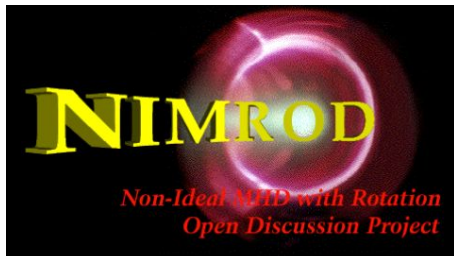# GPU infrastructure code and performance for NIMROD

Jacob King (Tech-X Corporation)

With contributions from
Eric Howell (Tech-X Corporation)
Brian Cornille (Univ. of Wisconsin)

April 3rd 2022
CTTS SciDAC Meeting

# Next gen computing: SciDAC codes need to exploit GPUs

- Perlmutter (NERSC), Frontier (OLCF) and Aurora (ALCF) contain GPUs
  - Three different hardware vendors (Nvidia, AMD, Intel)
  - Three different software eco-systems
- Dust yet to settle on unified programming paradigm
  - CUDA, HIP, SYCL, OpenCL, OpenACC, OpenMP 5, language standards
  - We've explored OpenACC for NIMROD – support by nvhpc / gcc
  - May switch to OpenMP 5 or language standards as dust settles

# Effort splits NIMROD code into infrastructure and physics repos

- Physics: anything with units
- Infrastructure: FE evaluation & integration, linear algebra, parallel decomp, FFTs, inverse mapping, etc.
- Infrastructure repo now open-source, open-access
  - See https://gitlab.com/NIMRODteam/nimrod-abstract
  - Allows access by ASCR/computing facility partners and compiler teams
  - Useful for hackathons and diagnosing performance and compiler issues
- Using gitlab enables modern software tools
  - Code review
  - Issue tracking
  - Continuous integration testing
  - Code coverage reports

# Infrastructure repo close to maturity

- Relies on modern Fortran abstract types to enable flexibility
  - Define interface at abstract (virtual) level
  - Now: working on version with real types
- Investigating API on GPU before expanding concrete base classes
- Done: I/O, seams, MPI with mpi_f08, timing, FEM types, linear algebra, Dirichlet boundary conditions, regularity
- Remaining big tasks: preconditioning, FFTs, surface integrals
- Next: report on progress on GPU by overviewing example application
- Example application: solve time-dependent Laplace equation in a periodic box

Thanks to the full GPU hackathon team!

For the openACC material here, I'd specifically like to thank Eric Howell, Brian Cornille, Torrin Bechtel, Robbie Searles (Nvidia) and Vassillios Mewes (OLCF)

```
 1   !----------------------------------------------------------------
 2   !! run a time-dependent calculation of a Laplacian
 3   !----------------------------------------------------------------
 4   #include "config.f"
 5   !----------------------------------------------------------------
 6   !* run a time-dependent calculation of a Laplacian
 7   !----------------------------------------------------------------
 8   PROGRAM ex1
 9     USE local
10     USE io
11     USE timer_mod
12     USE pardata_mod
13     USE seam_mod
14     USE gblock_mod
15     USE nodal_mod
16     USE quadrature_mod
17     USE vector_mod
18     USE matrix_mod
19     USE dump_mod
20     USE finite_element_mod
21     USE create_linalg_mod
22     USE xfer_vector_to_fem_mod
23     USE iter_cg_mod
24     USE iterdata_mod
25     USE integrand_mod
26     USE fem_utils_mod
27     IMPLICIT NONE
28
29     INTEGER(i4) :: nstep=10,maxit=100,ngr=2
30     REAL(r8) :: dt=1._r8,diff=1._r8,tol=1.e-8,theta=0.5
31     LOGICAL :: static_condensation=.FALSE.
32     CHARACTER(64) :: dump_file='dump.00000.h5'
33     CHARACTER(8) :: int_formula='gaussian' ! or 'lobatto'
34     CHARACTER(64) :: opt_name,opt_val
35     INTEGER(i4) :: opt_ival,ii,iarg,nargs,istep0,istep,fac,max_nq,ibl
36     REAL(r8) :: opt_rval,time,t0
```

ex1.f  13 KB

- All module import ABSTRACT base types at the example level

5

```
36    REAL(r8) :: opt_rval,time,t0
37    TYPE(block_storage), ALLOCATABLE :: blk(:)
38    TYPE(seam_type) :: seam
39    TYPE(rfem_storage), ALLOCATABLE, TARGET :: field(:)
40    TYPE(qp_real), ALLOCATABLE :: qpfield(:)
41    TYPE(field_list_pointer) :: io_field_list(1)
42    TYPE(rvec_storage), ALLOCATABLE :: rhs(:),sln(:)
43    TYPE(rmat_storage), ALLOCATABLE :: mat(:)
44    REAL(r8), ALLOCATABLE :: rhs_norm(:,:)
45    REAL(r8), PARAMETER :: ksq = 2*twopi**2
46    TYPE(iterdata) :: itdat
47
48    INTEGER(i4) :: idepth
49    INTEGER(i4), SAVE :: iftn=-1
50
51    CALL par%init
52    CALL timer%init
53    CALL timer%start_timer_l0('ex1','ex1',iftn,idepth)
54    CALL fch5init
55  !-----------------------------------------------------------------
56  ! determine configuration from the command line
57  !-----------------------------------------------------------------
58    nargs=command_argument_count()
59    IF (MOD(nargs,2) /= 0) THEN
60      CALL print_usage
61      CALL nim_stop('Argument error')
62    ENDIF
63    DO ii=1,nargs/2
64      iarg=(ii-1)*2+1
65      CALL get_command_argument(iarg,opt_name)
66      CALL get_command_argument(iarg+1,opt_val)
67      SELECT CASE(opt_name)
68      CASE("nstep","-nstep","--nstep")
69        READ(opt_val,'(i8)') opt_ival
70        nstep=opt_ival
71      CASE("dt","-dt","--dt")
72        READ(opt_val,'(f8.0)') opt_rval
73        dt=opt_rval
```

- All data types are only know at the abstract level
- Storage arrays allow different concrete instantiations for different blocks
- Timer and parallel types require initialization of singleton objects
- Example takes command line input (no namelist)

6

```
74     CASE("diff","-diff","--diff")
75       READ(opt_val,'(f8.0)') opt_rval
76       diff=opt_rval
77     CASE("maxit","-maxit","--maxit")
78       READ(opt_val,'(i8)') opt_ival
79       maxit=opt_ival
80     CASE("tol","-tol","--tol")
81       READ(opt_val,'(f8.0)') opt_rval
82       tol=opt_rval
83     CASE("theta","-theta","--theta")
84       READ(opt_val,'(f8.0)') opt_rval
85       theta=opt_rval
86     CASE("int_formula","-int_formula","--int_formula")
87       int_formula=opt_val(1:8)
88     CASE("stat_con","-stat_con","--stat_con")
89       READ(opt_val,'(i8)') opt_ival
90       IF (opt_ival==0) THEN
91         static_condensation=.FALSE.
92       ELSE
93         static_condensation=.TRUE.
94       ENDIF
95     CASE("dump","-dump")
96       dump_file=opt_val
97     CASE DEFAULT
98       CALL print_usage
99       CALL nim_stop('Argument error')
100    END SELECT
101   ENDDO
102   ALLOCATE(rhs_norm(3,nstep))
103  !----------------------------------------------------------
104  ! read dump file which initializes the parallel decomposition
105  !----------------------------------------------------------
106   io_field_list(1)%name="field"
107   io_field_list(1)%type="rfem_storage"
108   CALL dump_read(TRIM(dump_file),t0,istep0,blk,seam,io_field_list)
109   time=t0
110  !----------------------------------------------------------
111  ! set up blocks
112  !----------------------------------------------------------
113   DO ibl=1,par%nbl
114     CALL blk(ibl)%b%block_intg_formula_set(ngr,int_formula,int_formula)
115     CALL blk(ibl)%b%block_metric_set('lin')
116   ENDDO
```

- Mesh and block decomposition set by initialization
- "dump_read" also initializes parallel block decomposition
- Blocks govern FEM integration and setting quadrature rules and weights mirrors non-abstract code
- This structure allows for natural separate of blocks (e.g. CK type bound to a CK block with CK integration rules)

# Example overview

```
117  !-----------------------------------------------------
118  ! copy out field list and reset as pointer
119  !-----------------------------------------------------
120    ALLOCATE(field(par%nbl))
121  #ifdef __gfortran
122    DO ibl=1,par%nbl
123      SELECT TYPE(iofield=>io_field_list(1)%p)
124      TYPE IS (rfem_storage)
125        CALL copy_field(field(ibl)%f,iofield(ibl)%f)
126      END SELECT
127    ENDDO
128  #else
129    DO ibl=1,par%nbl
130      SELECT TYPE(iofield=>io_field_list(1)%p(ibl))
131      TYPE IS (rfem_storage)
132        CALL copy_field(field(ibl)%f,iofield%f)
133      END SELECT
134    ENDDO
135  #endif
136    io_field_list(1)%p=>field
137  !-----------------------------------------------------
138  ! initialize seams
139  !-----------------------------------------------------
140    max_nq=field(1)%f%nqty*MAX(1_i4,field(1)%f%pd-1_i4)
141    CALL seam%init(max_nq,max_nq)
142  !-----------------------------------------------------
143  ! set up linear algebra and quadrature point structures
144  !-----------------------------------------------------
145    ALLOCATE(rhs(par%nbl),sln(par%nbl),mat(par%nbl),qpfield(par%nbl))
146    DO ibl=1,par%nbl
147      on_gpu=.TRUE.
148      CALL create_vector_for_fem(rhs(ibl)%v,field(ibl)%f)
149      CALL rhs(ibl)%v%alloc_with_mold(sln(ibl)%v)
150      CALL sln(ibl)%v%zero
151      CALL rhs(ibl)%v%set_edge_vars(seam%s(ibl))
152      CALL create_matrix_for_fem(mat(ibl)%m,field(ibl)%f)
153      on_gpu=.FALSE.
154      CALL field(ibl)%f%qp_alloc(blk(ibl)%b%ng,qpfield(ibl))
155      CALL field(ibl)%f%init_basis_ftn(blk(ibl)%b%xg)
156      CALL compute_alpha_real(field(ibl)%f,blk(ibl)%b%metric,blk(ibl)%b%ng)
157    ENDDO
158    fac=1_i4 ! dummy factor structure
```

- Minor penalty for not binding fields to blocks: one extra copy after read
- #ifdefs based on compiler demonstrate difficulty with modern Fortran
- Vectors and matrices are initialized from FE field
- Each field (or combination of fields) needs to know how to map to appropriate linear algebra structures
  - See create_*_for_fem
- Basis functions are precomputed based on FE field

8

# Example overview

```
159  !-----------------------------------------------------------------
160  ! set iterdata input
161  !-----------------------------------------------------------------
162    itdat%maxit=maxit
163    itdat%tol=tol
164  !-----------------------------------------------------------------
165  ! time-step loop to solve the diffusion equation
166  ! while we could form this linear operator once outside the loop, do it inside
167  ! to better mirror the timing behavior with a nonlinear operator
168  !-----------------------------------------------------------------
169    DO istep=istep0+1,istep0+nstep
170      DO ibl=1,par%nbl
171        CALL field(ibl)%f%qp_update(qpfield(ibl),blk(ibl)%b%metric)
172      ENDDO
173      IF (static_condensation) THEN
174        CALL create_vector(blk,rhs,seam,diff_rhs,.FALSE.,'none',elim_matrix=mat)
175        CALL create_matrix(blk,mat,fac,seam,diff_op,.FALSE.,'none',.TRUE.)
176      ELSE
177        CALL create_vector(blk,rhs,seam,diff_rhs,.FALSE.,'none')
178        CALL create_matrix(blk,mat,fac,seam,diff_op,.FALSE.,'none',.FALSE.)
179      ENDIF
180      CALL iter_cg_real_dir_solve(mat,1_i4,rhs,sln,seam,itdat)
181      ii=istep-istep0
182      rhs_norm(1,ii)=itdat%rhs_norm
183      rhs_norm(2,ii)=rhs_norm(1,1)*((1._r8-dt*(1._r8-theta)*diff*ksq)        &
184                                  /(1._r8+dt*theta*diff*ksq))**(ii-1)
185      rhs_norm(3,ii)=rhs_norm(1,1)*EXP(-(time-t0)*diff*ksq)
186      IF (itdat%converged) THEN
187        IF (par%node==0) THEN
188          WRITE(nim_wr,'(a,i4,a,es10.3,a,i6)')                               &
189            'Step=',istep,': solved system to err=',itdat%err,' in ',        &
190            itdat%its,' iterations'
191          WRITE(nim_wr,'(a,es10.3,2a)')                                      &
192            '  rhs_norm=',itdat%rhs_norm,' seed=',TRIM(itdat%seed)
193        ENDIF
```

- Time step loop:
    - update quadrature storage
    - Create matrix/RHS
    - Call iterative solver
    - Use RHS norm (inf_norm in this case) to test solution
    - Check convergence

```
194    DO ibl=1,par%nbl
195      IF (static_condensation) THEN
196        !TODO: edit postsolve doc string to be correct
197        CALL mat(ibl)%m%elim_postsolve(rhs(ibl)%v,sln(ibl)%v)
198        CALL xfer_vector_to_fem(sln(ibl)%v,field(ibl)%f)
199      ELSE
200        CALL xfer_vector_to_fem(sln(ibl)%v,field(ibl)%f)
201      ENDIF
202    ENDDO
203  ELSE
204    IF (par%node==0) THEN
205      WRITE(nim_wr,'(a,i4,a,es10.3,a,i6,2a)')                &
206        'Step=',istep,': failed to converge with err=',itdat%err,' in ',  &
207        itdat%its,' iterations'
208      WRITE(nim_wr,'(a,es10.3,2a)')                          &
209        '  rhs_norm=',itdat%rhs_norm,' seed=',TRIM(itdat%seed)
210    ENDIF
211    EXIT
212  ENDIF
213  time=time+dt
214  ENDDO
215 !--------------------------------------------------------------------
216 ! write dump file and exit
217 !--------------------------------------------------------------------
218  IF (itdat%converged) THEN
219    CALL dump_write(time,istep0+nstep,blk,seam,io_field_list)
220  ENDIF
221 !--------------------------------------------------------------------
222 ! display timings and expected result
223 !--------------------------------------------------------------------
224  CALL timer%end_timer_l0(iftn,idepth)
225  IF (par%node==0) THEN
226    CALL timer%report
227    WRITE(nim_wr,'(2a)') '  rhs_norm              expected',      &
228      '                     exact'
229    WRITE(nim_wr,*) rhs_norm
230  ENDIF
231  CALL timer%finalize
232  CALL par%finalize
233 CONTAINS
```

- After time step loop, we test solution

Wait! Where's the GPU stuff?!?

# Example overview

```fortran
249    !----------------------------------------------------------------
250    ! Time-dependent diffusion operator integrand routine
251    !----------------------------------------------------------------
252    SUBROUTINE diff_op(bl, integrand)
253      USE local
254      USE gblock_mod
255      CLASS(gblock), INTENT(IN) :: bl
256      REAL(r8), CONTIGUOUS, INTENT(OUT) :: integrand(:,:,:,:,:)
257
258      INTEGER(i4) :: idepth
259      INTEGER(i4), SAVE :: iftn=-1
260
261      CALL timer%start_timer_l1('diff_op','diff_op',iftn,idepth)
262
263      ASSOCIATE (nvert=>field(bl%ibl)%f%nbasis, nel=>bl%nel, ng=>bl%ng,      &
264                 nq=>field(bl%ibl)%f%nqty,                                   &
265                 alpha=>field(bl%ibl)%f%qab%alf,                            &
266                 grad_alpha=>field(bl%ibl)%f%qab%aldf, wdetj=>bl%wdetj)
267      BLOCK
268        INTEGER(i4) :: iq, iv, jv, ie, ig
269        REAL(r8) :: tmpsum, fac
270
271        fac=dt*diff*theta
272        !$acc parallel present(integrand,alpha,grad_alpha,wdetj) if(on_gpu) &
273        !$acc copyin(fac,nel,nvert,nq,ng)
274        !$acc loop gang worker
275        DO ie = 1, nel
276          !!$acc cache(wdetj(:,ie))
277          !$acc loop vector collapse(2) independent
278          DO iv = 1, nvert
279            DO jv = 1, nvert
280              tmpsum = 0._r8
281              !$acc loop seq
282              DO ig = 1, ng
283                tmpsum = tmpsum                                             &
284                       + wdetj(ig,ie)*(alpha(ig,jv,ie,1)*alpha(ig,iv,ie,1)   &
285                       + fac*SUM(grad_alpha(ig,jv,ie,:)*grad_alpha(ig,iv,ie,:)))
286              ENDDO
287              DO iq=1,nq
288                integrand(iq,iq,ie,jv,iv) = tmpsum
289              ENDDO
290            ENDDO
291          ENDDO
292        ENDDO
293        !$acc end parallel
294      END BLOCK
295      END ASSOCIATE
296
297      CALL timer%end_timer_l1(iftn,idepth)
298    END SUBROUTINE diff_op
```

- Physics application is not completely exempt from being GPU aware, openACC statements required in integrand routines
- Parallel clause create GPU kernel
- present/copyin clauses manage data movement from host to device
- Structure of acc loops follow pre-set pattern
- With full time step loop on device, minimal copyin statements will be needed (more on this later)

# Example overview

```fortran
300  !-----------------------------------------------------------------
301  ! Time-dependent diffusion RHS integrand routine
302  !-----------------------------------------------------------------
303    SUBROUTINE diff_rhs(bl, integrand)
304      USE local
305      USE gblock_mod
306      IMPLICIT NONE
307
308      CLASS(gblock), INTENT(IN) :: bl
309      REAL(r8), CONTIGUOUS, INTENT(OUT) :: integrand(:,:,:)
310
311      INTEGER(i4) :: idepth
312      INTEGER(i4), SAVE :: iftn=-1
313
314      CALL timer%start_timer_l1('diff_rhs','diff_rhs',iftn,idepth)
315
316      ASSOCIATE (nvert=>field(bl%ibl)%f%nbasis, nel=>bl%nel, ng=>bl%ng,      &
317                 nq=>field(bl%ibl)%f%nqty, qfield=>qpfield(bl%ibl)%qpf,       &
318                 dqfield=>qpfield(bl%ibl)%qpdf,                               &
319                 alpha=>field(bl%ibl)%f%qab%alf,                             &
320                 grad_test=>field(bl%ibl)%f%qab%aldf, wdetj=>bl%wdetj)
321      BLOCK
322        INTEGER(i4) :: iq, iv, ie, ig
323        REAL(r8) :: tmpsum, fac
324
325        fac = dt*diff*(1._r8-theta)
326        !$acc parallel present(integrand,alpha,grad_test,qfield,dqfield,wdetj) &
327        !$acc copyin(fac,nel,nvert,nq,ng) if(on_gpu)
328        !$acc loop gang worker
329        DO ie = 1, nel
330          !!$acc cache(wdetj(:,ie))
331          !$acc loop vector collapse(2) independent private(tmpsum)
332          DO iv = 1, nvert
333            DO iq= 1, nq
334              tmpsum = 0._r8
335              !$acc loop seq
336              DO ig = 1, ng
337                tmpsum = tmpsum + wdetj(ig,ie)*                                &
338                         (qfield(ig,ie,1,iq)*alpha(ig,iv,ie,1) -               &
339                          fac*SUM(dqfield(ig,ie,:,iq)*grad_test(ig,iv,ie,:)))
340              ENDDO
341              integrand(iq,ie,iv)=tmpsum
342            ENDDO
343          ENDDO
344        ENDDO
345        !$acc end parallel
346      END BLOCK
347      END ASSOCIATE
348
349      CALL timer%end_timer_l1(iftn,idepth)
350    END SUBROUTINE diff_rhs
351  END PROGRAM ex1
```

- For completeness, here's the RHS routine
- It is pretty similar, but quadrature point data from field is used

- Majority of OpenACC code in infrastructure repo.
- Management of data locality is key (next slide)

12

# Device data management

```
109   !-----------------------------------------------------------
110   !* call block-wise finite element computations to create a vector
111   !-----------------------------------------------------------
112     SUBROUTINE create_vector_real(bl,vector,seam,integrand_func,dirichlet_bc,  &
113                          bc_component,elim_matrix)
114       USE matrix_mod
115       USE vector_mod
116       IMPLICIT NONE
117
118       !> block storage
119       TYPE(block_storage), INTENT(IN) :: bl(:)
120       !> vector to create stored in array by block
121       TYPE(rvec_storage), INTENT(INOUT) :: vector(:)
122       !* seam
123       TYPE(seam_type), INTENT(INOUT) :: seam
124       !> integrand procedure to use
125       PROCEDURE(integrand_vec_real) :: integrand_func
126       !> if true a dirichlet BC is applied to bc_component
127       LOGICAL, INTENT(IN) :: dirichlet_bc
128       !> components to use for dirichlet BC
129       CHARACTER(*), INTENT(IN) :: bc_component
130       !> matrix to use to eliminate interior DOFs stored in array by block
131       TYPE(rmat_storage), OPTIONAL, INTENT(IN) :: elim_matrix(:)
132
133       REAL(r8), ALLOCATABLE :: int_arr(:,:,:)
134       INTEGER(i4) :: ibl
135       CLASS(rvector), ALLOCATABLE :: temp_vec
136       INTEGER(i4) :: idepth
137       INTEGER(i4), SAVE :: iftn=-1
138
139       CALL timer%start_timer_l1(mod_name,'create_vector_real',iftn,idepth)
140   !-----------------------------------------------------------
141   !   loop over blocks and integrate
142   !-----------------------------------------------------------
143       DO ibl=1,SIZE(bl)
144         ASSOCIATE (vec=>vector(ibl)%v)
145           on_gpu=.TRUE.
146           ALLOCATE(int_arr(vec%nqty,vec%nel,vec%u_ndof))
147           !$acc data create(int_arr) if(on_gpu)
148           !$acc kernels async(vec%id) present(int_arr) if(on_gpu)
149           int_arr=0._r8
150           !$acc end kernels
151           CALL integrand_func(bl(ibl)%b,int_arr)
152           CALL vec%zero
153           CALL vec%assemble(int_arr)
154           DEALLOCATE(int_arr)
155           !$acc end data
156           on_gpu=.FALSE.
157   !-----------------------------------------------------------
158   !      call the regularity condition and the boundary condition routines.
159   !-----------------------------------------------------------
160           IF (bl(ibl)%b%r0block) CALL vec%regularity(seam%s(ibl),'all')
161           IF (dirichlet_bc) CALL vec%dirichlet_bc(bc_component,seam%s(ibl))
```

- Unstructured data blocks:
  - persistent data on GPU
  - "!$acc enter data create"
  - "!$acc exit data delete"
- Structured data blocks:
  - GPU data in local scope
  - "!$acc data create"
  - "!$acc end data"

```
4    !-----------------------------------------------------------
5    !* allocate vector
6    !-----------------------------------------------------------
7      SUBROUTINE alloc_real(rvec,poly_degree,mx,my,nqty,id)
8        IMPLICIT NONE
9
10       !> vector to allocate
11       CLASS(vec_rect_2D_real_acc), INTENT(INOUT) :: rvec
12       !> polynomial degree
13       INTEGER(i4), INTENT(IN) :: poly_degree
14       !> number of elements in the horizontal direction
15       INTEGER(i4), INTENT(IN) :: mx
16       !> number of elements in the vertical direction
17       INTEGER(i4), INTENT(IN) :: my
18       !> number of quantities
19       INTEGER(i4), INTENT(IN) :: nqty
20       !> ID for parallel streams
21       INTEGER(i4), INTENT(IN) :: id
22
23       INTEGER(i4) :: idepth
24       INTEGER(i4), SAVE :: iftn=-1
25
26       CALL timer%start_timer_l2(mod_name,'alloc_real',iftn,idepth)
27   !-----------------------------------------------------------
28   !   store grid and vector dimensions
29   !-----------------------------------------------------------
30       rvec%nqty=nqty
31       rvec%mx=mx
32       rvec%my=my
33       rvec%n_side=poly_degree-1
34       rvec%n_int=(poly_degree-1)**2
35       rvec%u_ndof=(poly_degree+1)**2
36       rvec%pd=poly_degree
37       rvec%nel=mx*my
38       rvec%ndim=2
39       rvec%id=id
40       !$acc enter data copyin(rvec)
41   !-----------------------------------------------------------
42   !   allocate space according to the basis functions needed.
43   !-----------------------------------------------------------
44       SELECT CASE(poly_degree)
45       CASE(1)  !  linear elements
46         ALLOCATE(rvec%arr(nqty,0:mx,0:my))
47         !$acc enter data create(rvec%arr) if(on_gpu)
48         NULLIFY(rvec%arri,rvec%arrh,rvec%arrv)
49       CASE(2:)  !  higher-order elements
50         ALLOCATE(rvec%arr(nqty,0:mx,0:my))
51         !$acc enter data create(rvec%arr) if(on_gpu)
52         ALLOCATE(rvec%arrh(nqty,poly_degree-1,1:mx,0:my))
53         !$acc enter data create(rvec%arrh) if(on_gpu)
54         ALLOCATE(rvec%arrv(nqty,poly_degree-1,0:mx,1:my))
55         !$acc enter data create(rvec%arrv) if(on_gpu)
56         ALLOCATE(rvec%arri(nqty,(poly_degree-1)**2,1:mx,1:my))
57         !$acc enter data create(rvec%arri) if(on_gpu)
58       END SELECT
```

13

# GPU port progress

- Integrand, vector assemble kernels on GPU
- Left to port: matvec, qp_update, matrix assemble
- Testing on Ascent at ORNL
- Kernels on 1x V100 GPU comparable or 2x faster than 32 POWER9 cores
  - CPU performance measured with timer
  - GPU performance measured with Nsight Compute
  - 64x64 FE mesh with bi-quintic basis functions
- But there's 6x GPU/node but only 42 POWER9 cores
- Overall application performance degraded on GPU
  - Until full time-step loop on GPU additional host-device memory transfers
  - After time-step loop on GPU plan optimization