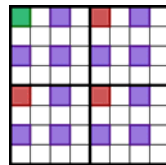


# A Brief Introduction to the NAS MG Benchmark

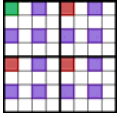
LaR 2005

November 12, 2005



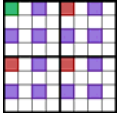
Brad Chamberlain

Cray Inc.



# NAS Parallel Benchmarks (NPB)

- A set of benchmarks developed...
  - ...in the early `90's by NASA's Advanced Supercomputing division
  - ...to model computations and data access patterns from CFD codes
  - ...originally released as paper & pencil benchmarks (v1.x)
  - ...then as MPI reference implementations (v2.x)
  - ...now versions available in a variety of languages
    - Java, OpenMP, HPF (v3.x)
    - UPC, Co-Array Fortran, Titanium, ZPL, ... (by respective groups)
- Among the most useful benchmark suites in HPC
  - well-designed and -maintained
  - good variety of data access patterns, communication requirements
  - open-source
  - well-understood, -used



# NAS Parallel Benchmarks (NPB)

- 8 Benchmarks:

- 5 kernels:

- **EP**: embarrassingly parallel
- **MG**: multigrid
- **CG**: conjugate gradient
- **FT**: Fourier transform
- **IS**: integer sort

- 3 pseudo-applications

- **BT**: block transpose
- **LU**: LU factorization
- **SP**: pentadiagonal

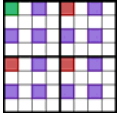
- Though useful, also domain-specific

- focus on CFD algorithms is good, but restrictive
- other HPC application areas would do well to create similar suites

- Often difficult to understand from the code

- terse variable names
- SPMD-style programming details

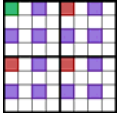
<http://www.nas.nasa.gov/Software/NPB/>



# NPB Basics

---

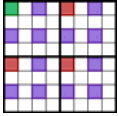
- Each benchmark supports a number of problem sizes:
  - **S**: Small
  - **W**: Workstation
  - **A, B, C, D**: “Production” sizes (added as machines grow)
- NPB codes use common...
  - ...timing routines
  - ...pseudo-random number generators
  - ...output routines
- NAS implementations are simply “one way to write the algorithm”
  - typically written to work well with both caches and vectors
  - how closely they should be followed can be a tricky question



# The NAS MG Benchmark

---

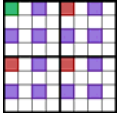
**Mathematically:** use a 3D multigrid method to find an approximate solution to a discrete Poisson problem ( $\nabla^2 u = v$ )



# Disclaimers

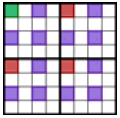
---

- I tend to take a very operational approach to describing algorithms
  - what are data structures?
  - what operations are performed on them?
  - what data distribution is used for parallel version?
  - what communication is required?
- I'm not a multigrid expert
  - More of a parallel algorithm/language person
  - A good deal of experience with NAS MG



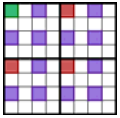
# NAS MG: Operational View

- Data structures:
  - 3D arrays & 3D hierarchical arrays (2D in my pictures)
  - 3D sparse arrays can also be useful
- 4 primary kernels:
  - each computed using 27-point stencils
    - **resid**: compute residual
    - **psinv**: compute approximate inverse
    - **rprj3**: projection from fine grid to coarse
    - **interp**: interpolation from coarse grid to fine
  - periodic boundary conditions
- computation of approximate norms
  - **norm2u3**: approximate L2 & uniform norms
- initialization, output

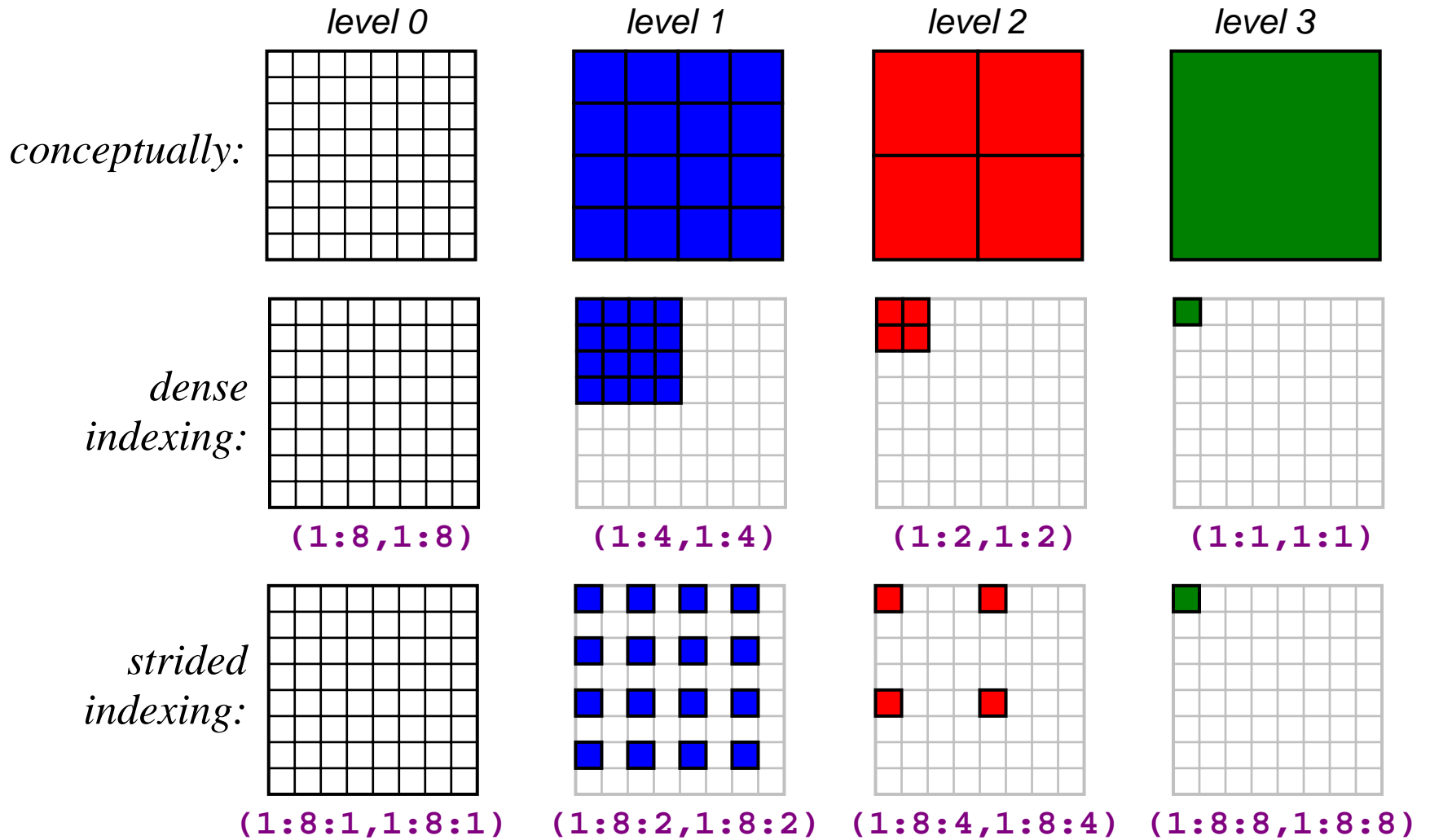


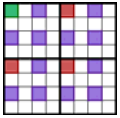
# NAS MG: Parallel Implementation

- Arrays typically use 3D block distributions
  - optimal load balance
  - good surface-to-volume ratio
  - ghost cells allocated for caching neighbors' values
- Communication Idioms:
  - 4 kernels require point-to-point communication
  - toroidal communication required for boundaries
  - global reductions required to compute norms
  - reductions useful during initialization as well

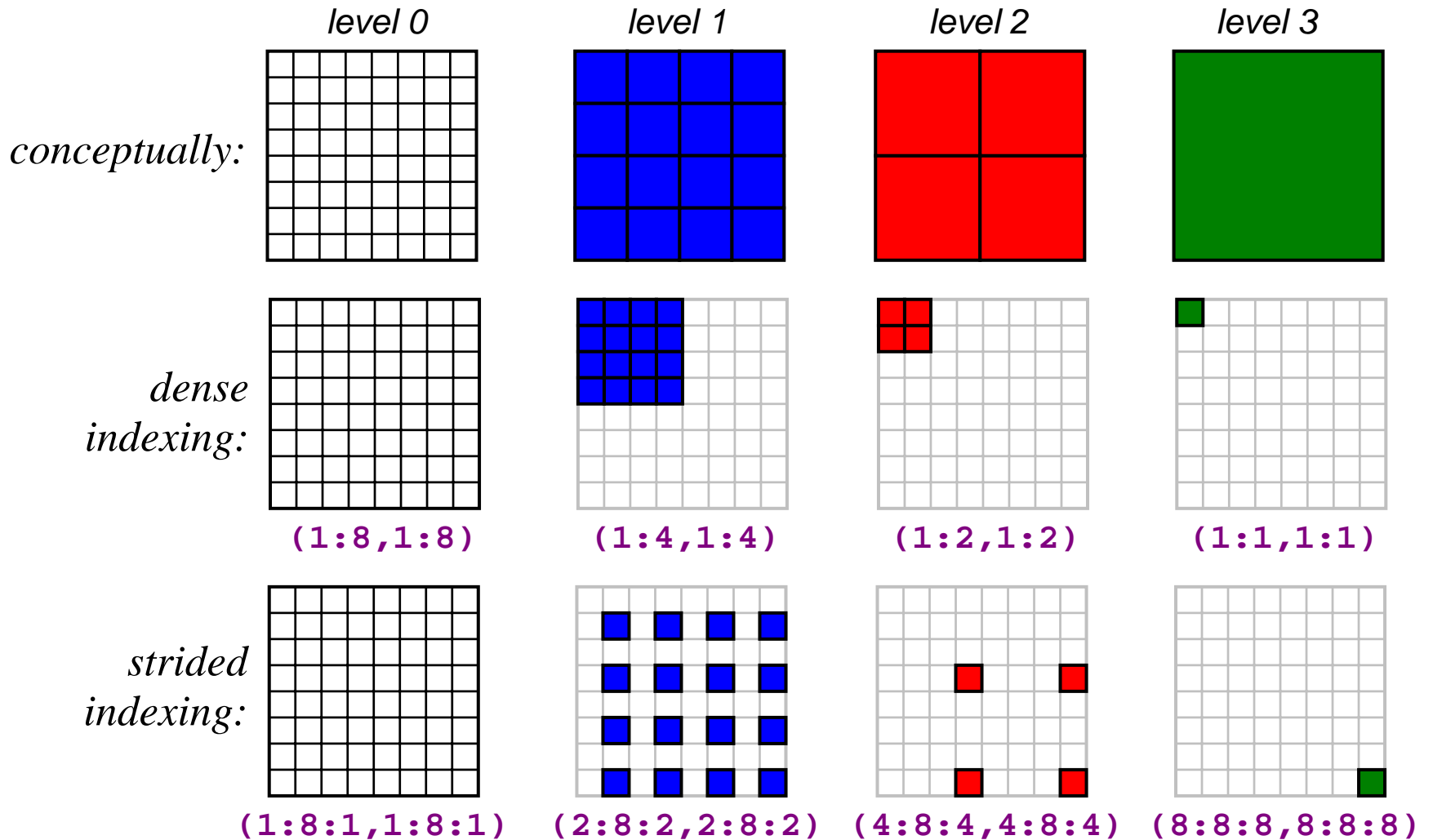


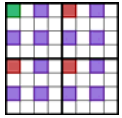
# Hierarchical Arrays



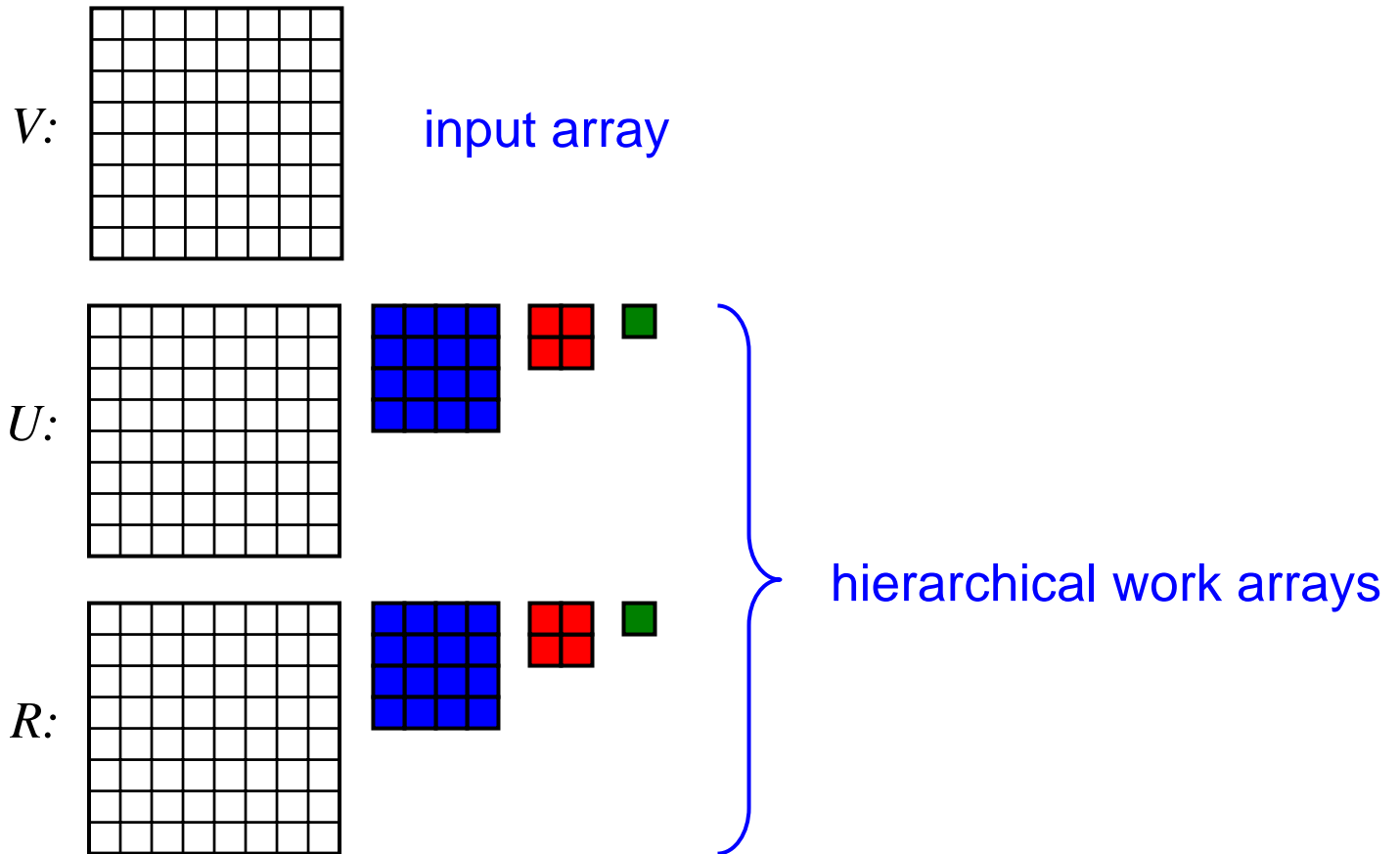


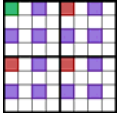
# Hierarchical Arrays



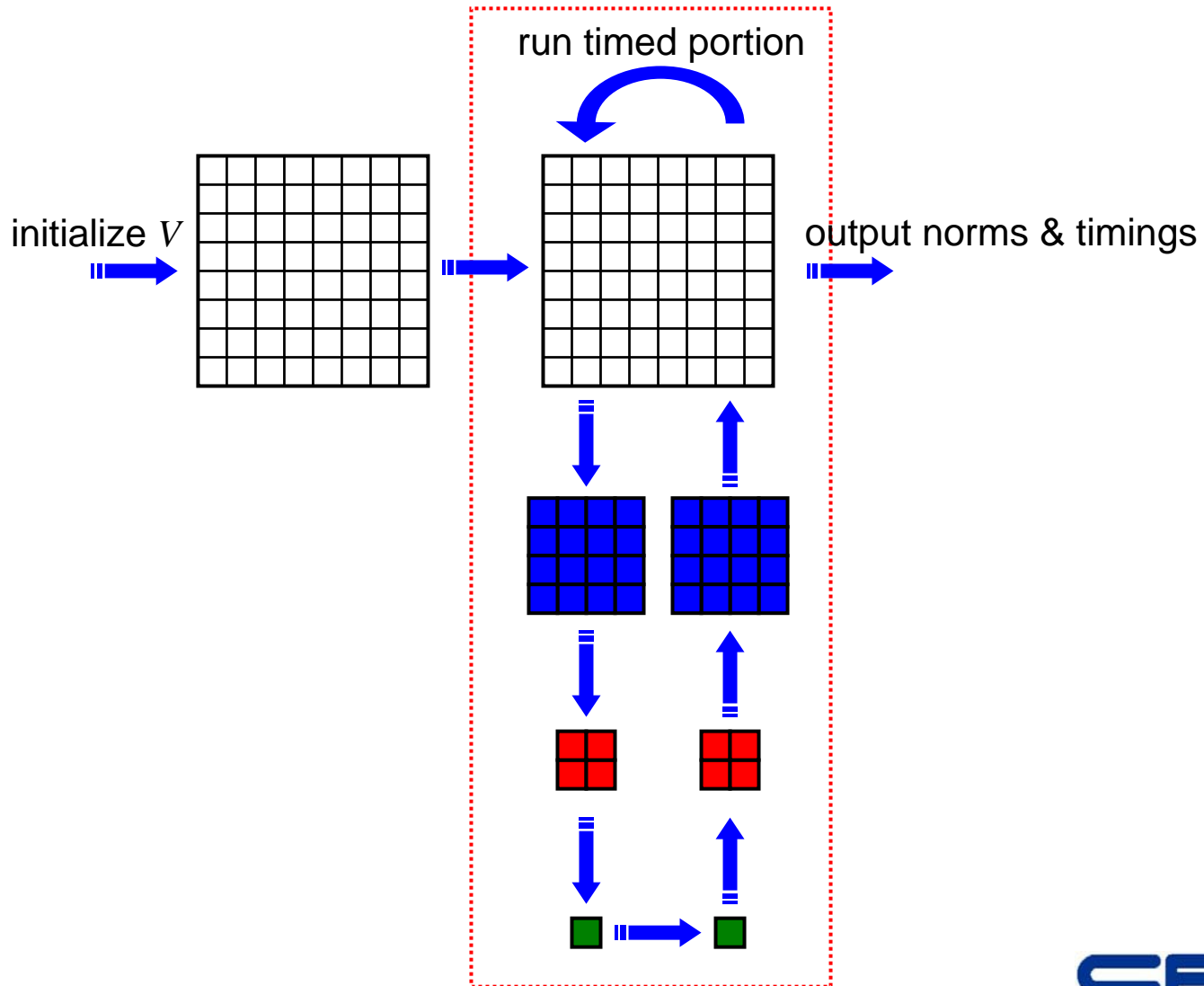


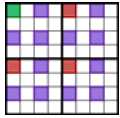
# MG's arrays



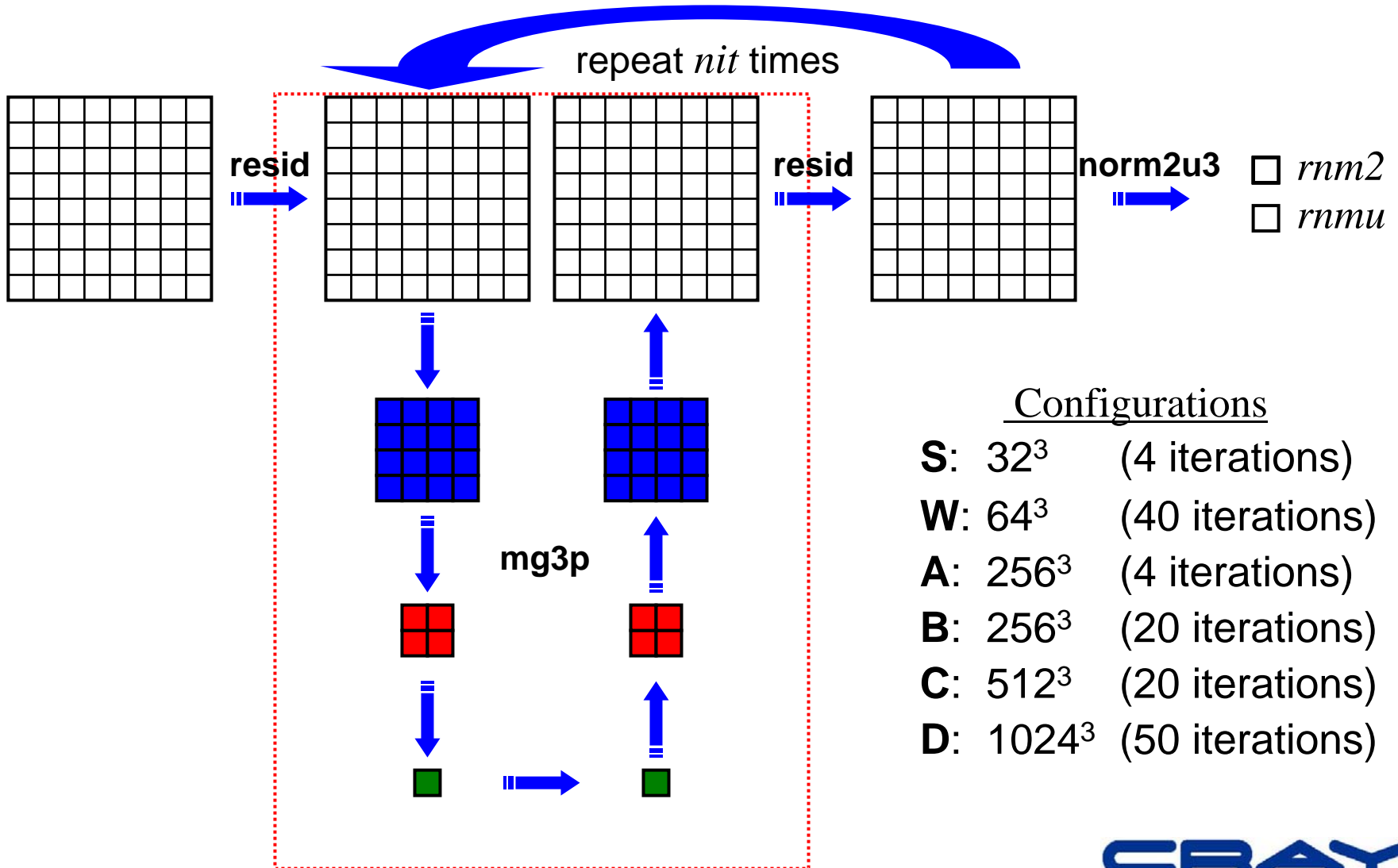


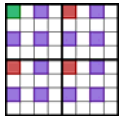
# Overview of MG



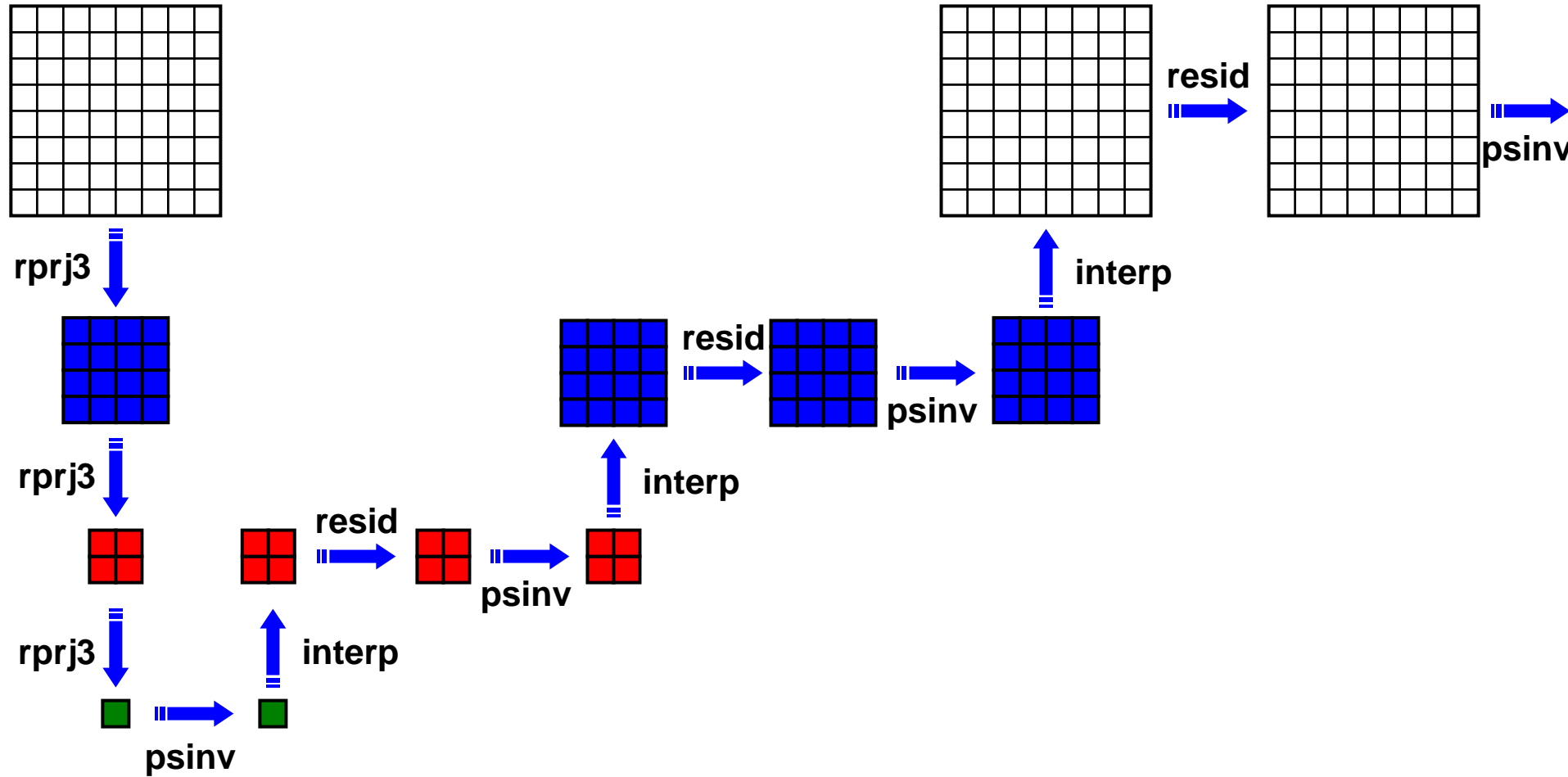


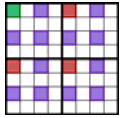
# MG's Timed Portion



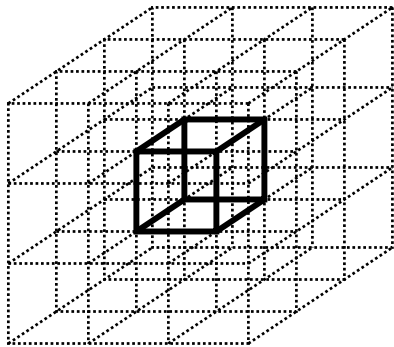


# MG's Guts (*mg3P*)

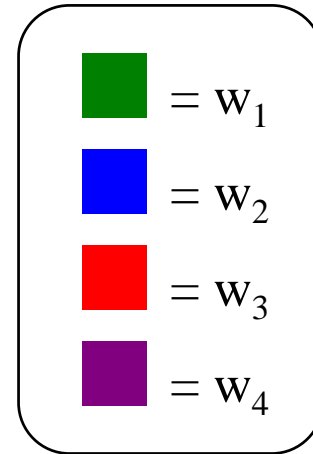
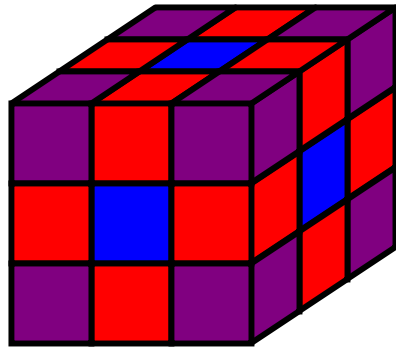




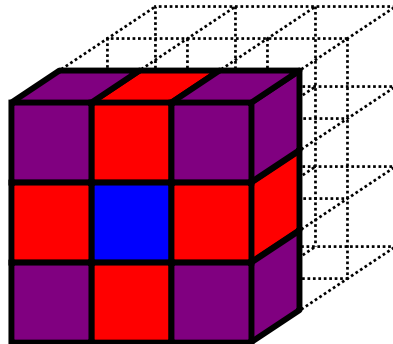
# 27-point stencils



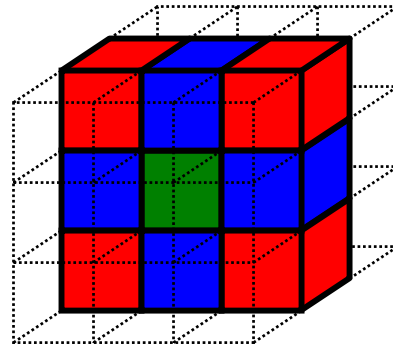
=



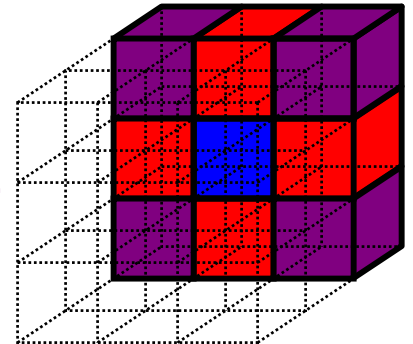
=

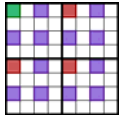


+



+





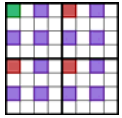
# Our First Stencil: $resid(R, V, U)$

$$R = V - \text{convolve}(\begin{array}{|c|c|c|} \hline \color{purple} & \color{white} & \color{purple} \\ \hline \color{purple} & \color{blue} & \color{white} \\ \hline \color{purple} & \color{white} & \color{purple} \\ \hline \end{array}, U)$$

The diagram shows a 5x5 grid labeled  $R$  on the left, followed by an equals sign, a 5x5 grid labeled  $V$  in the middle, a minus sign, the word "convolve" in parentheses, a 3x3 stencil kernel with a blue center and purple neighbors, a comma, another 5x5 grid labeled  $U$ , and a closing parenthesis.

$$R = V - \sum U$$

The diagram shows a 5x5 grid labeled  $R$  on the left, followed by an equals sign, a 5x5 grid labeled  $V$  in the middle, a minus sign, a summation symbol  $\Sigma$ , and a 5x5 grid labeled  $U$  on the right. The grid  $U$  contains a 3x3 stencil kernel with a blue center and purple neighbors.



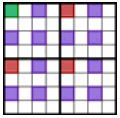
# Periodic Boundary Conditions

$$R = V - \sum U$$

*R*      *V*      *U*

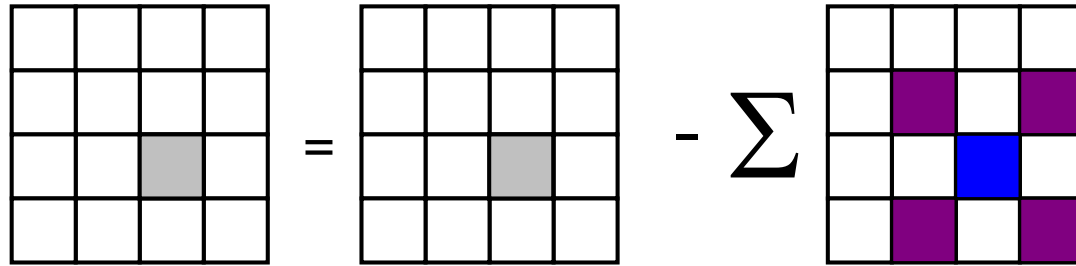
$$R = V - \sum U$$

*R*      *V*      *U*

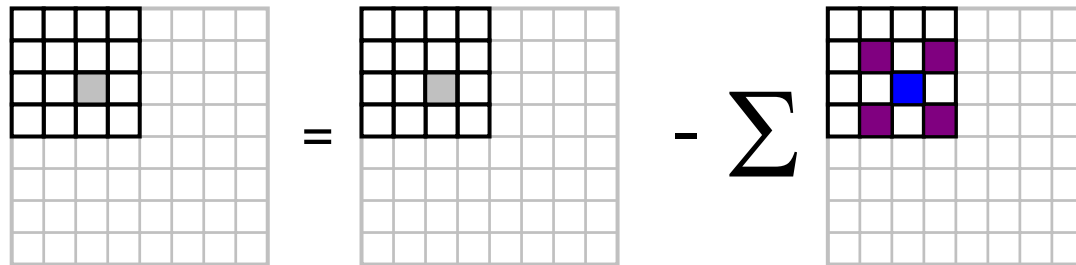


# At Other Levels of the Hierarchy

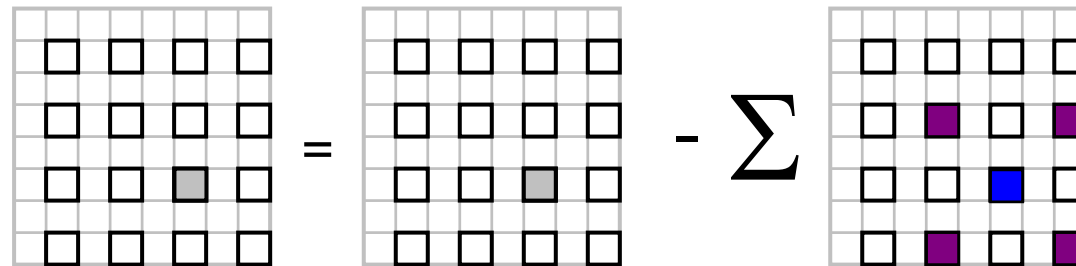
*conceptually:*



*dense indexing:*



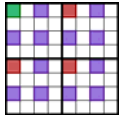
*strided indexing:*



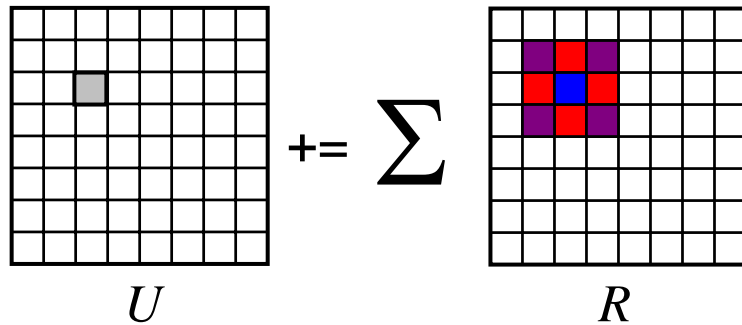
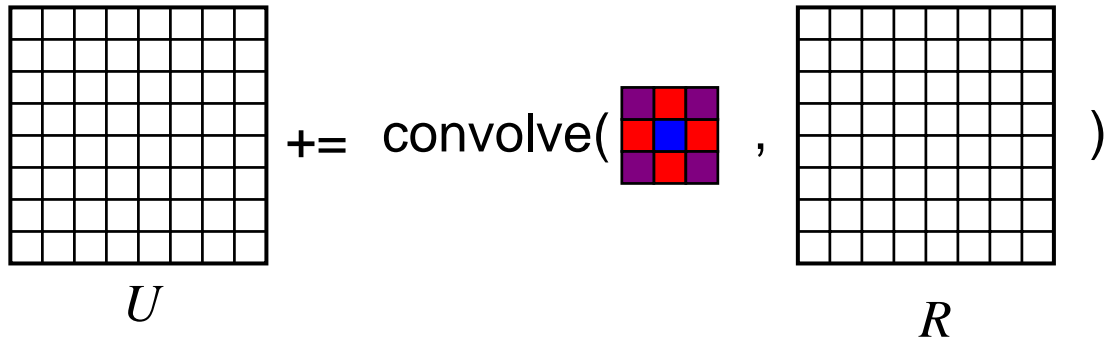
*R*

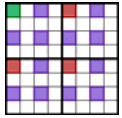
*V*

*U*

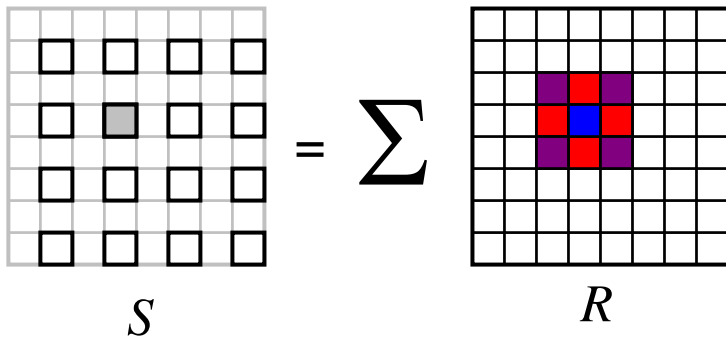
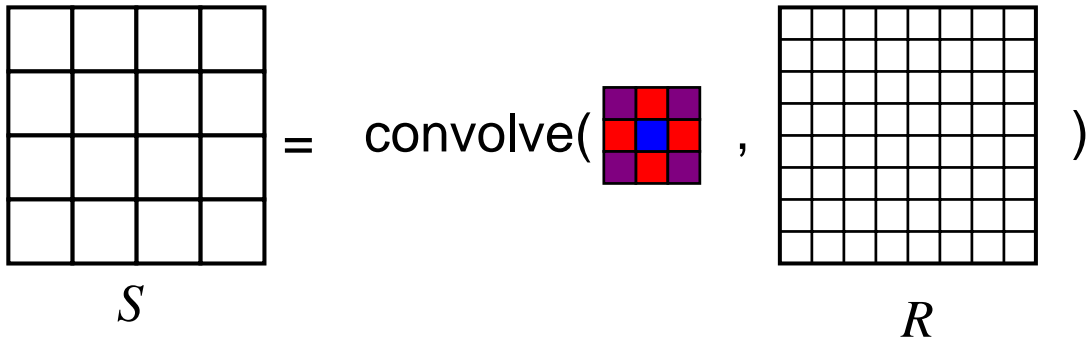


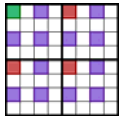
# $psinv(U, R)$



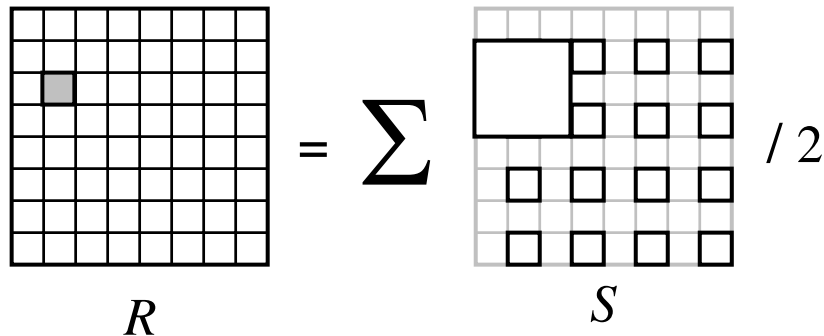
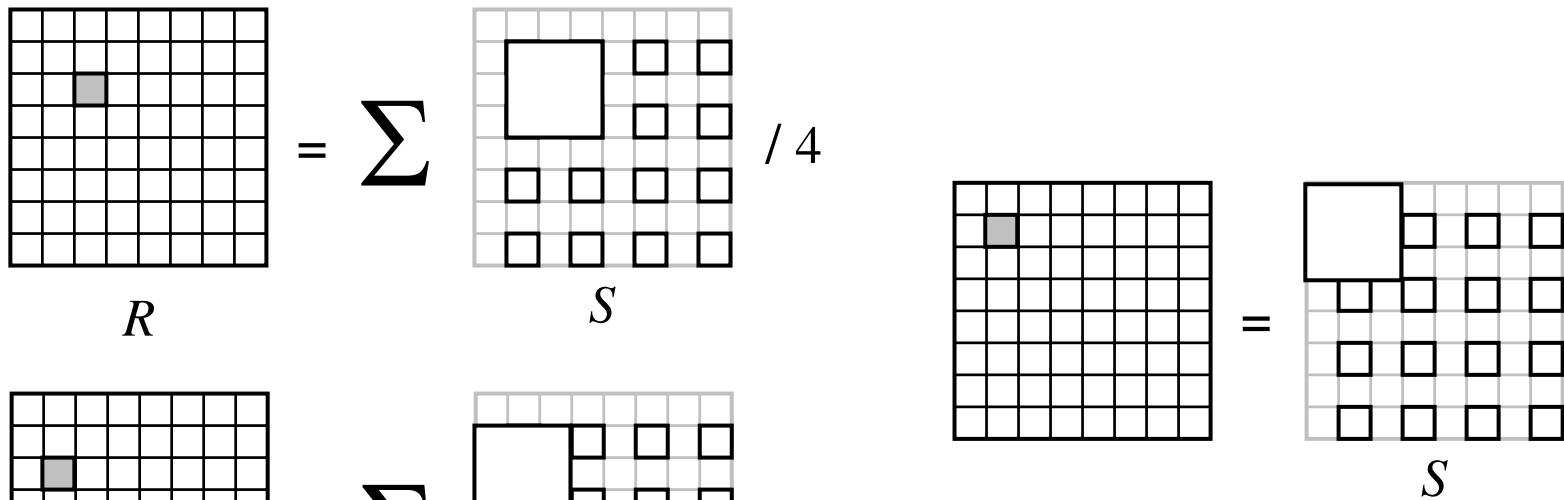
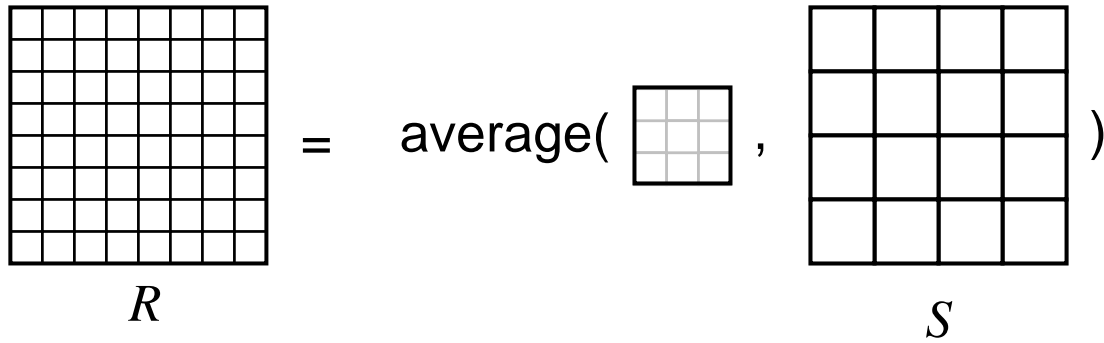


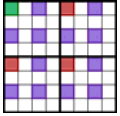
# $rprj3(S, R)$





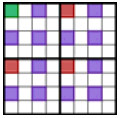
# *interp(R, S)*





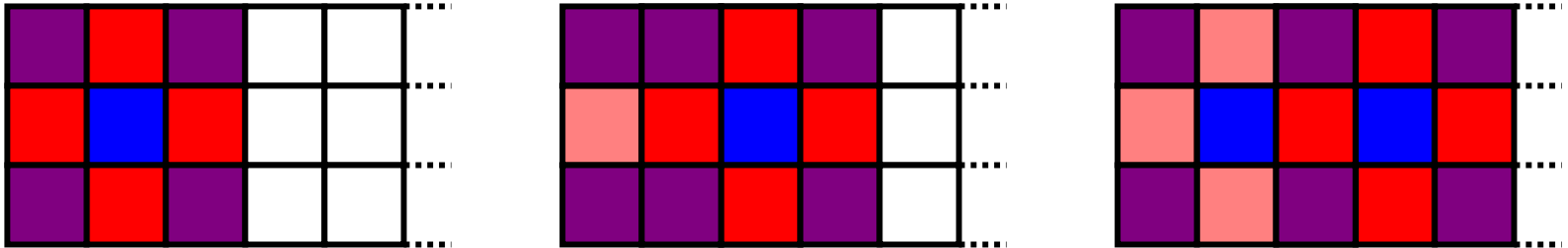
# rprj3 in Fortran

```
do j3=2,m3j-1
  i3 = 2*j3-d3
  do j2=2,m2j-1
    i2 = 2*j2-d2
    do j1=2,m1j-1
      i1 = 2*j1-d1
      s(j1,j2,j3) =
>         0.5D0 * r(i1,i2,i3)
>         + 0.25D0 * (r(i1-1,i2,i3) + r(i1+1,i2,i3)
>                   + r(i1, i2-1,i3 ) + r(i1, i2+1,i3 )
>                   + r(i1, i2, i3-1) + r(i1, i2, i3+1))
>         + 0.125D0 * (r(i1-1,i2-1,i3 ) + r(i1-1,i2+1,i3 )
>                   + r(i1-1,i2, i3-1) + r(i1-1,i2 ,i3+1)
>                   + r(i1+1,i2-1,i3 ) + r(i1+1,i2+1,i3 )
>                   + r(i1+1,i2, i3-1) + r(i1+1,i2 ,i3+1))
>                   + r(i1, i2-1,i3-1) + r(i1, i2-1,i3+1)
>                   + r(i1, i2+1,i3-1) + r(i1, i2+1,i3+1))
>         + 0.0625D0 * (r(i1-1,i2-1,i3-1) + r(i1-1,i2-1,i3+1)
>                   + r(i1-1,i2+1,i3-1) + r(i1-1,i2+1,i3+1)
>                   + r(i1+1,i2-1,i3-1) + r(i1+1,i2-1,i3+1)
>                   + r(i1+1,i2+1,i3-1) + r(i1+1,i2+1,i3+1))
      enddo
    enddo
  enddo
```

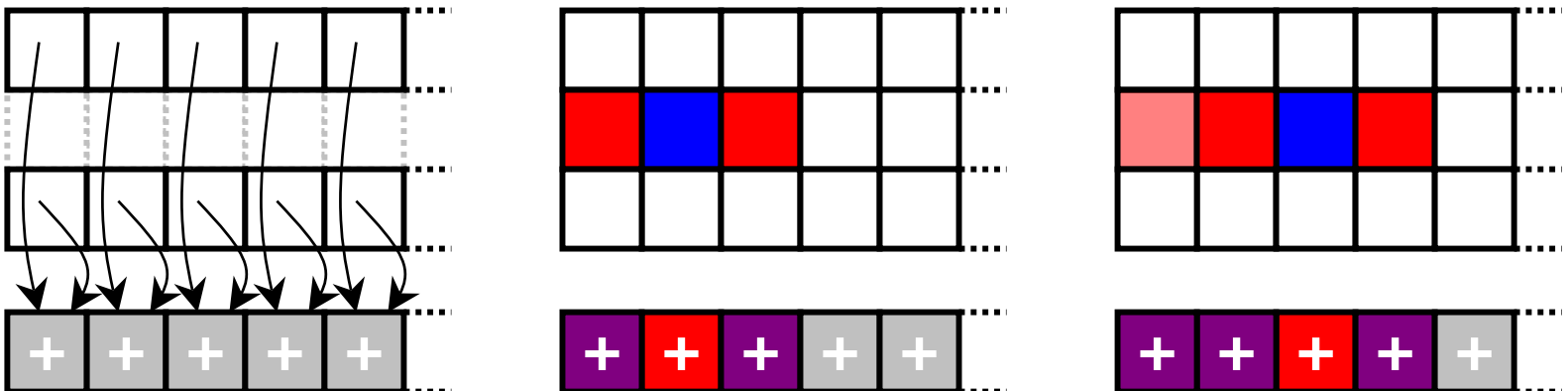


# Stencil Optimization (2D)

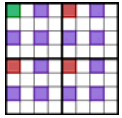
- Adjacent stencils use common subexpressions:



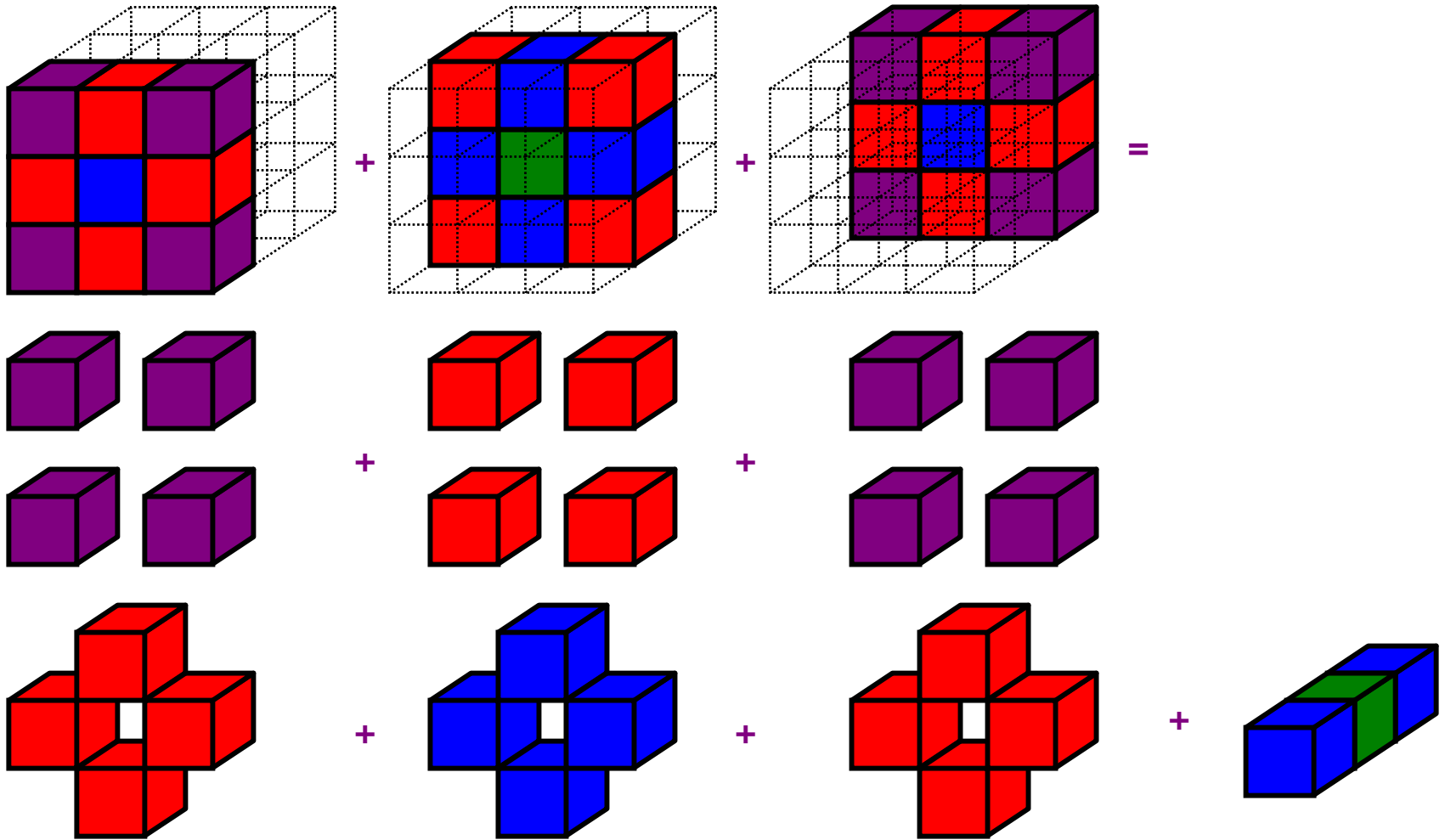
- Observation: Cache partial sums for reuse...

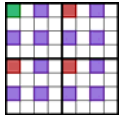


- Benefits are greater for 3D stencils...

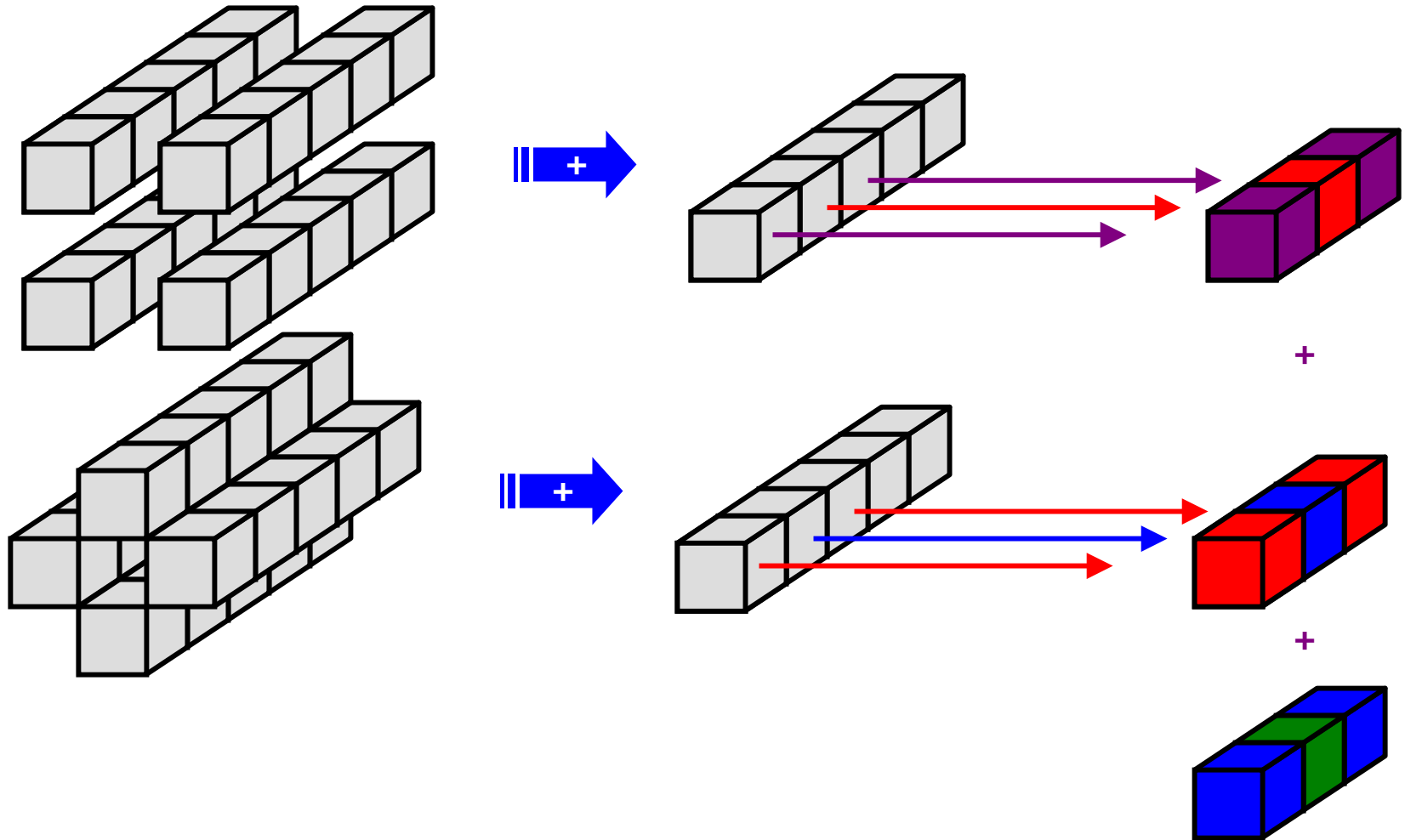


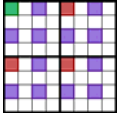
# MG Stencil Optimization





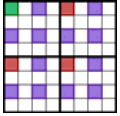
# MG Stencil Optimization





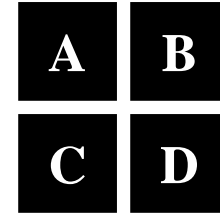
# *rprj3* in Fortran with stencil opt.

```
do j3=2,m3j-1
  i3 = 2*j3-d3
  do j2=2,m2j-1
    i2 = 2*j2-d2
    do j1=2,m1j
      i1 = 2*j1-d1
      x1(i1-1) = r(i1-1,i2-1,i3  ) + r(i1-1,i2+1,i3  )
>          + r(i1-1,i2,  i3-1) + r(i1-1,i2,  i3+1)
      y1(i1-1) = r(i1-1,i2-1,i3-1) + r(i1-1,i2-1,i3+1)
>          + r(i1-1,i2+1,i3-1) + r(i1-1,i2+1,i3+1)
    enddo
  do j1=2,m1j-1
    i1 = 2*j1-d1
    y2 = r(i1,  i2-1,i3-1) + r(i1,  i2-1,i3+1)
>      + r(i1,  i2+1,i3-1) + r(i1,  i2+1,i3+1)
    x2 = r(i1,  i2-1,i3  ) + r(i1,  i2+1,i3  )
>      + r(i1,  i2,  i3-1) + r(i1,  i2,  i3+1)
    s(j1,j2,j3) =
>      0.5D0 * r(i1,i2,i3)
>      + 0.25D0 * (r(i1-1,i2,i3) + r(i1+1,i2,i3) + x2)
>      + 0.125D0 * ( x1(i1-1) + x1(i1+1) + y2)
>      + 0.0625D0 * ( y1(i1-1) + y1(i1+1) )
  enddo
enddo
enddo
```

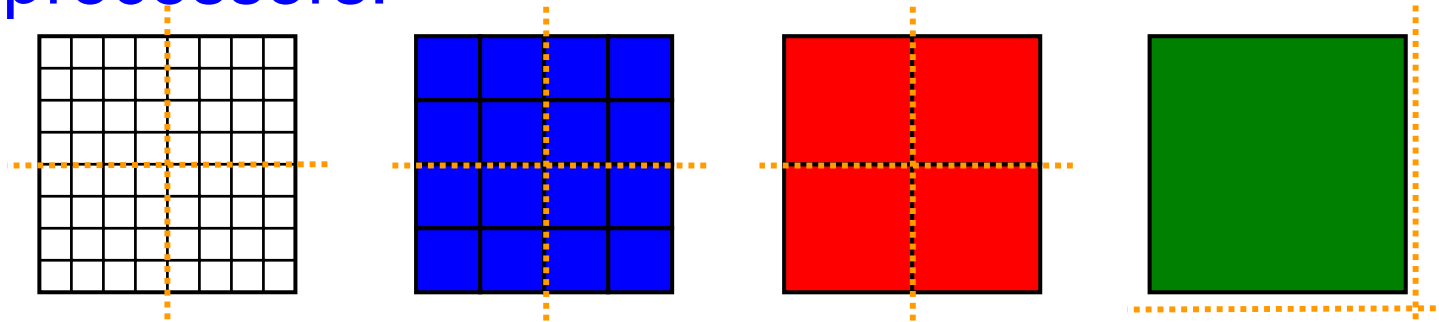


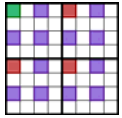
# Parallel Data Distribution

- Given a virtual processor grid...



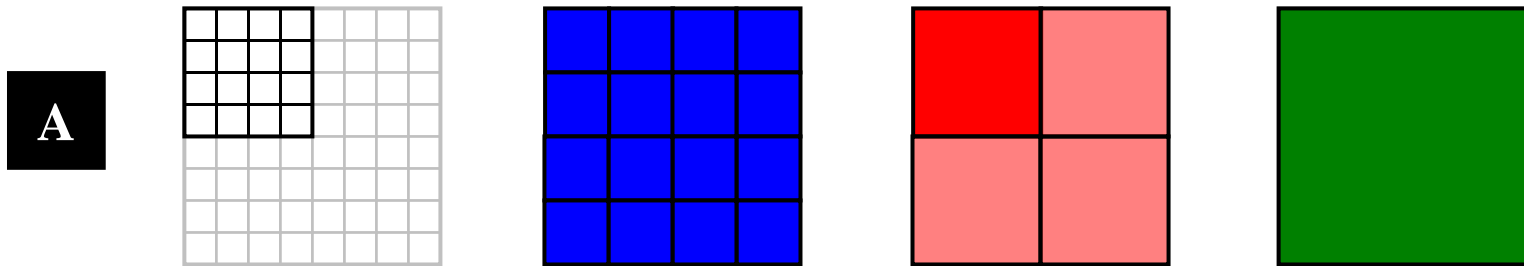
...arrays are block-distributed between processors:



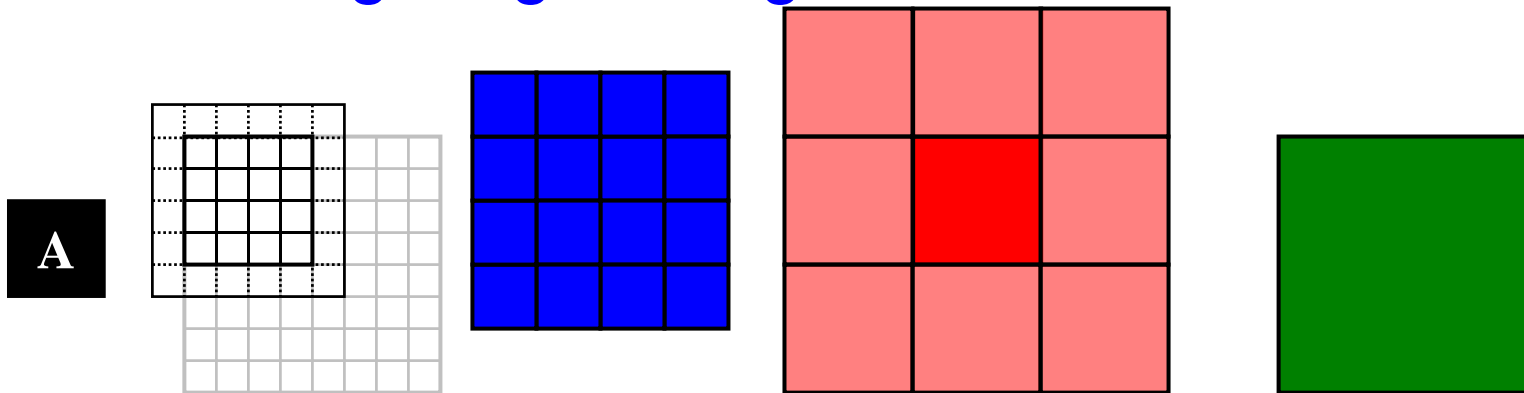


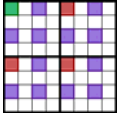
# Per-processor Data Allocation

- In addition to its local block of values...



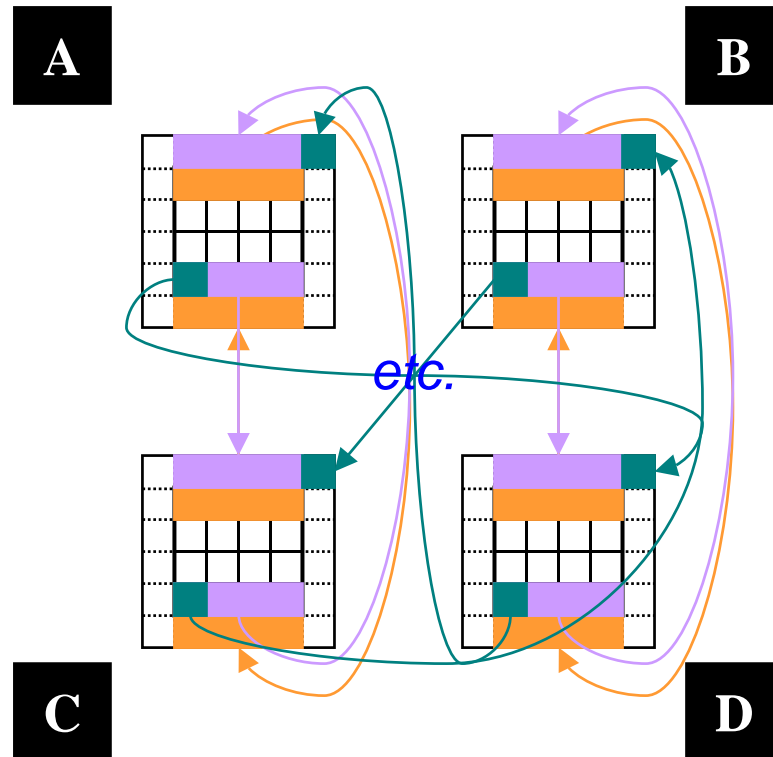
...each processor allocates ghost cells for caching neighboring values





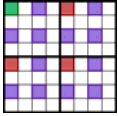
# Stencil Communication

Prior to computing a stencil, communication is typically required to refresh the ghost cells



Notes:

- Lots of optimization opportunities
- Have to eventually start skipping processors for coarser levels



# Distributed *rprj3* in Fortran

```
subroutine rprj3(r,m1k,m2k,m3k,s,m1j,m2j,m3j,k)
implicit none
include 'cafnpb.h'
include 'globals.h'

integer m1k, m2k, m3k, m1j, m2j, m3j,k

double precision r(m1k,m2k,m3k), s(m1j,m2j,m3j) >
integer j3, j2, j1, i3, i2, i1, d1, d2, d3, j >
double precision x1(m), y1(m), x2,y2 >

if(m1k.eq.3)then
  d1 = 2
else
  d1 = 1
endif

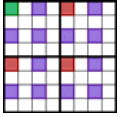
if(m2k.eq.3)then
  d2 = 2
else
  d2 = 1
endif

if(m3k.eq.3)then
  d3 = 2
else
  d3 = 1
endif

do j3=2,m3j-1
  i3 = 2*j3-d3
  do j2=2,m2j-1
    i2 = 2*j2-d2
    do j1=2,m1j
      i1 = 2*j1-d1
      x1(i1-1) = r(i1-1,i2-1,i3 ) + r(i1-1,i2+1,i3 )
                + r(i1-1,i2, i3-1) + r(i1-1,i2, i3+1)
      y1(i1-1) = r(i1-1,i2-1,i3-1) + r(i1-1,i2-1,i3+1)
                + r(i1-1,i2+1,i3-1) + r(i1-1,i2+1,i3+1)
    enddo
  do j1=2,m1j-1
    i1 = 2*j1-d1
    y2 = r(i1, i2-1,i3-1) + r(i1, i2-1,i3+1)
          + r(i1, i2+1,i3-1) + r(i1, i2+1,i3+1)
    x2 = r(i1, i2-1,i3 ) + r(i1, i2+1,i3 )
          + r(i1, i2, i3-1) + r(i1, i2, i3+1)
    s(j1,j2,j3) =
      0.5D0 * r(i1,i2,i3)
      + 0.25D0 * (r(i1-1,i2,i3) + r(i1+1,i2,i3) + x2)
      + 0.125D0 * ( x1(i1-1) + x1(i1+1) + y2)
      + 0.0625D0 * ( y1(i1-1) + y1(i1+1) )
  enddo
enddo
enddo
j = k-1
call comm3(s,m1j,m2j,m3j,j)
return
end
```

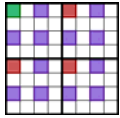




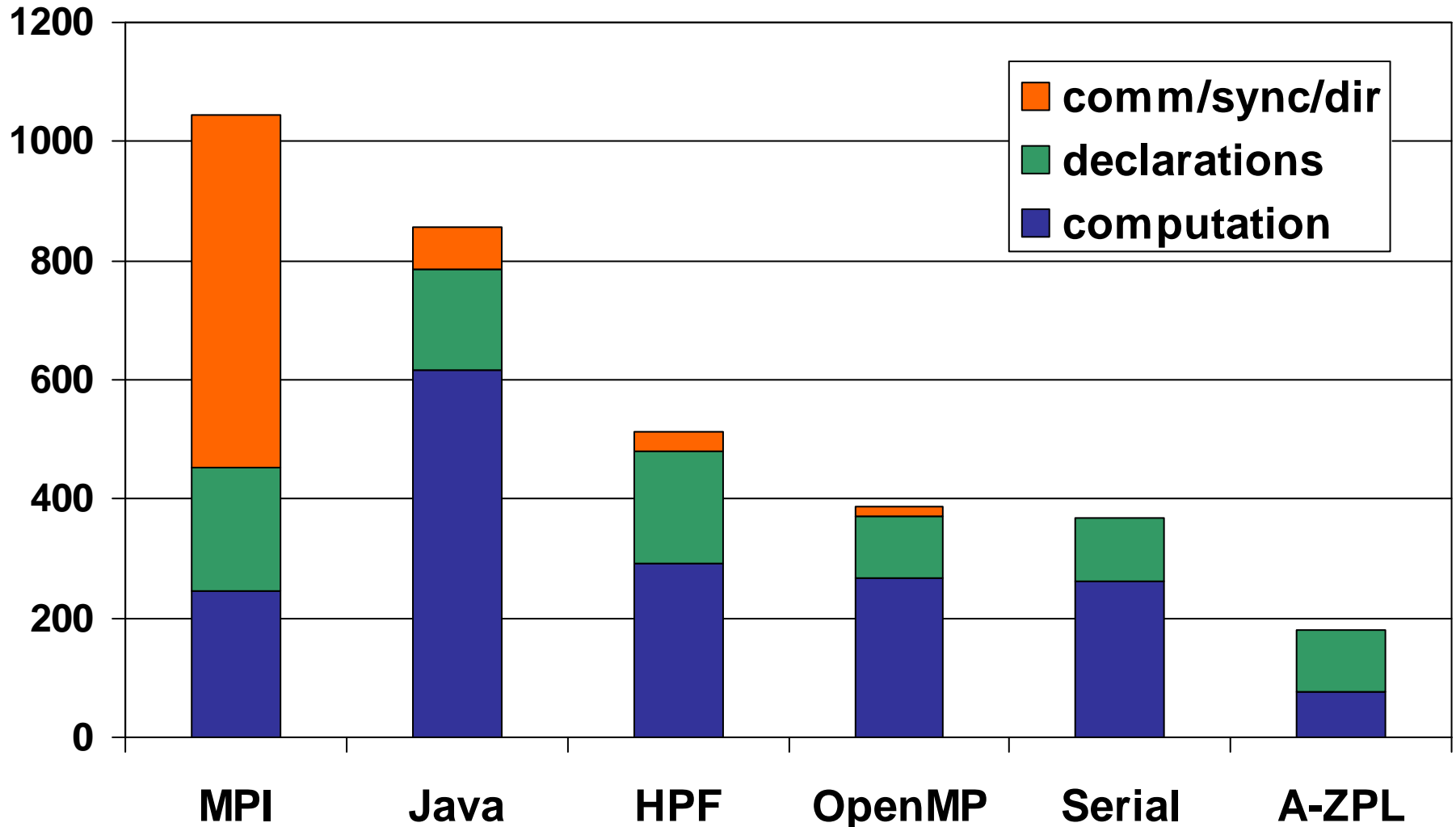


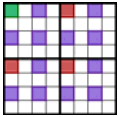
# rprj3 in A-ZPL

```
procedure rprj3(var S,R: [, ,] double;
                d: array [] of direction);
begin
  S := 0.5000 * R +
    0.2500 * (R@^d[ 1, 0, 0] + R@^d[ 0, 1, 0] + R@^d[ 0, 0, 1] +
              R@^d[-1, 0, 0] + R@^d[ 0,-1, 0] + R@^d[ 0, 0,-1] +
    0.1250 * (R@^d[ 1, 1, 0] + R@^d[ 1, 0, 1] + R@^d[ 0, 1, 1] +
              R@^d[ 1,-1, 0] + R@^d[ 1, 0,-1] + R@^d[ 0, 1,-1] +
              R@^d[-1, 1, 0] + R@^d[-1, 0, 1] + R@^d[ 0,-1, 1] +
              R@^d[-1,-1, 0] + R@^d[-1, 0,-1] + R@^d[ 0,-1,-1]) +
    0.0625 * (R@^d[ 1, 1, 1] + R@^d[ 1, 1,-1] +
              R@^d[ 1,-1, 1] + R@^d[ 1,-1,-1] +
              R@^d[-1, 1, 1] + R@^d[-1, 1,-1] +
              R@^d[-1,-1, 1] + R@^d[-1,-1,-1]);
end;
```



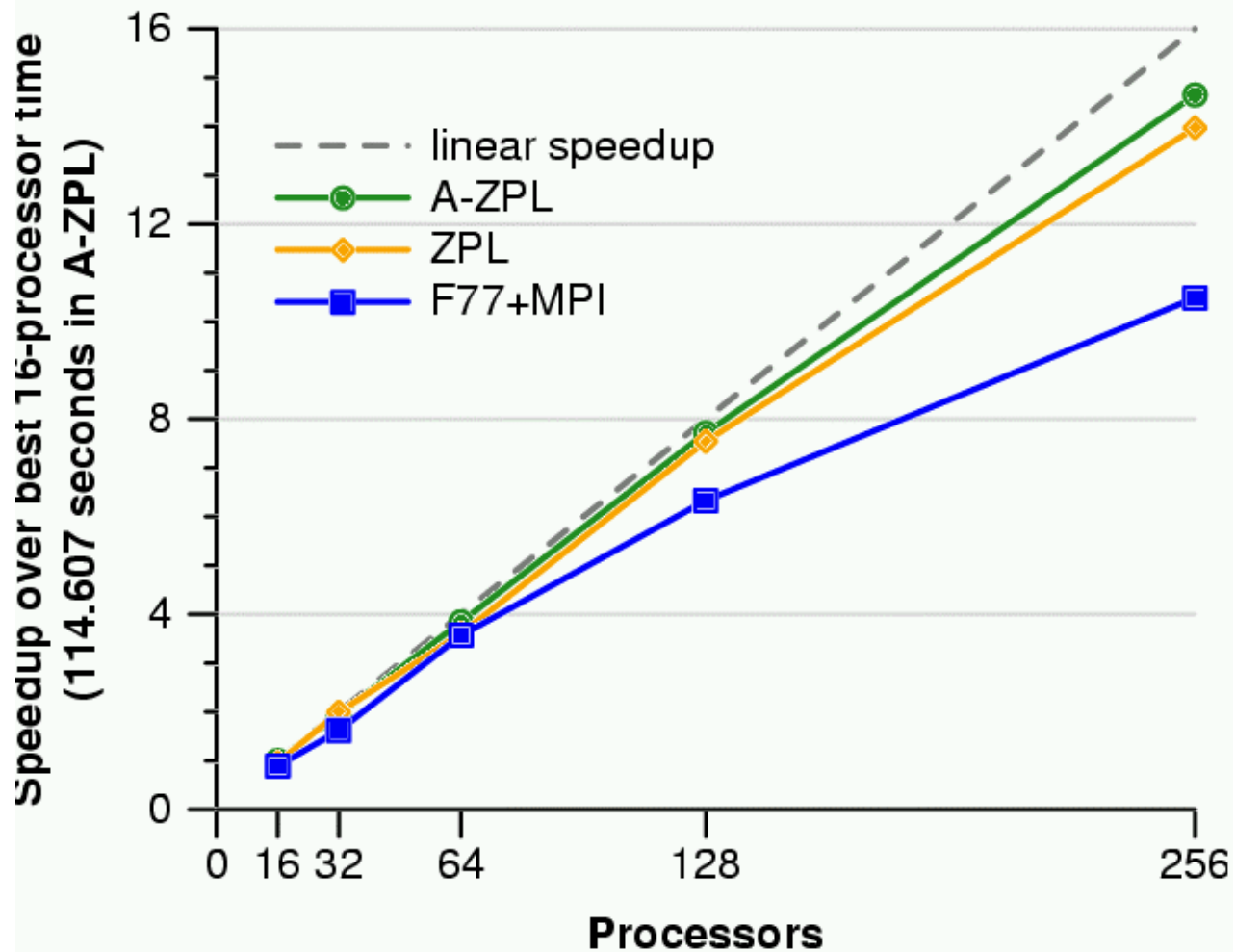
# NAS MG Linecounts

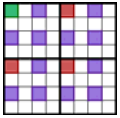




# MG Performance

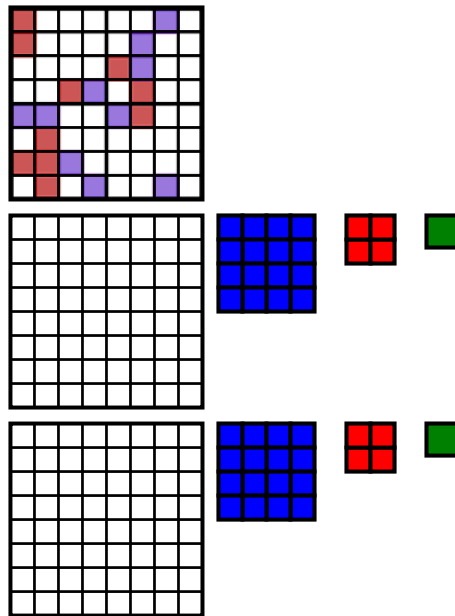
## MG Class C -- Cray T3E



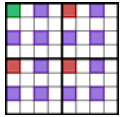


# Sparse Arrays in MG

- MG's input array,  $V$ , has 20 non-zeroes regardless of problem size

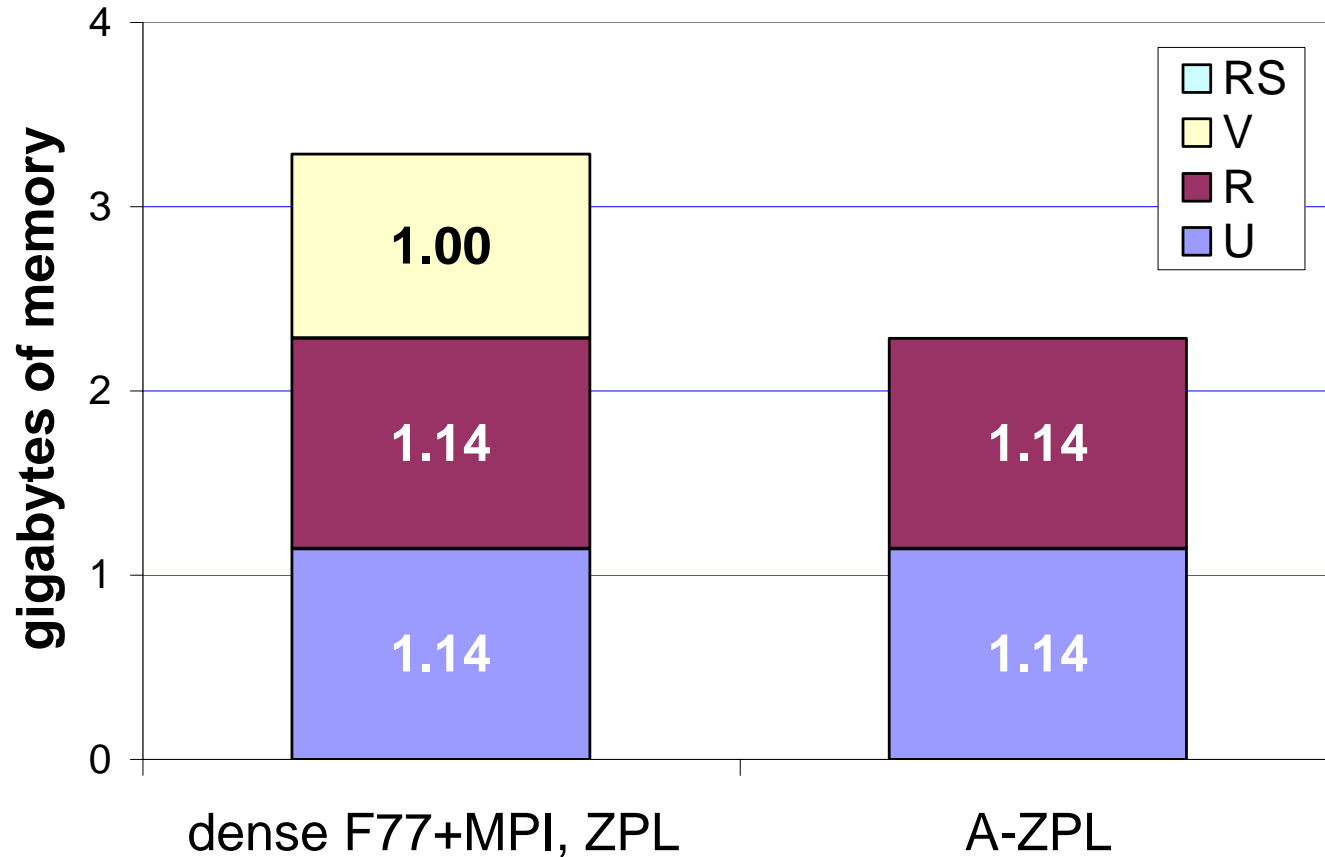


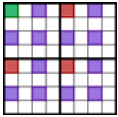
- Computes resid against  $V$  twice per iteration
- Wasted time and space



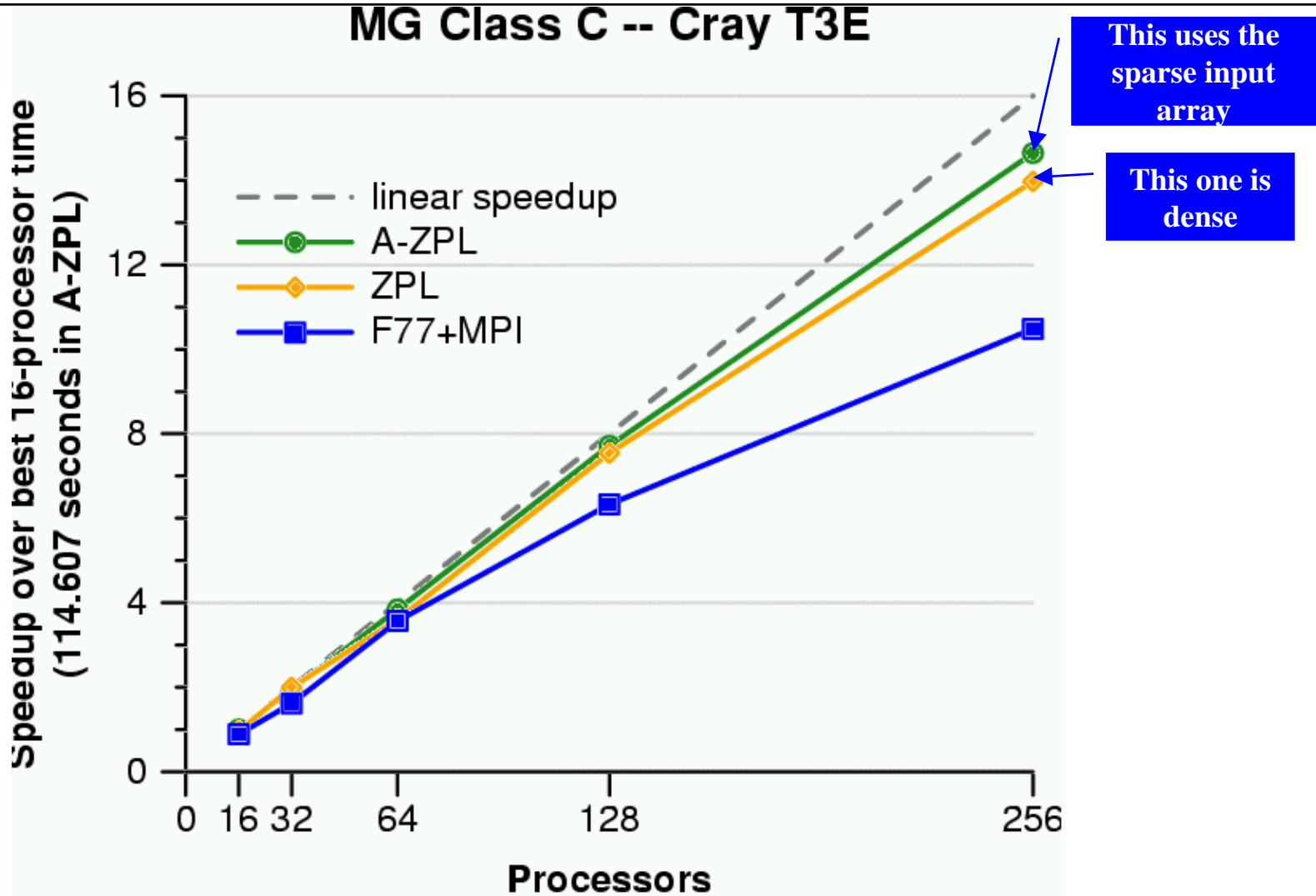
# NAS MG Memory Usage

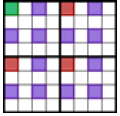
## MG Class C -- memory usage





# MG Performance





# Summary

---

- MG is a useful application
  - models real-world computations, parallel idioms
  - small enough to be easy to study
- Not so hard to describe in pictures and words
- Coding up a parallel version of it tends to be a different matter
- Surely we can do better than the current crop of parallel languages