LIFE SCIENCE EXAMPLES: HOW DO WE CONCRETELY PARALLELIZE **AND ACCELERATE SCIENTIFIC PROBLEMS?** (AND HOW DO YOU MAKE IT INTO A CAREER?)

Thomas E Cheatham III, University of Utah, Salt Lake City, UT Erik Lindahl, Stockholm University, Stockholm tec3@utah.edu erik.lindahl@scilifelab.se





Two Stories:



Cryo-EM!



With $\Delta t \sim 1 fs$ and μs to s timescales of interest, we need 10⁹-10¹⁵ 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

```
for(k=nj0; (k<nj1); k++)</pre>
    /* Get j neighbor index. and coordinate index */
                   = jjnr[k]; >
    jnr
    j3
                      >*INC;
    /* load j atom coordinates */
                     = pos[j3+0];
    jx1
                     = pos[j3+1];
   jy1
                     = pos[j3+2];
    jz1
    /* Calculate distance */
                     = ix1 - jx1;
    dx11
   dy11
                     = iy1 - jy1;
                     = iz1 - jz1;
    dz11
                     = dx11*dx11+dy11*dy11+dz11*dz11;
    rsq11
                     and 1/r2 */
    /* Calculate 1
                     = 1.0/sqrt(rsq11);
    rinv11
    /* Load paramete.
                       for i atom \star
                     = ig*charge[jnr];
    qq
                     = nti+2*type[jnr];
    tj
                     = vdwparam[tj];
    c6
    c12
                     = vdwparam[tj+1];
                     = rinv11*rinv11;
    rinvsq
    /* Coulomb interaction */
                     = qq*rinv11;
    vcoul
    vctot
                     = vctot+vcoul;
    /* Lennard-Jones interaction */
    rinvsix
                     = rinvsq*rinvsq*rinvsq;
    Vvdw6
                     = c6*rinvsix;
    Vvdw12
                     = c12*rinvsix*rinvsix;
    Vvdwtot
                     = Vvdwtot+Vvdw12-Vvdw6;
                     = (vcoul+12.0*Vvdw12-6.0*Vvdw6)*rinvsq;
    fscal
   /* Calculate temporary vectorial force */
                     = fscal*dx11;
    tx
                     = fscal*dy11;
    ty
    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 μ s 6666888888866688866688

To-Do Thursday 10:15:48.004.100 (simple version)

Adjust domain decomposition [communicate] Communicate coordinates to/from 26 neighbor nodes Find atoms in proximity **Ccommunicate**] 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

mmunica 00 <u></u> ന



What does a modern CPU core look like?



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



What does a modern CPU die look like? AMD Ryzen 1600x



The inside of a modern node looks like a cluster with advanced network topology!



Explicit Data Parallelism

Stream=your data Kernel=algoritm

Without dependencies, all these operations could be done in parallel if enough hardware was available!





NVLink NVLink	NVLink



Acceleration Approaches

	GPU libraries	OpenAcc	Pure CUDA	Heterogeneou 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 wo less CUDA, can use both CPU & GPU



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



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



Strategy: 1. Profile 2. Offload largest bottleneck to GF



From neighborlists to cluster proximity lists: Revisit everything





9 10 11 12 13 19 20 ...

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





Organize as tiles with all-vs-all interactions:





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







This creates a new problem: Tiling circles is difficult serial computing stream computing

• You need a lot of cubes to cover a sphere

All interactions beyond cutoff need to be zero

You want to calculate interactions with red neighbors



Lots of wasted FLOPS!



The big gain of heterogeneous acceleration: n_cuda lindahl\$ ls –ltar Very little CUDA required

•	0 0		c cuda_kernel.c
	< ▶	cuda_kernel.c > No Selection	on
351 352 353	4:64-6		<pre>/* load the rest of the i-atom parameters */ qi = xqbuf.w;</pre>
355	#11001	TATTPE_SHINEH	<pre>typei = atib[i * CL_SIZE + tidxi];</pre>
350 357 358 250	#etse #endif		<pre>typei = atom_types[ai];</pre>
360 361	#ifdef	USE TEXOBJ	/* LJ 6*C6 and 12*C12 */
362 363	#=1		<pre>c6 = tex1Dfetch<float>(nbparam.nbfp_texobj, 2 * (ntypes * typei + typej)); c12 = tex1Dfetch<float>(nbparam.nbfp_texobj, 2 * (ntypes * typei + typej) + 1);</float></float></pre>
304 365 366 367 269	#etse #endif		<pre>c6 = tex1Dfetch(nbfp_texref, 2 * (ntypes * typei + typej)); c12 = tex1Dfetch(nbfp_texref, 2 * (ntypes * typei + typej) + 1); /* USE_TEX0BJ */</pre>
369 370 371 372			<pre>/* avoid NaN for excluded pairs at r=0 */ r2 += (1.0f - int_bit) * NBNXN_AVOID_SING_R2_INC;</pre>
373 374 375 376	#if det	fined EXCLUSION_FORCES	<pre>inv_r = rsqrt(r2); inv_r2 = inv_r * inv_r; inv_r6 = inv_r2 * inv_r2 * inv_r2;</pre>
377 378 379 380	#endif		<pre>/* We could mask inv_r2, but with Ewald * masking both inv_r6 and F_invr is faster */ inv_r6 *= int_bit; /* EXCLUSION FORCES */</pre>
381 382 383 384	#if det	fined CALC_ENERGIES de	<pre>F_invr = inv_r6 * (c12 * inv_r6 - c6) * inv_r2; fined LJ_POT_SWITCH E li p = int bit * (c12 * (inv r6 * inv r6 + nbparam.repulsion shift.cpot)*ONE TWELVETH E</pre>
385 386	#endif		<pre>c6 * (inv_r6 + nbparam.dispersion_shift.cpot)*ONE_SIXTH_F);</pre>
387 388 389 390	#ifdef #ifdef	LJ_FORCE_SWITCH CALC_ENERGIES	calculate_force_switch_F_E(nbparam, c6, c12, inv_r, r2, &F_invr, &E_lj_p);
391 392 393 394	#else #endif #endif	/* CALC_ENERGIES */ /* LJ_FORCE_SWITCH */	calculate_force_switch_F(nbparam, c6, c12, inv_r, r2, &F_invr);
302 302 304 303	#endif #endif	/* CALC_ENERGIES */ /* LJ_FORCE_SWITCH */	



A total of ~3000 lines of CUDA, compared to 3M lines of C++

. 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 you have less than 100% utilization

- 4. NVML "application clocks" effectively overclock the GPU on-the-fly when
- Think of a node as a collection of compute & communication devices use them all!



Kernel latency is key for heterogeneous acceleration

CUDA runtime overheads

- Horrible ECC + cudaStreamSynchronize overhead gone (mostly)!
- But observed:
 - single stream (without DD): 35-60 us
 - two streams (with DD): 50-100 us
 - cudaStreamSynchronize ~30% faster than cudaEventSynchronize

CUDA RT call	single-stream (us/step)	dual-stream (us/step)
H2D pair list	0.3	1.1
H2D x+q	7.2	9.6
NB kernel	9.8	14.8
D2H forces	5.2	9.2
D2H E+shift F	1.4	1.4
cudaStreamSync	3.4	2.3
Clear kernel	9.2	8.8
Total	37.2	49.5

Hw/sw setup:

- CUDA 5.0 + driver API 5.0: 304.54, 304.60
- Ubuntu 12.04 AMD64
- Intel SB + GTX 680/K20
- On Titan-XP, our force kernel completes in 22ns for 1000 atoms
- But then the efficiency is less than 50%
 From 24k atoms, we have full efficiency
- If we could reach this for the smallest system, that kernel would run 2.5x faster
- This limits parallelization (strong scaling)



Exploiting multiple & high-priority streams



immediately so it can be returned faster

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



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





PME offload into separate stream

Much harder, but important for multi-GPU scaling



Multiple streams, neighbor list pruning, PME offload: 2.7X better! 8-10X faster than hand tuned SIMD assembly running on CPU



95.6k atoms in dodecahedron box



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 dual-pair list buffer:

- reduces overhead

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



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

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



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

<1 millisecond





Cray XC System Building Blocks

...

CRA





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

Intra-rank parallelisation: OpenMP today, future ? Efficient current parallelization of all algorithms using MPI + OpenMP

- We are convinced we are moving to a world where latency- and throughputoptimized 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.



Computational Bottleneck

Computing is the new experiment:

Cracking the Cryo-EM Image Reconstruction





Structural biology: Determine the structure of biomolecules from electron microscopy.

Take thousands of images each containing 100-1000 particles and use computers to reconstruct 3D electron density



Cryo-EM Image Reconstruction





RELION A Bayesian approach

p(A|B)p(B) = p(B|A)p(B) $p(A|B) = \frac{p(B|A)p(B)}{p(B)}$ Model Image data

This way, we can include all microscope CTF, noise estimates and possible orientational bias in the statistics, as well as the assumption of smoothly varying density



Particle picking 2D classification 3D classification 3D refinement ("autorefine")





Picking & extraction



Alignment & Classification

Inverse FFT & iterate

Parallelization



We scan many independent Views (100,000s) Objects (10) Pixels (0.25 Mpix)





34535	movaps xmm0, xmm4
34536	movaps xmm1, xmm4
34537	
34538	shufps xmm2, xmm6, <mark>0b10001000</mark>
34539	
34540	shufps xmm0, xmm5, <mark>0b10001000</mark>
34541	shufps xmm1, xmm5, <mark>0b11011101</mark>
34542	
34543	/* move ixO-izO to xmm4-xmm6 */
34544	movaps xmm4, [esp + _ix0]
34545	movaps xmm5, [esp + _iy0]
34546	movaps xmm6, [esp + _iz0]
34547	
34548	/* calc dr */
34549	subps xmm4, xmm0
34550	subps xmm5, xmm1
34551	subps xmm6, xmm2
34552	
34553	/* store dr */
34554	movaps [esp + _dx0], xmm4

ar and a second second





Good Kernels Express Parallelism, not Architecture

```
template <typename T>
_____global___ void cuda translate3D( T * g image in,
                                   T * g image out,
                                   int image size,
                                   int xdim,
                                   int ydim,
                                   int zdim,
                                   int dx,
                                   int dy,
                                   int dz)
   int tid = threadIdx.x;
   int bid = blockIdx.x;
   int x,y,z,xp,yp,zp,xy;
   int voxel=tid + bid*BLOCK SIZE;
   int new voxel;
   int xydim = xdim*ydim;
   if(voxel<image size)</pre>
      z = voxel / xydim;
      zp = z + dz;
      xy = voxel % xydim;
      y = xy / xdim;
      yp = y + dy;
      x = xy \% xdim;
      xp = x + dx;
      if (zp \ge 0 \& yp \ge 0 \& xp \ge 0 \& zp < zdim \& yp < ydim \& xp < xdim)
         new voxel = zp * xydim + yp * xdim + xp;
         if(new voxel >= 0 && new voxel < image size)</pre>
            g image out[new voxel] = g image in[voxel];
```

```
template <typename T>
void cpu translate3D(T * g_image_in,
                T * g image out,
                int
                         image size,
                int
                         xdim,
                int
                         ydim,
                int
                         zdim,
                int
                         dx,
                int
                         dy,
                int
                         dz)
   int x,y,z,xp,yp,zp,xy;
   int new voxel;
   for(int voxel=0; voxel<image size; voxel++)</pre>
      int xydim = xdim*ydim;
      z = voxel / xydim;
      zp = z + dz;
      xy = voxel % xydim;
      y = xy / xdim;
      yp = y + dy;
      x = xy % xdim;
      xp = x + dx;
      if (zp \ge 0 \& yp \ge 0 \& xp \ge 0 \& zp < zdim \& yp < ydim \& xp < xdim)
         new voxel = zp*xydim + yp*xdim + xp;
         if(new voxel >= 0 && new voxel < image size)</pre>
            g image out[new voxel] = g image in[voxel];
   }
```



"Ninja Coding": GROMACS vs. RELION

GROMACS

RELION

- tweaking, but if the compiler performs this well with only hints, why bother?
- anybody will now be x86-accelerated, even if they do not understand SIMD

Lines of raw SIMD code ~600,000

Used to be raw assembly, but now intrinsics

It is extremely likely that we could get even better speed-ups in RELION with manual

There is a huge advantage in only having a single code path. All modifications introduced by



Use GPU-specific features when you can

When we have switched to single precision, we realized we can use GPU textures to handle interpolation during image rotation.





When handling things dynamically, cudaMalloc() became bottleneck

We (Dari Kimanius) had to write a custom memory allocator for CUDA

in Queue

Allocate a large chunk of memory and handle it ourselves

Unified memory could theoretically be an alternative, but presently it is just as slow as CUDA handling memory



Single precision

- Most modern CPUs provide twice the throughput in single compared to double
- This also saves cache, memory bandwidth and RELION is quite memory-hungry
- Do you need more than 7 valid digits in the output?
- Converting scientific code to single-precision is not trivial: Most codes fail if you just do search & replace, but it can be made to work:
 - Identify sensitive parts (maximization step), and leave those in double
 - Only perform a few operations (exponentials) in double
 - Sum small numbers first, or use tree summation instead of linear order
 - Use strength-reduction algorithms (check open source math libraries)
- Beware of double-single-double conversions inside critical code paths





So, how did we do performance-wise? 14-core Xeon E5-2960v4 CPUs (est. cost \$85,000)

RELION-3: New x86 acceleration makes CPU competitive again!

A single quad-GPU workstation was on par with 10 nodes with dual ... and we were still often limited by our CPU code (a bit embarrassing...)

Preprocessing steps are even more impressive - 600x To find particles in the images, we rely on FFTs where we can use cuFFT

Picking quality

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 codes CUDA/AVX512/OpenCL 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

GROMACS: Berk Hess, Szilard Pall, Mark Abraham, Aleksei liupinov, John Eblen, Roland Shultz, Christian Wennberg, Viveca Lindahl RELION: Dari Kimanius, Björn Forsberg, Sjors Scheres, Alexey Amunts, Marta Carroni, Shintaro Aibara NVIDIA: Mark Berger, Duncan Poole, Julia Levites, Jiri Kraus, Nikolay Markovskiy **INTEL:** Charles Congdon, Sheng Fu, Kristina Kermanshahche, Yuping Zhao CSCS: Thomas Schulthess, Victor Holanda, Prashant Kanduri PDC: Erwin Laure

