# Optimisation : sgemm

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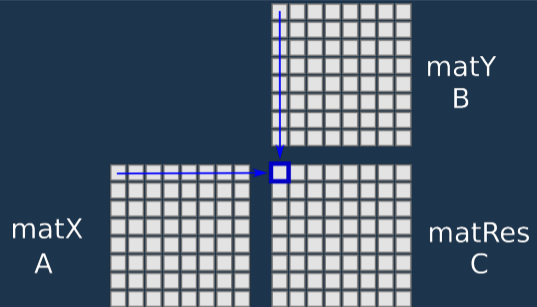
$$C_{i,j} \;=\; \sum_{k=1}^{N} A_{i,k} \cdot B_{k,j}, \quad \forall i,j \in 1, N$$

matY
B

matX
A

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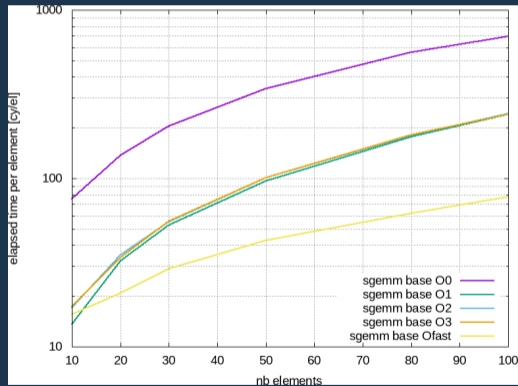
```
void sgemm(float* matOut, const float * matX, const float* matY, long unsigned int size){
        for(long unsigned int i(0lu); i < size; ++i){
                for(long unsigned int j(0lu); j < size; ++j){
                        float res(0.0f);
                        for(long unsigned int k(0lu);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

```
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;
        }
    }
}
```

```
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;
        }
    }
}
```

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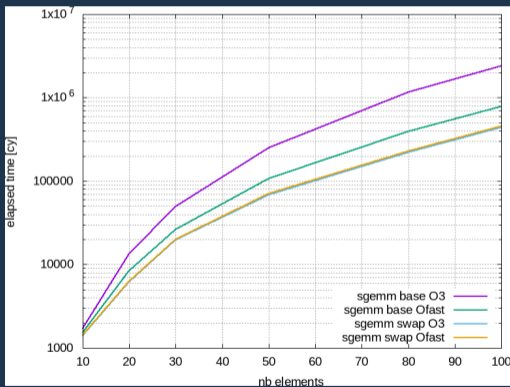
```cpp
void sgemm(float* matOut, const float * matX, const float* matY, long unsigned int size){
        memset(matOut, 0, sizeof(float)*size*size);
        for(long unsigned int i(0lu); i < size; ++i){
                for(long unsigned int k(0lu); k < size; ++k){
                        for(long unsigned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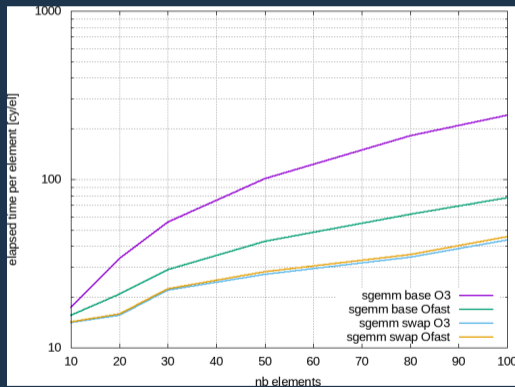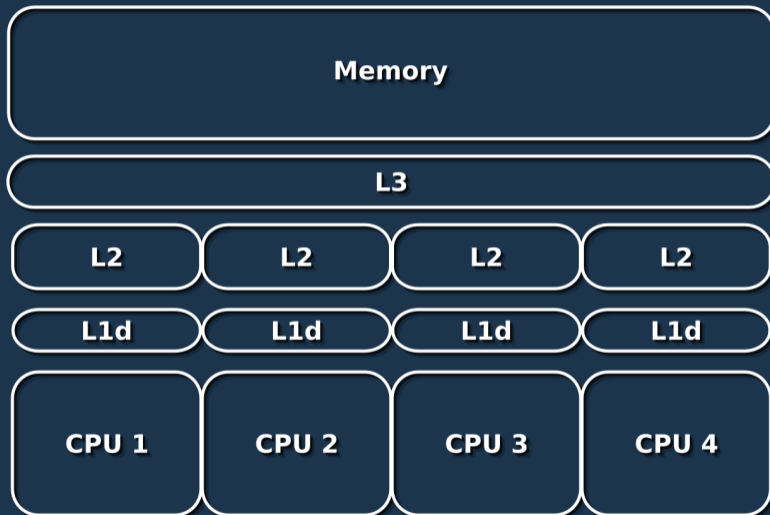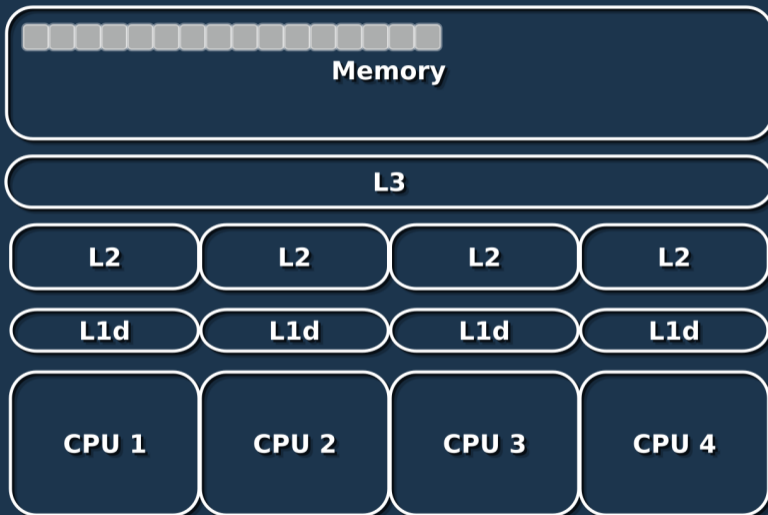
```
void sgemm(float* matOut, const float * matX, const float* matY, long unsigned int size){
    for(long unsigned int i(0lu); i < size; ++i){
        for(long unsigned int j(0lu); j < size; ++j){
            float res(0.0f);
            for(long unsigned int k(0lu); k < size; ++k){
               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

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

```
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;
        }
    }
}
```

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

Total Elapsed Time (cy)

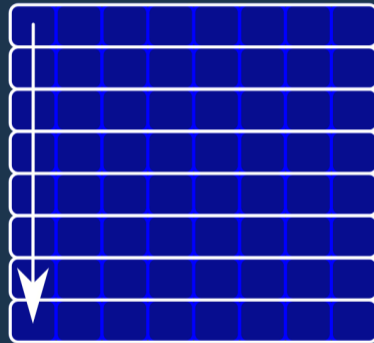Elapsed Time per element (cy/el)

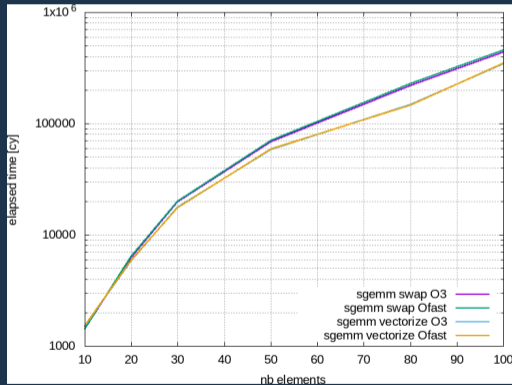**Data Fetch Size 8 elements**
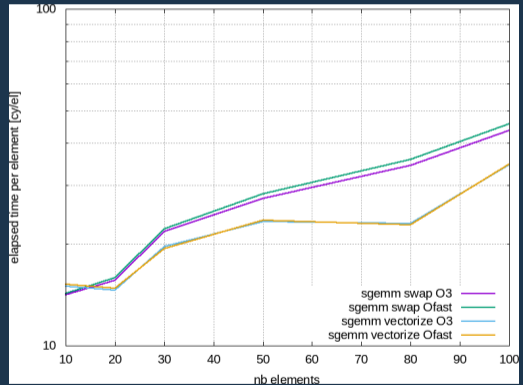


**Read 8 elements with one fetch**
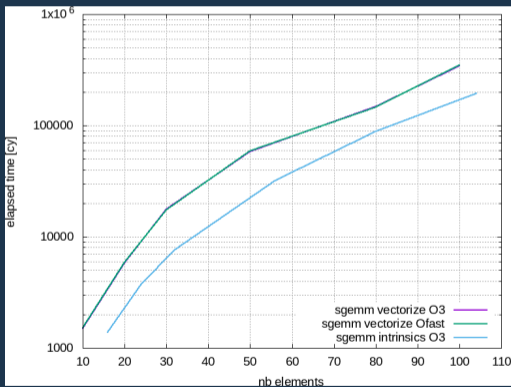
**Read 8 elements with 8 fetchs
=> fecth 64 elements instead of 8**

Total Elapsed Time (cy)

Elapsed Time per element (cy/el)
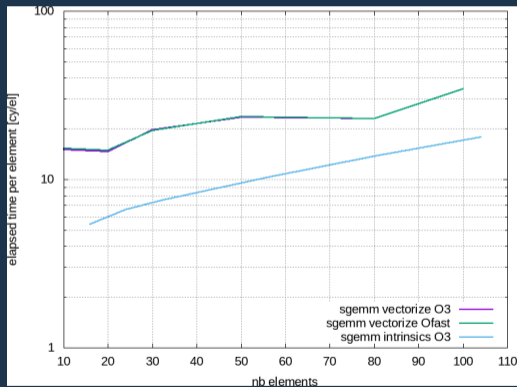
Total Elapsed Time (cy)

Elapsed Time per element (cy/el)

**For our intrinsics version :**
    **The number of columns has to be a multiple of 8**

**For our intrinsics version :**
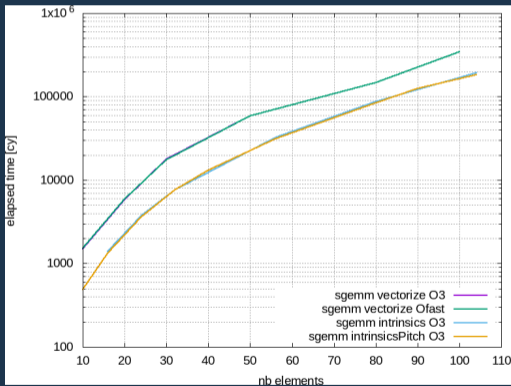   **The number of columns has to be a multiple of 8**

**If is it not the case :**



**Vectorial register size**

**For our intrinsics version :**
   **The number of columns has to be a multiple of 8**

**If is it not the case :**



Pitch (Matrix)
or
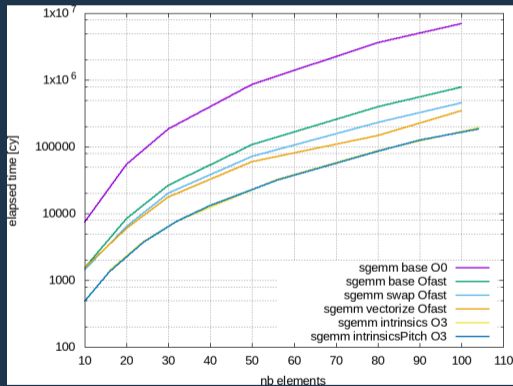Padding (Vector)

**Vectorial register size**
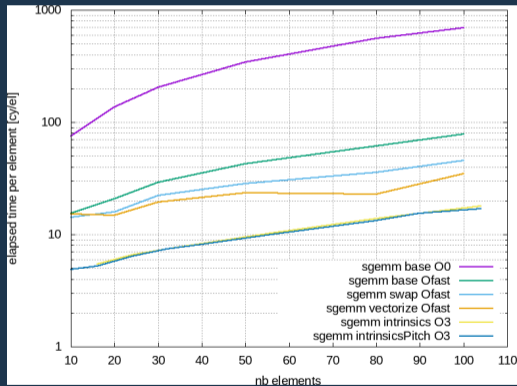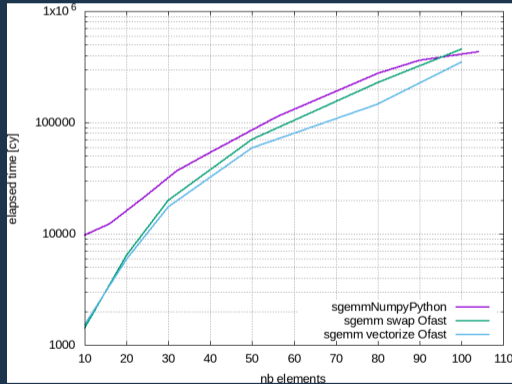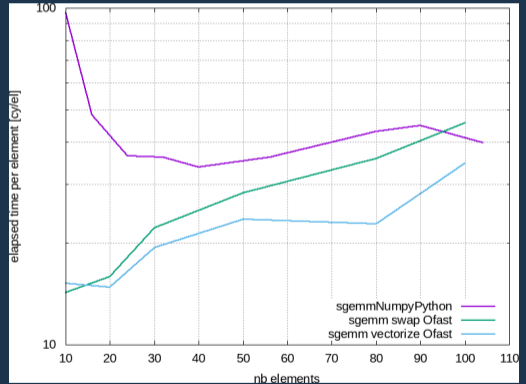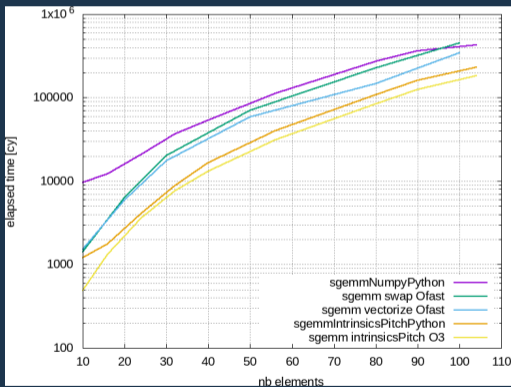
Total Elapsed Time (cy)

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)