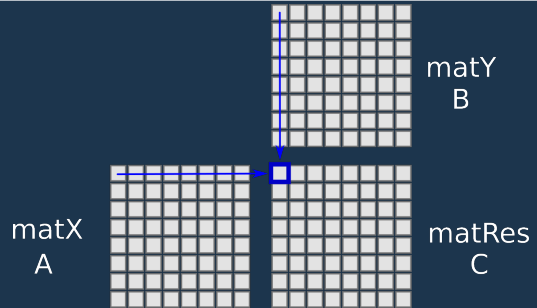


Optimisation : sgemm

Pierre Aubert

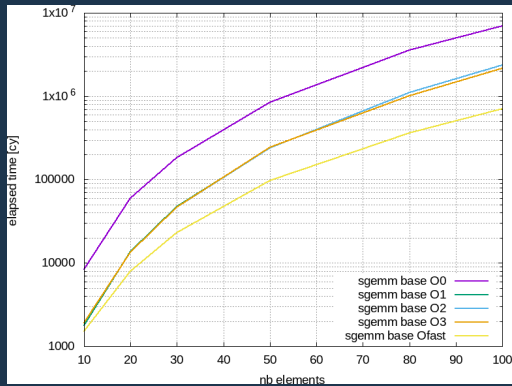


$$C_{i,j} = \sum_{k=1}^N A_{i,k} \cdot B_{k,j}, \quad \forall i, j \in 1, N$$

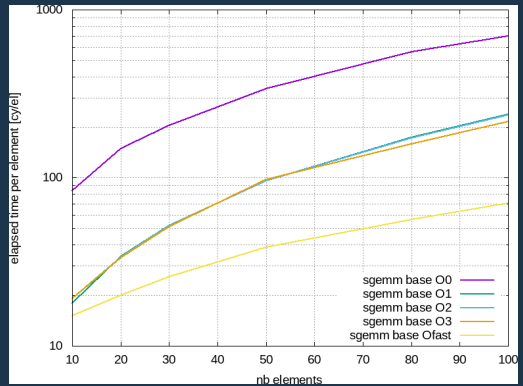


```
void sgemm(float* matOut, const float * matX, const float* matY, long unsigned int size){
»   for(long unsigned int i(0lu); i < size; ++i){
»       for(long unsigned int j(0lu); j < size; ++j){
»           float res(0.0f);
»           for(long unsigned int k(0lu); k < size; ++k){
»               res += matX[i*size + k]*matY[k*size + j];
»           }
»           matOut[i*size + j] = res;
»       }
»   }
}
```

Total Elapsed Time (cy)



Elapsed Time per element (cy/el)



Let's swap loops over j and k



A red rectangular box containing a white 'X' symbol, indicating a warning or error.

**Linked
Image
Not Found**

Let's swap loops over j and k



A red rectangular box containing a white 'X' symbol, indicating a warning or error.

**Linked
Image
Not Found**

Let's swap loops over j and k

X

**Linked
Image
Not Found**

X

**Linked
Image
Not Found**

Let's swap loops over j and k

X

**Linked
Image
Not Found**

X

**Linked
Image
Not Found**

Let's swap loops over j and k

X

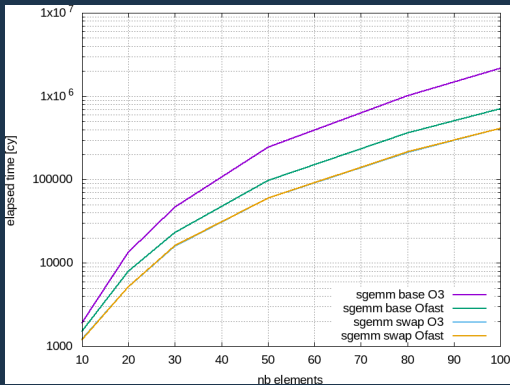
**Linked
Image
Not Found**

X

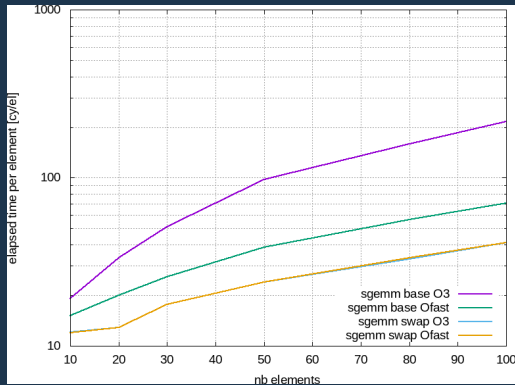
**Linked
Image
Not Found**

SGEMM : Swap loop performances

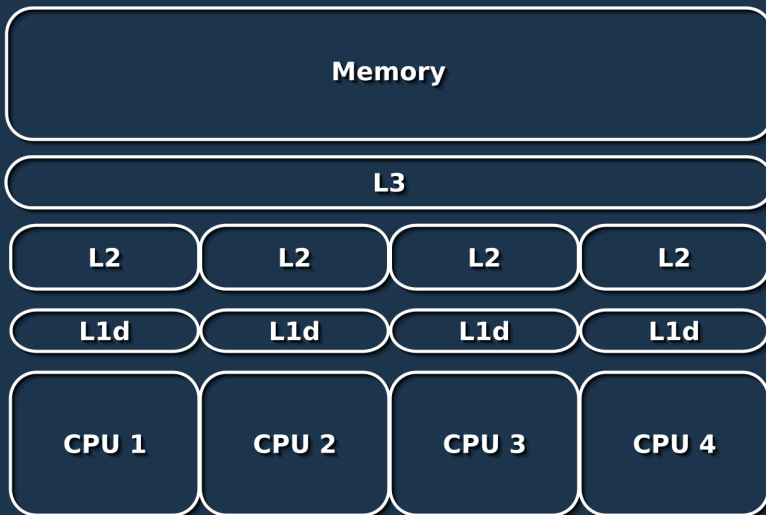
Total Elapsed Time (cy)



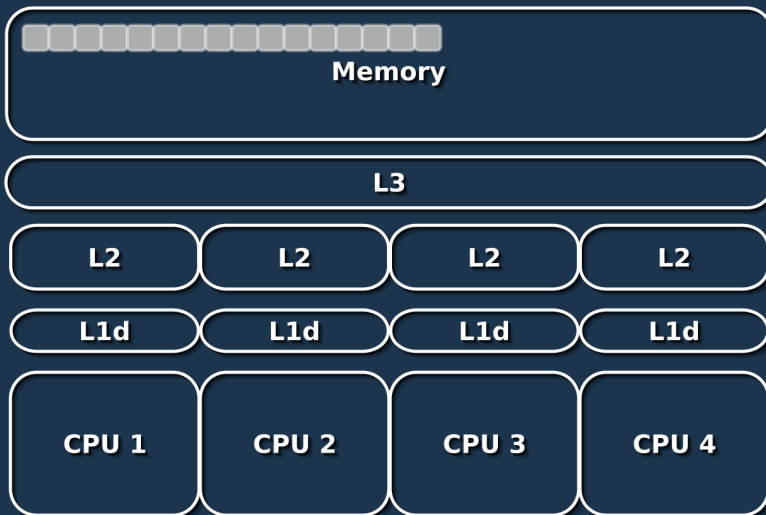
Elapsed Time per element (cy/el)



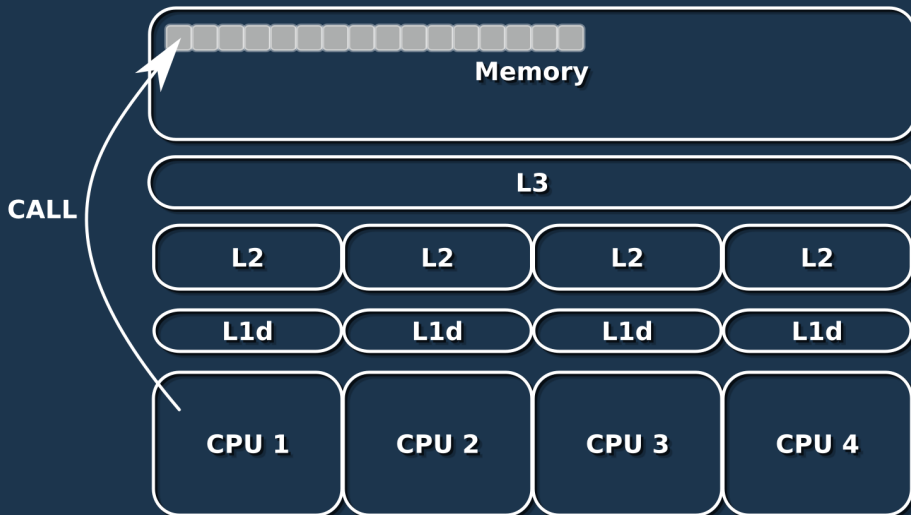
Data Fetching and pre-fetching



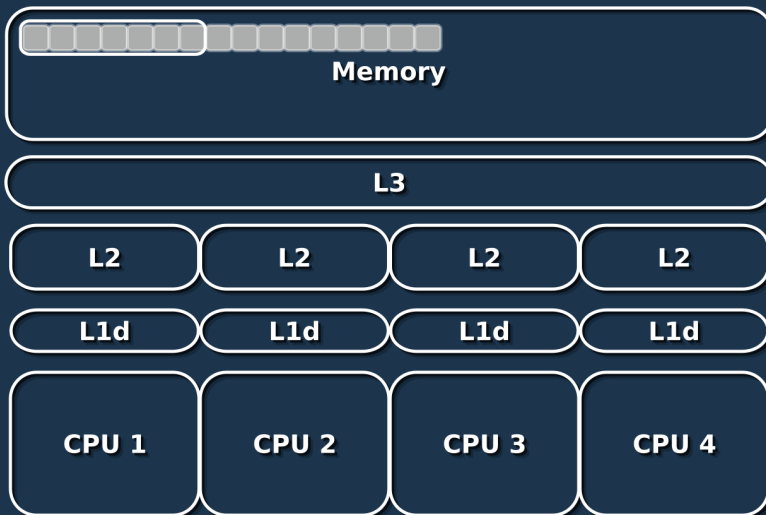
Data Fetching and pre-fetching



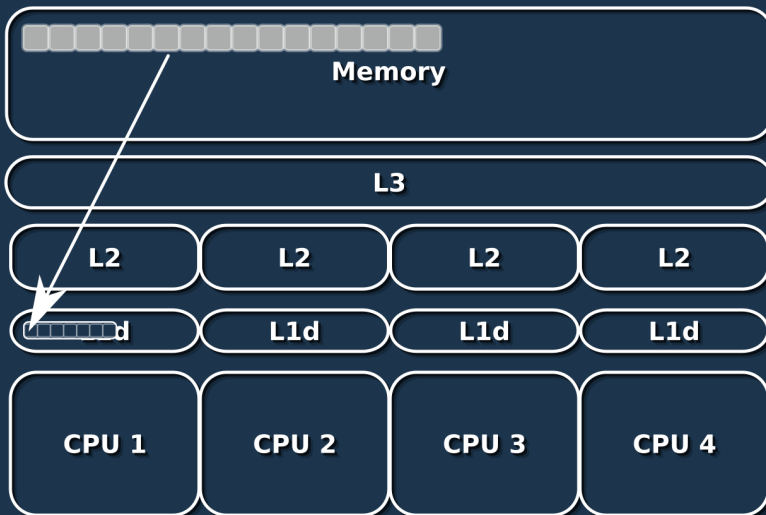
Data Fetching and pre-fetching



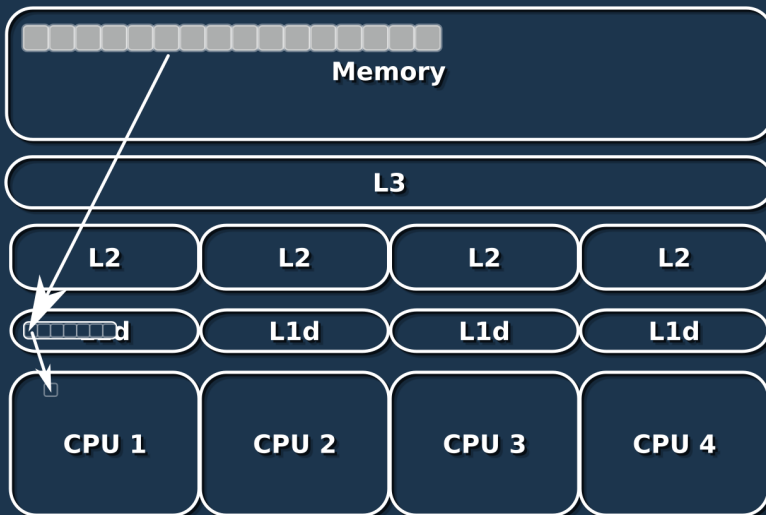
Data Fetching and pre-fetching



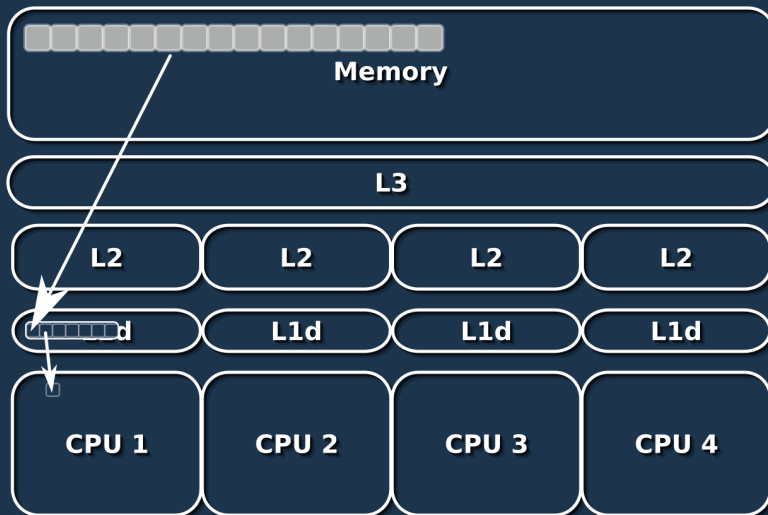
Data Fetching and pre-fetching



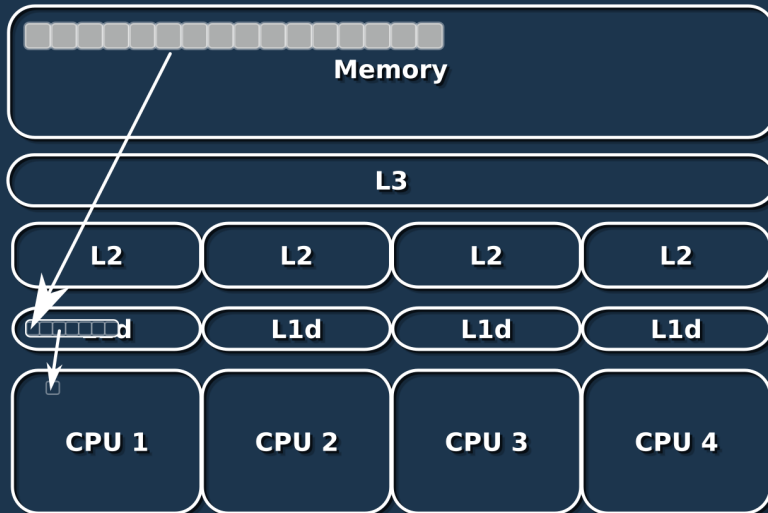
Data Fetching and pre-fetching



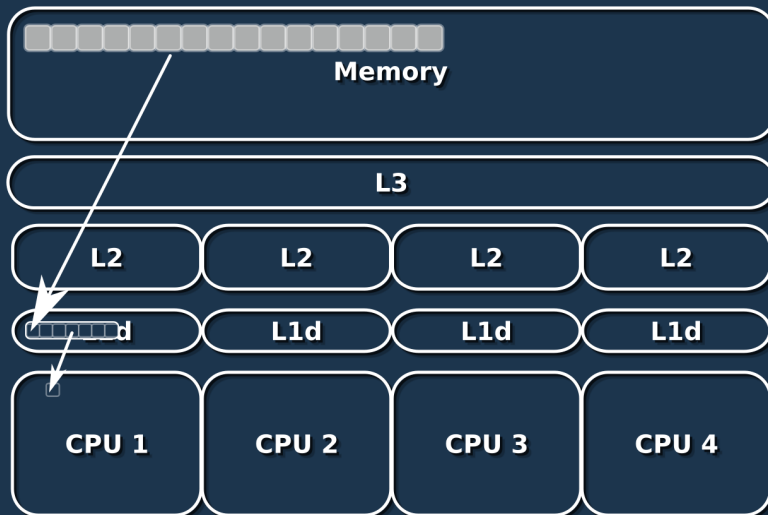
Data Fetching and pre-fetching



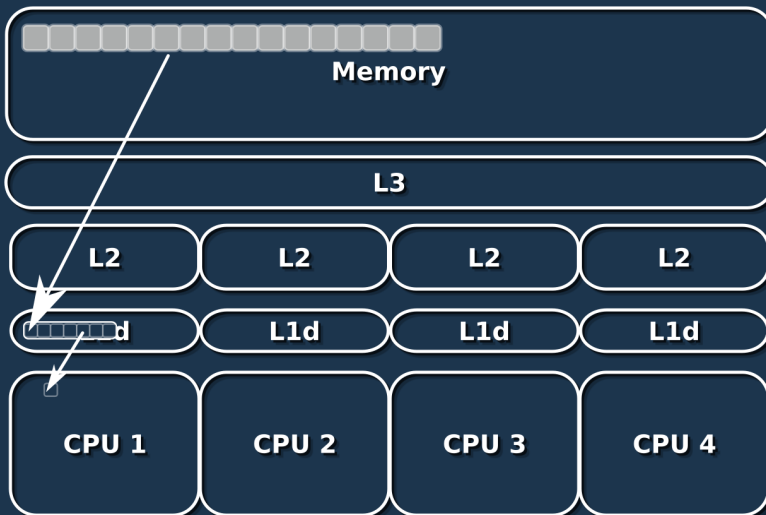
Data Fetching and pre-fetching



Data Fetching and pre-fetching

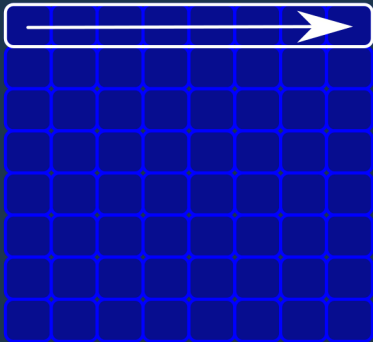


Data Fetching and pre-fetching

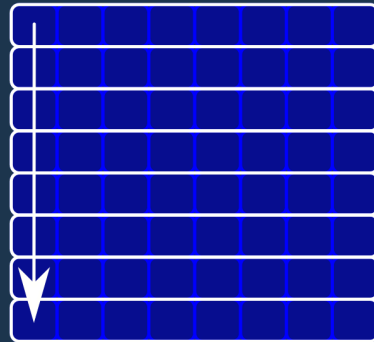


Data locality

Data Fetch Size 8 elements

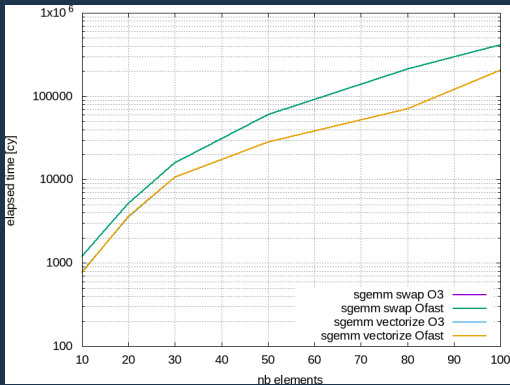


Read 8 elements with one fetch

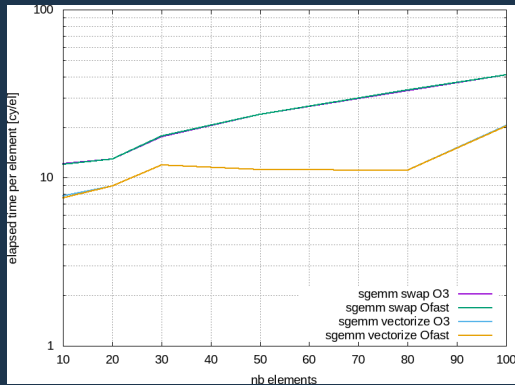


Read 8 elements with 8 fetchs
=> fetch 64 elements instead of 8

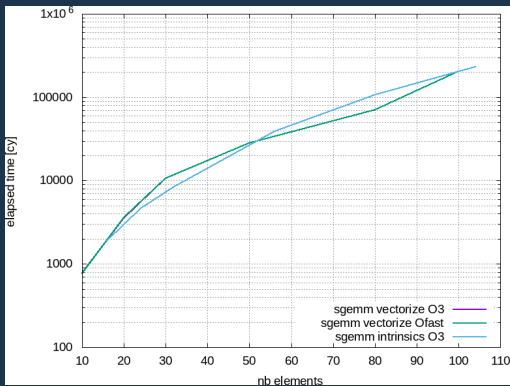
Total Elapsed Time (cy)



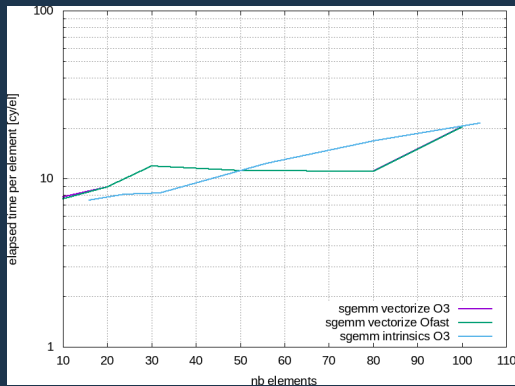
Elapsed Time per element (cy/el)



Total Elapsed Time (cy)



Elapsed Time per element (cy/el)



What if matrices don't have proper size ?

For our intrinsics version :

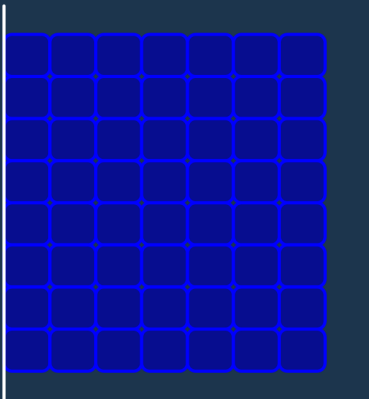
The number of columns has to be a multiple of 8

What if matrices don't have proper size ?

For our intrinsics version :

The number of columns has to be a multiple of 8

If is it not the case :



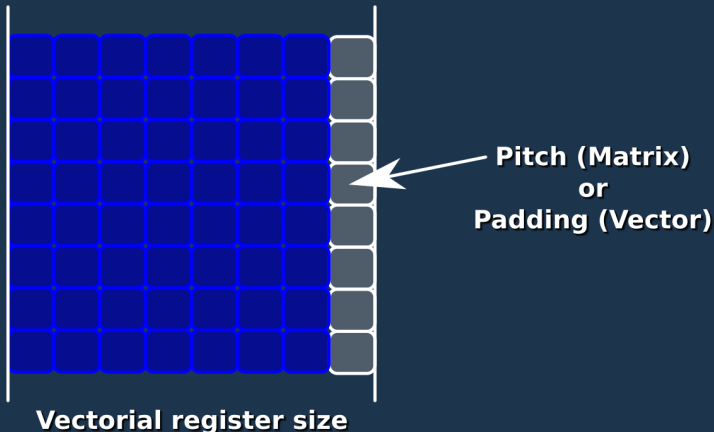
Vectorial register size

What if matrices don't have proper size ?

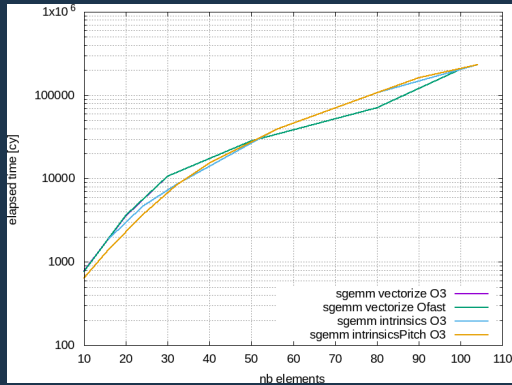
For our intrinsics version :

The number of columns has to be a multiple of 8

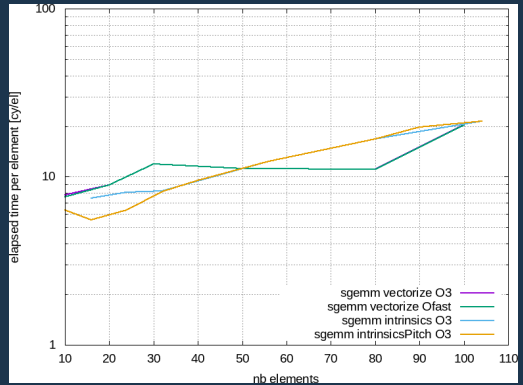
If is it not the case :



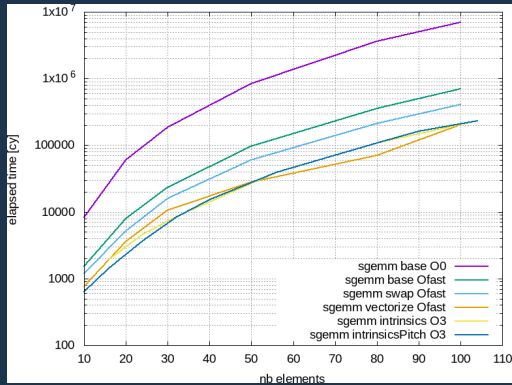
Total Elapsed Time (cy)



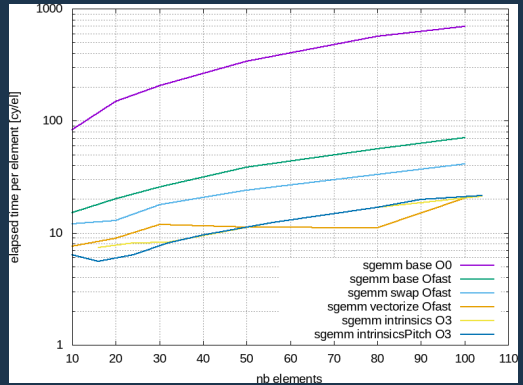
Elapsed Time per element (cy/el)



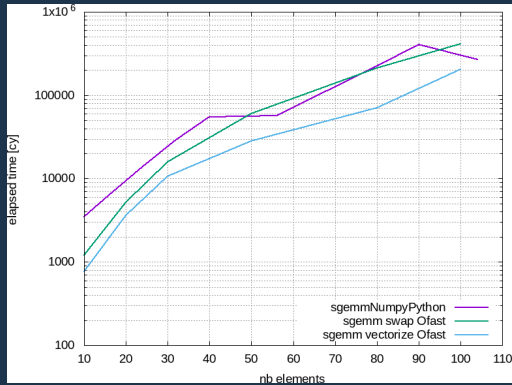
Total Elapsed Time (cy)



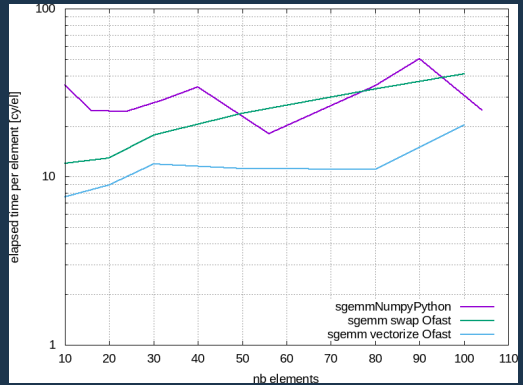
Elapsed Time per element (cy/el)



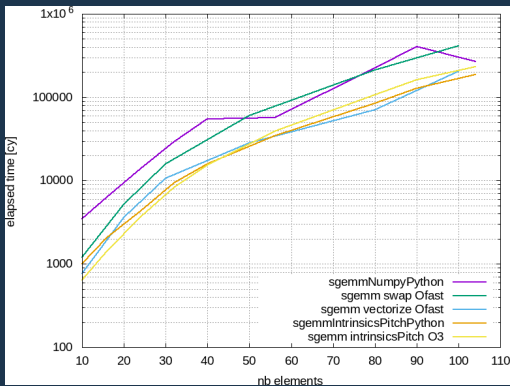
Total Elapsed Time (cy)



Elapsed Time per element (cy/el)



Total Elapsed Time (cy)



Elapsed Time per element (cy/el)

