

- For global version, radial wavevectors replaced by finite differences and multi-point gyro-average: $G = G^{\dagger} = J_0(k_{\perp}(\mu)\rho_s)$
local spectral
- Significant changes to field solve & gyroaverage $(G(\mathbf{x}, \mu, \text{sp})f(\mathbf{x}))_i = G_{ij}f_j$
non-spectral

$$\sum_{\text{sp}} \int d^3\mathbf{v} \left[Z_s e G f(\mathbf{x}) + F_M \frac{Z_s^2 e^2}{T_s} (G G^{\dagger} - 1) \phi(\mathbf{x}) \right] = 0$$

$$I(k_y, x_i, s) + P_{ij} \phi(k_y, x_j, s) = 0$$

- Advantages: Global profiles,
another direction for domain decomposition:

- 6D distribution
- 3D fields
- 5D gyro-averaged fields
- 4D collisions conservation field
- 5D domain decomposition
- 5D MPI Cartesian communicator

$$f(k_y, x, s, \mu, v_{\parallel}, \text{sp})$$

$$\phi(k_y, x, s)$$

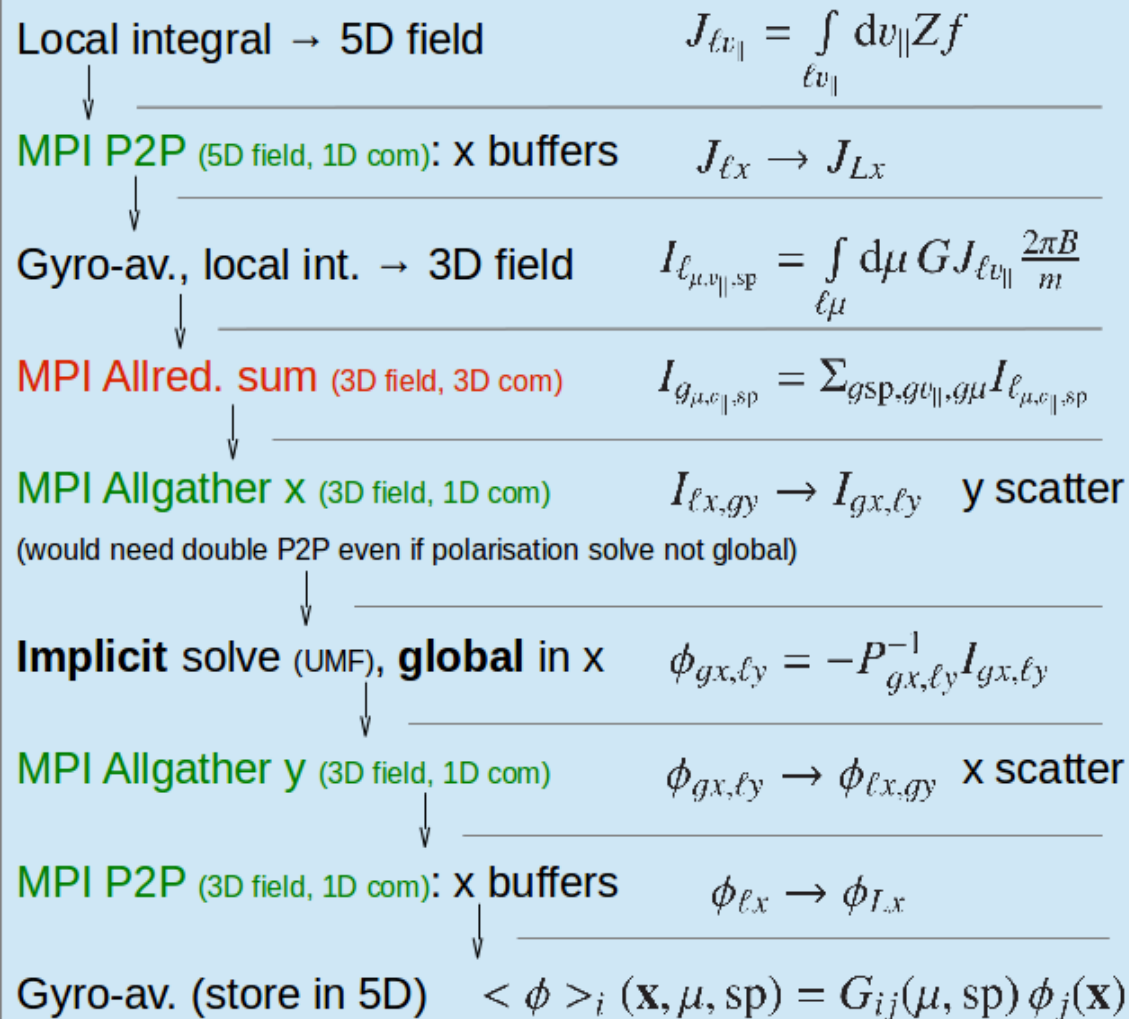
$$\langle \phi \rangle (k_y, x, s, \mu, \text{sp})$$

$$x, s, \mu, v_{\parallel}, \text{sp}$$

$$x, s, \mu, v_{\parallel}, \text{sp}$$

Field solve becomes complex, multi-stage MPI

Non-spectral parallel field solve



- Communicate radial buffer cells before each gyroav.
- Polarization solve implicit
 - Matrix $\sim Nx^2$
 - Uses UMFPack library
 - No radial decomposition
 - Instead parallelise over binormal modes, using Cartesian direction usually used for velocity
 - Allows good scaling
 - MPI Allgather / AllScatter
- 5 MPI communications total
 - Only Allreduce is costly
- Limitation: Gyroradii may cross ≤ 1 proc. boundary

Field solve becomes complex, multi-stage MPI

Non-spectral parallel field solve		28.9
Local integral → 5D field	$J_{\ell v_{\parallel}} = \int_{\ell v_{\parallel}} dv_{\parallel} Z f$	2.9
MPI P2P (5D field, 1D com): x buffers	$J_{\ell x} \rightarrow J_{Lx}$	1.3
Gyro-av., local int. → 3D field	$I_{\ell \mu, v_{\parallel}, sp} = \int_{\ell \mu} d\mu G J_{\ell v_{\parallel}} \frac{2\pi B}{m}$	3.8
MPI Allred. sum (3D field, 3D com)	$I_{g \mu, c_{\parallel}, sp} = \sum_{gsp, gv_{\parallel}, g\mu} I_{\ell \mu, c_{\parallel}, sp}$	5.4
MPI Allgather x (3D field, 1D com) (would need double P2P even if polarisation solve not global)	$I_{\ell x, gy} \rightarrow I_{gx, ly}$ y scatter	4.5
Implicit solve (UMF), global in x	$\phi_{gx, ly} = -P_{gx, ly}^{-1} I_{gx, ly}$	7.3
MPI Allgather y (3D field, 1D com)	$\phi_{gx, ly} \rightarrow \phi_{\ell x, gy}$ x scatter	1.3
MPI P2P (3D field, 1D com): x buffers	$\phi_{\ell x} \rightarrow \phi_{Ix}$	0.3
Gyro-av. (store in 5D)	$\langle \phi \rangle_i(\mathbf{x}, \mu, sp) = G_{ij}(\mu, sp) \phi_j(\mathbf{x})$	0.2

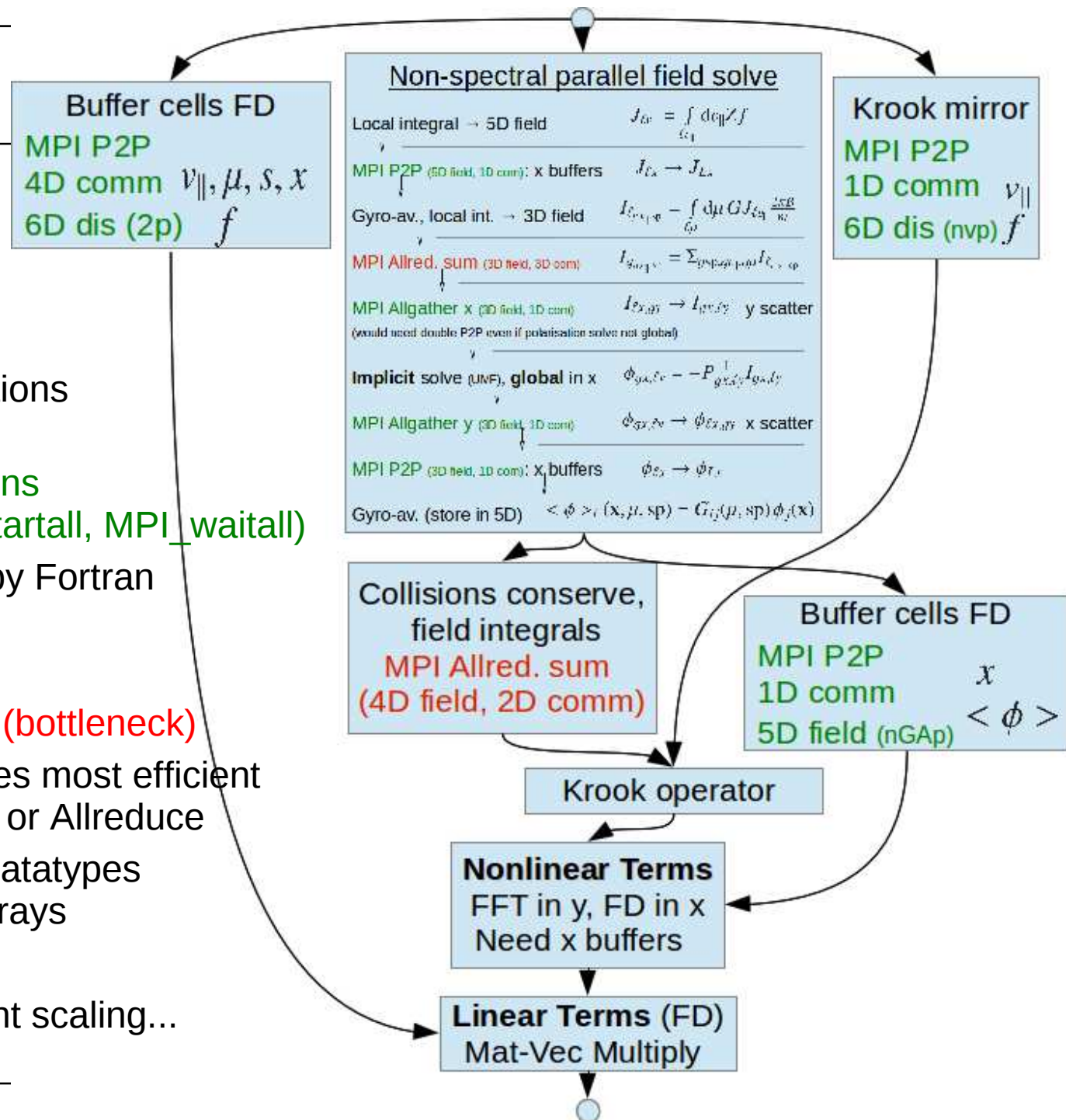
Percentages of total runtime for large electromagnetic collisional case on 32,768 processes

Important for good MPI scaling that the runtime fraction of the field solve does not rise

All MPI ops

Single Explicit Step

- Overlap communications with computation
- 7 P2P communications nonblocking (MPI_Startall, MPI_waitall)
 - Each managed by Fortran structure
- 2 MPI Allgather ops
- 2 MPI Allreduce ops (bottleneck)
- Parallel direction gives most efficient scaling: No Integrals or Allreduce
- 12 x 2 MPI derived datatypes for different buffer arrays
- All this allows efficient scaling...

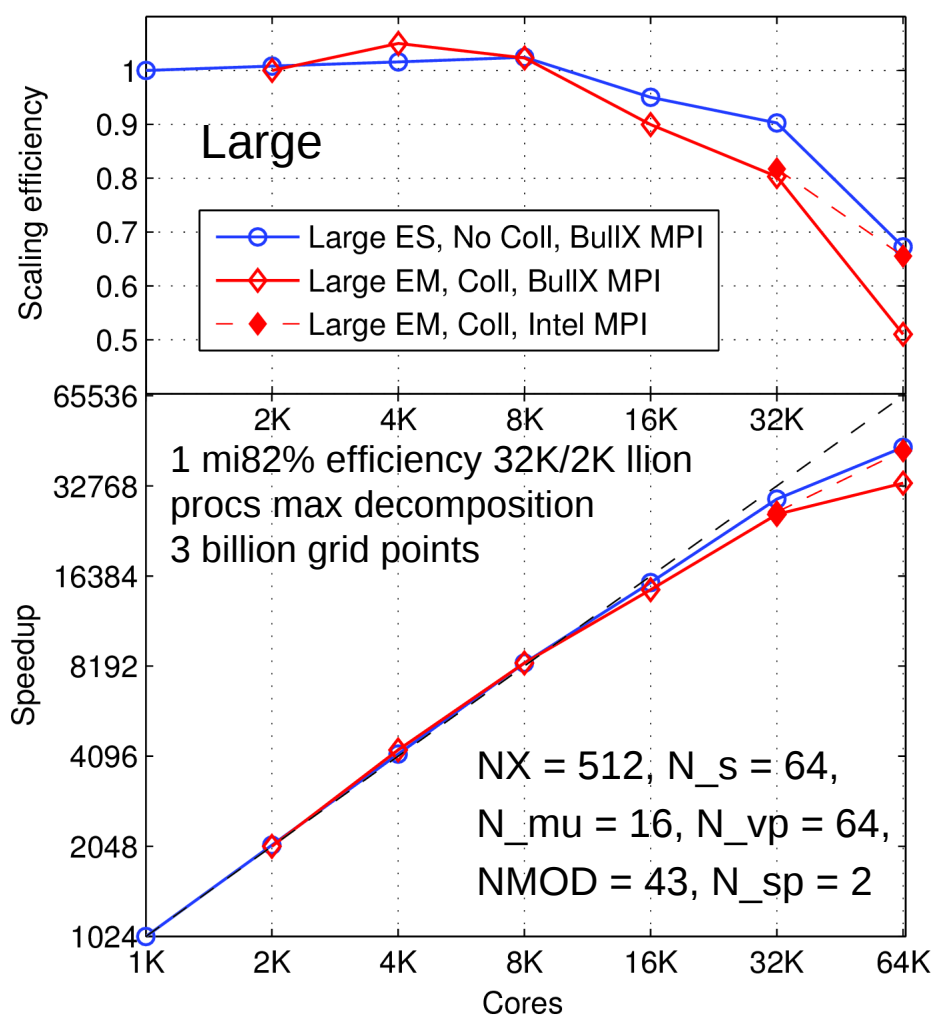
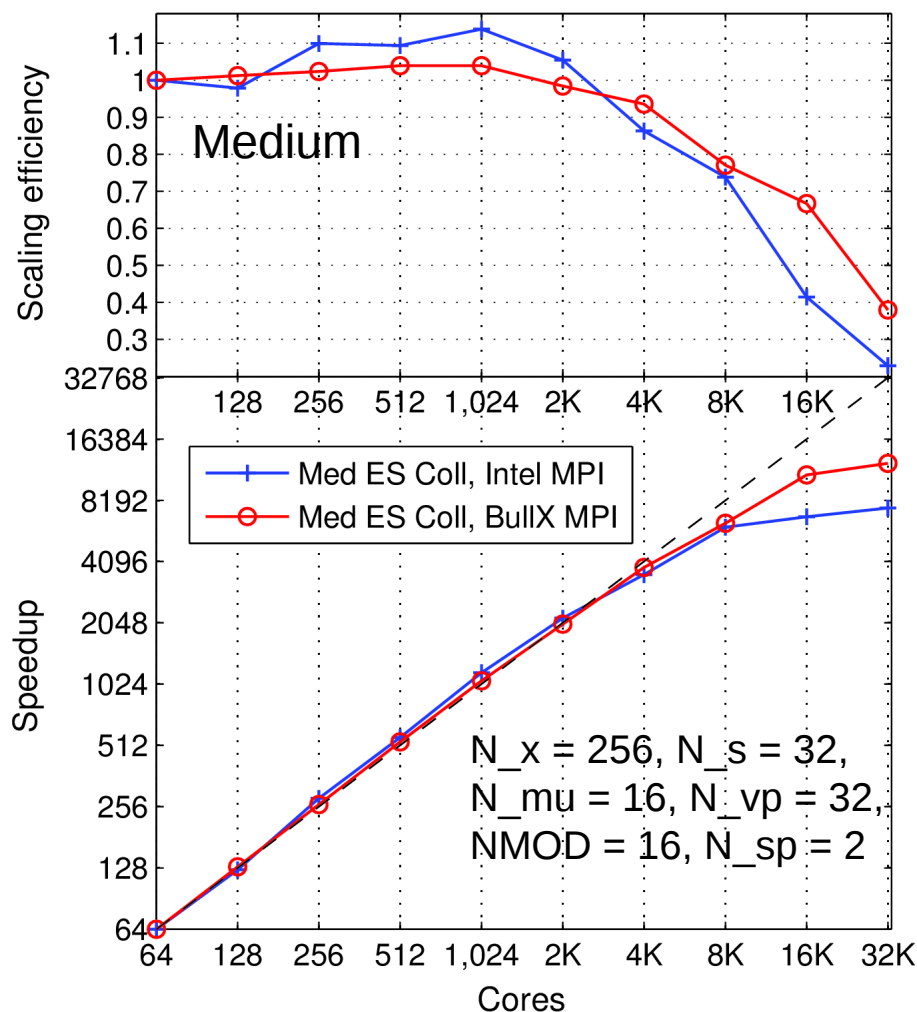


Strong Scaling Performance (HELIOS)

For typical physics problems: electromagnetic collisions (no conservation, no Krook)

95% efficiency 4K / 64

82% efficiency 32K / 2K



Comments on the HELIOS scaling

- Omit timing of the first loop (order of seconds for first MPI calls)
- Scaling breaks down when field solve fraction of runtime increases
 - Inefficient Collectives: large **Allreduce** and **Allgather**
 - Use MPI_STATUSES_IGNORE, MPI_IN_PLACE

- Intel MPI and BullX MPI perform differently
 - Intel MPI manages cache better affects memory bandwidth ?
 - Intel MPI better at largest collectives
 - BullX good at smaller # of cores
 - Many environment variables to tweak for largest jobs
 - No general optimal settings
 - Each case / code is different
 - Intel MPI tune might help ?
 - OMP_COLL_OPT

Large, time (s): 1600 iterations	Bull MPI		Intel MPI	
	32K	64K	32K	64K
Field Solve	21	25	26	22
Non Linear	13	07	12	06
Buffer cells f	09	11	10	05
Linear Term	21	11	17	08
Total	70	56	68	43
Efficiency	0.81	0.52	0.82	0.66

- Explicit Eulerian gyro-kinetics with finite differences on appropriate grids allows parallel domain decomposition in 5 dimensions
 - Permits efficient exploitation of the largest present day supercomputers
 - Limited by available resources, not code capability
 - Requires complex multi-dimensional MPI
 - multiple communicators
 - multiple derived datatypes
 - persistent non blocking communications
(some person-years of development and testing)
 - Allows us to do new and interesting physics
(e.g. NTM simulations in first part of the talk)
-