



南開大學  
Nankai University

计算机学院  
并行程序设计实验报告

**SIMD 编程实验—ANNS**

姓名：杨宇翔

学号：2312506

专业：计算机科学与技术

2025 年 4 月 24 日

# 目录

<b>1</b>	<b>算法设计</b>	<b>2</b>
1.1	暴力算法 flat	2
1.2	SQ	2
1.2.1	标量量化	2
1.2.2	欧式距离并行加速	4
1.2.3	SQ 搜索	5
1.3	PQ	5
1.3.1	切分量化	5
1.3.2	PQ 搜索	6
1.4	IVF-PQ	7
1.4.1	IVF 粗聚类	8
1.4.2	欧氏距离并行加速	8
1.4.3	IVF-PQ 搜索	9
<b>2</b>	<b>性能分析</b>	<b>11</b>
2.1	SDC 和 ADC	11
2.2	k_top	12
2.3	IVF-PQ 候选簇数量	12
2.4	其他参数	13
<b>3</b>	<b>总结</b>	<b>13</b>

# 1 算法设计

## 1.1 暴力算法 flat

直接计算查询向量和其他所有向量的内积距离，取最小的  $k$  个作为返回结果。内积距离的计算我们可以用 neon 并行进行优化。将原向量内积拆成 4 个部分，同时计算 4 个内积和，最后水平求和，得到总的内积和。

向量内积距离

```

1 float get_query_dis(float *query1, float *query2, int vecdim){
2     float32x4_t sum_vec = vdupq_n_f32(0);
3     for(int k=0;k<vecdim;k+=4){
4         float32x4_t vec_1 = vld1q_f32(query1+k);
5         float32x4_t vec_2 = vld1q_f32(query2+k);
6         float32x4_t mul_vec = vmulq_f32(vec_1,vec_2);
7         sum_vec = vaddq_f32(sum_vec,mul_vec);
8     }
9     float tmp[4];
10    vst1q_f32(tmp,sum_vec);
11    float dis = tmp[0]+tmp[1]+tmp[2]+tmp[3];
12    dis = 1-dis;
13    return dis;
14 }
```

## 1.2 SQ

我们已经把原始的内积和拆成 4 个部分并行计算了，能不能拆得再多一点呢？注意到原本一个 float 占 32 位，一个 128 位的向量寄存器只能存 4 个，即只能 4 位并行。于是我们可以考虑把 float 变成只占 8 位的 uint8\_t 类型，这样子一个向量寄存器就可以存 16 个，即 16 位并行，进一步加速计算。尽管这样会丢失精度，但在接受范围内。

### 1.2.1 标量量化

将原本向量里的所有值都映射到 [0,255] 这个范围内，这样就可以用 uint8\_t 来装了。具体映射方法就是确定所有向量原始值所在的区间，即确定最大值和最小值，用 255 除以这个区间长度从而得到缩放因子 scale。同时减去最小值得使负数映射到正区间。映射函数为  $y = \text{scale} * (x - \text{minval})$ 。

这个过程同样可以用 neon 并行大幅加速计算。

标量量化

```

1 void sq_compress_base(float* base, uint8_t* compressed_base, size_t base_number,
2     size_t vecdim){
3     int byte_number = base_number*vecdim;
4
5     float32x4_t min_vec = vdupq_n_f32(INFINITY);
6     float32x4_t max_vec = vdupq_n_f32(-INFINITY);
7
8     for(int i=0;i+4<=byte_number;i+=4){
```

```

8         float32x4_t vec = vld1q_f32(base+i);
9         min_vec = vminq_f32(min_vec, vec);
10        max_vec = vmaxq_f32(max_vec, vec);
11    }
12    float32x4_t temp1 = vminq_f32(min_vec, vrev64q_f32(min_vec));
13    float32x2_t temp2 = vpmin_f32(vget_low_f32(temp1), vget_high_f32(temp1));
14    float min_val = vget_lane_f32(temp2, 0);
15    temp1 = vmaxq_f32(max_vec, vrev64q_f32(max_vec));
16    temp2 = vpmax_f32(vget_low_f32(temp1), vget_high_f32(temp1));
17    float max_val = vget_lane_f32(temp2, 0);
18    for(int i=(byte_number/4)*4; i<byte_number; i++){
19        if(base[i]<min_val) min_val = base[i];
20        if(base[i]>max_val) max_val = base[i];
21    }
22
23    if(min_val==max_val){
24        uint8x16_t zero_vec = vdupq_n_u8(0);
25        for(int i=0; i<=byte_number-16; i+=16){
26            vst1q_u8(compressed_base+i, zero_vec);
27        }
28        for(int i=(byte_number/16)*16; i<byte_number; i++){
29            compressed_base[i] = 0;
30        }
31        return;
32    }
33
34    float scale = 255.0f/(max_val-min_val);
35    float32x4_t vmin = vdupq_n_f32(min_val);
36    float32x4_t vscale = vdupq_n_f32(scale);
37
38    for(int i=0; i+16<=byte_number; i+=16){
39        float32x4_t x0 = vld1q_f32(base+i);
40        float32x4_t x1 = vld1q_f32(base+i+4);
41        float32x4_t x2 = vld1q_f32(base+i+8);
42        float32x4_t x3 = vld1q_f32(base+i+12);
43
44        x0 = vsubq_f32(x0, vmin);
45        x1 = vsubq_f32(x1, vmin);
46        x2 = vsubq_f32(x2, vmin);
47        x3 = vsubq_f32(x3, vmin);
48        x0 = vmulq_f32(x0, vscale);
49        x1 = vmulq_f32(x1, vscale);
50        x2 = vmulq_f32(x2, vscale);
51        x3 = vmulq_f32(x3, vscale);
52        x0 = vmaxq_f32(vminq_f32(x0, vdupq_n_f32(255.0f)), vdupq_n_f32(0.0f));
53        x1 = vmaxq_f32(vminq_f32(x1, vdupq_n_f32(255.0f)), vdupq_n_f32(0.0f));
54        x2 = vmaxq_f32(vminq_f32(x2, vdupq_n_f32(255.0f)), vdupq_n_f32(0.0f));
55        x3 = vmaxq_f32(vminq_f32(x3, vdupq_n_f32(255.0f)), vdupq_n_f32(0.0f));
56        int32x4_t y0 = vcvtqnq_s32_f32(x0);

```

```

57     int32x4_t y1 = vcvtmq_s32_f32(x1);
58     int32x4_t y2 = vcvtmq_s32_f32(x2);
59     int32x4_t y3 = vcvtmq_s32_f32(x3);
60     int16x4_t y0_16 = vmovn_s32(y0);
61     int16x4_t y1_16 = vmovn_s32(y1);
62     int16x4_t y2_16 = vmovn_s32(y2);
63     int16x4_t y3_16 = vmovn_s32(y3);
64     int16x8_t y01 = vcombine_s16(y0_16, y1_16);
65     int16x8_t y23 = vcombine_s16(y2_16, y3_16);
66     uint8x8_t y01_u8 = vqmovun_s16(y01);
67     uint8x8_t y23_u8 = vqmovun_s16(y23);
68     vst1q_u8(compressed_base+i, vcombine_u8(y01_u8, y23_u8));
69 }
70 for (int i=(byte_number/16)*16; i<byte_number; i++){
71     float x = base[i];
72     float y = (x-min_val)*scale;
73     y = std::max(0.0f, std::min(y, 255.0f));
74     compressed_base[i] = static_cast<uint8_t>(y+0.5f);
75 }
76 }

```

### 1.2.2 欧式距离并行加速

由于原向量经过了标量量化，显然没有办法再用内积距离。所以我们直接用欧式距离计算。这个过程用 neon 并行算法优化。

#### 欧氏距离并行加速

```

1 float sq_l2_distance(uint8_t* query1, uint8_t* query2, size_t vecdim){//平方距离
2     uint32x4_t sum_vec = vdupq_n_u32(0);
3     for (int i=0; i<vecdim; i+=16){
4         uint8x16_t x = vld1q_u8(query1+i);
5         uint8x16_t y = vld1q_u8(query2+i);
6         uint8x16_t diff = vabdq_u8(x, y);
7         uint8x8_t diff_lo = vget_low_u8(diff);
8         uint8x8_t diff_hi = vget_high_u8(diff);
9         uint16x8_t diff2_lo = vmull_u8(diff_lo, diff_lo);
10        uint16x8_t diff2_hi = vmull_u8(diff_hi, diff_hi);
11
12        sum_vec = vaddq_u32(sum_vec, vpaddlq_u16(diff2_lo));
13        sum_vec = vaddq_u32(sum_vec, vpaddlq_u16(diff2_hi));
14    }
15    uint32_t sum = vaddvq_u32(sum_vec);
16    return sum;
17 }

```

### 1.2.3 SQ 搜索

本质和暴力算法相同，还是遍历所有向量算出距离，得到  $k\_top$  个向量。但速度有显著的提升。

#### 平凡算法

```

1  std::priority_queue<std::pair<float, uint32_t>> sq_search(uint8_t* base, uint8_t*
    query, size_t base_number, size_t vecdim, size_t k) {
2      std::priority_queue<std::pair<float, uint32_t>> q;
3
4      for(int i = 0; i < base_number; ++i) {
5          float dis = 0;
6          uint8_t* base_query = base+i*vecdim;
7
8          dis = sq_l2_distance(query, base_query, vecdim);
9
10         if(q.size() < k) {
11             q.push({dis, i});
12         } else {
13             if(dis < q.top().first) {
14                 q.push({dis, i});
15                 q.pop();
16             }
17         }
18     }
19     return q;
20 }
```

## 1.3 PQ

虽然 SQ 将原向量压缩成 `uint8_t` 后并行效率得到了很大提高，但遍历所有向量的做法还是太暴力了。如果向量规模继续扩大，遍历千万甚至亿的数量级显然不现实。于是 PQ 提出了将原向量切割成多个子空间，分别聚类，把向量替换成其所属簇中心的方法。

我们将原向量空间切成 4 个子空间，各自聚类得到 256 个类中心，总共 1024 个中心。当我们要计算查询向量与某个向量的距离时，我们只用直接计算到它所属中心的距离就可以了。

文档中提出的对查询向量进行同样的切分和量化编码，于是查询距离被替换成中心与中心的距离，即 SDC 的做法，显然是没道理的。首先查询向量没办法提前切分量化，而量化的过程又避免不了遍历一遍所有聚类中心并计算查询向量与聚类中心的距离。那为什么不直接拿这个距离用呢？还能避免查询向量量化产生的误差。所以我使用的是直接构建查询向量与聚类中心的距离表的方法，即 ADC。当然 SDC 的代码我也保留了，提前计算了类中心之间的距离，后续性能分析上会有比较。

### 1.3.1 切分量化

`kmeans` 函数省略，切分量化函数也有部分省略。由于是离线计算，所以没有特地做并行优化。

#### 切分量化

```

1 void kmeans_compress_base(float *base, size_t base_number, size_t vecdim){
2     const int K = 256;
3     uint8_t *compressed_base = new uint8_t[base_number*4];
4     float *centers = new float[vecdim*K];
5
6     for(int i=0;i<4;i++){
7         float *tmp_base = new float[base_number*vecdim/4];
8         for(int j=0;j<base_number;j++){
9             for(int k=0;k<vecdim/4;k++){
10                 tmp_base[j*vecdim/4+k] = base[i*vecdim/4+j*vecdim+k];
11             }
12         }
13
14         uint8_t *tmp_labels = new uint8_t[base_number];
15         float* tmp_centers = new float[vecdim/4*K];
16
17         kmeans(tmp_base, tmp_labels, tmp_centers, base_number, vecdim/4, K);
18
19         for(int j=0;j<base_number;j++){
20             compressed_base[j*4+i] = tmp_labels[j];
21         }
22
23         for(int j=0;j<K;j++){
24             for(int k=0;k<vecdim/4;k++){
25                 centers[(i*K + j)*vecdim/4 + k] = tmp_centers[j*vecdim/4 + k];
26             }
27         }
28         delete[] tmp_base;
29         delete[] tmp_labels;
30         delete[] tmp_centers;
31     }
32 }

```

### 1.3.2 PQ 搜索

注释掉的即为 SDC。在计算查询向量与聚类中心的距离时进行了 neon 并行优化。

#### PQ 搜索

```

1 std::priority_queue<std::pair<float, uint32_t>> pq_search(uint8_t* compressed_base,
2     float *centers, float *centers_dis, float* query, size_t base_number, size_t
3     vecdim, size_t k_top) {
4     const int K=256;
5     // int query_labels[4];
6     float query_centers_dis[4][256];
7
8     for(int i=0;i<4;i++){
9         float min_dis = INFINITY;

```

```

8      for(int j=0;j<256;j++){
9          float *center = centers+i*K*vecdim/4+j*vecdim/4;
10         float *sub_query = query+i*vecdim/4;
11         float32x4_t sum_vec = vdupq_n_f32(0);
12         for(int k=0;k<vecdim/4;k+=4){
13             float32x4_t vec_1 = vld1q_f32(center+k);
14             float32x4_t vec_2 = vld1q_f32(sub_query+k);
15             float32x4_t diff = vsubq_f32(vec_1, vec_2);
16             float32x4_t squared_diff = vmulq_f32(diff, diff);
17             sum_vec = vaddq_f32(sum_vec, squared_diff);
18         }
19         float sum = vaddvq_f32(sum_vec);
20         // if(dis<min_dis){
21         //     min_dis = dis;
22         //     query_labels[i] = j;
23         // }
24         query_centers_dis[i][j] = sum;
25     }
26 }
27
28 std::priority_queue<std::pair<float, uint32_t>> q;
29
30 for(int i = 0; i < base_number; i++) {
31     float dis = 0;
32     for(int j=0;j<4;j++){
33         // dis+=centers_dis[j*256*256+compressed_base[i*4+j]*256+query_labels[j]];
34         dis+=query_centers_dis[j][compressed_base[i*4+j]];
35     }
36     if(q.size() < k_top) {
37         q.push({dis, i});
38     } else {
39         if(dis < q.top().first) {
40             q.push({dis, i});
41             q.pop();
42         }
43     }
44 }
45 return q;
46 }

```

## 1.4 IVF-PQ

在原本的 PQ 基础上，我们加上了一层粗聚类，即 IVF 操作。通过粗聚类和倒排表，我们可以快速锁定查询向量相关的搜索范围，而不用像 PQ 一样遍历所有的聚类中心。

除此之外，IVF-PQ 在进行 PQ 量化之前，还将所有向量减去其所属的粗聚类中心，得到残差，最后对这个残差进行 PQ 量化。残差的动态范围更小，能更精确地表示向量间的局部差异，显著提高了



正确率。

#### 1.4.1 IVF 粗聚类

IVF 就是在切分量化之前先做一遍聚类，构建倒排表并计算残差。这些全部保存成二进制文件供后续处理。

##### IVF 粗聚类

```

1 void ivfpq_compress_base(float *base, uint8_t* ivf_base, float* ivf_centers, size_t
   base_number, size_t vecdim){
2     const int K = 256; // 聚类中心数
3     uint8_t *compressed_base = new uint8_t[base_number*4];
4     float *centers = new float[vecdim*K];
5
6     for(int i=0; i<4; i++){
7         float *tmp_base = new float[base_number*vecdim/4];
8         for(int j=0; j<base_number; j++){
9             for(int k=0; k<vecdim/4; k++){
10                tmp_base[j*vecdim/4+k] =
11                    base[j*vecdim+i*(vecdim/4)+k]-ivf_centers[ivf_base[j]*vecdim+i*vecdim/4+k];
12            }
13        }
14        uint8_t *tmp_labels = new uint8_t[base_number];
15        float* tmp_centers = new float[vecdim/4*K];
16
17        kmeans(tmp_base, tmp_labels, tmp_centers, base_number, vecdim/4, K);
18
19        for(int j=0; j<base_number; j++){
20            compressed_base[j*4+i] = tmp_labels[j];
21        }
22
23        for(int j=0; j<K; j++){
24            for(int k=0; k<vecdim/4; k++){
25                centers[(i*K + j)*vecdim/4 + k] = tmp_centers[j*vecdim/4 + k];
26            }
27        }
28        delete[] tmp_base;
29        delete[] tmp_labels;
30        delete[] tmp_centers;
31    }
32 }
```

#### 1.4.2 欧氏距离并行加速

虽然之前写过一个欧式距离了，但那个是针对 uint8\_t 类型的并行加速，float 类型还是得重新写一个。

## 欧式距离并行加速

```

1 float get_query_dis(float *query1, float *query2, int vecdim){
2     float32x4_t sum_vec = vdupq_n_f32(0);
3     for(int k=0;k<vecdim;k+=4){
4         float32x4_t vec_1 = vld1q_f32(query1+k);
5         float32x4_t vec_2 = vld1q_f32(query2+k);
6         float32x4_t diff = vsubq_f32(vec_1, vec_2);
7         float32x4_t squared_diff = vmulq_f32(diff, diff);
8         sum_vec = vaddq_f32(sum_vec, squared_diff);
9     }
10    float sum = vaddvq_f32(sum_vec);
11    return sum;
12 }

```

## 1.4.3 IVF-PQ 搜索

中间经过一系列处理后，我们最后拿到了粗聚类的 ivf\_base 信息，倒排表，以及残差 PQ 量化后的 pqivf\_base 信息。有了这些后我们就可以进行搜索了。拿到查询向量后遍历一遍粗聚类中心，得到 m 个相距最近的簇，然后遍历这些簇中的向量，计算查询向量与簇中心的残差到候选向量的残差的距离，从而得到 k\_top 的向量。

其中计算查询向量与簇中心的残差部分也可以做 neon 并行加速。

## IVFPQ 搜索

```

1 std::priority_queue<std::pair<float, uint32_t>> ivfpq_search(uint8_t*
2     compressed_base, float *centers, float *ivf_centers, std::vector<int>
3     *inverted_kmeans, float* query, size_t base_number, size_t vecdim, size_t k_top) {
4     const int K=256;
5
6     int m=16;
7     int *query_ivf_labels = new int[m];
8
9     std::priority_queue<std::pair<float, uint32_t>> ivf_q;
10
11     for(int i = 0; i < 256; i++) {
12         float dis = get_query_dis(ivf_centers+i*vecdim, query, vecdim);
13         if(ivf_q.size() < m) {
14             ivf_q.push({dis, i});
15         } else {
16             if(dis < ivf_q.top().first) {
17                 ivf_q.push({dis, i});
18                 ivf_q.pop();
19             }
20         }
21     }
22
23     for(int i=0;i<m;i++){
24         query_ivf_labels[i] = ivf_q.top().second;
25     }
26 }

```

```

23     ivf_q.pop();
24 }
25
26 int *size = new int[m];
27 int **candidate_base = new int*[m];
28 float **residual = new float*[m];
29 for(int i=0; i<m; i++){
30     size[i] = inverted_kmeans[query_ivf_labels[i]].size();
31     candidate_base[i] = inverted_kmeans[query_ivf_labels[i]].data();
32     residual[i] = new float[vecdim];
33
34     float* center_ptr = &ivf_centers[query_ivf_labels[i]*vecdim];
35     float* residual_ptr = residual[i];
36
37     for(int j=0; j<vecdim; j+=4){
38         float32x4_t center = vld1q_f32(center_ptr + j);
39         float32x4_t qvec = vld1q_f32(query + j);
40         float32x4_t residual_vec = vsubq_f32(qvec, center);
41         vst1q_f32(residual_ptr + j, residual_vec);
42     }
43 }
44
45 std::priority_queue<std::pair<float, uint32_t>> q;
46
47 for(int i = 0; i < m; i++) {
48     for(int j=0;j<size[i];j++){
49         int candidate_query = candidate_base[i][j];
50         float dis = 0;
51         for(int k=0;k<4;k++){
52             dis+=get_query_dis(residual[i]+k*vecdim/4,centers+k*256*vecdim/4+compressed_base[candidate_query]*vecdim/4);
53         }
54         if(q.size() < k_top) {
55             q.push({dis, candidate_query});
56         } else {
57             if(dis < q.top().first) {
58                 q.push({dis, candidate_query});
59                 q.pop();
60             }
61         }
62     }
63 }
64 return q;
65 }

```

## 2 性能分析

所有实验结果都是在服务器上取得的。其中 PQ 的正确率只有 30%，IVFPQ 也只有 40%，说实话我也觉得太低了。但我已经被这个正确率硬控 5 天了，代码怎么也找不出来哪里出了问题，我只能默认这个结果是正常的了。

latency 取 10 次实验结果的平均值。默认 test\_number 为 2000。k\_top 默认为 100。PQ 默认切割 4 个子空间，每个子空间 256 个聚类中心，采用 ADC。IVF-PQ 默认 256 个粗聚类中心，残差 PQ 结构与常规 PQ 结构相同，默认候选簇数量为 8。

从测试结果上可以看到，SQ 因为被查询向量的标量量化操作拖慢了速度，并没有和经过 neon 并行优化后的 flat 拉开太大的差距。PQ 虽然速度上有了很大的提高，但准确率大大降低。IVF-PQ 则是在速度进一步提升的同时更好地维持了准确率。

表 1: 各算法性能

指标	flat	SQ	PQ	IVF-PQ
recall@100	1.000000	0.809763	0.315720	0.40342
latency ( $\mu s$ )	6300.170	3022.926	571.951	248.692

```
load data /anndata/DEEP100K.query.fbin
dimension: 96 number:10000 size_per_element:4
load data /anndata/DEEP100K.gt.query.100k.top100.bin
dimension: 100 number:10000 size_per_element:4
load data /anndata/DEEP100K.base.100k.fbin
dimension: 96 number:100000 size_per_element:4
average recall: 1
average latency (us): 5425.13
```

图 2.1: flat

```
load data /anndata/DEEP100K.query.fbin
dimension: 96 number:10000 size_per_element:4
load data /anndata/DEEP100K.gt.query.100k.top100.bin
dimension: 100 number:10000 size_per_element:4
load data /anndata/DEEP100K.base.100k.fbin
dimension: 96 number:100000 size_per_element:4
成功读入sq_compressed_base.bin!
average recall: 0.809763
average latency (us): 3031.12
```

图 2.2: SQ

```
load data /anndata/DEEP100K.query.fbin
dimension: 96 number:10000 size_per_element:4
load data /anndata/DEEP100K.gt.query.100k.top100.bin
dimension: 100 number:10000 size_per_element:4
load data /anndata/DEEP100K.base.100k.fbin
dimension: 96 number:100000 size_per_element:4
成功读入pq_kmeans_base.bin!
成功读入pq_kmeans_centers.bin!
成功读入pq_kmeans_centers_distance.bin!
average recall: 0.31572
average latency (us): 571.253
```

图 2.3: PQ

```
load data /anndata/DEEP100K.query.fbin
dimension: 96 number:10000 size_per_element:4
load data /anndata/DEEP100K.gt.query.100k.top100.bin
dimension: 100 number:10000 size_per_element:4
load data /anndata/DEEP100K.base.100k.fbin
dimension: 96 number:100000 size_per_element:4
成功读入inverted_kmeans.bin!
成功读入ivf_centers.bin!
成功读入ivf_base.bin!
成功读入ivfpq_kmeans_base.bin!
成功读入ivfpq_kmeans_centers.bin!
average recall: 0.40342
average latency (us): 258.848
```

图 2.4: IVFPQ

### 2.1 SDC 和 ADC

计算查询向量与原始向量距离时，SDC 计算的是两个类中心之间的距离，而 ADC 计算的是查询向量与类中心的距离。显然我们这次实验的测试查询向量没办法提前量化，量化的时间应该纳入 latency 中。所以使用了 ADC。二者的性能我也做了对比。

可以看到 SDC 在正确率上远低于 ADC。甚至反而因为 SDC 需要大量偏移访问 pq\_kmeans\_centers 里的数据，速度严重落后于 ADC。

```
load data /anndata/DEEP100K.query.fbin
dimension: 96 number:10000 size_per_element:4
load data /anndata/DEEP100K.gt.query.100k.top100.bin
dimension: 100 number:10000 size_per_element:4
load data /anndata/DEEP100K.base.100k.fbin
dimension: 96 number:100000 size_per_element:4
成功读入pq_kmeans_base.bin!
成功读入pq_kmeans_centers.bin!
成功读入pq_kmeans_centers_distance.bin!
average recall: 0.22845
average latency (us): 1417.8
```

图 2.5: SDC

```
load data /anndata/DEEP100K.query.fbin
dimension: 96 number:10000 size_per_element:4
load data /anndata/DEEP100K.gt.query.100k.top100.bin
dimension: 100 number:10000 size_per_element:4
load data /anndata/DEEP100K.base.100k.fbin
dimension: 96 number:100000 size_per_element:4
成功读入pq_kmeans_base.bin!
成功读入pq_kmeans_centers.bin!
成功读入pq_kmeans_centers_distance.bin!
average recall: 0.31572
average latency (us): 571.253
```

图 2.6: ADC

## 2.2 k\_top

k\_top 的高低在 latency 上影响较小, 因为它只会影响维护最大堆消耗的时间。而无论是 SQ、PQ 还是 IVF-PQ, 它们遍历向量的数量都和 k\_top 无关。所以 latency 就不画图了。

k\_top 在 recall@k 上的影响就非常大了。毕竟 k\_top 越小, 搜索结果的范围越来越小, 召回率自然下降。

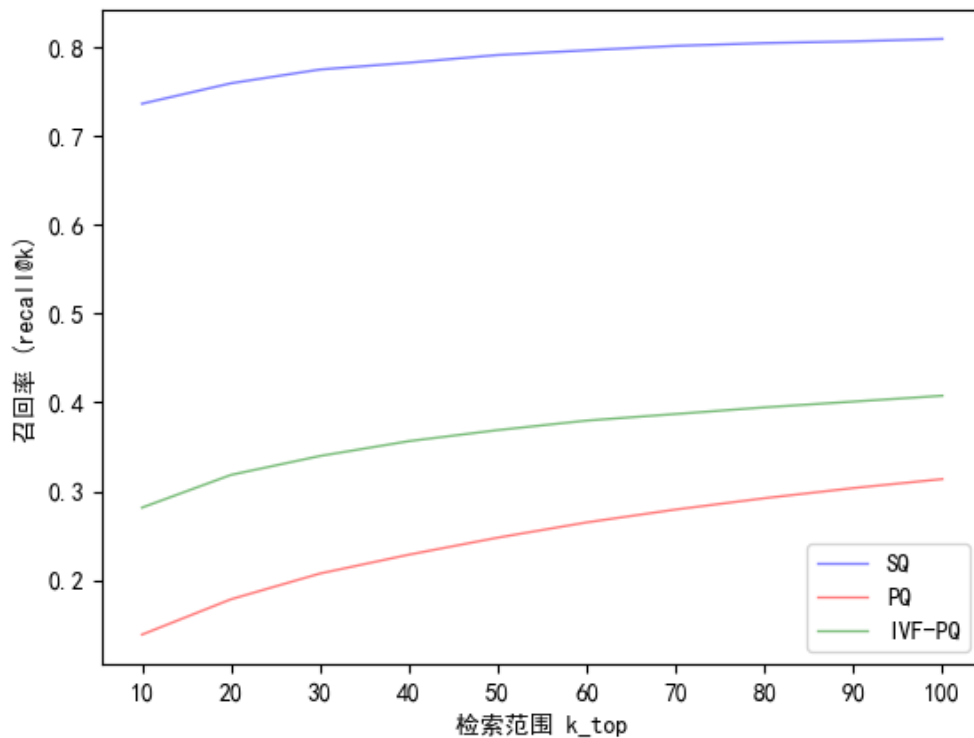


图 2.7: k\_top

## 2.3 IVF-PQ 候选簇数量

虽然 IVF-PQ 通过粗簇筛选掉了大量冗余数据, 但这同时也可能会漏掉真实的最近邻向量。假设搜索时检索 m 个候选簇, 对于这 m 个候选簇我们不但要遍历其中的所有向量分别计算距离, 而且由于我们 PQ 量化的是残差, 我们还要先分别算出查询向量和各候选簇的残差才能用来计算距离。这个

计算开销是巨大的。

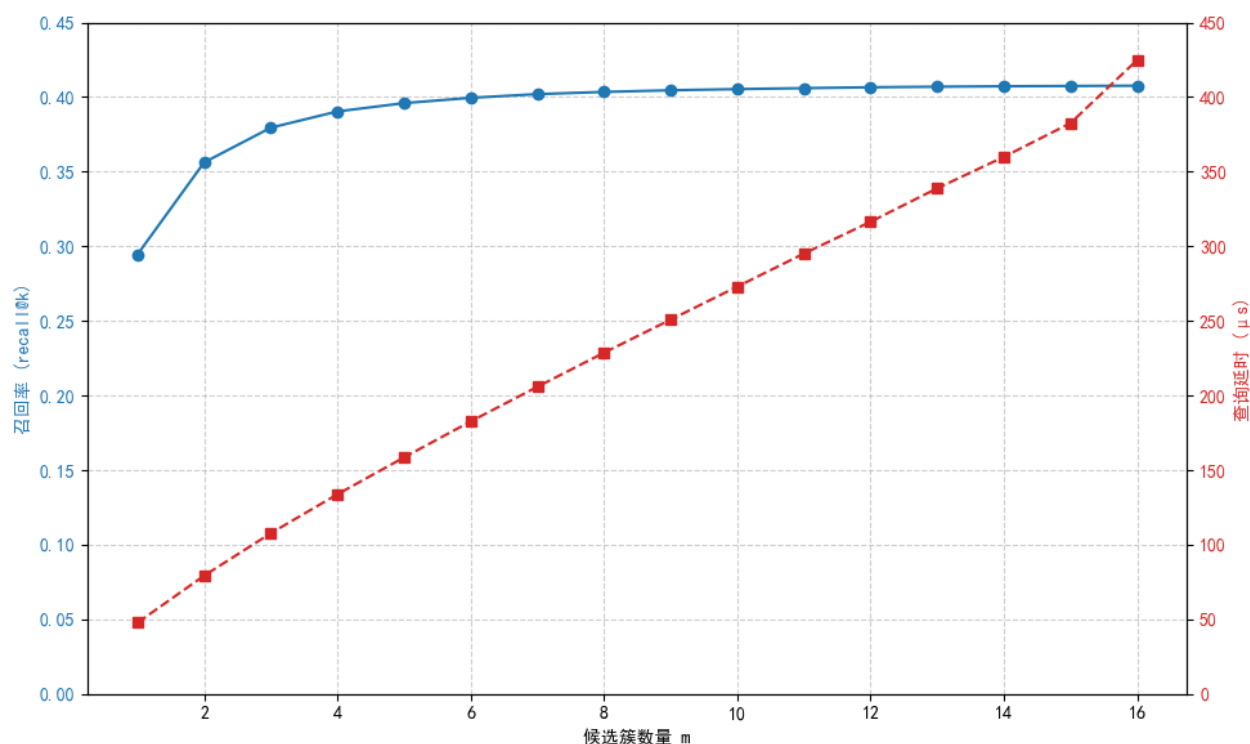


图 2.8: 候选簇数量  $m$

可以看到召回率在候选簇数量到达 4 时涨幅就不明显了，但查询时间却是线性地在增长。如果追求速度的话候选簇数量  $m$  可以设置为 4。但考虑到正确率，我选择设置的候选簇数量为 8。

## 2.4 其他参数

还有一些其他的参数比如 `test_number`，PQ 的聚类中心数量，IVF 粗聚类中心数量等需要考虑。但尝试修改测试过后发现影响很小。例如 IVF 粗聚类中心设置为 256 和 1024，对应的设置候选簇数量为 16 和 256 时，二者的 `recall@k` 几乎没有区别，后者反而因为计算开销导致 `latency` 升高。

## 3 总结

这次实验虽然是 SIMD 编程实验，但我几乎全都在忙着实现各种算法。尤其是在正确率上纠结了很长时间，期间换过距离计算公式，对聚类结果进行标准化，但完全没用。结果就是更多探索已经没时间实现了。考虑到这次实验的重点是 SIMD 编程，应该已经完成任务了吧 (?)

我用 `neon` 进行 SIMD 编程，对欧式距离、内积距离的计算，标量量化过程，搜索时的查询向量归类过程做了并行优化，实现了 `arm` 平台上查询速度的显著提高。

更多的 ANNS 算法实现和并行优化就留到期末报告里实现罢。