

Optimisation : Reduction

Pierre Aubert



LISTIC



UNIVERSITÉ
SAVOIE
MONT BLANC



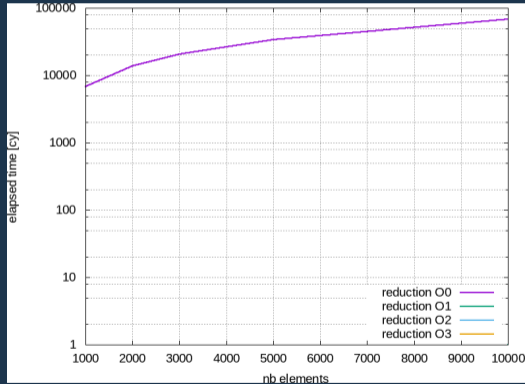
The Reduction (sum)

$$\alpha = \sum_{i=1}^N x_i$$

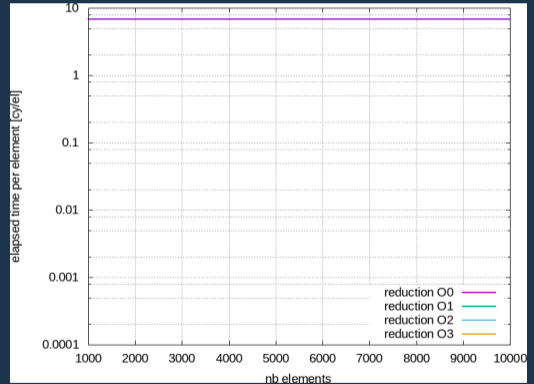


```
float reduction(const float * tabValue, long unsigned int nbElement){
    >> float res(0.0f);
    >> for(long unsigned int i(0lu); i < nbElement; ++i){
    >>     >> res += tabValue[i];
    >> }
    >> return res;
}
```

Total Elapsed Time (cy)

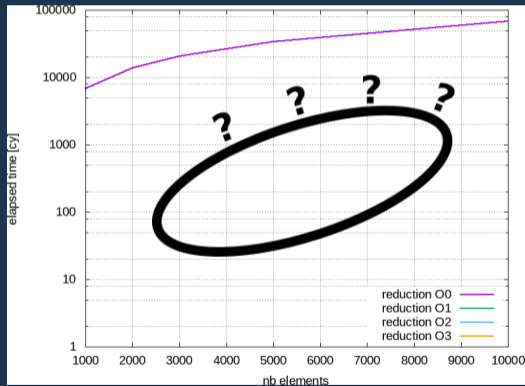


Elapsed Time per element (cy/el)

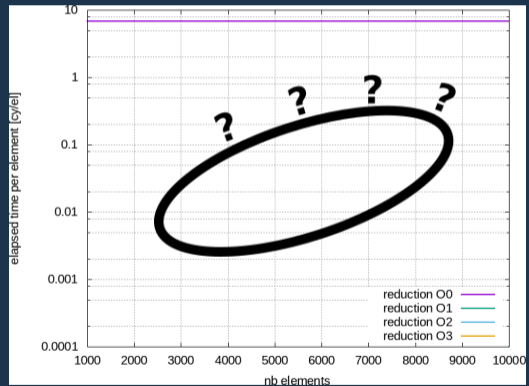


The reduction : first performances

Total Elapsed Time (cy)



Elapsed Time per element (cy/el)



- ▶ Performances **-O0** : slow but reasonable
- ▶ Other performances (**-O1**, **-O2**, **-O3**, **-Ofast**) are too fast (non sense)

GCC is smart of guileful depending on the points of view.

- ▶ GCC noticed you **do not** use the result of the **reduction** function.
- ▶ The call to **reduction** is considered as dead code (or never called code).

To avoid that, you have to compile the **reduction** function in an other file.

Dead code for compiler

```

float reduction(const float * tabValue, long unsigned int nbElement){
    float res(0.0f);
    for(long unsigned int i(0lu); i < nbElement; ++i){
        res += tabValue[i];
    }
    return res;
}

void evaluateReduction(long unsigned int nbElement, long unsigned int nbRepetition){
    float * tabValue = (float*)asterics_malloc(sizeof(float)*nbElement);
    for(long unsigned int i(0lu); i < nbElement; ++i){
        tabValue[i] = (float)(i*32lu%17lu);
    }

    long unsigned int beginTime(rdtsc());
    for(long unsigned int i(0lu); i < nbRepetition; ++i){
        reduction(tabValue, nbElement);
    }
    long unsigned int elapsedTime((double)(rdtsc() - beginTime)/((double)nbRepetition));

    double cyclePerElement(((double)elapsedTime)/((double)nbElement));
    cout << "evaluateReduction : nbElement = " << nbElement << ", cyclePerElement = " << cyclePerElement << " cy/
    el, elapsedTime = " << elapsedTime << " cy" << endl;
    cerr << nbElement << "\t" << cyclePerElement << "\t" << elapsedTime << endl;

    asterics_free(tabValue);
}

```

Dead code for compiler

```

float reduction(const float * tabValue, long unsigned int nbElement){
    float res(0.0f);
    for(long unsigned int i(0lu); i < nbElement; ++i){
        res += tabValue[i];
    }
    return res;
}

void evaluateReduction(long unsigned int nbElement, long unsigned int nbRepetition){
    float * tabValue = (float*)asterics_malloc(sizeof(float)*nbElement);
    for(long unsigned int i(0lu); i < nbElement; ++i){
        tabValue[i] = (float)(i*32lu%17lu);
    }

    long unsigned int beginTime(rdtsc());
    for(long unsigned int i(0lu); i < nbRepetition; ++i){
        reduction(tabValue, nbElement); ← Unused result
    }
    long unsigned int elapsedTime((double)(rdtsc() - beginTime)/((double)nbRepetition));

    double cyclePerElement(((double)elapsedTime)/((double)nbElement));
    cout << "evaluateReduction : nbElement = " << nbElement << ", cyclePerElement = " << cyclePerElement << " cy/
    el, elapsedTime = " << elapsedTime << " cy" << endl;
    cerr << nbElement << "\t" << cyclePerElement << "\t" << elapsedTime << endl;

    asterics_free(tabValue);
}

```

Dead code for compiler

```

float reduction(const float * tabValue, long unsigned int nbElement){
    float res(0.0f);
    for(long unsigned int i(0lu); i < nbElement; ++i){
        res += tabValue[i];
    }
    return res;
}

void evaluateReduction(long unsigned int nbElement, long unsigned int nbRepetition){
    float * tabValue = (float*)asterics_malloc(sizeof(float)*nbElement);
    for(long unsigned int i(0lu); i < nbElement; ++i){
        tabValue[i] = (float)(i*32lu%17lu);
    }

    long unsigned int beginTime(rdtsc());
    for(long unsigned int i(0lu); i < nbRepetition; ++i){
        reduction(tabValue, nbElement); ← Unused result
    }
    long unsigned int elapsedTime((double)(rdtsc() - beginTime)/((double)nbRepetition));

    double cyclePerElement(((double)elapsedTime)/((double)nbElement));
    cout << "evaluateReduction : nbElement = " << nbElement << ", cyclePerElement = " << cyclePerElement << " cy/
el, elapsedTime = " << elapsedTime << " cy" << endl;
    cerr << nbElement << "\t" << cyclePerElement << "\t" << elapsedTime << endl;

    asterics_free(tabValue);
}

```


Dead code for compiler

```
float reduction(const float * tabValue, long unsigned int nbElement){
    float res(0.0f);
    for(long unsigned int i(0lu); i < nbElement; ++i){
        res += tabValue[i];
    }
    return res;
}
```

Useless function

```
void evaluateReduction(long unsigned int nbElement, long unsigned int nbRepetition){
    float * tabValue = (float*)asterics_malloc(sizeof(float)*nbElement);
    for(long unsigned int i(0lu); i < nbElement; ++i){
        tabValue[i] = (float)(i*32lu%17lu);
    }

    long unsigned int beginTime(rdtsc());
    for(long unsigned int i(0lu); i < nbRepetition; ++i){
        reduction(tabValue, nbElement);
    }
    long unsigned int elapsedTime(((double)rdtsc() - beginTime)/((double)nbRepetition));

    double cyclePerElement(((double)elapsedTime)/((double)nbElement));
    cout << "evaluateReduction : nbElement = " << nbElement << ", cyclePerElement = " << cyclePerElement << " cy/
    el, elapsedTime = " << elapsedTime << " cy" << endl;
    cerr << nbElement << "\t" << cyclePerElement << "\t" << elapsedTime << endl;

    asterics_free(tabValue);
}
```

Unused result

Dead code for compiler

```
float reduction(const float * tabValue, long unsigned int nbElement){
    float res(0.0f);
    for(long unsigned int i(0lu); i < nbElement; ++i){
        res += tabValue[i];
    }
    return res;
}
```

Useless function

```
void evaluateReduction(long unsigned int nbElement, long unsigned int nbRepetition){
    float * tabValue = (float*)asterics_malloc(sizeof(float)*nbElement);
    for(long unsigned int i(0lu); i < nbElement; ++i){
        tabValue[i] = (float)(i*32lu%17lu);
    }
}
```

Useless loop

```
long unsigned int beginTime(rdtsc());
for(long unsigned int i(0lu); i < nbRepetition; ++i){
    reduction(tabValue, nbElement);
}
```

Unused result

```
long unsigned int elapsedTime((double)(rdtsc() - beginTime)/((double)nbRepetition));

double cyclePerElement(((double)elapsedTime)/((double)nbElement));
cout << "evaluateReduction : nbElement = " << nbElement << ", cyclePerElement = " << cyclePerElement << " cy/
el, elapsedTime = " << elapsedTime << " cy" << endl;
cerr << nbElement << "\t" << cyclePerElement << "\t" << elapsedTime << endl;

asterics_free(tabValue);
}
```

Dead code for compiler

```
float reduction(const float * tabValue, long unsigned int nbElement){
    float res(0.0f);
    for(long unsigned int i(0lu); i < nbElement; ++i){
        res += tabValue[i];
    }
    return res;
}
```

Useless function

```
void evaluateReduction(long unsigned int nbElement, long unsigned int nbRepetition){
    float * tabValue = (float*)asterics_malloc(sizeof(float)*nbElement);
    for(long unsigned int i(0lu); i < nbElement; ++i){
        tabValue[i] = (float)(i*32lu%17lu);
    }
}
```

Nothing between timers

```
long unsigned int beginTime(rdtsc());
for(long unsigned int i(0lu); i < nbRepetition; ++i){
    reduction(tabValue, nbElement);
}
long unsigned int elapsedTime(((double)rdtsc() - beginTime)/((double)nbRepetition));
```

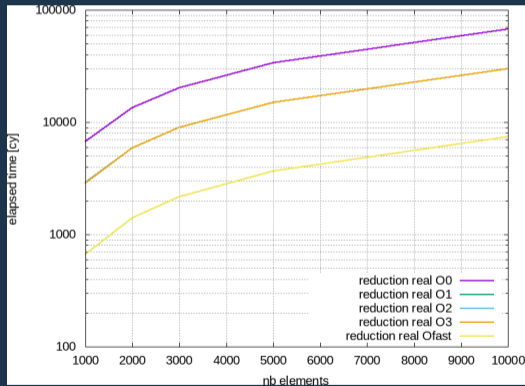
Useless loop

Unused result

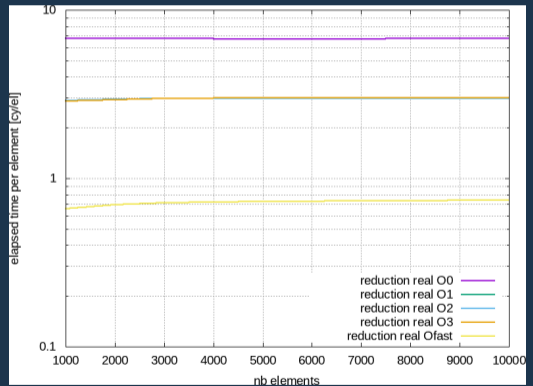
```
double cyclePerElement(((double)elapsedTime)/((double)nbElement));
cout << "evaluateReduction : nbElement = " << nbElement << ", cyclePerElement = " << cyclePerElement << " cy/
el, elapsedTime = " << elapsedTime << " cy" << endl;
cerr << nbElement << "\t" << cyclePerElement << "\t" << elapsedTime << endl;

asterics_free(tabValue);
}
```

Total Elapsed Time (cy)



Elapsed Time per element (cy/el)



▶ Data alignment :

- ▶ All the data to be aligned on vectorial registers size.
- ▶ Change **new** or **malloc** to **memalign** or **posix_memalign**

You can use **asterics_malloc** to have LINUX/MAC compatibility (in **evaluateReduction**):

```
(float*)asterics_malloc(sizeof(float)*nbElement);
```

The **__restrict__** keyword (arguments of **reduction** function):

```
const float * __restrict__ ptabValue
```

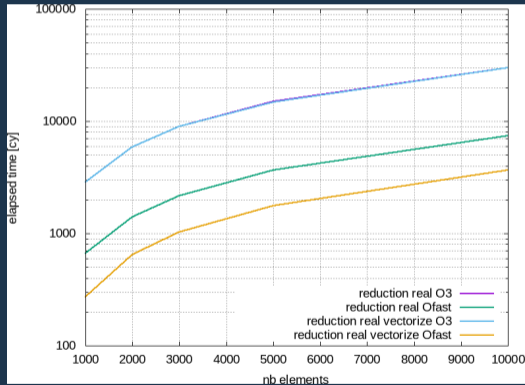
The **__builtin_assume_aligned** function call (in **reduction** function):

```
const float* tabValue = (const float*)__builtin_assume_aligned(ptabValue, VECTOR_ALIGNEMENT);
```

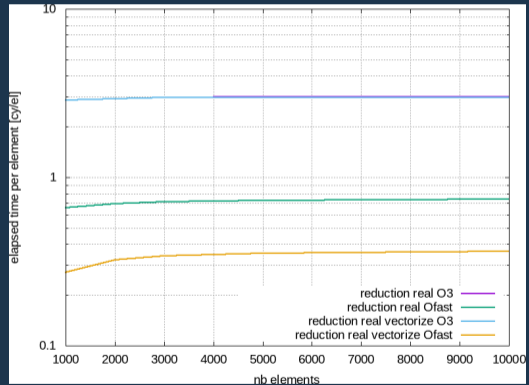
▶ The Compilation Options become :

- ▶ **-O3 -ftree-vectorize -march=native -mtune=native -mavx2**

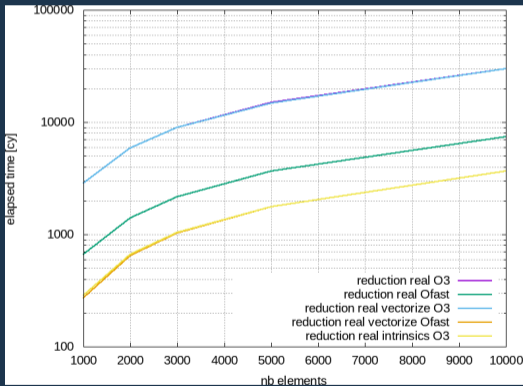
Total Elapsed Time (cy)



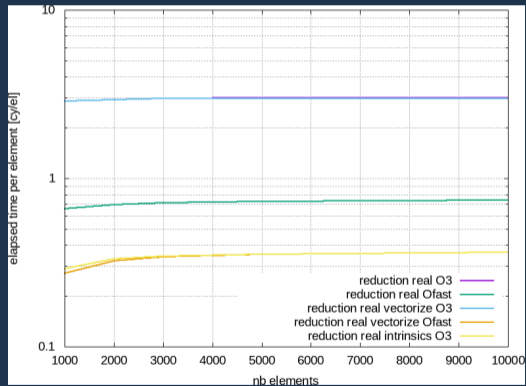
Elapsed Time per element (cy/el)



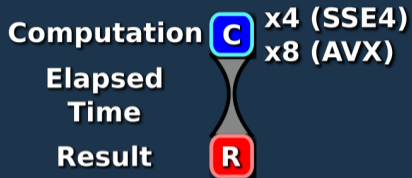
Total Elapsed Time (cy)



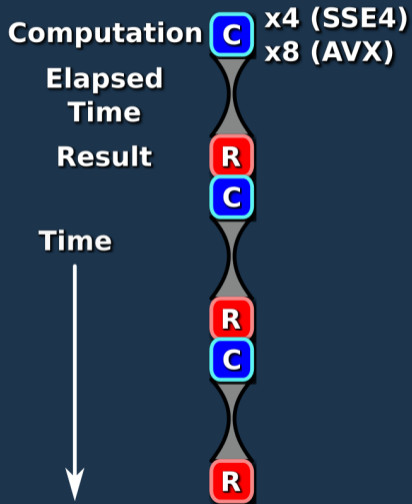
Elapsed Time per element (cy/el)





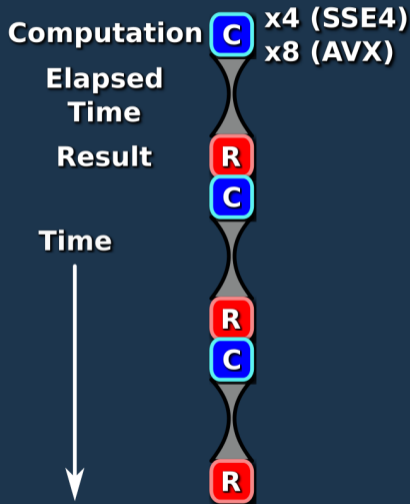


1 accumulator

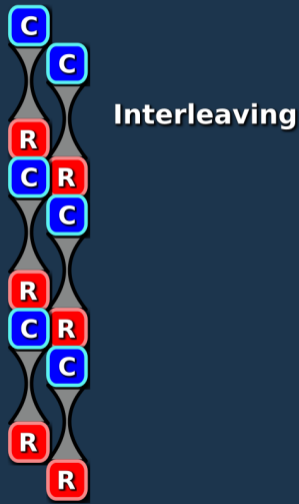


Reduction optimisation

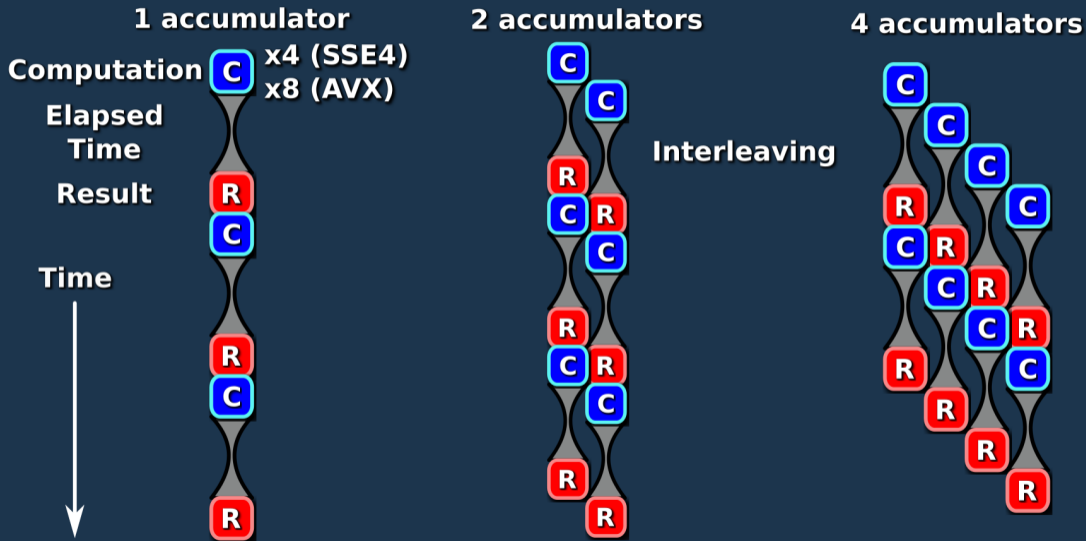
1 accumulator



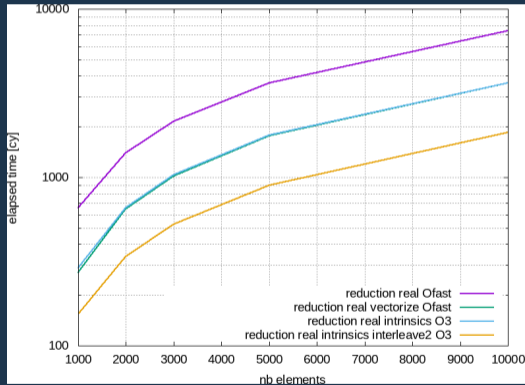
2 accumulators



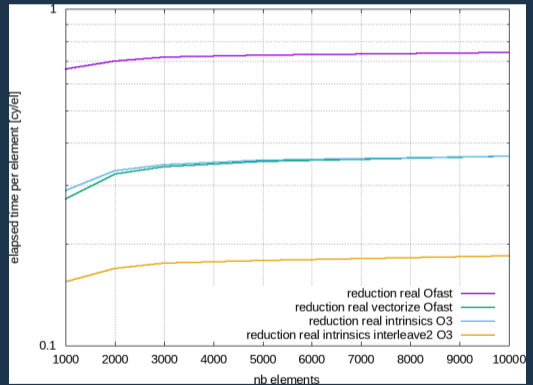
Reduction optimisation



Total Elapsed Time (cy)

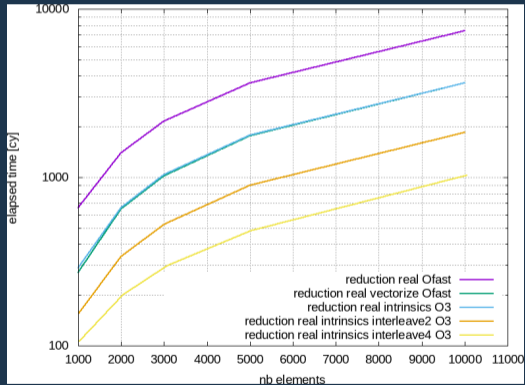


Elapsed Time per element (cy/el)

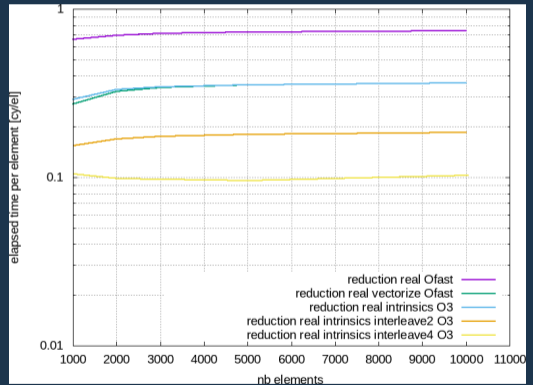


The reduction : intrinsics interleaved 4

Total Elapsed Time (cy)

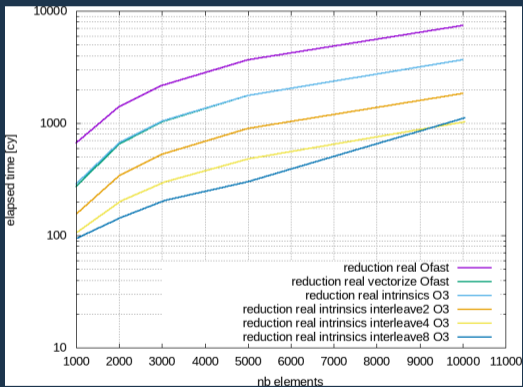


Elapsed Time per element (cy/el)

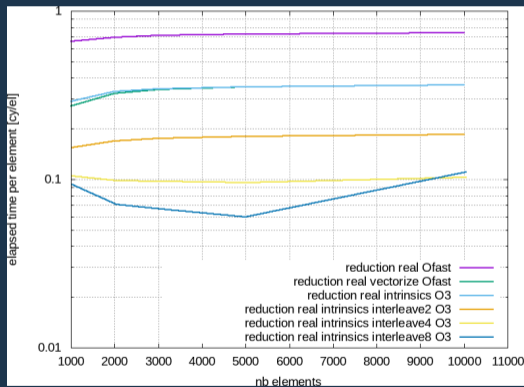


The reduction : intrinsics interleaved 8

Total Elapsed Time (cy)

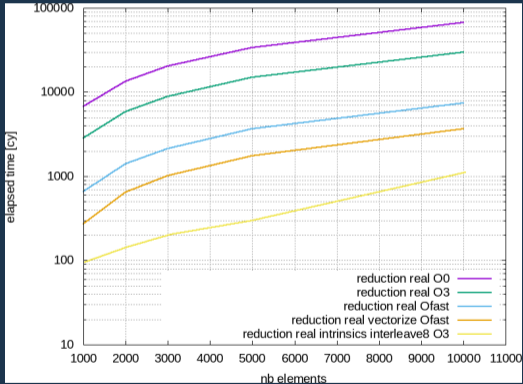


Elapsed Time per element (cy/el)

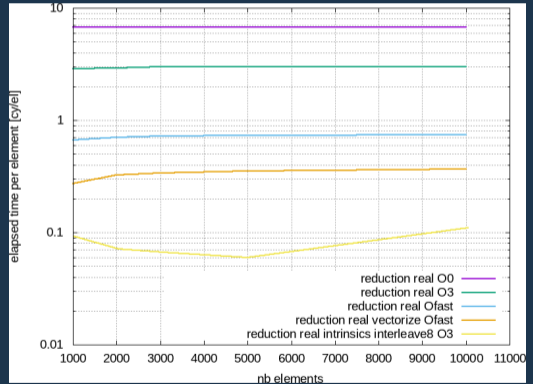


The reduction : summary

Total Elapsed Time (cy)



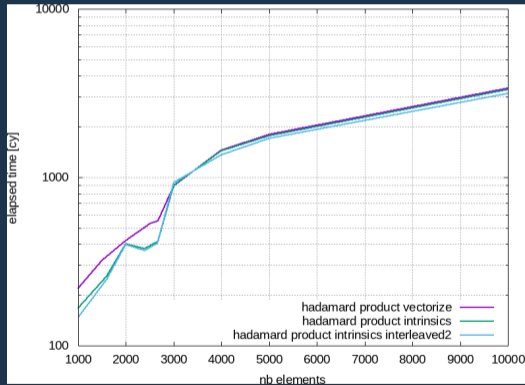
Elapsed Time per element (cy/el)



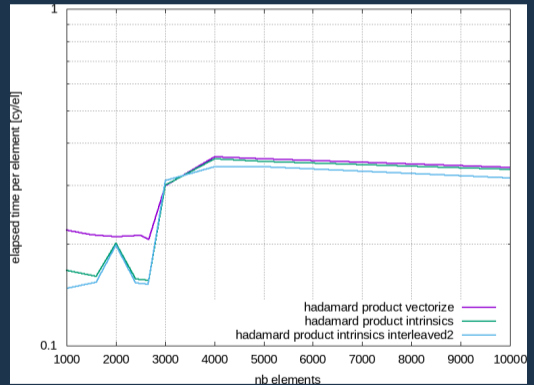
5000 elements, Intrinsics is **166** times faster than **-O0** and **7** times faster than **-Ofast** vectorized

What about the Hadamard product ?

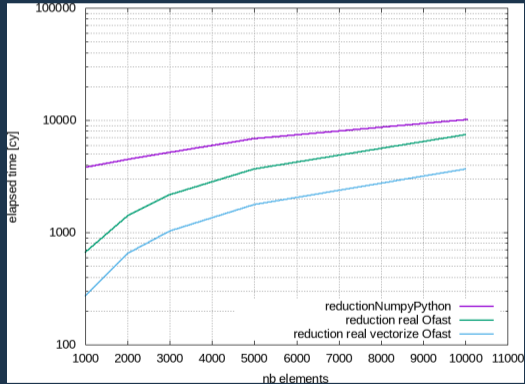
Total Elapsed Time (cy)



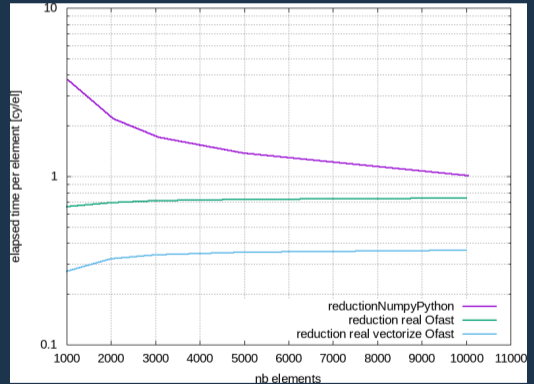
Elapsed Time per element (cy/el)



Total Elapsed Time (cy)

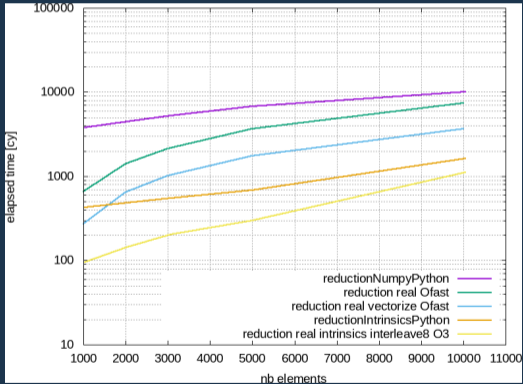


Elapsed Time per element (cy/el)

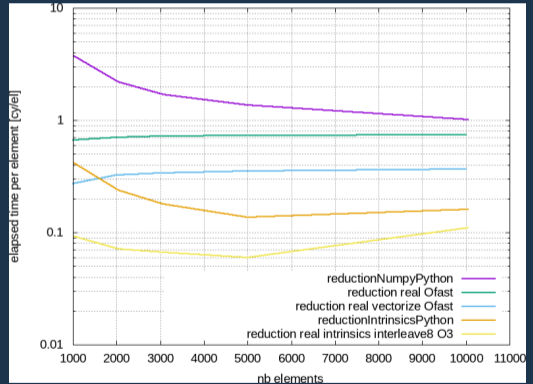


1000 elements, GCC vectorized version is **13** times faster than **numpy** sum

Total Elapsed Time (cy)



Elapsed Time per element (cy/el)



1000 elements, our python reduction is **10** times faster than **numpy** sum