

Hardware-aware Performance Engineering

Jiří Filipovič

fall 2019

Focus of the Lecture

We will learn how to optimize C/C++ code to get more performance from contemporary processors

- maximizing benefit from cache architecture
- writing code taking advantage of compiler auto-vectorization
- using multiple cores efficiently

We will not cover all interesting topics...

- only basic optimizations from each category
- no language-specific optimizations (inlining, proper usage of virtual functions etc.)
- no hardcore, assembly-level optimizations

Focus of the Lecture

All optimization techniques will be demonstrated on two examples

- so we will see how to change simple code in multiple steps, getting more and more speedup

Timing presented in this lecture will be obtained by using Intel C++

- it has more advanced loop optimization (mainly autovectorization) comparing to competitors
- described optimization methods are of course usable also with different compilers, but it may need some compiler tweaking or writing more complex code
- students can get free license from Intel, anybody from academia can access license via METACentrum

Motivation

Optimizations will target HW properties, which may change.
However, in world of x86 processors:

- caches appear in 80486 (1989), even 80386 had optional off-chip cache
- vector instructions appear in Pentium MMX (1996)
- multiple cores appear in Pentium D (2005), multi-socket configurations much earlier

So, knowledge gathered from this lecture should have sufficiently long lifetime :-).

Motivation

What if we ignore HW properties and just write good algorithms?

- suppose we have Core i7-5960X processor: 8 cores, Skylake architecture (AVX2 + FMA3 instructions)
- L1 cache latency: 4 cycles, L2: 12 cycles, L3: 42 cycles, RAM about 200 cycles
- vectorized code finishes up to 32 single-precision operations per cycle
- parallelized code take advantage of 8 (or 16 virtual) cores

You can get a lot of speedup by hardware-aware programming.

Demonstration Examples

We will demonstrate optimization methods using two examples

- I have tried to find as simple as possible computational problems, which still expose a lot of opportunity for various optimization techniques

The code is not very abstract or generic

- in productivity-optimized programming, we want to hide how are algorithms performed, how are data stored etc.
- however, when optimizing code, we have to focus on implementation details, thus, code looks more "old school"
- in practice, usually very small fraction of source code is performance-critical, so different programming style for optimized code is not a problem

Electrostatic Potential Map

Important problem from computational chemistry

- we have a molecule defined by position and charges of its atoms
- the goal is to compute charges at a 3D spatial grid around the molecule

In a given point of the grid, we have

$$V_i = \sum_j \frac{w_j}{4\pi\epsilon_0 r_{ij}}$$

Where w_j is charge of the j -th atom, r_{ij} is Euclidean distance between atom j and the grid point i and ϵ_0 is vacuum permittivity.

Electrostatic Potential Map

Initial implementation

- suppose we know nothing about HW, just know C++
- algorithm needs to process 3D grid such that it sums potential of all atoms for each grid point
- we will iterate over atoms in outer loop, as it allows to precompute positions of grid points and minimizes number of accesses into input/output array

Electrostatic Potential Map

```

void coulomb(const sAtom* atoms, const int nAtoms,
             const float gs, const int gSize, float *grid) {
    for (int a = 0; a < nAtoms; a++) {
        sAtom myAtom = atoms[a];
        for (int x = 0; x < gSize; x++) {
            float dx2 = powf((float)x * gs - myAtom.x, 2.0f);
            for (int y = 0; y < gSize; y++) {
                float dy2 = powf((float)y * gs - myAtom.y);
                for (int z = 0; z < gSize; z++) {
                    float dz = (float)z * gs - myAtom.z;
                    float e = myAtom.w / sqrtf(dx2 + dy2 + dz*dz);
                    grid[z*gSize*gSize + y*gSize + x] += e;
                }
            }
        }
    }
}

```

Histogram

Used in many scientific applications

- computes a frequency of input values occurrence in defined intervals
- in our example, we will compute histogram of population age in uniformly-sized intervals
- input is vector of ages (floating point) and interval size, output is histogram

Histogram

```
void hist(const float* age, int* const hist, const int n,
          const float group_width, const int m) {
    for (int i = 0; i < n; i++) {
        const int j = (int) ( age[i] / group_width );
        hist[j]++;
    }
}
```

Benchmarking

We will benchmark codes on pretty average desktop system

- 4 cores
- AVX2 (256-bit vectors), no FMA

Guess speedup of original codes :-).

Cache Memories

Why we have cache memories in modern processors?

- main memory is too slow (both latency and bandwidth) comparing to compute cores
- we can build much faster, but also more expensive memory
- cache is fast memory, which temporary keeps parts of larger and slower memories

Cache Implementation

How it is working?

- multiple levels (usually L1 and L2 private for core, L3 shared)
- accessed by cache lines (64 bytes on Intel architectures)
- when data are accessed, they are stored in cache and kept until cache line is needed for another data
- limited associativity (each cache line may cache only defined parts of main memory)
- parallel access into memory – cache lines must be somehow synchronized (broadcast, invalidation)

Optimization for Cache

Optimization for spatial locality

- access consequent elements
- align data to a multiple of cache line size
- otherwise only part of transferred data is used

Optimization for temporal locality

- when data element needs to be accessed multiple times, perform accesses in a short time
- otherwise it may be removed from cache due to its limited capacity or associativity

Omit inefficient usage

- conflict misses
- false sharing

Electrostatic Potential Map

```

void coulomb(const sAtom* atoms, const int nAtoms,
             const float gs, const int gSize, float *grid) {
    for (int a = 0; a < nAtoms; a++) {
        sAtom myAtom = atoms[a];
        for (int x = 0; x < gSize; x++) {
            float dx2 = powf((float)x * gs - myAtom.x, 2.0f);
            for (int y = 0; y < gSize; y++) {
                float dy2 = powf((float)y * gs - myAtom.y, 2.0f);
                for (int z = 0; z < gSize; z++) {
                    float dz = (float)z * gs - myAtom.z;
                    float e = myAtom.w / sqrtf(dx2 + dy2 + dz*dz);
                    grid[z*gSize*gSize + y*gSize + x] += e;
                }
            }
        }
    }
}

```

Evaluation

We have compiled the code above with vectorization switched off (as we are interested in effects of memory access only)

- 31.6 millions of atoms evaluated per second (MEvals/s) using $256 \times 256 \times 256$ grid and 4096 atoms
- by changing grid size to $257 \times 257 \times 257$, performance changes to 164.7 Mevals/s

Interpretation

- strong dependence on input size indicates problems with cache associativity
- even 164.7 Mevals/s is not very good result, considering 8 floating point operations are performed in innermost loop

Spatial Locality

We are interested in the innermost loop

- it defines memory access pattern (i.e. which elements are accessed consequently)
- the innermost loop runs over z , which creates large memory strides in accessing grid
- when grid size is power of two, columns hits the same associativity region

Optimization

- we need to rearrange loops: the innermost loop should iterate through x

Spatial Locality

```

void coulomb(const sAtom* atoms, const int nAtoms,
             const float gs, const int gSize, float *grid) {
    for (int a = 0; a < nAtoms; a++) {
        sAtom myAtom = atoms[a];
        for (int z = 0; z < gSize; z++) {
            float dz2 = powf((float)z * gs - myAtom.z, 2.0f);
            for (int y = 0; y < gSize; y++) {
                float dy2 = powf((float)y * gs - myAtom.y, 2.0f);
                for (int x = 0; x < gSize; x++) {
                    float dx = (float)x * gs - myAtom.x;
                    float e = myAtom.w / sqrtf(dx*dx + dy2 + dz2);
                    grid[z*gSize*gSize + y*gSize + x] += e;
                }
            }
        }
    }
}

```

Evaluation

Performance measurement

- 371.8 Mevals/s using $256 \times 256 \times 256$ grid and 4096 atoms
- no sensitivity to changing grid size (no cache associativity problem)
- much better spatial locality

Analysis of cache pattern

- each atom is applied to the whole grid
- poor temporal locality (grid is too large structure)

Temporal Locality

Atoms array is much smaller than grid

- we can rearrange loops to iterate over atoms in the innermost loop: z-y-x-a
- alternatively, we may apply atom forces per rows of a grid, creating iteration order z-y-a-x
- or tiling may be used

Memory tiling

- we break some loop into nested loops, such that outer loop iterates with step $s > 1$ and those steps are performed in some inner loop
- multiple loops may be tiled
- in our example, we will tile loop running over atoms

Tiled Algorithm

```

const int TILE = 16;
sAtom myAtom[TILE]; float dy2[TILE], dz2[TILE];
for (int a = 0; a < nAtoms; a+=TILE) {
    myAtom[0:TILE] = atoms[a:a+TILE];
    for (int z = 0; z < gSize; z++) {
        for (int aa = 0; aa < TILE; aa++)
            dz2[aa] = powf((float)z * gs - myAtom[aa].z, 2.0f);
        for (int y = 0; y < gSize; y++) {
            for (int aa = 0; aa < TILE; aa++)
                dy2[aa] = powf((float)y * gs - myAtom[aa].y, 2.0f);
            for (int x = 0; x < gSize; x++) {
                float e = 0.0f;
                for (int aa = 0; aa < TILE; aa++) {
                    float dx = (float)x * gs - myAtom[aa].x;
                    e += myAtom[aa].w / sqrtf(dx*dx + dy2[aa] + dz2[aa]);
                }
                grid[z*gSize*gSize + y*gSize + x] += e;
            }
        }
    }
}

```

Evaluation

Note that autovectorization is switched off for all implementations.

Implementation	Performance	speedup
Naive (grid 257)	164.7	n/a
Spatial loc.	371.8	2.26×
ZYXA	359.7	2.18×
ZYAX	382.2	2.32×
Tiled	476.9	2.9×

Temporal locality brings only minor improvement, but it may change when instructions are optimized/code is parallelized.

Histogram

```
void hist(const float* age, int* const hist, const int n,
          const float group_width, const int m) {
    for (int i = 0; i < n; i++) {
        const int j = (int) ( age[i] / group_width );
        hist[j]++;
    }
}
```

Memory access is already consequent in `age`. The random access into `hist` cannot be omitted. So, nothing to optimize so far...

Vector Instructions

Modern processors have complex logic preparing instructions

- arithmetical units are relatively cheap
- when instruction is to be executed, it may process multiple data elements in parallel

Data-parallel programming

- the same instruction is applied onto multiple data (SIMD model)
- explicit usage: we need to generate vector instructions

Vector Instructions

Vector instructions

- the same operation is applied to a short vector
- mainly arithmetic operations, may be masked, may contain support for reduction, binning etc.
- vector length depends on data type and instruction set, e.g. AVX2 works with vector of size 256 bytes, so 8 32-bit numbers or 4 64-bit numbers are processed in parallel

Vectorization in C/C++

- explicit: inline assembly or intrinsics
- implicit: compiler generates vector instructions automatically

Automatic Vectorization

Better portability

- the code can be compiled for any vector instruction set

Supported in modern compilers

- however, it is difficult task, so allowing compiler to vectorize code needs programmer assist

Automatic Vectorization

Current limitations

- only innermost for loops are vectorized
- number of iterations must be known when loop is entered, or (preferably) at compilation time
- memory access must be regular, ideally with unit stride (i.e. consequent elements are accessed in vector instructions)
- vector dependence usually disallows vectorization

Vector Dependence

Vector dependence

- the for loop cannot be vectorized, if there is flow dependence between iterations
- however, compiler may wrongly assume vector dependence (it must be conservative to generate correct code)
- `#pragma ivdep` (Intel) or `#pragma GCC ivdep` (gcc) instruct compiler to ignore assumed vector dependences (true dependence still disallows vectorization)

Vector Dependence

```
float a[n], b[n];
for (int i = 0; i < n; i++)
    a[i] = b[i]*2.0f;
```

No vector dependence, the code is vectorized.

```
void foo(float* a, const float* b, int n) {
    for (int i = 0; i < n; i++)
        a[i] = b[i]*2.0f;
}
```

The compiler must generate correct code also for pointer aliasing (i.e. when a and b overlaps): it generates vectorized and non-vectorized code with runtime check, or not vectorize at all. We may help the compiler using `restrict` quantifier with a and b, or use `ivdep` pragma.

Contiguous Memory Access

```

struct vec{
    float x,y;
};
vec v[n];
for (int i = 0; i < n; i++)
    v[i].x *= 2.0f;

```

The loop is vectorized, however, access into `v` is strided. Typical optimization is transferring array of structures (AoS) to structure of arrays (SoA).

```

struct vec{
    float *x;
    float *y;
};
vec v;
// allocation ...
for (int i = 0; i < n; i++)
    v.x[i] *= 2.0f;

```

Loop Strip-mining

If part of the code within a loop cannot be vectorized

- we split loop into two nested loops (similarly to tiling)
- we divide the inner loop according to vectorization possibility into vectorizable loop(s) and non-vectorizable loop(s)

Electrostatic Potential Map

Naive implementation has assumed dependence, which needs to be manually fixed.

```

for (int a = 0; a < nAtoms; a++) {
    sAtom myAtom = atoms[a];
    for (int x = 0; x < gSize; x++) {
        float dx2 = powf((float)x * gs - myAtom.x, 2.0f);
        for (int y = 0; y < gSize; y++) {
            float dy2 = powf((float)y * gs - myAtom.y);
            for (int z = 0; z < gSize; z++) {
                float dz = (float)z * gs - myAtom.z;
                float e = myAtom.w / sqrtf(dx2 + dy2 + dz*dz);
                #pragma ivdep
                grid[z*gSize*gSize + y*gSize + x] += e;
            }
        }
    }
}

```

Electrostatic Potential Map

AZYX and ZYAX Innermost Loop

```
for (int x = 0; x < gSize; x++) {  
    float dx = (float)x * gs - myAtom.x;  
    float e = myAtom.w / sqrtf(dx*dx + dy2 + dz2);  
    grid[z*gSize*gSize + y*gSize + x] += e;  
}
```

The loop is automatically vectorized without problems.

Electrostatic Potential Map

ZYXA implementation

```

for (int z = 0; z < gSize; z++) {
    for (int y = 0; y < gSize; y++) {
        for (int x = 0; x < gSize; x++) {
            float e = 0.0f;
            for (int a = 0; a < nAtoms; a++) {
                sAtom myAtom = atoms[a];
                float dx = (float)x * gs - myAtom.x;
                float dy = (float)y * gs - myAtom.y;
                float dz = (float)z * gs - myAtom.z;
                e += myAtom.w / sqrtf(dx*dx + dy*dy + dz*dz);
            }
            grid[z*gSize*gSize + y*gSize + x] += e;
        }
    }
}

```

ZYXA and Tiled

The innermost loop is difficult to vectorize

- strided memory access into atoms elements
- reduction

Two possible solutions

- AoS to SoA optimization
- vectorization of outer loop running over x

SoA

```

for (int z = 0; z < gSize; z++) {
    for (int y = 0; y < gSize; y++) {
        for (int x = 0; x < gSize; x++) {
            float e = 0.0f;
            for (int a = 0; a < nAtoms; a++) {
                float dx = (float)x * gs - atoms.x[a];
                float dy = (float)y * gs - atoms.y[a];
                float dz = (float)z * gs - atoms.z[a];
                e += atoms.w[a] / sqrtf(dx*dx + dy*dy + dz*dz);
            }
            grid[z*gSize*gSize + y*gSize + x] += e;
        }
    }
}

```

Outer-loop Vectorization

```

for (int z = 0; z < gSize; z++) {
    for (int y = 0; y < gSize; y++) {
        #pragma simd
        for (int x = 0; x < gSize; x++) {
            float e = 0.0f;
            for (int a = 0; a < nAtoms; a++) {
                sAtom myAtom = atoms[a];
                float dx = (float)x * gs - myAtom.x;
                float dy = (float)y * gs - myAtom.y;
                float dz = (float)z * gs - myAtom.z;
                e += myAtom.w / sqrtf(dx*dx + dy*dy + dz*dz);
            }
            grid[z*gSize*gSize + y*gSize + x] += e;
        }
    }
}

```

All implementations

We will use `restrict` quantifier

- otherwise, compiler may expect aliasing even between atoms and `grid` and give up vectorization

Evaluation

Implementation	Performance	speedup	(vect. speedup)
Naive (grid 257)	164.7	n/a	n/a
Naive vec. (grid 257)	330.6	2.01×	2.01×
Spatial loc.	1838	11.2×	4.94×
ZYXA outer	2189	13.3×	6.09×
ZYXA SoA	2203	13.4×	6.12×
ZYAX	2197	13.3×	5.75×
Tiled outer	2577	15.6×	5.4×
Tiled SoA	2547	15.5×	5.34×

Histogram

```
void hist(const float* age, int* const hist, const int n,
          const float group_width, const int m) {
    for (int i = 0; i < n; i++) {
        const int j = (int) ( age[i] / group_width );
        hist[j]++;
    }
}
```

The loop cannot be vectorized due to dependency in `hist`. We will use strip-mining.

Histogram

```

void hist(const float* restrict age, int* const restrict hist,
         const int n, const float group_width, const int m) {
    const int vecLen = 16;
    //XXX: this algorithm assumes n%vecLen == 0.
    for (int ii = 0; ii < n; ii += vecLen) {
        int histIdx[vecLen];
        for (int i = ii; i < ii + vecLen; i++)
            histIdx[i-ii] = (int) ( age[i] / group_width );
        for (int c = 0; c < vecLen; c++)
            hist[histIdx[c]]++;
    }
}

```

Division is heavy-weight operation, we will remove it.

Histogram

```
void hist(const float* restrict age, int* const restrict hist,
         const int n, const float group_width, const int m) {
    const int vecLen = 16;
    const float invGroupWidth = 1.0f/group_width;
    //XXX: this algorithm assumes n%vecLen == 0.
    for (int ii = 0; ii < n; ii += vecLen) {
        int histIdx[vecLen];
        for (int i = ii; i < ii + vecLen; i++)
            histIdx[i-ii] = (int) ( age[i] * invGroupWidth );
        for (int c = 0; c < vecLen; c++)
            hist[histIdx[c]]++;
    }
}
```

Evaluation

Implementation	Performance	speedup
Naive	1020 MB/s	n/a
Vectorized	2455 MB/s	2.41×
Removed div.	4524 MB/s	4.44×

Parallelization

Why we have multicore processors?

- processors frequency is no longer substantially improved due to energy requirements
- however, with new manufacturing processes, it is possible to build smaller cores, thus, multiple cores can be integrated into a die

Programming multiple cores

- coarse-grained parallelism (compared to vectorization)
- threads are asynchronous by default (MIMD model), synchronization is explicit and relatively expensive

Parallelization in C/C++

Thread-level parallelism in C/C++

- many possible ways to parallelize a code: pthreads, Boost threads, TBB etc.
- we will use OpenMP in our examples, as it broadly-supported standard and it requires only small changes in our code
- however, optimization principles are general and can be used with any parallelization interface

OpenMP

OpenMP standard

- for shared-memory parallelism
- uses pragmas to declare, which parts of the code runs in parallel
- very easy to use, but writing efficient code may be challenging (much like in other interfaces)
- implements fork-join model
- standard, implemented in all major C/C++ compilers

OpenMP

The parallel region of the code is declared by `#pragma omp parallel`

```
//serial code
const int n = 100;
#pragma omp parallel
{
    //parallel code
    printf("Hello from thread %d\n", omp_get_thread_num());
    //parallel loop, iterations order is undefined
    #pragma omp for
    for (int i = 0; i < n; i++) {
        //iteration space is distributed across all threads
        printf("%d ", i);
    }
}
//serial code
```

OpenMP

We can define private and shared variables

- `#pragma omp parallel for private(a) shared(b)`
- variables declared before parallel block are shared by default
- `private` statement creates private copy for each thread

Thread synchronization

- we can define critical section by `#pragma omp critical`
- or use lightweight atomic operations, which are restricted to simple scalar operations, such as `+` `-` `*` `/`

Electrostatic Potential Map

Which loop can be parallelized?

- AZYX: loop running over atoms would need synchronization, so we prefer to parallelize loop running over Z, Y or X
- ZYXA: we can parallelize up to three outermost loops
- ZYAX: we can parallelize up to two outermost loops
- tiled: we can parallelize loop running over Z, Y and X

Which loop to parallelize?

- enter and exit of the loop is synchronized
- we want to minimize number of synchronizations, so we will parallelize loops performing more work
- to scale better, we may collapse n perfectly-nested loops using `#pragma omp for collapse(n)`

ZYXA Example

```

#pragma omp parallel for
for (int z = 0; z < gSize; z++) {
    for (int y = 0; y < gSize; y++) {
        #pragma simd
        for (int x = 0; x < gSize; x++) {
            float e = 0.0f;
            for (int a = 0; a < nAtoms; a++) {
                sAtom myAtom = atoms[a];
                float dx = (float)x * gs - myAtom.x;
                float dy = (float)y * gs - myAtom.y;
                float dz = (float)z * gs - myAtom.z;
                e += myAtom.w / sqrtf(dx*dx + dy*dy + dz*dz);
            }
            grid[z*gSize*gSize + y*gSize + x] += e;
        }
    }
}

```

ZYXA Example

```

#pragma omp parallel for collapse(2)
for (int z = 0; z < gSize; z++) {
    for (int y = 0; y < gSize; y++) {
        #pragma simd
        for (int x = 0; x < gSize; x++) {
            float e = 0.0f;
            for (int a = 0; a < nAtoms; a++) {
                sAtom myAtom = atoms[a];
                float dx = (float)x * gs - myAtom.x;
                float dy = (float)y * gs - myAtom.y;
                float dz = (float)z * gs - myAtom.z;
                e += myAtom.w / sqrtf(dx*dx + dy*dy + dz*dz);
            }
            grid[z*gSize*gSize + y*gSize + x] += e;
        }
    }
}

```

Evaluation

Implementation	Performance	speedup	(par. speedup)
Naive (grid 257)	164.7	n/a	n/a
Spatial loc.	2272	13.8×	1.24×
ZYXA outer	7984	48.5×	3.62×
ZYAX	8092	49.1×	3.68×
Tiled SoA	9914	60.2×	3.92×

Histogram

```

void hist(const float* restrict age, int* const restrict hist,
          const int n, const float group_width, const int m) {
    const int vecLen = 16;
    const float invGroupWidth = 1.0f/group_width;
    for (int ii = 0; ii < n; ii += vecLen) {
        int histIdx[vecLen];
        for (int i = ii; i < ii + vecLen; i++)
            histIdx[i-ii] = (int) ( age[i] * invGroupWidth );
        for (int c = 0; c < vecLen; c++)
            hist[histIdx[c]]++;
    }
}

```

We can parallelize the outer loop and atomically update hist.

Histogram

```
void hist(const float* restrict age, int* const restrict hist,
         const int n, const float group_width, const int m) {
    const int vecLen = 16;
    const float invGroupWidth = 1.0f/group_width;
    #pragma omp parallel for
    for (int ii = 0; ii < n; ii += vecLen) {
        int histIdx[vecLen];
        for (int i = ii; i < ii + vecLen; i++)
            histIdx[i-ii] = (int) ( age[i] * invGroupWidth );
        for (int c = 0; c < vecLen; c++)
            #pragma omp atomic
            hist[histIdx[c]]++;
    }
}
```

Evaluation

Implementation	Performance	speedup
Naive	1020 MB/s	n/a
Vectorized	2455 MB/s	2.41×
Removed div.	4524 MB/s	4.44×
Parallel	290.2 MB/s	0.28×

So, overhead of atomic operations is too high...

Histogram

```

void hist(const float* restrict age, int* const restrict hist,
         const int n, const float group_width, const int m) {
    const int vecLen = 16;
    const float invGroupWidth = 1.0f/group_width;
    #pragma omp parallel
    {
        int histPriv[m];
        histPriv[:] = 0;
        int histIdx[vecLen];
        #pragma omp for
        for (int ii = 0; ii < n; ii += vecLen) {
            for (int i = ii; i < ii + vecLen; i++)
                histIdx[i-ii] = (int) ( age[i] * invGroupWidth );
            for (int c = 0; c < vecLen; c++)
                histPriv[histIdx[c]]++;
        }
        for (int c = 0; c < m; c++)
            #pragma omp atomic
            hist[c] += histPriv[c];
    }
}

```

Evaluation

Implementation	Performance	speedup
Naive	1020 MB/s	n/a
Vectorized	2455 MB/s	2.41×
Removed div.	4524 MB/s	4.44×
Parallel	290.2 MB/s	0.28×
Parallel opt.	20086 MB/s	19.7×

Histogram

This implementation is OK on our system

- however, combination of small m (wide groups, for which the histogram is computed) and highly-parallel system decreases performance significantly
- false sharing issue!

False sharing

- array `histPriv` is created by all threads
- if the array is small, multiple arrays may share the same cache line
- so write access to the independent array causes frequent synchronization of the cache
- very simple optimization: padding `histPriv` array

Conclusion

We have demonstrated basic hardware-aware optimization methods

- there is still a lot of uncovered topics
- however, knowledge of basic optimization methods can still make a big difference in performance

We have demonstrated optimization on two examples

- electrostatic potential map: up to $60\times$ speedup
- histogram: up to $20\times$ speedup
- this is much more, than people usually expect...

More info

- Intel/AMD optimization manuals
- Colfax Research courses