

# BONE ANISOTROPY MAPPING

Jarunan Panyasantisuk, Joao Rivera, Rajan Gill, Ryan Cherifa

Department of Computer Science  
ETH Zürich  
Zürich, Switzerland

## ABSTRACT

Todo

## 1. INTRODUCTION

Bone fabric anisotropy or microstructure orientation was recently included in finite element (FE) models to improve the accuracy in predicting bone stiffness and strength. To save computing cost, an FE model of bone is generated from a clinical computer tomography (CT) scanned image with low resolution (1-3mm). However, the bone microstructure details can be obtained only by high resolution peripheral CT with the resolution of 60-82  $\mu\text{m}$ . Therefore, bone anisotropy mapping methodology was required to map the bone microstructure orientation from the high resolution image onto the low resolution image.

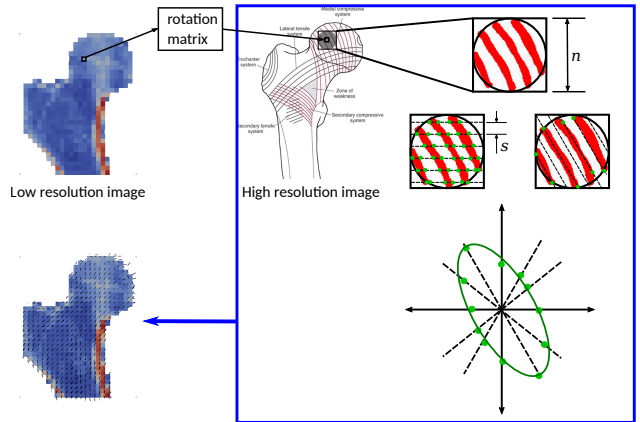
Bone anisotropy mapping methodology includes coordinate mapping between a low and a high-resolution images, region extraction, the mean intercept length (MIL) method for quantification of the microstructure orientation, ellipsoid fitting of MIL and eigendecomposition to obtain the major direction of the microstructure.

**Motivation.** In the recent study, bone anisotropy mapping algorithms are performed for all low resolution image voxels and for all pairs ( $n=71$ ) of low and high resolution images. This preprocessing step consumes a significant amount of computing time. Moreover, researchers expect larger dataset to create a more general FE models of bone. Therefore, the performance of these algorithms are crucial. Two software packages for image processing which include MIL calculation are Medtool, a commercialized PYTHON package, and BoneJ, an open-source JAVA plugin for ImageJ. The external packages needed to be integrated to the computation pipeline and optimization is not straightforward.

In this paper, we are presenting an integrated and optimized methodology of bone anisotropy mapping. At our best knowledge, this is the first paper to study and optimize the performance of MIL calculation depending on the extracted region size.

## 2. BACKGROUND

The methodology is shown in Fig. 1. Coordinates from a low resolution image were mapped to its own high resolution image. Then, a sphere region is extracted and centered at the mapped coordinate in the high resolution image. Subsequently, the anisotropy of the extracted bone region is quantified by using the mean intercept length (MIL) method which imposed direction vectors on the regions. The mean length of each vector is the sum of the length inside the bone region divided by the number of intercepts which intersect with bone/non-bone transition. The MIL values can be plots as a cloud of points in the direction vector space and an ellipsoid can be fitted to obtain a representative two dimensional tensor, for which three eigenvalues and three eigenvectors are calculated. The eigenvector associated with the minimum eigenvalue is the major direction of that bone region.



**Fig. 1.** Bone anisotropy mapping methodology includes coordinate mapping from low to high resolution image, region extraction from the high resolution image, MIL calculation, ellipsoid fitting to obtain a fabric tensor and eigendecomposition of the fabric tensor from which the major eigenvector can be visualized.  $n$  is the extracted region dimension in one direction.  $s$  is the stride between parallel direction vectors.

**Region extraction.** This algorithm applies a sphere mask on the high resolution image and copies the extraction region to a separate array. A multiplication was performed for each image voxel.

**Mean Intercept Length (MIL) method.**

The mean intercept length of a vector  $v$  can be expressed as

$$MIL(v) = \frac{h(v)}{C(v)}, \quad (1)$$

where  $h(v)$  is the summation of the bone intensities “touched” by all the rays formed by direction vector  $v$ , and  $C(v)$  is the number of intercepts of the vector  $v$ . In general, each ray is separated by a stride  $s = 2$ . Thus, each direction vector touches only  $\frac{n^3}{s^2}$  voxels in the region for  $n$  even. We consider only  $N_{vec} = 13$  direction vectors.

As an example, consider the (2-D) extracted region in Fig 1 with horizontal vector  $\vec{v} = (1, 0, 0)$ . In this region,  $C(v)$  is the number of green dots, whereas  $h(v)$  is calculated summing up all voxels touched by the horizontal dotted lines.

The algorithm includes floating point additions to calculate  $h(v)$  and comparisons to detect interceptions. Note that additions performed to calculate  $C(v)$  are not considered in the cost analysis because those are performed over integers. In addition, one division is performed per direction vector. The total cost of the MIL algorithm per extracted region is therefore  $C(n) = N_{vec}(\frac{2n^3}{s^2} + 1)$ . Table 1 provides a detailed breakdown of the cost analysis.

**Ellipsoid fitting. Todo**

From the basic implementation, MIL calculation, ellipsoid fitting and region extraction consume approximately 80%, 15% and 5 % of the overall computing time, respectively. Therefore, we focused on these three algorithms for the optimization. The run time was measured by using time stamp counter (TSC). We define our cost measure as the number of floating point operations (flops) performed by the algorithms. The cost analysis for each algorithm is shown in Table 1, where  $n$  is the size of the region extracted during RE. Note that each algorithm is performed once per non-zero voxel in the low resolution image. (This can be omitted) Thus, the cost of the whole program is total  $\cdot N_{nz}$ , where  $N_{nz}$  is the number of non-zero voxels in the low resolution image which represent the bone voxel.

**Table 1.** Cost analysis per low-resolution bone voxel.

	$N_{add}$	$N_{mul}$	$N_{div}$	$N_{cmp}$	Total
RE	-	$n^3$	-	-	$n^3$
MIL	$\frac{13n^3}{4}$	-	13	$\frac{13n^3}{4}$	$13(\frac{n^3}{2} + 1)$
EF					

### 3. METHODS

This section explains the baseline implementation and optimization strategies of each focused algorithm.

**Region extraction.** The baseline implementation loops over all voxels in the sphere mask with the size of  $n^3$ . For each sphere mask voxel, if a corresponding voxel in the high resolution image is found, the extracted region voxel is set as the multiplication of the sphere mask voxel and the high resolution image voxel. Therefore, loop unrolling with four parallel multiplications and scalar replacement were applied to speed up the computation and avoid aliasing.

**MIL.** The following C-like pseudo code summarizes the baseline implementation of MIL algorithm:

```

1 void mil_base(double* reg, int n, double* mil){
2   double h[13] = 0; int C[13] = 0;
3   // Core of MIL
4   for (int v = 0; v < 13; ++v)
5     for every ray r of v do
6       {k,j,i} = get_start_of_ray(r);
7       while ( {k,j,i} < n && {k,j,i} > 0 )
8         h[v] += reg[k][j][i];
9         curr = reg[k][j][i] > 0.5; //Is it bone?
10        C[v] += curr ^ prev; // interception
11        prev = curr;
12        {k,j,i} = next_voxel_indices({k,j,i}, v);
13        mil[v] = h[v] / C[v];
14  }
```

As can be seen, the baseline implementation iterates through all rays of all 13 direction vectors (see the `for` loop in lines 4-5). For each ray, we first get its starting voxel to begin iterating from (line 6). Finally, the inner loop iterates through the whole ray and computes  $h(v)$  and  $C(v)$ . Note that  $C(v)$  is an integer array; thus, the only floating point operations in this loop are the addition and comparison in lines 8-9.

The first performance bottleneck in the baseline implementation is due to the non-associativity of floating point additions in line 8. This creates inter-loop dependencies between addition operations that limit the instruction level parallelism (ILP) of the implementation. Assuming a latency of floating point additions of 4 cycles, an upper bound on the performance due to this dependency is 0.5 flops/cycle (1 add + 1 cmp every 4 cycles).

In order to improve ILP, we used accumulators with loop-unrolling. We first unrolled the loop in line 5, to get four different rays of the same vector. We created an accumulator for each ray; thus, now we can perform 4 additions and 4 comparisons in parallel. Note that the rays for each accumulator should be of the same length to avoid handling leftovers at the end. For 1-Dimensional direction vectors, e.g. (1,0,0), all rays have the same length (equal to  $n$ ); however, this is not the case for 2-D and 3-D vectors. Thus, we split the implementation for 1-D, 2-D and 3-D direction vectors in order to select a suitable set of rays of the same length for the inner most loop in each case. If space, explain

that the set of rays were chosen close to each other on the  $x$  dimension to try to exploit some spacial locality.

The second performance blocker are memory accesses. Assume that the extracted region is much larger than the size of the last level cache, i.e.  $8n^3 \gg N_{L3}$ . For the first horizontal direction vector  $v_1 = (1, 0, 0)$ , a lower bound for the data read from memory is  $Q_1(n) \geq \frac{n^3}{4}$  doubles. Recall that rays are spaced by a stride  $s = 2$ ; thus, not all data in the region is accessed by a vector. The bound of  $Q_1(n)$  is close to tight for the horizontal vector because data is accessed sequentially for each ray. Thus, when a (mandatory) cache miss occurs, the data in the whole cache line that is brought to cache will be used in the following iterations taking advantage of spacial locality (except for the first or last cache line in the ray in case of misalignment).

The second direction vector,  $v_2 = (0, 1, 0)$ , iterates through the extracted region vertically. Since we assumed a very large data-set, there is no reuse of data from the previous vector. Further, note that for the vertical vector, only half of the data in a cache line will be used in the best case because we have a stride of  $s = 2$  for consecutive rays. Since unused data in a cache line still have to be read from memory, a lower bound for the data read for the vertical vector is  $Q_2(n) \geq \frac{n^3}{2}$  doubles. The remaining vectors will behave similarly as the vertical vector. Hence, a lower bound on the total amount of data read from memory is  $Q(n) \geq \frac{n^3}{4} + \frac{12n^3}{2} = 6.25n^3$ . Recall from the cost analysis of the previous section that the total number of flops performed by MIL is  $W(n) = 13(\frac{n^3}{2} + 1) \approx 6.5n^3$ . An upper bound on the operation intensity of MIL for a large  $n$  is therefore  $I(n) \leq \frac{6.5n^3}{6.25n^3} = 1.04$  flops/double. We conclude that the implementation of the algorithm is memory bound when the data does not longer fit in the cache.

In order to improve the operational intensity and cache locality, the next optimization applied was blocking. Instead of calculating

Observe that all vectors iterate through the same region in a predefined pattern.

**Ellipsoid fitting.** moving loop invariant code, SIMD

#### 4. EXPERIMENTAL RESULTS

The experimental results are presenting here according to the algorithms.

**Region extraction.** *Experimental setup* The tests were performed on Intel i7 U7600 (Kaby Lake), 2.8 GHz without TurboBoost. The L1, L2 and L3 cache size were 32KB, 256KB, and 8MB, respectively. The highest optimization flag was used without vectorization (`-O3 -fno-vectorize-tree`). The sphere region size were ranged from  $n = 8^3$  to  $128^3$ .

*Results* As seen in Fig. 2, the loop unrolling and scalar replacement did not improve the performance for the region extraction algorithms. The compiler might already perform

well in alias checking and parallelizing the multiplications. When the extracted region size fitted in the cache, the performance was bound to 1 flops/cycle due to the instruction combinations. On the other hand, when the extracted region size was too large for the cache, the performance was bound to 0.5 flops/cycle due to the memory. The results shows that the algorithm reached around 0.8 flops/cycle in the cache and stayed around 0.2 flops/cycle in the memory.

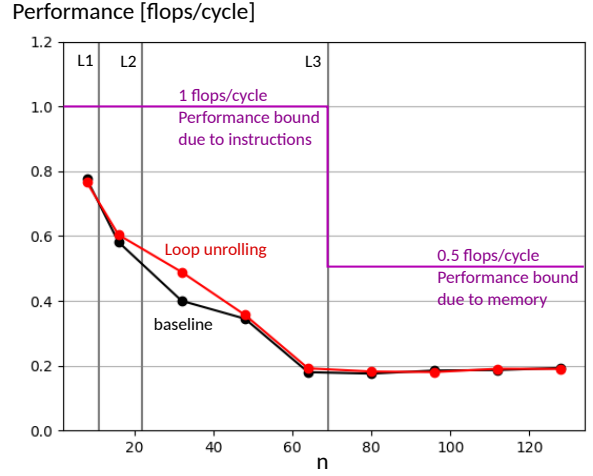


Fig. 2. Performance of region extraction

#### MIL calculation. *Experimental setup*

##### Results

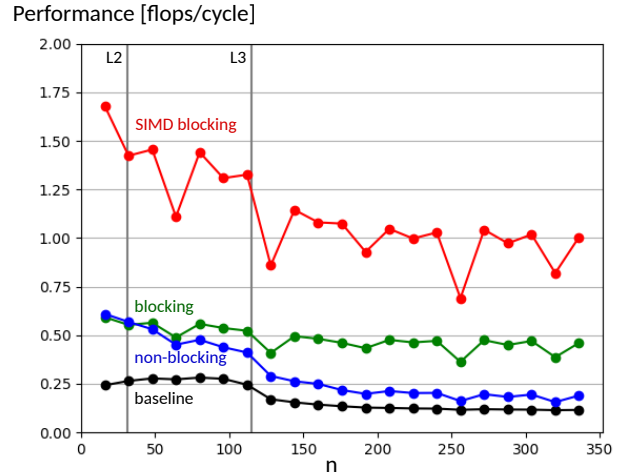
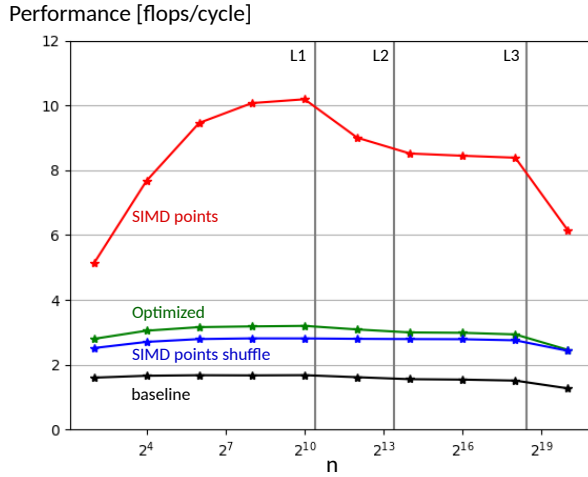


Fig. 3. Performance of MIL

#### Ellipsoid fitting. *Experimental setup*

##### Results



**Fig. 4.** Performance of ellipsoid

**Overall performance.**

## 5. CONCLUSIONS

The bone anisotropy mapping was integrated and optimized. MIL calculation consumes the large percentage of the overall computation time, followed by ellipsoid fitting and region extraction.

The performance of region extraction was bounded by the instructions mix when it was in cache and memory bound when in memory. The loop unrolling did not improve the performance and the compiler already optimized the alias checking in the basic implementation.

MIL calculation

Ellipsoid fitting

## 6. FURTHER COMMENTS

Here we provide some further tips.

**Further general guidelines.**