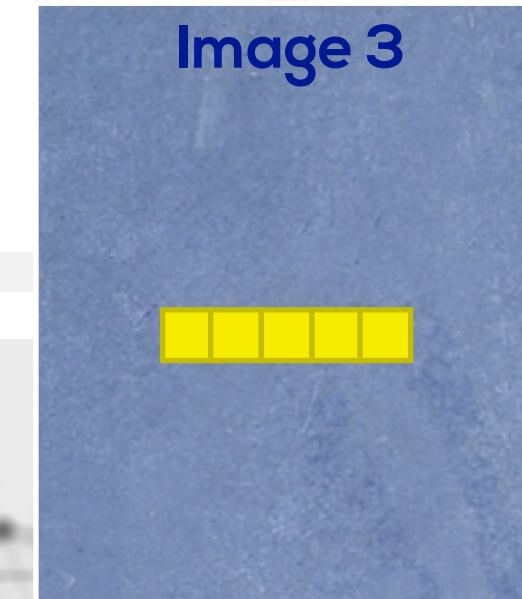
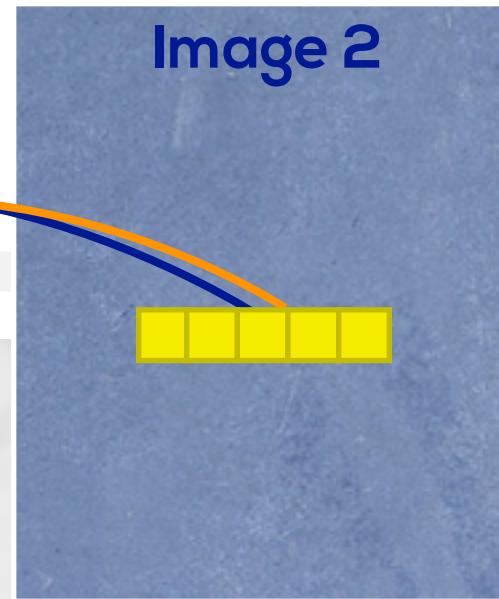
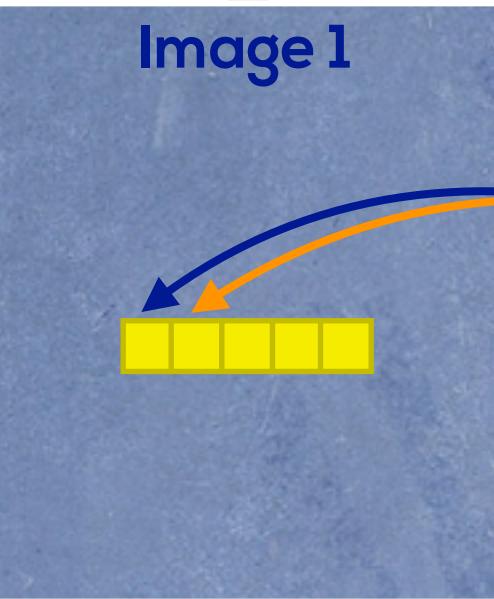
**Image 3**



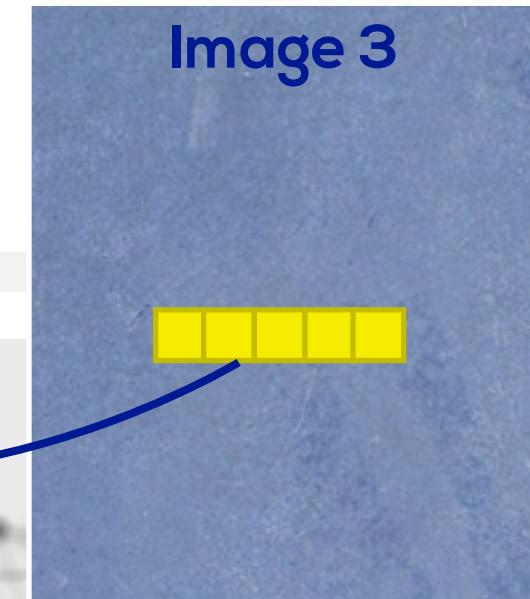
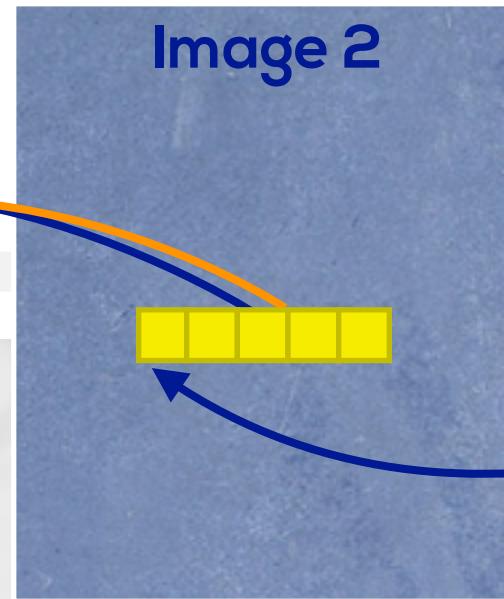
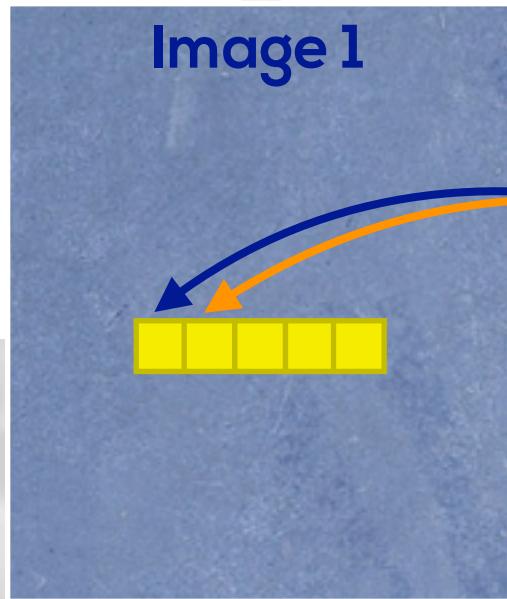
`if (me<n) a(1:2) = a(2:3)[me+1]`



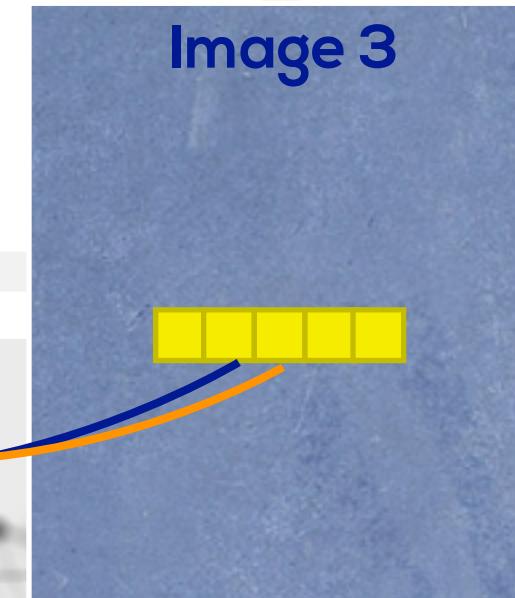
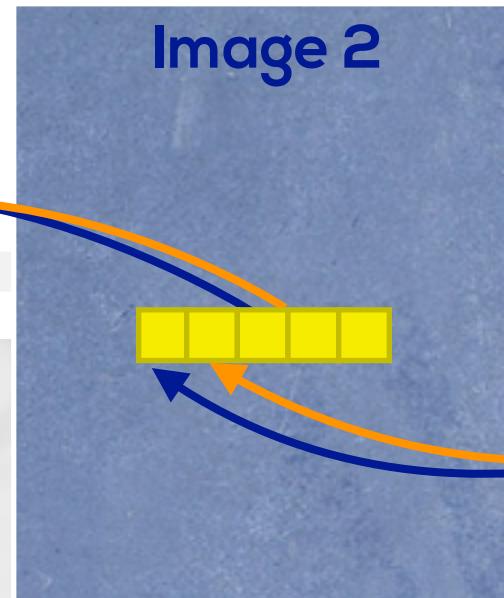
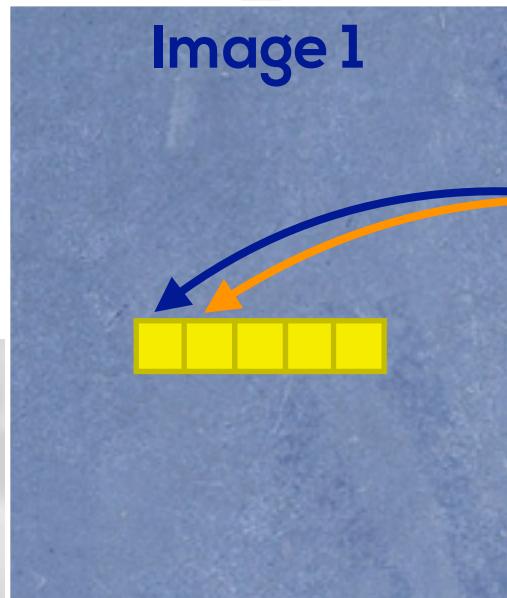
`if (me < n) a(1:2) = a(2:3) [me + 1]`



`if (me<n) a(1:2) = a(2:3)[me+1]`

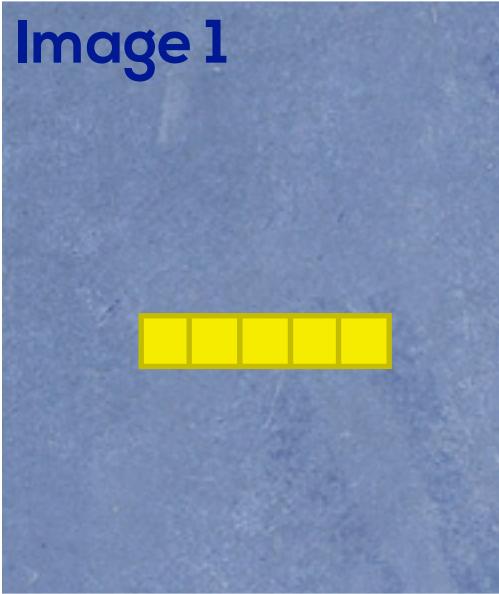


`if (me < n) a(1:2) = a(2:3) [me+1]`

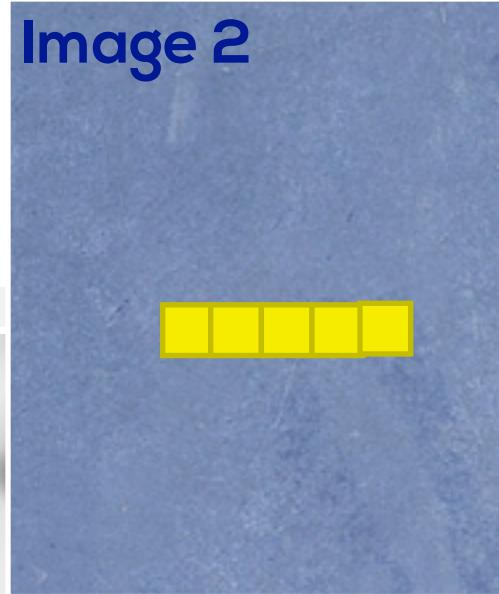


**if (me==1) a(4)[2] = a(5)[3]**

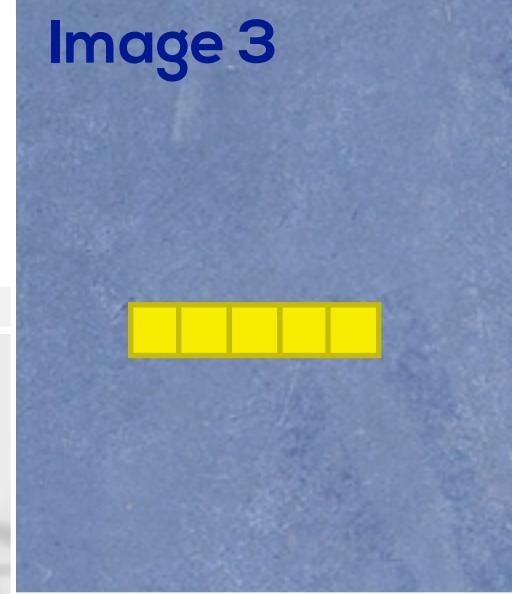
**Image 1**



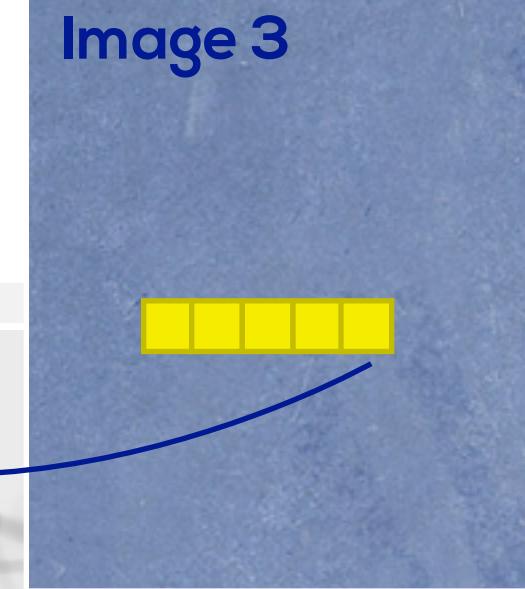
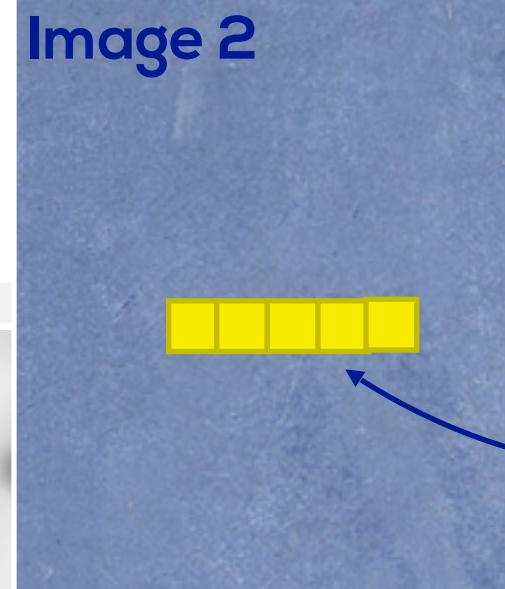
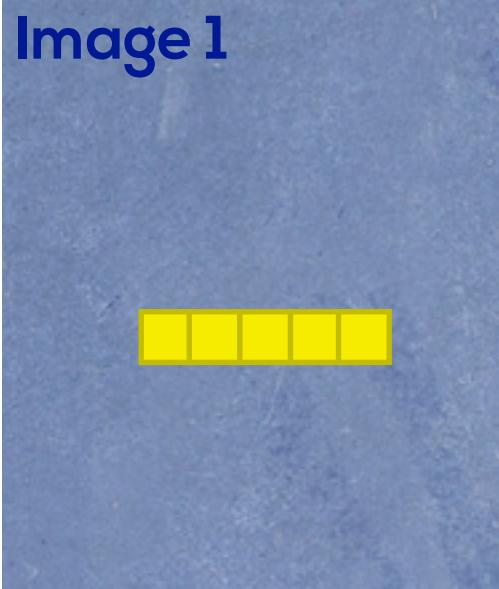
**Image 2**



**Image 3**



**if (me==1) a(4)[2] = a(5)[3]**



# Segment Ordering:

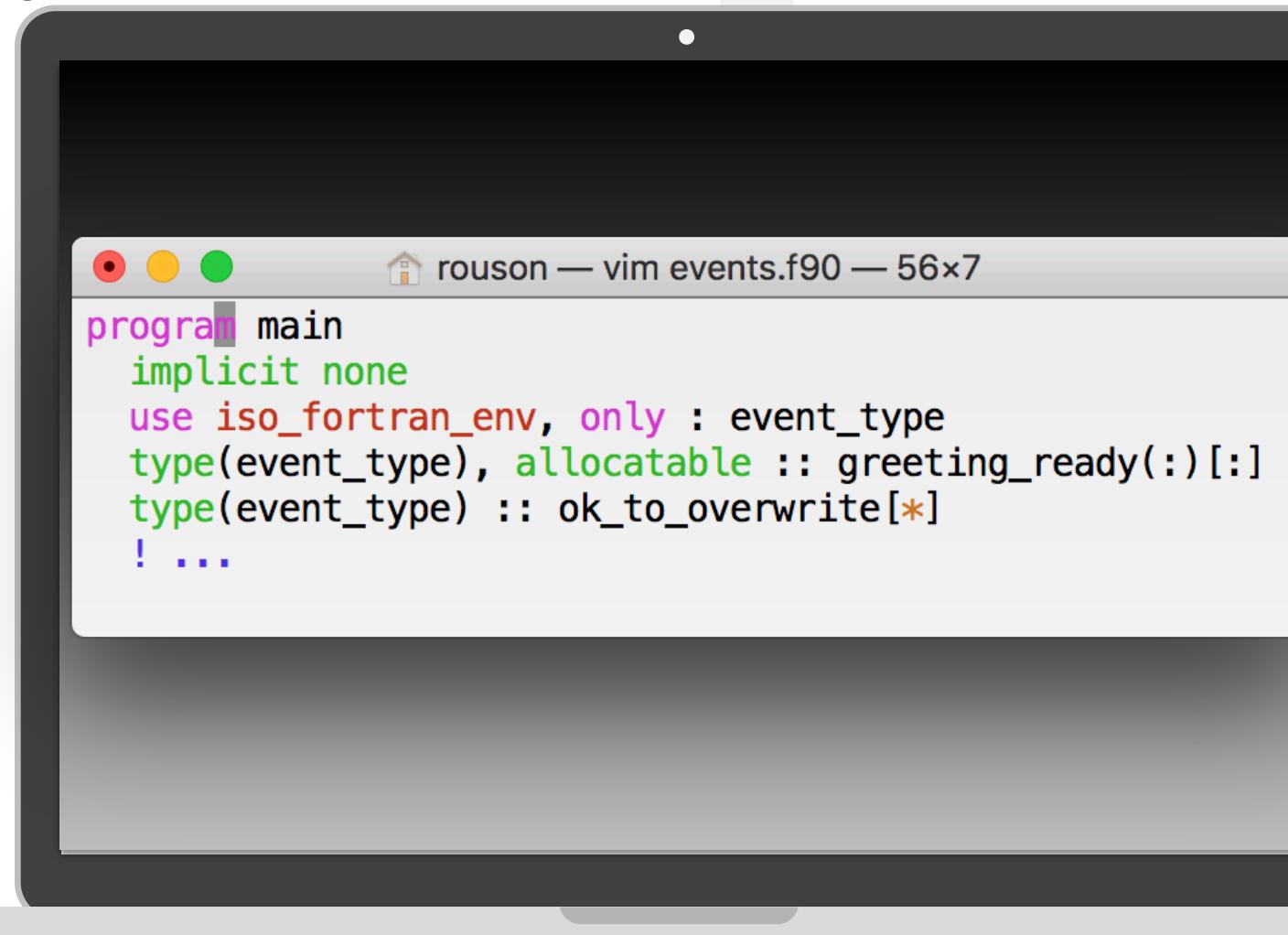
## Events

An intrinsic module provides the derived type `event_type`, which encapsulates an `atomic_int_kind` integer component default-initialized to zero.

 An image increments the event count on a remote image by executing `event_post`.

 The remote image obtains the post count by executing `event_query`.

	Image Control	Side Effect
<code>event post</code>	<input checked="" type="checkbox"/>	<code>atomic_add 1</code>
<code>event_query</code>		defines count
<code>event wait</code>	<input checked="" type="checkbox"/>	<code>atomic_add -1</code>

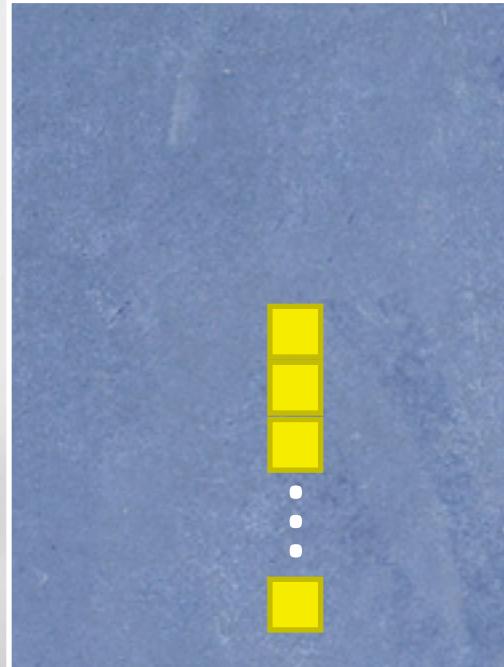


A screenshot of a laptop screen displaying a Vim editor window. The title bar shows "rouson — vim events.f90 — 56x7". The code in the editor is:

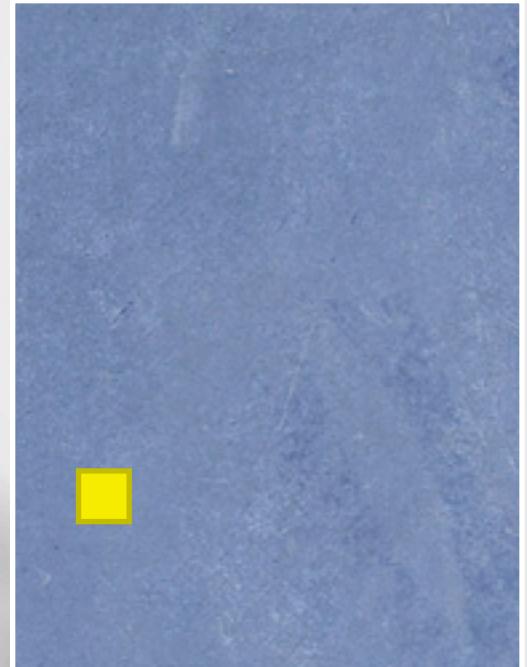
```
program main
    implicit none
    use iso_fortran_env, only : event_type
    type(event_type), allocatable :: greeting_ready(:)[:]
    type(event_type) :: ok_to_overwrite[*]
!
```

# Events

Hello, world!



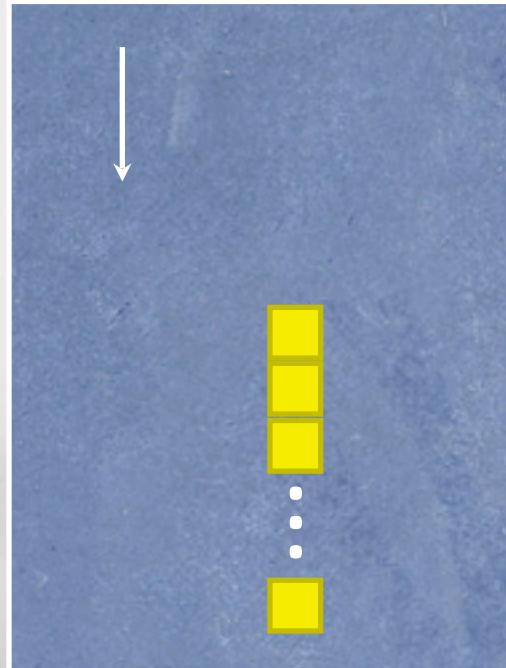
greeting\_ready(2:n) [1]



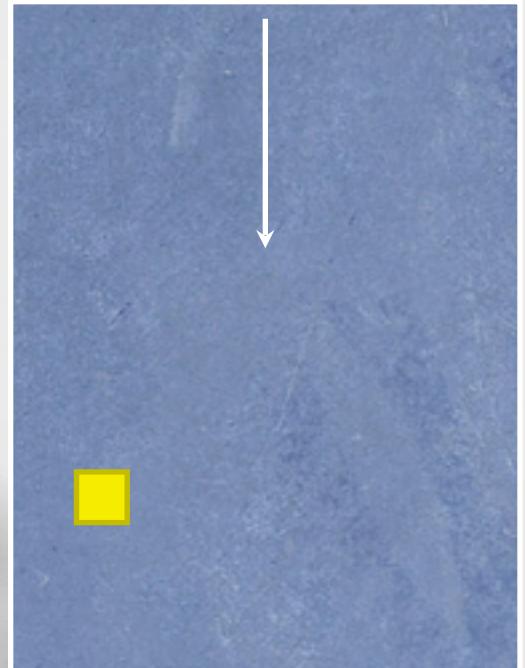
ok\_to\_overwrite[3]

# Events

Hello, world!



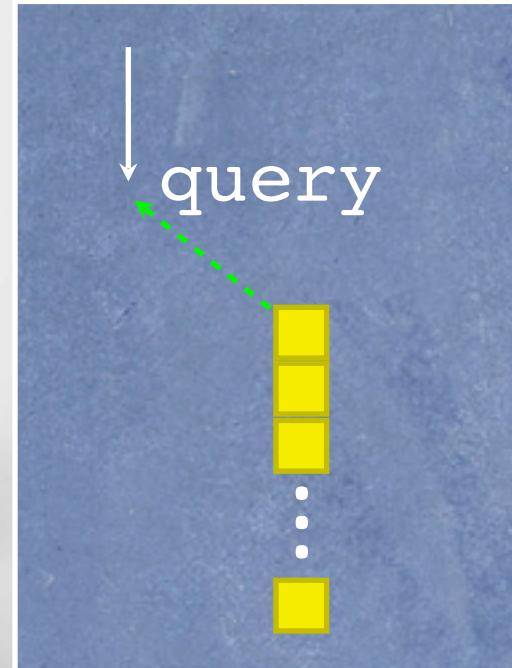
greeting\_ready(2:n) [1]



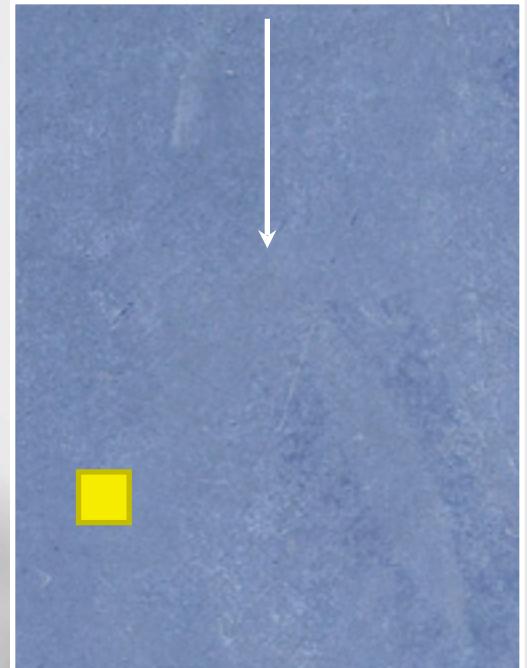
ok\_to\_overwrite[3]

# Events

Hello, world!



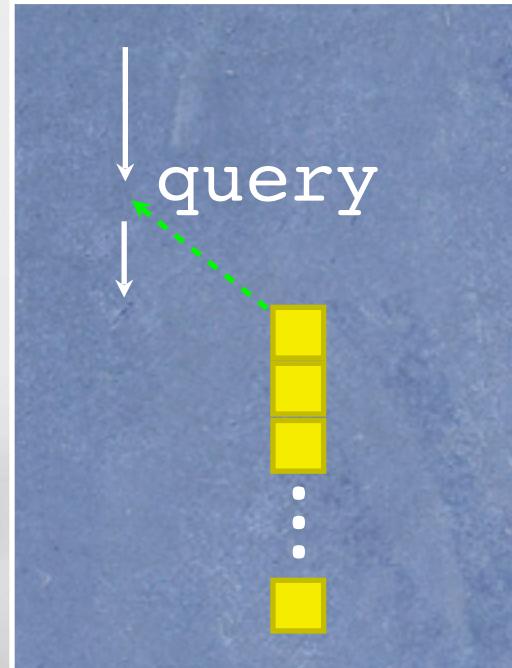
greeting\_ready(2:n) [1]



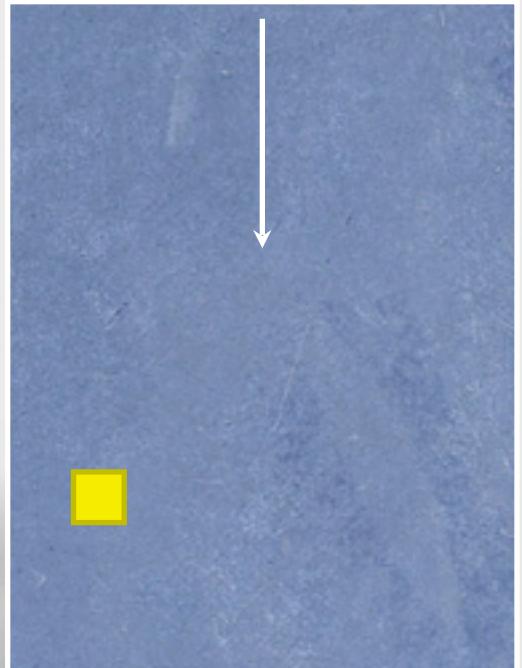
ok\_to\_overwrite[3]

# Events

Hello, world!



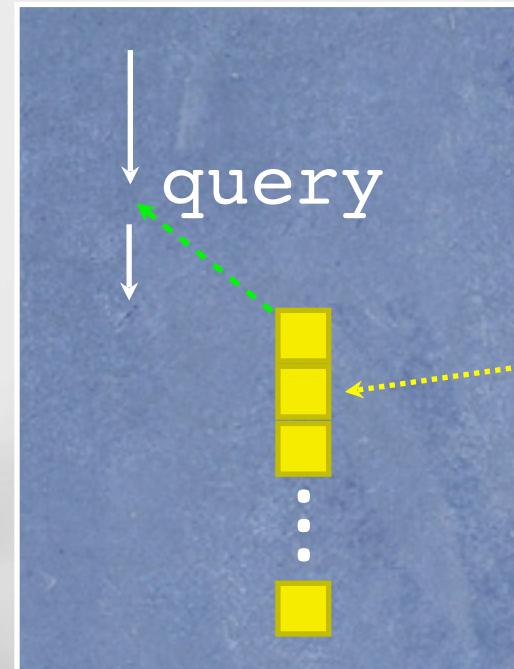
greeting\_ready(2:n) [1]



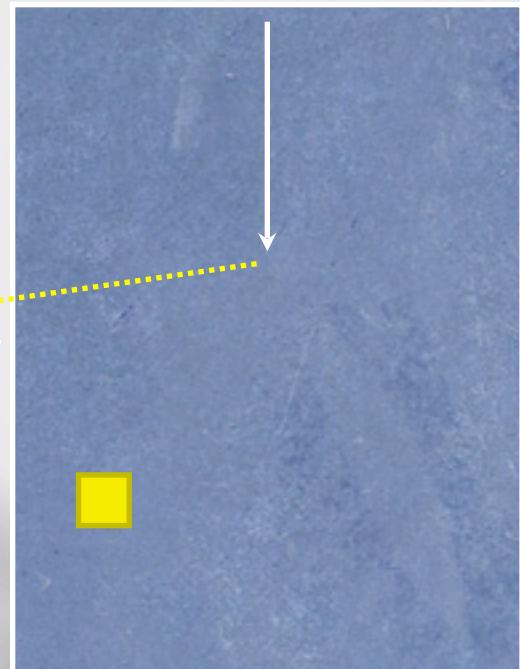
ok\_to\_overwrite[3]

# Events

Hello, world!



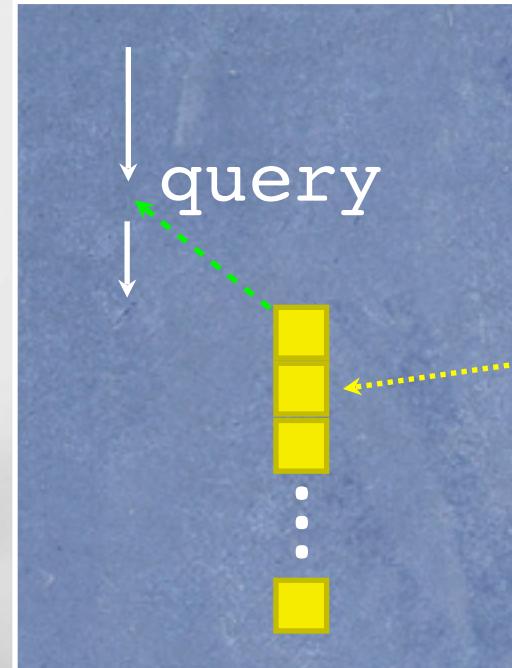
greeting\_ready(2:n) [1]



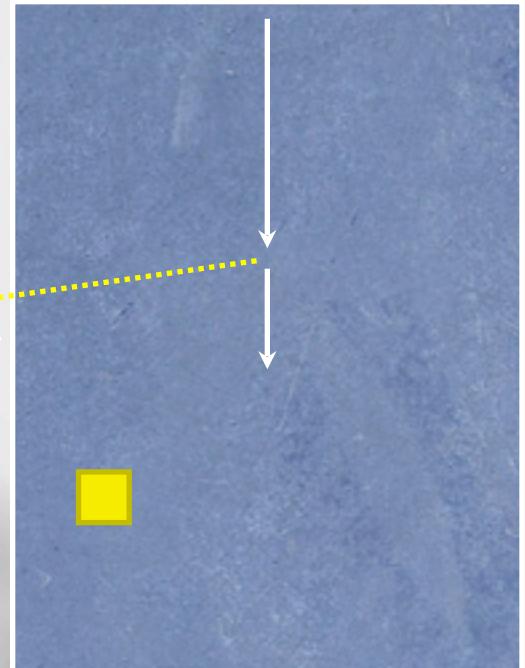
ok\_to\_overwrite[3]

# Events

Hello, world!



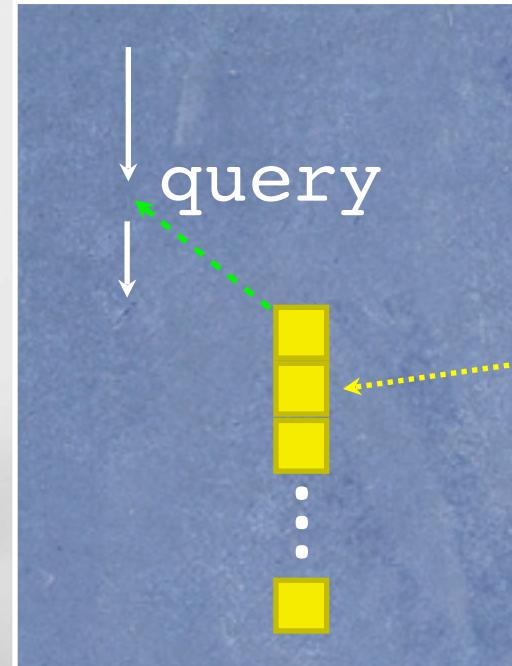
greeting\_ready(2:n) [1]



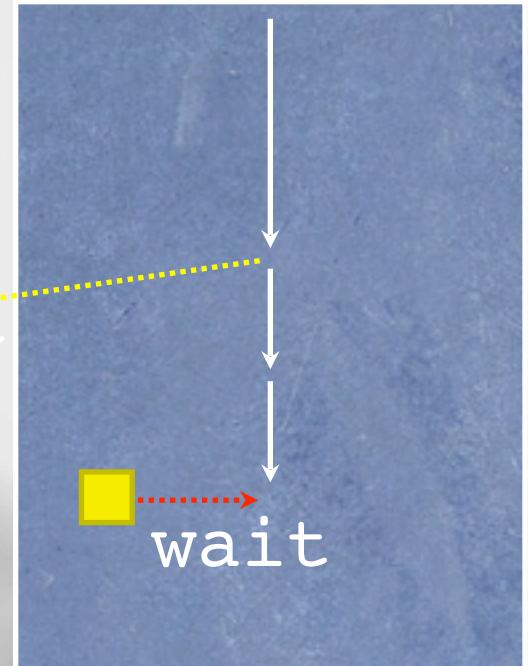
ok\_to\_overwrite[3]

# Events

Hello, world!



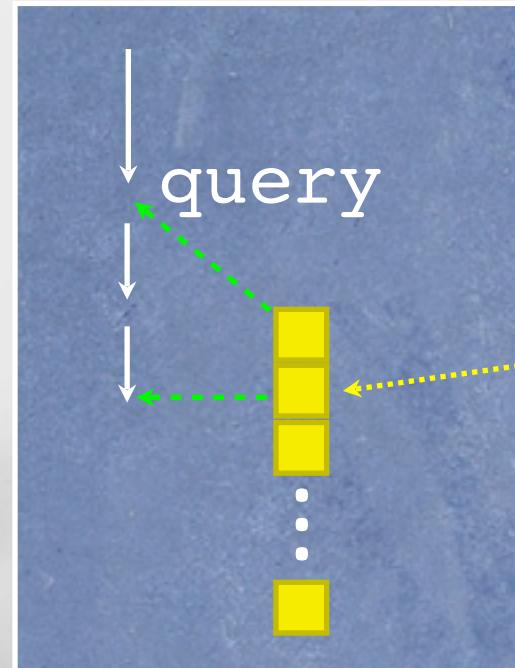
greeting\_ready(2:n) [1]



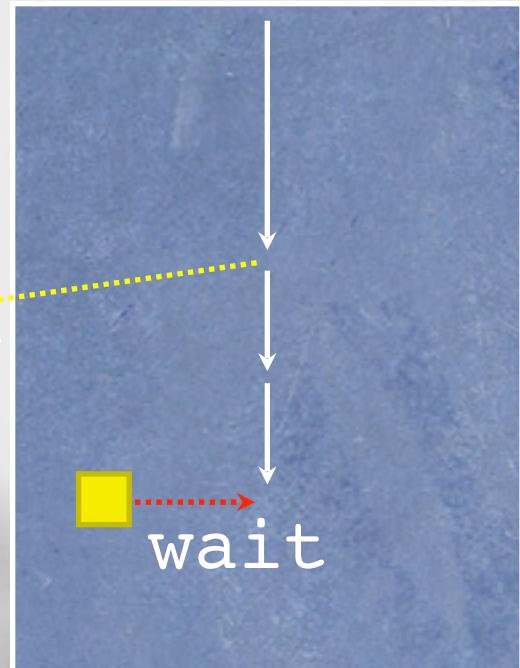
ok\_to\_overwrite[3]

# Events

Hello, world!



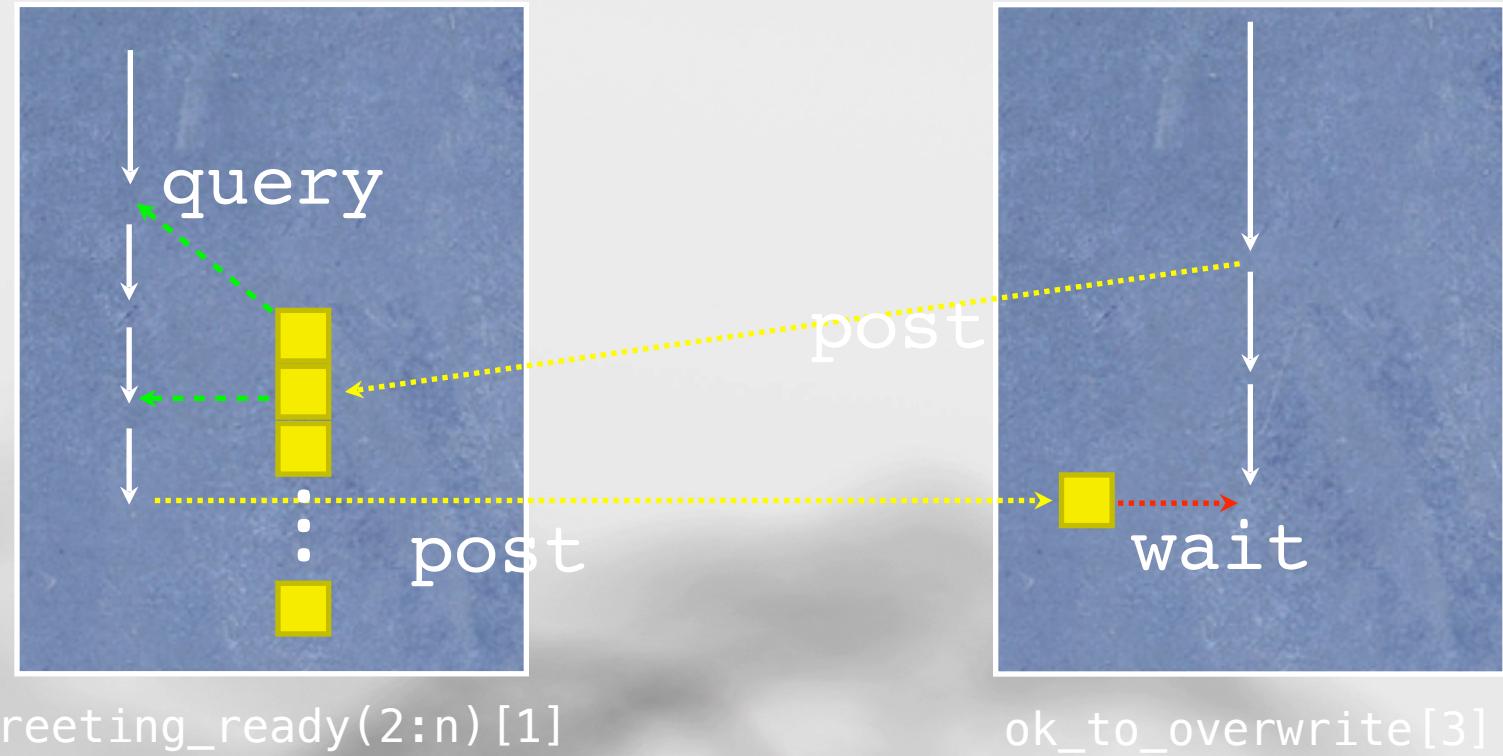
greeting\_ready(2:n) [1]



ok\_to\_overwrite[3]

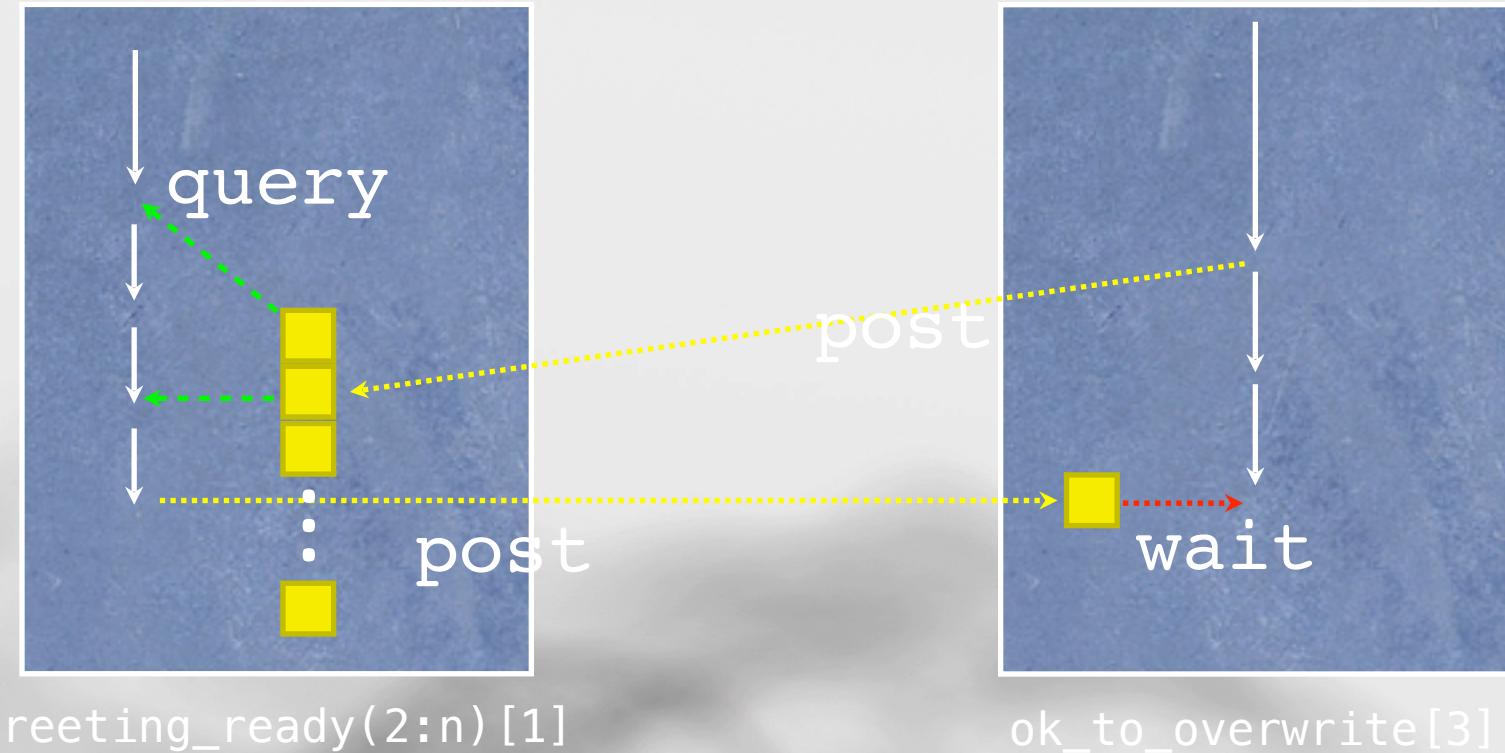
# Events

Hello, world!



# Events

## Hello, world!



greeting\_ready(2:n) [1]

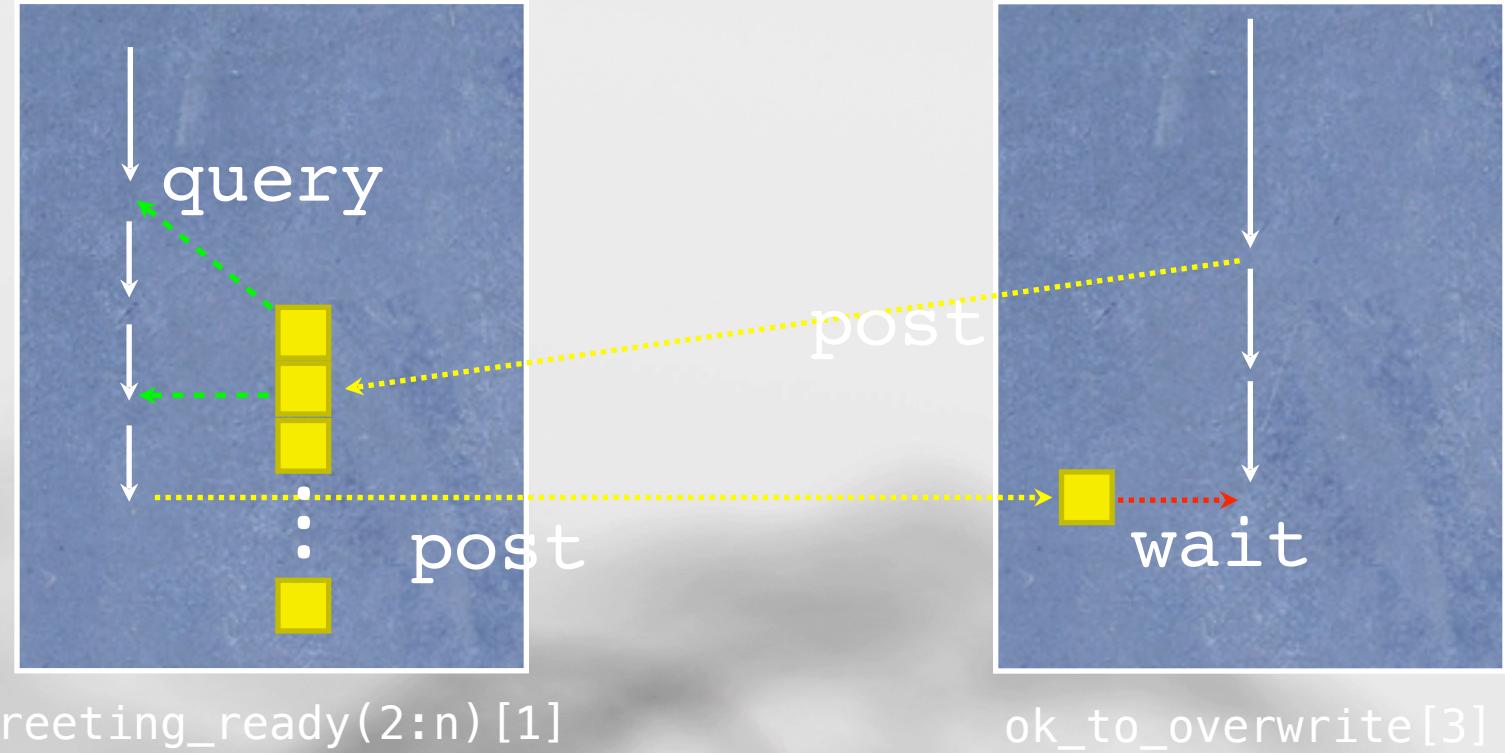
ok\_to\_overwrite[3]

Performance-oriented constraints:

- Query and wait must be local.
- Post and wait are disallowed in do concurrent constructs.

# Events

## Hello, world!



Performance-oriented constraints:

- Query and wait must be local.
- Post and wait are disallowed in do concurrent constructs.

Pro tips:

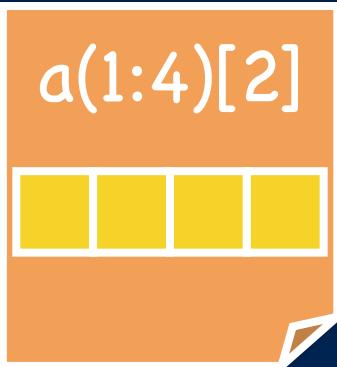
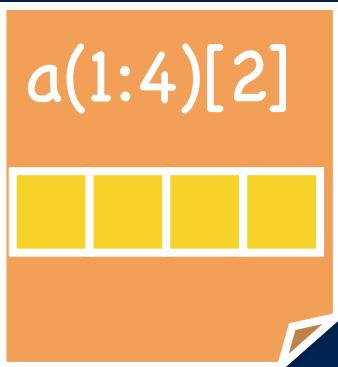
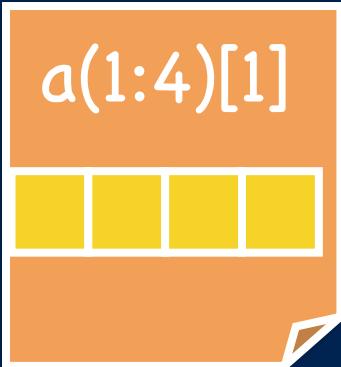
- Overlap communication and computation.
- Wherever safety permits, query without waiting.
- Write a spin-query-work loop & build a logical mask describing the remaining work.

# TEAM

A set of images that can readily execute independently of other images

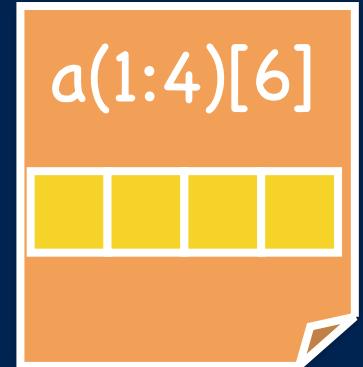
**Team 1**

Image 1    Image 2    Image 3



**Team 2**

Image 4    Image 5    Image 6



# Collective Subroutines



Each non-failed image of the current team must invoke the collective.



After parallel calculation/communication, the result is placed on one or all images.



Optional arguments: `stat`, `errmsg`, `result_image`, `source_image`



All collectives have `intent(inout)` argument `A` holding the input data and may hold the result on return, depending on `result_image`



No implicit synchronization at beginning/end, which allows for overlap with other actions.



No image's data is accessed before that image invokes the collective subroutine.



# Extensible Set of Collective Subroutines

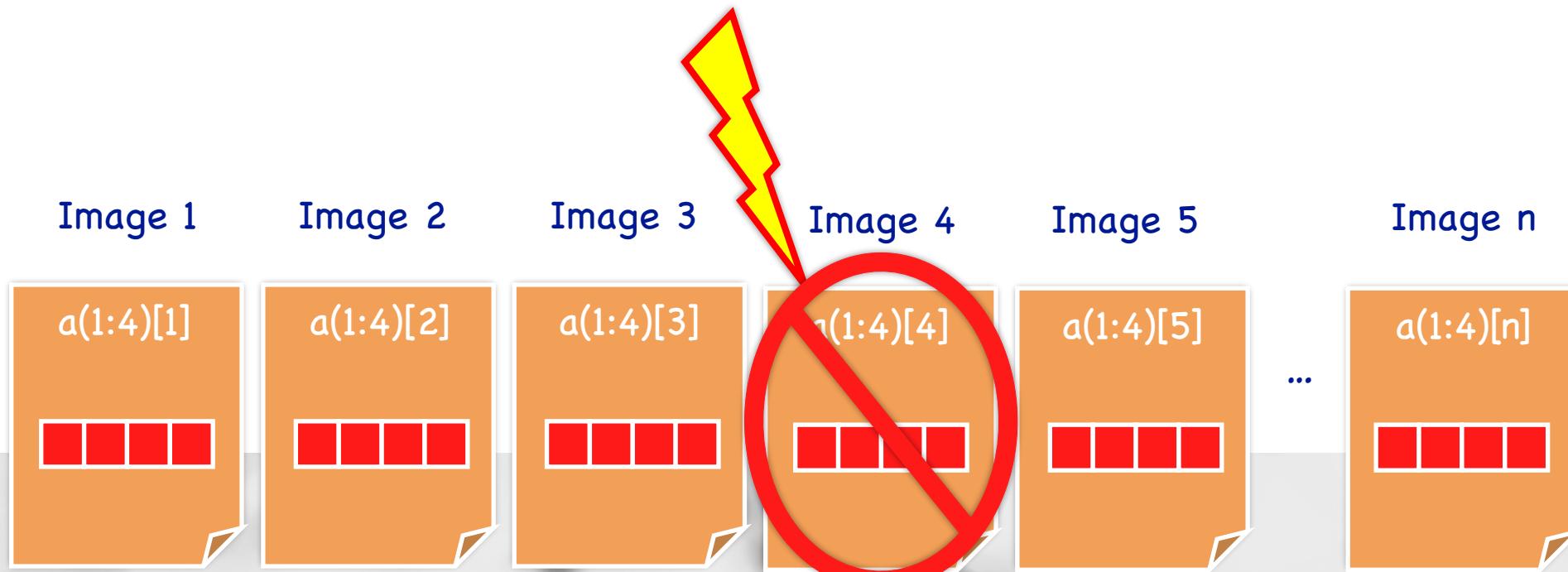
```
intel-hpc-developer-conference — vim collective-subroutines.f90...
call co_broadcast (a, source_image [, stat, errmsg])
call co_max (a [, result_image, stat, errmsg])
call co_min (a [, result_image, stat, errmsg])
call co_sum (a [, result_image, stat, errmsg])
call co_reduce (a, operation [, result_image, stat, errmsg])
~
```

# Failure Detection



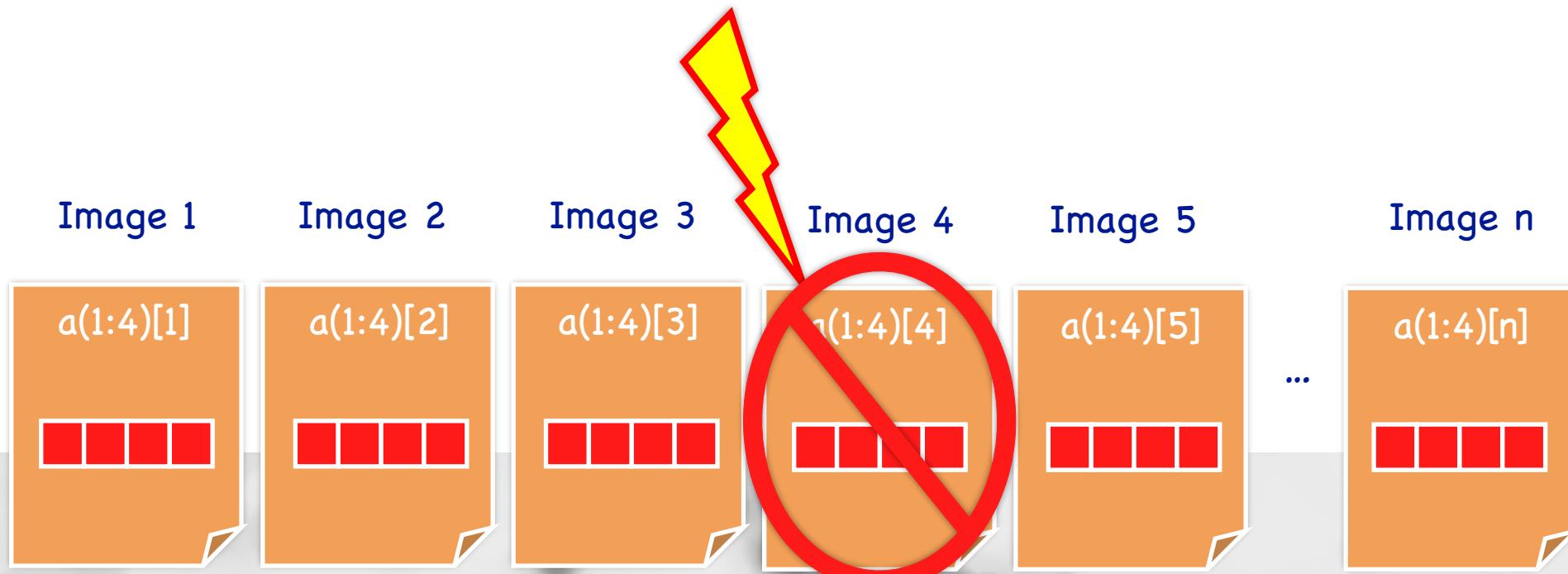
```
use iso_fortran_env, only : STAT_FAILED_IMAGE
integer :: status
sync all(stat==status)
if (status==STAT_FAILED_IMAGE) call fault_tolerant_algorithm()
```

# Failure Detection



```
use iso_fortran_env, only : STAT_FAILED_IMAGE
integer :: status
sync all(stat==status)
if (status==STAT_FAILED_IMAGE) call fault_tolerant_algorithm()
```

# Failure Detection



# FORTRAN 2018

## Failed-Image Detection

- FAIL IMAGE (simulates a failures)
- IMAGE\_STATUS (checks the status of a specific image)
- FAILED\_IMAGES (provides the list of failed images)
- Coarray operations may fail. STAT= attribute used to check correct behavior.

# Overview



Fortran 2018 in a Nutshell



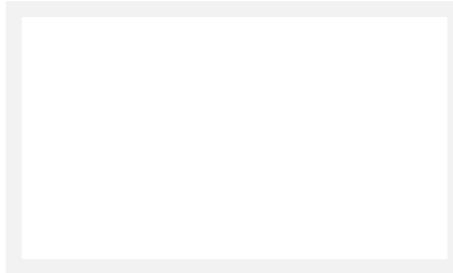
**ICAR & Coarray ICAR**



Results

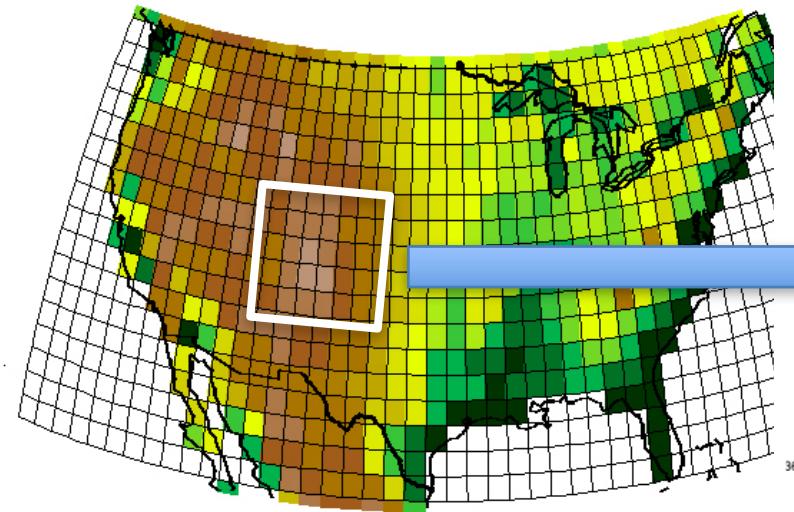


Conclusions

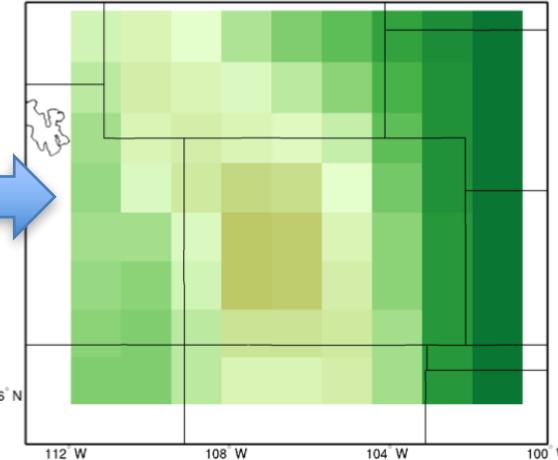


# The Climate Downscaling Problem

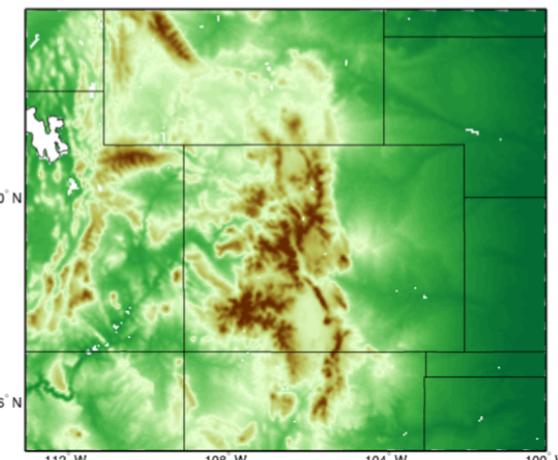
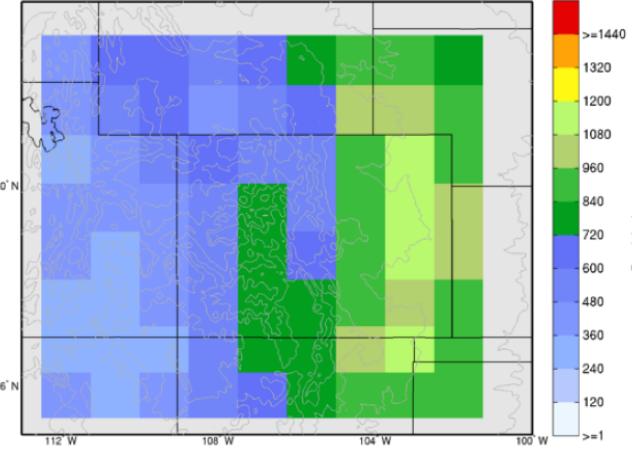
Climate model (top row) and application needs (bottom row)



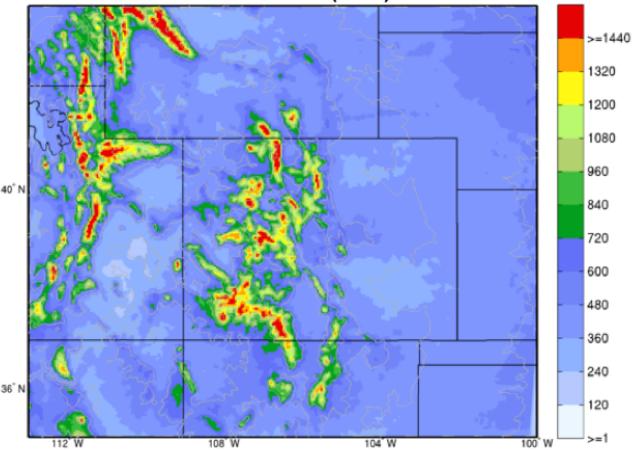
Topography



Precipitation



WRF Ens002 (4km)



## Computational Cost

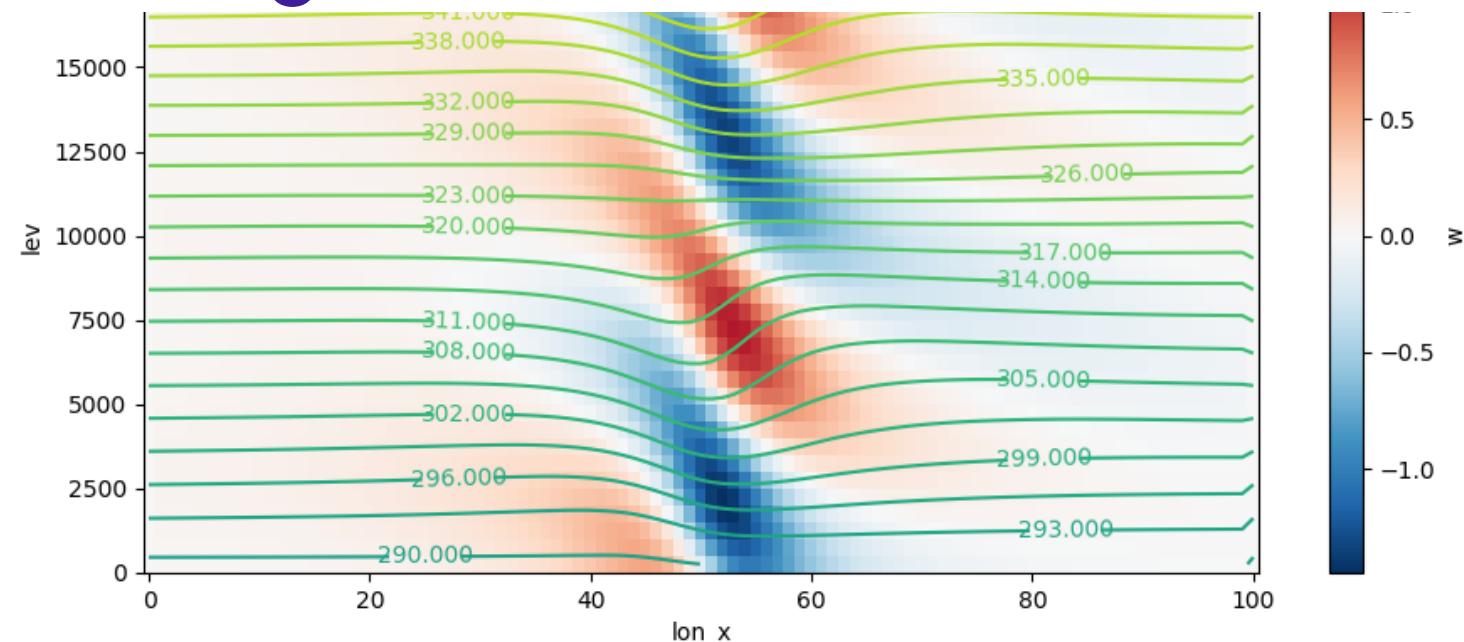
High-res regional model  
>10 billion core hours

(CONUS 4km 150yrs 40 scenarios)

# The Intermediate Complexity Atmospheric Research Model (ICAR)

- Analytical solution for flow over topography (right)
- 90% of the information for 1% of the computational cost
- Core Numerics:
  - 90% of cost = Cloud physics
    - grid-cells independent
  - 5% of cost = Advection
    - fully explicit, requires local neighbor communication

## ICAR Wind Field over Topography



# ICAR

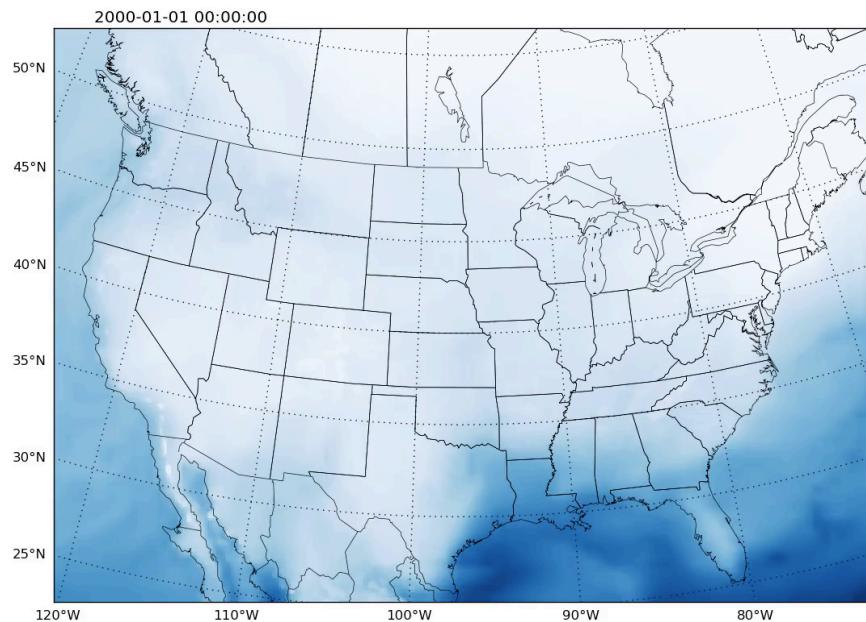
## Intermediate Complexity Atmospheric Research

Animation of  
Water vapor (blue) and  
Precipitation (Green to Red)  
over the contiguous United States

Output timestep : 1hr  
Integration timestep : Variable (~30s)  
Boundary Conditions: ERA-interim (historical)

Run on 16 cores with OpenMP

Limited scalability



# ICAR

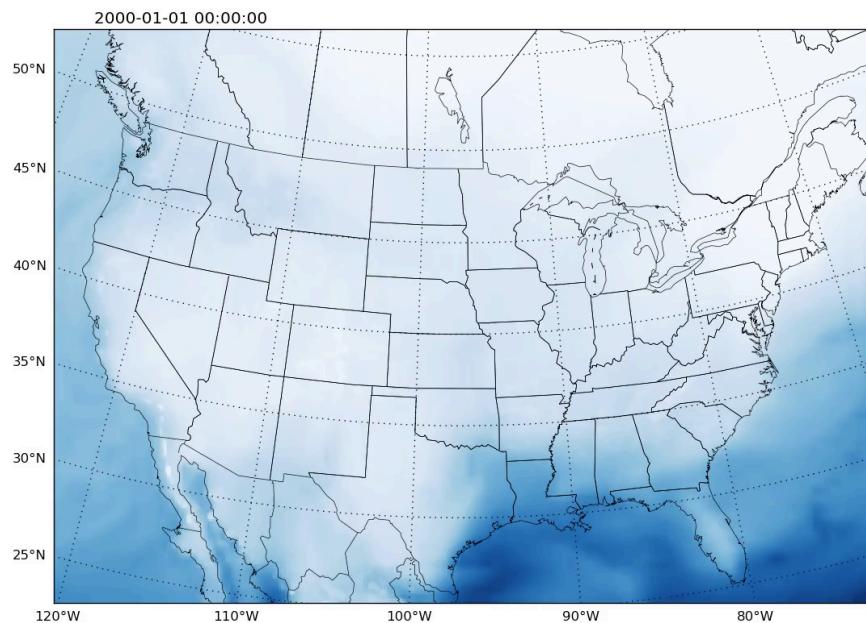
## Intermediate Complexity Atmospheric Research

Animation of  
Water vapor (blue) and  
Precipitation (Green to Red)  
over the contiguous United States

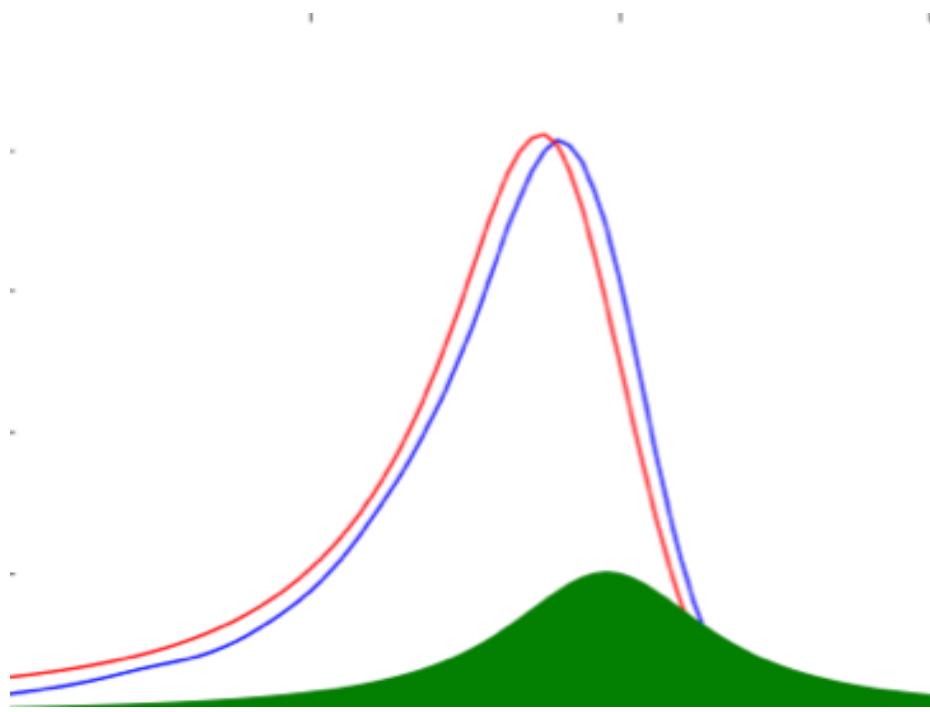
Output timestep : 1hr  
Integration timestep : Variable (~30s)  
Boundary Conditions: ERA-interim (historical)

Run on 16 cores with OpenMP

Limited scalability

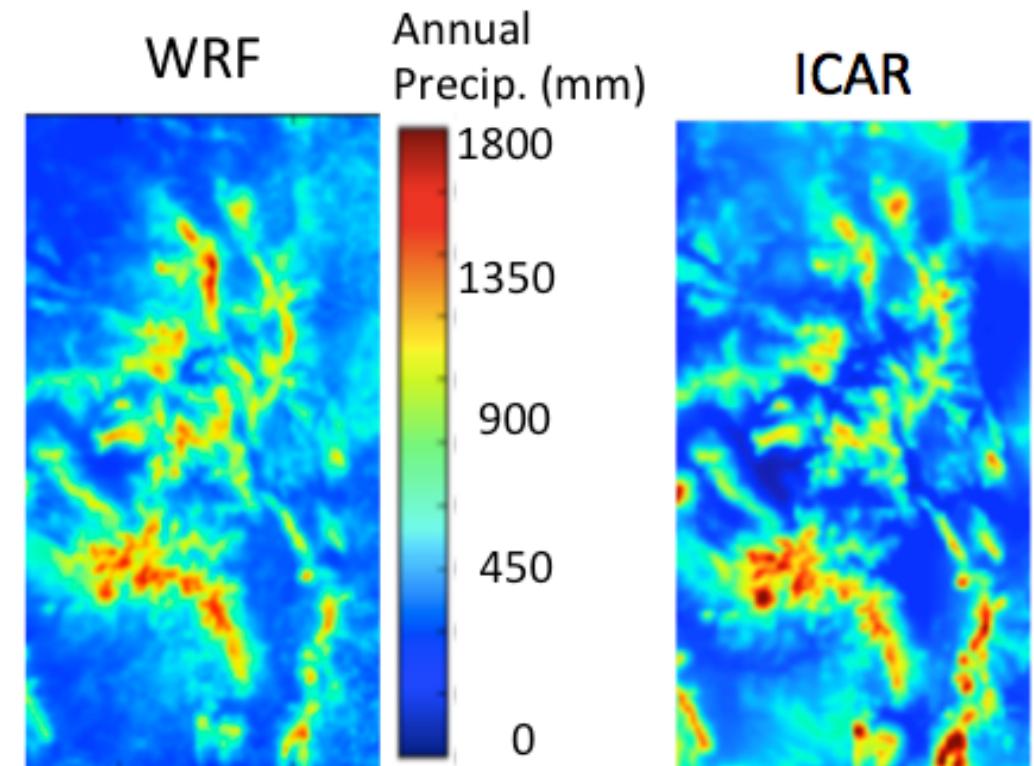


# ICAR comparison to “Full-physics” atmospheric model (WRF)



## Ideal Hill case

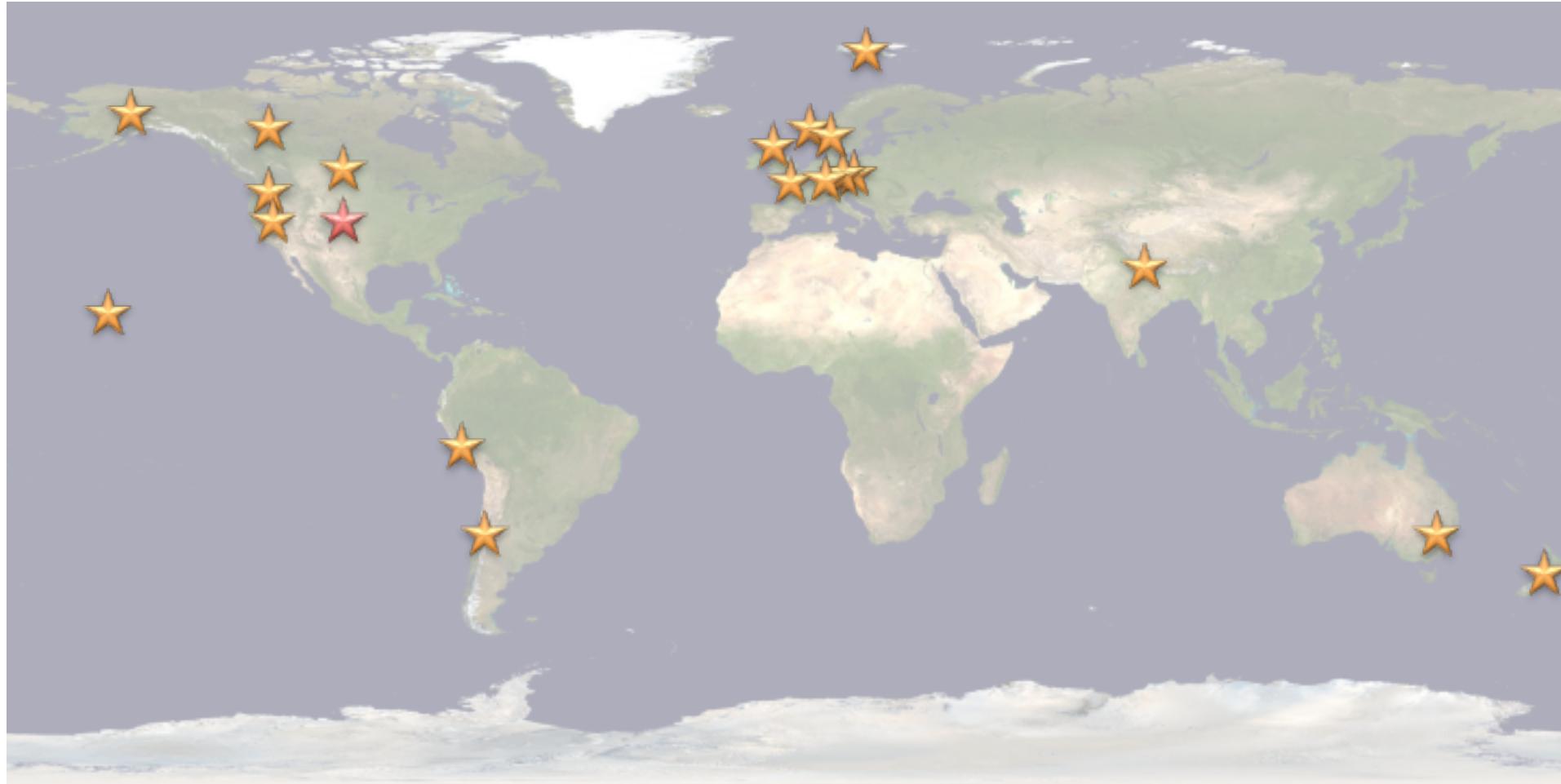
ICAR (red) and WRF (blue)  
precipitation over an idealized  
hill (green)



## Real Simulation

WRF (left) and ICAR (right)  
Season total precipitation over  
Colorado Rockies

# ICAR Users and Applications



<http://github.com/NCAR/icar.git>

# Coarray ICAR Mini-App



Object-oriented design



Collective broadcast of initial data



Overlaps communication & computation via  
one-sided “puts.”



Coarray halo exchanges with 3D, 1st-order upwind  
advection (~2000 lines of new code)



Cloud microphysics (~5000 lines of pre-  
existing code)



# Overview



Fortran 2018 in a Nutshell



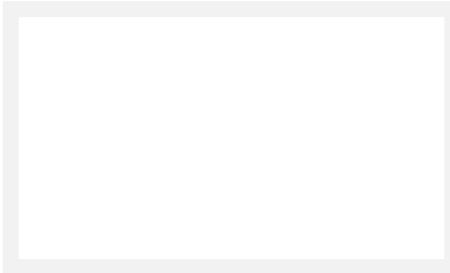
ICAR & Coarray ICAR



**Results**



Conclusions

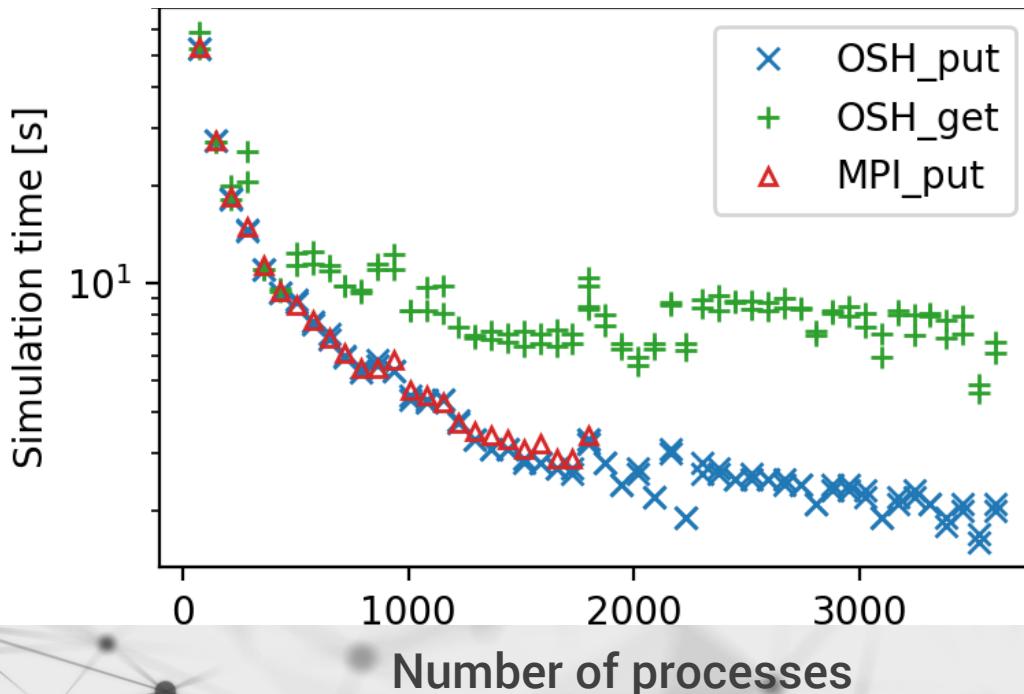


# Coarray ICAR

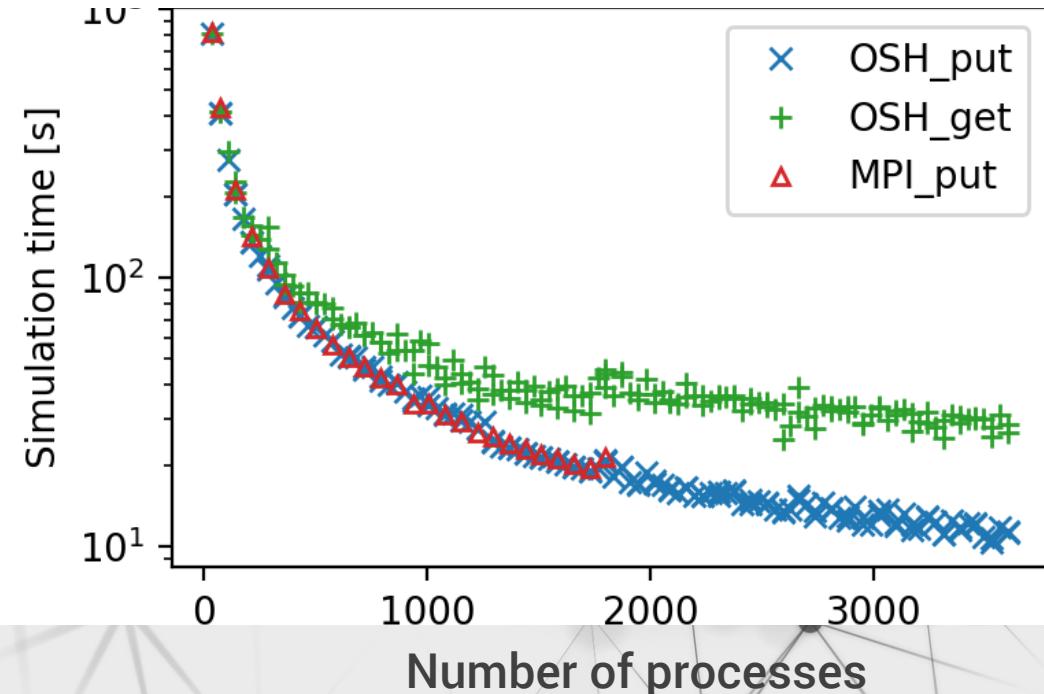
## Simulation Time

- OpenSHMEM vs MPI
- Puts vs gets

500 x 500 x 20 grid

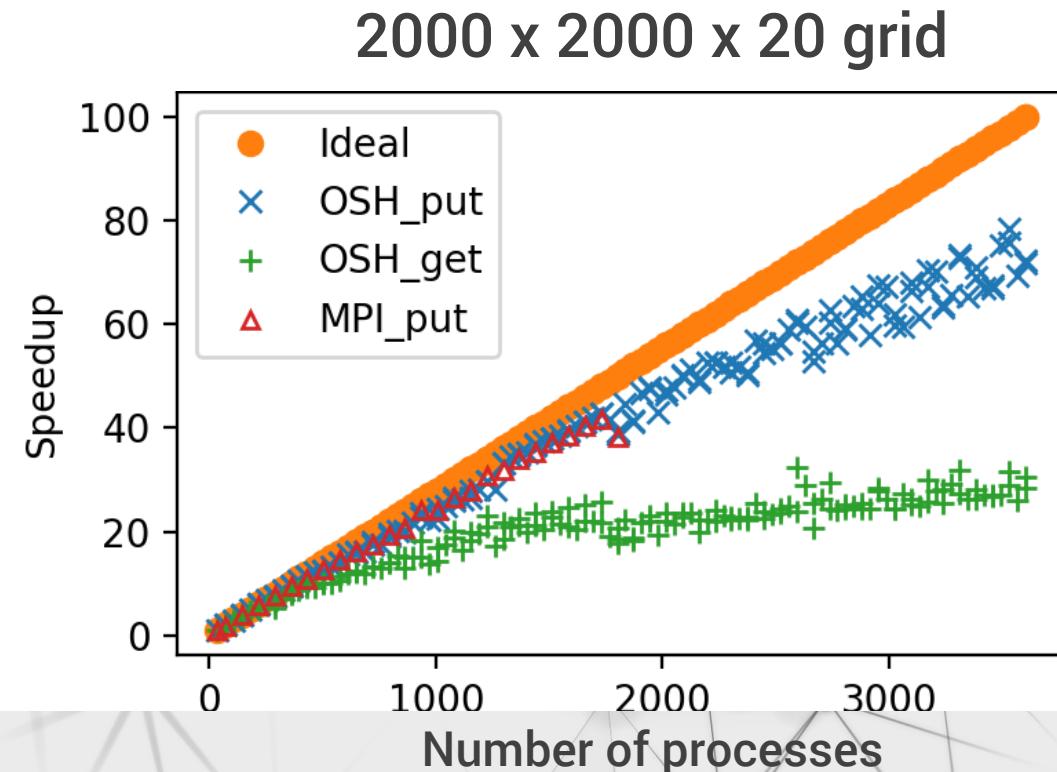
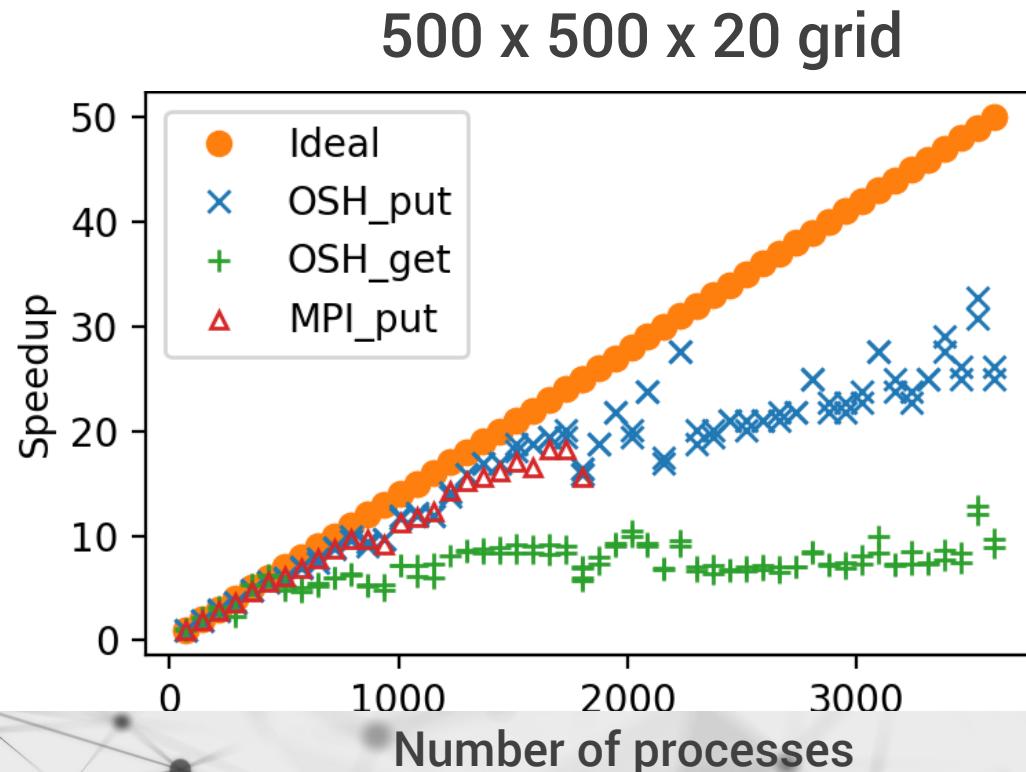


2000 x 2000 x 20 grid



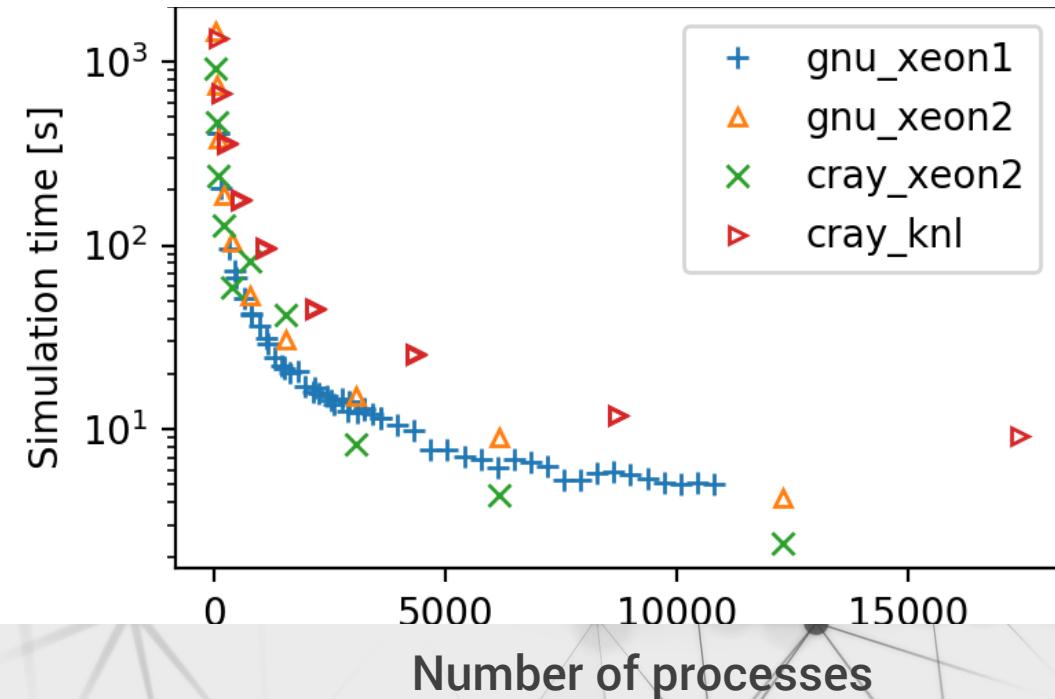
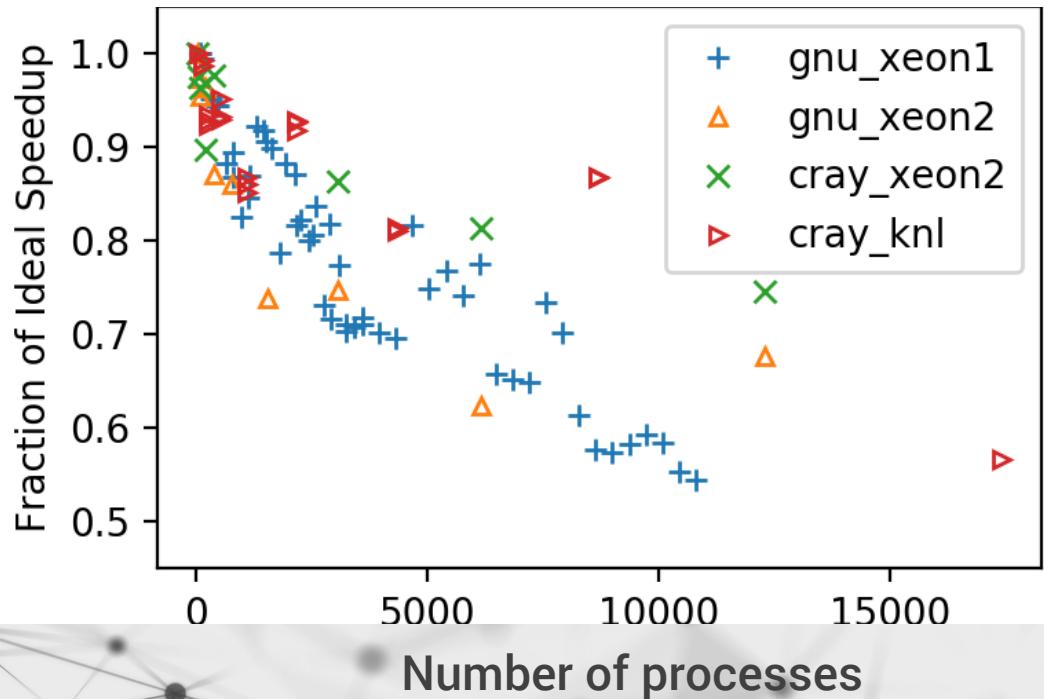
# Coarray ICAR

## Speedup



# Coarray ICAR

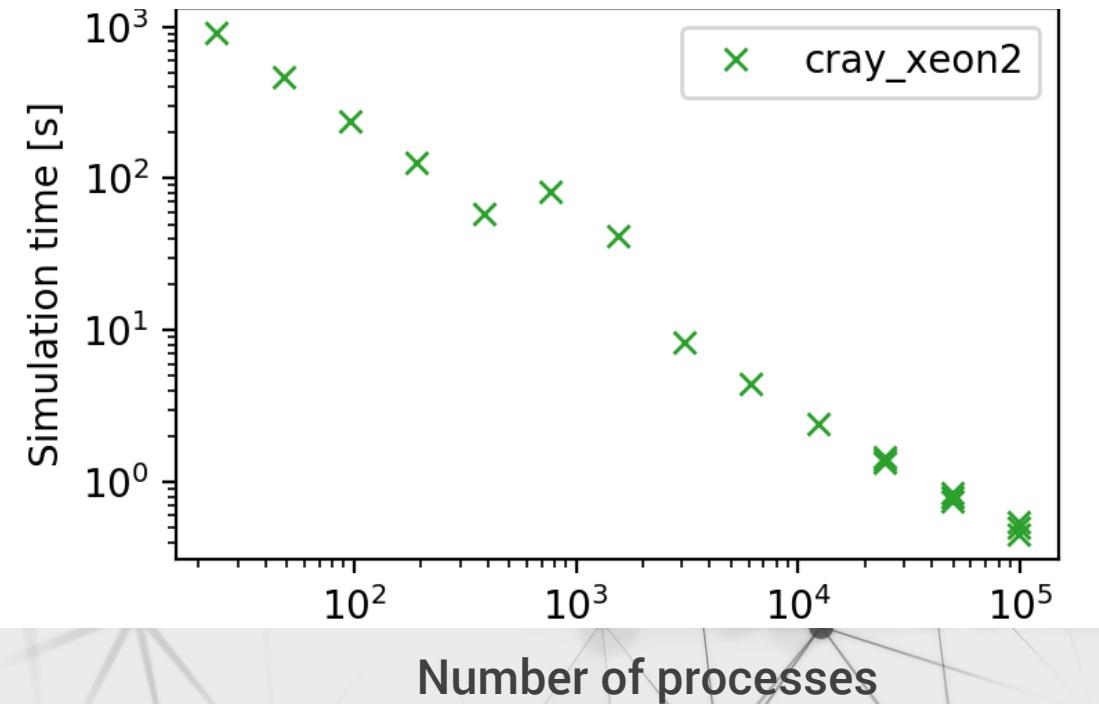
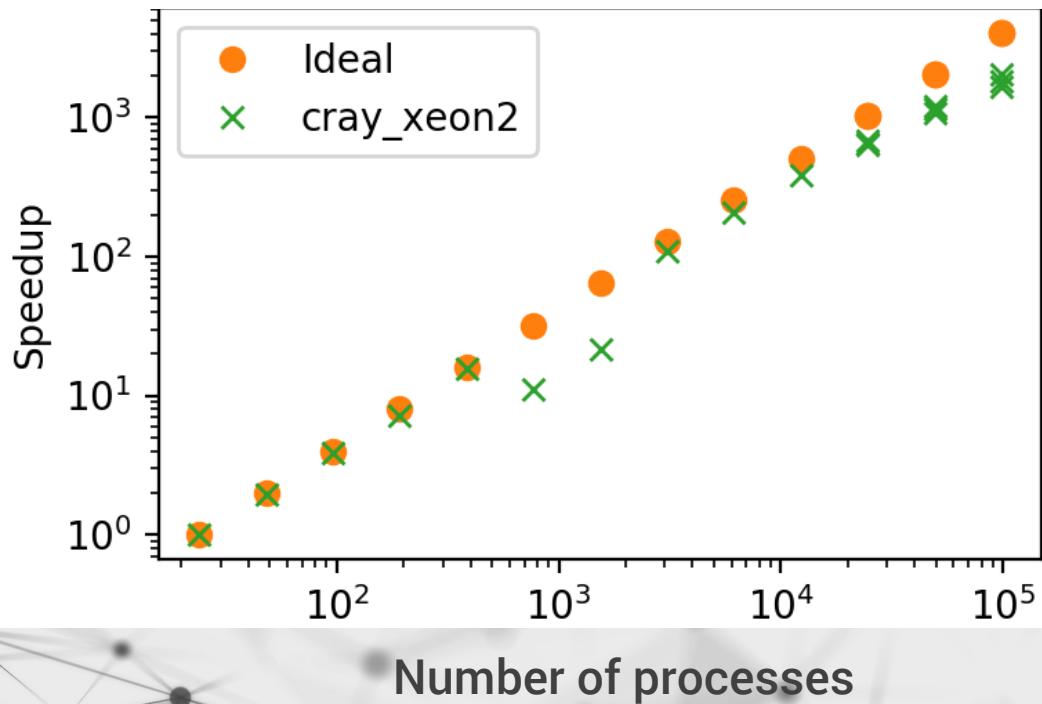
## Xeon vs KNL



Compilers & platforms: GNU on Cheyenne; Cray compiler on Edison (Broadwell), Cori (KNL).

# Coarray ICAR

## At Scale



Cray compiler on Edison.

# Conclusions

-  Fortran 2018 is a PGAS language supporting SPMD programming at scale.
-  Programming-model agnosticism is a life-saver.
-  High productivity pays off: from shared-memory parallelism to 100,000 cores in ~100 person-hours.

