# Progamming the OpenMP API
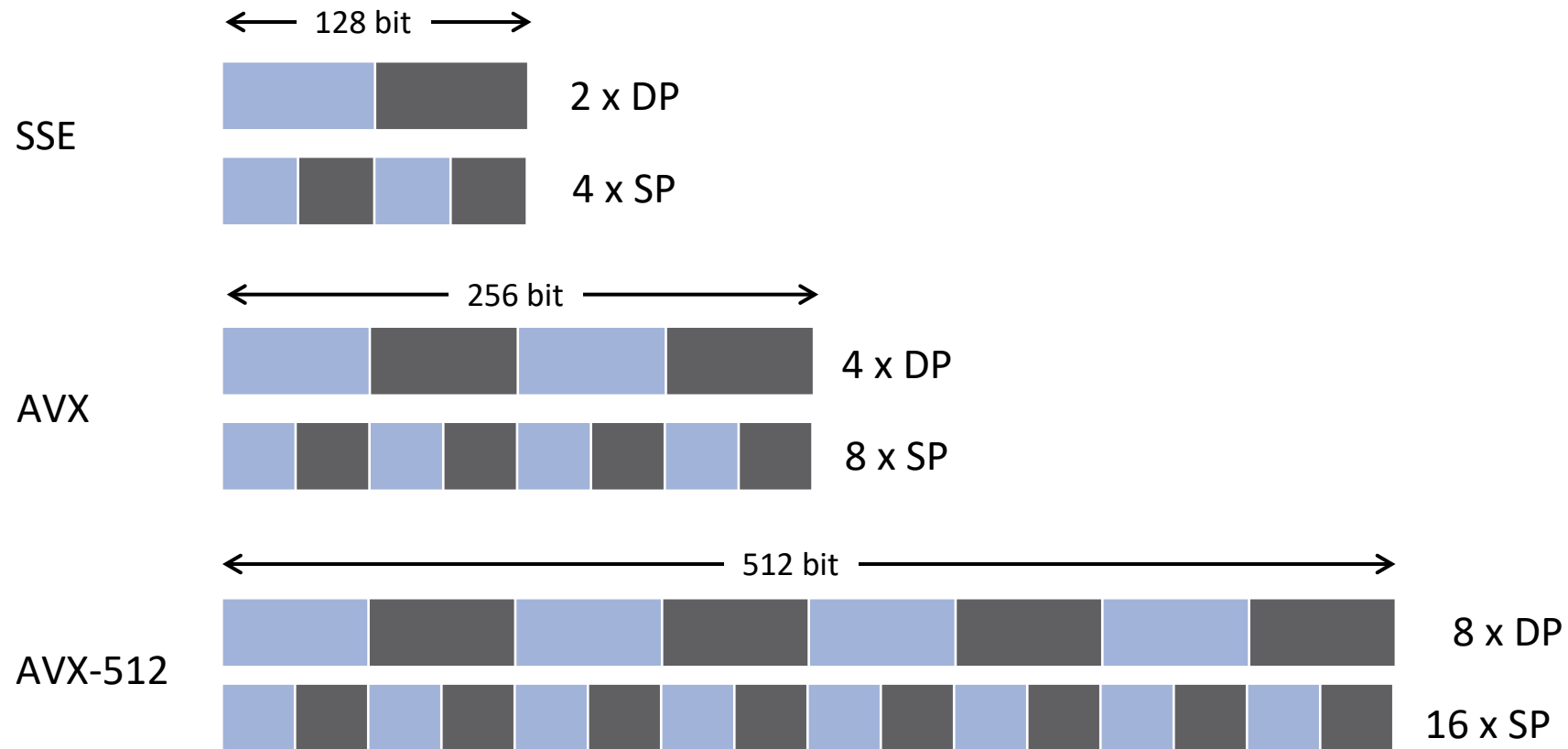
## *Vectorization*

# Topics

- Exploiting SIMD parallelism with OpenMP
- Using SIMD directives with loops
- Creating SIMD functions
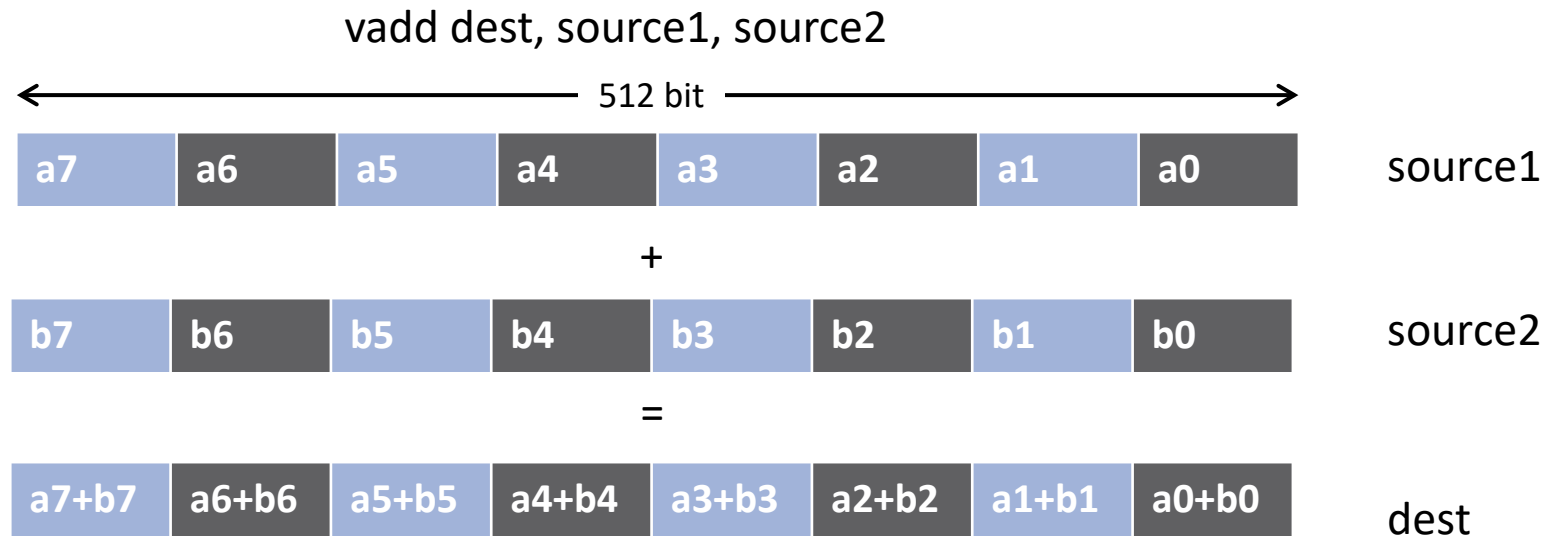
# SIMD on x86 Architectures

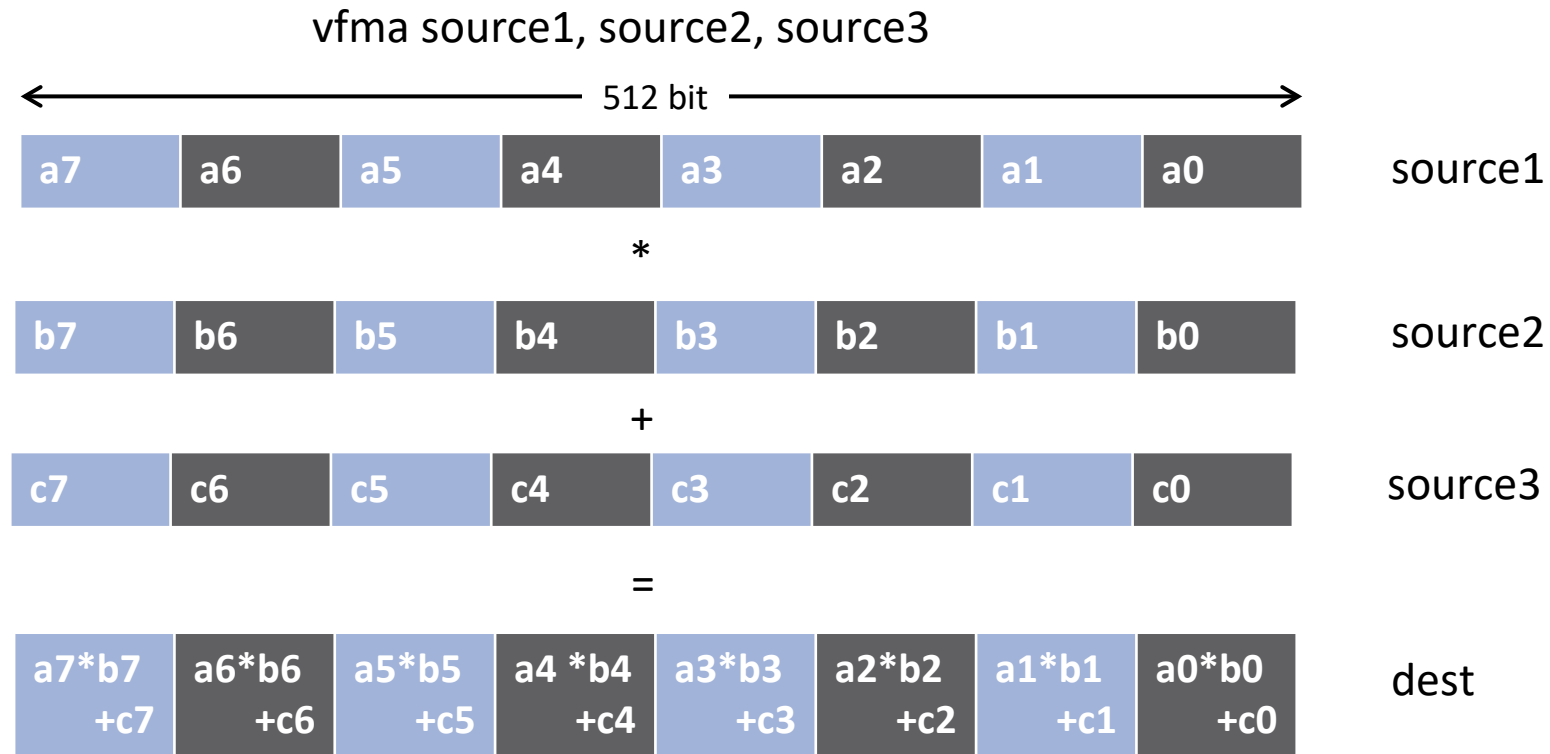- Width of SIMD registers has been growing in the past:



SSE
128 bit
2 x DP
4 x SP

AVX
256 bit
4 x DP
8 x SP

AVX-512
512 bit
8 x DP
16 x SP

**Programming the OpenMP API**
**SIMD**

# More Powerful SIMD Units

■ SIMD instructions become more powerful

vadd dest, source1, source2

| a7 | a6 | a5 | a4 | a3 | a2 | a1 | a0 | source1 |

+

| b7 | b6 | b5 | b4 | b3 | b2 | b1 | b0 | source2 |

=

| a7+b7 | a6+b6 | a5+b5 | a4+b4 | a3+b3 | a2+b2 | a1+b1 | a0+b0 | dest |

512 bit

# More Powerful SIMD Units

- SIMD instructions become more powerful

vfma source1, source2, source3

← 512 bit →

| a7 | a6 | a5 | a4 | a3 | a2 | a1 | a0 | source1 |

\*

| b7 | b6 | b5 | b4 | b3 | b2 | b1 | b0 | source2 |

+

| c7 | c6 | c5 | c4 | c3 | c2 | c1 | c0 | source3 |

=

| a7*b7 +c7 | a6*b6 +c6 | a5*b5 +c5 | a4 *b4 +c4 | a3*b3 +c3 | a2*b2 +c2 | a1*b1 +c1 | a0*b0 +c0 | dest |

# More Powerful SIMD Units

- SIMD instructions become more powerful

vadd dest{k1}, source2, source3

← 512 bit →

| a7 | a6 | a5 | a4 | a3 | a2 | a1 | a0 | source1 |

+

| b7 | b6 | b5 | b4 | b3 | b2 | b1 | b0 | source2 |

| 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | mask |

=

| a7+b7 | d6 | a5+b5 | d4 | d3 | a2+b2 | d1 | a0+b0 | dest |

# More Powerful SIMD Units

■ SIMD instructions become more powerful

# Auto-vectorization

- Compilers offer auto-vectorization as an optimization pass
    - → Usually, part of the general loop optimization passes
    - → Code analysis detects code properties that inhibit SIMD vectorization **?**
    - → Heuristics determine if SIMD execution might be beneficial
    - → If all goes well, the compiler will generate SIMD instructions

- Example: clang/LLVM      GCC      Intel Compiler
    - → -fvectorize      -ftree-vectorize      -vec (enabled w/ -O2)
    - → -Rpass=loop-.\*      -ftree-loop-vectorize      -qopt-report=vec
    - → -mprefer-vector-width=*<width>*    -fopt-info-vec-all

# Why Auto-vectorizers Fail

- Data dependencies

- Other potential reasons
  - → Alignment
  - → Function calls in loop block
  - → Complex control flow / conditional branches
  - → Loop not "countable"
    - → e.g., upper bound not a runtime constant
  - → Mixed data types
  - → Non-unit stride between elements
  - → Loop body too complex (register pressure)
  - → Vectorization seems inefficient

- Many more … but less likely to occur

**Programming the OpenMP API**
**SIMD**

# Data Dependencies

- Suppose two statements S1 and S2
- S2 depends on S1, iff S1 must execute before S2
  - → Control-flow dependence
  - → Data dependence
  - → Dependencies can be carried over between loop iterations
- Important flavors of data dependencies
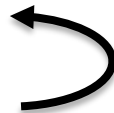
FLOW

```
s1: a = 40

    b = 21
s2: c = a + 2
```

ANTI

```
    b = 40

s1: a = b + 1
s2: b = 21
```

# Loop-Carried Dependencies

■ Dependencies may occur across loop iterations
- →Loop-carried dependency

■ The following code contains such a dependency:

```
void lcd_ex(float* a, float* b, size_t n, float c1, float c2)
{
    size_t i;
    for (i = 0; i < n; i++) {
        a[i] = c1 * a[i + 17] + c2 * b[i];
    }
}
```

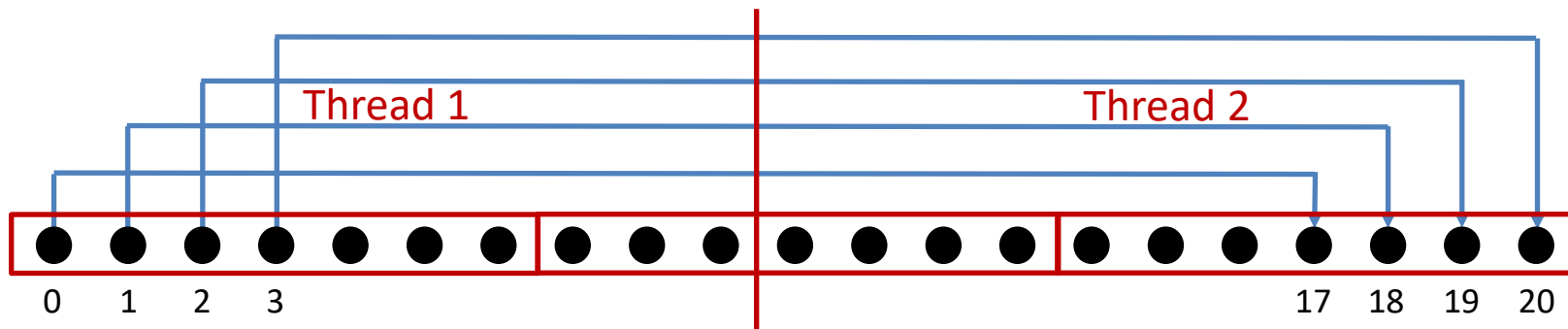Loop-carried dependency for a[i] and a[i+17]; distance is 17.

■ Some iterations of the loop have to complete before the next iteration can run
- →Simple trick: Can you reverse the loop w/o getting wrong results?

# Loop-carried Dependencies

■ Can we parallelize or vectorize the loop?

```
void lcd_ex(float* a, float* b, size_t n, float c1, float c2) {
    for (int i = 0; i < n; i++) {
        a[i] = c1 * a[i + 17] + c2 * b[i];
}    }
```



→ Parallelization: no
(except for very specific loop schedules)

→ Vectorization: yes
(iff vector length is shorter than any distance of any dependency)

# In a Time Before OpenMP 4.0

■ Support required vendor-specific extensions
→ Programming models (e.g., Intel® Cilk Plus)
→ Compiler pragmas (e.g., `#pragma vector`)
→ Low-level constructs (e.g., `_mm_add_pd()`)

```
#pragma omp parallel for
#pragma vector always
#pragma ivdep
for (int i = 0; i < N; i++) {
    a[i] = b[i] + ...;
}
```

You need to trust your compiler to do the "right" thing.

# SIMD Loop Construct

- Vectorize a loop nest
  - →Cut loop into chunks that fit a SIMD vector register
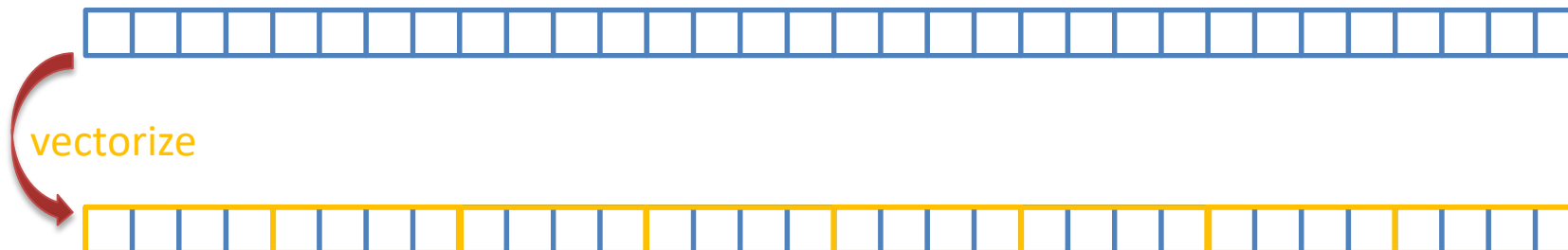  - →No parallelization of the loop body

- Syntax (C/C++)
  ```
  #pragma omp simd [clause[[,] clause],…]
  for-loops
  ```

- Syntax (Fortran)
  ```
  !$omp simd [clause[[,] clause],…]
  do-loops
  [!$omp end simd]
  ```

# Example

```
float sprod(float *a, float *b, int n) {
    float sum = 0.0f;
#pragma omp simd reduction(+:sum)
    for (int k=0; k<n; k++)
        sum += a[k] * b[k];
    return sum;
}
```



vectorize

**Programming the OpenMP API**
**SIMD**

# Data Sharing Clauses

- `private(`*`var-list`*`):`
  Uninitialized vectors for variables in *var-list*

  x: `42` ⟶ `? ? ? ?`

- `firstprivate(`*`var-list`*`):`
  Initialized vectors for variables in *var-list*

  x: `42` ⟶ `42 42 42 42`

- `reduction(`*`op`*`:`*`var-list`*`):`
  Create private variables for *var-list* and apply reduction operator *op* at the end of the construct

  `12 5 8 17` ⟶ x: `42`

**Programming the OpenMP API**
**SIMD**

# SIMD Loop Clauses

- `safelen (`*`length`*`)`
  - → Maximum number of iterations that can run concurrently without breaking a dependence
  - → In practice, maximum vector length
- `linear (`*`list`*`[:`*`linear-step`*`])`
  - → The variable's value is in relationship with the iteration number
    - → $x_i = x_{orig} + i * \text{linear-step}$
- `aligned (`*`list`*`[:`*`alignment`*`])`
  - → Specifies that the list items have a given alignment
  - → Default is alignment for the architecture
- `collapse (`*`n`*`)`

# SIMD Worksharing Construct

- **Parallelize and vectorize a loop nest**
  - → Distribute a loop's iteration space across a thread team
  - → Subdivide loop chunks to fit a SIMD vector register

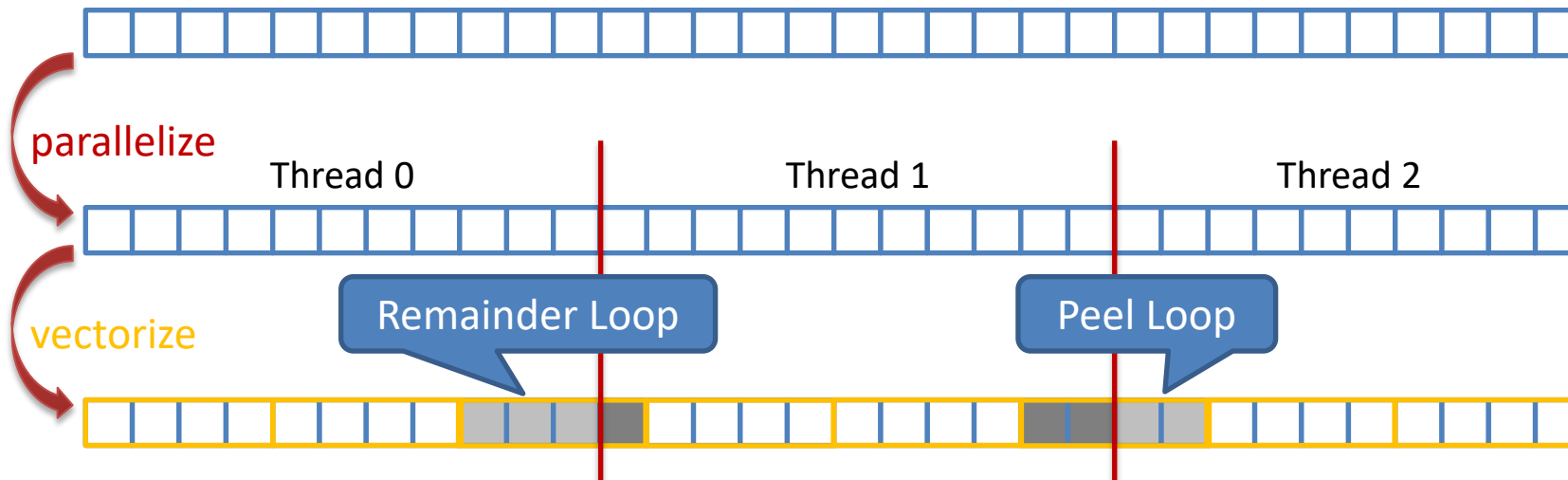- **Syntax (C/C++)**
  ```
  #pragma omp for simd [clause[[,] clause],…]
  for-loops
  ```

- **Syntax (Fortran)**
  ```
  !$omp do simd [clause[[,] clause],…]
  do-loops
  [!$omp end do simd [nowait]]
  ```

```
float sprod(float *a, float *b, int n) {
  float sum = 0.0f;
#pragma omp for simd reduction(+:sum)
  for (int k=0; k<n; k++)
    sum += a[k] * b[k];
  return sum;
}
```

# Be Careful What You Wish For...

```
float sprod(float *a, float *b, int n) {
  float sum = 0.0f;
#pragma omp for simd reduction(+:sum) \
                        schedule(static, 5)

  for (int k=0; k<n; k++)
    sum += a[k] * b[k];
  return sum;
}
```

- You should choose chunk sizes that are multiples of the SIMD length
  - → Remainder loops are not triggered
  - → Likely better performance
- In the above example …
  - → and AVX2, the code will only execute the remainder loop!
  - → and SSE, the code will have one iteration in the SIMD loop plus one in the remainder loop!

# OpenMP 4.5 Simplifies  SIMD Chunks

```c
float sprod(float *a, float *b, int n) {
  float sum = 0.0f;
#pragma omp for simd reduction(+:sum) \
                       schedule(simd: static, 5)
  for (int k=0; k<n; k++)
    sum += a[k] * b[k];
  return sum;
}
```

- Chooses chunk sizes that are multiples of the SIMD length
  - First and last chunk may be slightly different to fix alignment and to handle loops that are not exact multiples of SIMD width
  - Remainder loops are not triggered
  - Likely better performance

# SIMD Function Vectorization

```c
float min(float a, float b) {
    return a < b ? a : b;
}


float distsq(float x, float y) {
    return (x - y) * (x - y);
}

void example() {
#pragma omp parallel for simd
    for (i=0; i<N; i++) {
        d[i] = min(distsq(a[i], b[i]), c[i]);
}    }
```

**Programming the OpenMP API**
**SIMD**

# SIMD Function Vectorization

- Declare one or more functions to be compiled for calls from a SIMD-parallel loop

- Syntax (C/C++):

```
#pragma omp declare simd [clause[[,] clause],…]
[#pragma omp declare simd [clause[[,] clause],…]]
[…]
function-definition-or-declaration
```

- Syntax (Fortran):

```
!$omp declare simd (proc-name-list)
```

# SIMD Function Vectorization

```
#pragma omp declare simd
float min(float a, float b) {
    return a < b ? a : b;
}
```

```
_ZGVZN16vv_min(%zmm0, %zmm1):
    vminps %zmm1, %zmm0, %zmm0
    ret
```

```
#pragma omp declare simd
float distsq(float x, float y)
    return (x - y) * (x - y);
}
```

```
_ZGVZN16vv_distsq(%zmm0, %zmm1):
    vsubps %zmm0, %zmm1, %zmm2
    vmulps %zmm2, %zmm2, %zmm0
    ret
```

```
void example() {
#pragma omp parallel for simd
    for (i=0; i<N; i++) {
        d[i] = min(distsq(a[i], b[i]), c[i]);
    }   }
```

```
vmovups (%r14,%r12,4), %zmm0
vmovups (%r13,%r12,4), %zmm1
call _ZGVZN16vv_distsq
vmovups (%rbx,%r12,4), %zmm1
call _ZGVZN16vv_min
```

**Programming the OpenMP API**
**SIMD**

# SIMD Function Vectorization

- `simdlen (length)`
  - → generate function to support a given vector length
- `uniform (argument-list)`
  - → argument has a constant value between the iterations of a given loop
- `inbranch`
  - → function always called from inside an if statement
- `notinbranch`
  - → function never called from inside an if statement
- `linear (argument-list[:linear-step])`
- `aligned (argument-list[:alignment])`

# inbranch & notinbranch

```
#pragma omp declare simd inbranch
float do_stuff(float x) {
    /* do something */
    return x * 2.0;
}

void example() {
#pragma omp simd
    for (int i = 0; i < N; i++)
        if (a[i] < 0.0)
            b[i] = do_stuff(a[i]);
}
```
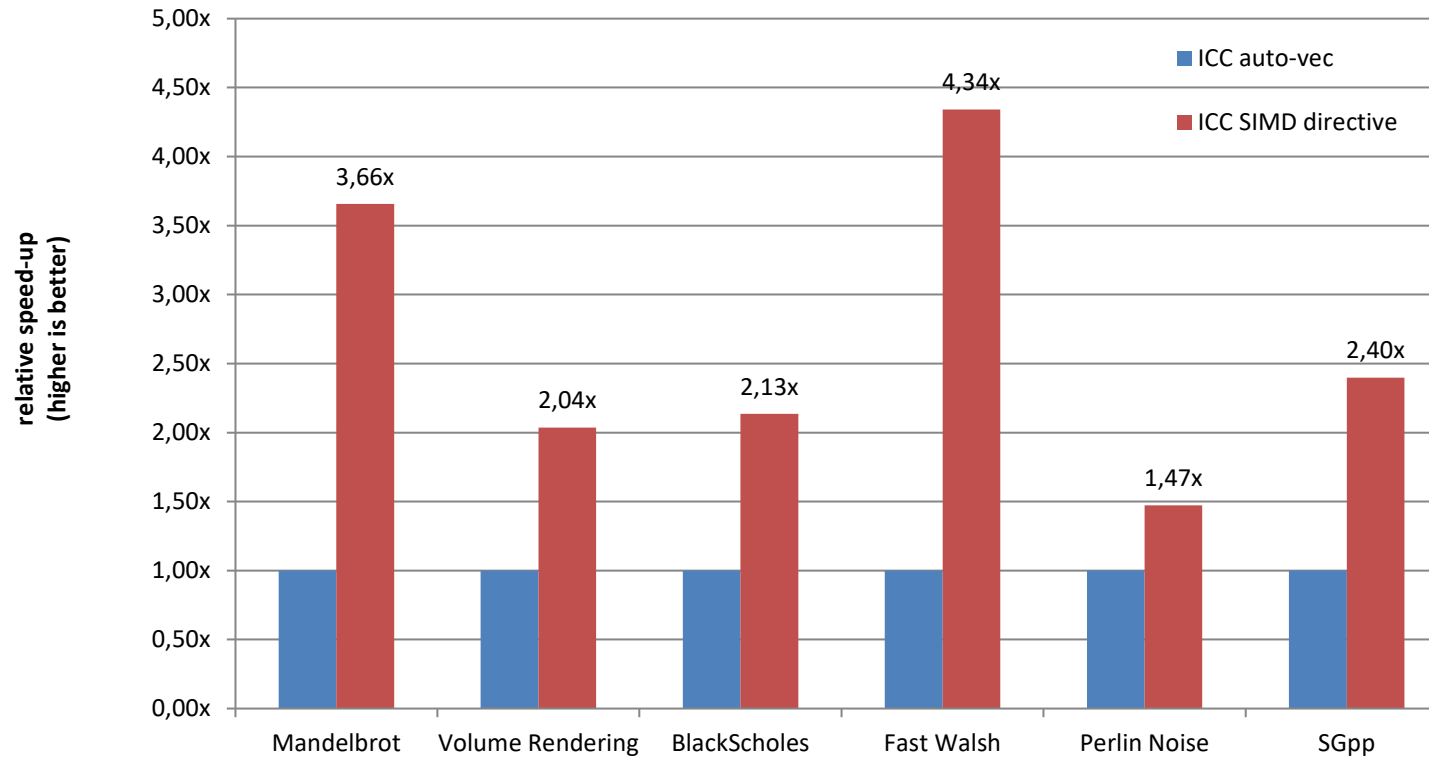
```
vec8 do_stuff_v(vec8 x, mask m) {
    /* do something */
    vmulpd x{m}, 2.0, tmp
    return tmp;
}
```

```
for (int i = 0; i < N; i+=8) {
    vcmp_lt &a[i], 0.0, mask
    b[i] = do_stuff_v(&a[i], mask);
}
```

**Programming the OpenMP API**
**SIMD**

# SIMD Constructs & Performance



M.Klemm, A.Duran, X.Tian, H.Saito, D.Caballero, and X.Martorell. Extending OpenMP with Vector Constructs for Modern Multicore SIMD Architectures. In Proc. of the Intl. Workshop on OpenMP, pages 59-72, Rome, Italy, June 2012. LNCS 7312.