

Programming Techniques for Supercomputers Tutorial

Erlangen National High Performance Computing Center

Department of Computer Science

FAU Erlangen-Nürnberg

Sommersemester 2025



Assignment 3 – Task 1

Max in-core performance a) AVX512

```
for(int i=0; i<N; ++i)
  for(int j=1; j<M-1; ++j)
    a[i][j] = 0.2*(b[i][j] +
      b[i][j-1] + b[i][j+1]);
```

- Double precision arrays, 1.9 [GHz]

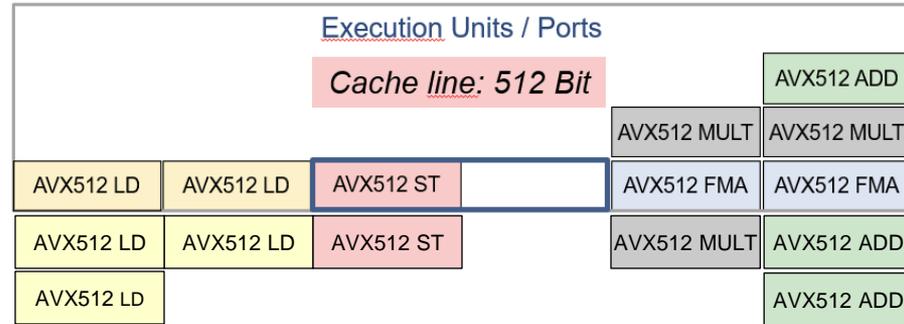
$$P_{max} = 12 \left[\frac{flops}{cy} \right] * 1.9 \left[\frac{Gcy}{s} \right] = 22.8 \left[\frac{GF}{s} \right]$$

$$P_{peak} = n_{super}^{FP} * n_{FMA} * n_{SIMD} * f =$$

$$\underbrace{2 \left[\frac{instr}{cy} \right] * 2 * 8 \left[\frac{flops}{instr} \right]}_{32 \left[\frac{flops}{cy} \right]} * 1.9 \left[\frac{Gcy}{s} \right] = 60.8 \left[\frac{GF}{s} \right]$$

$$\Rightarrow P_{max} = \frac{3}{8} P_{peak}$$

$$7[instr]/2[cy]$$



Assignment 3 – Task 1

Max in-core performance a) AVX2

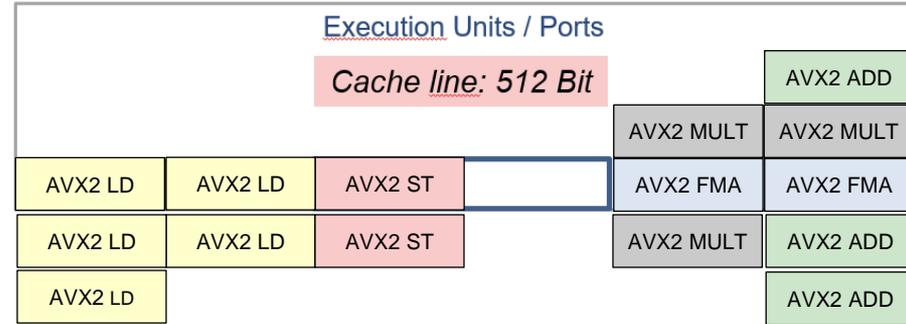
```
for(int i=0; i<N; ++i)
  for(int j=1; j<M-1; ++j)
    a[i][j] = 0.2*(b[i][j] +
      b[i][j-1] + b[i][j+1]);
```

- Double precision arrays, 1.9 [GHz]

$$P_{max} = 6 \left[\frac{flops}{cy} \right] * 1.9 \left[\frac{Gcy}{s} \right] = 11.4 \left[\frac{GF}{s} \right]$$

$$\Rightarrow P_{max} = \frac{3}{16} P_{peak}$$

$$7[instr]/2[cy]$$



Assignment 3 – Task 1

Max in-core performance a) Scalar

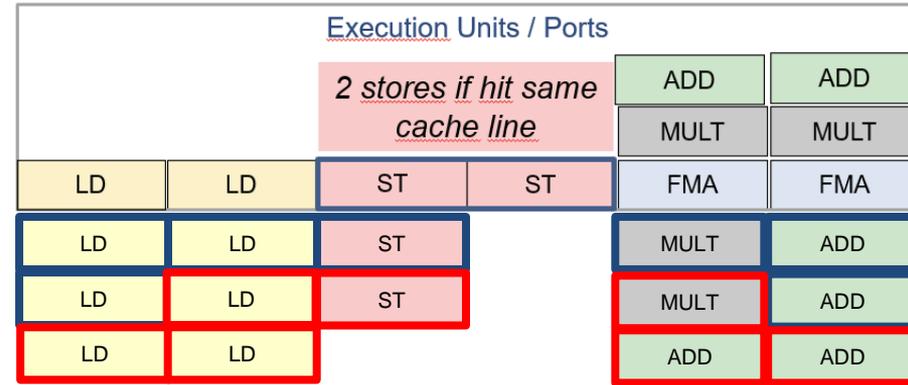
```
for(int i=0; i<N; ++i)
  for(int j=1; j<M-1; ++j)
    a[i][j] = 0.2*(b[i][j] +
      b[i][j-1] + b[i][j+1]);
```

- Double precision arrays, 1.9 [GHz]

$$P_{max} = 2 \left[\frac{flops}{cy} \right] * 1.9 \left[\frac{Gcy}{s} \right] = 3.8 \left[\frac{GF}{s} \right]$$

$$\Rightarrow P_{max} = \frac{1}{16} P_{peak}$$

$$14[instr]/3[cy]$$



Assignment 3 – Task 1

Max in-core performance b) Code Balance

```
for(int i=0; i<N; ++i)
  for(int j=1; j<M-1; ++j)
    a[i][j] = 0.2*(b[i][j] +
                  b[i][j-1] + b[i][j+1]);
```

RHS: 1 RD
LHS: 1 WR (+ 1RD for W/A)

- Double precision arrays
- Nontemporal stores or cache line claim cannot be used
- Assume N and M are “large” (arrays will not fit into cache)
- In each iteration, $b[i][j+1]$ is read from memory
 - This element will be used from cache in the next two iterations
 - **No main memory transfers are necessary** for $b[i][j]$ and $b[i][j-1]$

$$\Rightarrow B_c = \frac{24[B/iter]}{3[flops/iter]} = 8[B/F]$$

Assignment 3 – Task 2

Latency and bandwidth

$$b_s = 250 \text{ [GB/s]}$$

$$T_l = 130 \text{ [ns]}$$

a) $N = 4096 \text{ [B]}$

$$\text{Hockney Model: } T_{trans} = T_l + \frac{N}{b} = 130 \times 10^{-9} \text{ [s]} + \frac{4096 \text{ [B]}}{250 \times 10^9 \text{ [B/s]}} = \mathbf{146 \text{ [ns]}}$$

$$B_{eff} = \frac{N}{T_{trans}} = 4096 \text{ [B]} / 146 \text{ [ns]} = \mathbf{28 \text{ [GB/s]}}$$

b) $B_{eff} = b_s/2$ with $B_{eff} = \frac{N}{B} \rightarrow N_{1/2} = T_l b_s$

c) $P = 1 + \frac{T_l}{N/b} = 1 + \frac{130 \text{ [ns]}}{64 \text{ [B]} / 250 \text{ [GB/s]}} = 508.8125 \rightarrow \mathbf{509}$ prefetches needed

$\approx \mathbf{32.6 \text{ kB}}$ in-flight

Who is the major contributor to runtime?

Assignment 3 – Task 3

Strided access

```
for(i=0; i<N; i+=M)  
  a[i] = s * a[i];
```

Only every Mth
element touched

Cache thrashing!

