

Accurate and Efficiently Vectorized Sums and Dot Products in Julia

Julia Paris Meetup

October 17th, 2019

François Févotte – TriScale innov

Chris Elrod – Baylor University



Summation & dot product

Summation & dot product

Algorithm 1 Naive summation

$\sigma \leftarrow 0$

for $i = 1 \dots n$ **do**

$\sigma \leftarrow \sigma + x_i$

end for

return σ

Algorithm 2 Naive dot product

$\sigma \leftarrow 0$

for $i = 1 \dots n$ **do**

$\sigma \leftarrow \sigma + x_i y_i$

end for

return σ

Accuracy

Accuracy of summation algorithms

Naive summation

Algorithm 3 Naive summation

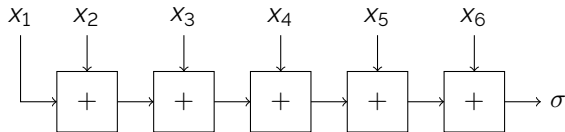
$\sigma \leftarrow 0$

for $i = 1 \dots n$ **do**

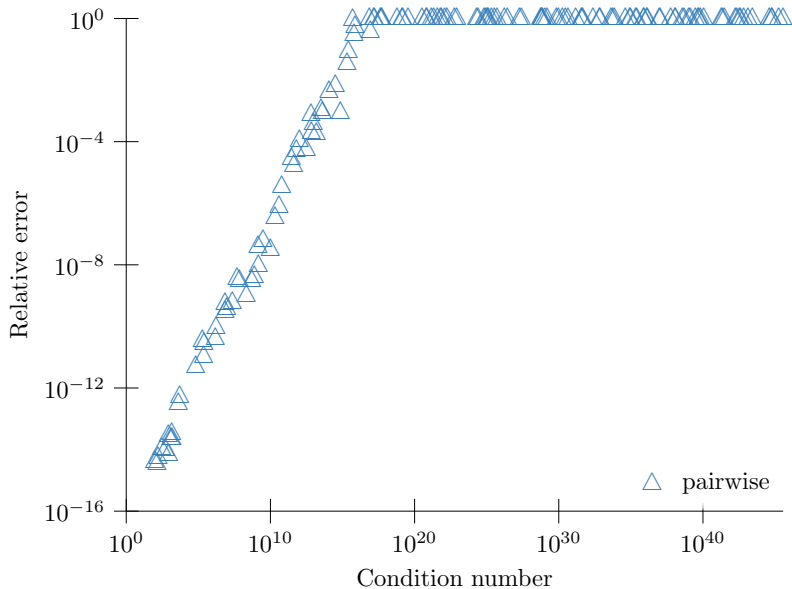
$\sigma \leftarrow \sigma + X_i$

end for

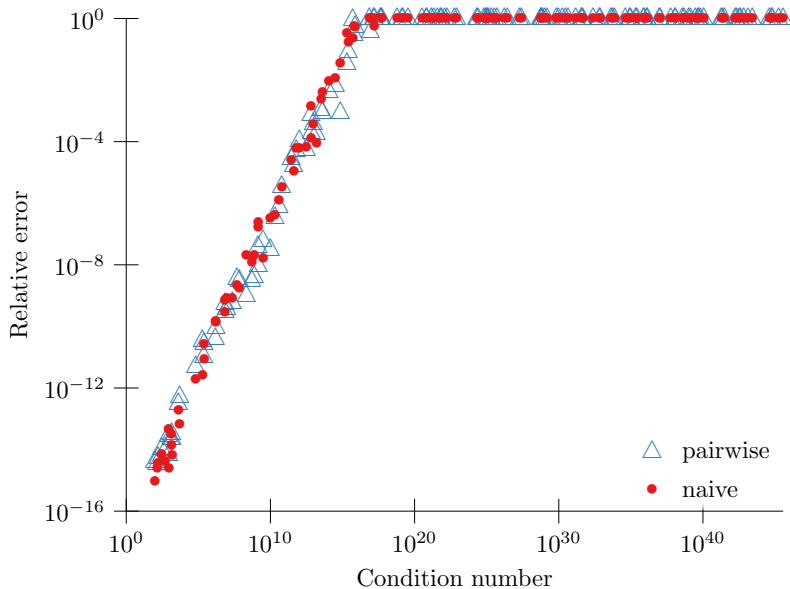
return σ



Error of summation algorithms



Error of summation algorithms



Accuracy of summation algorithms

Compensated algorithms

- ▶ For some classical algorithms, there exist “compensated” variants :
 - ▶ summation
 - ▶ dot product
 - ▶ polynomial evaluation
 - ▶ ...
- ▶ these “compensated algorithms” are based on Error Free Transforms (EFTs) :

$$x \underset{\text{EFT}}{\circ} y = (r, \delta) \quad \left(\forall \circ \in \{+, -, \times\} \right)$$

such that

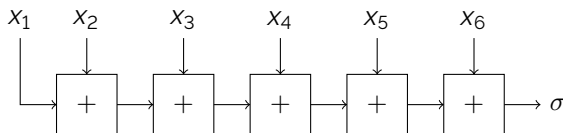
$$\begin{aligned} r &= \lfloor x \circ y \rfloor \\ r + \delta &= x \circ y \end{aligned}$$

Accuracy of summation algorithms

Compensated summation

Algorithm 4 Naive summation

```
 $\sigma \leftarrow 0$   
for  $i = 1 \dots n$  do  
   $\sigma \leftarrow \sigma + X_i$   
end for  
return  $\sigma$ 
```



Accuracy of summation algorithms

Compensated summation

Algorithm 7 Compensated summation

$\sigma \leftarrow 0$

$\pi \leftarrow 0$

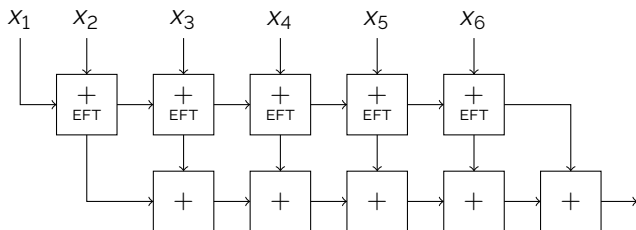
for $i = 1 \dots n$ **do**

$(\sigma, e) \leftarrow \sigma + X_i$
EFT

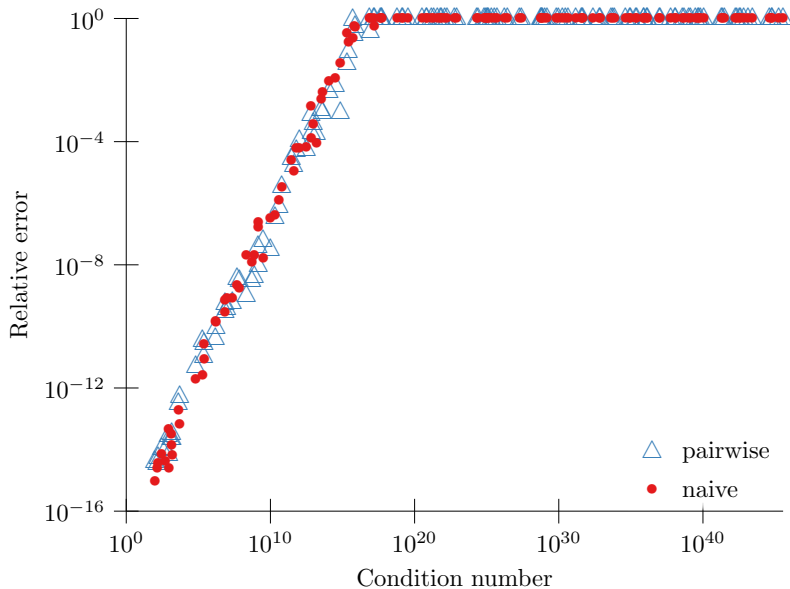
$\pi \leftarrow \pi + e$

end for

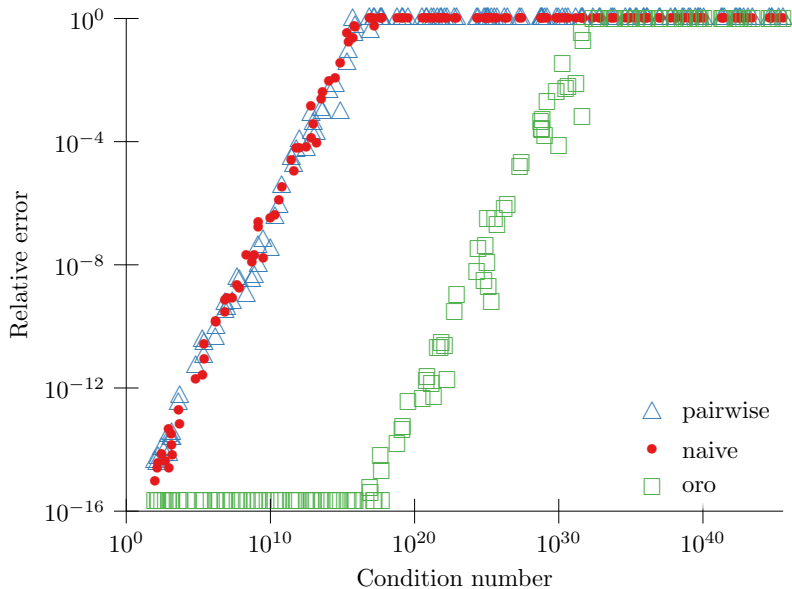
return $\sigma + \pi$



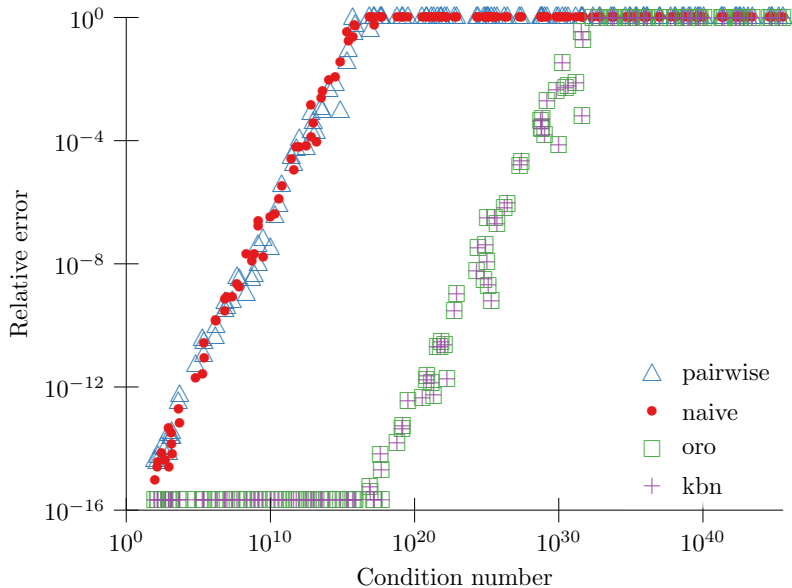
Error of summation algorithms



Error of summation algorithms



Error of summation algorithms

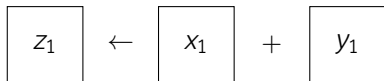


Performance & Vectorization (SIMD)

Performance of summation algorithms

Vectorization (SIMD)

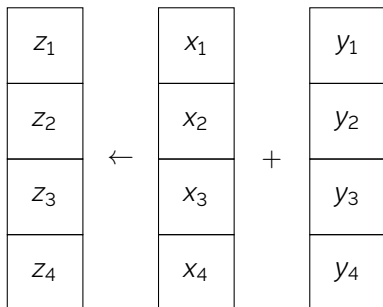
► Scalar operation



Performance of summation algorithms

Vectorization (SIMD)

- ▶ Vector operation (Single Instruction Multiple Data)

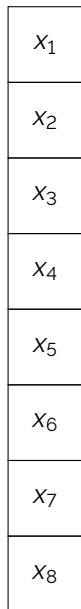
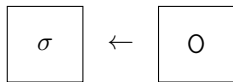


- ▶ Hardware generations :

- ▶ SSE
- ▶ AVX
- ▶ AVX-2
- ▶ AVX-512

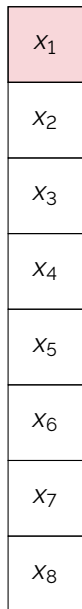
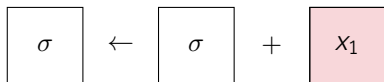
Performance of summation algorithms

Naive summation, scalar version



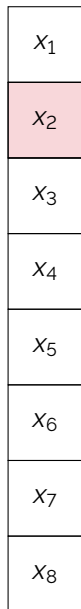
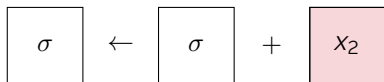
Performance of summation algorithms

Naive summation, scalar version



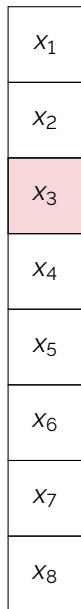
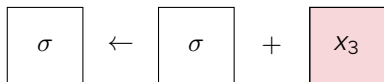
Performance of summation algorithms

Naive summation, scalar version



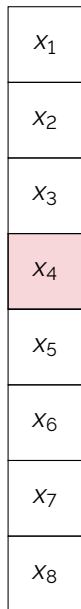
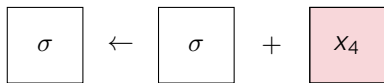
Performance of summation algorithms

Naive summation, scalar version



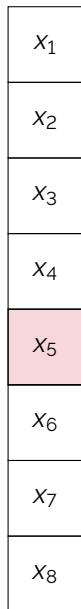
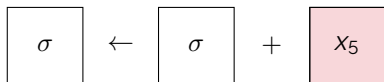
Performance of summation algorithms

Naive summation, scalar version



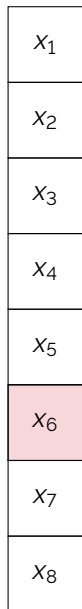
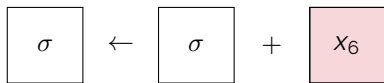
Performance of summation algorithms

Naive summation, scalar version



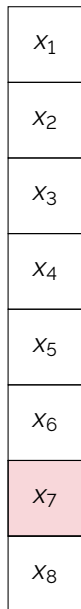
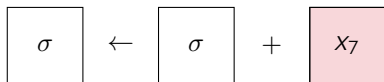
Performance of summation algorithms

Naive summation, scalar version



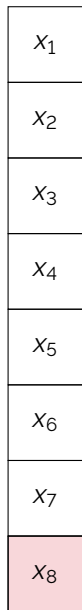
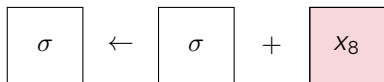
Performance of summation algorithms

Naive summation, scalar version



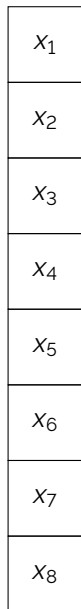
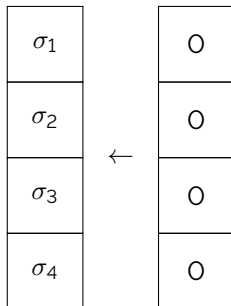
Performance of summation algorithms

Naive summation, scalar version



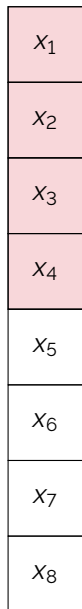
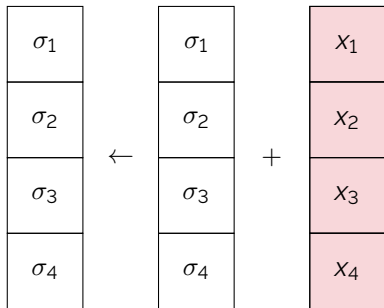
Performance of summation algorithms

Naive summation, vectorized version



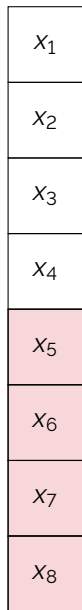
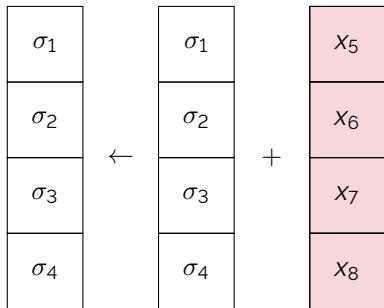
Performance of summation algorithms

Naive summation, vectorized version



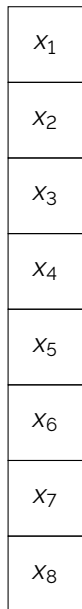
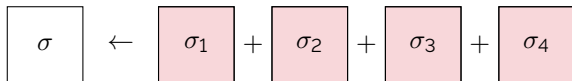
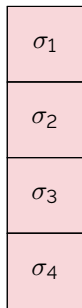
Performance of summation algorithms

Naive summation, vectorized version



Performance of summation algorithms

Naive summation, vectorized version



Performance of summation algorithms

Vectorized summation

Algorithm 8 Naive summation

{Initialization :}

$a \leftarrow 0$

{Loop on vector elements :}

for $i \in 1 : N$ **do**

$a \leftarrow a \oplus x_i$

end for

return a

► Naive

variant	ns/el.	speedup
scalar	1.35	

Performance of summation algorithms

Vectorized summation

Algorithm 12 Vectorized summation

```
1: {Initialization :}  
2: a  $\leftarrow$  0  
  
3: {Loop on full packs :}  
4: for  $j \in 1 : \lfloor \frac{N}{W} \rfloor$  do  
5:    $i \leftarrow W(j - 1) + 1$   
6:   p  $\leftarrow (x_i, x_{i+1}, x_{i+2}, x_{i+3})$   
7:   a  $\leftarrow \mathbf{a} \oplus \mathbf{p}$   
8: end for  
  
9: {Reduction of SIMD accumulator :}  
10:  $a \leftarrow \mathbf{vsum}(\mathbf{a})$   
  
11: {Loop on remaining elements :}  
12: for  $j \in W \lfloor \frac{N}{W} \rfloor + 1 : N$  do  
13:    $a \leftarrow a \oplus x_j$   
14: end for  
  
15: return  $a$ 
```

► Naive

variant	ns/el.	speedup
scalar	1.35	
vector	0.17	$\times 8$

Performance of summation algorithms

Vectorized summation

Algorithm 16 Unrolled vectorized summation

```
1: [Initialization :]
2:  $\mathbf{a}_1 \leftarrow 0$ 
3:  $\mathbf{a}_2 \leftarrow 0$ 

4: [Loop on full packs, unrolled twice :]
5: for  $j \in 1 : \lfloor \frac{N}{2W} \rfloor$  do
6:    $i_1 \leftarrow 2W(j-1) + 1$ 
7:    $\mathbf{p}_1 \leftarrow (X_{i_1}, X_{i_1+1}, X_{i_1+2}, X_{i_1+3})$ 
8:    $\mathbf{a}_1 \leftarrow \mathbf{a}_1 \oplus \mathbf{p}_1$ 
9:    $i_2 \leftarrow 2W(j-1) + W + 1$ 
10:   $\mathbf{p}_2 \leftarrow (X_{i_2}, X_{i_2+1}, X_{i_2+2}, X_{i_2+3})$ 
11:   $\mathbf{a}_2 \leftarrow \mathbf{a}_2 \oplus \mathbf{p}_2$ 
12: end for

13: [Loop on remaining full packs :]
14: for  $j \in 2 \lfloor \frac{N}{2W} \rfloor + 1 : \lfloor \frac{N}{W} \rfloor$  do
15:    $i \leftarrow W(j-1) + 1$ 
16:    $\mathbf{p} \leftarrow (X_i, X_{i+1}, X_{i+2}, X_{i+3})$ 
17:    $\mathbf{a}_1 \leftarrow \mathbf{a}_1 \oplus \mathbf{p}$ 
18: end for

19: [Reduction of SIMD accumulators :]
20:  $\mathbf{a}_1 \leftarrow \mathbf{a}_1 \oplus \mathbf{a}_2$ 
21:  $\alpha \leftarrow \mathbf{vsum}(\mathbf{a}_1)$ 

22: [Loop on remaining elements :]
23: for  $j \in W \lfloor \frac{N}{W} \rfloor + 1 : N$  do
24:    $\alpha \leftarrow \alpha \oplus X_j$ 
25: end for

26: return  $\alpha$ 
```

► Naive

variant	ns/el.	speedup
scalar	1.35	
vector	0.17	$\times 8$

► Unrolled

variant	ns/el.	speedup
scalar	0.35	
vector	0.09	$\times 3.8$

Implementation

AccurateArithmetic.jl

Implementation

Why Julia?

Algorithm

- ▶ Summation
- ▶ Dot product

Compensation

- ▶ Naive
- ▶ Compensated (ORO)
- ▶ Compensated (KBN)

Vectorization

- ▶ Scalar
- ▶ Vectorized
- ▶ Vectorized + unrolled

Julia implementation

Why Julia? Separation of concerns : multiple dispatch, splatting...

```
function sum(x)
    acc = zero(eltypes(x))
    for e in x
        acc += e
    end
    return acc
end
```

Julia implementation

Why Julia? Separation of concerns : multiple dispatch, splatting...

```
sum(x) = sum_(x, NaiveAcc)
```

```
function sum_(x, accType)
    acc = zero(accType)
    for e in x
        add!(acc, e)
    end
    return value(acc)
end
```

Julia implementation

Why Julia? Separation of concerns : multiple dispatch, splatting...

```
sum(x, y) = acc_((x,), NaiveSum)
dot(x, y) = acc_((x,y), NaiveDot)
```

```
function acc_(operands, accType)
    acc = zero(accType)
    for e in zip(operands...)
        add!(acc, e)
    end
    return value(acc)
end
```

Implementation

Why Julia?

Algorithm

- ▶ Summation
- ▶ Dot product

Compensation

- ▶ Naive
- ▶ Compensated (ORO)
- ▶ Compensated (KBN)

Vectorization

- ▶ Scalar
- ▶ Vectorized
- ▶ Vectorized + unrolled

```
sum1(x) = acc_((x,), NaiveSum)
```

Implementation

Why Julia?

Algorithm

- ▶ Summation
- ▶ Dot product

Compensation

- ▶ Naive
- ▶ Compensated (ORO)
- ▶ Compensated (KBN)

Vectorization

- ▶ Scalar
- ▶ Vectorized
- ▶ Vectorized + unrolled

```
sum2(x) = acc_((x,), CompSum{two_sum})
```

Implementation

Why Julia?

Algorithm

- ▶ Summation
- ▶ Dot product

Compensation

- ▶ Naive
- ▶ Compensated (ORO)
- ▶ Compensated (KBN)

Vectorization

- ▶ Scalar
- ▶ Vectorized
- ▶ Vectorized + unrolled

```
dot1(x, y) = \  
    acc_((x, y), CompDot{two_sum})
```


Implementation

Why Julia?

Algorithm

- ▶ Summation
- ▶ Dot product

Compensation

- ▶ Naive
- ▶ Compensated (ORO)
- ▶ Compensated (KBN)

Vectorization

- ▶ Scalar
- ▶ Vectorized
- ▶ Vectorized + unrolled

```
dot2(x, y) = \
    acc_vec((x, y), CompDot{two_sum}, Val(3))
```

Julia implementation

Why Julia? Textbook-like implementation

Algorithm 17 `two_sum` error-free transform

Require: $(a, b) \in \mathbb{F}^2$

Ensure: $x = a \oplus b$ and $x + e = a + b$

$x \leftarrow a \oplus b$

$y \leftarrow x \ominus a$

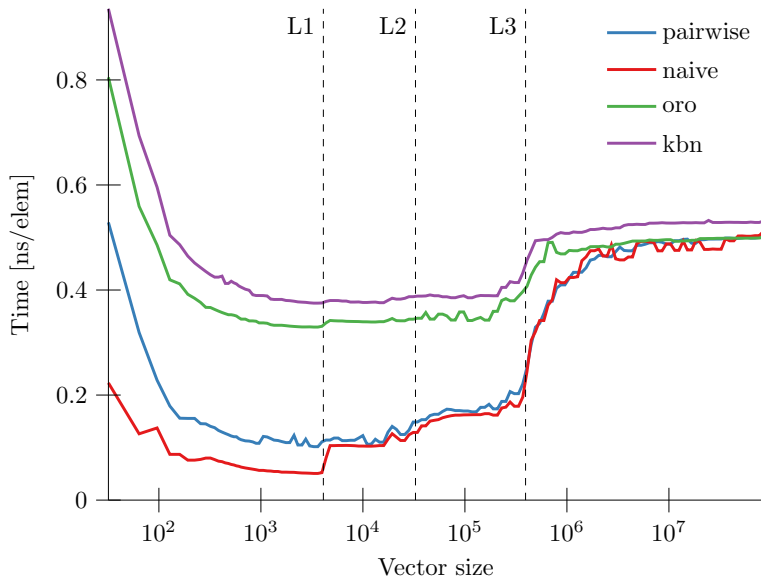
$e \leftarrow (a \ominus (x \ominus y)) \oplus (b \ominus y)$

```
1  function two_sum(a::T, b::T) where {T}
2      SIMDops.@explicit
3
4      x = a + b
5      y = x - a
6      e = (a - (x - y)) + (b - y)
7      return x, e
8  end
```

Julia implementation

Why Julia?

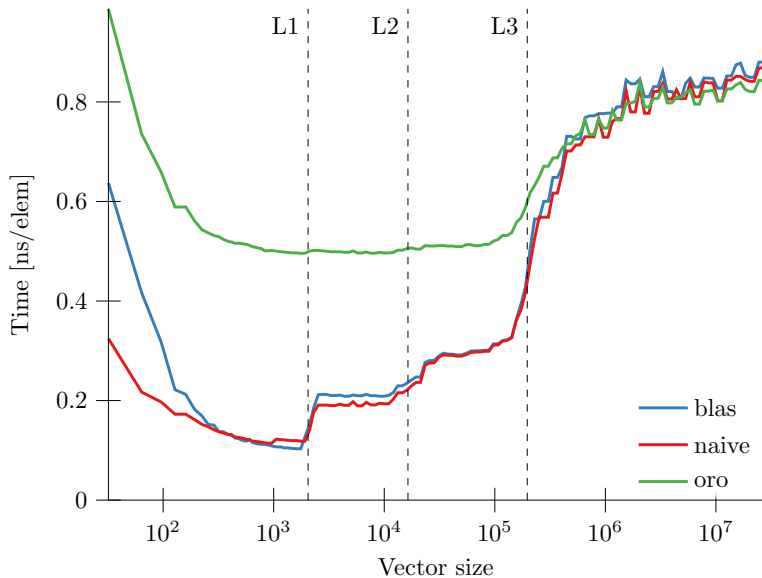
Because it's fast!



Julia implementation

Why Julia?

Because it's fast!



Conclusions

- ▶ Julia solves the 2-language problem :
 - ▶ (relatively) easy to implement (complex) algorithms
 - ▶ good performance (on par with OpenBLAS)
- ▶ Browse the sources of `AccurateArithmetic.jl` to see how this is really done
 - ▶ www.github.com/JuliaMath/AccurateArithmetic.jl
 - ▶ hal.archives-ouvertes.fr/hal-02265534

Thanks!

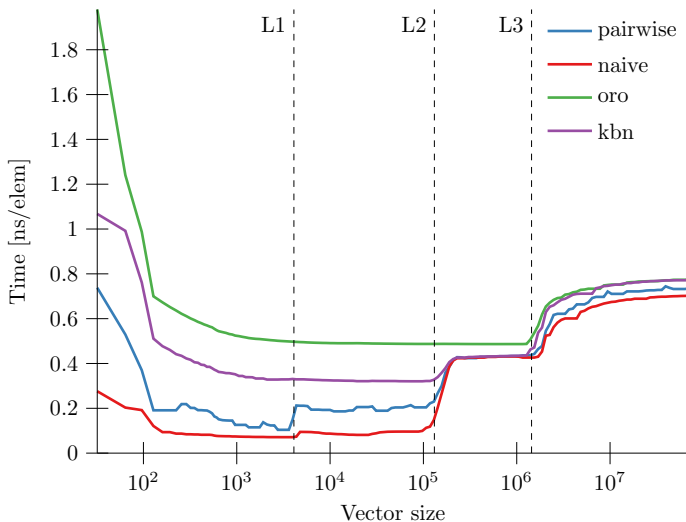
Questions?

Julia implementation

Why Julia?

Because it's fast!

Performance of summation implementations



Julia implementation

Why Julia?

Because it's fast!

Performance of dot product implementations

