# HPC IN LIFE SCIENCES: MOLECULAR DYNAMICS

1. CAVEATS FOR COMPLEX MODELS
2. PARALLELIZING & ACCELERATING REAL CODES
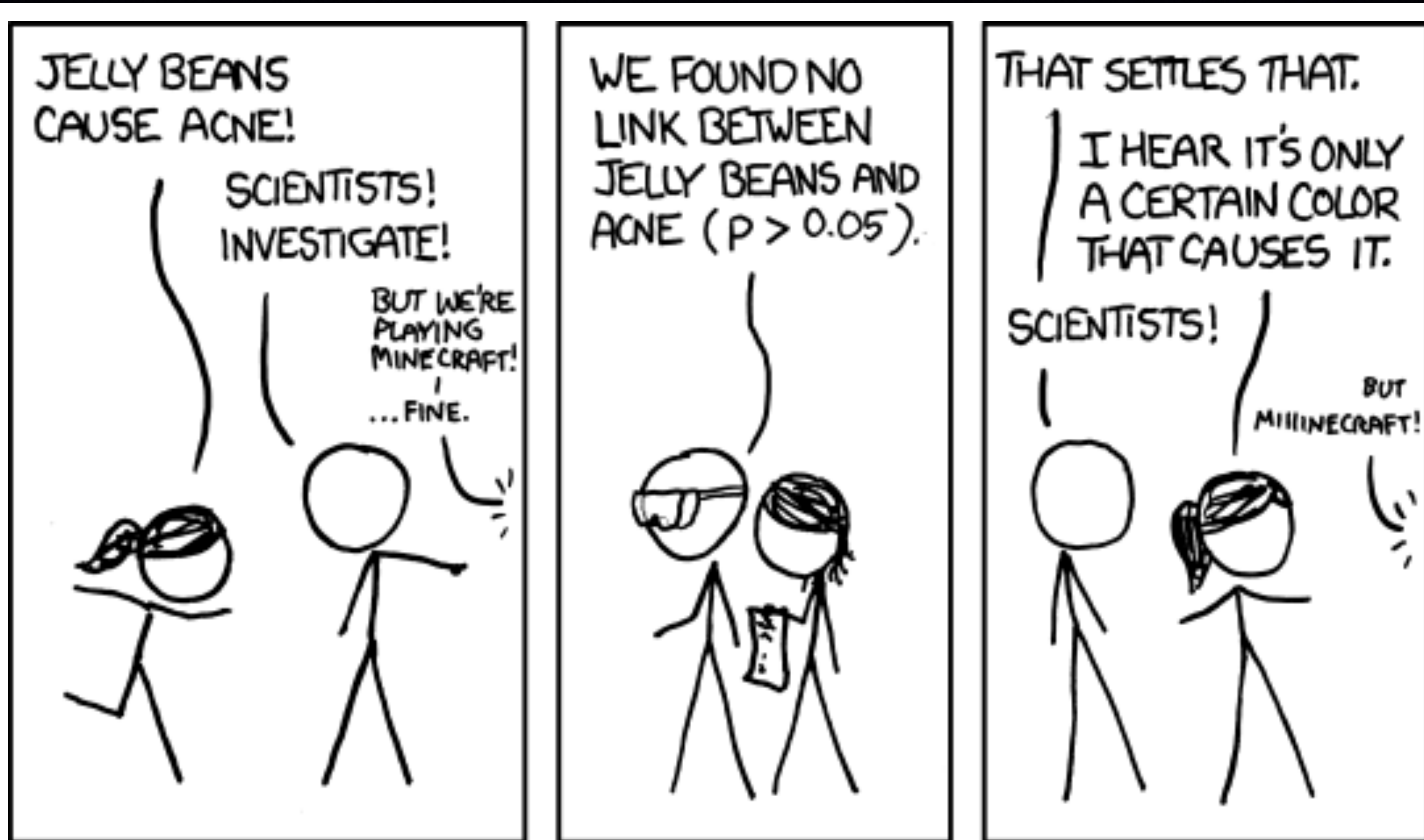3. HETEROGENEOUS (CPU+GPU) ACCELERATION

Erik Lindahl, Stockholm University, Stockholm University <erik.lindahl@scilifelab.se>

High accuracy

Low accuracy

Unreliable, but valid

Reliable & valid

Unreliable & not valid

Reliable, but not valid

Low precision

High precision

# P-value hacking
https://www.xkcd.com/882/

# P-value hacking - for real

## RESEARCH ARTICLE

**Open Access**

# Change in paternal grandmothers´ early food supply influenced cardiovascular mortality of the female grandchildren

### Abstract

**Background:** This study investigated whether large fluctuations in food availability during grandparents' early development influenced grandchildren's cardiovascular mortality. We reported earlier that changes in availability of food - from good to poor or from poor to good - during intrauterine development was followed by a double risk of sudden death as an adult, and that mortality rate can be associated with ancestors´ childhood availability of food. We have now studied transgenerational responses (TGR) to sharp differences of harvest between two consecutive years´ for ancestors of 317 people in Överkalix, Sweden.

**Results:** The confidence intervals were very wide but we found a striking TGR. There was no response in

Ceci n'est pas une pipe.

This is not a protein

$$\left[\frac{-\hbar^2}{2m}\nabla^2 + V\right]\Psi = i\hbar\frac{\partial}{\partial t}\Psi$$

F=ma

# Experiments

Efficient averaging

Less detail

**Where we need to be**

Physics — Chemistry — Biology

$10^{-15}$ s    $10^{-12}$ s    $10^{-9}$ s    $10^{-6}$ s    $10^{-3}$ s    $10^{0}$ s    $10^{3}$ s

**Where we are**

**Where we want to be**

# Simulations

Extreme detail

Sampling issues?

Parameter quality?

Another tRNA moves into the A site in order to add another amino acid to the peptide chain.

# Challenge: MD is intrinsically a *sequential* problem

Initial input data:
Interaction function V(**r**) - "force field"
coordinates **r**, velocities **v**

Compute potential V(**r**) and
forces $\mathbf{F}_i = \nabla_i V(r)$ on atoms

Update coordinates &
velocities according to
equations of motion

Collect statistics and write
energy/coordinates to
trajectory files

More steps?

Yes

No

Done!

Repeat for millions of steps

$$V(r) = \sum_{bonds} \frac{1}{2} k_{ij}^b \left( r_{ij} - r_{ij}^0 \right)^2$$

$$+ \sum_{angles} \frac{1}{2} k_{ijk}^\theta \left( \theta_{ijk} - \theta_{ijk}^0 \right)^2$$

$$+ \sum_{torsions} \left\{ \sum_n k_\theta \left[ 1 + \cos\left( n\phi - \phi_0 \right) \right] \right\}$$

$$+ \sum_{impropers} k_\xi \left( \xi_{ijkl} - \xi_{ijkl}^0 \right)$$

$$+ \sum_{i,j} \frac{q_i q_j}{4\pi\varepsilon_0 r_{ij}}$$

$$+ \sum_{i,j} \left[ \frac{C_{12}}{r_{ij}^{12}} - \frac{C_6}{r_{ij}^6} \right]$$

Costly, because these
terms involve all pairs

$$m_i \frac{\partial^2 r_i}{\partial t^2} = F_i \qquad i = 1..N$$

$$F_i = -\frac{\partial V(r)}{\partial r_i}$$

With $\Delta$t ~ 1fs and µs to s
timescales of interest, we
need $10^9$-$10^{15}$ steps.

The challenge:
- ~100,000 atoms
- Each has ~500 neighbors
  - Maintain a list of them, update ea 10 steps
- ~50M interactions/step
- ~2B FLOPS per step
- ~1ms real time per step

```c
for(k=nj0; (k<nj1); k++)
{

    /* Get j neighbor index, and coordinate index */
    jnr           = jjnr[k];
    j3            = 3*jnr;

    /* load j atom coordinates */
    jx1           = pos[j3+0];
    jy1           = pos[j3+1];
    jz1           = pos[j3+2];

    /* Calculate distance */
    dx11          = ix1 - jx1;
    dy11          = iy1 - jy1;
    dz11          = iz1 - jz1;
    rsq11         = dx11*dx11+dy11*dy11+dz11*dz11;

    /* Calculate 1/r and 1/r2 */
    rinv11        = 1.0/sqrt(rsq11);

    /* Load parameters for j atom */
    qq            = iq*charge[jnr];
    tj            = nti+2*type[jnr];
    c6            = vdwparam[tj];
    c12           = vdwparam[tj+1];
    rinvsq        = rinv11*rinv11;

    /* Coulomb interaction */
    vcoul         = qq*rinv11;
    vctot         = vctot+vcoul;

    /* Lennard-Jones interaction */
    rinvsix       = rinvsq*rinvsq*rinvsq;
    Vvdw6         = c6*rinvsix;
    Vvdw12        = c12*rinvsix*rinvsix;
    Vvdwtot       = Vvdwtot+Vvdw12-Vvdw6;
    fscal         = (vcoul+12.0*Vvdw12-6.0*Vvdw6)*rinvsq;

    /* Calculate temporary vectorial force */
    tx            = fscal*dx11;
    ty            = fscal*dy11;
    tz            = fscal*dz11;

    /* Increment i atom force */
    fix1          = fix1 + tx;
    fiy1          = fiy1 + ty;
    fiz1          = fiz1 + tz;

    /* Decrement j atom force */
    faction[j3+0]    = faction[j3+0] - tx;
    faction[j3+1]    = faction[j3+1] - ty;
    faction[j3+2]    = faction[j3+2] - tz;

    /* Inner loop uses 38 flops/iteration */

}
```
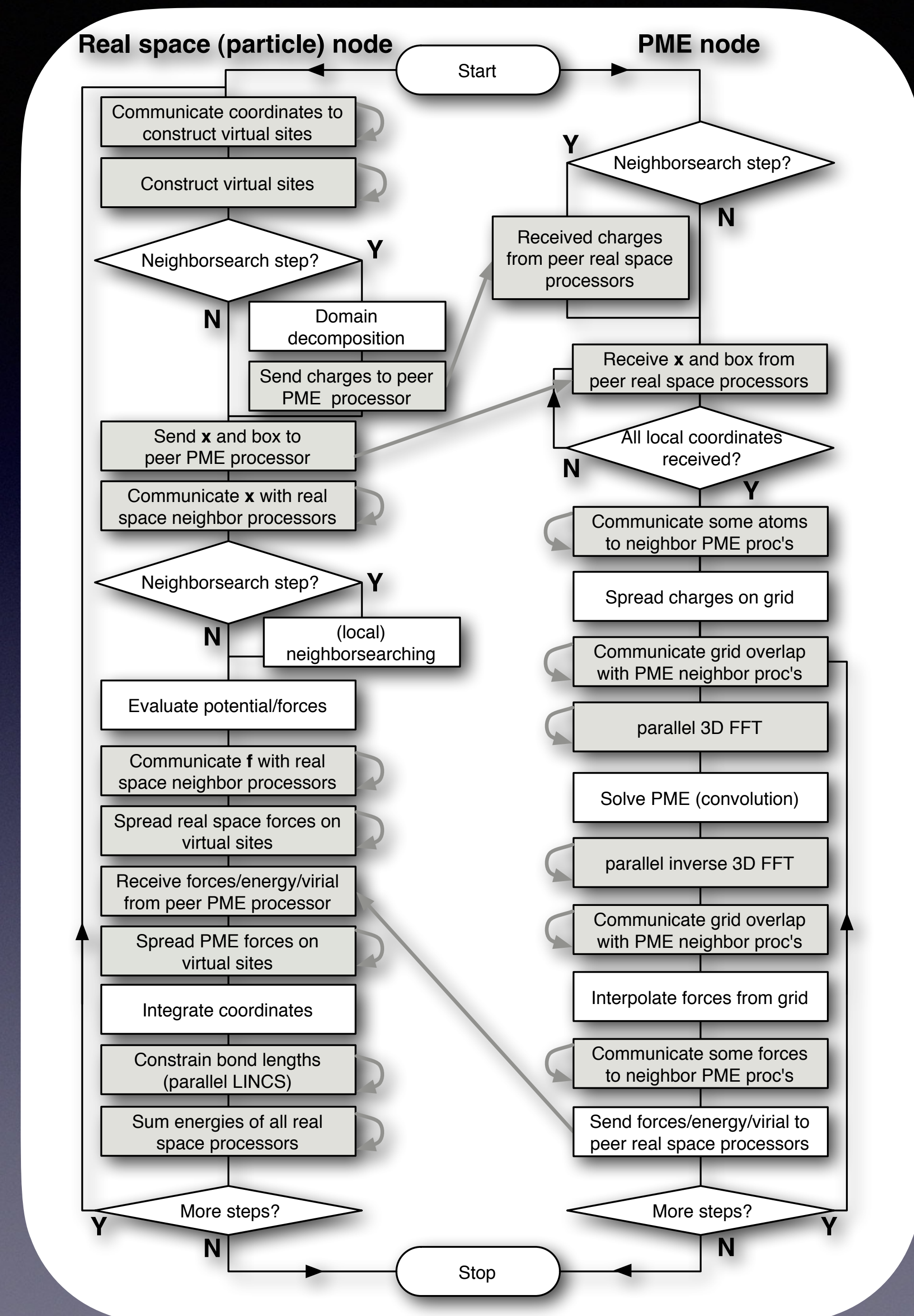
# The things we do every ~100 μs

**To-Do Monday 09:1 5:48.004.1 00 (simple version)**

Adjust domain decomposition [communicate]
Communicate coordinates to/from 26 neighbor nodes
Find atoms in proximity [communicate]
Change charges or parameters for free energy
Create local virtual particles [communicate]
Send coordinates to GPU
Calculate short-range electrostatics & VdW
Calculate bonds
Calculate angles
Calculate torsions
Perform long-range lattice summation [communicate]
Apply external fields/forces
Get forces back from GPU
Send forces to 26 neighbors [communicate]
Integrate new positions
Constrain bonds [communicate]
Update stats. (temperature, energy) [communicate]
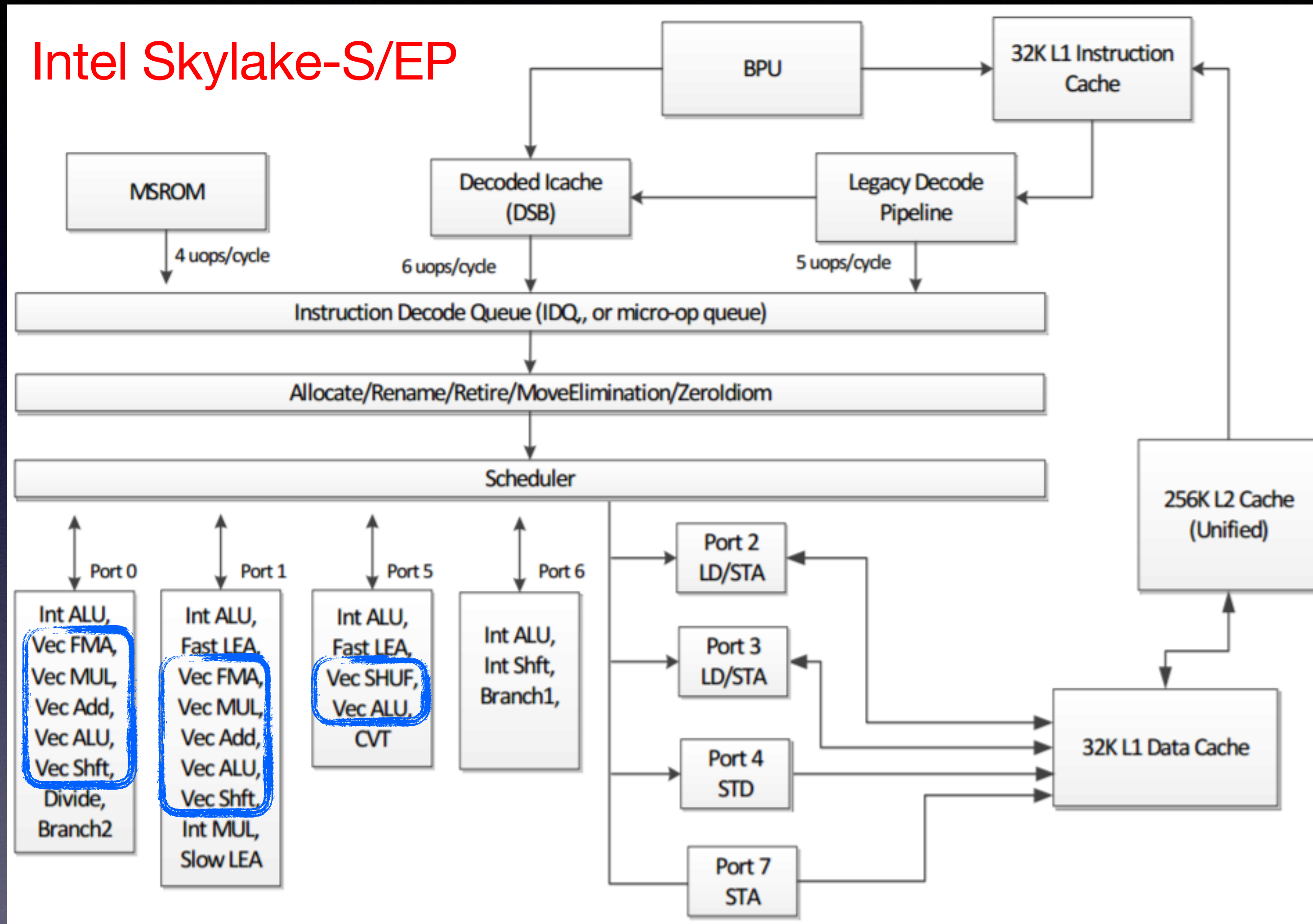Write coordinates/forces if necessary

A fairly typical HPC application - complex & fast

Every arrow is communication



**Real space (particle) node** — **PME node**

Start

Communicate coordinates to construct virtual sites
Construct virtual sites
Neighborsearch step? — N / Y
Domain decomposition
Send charges to peer PME processor
Send **x** and box to peer PME processor
Communicate **x** with real space neighbor processors
Neighborsearch step? — N / Y
(local) neighborsearching
Evaluate potential/forces
Communicate **f** with real space neighbor processors
Spread real space forces on virtual sites
Receive forces/energy/virial from peer PME processor
Spread PME forces on virtual sites
Integrate coordinates
Constrain bond lengths (parallel LINCS)
Sum energies of all real space processors
More steps? — Y / N

Neighborsearch step? — Y / N
Received charges from peer real space processors
Receive **x** and box from peer real space processors
All local coordinates received? — N / Y
Communicate some atoms to neighbor PME proc's
Spread charges on grid
Communicate grid overlap with PME neighbor proc's
parallel 3D FFT
Solve PME (convolution)
parallel inverse 3D FFT
Communicate grid overlap with PME neighbor proc's
Interpolate forces from grid
Communicate some forces to neighbor PME proc's
Send forces/energy/virial to peer real space processors
More steps? — N / Y

Stop

# What does a modern CPU core look like?



Intel Skylake-S/EP

SIMD: 128, 256, or 512 bit vectors

5 μops/core each cycle
SIMD + FMA: 64 flops @ single prec.

Up to 32 cores per chip

2 sockets per node

Theoretically:
320 instructions or 4096 flops
per *cycle* on each node.

latency is ~4 cycles on Skylake.
You need *256* independent FMA
flops (128 FMA operations) to
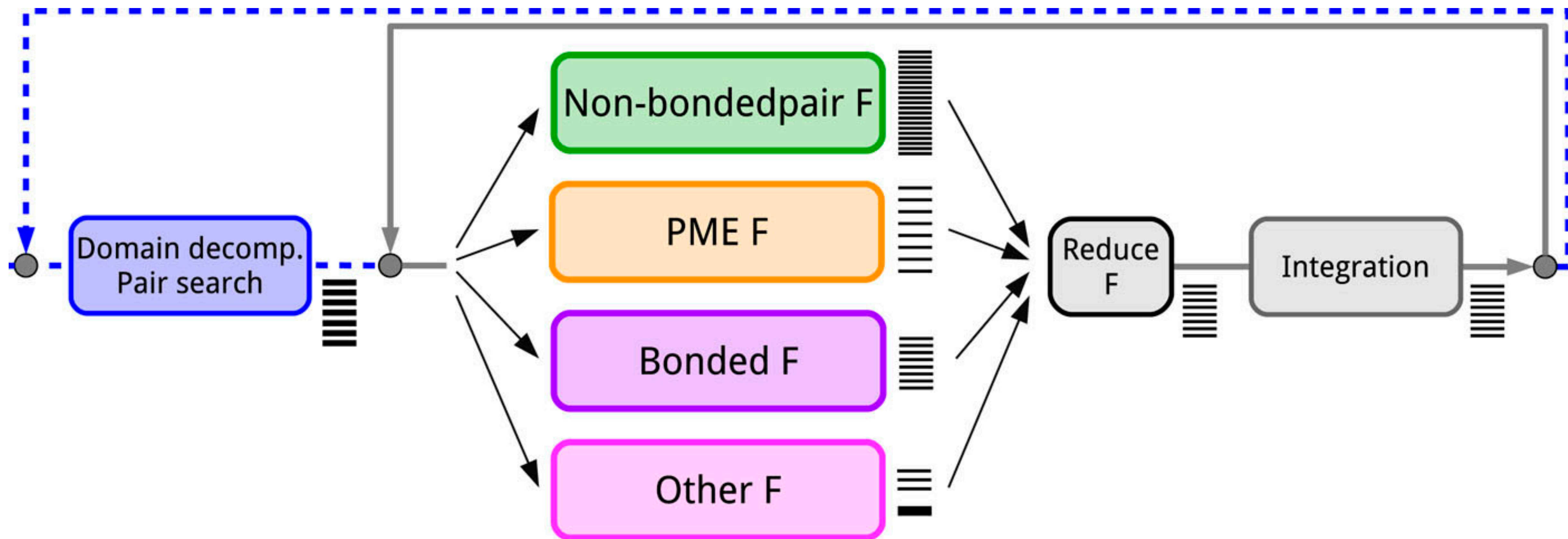saturate just a single core

# Acceleration Approaches

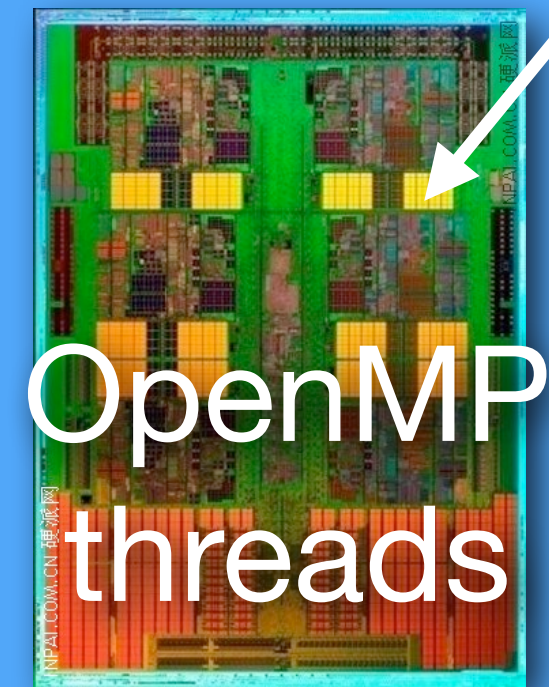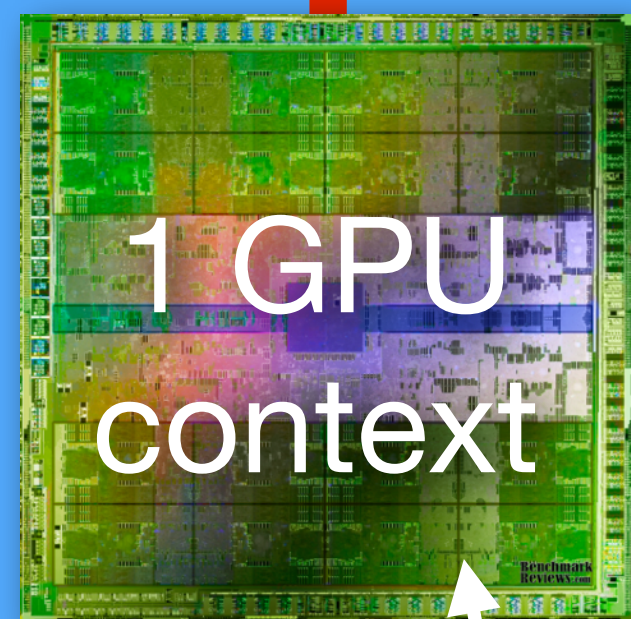| | GPU libraries | OpenAcc | Pure CUDA | Heterogeneous CPU/GPU |
|---|---|---|---|---|
| Initial effort / Expertise req. | 🟩 | 🟨 | 🟧 | 🟥 |
| Generality / Portability | 🟩 | 🟩 | 🟥 | 🟩 |
| Performance | 🟩→🟥 | 🟨 | 🟨 | 🟩 |
| Code maintainability | 🟩 | 🟩 | 🟥 | 🟨 |
| | Works if your code offloads to libraries | Always works, but success depends on you & compiler | Lots of work, assumes impl. can run entirely on GPU | Even more work, less CUDA, can use both CPU & GPU |

Explicit SIMD instructions on CPUs & Xeon Phi; each instruction does up to 32 flops
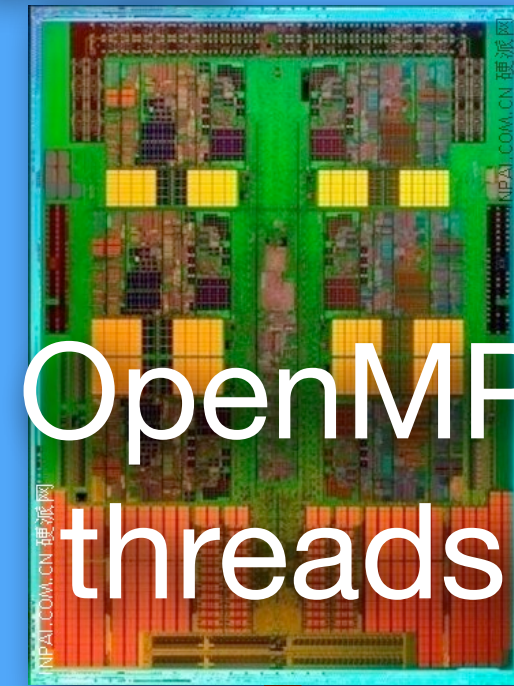
Node

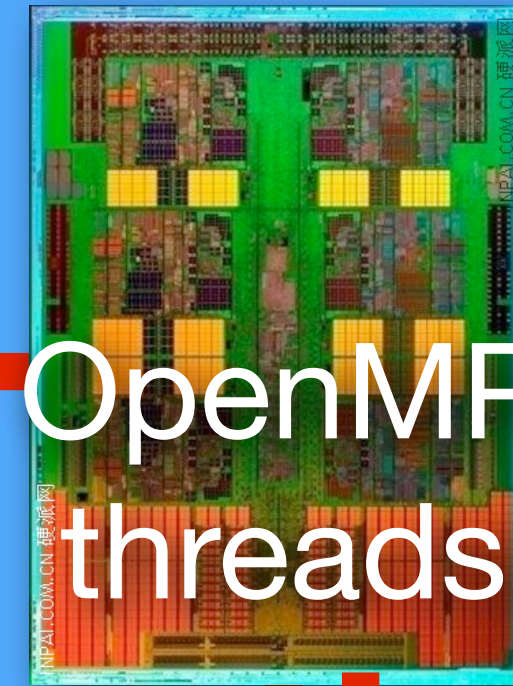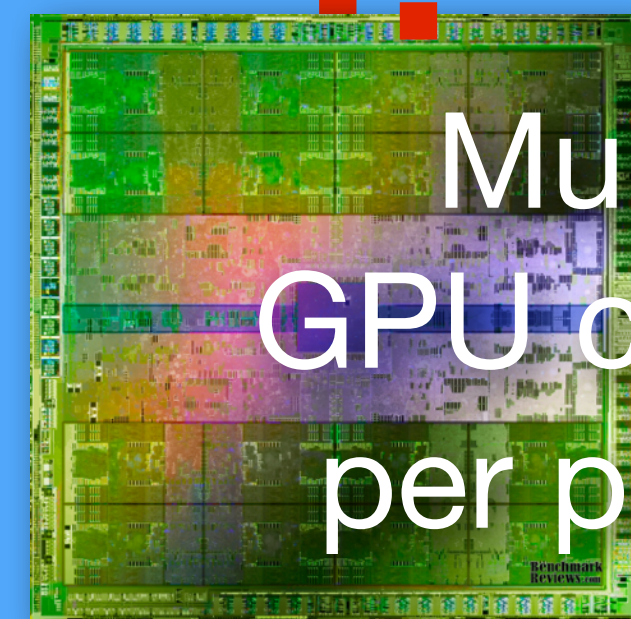OpenMP threads

Load balancing

1 GPU context

Node

OpenMP threads — OpenMP threads

Load balancing — Load balancing

Multiple GPU contexts per process

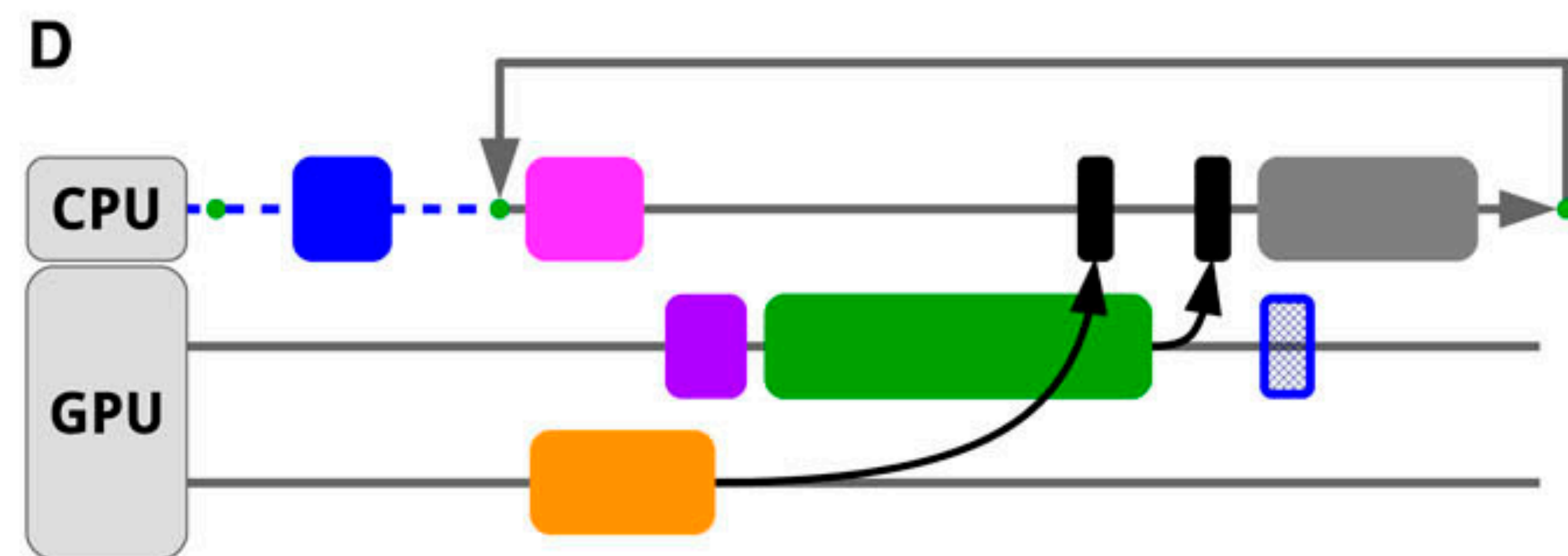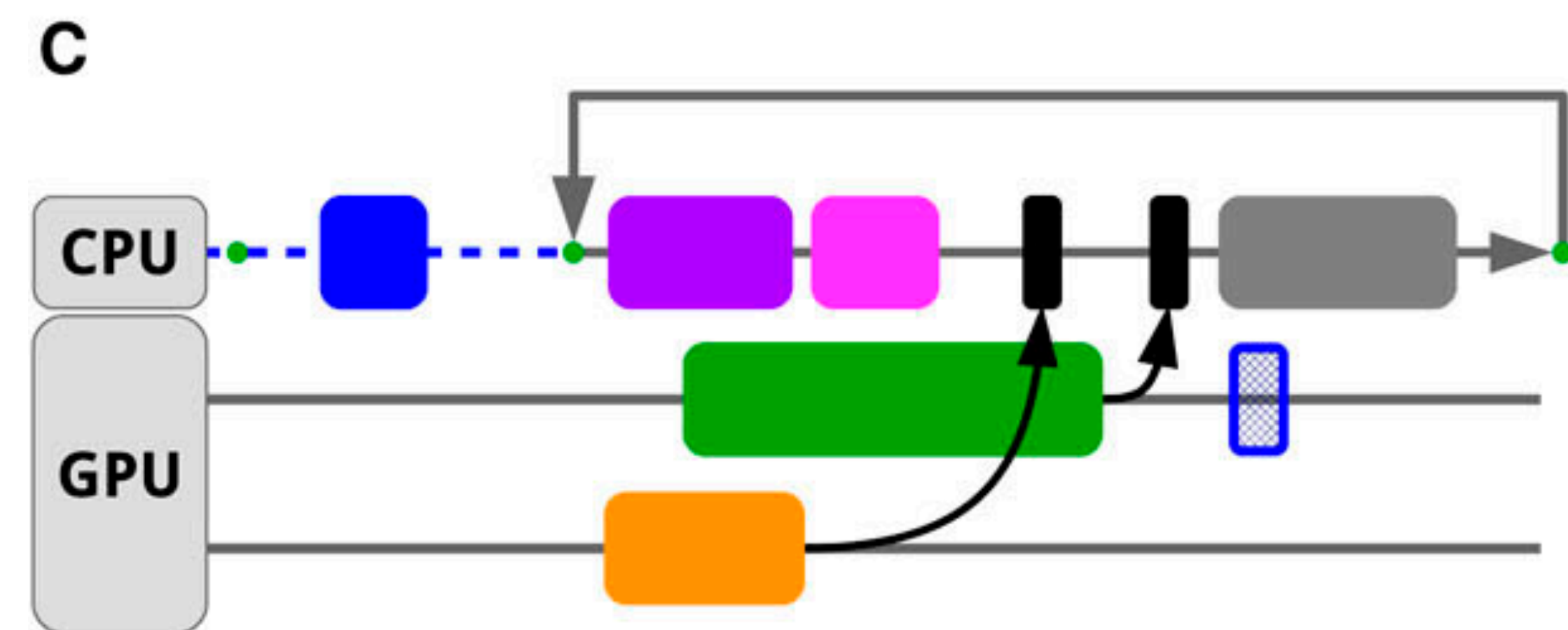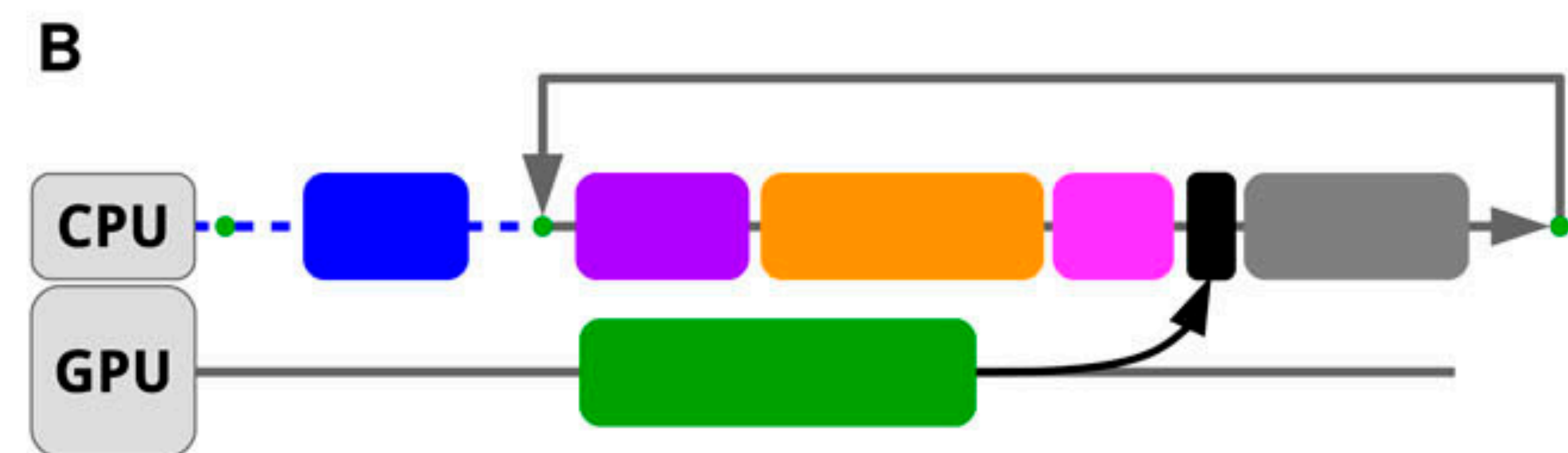CUDA kernels on NVIDIA GPUs, OpenCL for AMD/Intel GPUs

MPI MPI MPI

MPI MPI MPI

MPI MPI MPI

**A**

CPU — Pair Search — Bonded F — Non-bonded F — PME — Other Forces — Reduce Forces — Integration, Constraints

**B**

CPU / GPU

**C**

CPU / GPU

**D**

CPU / GPU

**E**

CPU / GPU

# From neighborlists to cluster proximity lists: Revisit algorithms

**x,y grid**
**z sort**
**z bin**

**x,y,z gridding**

**Cluster pairlist**

$r_{list}$

**Organize as tiles with all-vs-all interactions:**

| X | X | X | X |
|---|---|---|---|
| X | X | X | X |
| X | X | X | X |
| X | X | X | X |

*i=3:*  | 5 | 6 | 9 | 12 | 15 | 17 | 18 | 25 | 32 | … |

*i=4:*  | 7 | 8 | 9 | 11 | 12 | 15 | 17 | 25 | 32 | 43 | 54 | … |

*…*  | 8 | 9 | 10 | 11 | 12 | 13 | 19 | 20 | … |

*Tile interaction algorithms:*
*Load N atoms, compute N^2 forces*

*The Link-cell algorithm: Load 1 atom, calculate 1 interaction*
**Verlet, Phys Rev 159, 98-103 (1967)**

# This creates a new problem: Tiling circles is difficult

You want to calculate interactions with red neighbors



serial computing

stream computing

Lots of wasted FLOPS!

- You need a lot of cubes to cover a sphere

- All interactions beyond cutoff need to be zero

# Clusters on CPUs, Superclusters on GPUs

# Bringing the Performace back to the CPU

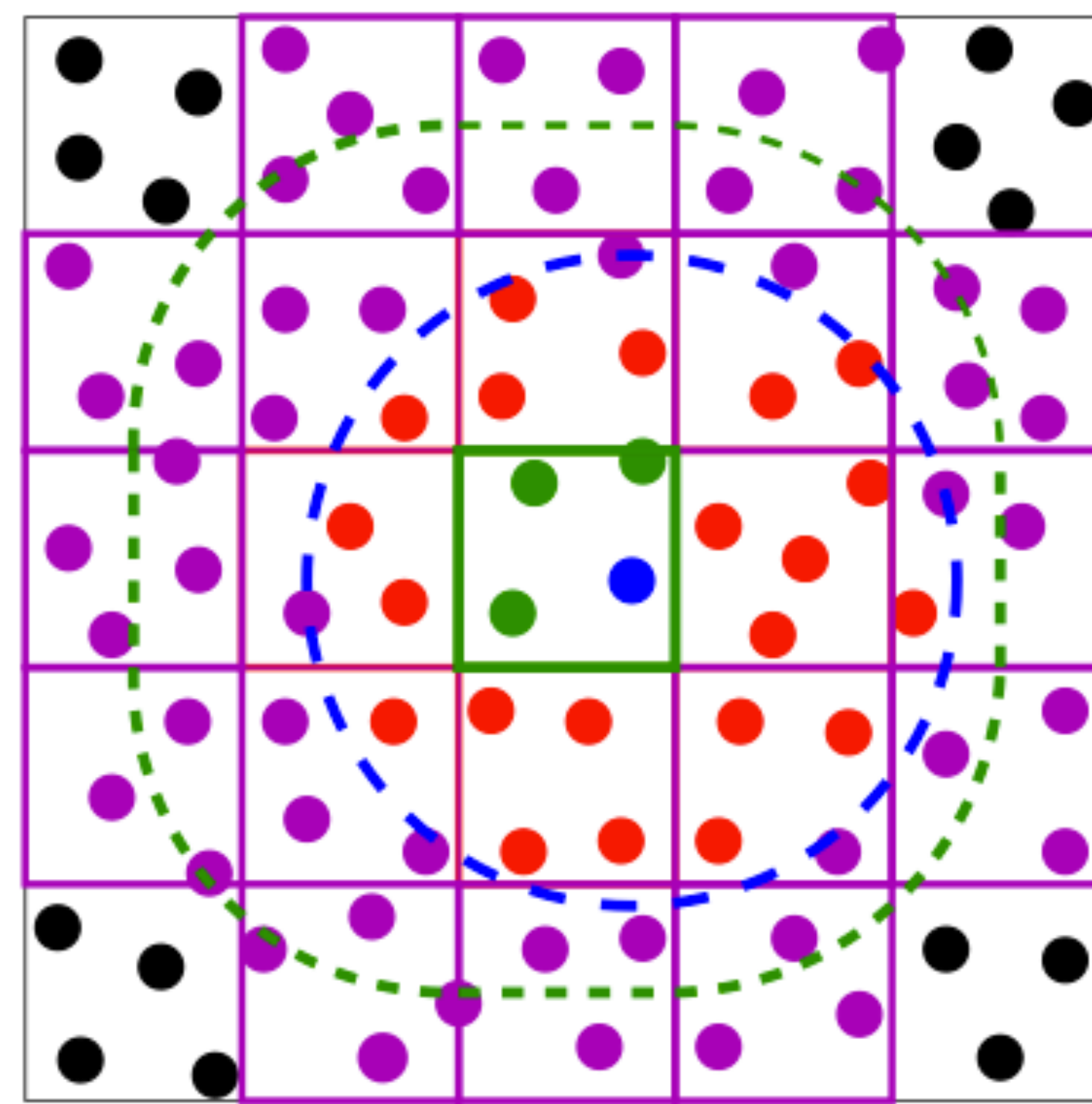# Unified GPU/CPU architecture - completely portable

CUDA
OpenCL
Intel MIC
x86 SSE2
x86 SSE4.1
x86 AVX
x86 AVX-128-FMA
x86 AVX2
x86 AVX2_128
x86 AVX-512F
x86 AVX-512ER
Arm Neon
Arm64 Asimd
IBM QPX
IBM VMX
IBM VSX
Fujitsu HPC-ACE

# Atom clustering and pair list buffering



We can
Adjust the
size of this buffer

Larger buffer
means more
calculations, but
we can update
the neighbor list
less frequently

New dual-pair list buffer:
- Use very large buffers, and prune it every few steps
- reduces overhead
- less sensitive to parameters

# The big gain of heterogeneous acceleration:
# Very little CUDA required

```
n_cuda lindahl$ ls -ltar

lindahl    staff    13012 Apr 18 15:10 nbnxn_cuda_types.h
lindahl    staff     9155 Apr 18 15:10 nbnxn_cuda_kernels.cuh
lindahl    staff    21576 Apr 18 15:10 nbnxn_cuda_kernel_utils.cuh
lindahl    staff    20945 Apr 18 15:10 nbnxn_cuda_kernel.cuh
lindahl    staff     1965 Apr 18 15:10 CMakeLists.txt
lindahl    staff    39049 Apr 18 15:10 nbnxn_cuda_data_mgmt.cu
lindahl    staff     3667 Apr 18 15:10 nbnxn_cuda.h
  ahl    staff    30920 May 22 09:13 nbnxn_cuda.cu
  ahl    staff     2686 May 22 09:13 ..
  ahl    staff      340 May 22 09:13 .
```

**A total of ~3000 lines of CUDA, compared to ~2M lines of C++**

```c
cuda_kernel.c

  cuda_kernel.c > No Selection

351                        {
352                            /* load the rest of the i-atom parameters */
353                            qi        = xqbuf.w;
354  #ifdef IATYPE_SHMEM
355                            typei     = atib[i * CL_SIZE + tidxi];
356  #else
357                            typei     = atom_types[ai];
358  #endif
359
360                            /* LJ 6*C6 and 12*C12 */
361  #ifdef USE_TEXOBJ
362                            c6        = tex1Dfetch<float>(nbparam.nbfp_texobj, 2 * (ntypes * typei + typej));
363                            c12       = tex1Dfetch<float>(nbparam.nbfp_texobj, 2 * (ntypes * typei + typej) + 1);
364  #else
365                            c6        = tex1Dfetch(nbfp_texref, 2 * (ntypes * typei + typej));
366                            c12       = tex1Dfetch(nbfp_texref, 2 * (ntypes * typei + typej) + 1);
367  #endif /* USE_TEXOBJ */
368
369
370                            /* avoid NaN for excluded pairs at r=0 */
371                            r2        += (1.0f - int_bit) * NBNXN_AVOID_SING_R2_INC;
372
373                            inv_r     = rsqrt(r2);
374                            inv_r2    = inv_r * inv_r;
375                            inv_r6    = inv_r2 * inv_r2 * inv_r2;
376  #if defined EXCLUSION_FORCES
377                            /* We could mask inv_r2, but with Ewald
378                             * masking both inv_r6 and F_invr is faster */
379                            inv_r6   *= int_bit;
380  #endif  /* EXCLUSION_FORCES */
381
382                            F_invr  = inv_r6 * (c12 * inv_r6 - c6) * inv_r2;
383  #if defined CALC_ENERGIES || defined LJ_POT_SWITCH
384                            E_lj_p  = int_bit * (c12 * (inv_r6 * inv_r6 + nbparam.repulsion_shift.cpot)*ONE_TWELVETH_F -
385                                                   c6 * (inv_r6 + nbparam.dispersion_shift.cpot)*ONE_SIXTH_F);
386  #endif
387
388  #ifdef LJ_FORCE_SWITCH
389  #ifdef CALC_ENERGIES
390                            calculate_force_switch_F_E(nbparam, c6, c12, inv_r, r2, &F_invr, &E_lj_p);
391  #else
392                            calculate_force_switch_F(nbparam, c6, c12, inv_r, r2, &F_invr);
393  #endif /* CALC_ENERGIES */
394  #endif /* LJ_FORCE_SWITCH */
```
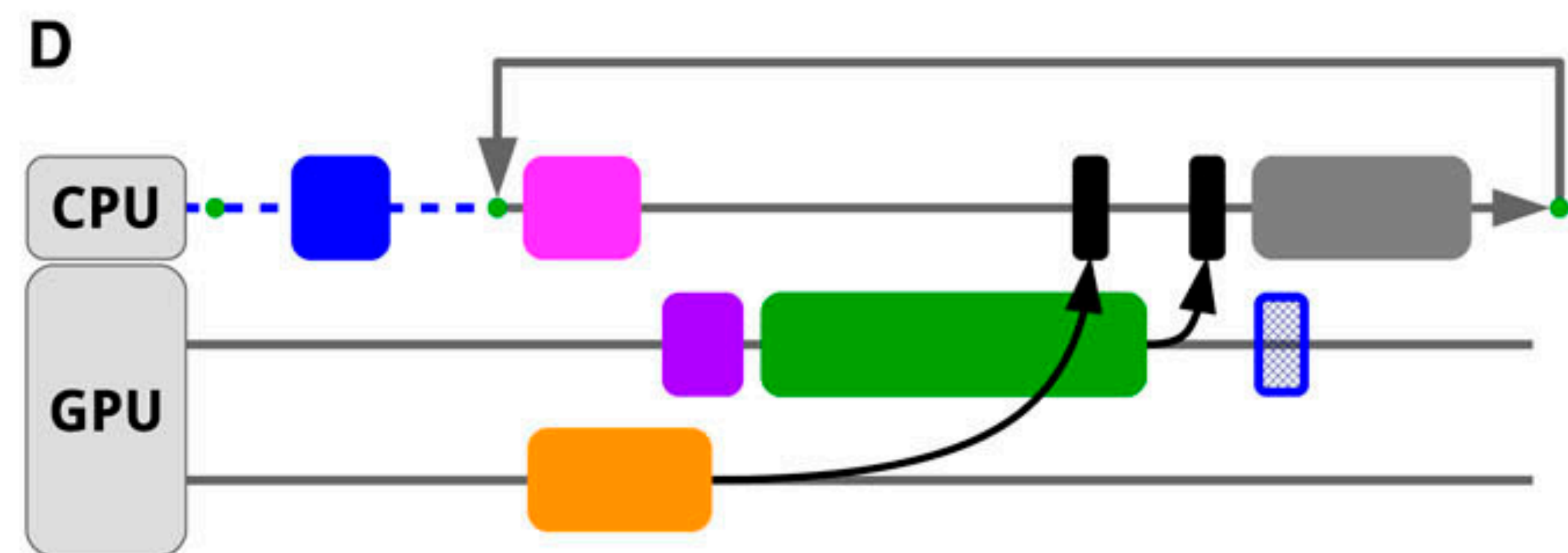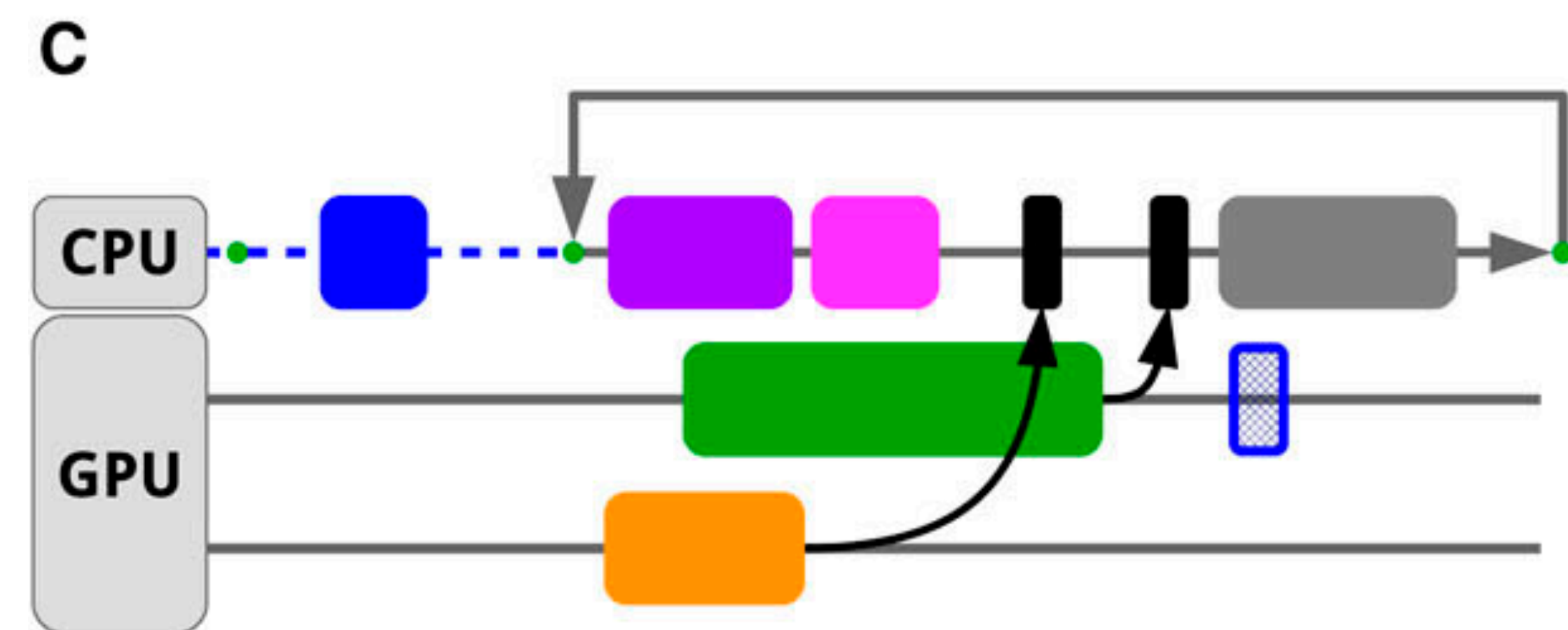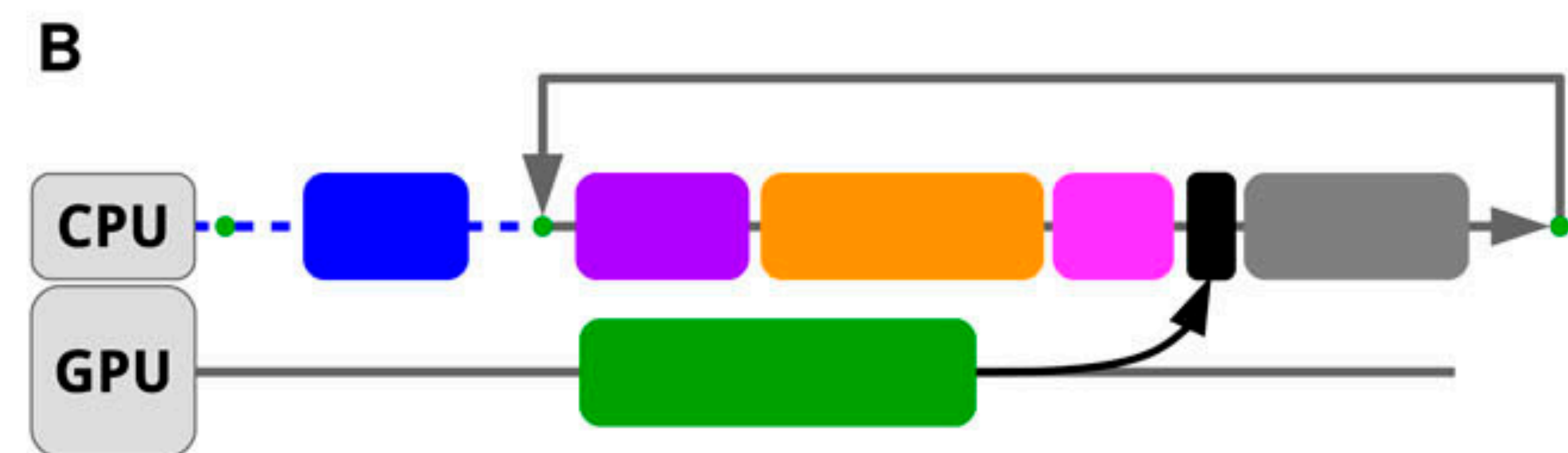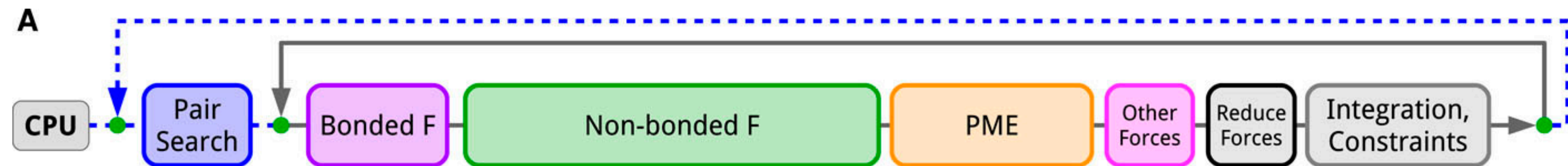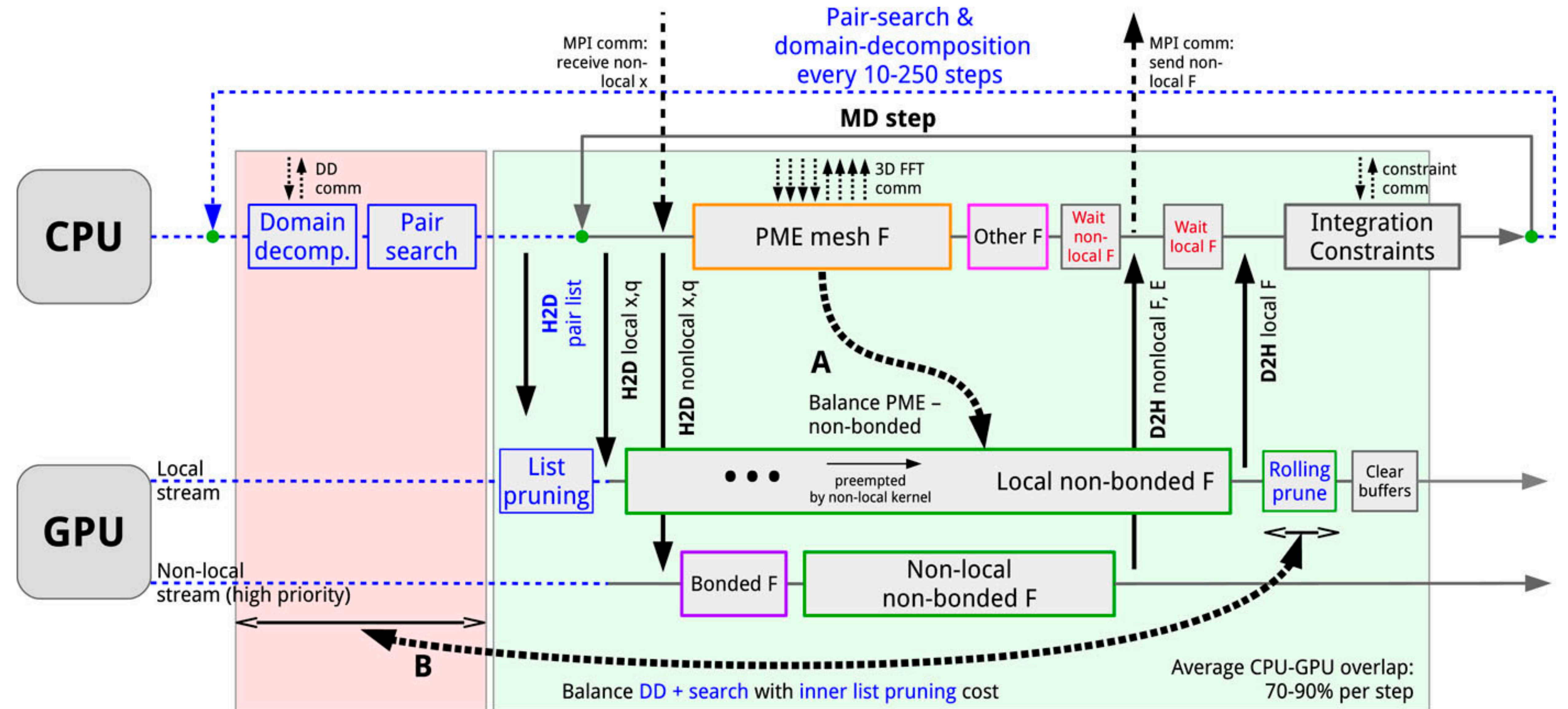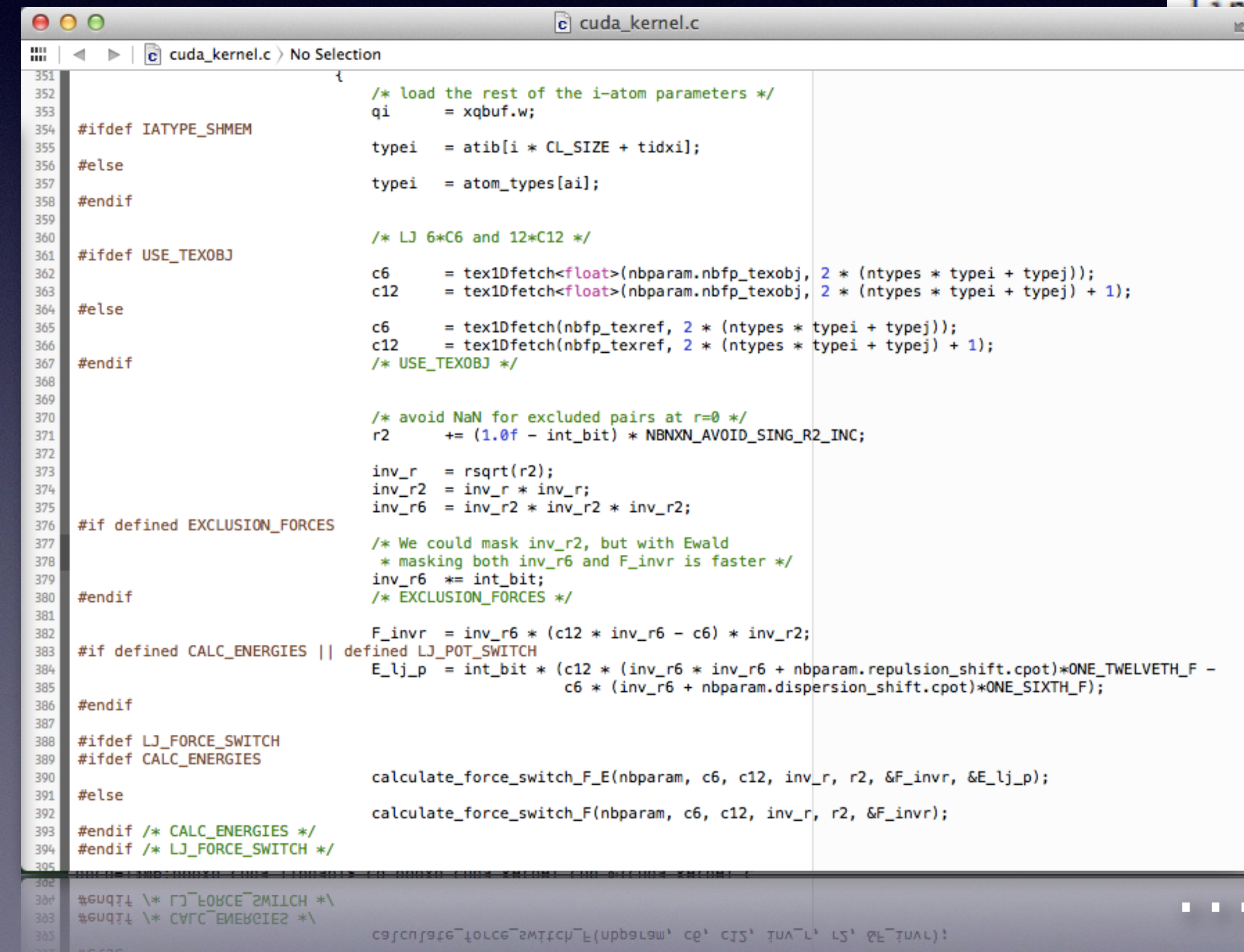
**… so we wrote OpenCL kernels too!**

# Make good use of both GPU & CPU

- Use the CPU to pre-calculate or optimize data structures, so there is less work for the GPU to do in your kernels
- Easier to implement more complex optimization on CPU
- Advanced multi-node domain decomposition easier on CPU
- Run some parts of the algorithm on the CPU (avoid wasting flops)

1. It's important to keep the GPU busy
2. … but it doesn't have to be busy 100% of the time!
3. A CUDA GPU running at 100% will get hot, and clock down
4. NVML "application clocks" effectively overclock the GPU on-the-fly when you have less than 100% utilization

*Think of a node as a collection of compute & communication devices - use them all!*

# Kernel latency is key for (heterogeneous) acceleration

**A**

CPU | Pair Search | Bonded F | Non-bonded F | PME | Other Forces | Reduce Forces | Integration, Constraints

**B**

CPU | GPU

**C**

CPU | GPU

**D**

CPU | GPU

**E**

CPU | GPU

# Exploiting multiple & high-priority streams

**GPU**
**CUDA**

Local stream

Non-local stream (high priority)

Stream that only needs local data can start directly, but can be preempted by the high-priority nonlocal data kernel

When remote data is delivered, handle it immediately so it can be returned faster

# Exploiting multiple & high-priority streams



Stream that only needs local data can start directly, but can be preempted by the high priority nonlocal data kernel
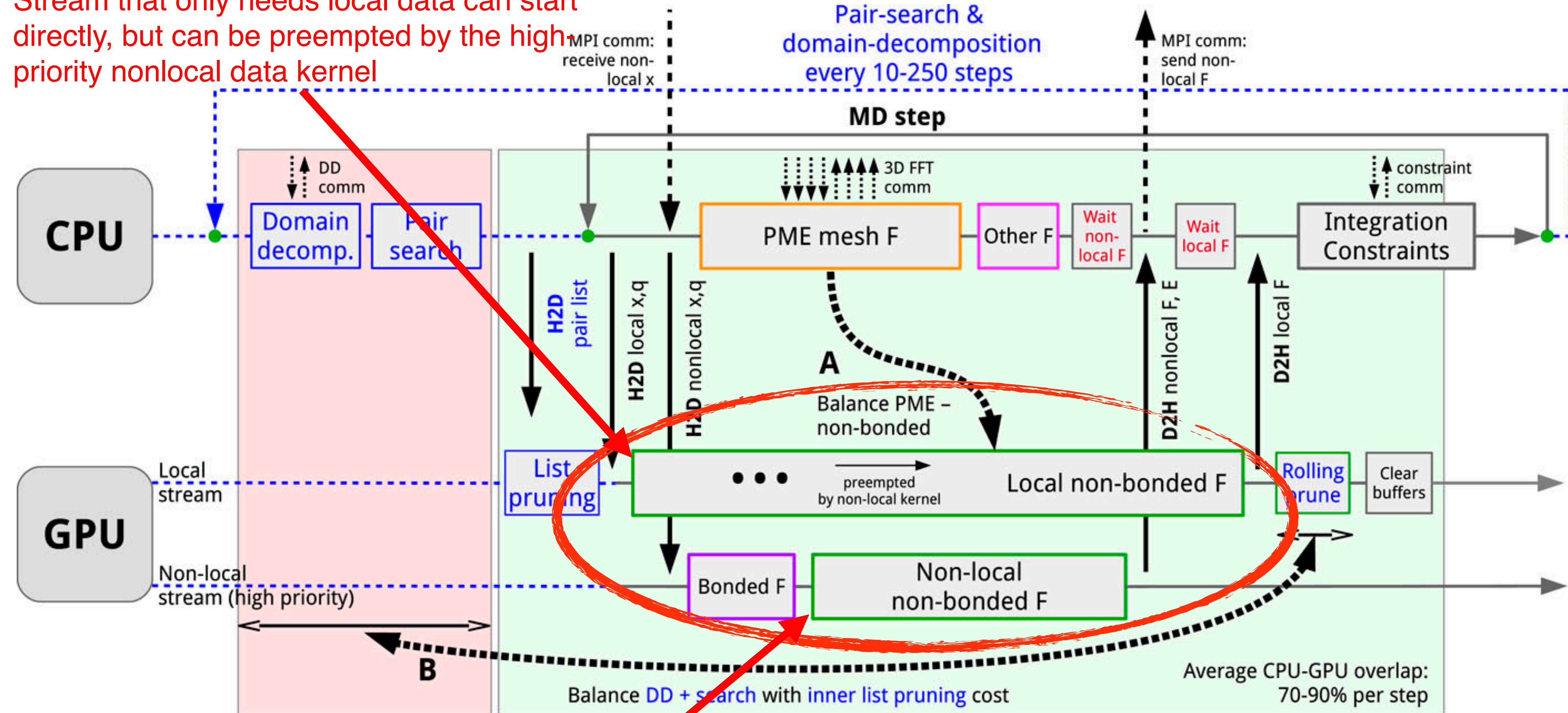
When remote data is delivered, handle it immediately so it can be returned faster

MPI comm: receive non-local x

Pair-search & domain-decomposition every 10-250 steps

MPI comm: send non-local F

MD step

**CPU**

DD comm

Domain decomp.

Pair search

3D FFT comm

PME mesh F

Other F

Wait non-local F

Wait local F

constraint comm

Integration Constraints

**GPU**

Local stream

Non-local stream (high priority)

H2D pair list

H2D local x,q

H2D nonlocal x,q

A

Balance PME – non-bonded

List pruning

preempted by non-local kernel

Local non-bonded F

D2H nonlocal F, E

D2H local F

Rolling prune

Clear buffers

Bonded F

Non-local non-bonded F

B

Balance DD + search with inner list pruning cost

Average CPU-GPU overlap: 70-90% per step

# Revisiting Amdahl's law - give GPU more work

*The least parallel part of the code (or at least slowest piece of hardware)
will eventually dominate execution completely and limit scaling*

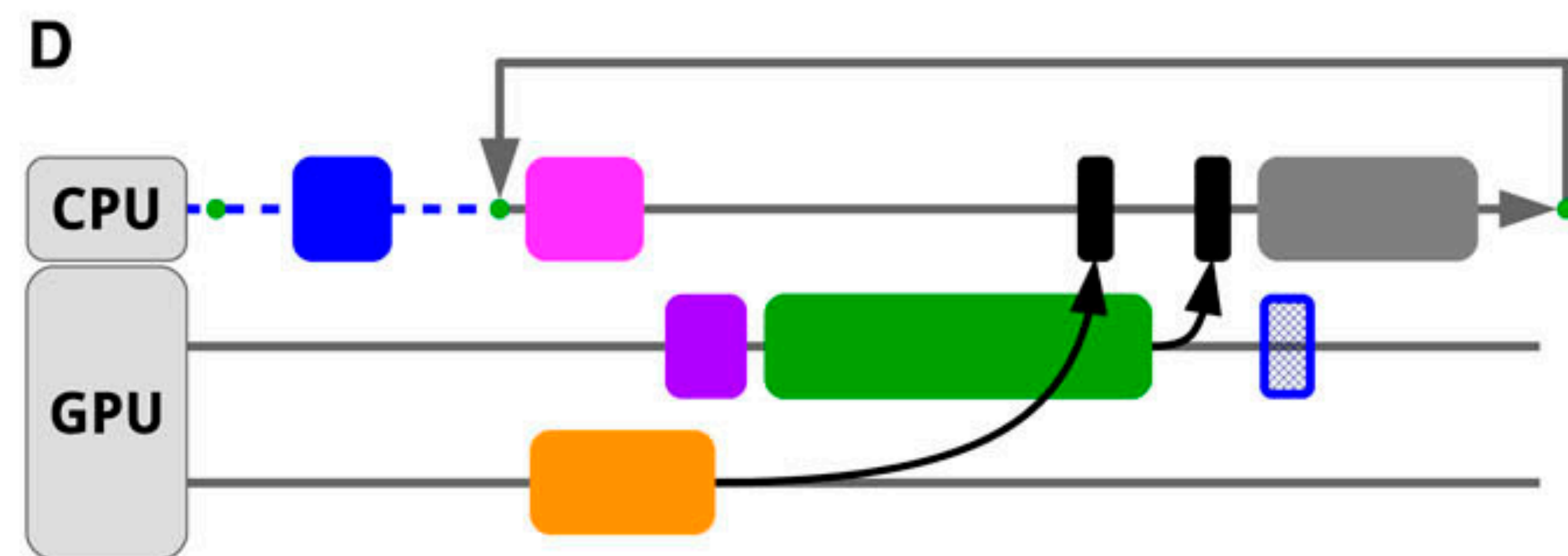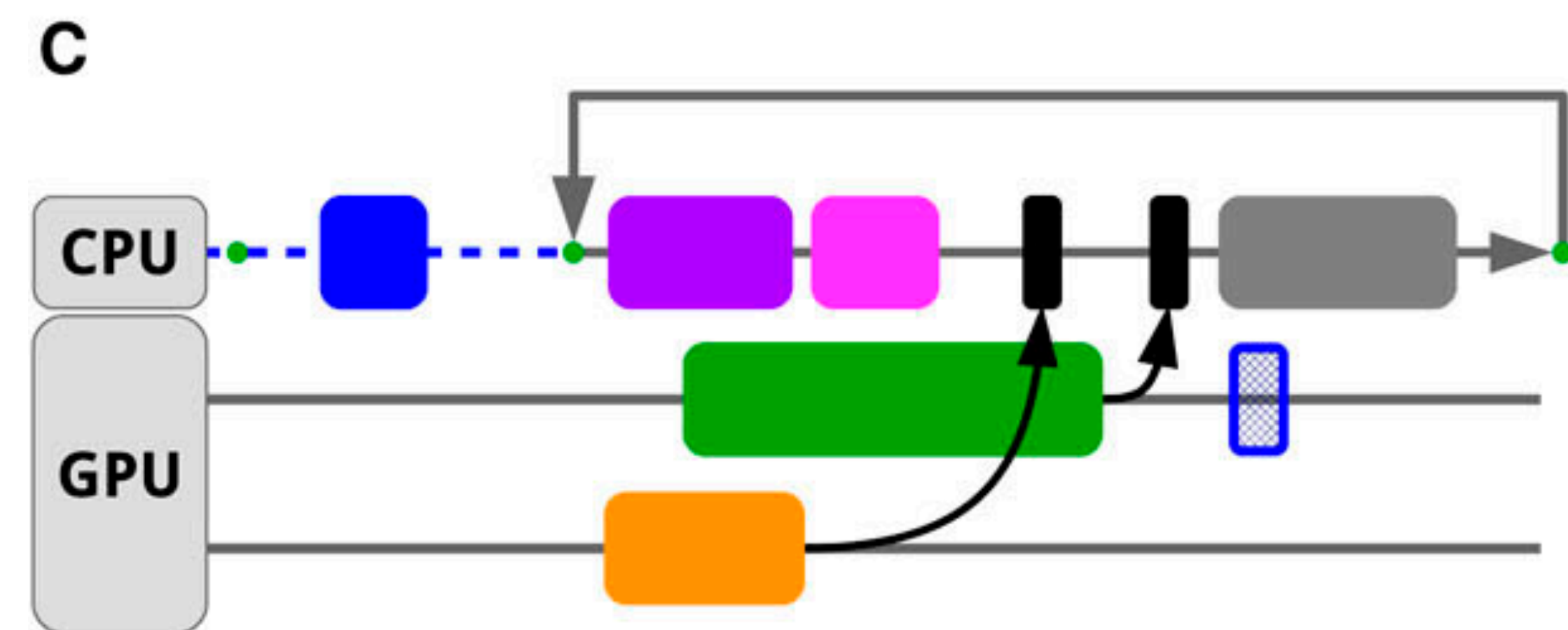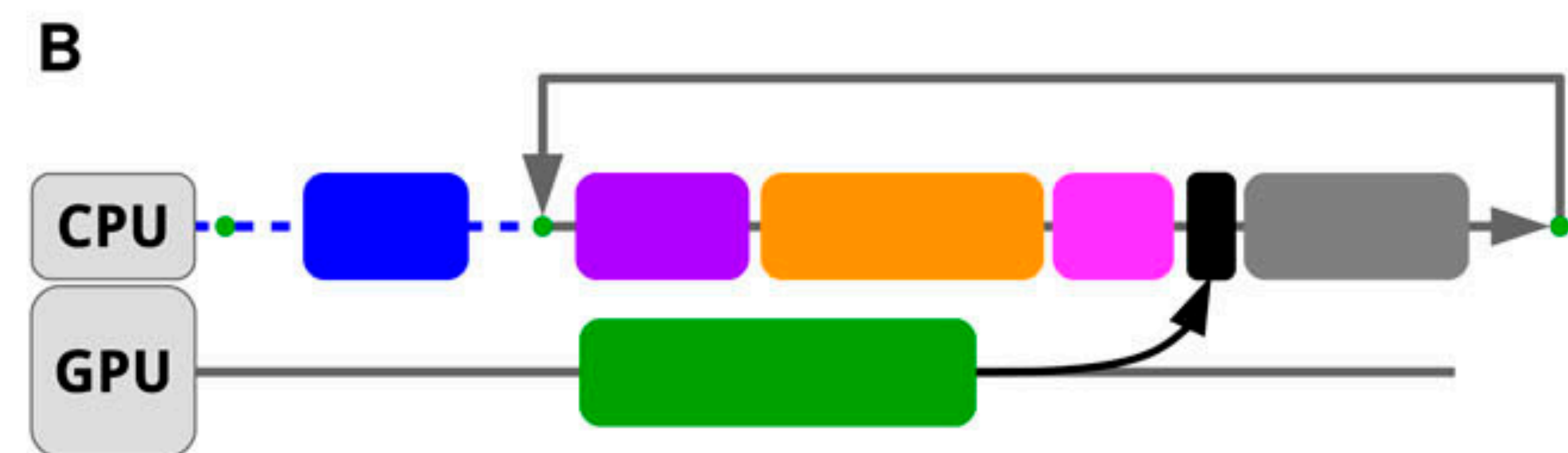Thanks to heterogeneous parallelism and efficient CPU-side algorithms, GROMACS frequently outperforms GPU-only implementations - and yet we only need a few thousand lines of CUDA.

But... GPU performance grows faster than CPU performance, and sometimes we want to put a high-end GPU in an old low-end CPU box

Our CPUs used to wait for the GPUs, now it's often the opposite

**A**

CPU — Pair Search — Bonded F — Non-bonded F — PME — Other Forces — Reduce Forces — Integration, Constraints

**B**

CPU / GPU

**C**

CPU / GPU

**D**

CPU / GPU

**E**

CPU / GPU

# The new bottleneck (for slow CPUs) is the PME algorithm
## 3D grid spreading, FFTs, convolution, iFFT, interpolation

PME offload into separate stream

# Load-balance the *algorithm* (not just work) between CPU & GPU

# GPU Timestep profile before/after PME offload (P100)

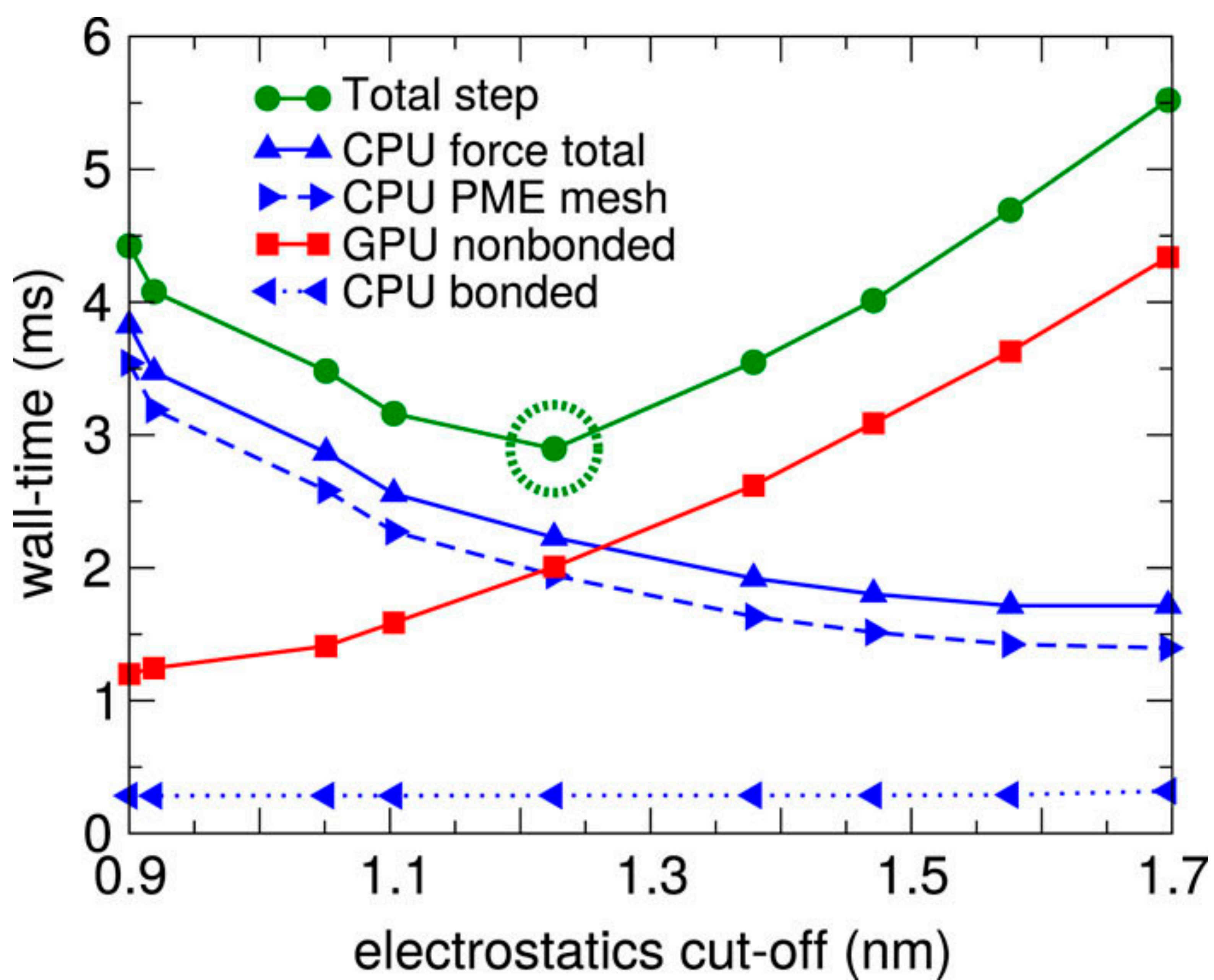Highly challenging small system (25k atoms), very fast iterations.
Much harder, but important for multi-GPU scaling



*450µs for a complete step - note x scale!*

Harder to retain full compute throughput than communication BW when scaling to small sys. Scheduling limited below - kernels should overlap. Lauch operations limiting us.

# Even codes that have been tuned for ~20 years on CPUs and ~10 years on GPUs can get great performance gains just from better algorithms & implementations

# We solved our issue of fast-CPU-dependency: Fast with a single core per GPU, even faster with many



Legend:
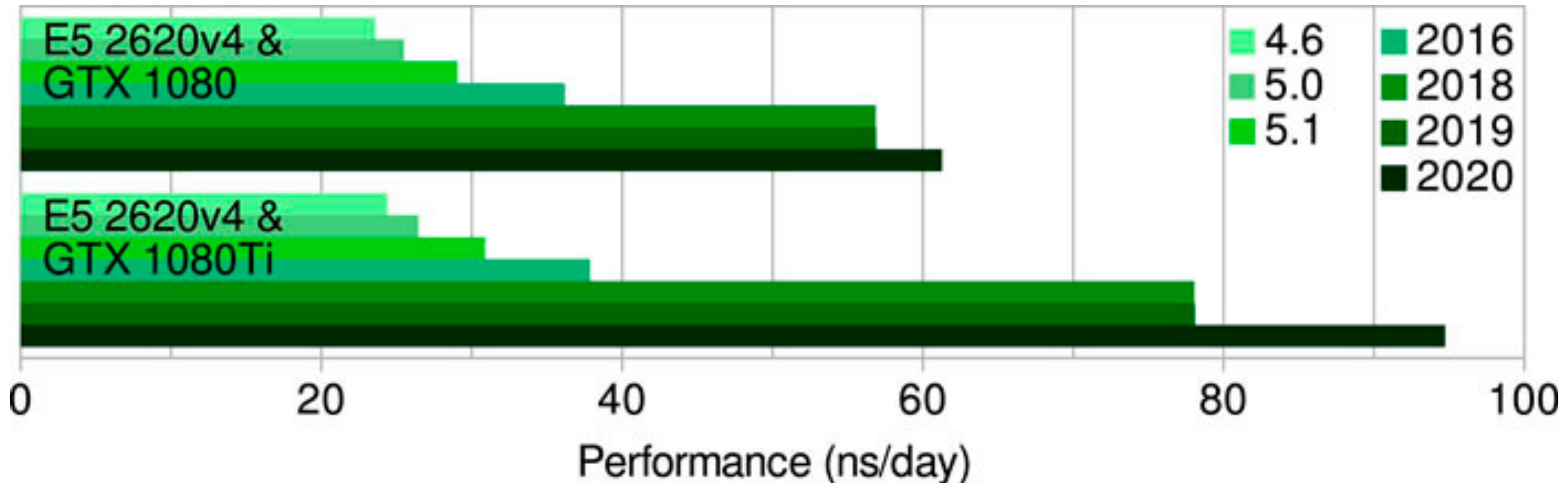- CPU
- Nonbonded
- Nonbonded and PME
- Nonbonded, bonded and PME
- Nonbonded, bonded, PME and update
- Nonbonded, PME and update

Panel A: Performance (ns/day) vs Number of CPU cores
Panel B: Performance (ns/day) vs Number of CPU cores

# Heterogeneous parallelisation provides good scaling even for the highly latency-sensitive algorithms in molecular dynamics

# Strong caling issues - challenges at 100µs per iteration

- The 3D-FFT in PME
- MPI overhead - we need MPI_Put_notify()
- OpenMP barriers take significant time

- Load imbalance
- CUDA API overhead can be 50% of CPU time
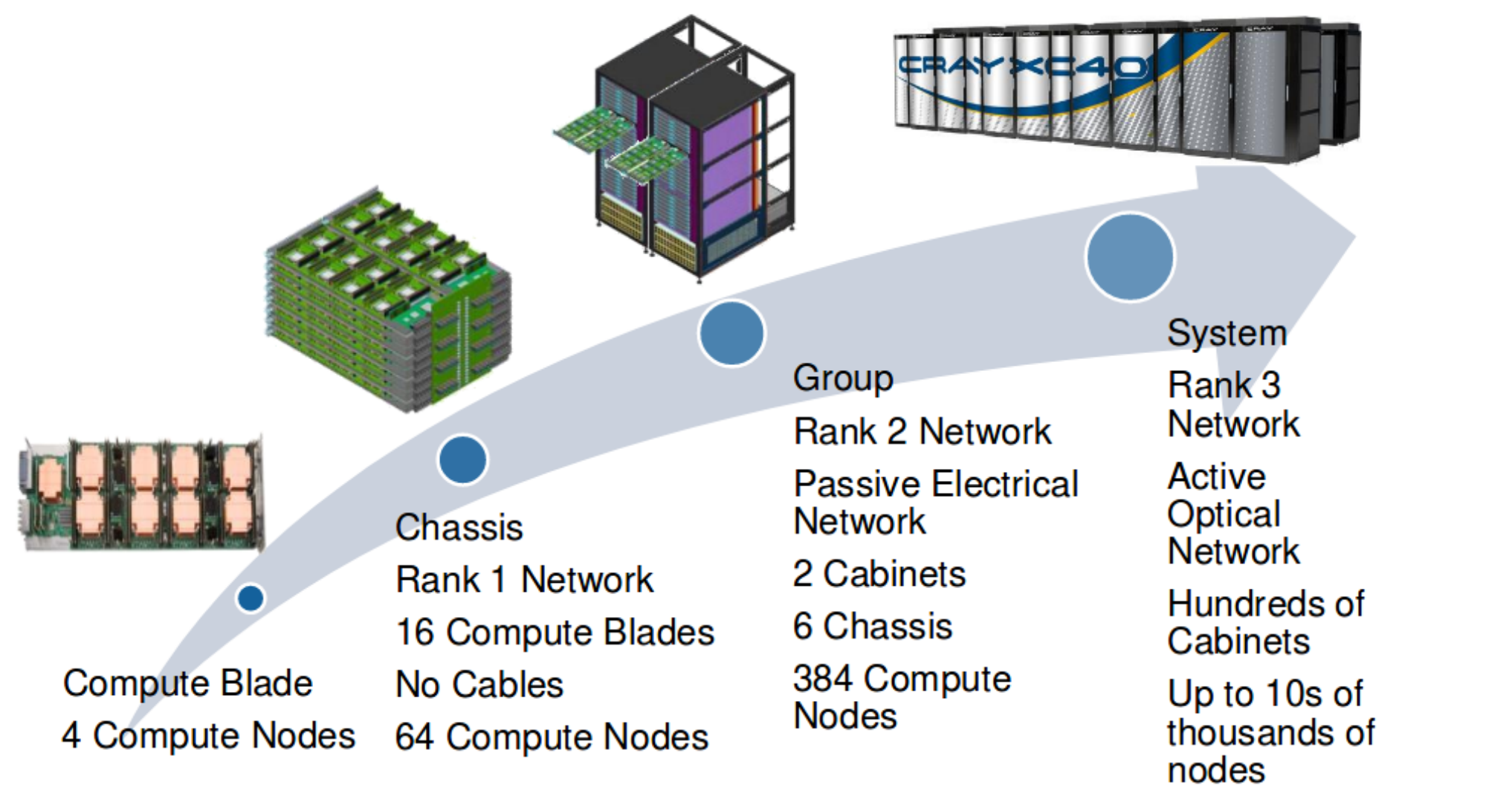- Too many GROMACS options to tweak manually

Large performance loss due to imbalance and network speed variation on Cray XC (interference from other jobs on the "smart" network)

<1 millisecond



**It would help a lot to have more control over node placement**

Cray XC System Building Blocks

Compute Blade
4 Compute Nodes

Chassis
Rank 1 Network
16 Compute Blades
No Cables
64 Compute Nodes

Group
Rank 2 Network
Passive Electrical Network
2 Cabinets
6 Chassis
384 Compute Nodes

System
Rank 3 Network
Active Optical Network
Hundreds of Cabinets
Up to 10s of thousands of nodes

COMPUTE | STORE | ANALYZE

Unconstrained
Constrained to groups

# Intra-rank parallelisation: OpenMP today, future ?

Efficient current parallelization of all algorithms using MPI + OpenMP

OpenMP is (performance) portable, but limited:

- No way to run parallel tasks next to each other
- No binding of threads to cores (cache locality)

Need for a better threading model, requirements:

- Extremely low overhead barriers (all-all, all-1, 1-all)
- Binding of threads to cores
- Portable

We are convinced we are moving to a world where latency- and throughput-optimized units converge into the same chip - the future is heterogeneous!

Urgent need for better, standardized and portable HPC-focused task parallelism frameworks. We are looking into both ArgoBots and home-grown solutions.

Spend time with your algorithms, not just code tuning.

A single Skylake-EP node has 4096-fold parallelism.
Your code likely doesn't.

Think accelerators - because a modern CPU looks like an accelerator,
and they will likely converge to multiple units on one die in the future.

Heterogeneous parallelism uses all resources and provides architecture portability.

Fast-iteration codes are very sensitive to node placement,
and they need task parallelism sooner rather than later.

Fast-iteration coding for CUDA/AVX512/OpenCL/SYCL isn't hard - but new algorithms are.

You can accomplish miracles with more codes than you think,
but it takes 6-12 months - not an afternoon.

Theory & Computation is the new experiment!

# Acknowledgments