# Code Optimization

Dalibor Djukanovic

# Generalities

This is an overview talk.
But technical.

If you have specific questions/comments concerning code performance issues contact me
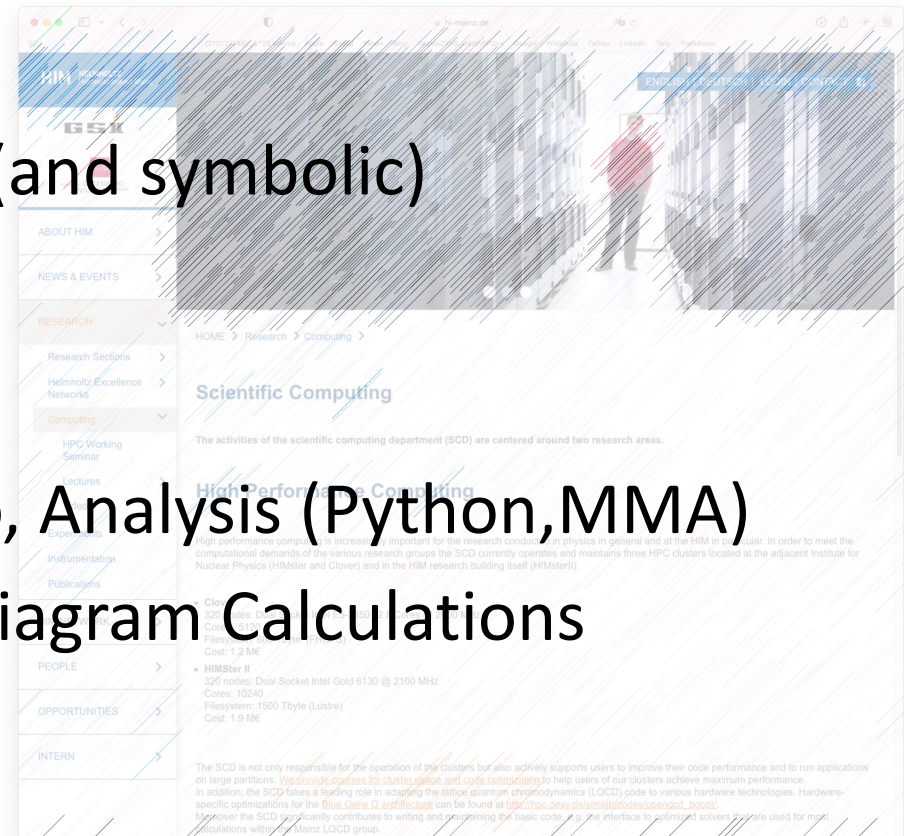
djukanov@uni-mainz.de

Reminder:
Fortnightly HPC seminar

HIM HELMHOLTZ
Helmholtz-Institut Mainz

# My Background

- Member of the Nuclear Theory Group

- Doing Calculations in:
  - Lattice QCD numerical (and symbolic)
  - EFT (mostly symbolic)

- Optimization interests:
  - Numerical: MonteCarlo, Analysis (Python,MMA)
  - Symbolical: Feynman Diagram Calculations (FORM,MMA)
  - Get results fast

- Not a performance architect:
  - Some of the thing I say might not even be wrong!

# Contents - I

- **Performance = Divide & Conquer**

- **Core Level Optimization:**
SIMD - Domain
Optimize Code to use Oncore vector units

- Node Level Optimization:
SIMT – Domain
Optimize Code to use threads working in parallel on data

- Application Level Optimization:
Algorithmic Improvements, Avoiding Data Transfer, Hiding Latencies

**HIM** HELMHOLTZ
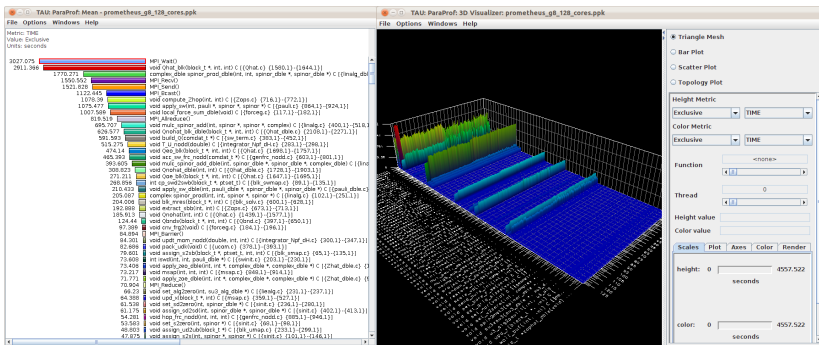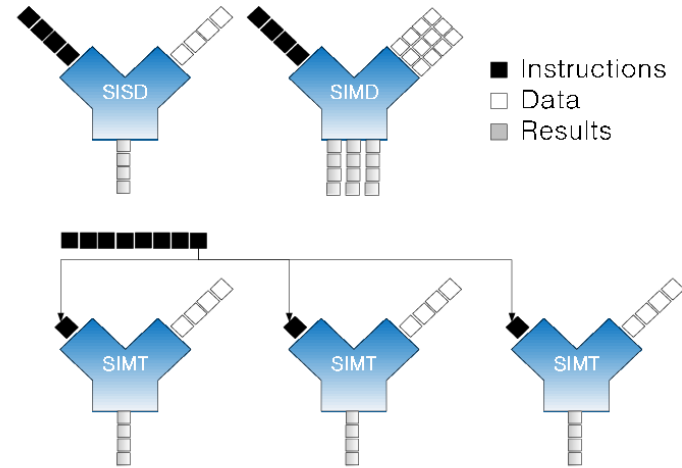Helmholtz-Institut Mainz

# Contents - II

- I will try to cover the following aspects:
  - How can you get the most out of a modern CPU
  - How do you know where your code spends most of the time
- Takeaways:
  - Be modest. 20x improvement will (probably) not happen.
  - Try to understand your codes performance ceiling and the bottlenecks.
  - Get used to tools (profilers, debuggers, … ).

HIM **HELMHOLTZ** Helmholtz-Institut Mainz

# Outline

- Performance Paradigms

- Profiling

- Assembler Coding

Bottleneck - the part of a bottle that slows the flow of liquid.
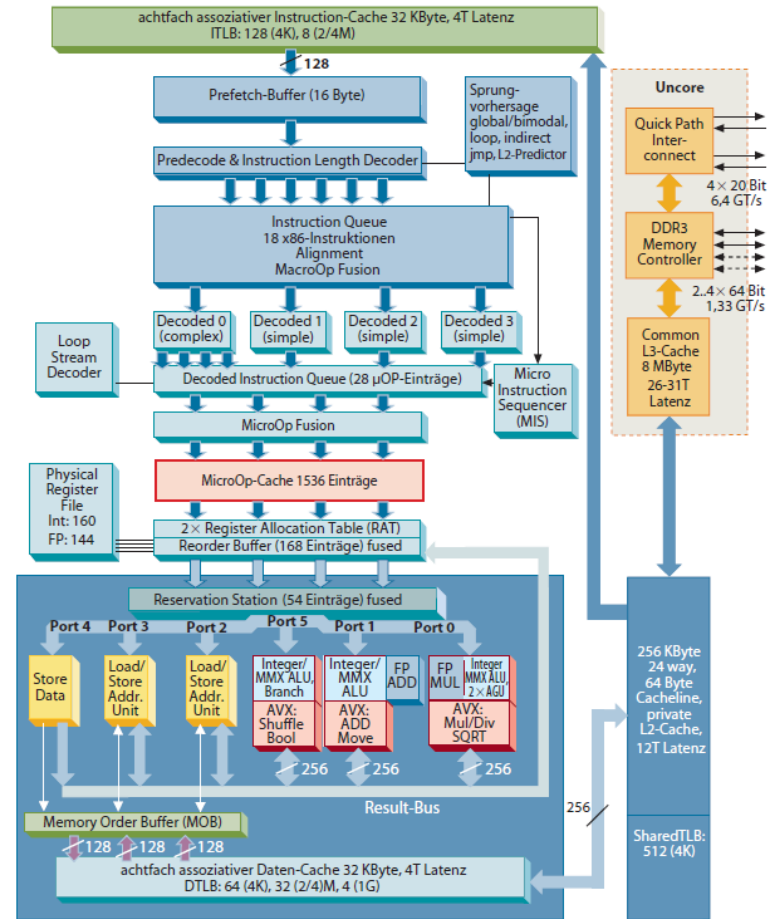
- ## Moore



- ## CPUs went multicore/vector
- ## RAM does not keep up

(von Neumann-Bottleneck)

- ## Latency Hierarchies

HIM  HELMHOLTZ
Helmholtz-Institut Mainz
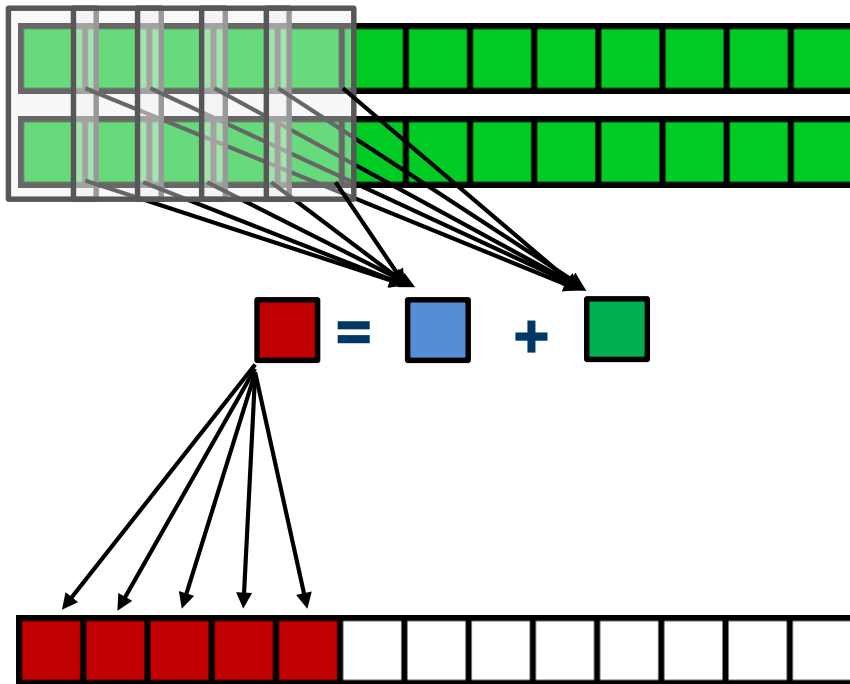
# Many Core

- Performance gains through parallelism
  - Thread Level:
    - More cores running tasks in parallel
  - Core Level:
    - Execution units are vector machines

- Having one set of instructions and data we distinguish



  - SIMT: Single Instruction Multiple Threads
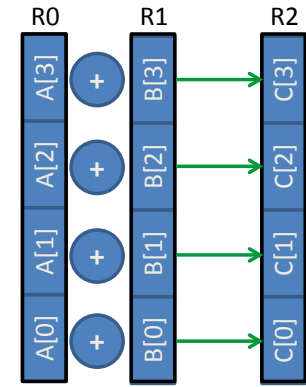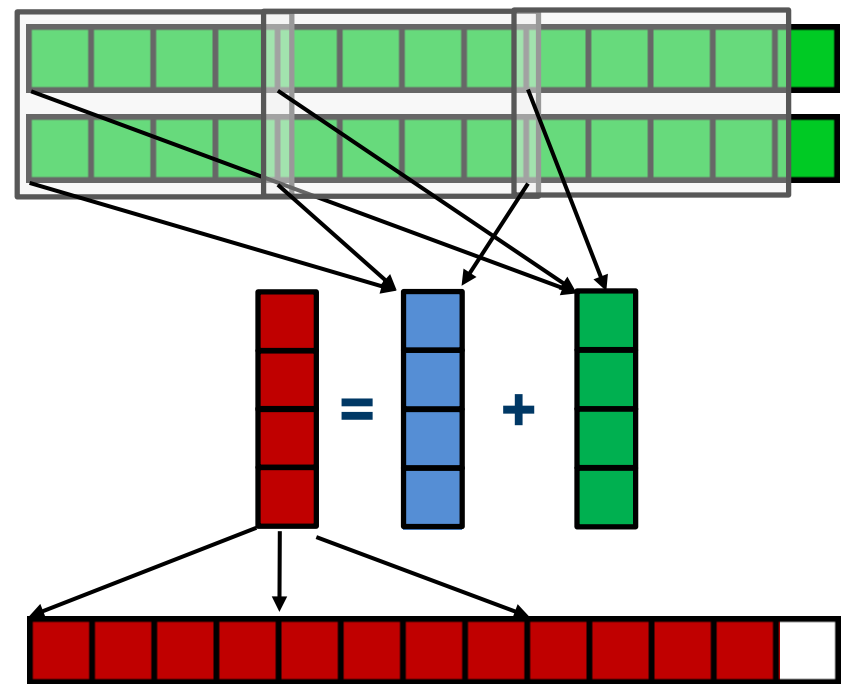  - SIMD: Single Instruction Multiple Data

HIM **HELMHOLTZ** Helmholtz-Institut Mainz

# SIMD

- Vector Execution of Scalar Code
- Works on registers

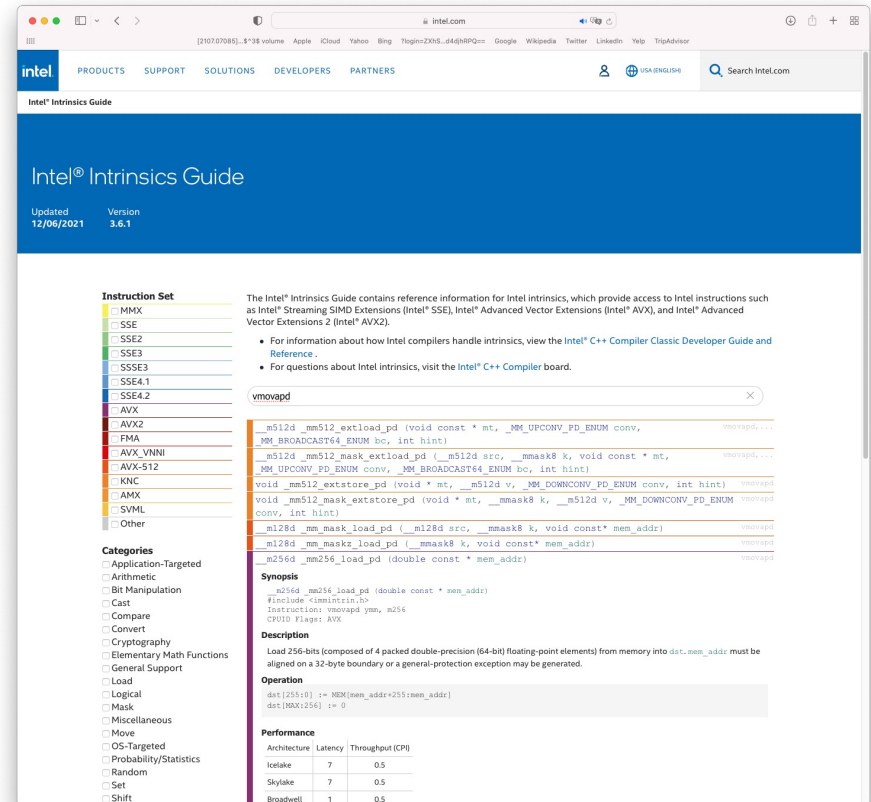**Scalar execution**

**SIMD execution**

# SIMD

- Width of registers depend on CPU
  - 128 Bit SSE
  - 256 Bit AVX
  - 512 Bit AVX512

- Instruction set also depends on CPU

- Theoretical Peak (Himster2)
  Intel(R) Xeon(R) Gold 6130 CPU @ 2.10GHz

- 16 (cores) x 8 (doubles in regs)  x 2 (fma)
  x 2 (fma units) x 2.1 GHz =  1000 Gflops/s

- 2 Tflops per node

HIM **HELMHOLTZ** Helmholtz-Institut Mainz

# What to Expect?

- MogonII/HIMster2 still in Top500 List

- Pos 364

- $R_{theo}$ = 2800 Tflops

- $R_{linpack}$ = 1967 Tflops

- @ 657 kWatt

- 70 % of peak looks promising

- Linpack not very common work load

HIM HELMHOLTZ
Helmholtz-Institut Mainz

# Performance Models

- ## Roofline Model

  - – Assume perfect overlap of data transfer and ops

  - – Latency is ignored

  - – Count amount of „work" against „data"

```
for(i=0;i<num;i++)
{
  c[i]=x[i]+y[i];
}
```

Data:  2 x load + 1 store = 24 byte (30)

Work: 1 x add  = 1 Flop

Intensity: Work/Data = 1 Flop/24 B

  - – Roofline Model:

    - Take minimum of applicable Peak perf vs BW of slowest data path (suppose 190 Gbyte/s)

    - $P_{roofline}$ = min(2 x 1075.2 Gflops, 1/24 x 190 Gflops)
       = 8 Gflops < 1 % peak

**HIM** **HELMHOLTZ** Helmholtz-Institut Mainz

# Roofline Model

- Ways to improve:
  - Increase Intensity
    Remember
    $I = Work/Data$



  - Avoid slow paths
    Keep data in caches

**HIM**  **HELMHOLTZ**
Helmholtz-Institut Mainz

# Roofline Model

- Simple model tells you „speed of light" for your kernel

- Refinements Possible:
  - Take into account cache hierarchy
  - Take into account latencies
  - Take into account pipeline design of CPUs

KERNCRAFT

- Usually I stop here at the simple model and try to get to „speed of light" according to Roofline

# Memory Alignment - Interlude

- **<u>Usually</u>** performance problem may be attributed to bad memory layout

- Ideal is stride-1 access pattern

- NUMA
  - Avoid "foreign memory"
  - Process pinning important
  - Topic for HPC-seminar?

Non-Uniform Memory Access (NUMA)

```
[djukanov@x0793 mystream]$ numactl -H
available: 2 nodes (0-1)
node 0 cpus: 0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47
node 0 size: 46750 MB
node 0 free: 44717 MB
node 1 cpus: 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63
node 1 size: 48339 MB
node 1 free: 47164 MB
node distances:
node   0   1
  0:  10  21
  1:  21  10
```

# Profiling

- ## I want to SIMD my code, where do I start?

  - Profile code

    - Manually using timings

    - Automatically using some tool

      - Instrumentation

        » Gprof, TAU …

      - Sampling

        » Intel Vtune, perf …

# Perf

- Perf is a profiling tool already installed on most Linux systems

- No need to recompile binary

- May attach to running processes

```
gcc -g <source> -o <bin>
perf record ./<bin>
perf report
```

**HIM**  HELMHOLTZ
Helmholtz-Institut Mainz

# Perf

- Attach to running process
  perf top –p <pid>

- Right:
  Sample of Jupyter
  Notebook run

  - HDF5 reading

  - Creating new Dict

  - Moving data
    around

  - Some kernel thread messing around at 6 %
    (Omnipath)

```
● ● ●        djukanov — djukanov@x0793:~ — ssh ‹ ssh -L 8889:x0793:8888 mogon2 — 80×24
Samples: 123  of event 'cycles:u', 4000 Hz, Event count (approx.): 10079838 lost
Overhead   Shared Object                    Symbol
   9,87%   libpython2.7.so.1.0              [.] PyDict_New
   6,02%   _objects.so                      [.] 0x0000000000011f17
   5,91%   [kernel]                         [k] 0xffffffffb2a00b87
   5,59%   libhdf5-9028dcc4.so.103.0.0      [.] memcpy@plt
   5,23%   libhdf5-9028dcc4.so.103.0.0      [.] H5SL_search
   5,03%   libpython2.7.so.1.0              [.] PyMethod_New
   4,94%   utils.so                         [.] PyObject_Size@plt
   4,65%   libhdf5-9028dcc4.so.103.0.0      [.] H5FD_read@plt
   4,63%   _multiarray_umath.so             [.] PyArray_DescrFromType
   4,54%   libc-2.28.so                     [.] __memmove_sse2_unaligned_erms
   4,31%   libhdf5-9028dcc4.so.103.0.0      [.] H5I_get_type
   4,31%   libhdf5-9028dcc4.so.103.0.0      [.] H5open
   3,61%   libhdf5-9028dcc4.so.103.0.0      [.] H5O_msg_read@plt
   3,51%   libpython2.7.so.1.0              [.] 0x00000000000815c0
   2,29%   libpython2.7.so.1.0              [.] PyObject_GetAttr
   2,16%   libpython2.7.so.1.0              [.] 0x00000000000aae18
   2,01%   libhdf5-9028dcc4.so.103.0.0      [.] H5S_extent_release
   1,90%   libhdf5-9028dcc4.so.103.0.0      [.] H5F_block_read
   1,79%   libhdf5-9028dcc4.so.103.0.0      [.] H5F__accum_read
   1,62%   libpthread-2.28.so               [.] __libc_read
   1,61%   libhdf5-9028dcc4.so.103.0.0      [.] H5F_get_eoa
For a higher level overview, try: perf top --sort comm,dso
```

HIM HELMHOLTZ
Helmholtz-Institut Mainz

# TAU

- Tuning and Analysis Tookit

- Application is instrumented in source code automatically by replacing CC with tau_cc.sh

  CC=tau_cc.sh -optTauSelectFile=./select.file

- Works with MPI

# Profiling General

- Helps you identify HotSpots

- Helps you identify bottlenecks (performance counters)

- Be careful:
  - Overhead due to instrumentation or sampling

- Notable alternative

**HIM**  HELMHOLTZ
Helmholtz-Institut Mainz

# Profiling Caveats

- For the experts:
  - Timing code not easy on OutOfOrder execution archs and different Clocks
  - rdtsc: Instruction getting time stamp counters
  - Use serializing muops like cpuid

```
static __inline__ unsigned long long rdtsc(void)
{
    unsigned hi, lo;
    __asm__ __volatile__ ("rdtsc" : "=a"(lo), "=d"(hi));
    return ( (unsigned long long)lo)|( ((unsigned long long)hi)<<32 );
}
```

- I haven't seen big differences!
  Maybe if you want cycle accurate measurements?

# Assembler Coding

- Pros assembler:
  - Control over what is done
  - You can do better than the compiler (most of the time)

- Cons assembler:
  - Hard to maintain
  - Hard to read
  - Have to take care of everything

- Only use:
  - For best possible perf
  - For small kernels

```
#define _load_cst(c) \
__asm__ __volatile__ ("vbroadcastss %0, %%ymm15 \n\t" \
                                    : \
                                    : \
                                    "m" (c) \
                                    : \
                                    "xmm15")

#define _mul_cst() \
__asm__ __volatile__ ("vmulps %%ymm15, %%ymm0, %%ymm0 \n\t" \
                       "vmulps %%ymm15, %%ymm1, %%ymm1 \n\t" \
                       "vmulps %%ymm15, %%ymm2, %%ymm2" \
                                    : \
                                    : \
                                    : \
                       "xmm0", "xmm1", "xmm2")
```

HIM **HELMHOLTZ** Helmholtz-Institut Mainz

# Assembler

- Goals:
  - Learn how to code Assembler
  - Write a simple programm squaring a vector

- Notable Omissions:
  - Will not discuss in detail pipelining or latency hiding

# Assembler

- Actually no need to write assembler to use vector units

  Intrinsics:

  – Look like ordinary function calls working on vector memory, e.g. __m256

  – We do not have to deal with registers

  – Could lead to spilling (or other things)

- We'll do the hard thing

# Assembler

- Example:
  We want to do

- Everything double

- Vectorize:
  We need to do

- In AVX can load 4 doubles

```
for(i=0;i<num;i++)
{
  c[i]=x[i]+y[i];
}
```

No Intraloop Dependency!

**SIMD execution**

How?

**HIM HELMHOLTZ** Helmholtz-Institut Mainz

# Assembler

- ## Have different options (NASM, etc.)

```
for(i=0;i<num;i++)
{
    c[i]=x[i]+y[i];
}
```

Do not delete

```
__asm__  __volatile__(\
            "vmovapd %1, %%ymm0\n"\
            "vmovapd %2, %%ymm1\n"\
            "vaddpd %%ymm0,%%ymm1,%%ymm0\n"\
            "vmovapd %%ymm0,%0\n"\
            :\
            "=m" (c[i*4])\
            :\
            "m" (x[i*4]),\
            "m" (y[i*4])\
            :\
            "ymm0", "ymm1");
```



Store operators

Load operators

Clobber

- Need 2 AVX Registers YMM0 & YMM1
- Need 2 Loads and 1 Store
- Need 1 Add on those registers
- Total 4 instructions to add 4+4 doubles
- Dependency:
  - Loads before add
  - Add before store
  - Store at the end

HIM **HELMHOLTZ** Helmholtz-Institut Mainz

# Assembler

- What we did is basically loop-unrolling

```
for(i=0;i<num;i++)
{
  c[i]=x[i]+y[i];
}
```

```
for(i=0;i<num/4;i++)
{
__asm__ __volatile__(\
            "vmovapd %1, %%ymm0\n"\
            "vmovapd %2, %%ymm1\n"\
            "vaddpd %%ymm0,%%ymm1,%%ymm0\n"\
            "vmovapd %%ymm0,%0\n"\
            :\
            "=m" (c[i*4])\
            :\
            "m" (x[i*4]),\
            "m" (y[i*4])\
            :\
            "ymm0", "ymm1");
}
```

- Compiler could do that in simple cases

- OpenMP Pragmas to tell compilers:
  - No Dependency:
    #pragma omp simd

- Need a lot of flags to get the above

HIM **HELMHOLTZ**
Helmholtz-Institut Mainz

# Assembler

- gcc ex2.c   -O3   -march=skylake-avx512   -S

```
.L5:
        vmovupd (%rdi,%rax), %ymm1
        vaddpd  (%rsi,%rax), %ymm1, %ymm0
        vmovupd %ymm0, (%rdx,%rax)
        addq    $32, %rax
        cmpq    %rax, %rcx
        jne     .L5
        movl    %r8d, %ecx
        andl    $-4, %ecx
        movl    %ecx, %eax
        cmpl    %ecx, %r8d
        je      .L23
        movl    %r8d, %r9d
        subl    %ecx, %r9d
        cmpl    $1, %r9d
        je      .L26
        vzeroupper
```

Not clear why. Frequency dropped @ AVX512
Note vmovupd vs vmovapd!

- ## Not using ZMM registers!

- gcc ex2.c   -O3   -march=skylake-avx512  -mprefer-vector-width=512 -S

```
.L5:
        vmovupd (%rsi,%rax), %zmm1
        vaddpd  (%rdi,%rax), %zmm1, %zmm0
        vmovupd %zmm0, (%rdx,%rax)
        addq    $64, %rax
        cmpq    %rcx, %rax
        jne     .L5
        movl    %r8d, %eax
        andl    $-8, %eax
        movl    %eax, %ecx
        cmpl    %eax, %r8d
        je      .L23
        movl    %r8d, %r9d
        subl    %eax, %r9d
        leal    -1(%r9), %r10d
        cmpl    $2, %r10d
        jbe     .L7
```

HIM **HELMHOLTZ**
Helmholtz-Institut Mainz

# Assembler

- Aligned Mem Access vmov**a**p(sd) Memory must be aligned at cache line boundary

- Avoids cache line splits



- Not a drastic penalty anymore

# Assembler

- Example with Intraloop dependency

```
double square(int num,double *a, double c) {
    int i;
    for(i=0;i<num;i++)
    {
        c=c+a[i];
    }
    return 0;
}
```

- Every iteration needs result of the one before

- Add accumulators

```
double square(int num,double *a, double c) {
    int i;
    double c0,c1,c2,c3;
    for(i=0;i<num/4;i++)
    {
        c0=c0+a[i*4];
        c1=c1+a[i*4+1];
        c2=c2+a[i*4+2];
        c3=c3+a[i*4+3];
    }
    return c0+c1+c2+c3;
}
```

- 4-Way vectorization final horizontal add

HIM HELMHOLTZ
Helmholtz-Institut Mainz

# Real World Example

- Lattice QCD Add 2 SU(3) Vectors

```
#define _vector_add(r,s1,s2) \
   (r).c1.re=(s1).c1.re+(s2).c1.re; \
   (r).c1.im=(s1).c1.im+(s2).c1.im; \
   (r).c2.re=(s1).c2.re+(s2).c2.re; \
   (r).c2.im=(s1).c2.im+(s2).c2.im; \
   (r).c3.re=(s1).c3.re+(s2).c3.re; \
   (r).c3.im=(s1).c3.im+(s2).c3.im
```

HIM    **HELMHOLTZ**
       Helmholtz-Institut Mainz

# Real World Example

- Lattice QCD Add 2 SU(3) Vectors Intrinsics

```
vector4double v1,v2;
#define _vector_add_qpx(r,s1,s2) \
    v1=vec_lda(0,&((s1).c1.re)); \
    v2=vec_lda(0,&((s2).c1.re)); \
    vec_sta(vec_add(v1,v2),0,&((r).c1.re)); \
    v1=vec_ld2a(0,&((s1).c3.re)); \
    v2=vec_ld2a(0,&((s2).c3.re)); \
    vec_st2a(vec_add(v1,v2),0,&((r).c3.re));
```

HIM **HELMHOLTZ**
Helmholtz-Institut Mainz

# Real World Example

- Profile to identify HotSpots

# Real World Example

- ## Can get complicated

```
static vector4double perm1={2.000000,2.250000,3.000000,3.250000};
static vector4double perm2={2.500000,2.750000,3.500000,3.750000};
.
.
.
v1=vec_ld2a(0,&((*m).u[6]));
v2=vec_lda(0,&((*m).u[8]));
v3=vec_lda(0,&((*m).u[12]));
v4=vec_lda(0,&((*m).u[16]));
v5=vec_lda(0,&((*m).u[20]));
/* First 4 Components */

v100=vec_perm(v10,v10,perm0011);           /*(u0,u0,u1,u1) */
s10=vec_sldw(s1,s2,2);                     /*s1.re, s1.im, s2.re, s2.im*/
s11=vec_sldw(s2,s1,2);                     /*s2.re, s2.im, s1.re, s1.im*/
v12=vec_perm(v2,v4,perm1);                 /*(u8,u9,u16,u17)*/
v13=vec_perm(v2,v4,perm2);                 /*(u10,u11,u18,u19)*/
v14=vec_perm(v3,v5,perm1);                 /*(u12,u13,u20,u21)*/
v15=vec_perm(v3,v5,perm2);                 /*(u14,u15,u22,u23)*/
r10=vec_xmul(v100,s10);
r11=vec_xxnpmadd(s11,vec_mul((vector4double)(1,1,1,-1),v1),vec_xmul(v1,s11));
r12=vec_xxnpmadd(v12,s3,vec_xmul(s3,v12));
r13=vec_xxnpmadd(v13,s4,vec_xmul(s4,v13));
r14=vec_xxnpmadd(v14,s5,vec_xmul(s5,v14));
r15=vec_xxnpmadd(v15,s6,vec_xmul(s6,v15));

r100=vec_add(r10,r11);
r101=vec_add(r12,r13);
r102=vec_add(r14,r15);
r110=vec_add(r100,r101);
r111=vec_add(r110,r102);
vec_sta(r111,0,&((*r).c1.c1.re));
```

Further Improvements:
Mix Loads and (F)MADs!
Order data dependency according to
Instruction latency. (See Bagel maybe?)
This can only be done in assembly, I guess
Not Sure how the compiler orders regs?

Bad Version not using Multiply-Add!
This was version 3. Changed in v4.

**HIM** **HELMHOLTZ**
Helmholtz-Institut Mainz

# Improvements - Version 0.4

- After 4 weeks of tuning
- Using timing programs from devel/dirac subdirectory in DD-HMC
- Local Lattice 16x16x16x16
- Processorgrid 4x2x2x2
- Results in Mflops/process
- Speedup for QPXD in brackets

| Routine | -O3 -qstrict | -O3 –qstrict -DQPXD | -O2 |
|---|---|---|---|
| Qhat | 688 (2.08x) | 1428 | 640 (2.23x) |
| Qhat_blk | 865 (2.15x) | 1859 | 691 (2.69x) |
| Qnohat | 562 (2.38x) | 1335 | 499 (2.68x) |
| Qhat_dble | 436 (1.78x) | 774 | 438 (1.77x) |
| Qhat_blk_dble | 576 (1.89x) | 1089 | 566 (1.92x) |
| Qnohat_dble | 403 (1.71x) | 690 | 392 (1.76x) |

- Using Prefetch by data cache block touch to preload __dcbt
- Set pipeline depth should improve DP code?