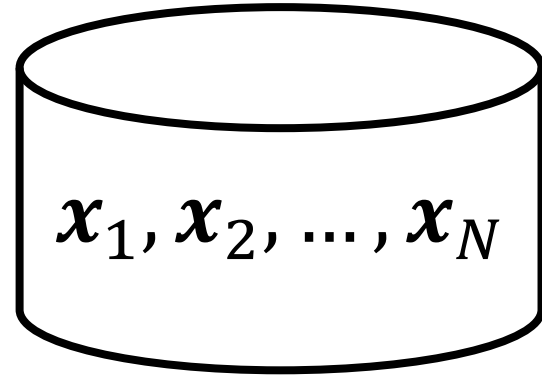


Billion-scale Approximate Nearest Neighbor Search

Yusuke Matsui
The University of Tokyo



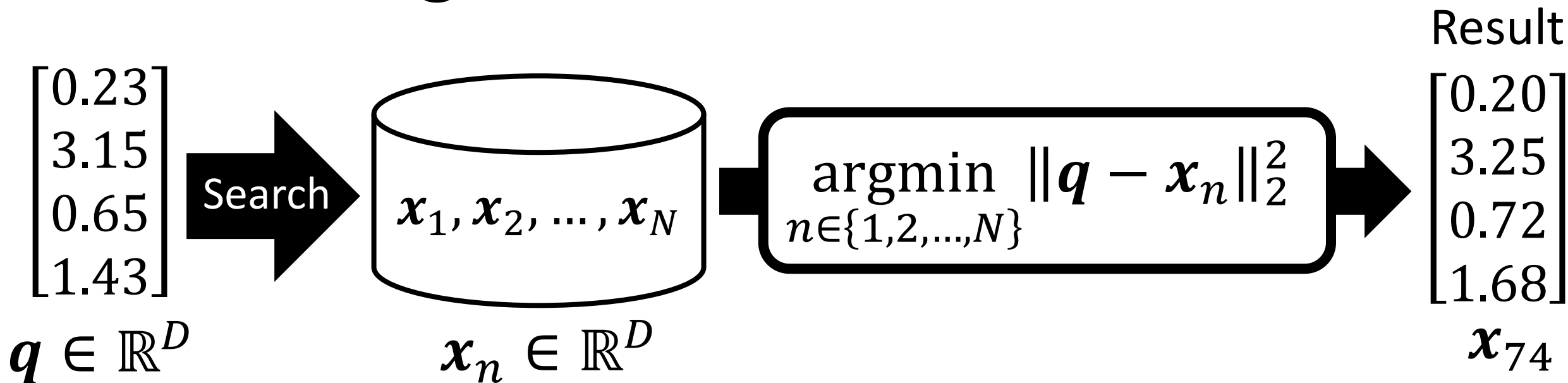
Nearest Neighbor Search; NN



$$x_n \in \mathbb{R}^D$$

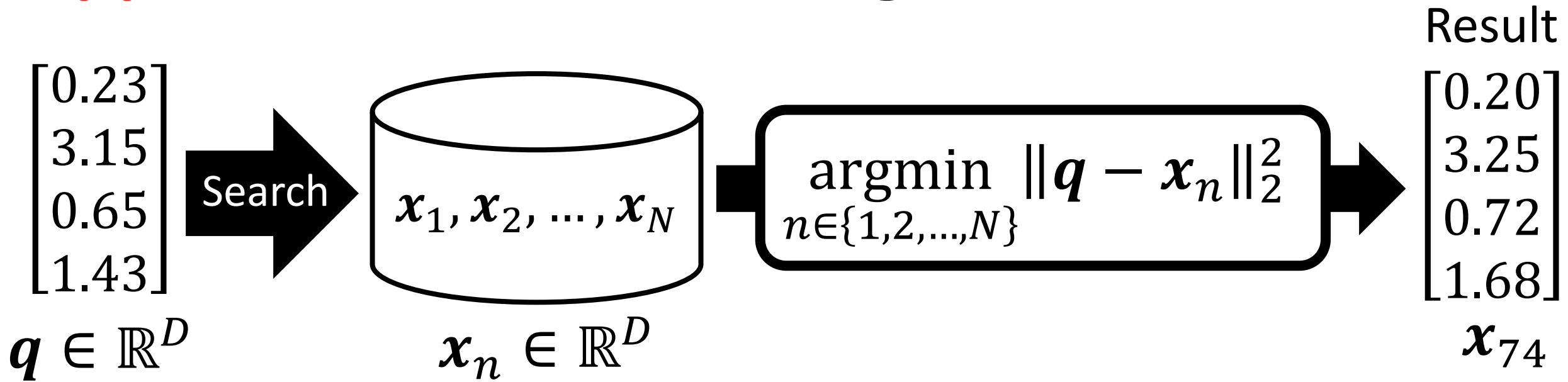
➤ N D -dim database vectors: $\{x_n\}_{n=1}^N$

Nearest Neighbor Search; NN



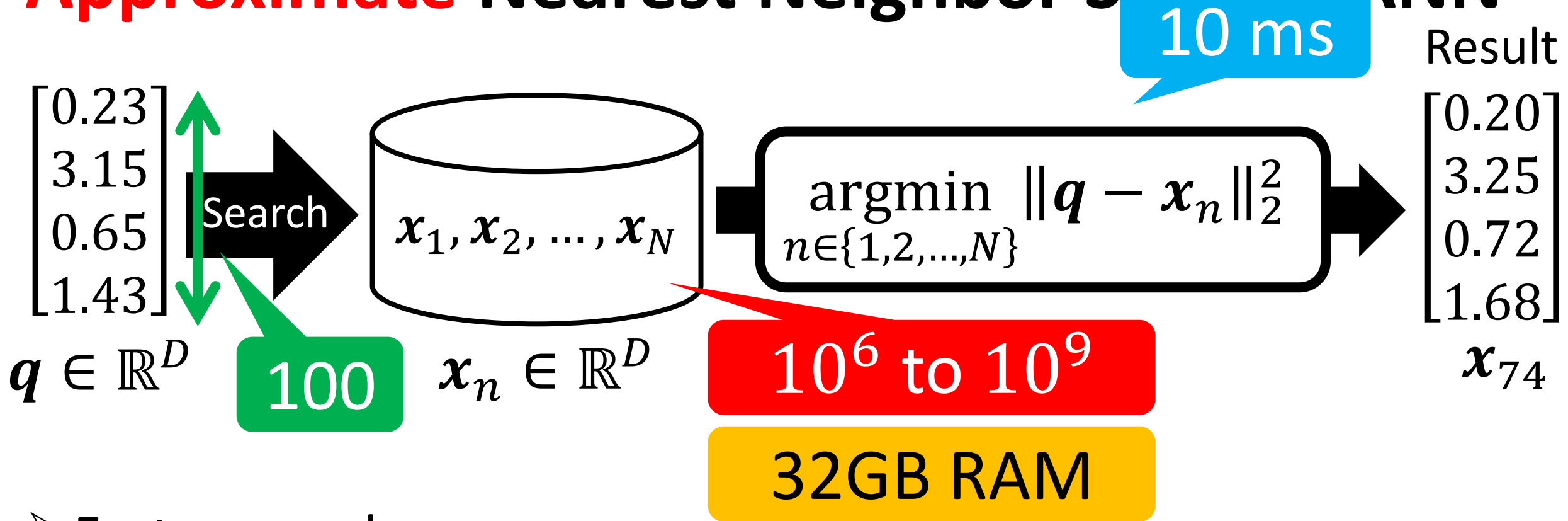
- N D -dim database vectors: $\{\mathbf{x}_n\}_{n=1}^N$
- Given a query \mathbf{q} , find the closest vector from the database
- One of the fundamental problems in computer science
- Solution: linear scan, $O(ND)$, slow 😞

Approximate Nearest Neighbor Search; ANN



- Faster search
- Don't necessarily have to be exact neighbors
- Trade off: runtime, accuracy, and memory-consumption

Approximate Nearest Neighbor Search: ANN



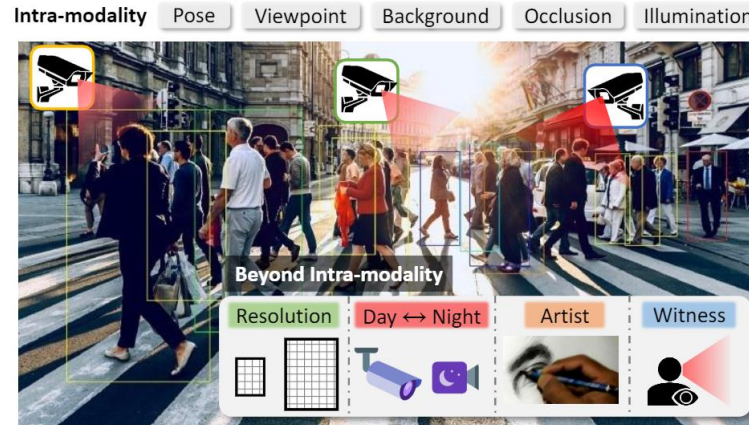
- Faster search
- Don't necessarily have to be exact neighbors
- Trade off: runtime, accuracy, and memory-consumption
- A sense of scale: billion-scale data on memory

NN/ANN for CV

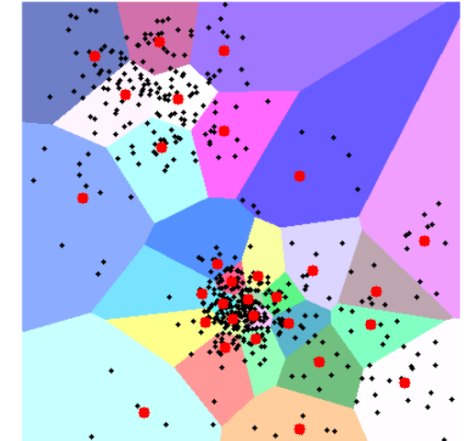
https://about.mercari.com/press/news/article/20190318_image_search/



Image retrieval

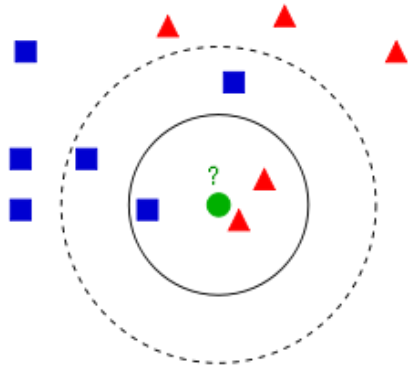


Person Re-identification



Clustering

https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm



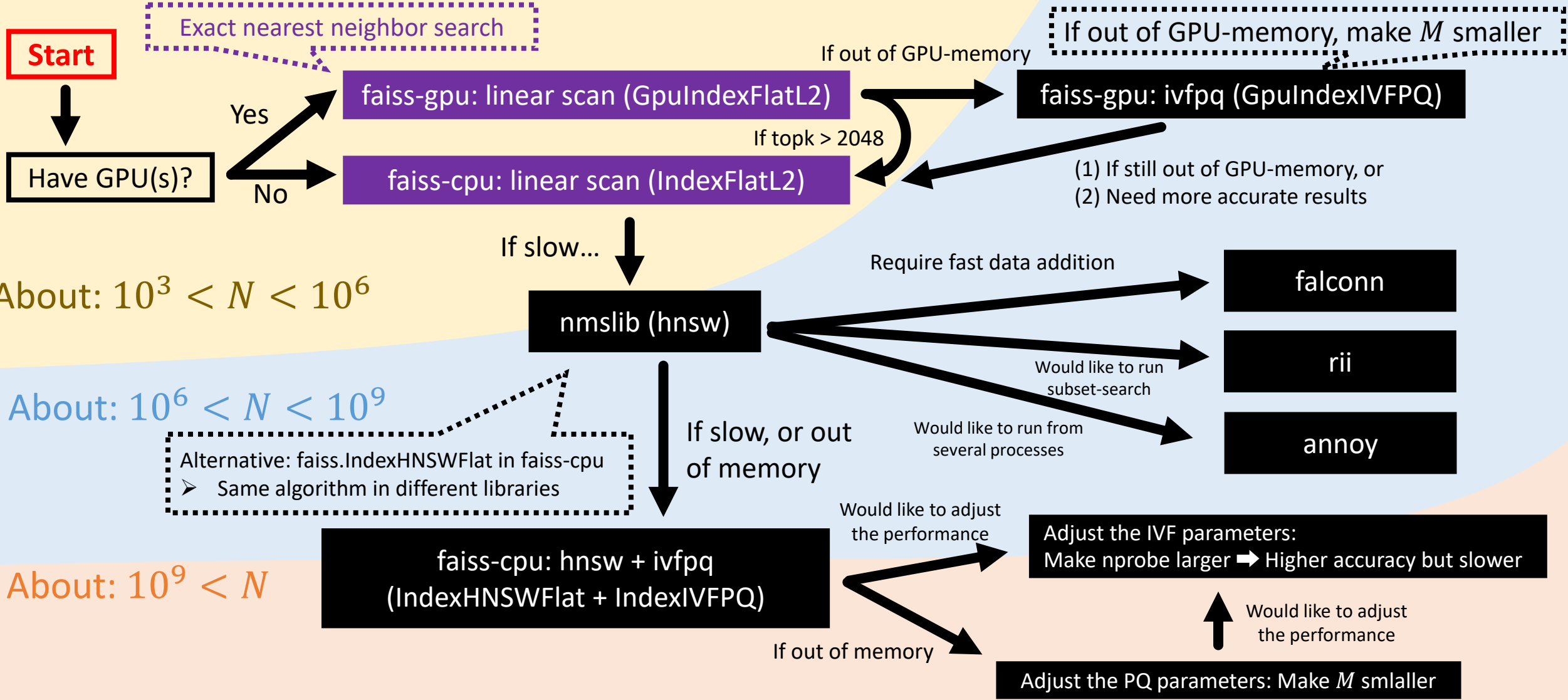
kNN recognition

<https://jp.mathworks.com/help/vision/ug/image-classification-with-bag-of-visual-words.html>



- Originally: fast construction of bag-of-features
- One of the benchmarks is still SIFT

cheat-sheet for ANN in Python (as of 2020. Can be installed by conda or pip)



About: $10^3 < N < 10^6$

About: $10^6 < N < 10^9$

About: $10^9 < N$

Note: Assuming $D \cong 100$. The size of the problem is determined by DN . If $100 \ll D$, run PCA to reduce D to 100

Part 1:

Nearest Neighbor Search

Part 2:

Approximate Nearest Neighbor Search

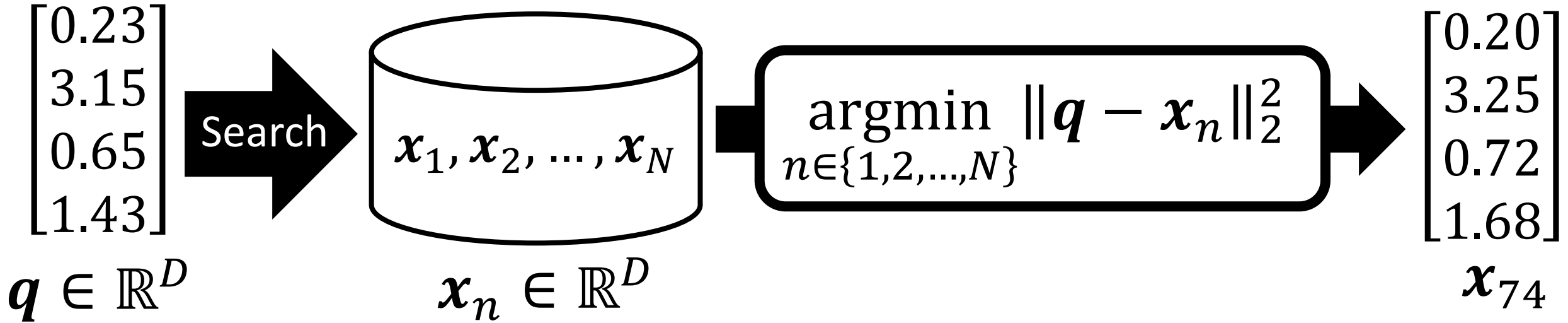
Part 1:

Nearest Neighbor Search

Part 2:

Approximate Nearest Neighbor Search

Nearest Neighbor Search



- Should try this first of all
- Introduce a naïve implementation
- Introduce a fast implementation
 - ✓ **Faiss** library from FAIR (you'll see many times today. CPU & GPU)
- Experience the drastic difference between the two impls

M D -dim query vectors $\mathcal{Q} = \{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_M\}$
 N D -dim database vectors $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ $M \ll N$
Task : Given $\mathbf{q} \in \mathcal{Q}$ and $\mathbf{x} \in \mathcal{X}$, compute $\|\mathbf{q} - \mathbf{x}\|_2^2$

M D -dim query vectors

$$Q = \{\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_M\}$$

N D -dim database vectors

$$X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \quad M \ll N$$

Task : Given $\mathbf{q} \in Q$ and $\mathbf{x} \in X$, compute $\|\mathbf{q} - \mathbf{x}\|_2^2$

Naïve impl.

```
def l2sqr(q, x):  
    diff = 0.0  
    for (d = 0; d < D; ++d):  
        diff += (q[d] - x[d])**2  
    return diff
```

```
parfor q in Q:
```

```
    for x in X:
```

```
        l2sqr(q, x)
```

Parallelize
query-side

Select min by heap,
but omit it now

M D -dim query vectors

$$Q = \{q_1, q_2, \dots, q_M\}$$

N D -dim database vectors

$$X = \{x_1, x_2, \dots, x_N\} \quad M \ll N$$

Task : Given $q \in Q$ and $x \in X$, compute $\|q - x\|_2^2$

Naïve impl.

```
def l2sqr(q, x):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (q[d] - x[d])**2
    return diff
```

```
parfor q in Q:
    for x in X:
        l2sqr(q, x)
```

Parallelize query-side

Select min by heap, but omit it now

faiss impl.

if $M < 20$:

compute $\|q - x\|_2^2$ by SIMD

else :

compute $\|q - x\|_2^2 = \|q\|_2^2 - 2q^T x + \|x\|_2^2$ by BLAS

M D -dim query vectors

$$Q = \{q_1, q_2, \dots, q_M\}$$

N D -dim database vectors

$$X = \{x_1, x_2, \dots, x_N\} \quad M \ll N$$

Task : Given $q \in Q$ and $x \in X$, compute $\|q - x\|_2^2$

Naïve impl.

```
def l2sqr(q, x):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (q[d] - x[d])**2
    return diff
```

```
parfor q in Q:
    for x in X:
        l2sqr(q, x)
```

Parallelize query-side

Select min by heap, but omit it now

faiss impl.

if $M < 20$:

compute $\|q - x\|_2^2$ by SIMD

else :

compute $\|q - x\|_2^2 = \|q\|_2^2 - 2q^T x + \|x\|_2^2$ by BLAS

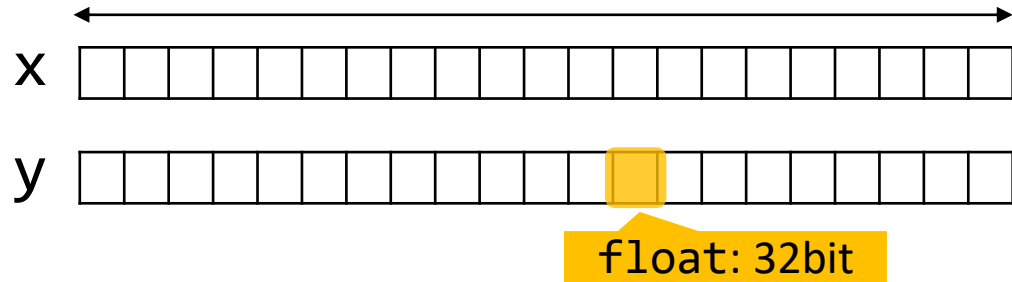
$\|x - y\|_2^2$ by SIMD

Rename variables for the sake of explanation

```
def l2sqr(x, y):  
    diff = 0.0  
    for (d = 0; d < D; ++d):  
        diff += (x[d] - y[d])**2  
    return diff
```

Ref.

D=31



```
float fvec_l2sqr (const float * x,  
                 const float * y,  
                 size_t d)  
{  
    __m256 msum1 = _mm256_setzero_ps();  
  
    while (d >= 8) {  
        __m256 mx = _mm256_loadu_ps (x); x += 8;  
        __m256 my = _mm256_loadu_ps (y); y += 8;  
        const __m256 a_m_b1 = mx - my;  
        msum1 += a_m_b1 * a_m_b1;  
        d -= 8;  
    }  
  
    __m128 msum2 = _mm256_extractf128_ps(msum1, 1);  
    msum2 += _mm256_extractf128_ps(msum1, 0);  
  
    if (d >= 4) {  
        __m128 mx = _mm_loadu_ps (x); x += 4;  
        __m128 my = _mm_loadu_ps (y); y += 4;  
        const __m128 a_m_b1 = mx - my;  
        msum2 += a_m_b1 * a_m_b1;  
        d -= 4;  
    }  
  
    if (d > 0) {  
        __m128 mx = masked_read (d, x);  
        __m128 my = masked_read (d, y);  
        __m128 a_m_b1 = mx - my;  
        msum2 += a_m_b1 * a_m_b1;  
    }  
  
    msum2 = _mm_hadd_ps (msum2, msum2);  
    msum2 = _mm_hadd_ps (msum2, msum2);  
    return _mm_cvtss_f32 (msum2);  
}
```

$\|x - y\|_2^2$ by SIMD

Rename variables for the sake of explanation

```
float fvec_L2sqr (const float * x,
                 const float * y,
                 size_t d)
{
```

```
    __m256 msum1 = _mm256_setzero_ps();
```

```
    while (d >= 8) {
        __m256 mx = _mm256_loadu_ps (x); x += 8;
        __m256 my = _mm256_loadu_ps (y); y += 8;
        const __m256 a_m_b1 = mx - my;
        msum1 += a_m_b1 * a_m_b1;
        d -= 8;
    }
```

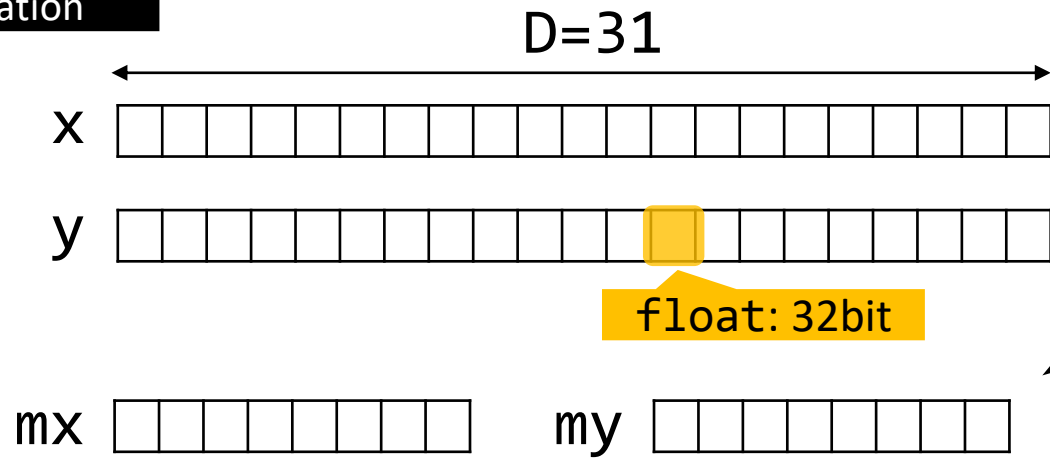
```
    __m128 msum2 = _mm256_extractf128_ps(msum1, 1);
    msum2 += _mm256_extractf128_ps(msum1, 0);
```

```
    if (d >= 4) {
        __m128 mx = _mm_loadu_ps (x); x += 4;
        __m128 my = _mm_loadu_ps (y); y += 4;
        const __m128 a_m_b1 = mx - my;
        msum2 += a_m_b1 * a_m_b1;
        d -= 4;
    }
```

```
    if (d > 0) {
        __m128 mx = masked_read (d, x);
        __m128 my = masked_read (d, y);
        __m128 a_m_b1 = mx - my;
        msum2 += a_m_b1 * a_m_b1;
    }
```

```
    msum2 = _mm_hadd_ps (msum2, msum2);
    msum2 = _mm_hadd_ps (msum2, msum2);
    return _mm_cvtss_f32 (msum2);
}
```

```
Ref.
def l2sqr(x, y):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (x[d] - y[d])**2
    return diff
```



- 256bit SIMD Register
- Process eight floats at once

Ref.

$\|x - y\|_2^2$ by SIMD

Rename variables for the sake of explanation

```
float fvec_L2sqr (const float * x,
                 const float * y,
                 size_t d)
```

```
{
    __m256 msum1 = _mm256_setzero_ps();
```

```
while (d >= 8) {
    __m256 mx = _mm256_loadu_ps (x); x += 8;
    __m256 my = _mm256_loadu_ps (y); y += 8;
    const __m256 a_m_b1 = mx - my;
    msum1 += a_m_b1 * a_m_b1;
    d -= 8;
}
```

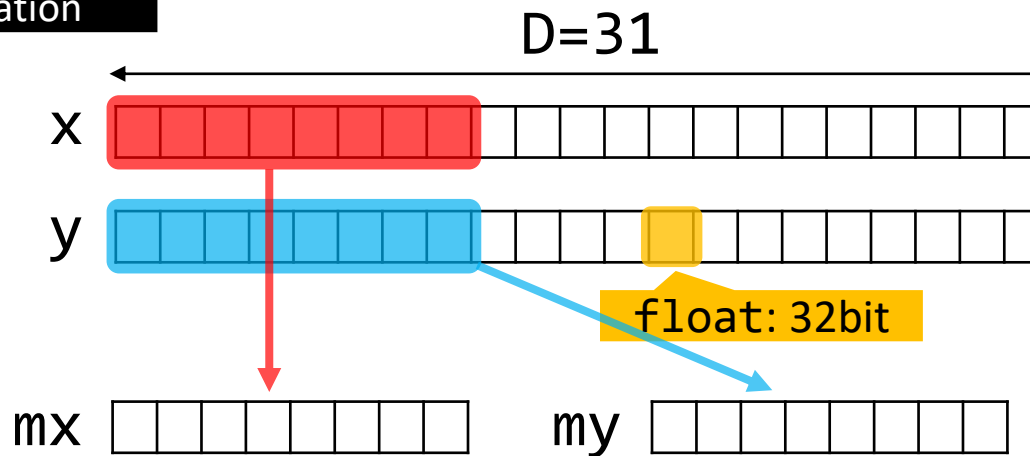
```
__m128 msum2 = _mm256_extractf128_ps(msum1, 1);
msum2 += _mm256_extractf128_ps(msum1, 0);
```

```
if (d >= 4) {
    __m128 mx = _mm_loadu_ps (x); x += 4;
    __m128 my = _mm_loadu_ps (y); y += 4;
    const __m128 a_m_b1 = mx - my;
    msum2 += a_m_b1 * a_m_b1;
    d -= 4;
}
```

```
if (d > 0) {
    __m128 mx = masked_read (d, x);
    __m128 my = masked_read (d, y);
    __m128 a_m_b1 = mx - my;
    msum2 += a_m_b1 * a_m_b1;
}
```

```
msum2 = _mm_hadd_ps (msum2, msum2);
msum2 = _mm_hadd_ps (msum2, msum2);
return _mm_cvtss_f32 (msum2);
}
```

```
Ref.
def l2sqr(x, y):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (x[d] - y[d])**2
    return diff
```



- 256bit SIMD Register
- Process eight floats at once

$\|x - y\|_2^2$ by SIMD

Rename variables for the sake of explanation

```
float fvec_L2sqr (const float * x,
                 const float * y,
                 size_t d)
{
```

```
    __m256 msum1 = _mm256_setzero_ps();
```

```
    while (d >= 8) {
        __m256 mx = _mm256_loadu_ps (x); x += 8;
        __m256 my = _mm256_loadu_ps (y); y += 8;
        const __m256 a_m_b1 = mx - my;
        msum1 += a_m_b1 * a_m_b1;
        d -= 8;
    }
```

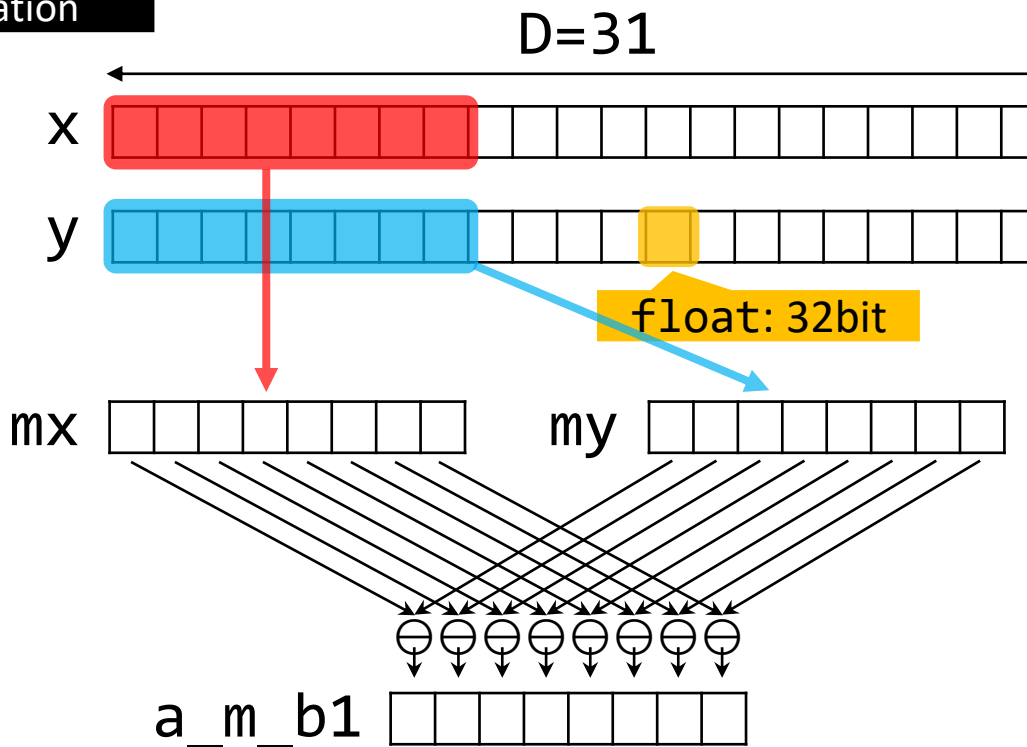
```
    __m128 msum2 = _mm256_extractf128_ps(msum1, 1);
    msum2 += _mm256_extractf128_ps(msum1, 0);
```

```
    if (d >= 4) {
        __m128 mx = _mm_loadu_ps (x); x += 4;
        __m128 my = _mm_loadu_ps (y); y += 4;
        const __m128 a_m_b1 = mx - my;
        msum2 += a_m_b1 * a_m_b1;
        d -= 4;
    }
```

```
    if (d > 0) {
        __m128 mx = masked_read (d, x);
        __m128 my = masked_read (d, y);
        __m128 a_m_b1 = mx - my;
        msum2 += a_m_b1 * a_m_b1;
    }
```

```
    msum2 = _mm_hadd_ps (msum2, msum2);
    msum2 = _mm_hadd_ps (msum2, msum2);
    return _mm_cvtss_f32 (msum2);
}
```

```
Ref.
def l2sqr(x, y):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (x[d] - y[d])**2
    return diff
```



➤ 256bit SIMD Register
➤ Process eight floats at once

$\|x - y\|_2^2$ by SIMD

Rename variables for the sake of explanation

```
float fvec_L2sqr (const float * x,
                 const float * y,
                 size_t d)
```

```
{
    __m256 msum1 = _mm256_setzero_ps();
```

```
while (d >= 8) {
    __m256 mx = _mm256_loadu_ps (x); x += 8;
    __m256 my = _mm256_loadu_ps (y); y += 8;
    const __m256 a_m_b1 = mx - my;
    msum1 += a_m_b1 * a_m_b1;
    d -= 8;
}
```

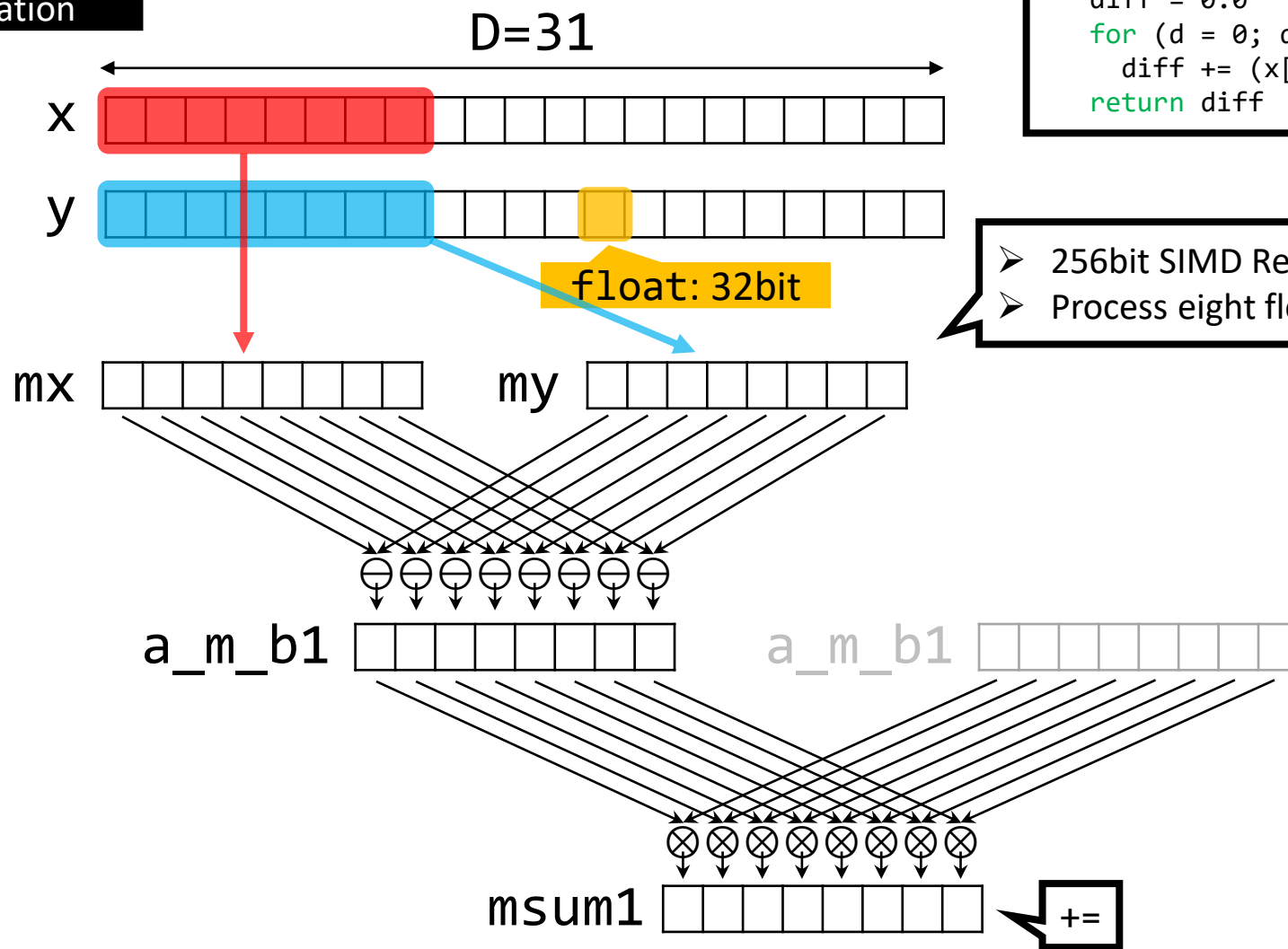
```
__m128 msum2 = _mm256_extractf128_ps(msum1, 1);
msum2 += _mm256_extractf128_ps(msum1, 0);
```

```
if (d >= 4) {
    __m128 mx = _mm_loadu_ps (x); x += 4;
    __m128 my = _mm_loadu_ps (y); y += 4;
    const __m128 a_m_b1 = mx - my;
    msum2 += a_m_b1 * a_m_b1;
    d -= 4;
}
```

```
if (d > 0) {
    __m128 mx = masked_read (d, x);
    __m128 my = masked_read (d, y);
    __m128 a_m_b1 = mx - my;
    msum2 += a_m_b1 * a_m_b1;
}
```

```
msum2 = _mm_hadd_ps (msum2, msum2);
msum2 = _mm_hadd_ps (msum2, msum2);
return _mm_cvtss_f32 (msum2);
}
```

```
Ref.
def l2sqr(x, y):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (x[d] - y[d])**2
    return diff
```



➤ 256bit SIMD Register
➤ Process eight floats at once

Ref.

$\|x - y\|_2^2$ by SIMD

Rename variables for the sake of explanation

```
float fvec_L2sqr (const float * x,
                 const float * y,
                 size_t d)
```

```
{
    __m256 msum1 = _mm256_setzero_ps();
```

```
while (d >= 8) {
    __m256 mx = _mm256_loadu_ps (x); x += 8;
    __m256 my = _mm256_loadu_ps (y); y += 8;
    const __m256 a_m_b1 = mx - my;
    msum1 += a_m_b1 * a_m_b1;
    d -= 8;
}
```

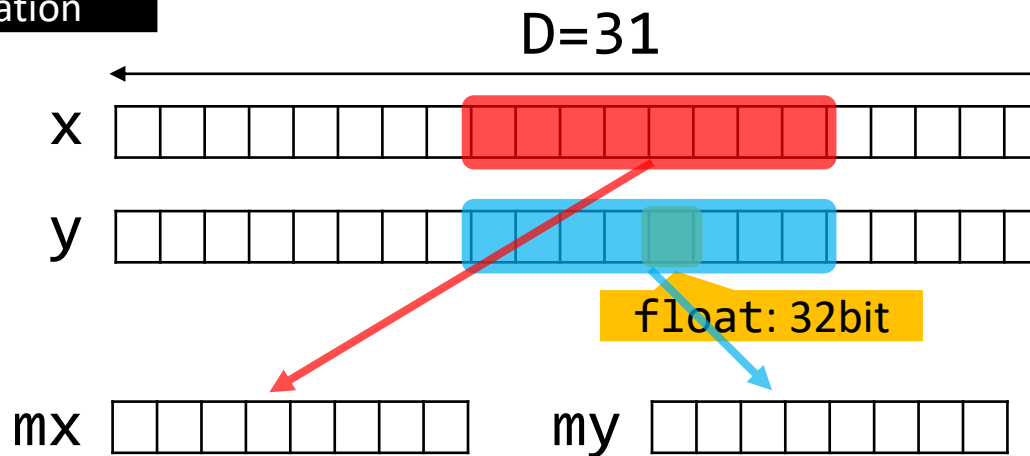
```
__m128 msum2 = _mm256_extractf128_ps(msum1, 1);
msum2 += _mm256_extractf128_ps(msum1, 0);
```

```
if (d >= 4) {
    __m128 mx = _mm_loadu_ps (x); x += 4;
    __m128 my = _mm_loadu_ps (y); y += 4;
    const __m128 a_m_b1 = mx - my;
    msum2 += a_m_b1 * a_m_b1;
    d -= 4;
}
```

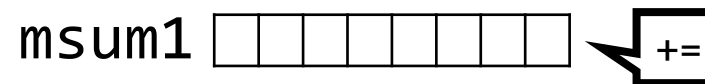
```
if (d > 0) {
    __m128 mx = masked_read (d, x);
    __m128 my = masked_read (d, y);
    __m128 a_m_b1 = mx - my;
    msum2 += a_m_b1 * a_m_b1;
}
```

```
msum2 = _mm_hadd_ps (msum2, msum2);
msum2 = _mm_hadd_ps (msum2, msum2);
return _mm_cvtss_f32 (msum2);
}
```

```
Ref.
def l2sqr(x, y):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (x[d] - y[d])**2
    return diff
```



➤ 256bit SIMD Register
➤ Process eight floats at once



$\|x - y\|_2^2$ by SIMD

Rename variables for the sake of explanation

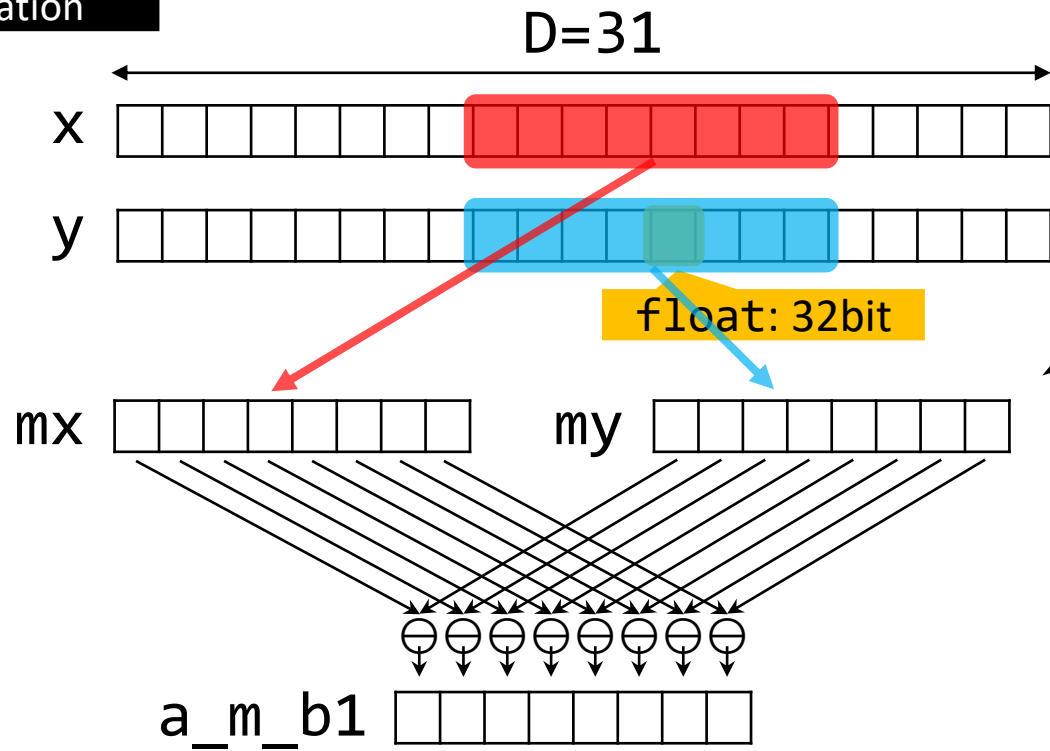
Ref.

```
def l2sqr(x, y):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (x[d] - y[d])**2
    return diff
```

```
float fvec_l2sqr (const float * x,
                 const float * y,
                 size_t d)
```

```
{
    __m256 msum1 = _mm256_setzero_ps();

    while (d >= 8) {
        __m256 mx = _mm256_loadu_ps (x); x += 8;
        __m256 my = _mm256_loadu_ps (y); y += 8;
        const __m256 a_m_b1 = mx - my;
        msum1 += a_m_b1 * a_m_b1;
        d -= 8;
    }
```



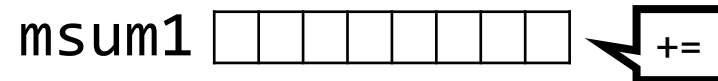
➤ 256bit SIMD Register
➤ Process eight floats at once

```
__m128 msum2 = _mm256_extractf128_ps(msum1, 1);
msum2 += _mm256_extractf128_ps(msum1, 0);
```

```
if (d >= 4) {
    __m128 mx = _mm_loadu_ps (x); x += 4;
    __m128 my = _mm_loadu_ps (y); y += 4;
    const __m128 a_m_b1 = mx - my;
    msum2 += a_m_b1 * a_m_b1;
    d -= 4;
}
```

```
if (d > 0) {
    __m128 mx = masked_read (d, x);
    __m128 my = masked_read (d, y);
    __m128 a_m_b1 = mx - my;
    msum2 += a_m_b1 * a_m_b1;
}
```

```
msum2 = _mm_hadd_ps (msum2, msum2);
msum2 = _mm_hadd_ps (msum2, msum2);
return _mm_cvtss_f32 (msum2);
}
```



$\|x - y\|_2^2$ by SIMD

Rename variables for the sake of explanation

```
float fvec_L2sqr (const float * x,
                 const float * y,
                 size_t d)
```

```
{
    __m256 msum1 = _mm256_setzero_ps();
```

```
while (d >= 8) {
    __m256 mx = _mm256_loadu_ps (x); x += 8;
    __m256 my = _mm256_loadu_ps (y); y += 8;
    const __m256 a_m_b1 = mx - my;
    msum1 += a_m_b1 * a_m_b1;
    d -= 8;
}
```

```
__m128 msum2 = _mm256_extractf128_ps(msum1, 1);
msum2 += _mm256_extractf128_ps(msum1, 0);
```

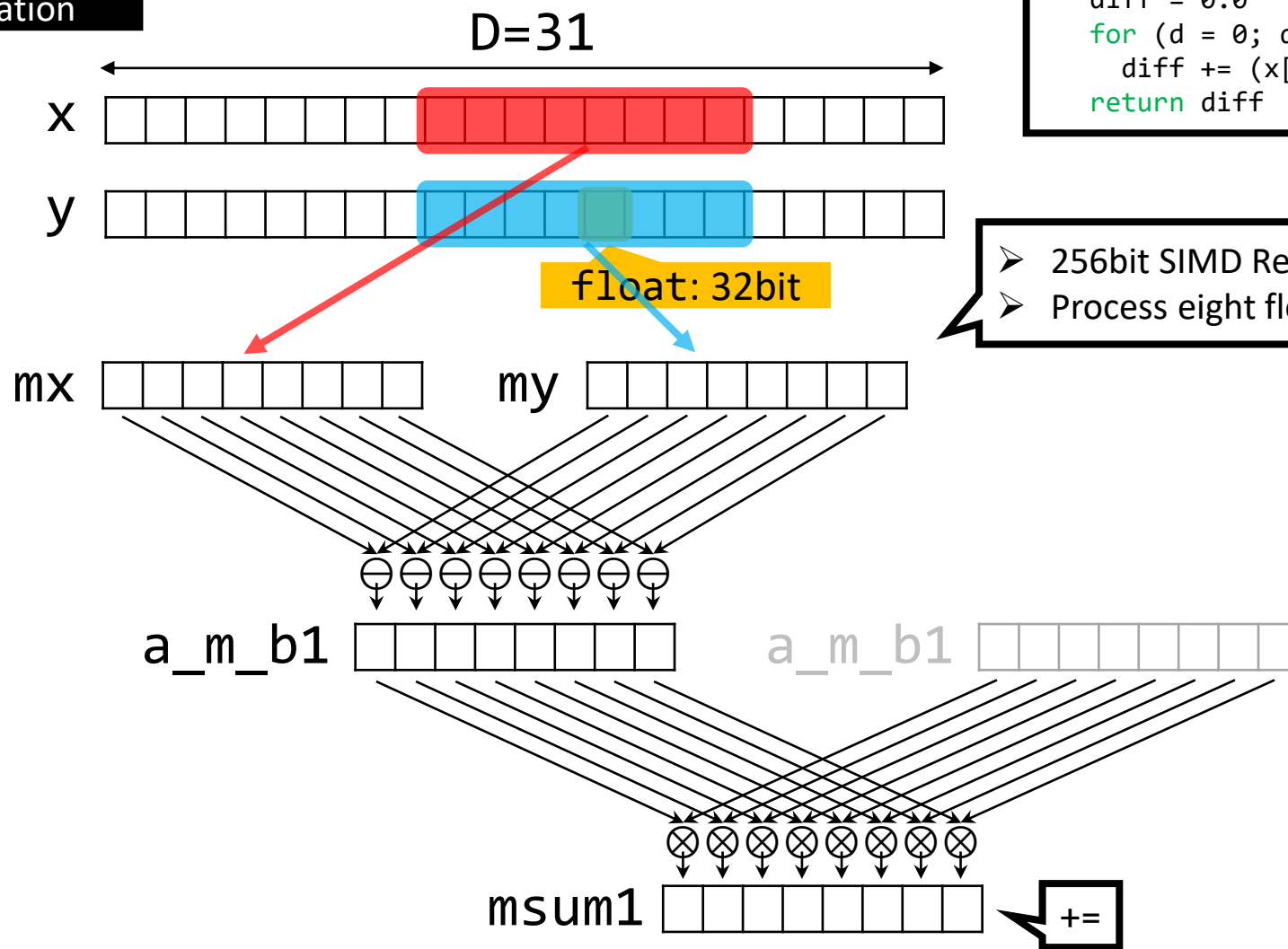
```
if (d >= 4) {
    __m128 mx = _mm_loadu_ps (x); x += 4;
    __m128 my = _mm_loadu_ps (y); y += 4;
    const __m128 a_m_b1 = mx - my;
    msum2 += a_m_b1 * a_m_b1;
    d -= 4;
}
```

```
if (d > 0) {
    __m128 mx = masked_read (d, x);
    __m128 my = masked_read (d, y);
    __m128 a_m_b1 = mx - my;
    msum2 += a_m_b1 * a_m_b1;
}
```

```
msum2 = _mm_hadd_ps (msum2, msum2);
msum2 = _mm_hadd_ps (msum2, msum2);
return _mm_cvtss_f32 (msum2);
}
```

Ref.

```
def l2sqr(x, y):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (x[d] - y[d])**2
    return diff
```



➤ 256bit SIMD Register
➤ Process eight floats at once

$\|x - y\|_2^2$ by SIMD

Rename variables for the sake of explanation

```

Ref.
def l2sqr(x, y):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (x[d] - y[d])**2
    return diff
    
```

```

float fvec_l2sqr (const float * x,
                 const float * y,
                 size_t d)
{
    __m256 msum1 = _mm256_setzero_ps();

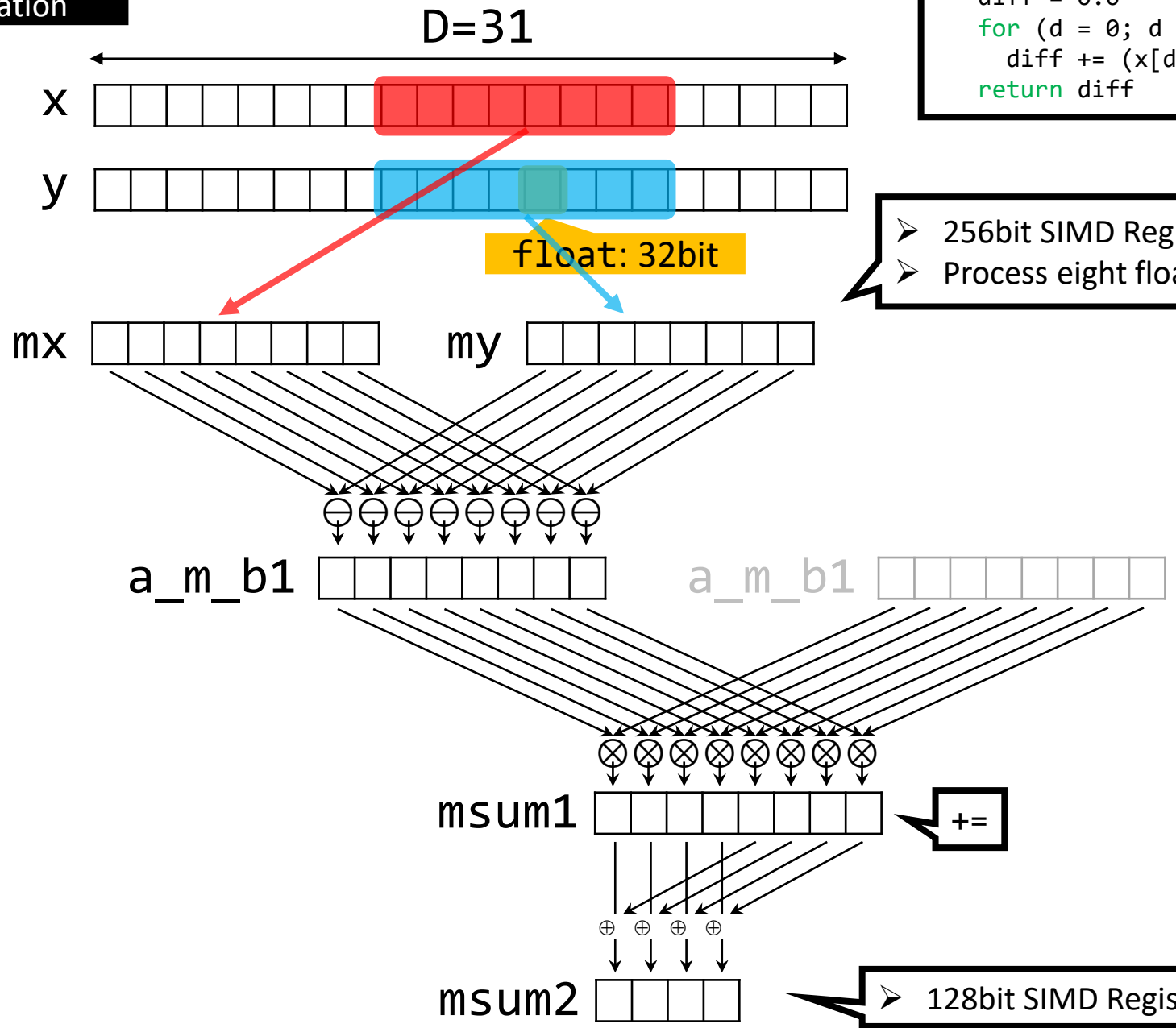
    while (d >= 8) {
        __m256 mx = _mm256_loadu_ps (x); x += 8;
        __m256 my = _mm256_loadu_ps (y); y += 8;
        const __m256 a_m_b1 = mx - my;
        msum1 += a_m_b1 * a_m_b1;
        d -= 8;
    }

    __m128 msum2 = _mm256_extractf128_ps(msum1, 1);
    msum2 += _mm256_extractf128_ps(msum1, 0);

    if (d >= 4) {
        __m128 mx = _mm_loadu_ps (x); x += 4;
        __m128 my = _mm_loadu_ps (y); y += 4;
        const __m128 a_m_b1 = mx - my;
        msum2 += a_m_b1 * a_m_b1;
        d -= 4;
    }

    if (d > 0) {
        __m128 mx = masked_read (d, x);
        __m128 my = masked_read (d, y);
        __m128 a_m_b1 = mx - my;
        msum2 += a_m_b1 * a_m_b1;
    }

    msum2 = _mm_hadd_ps (msum2, msum2);
    msum2 = _mm_hadd_ps (msum2, msum2);
    return _mm_cvtss_f32 (msum2);
}
    
```



➤ 256bit SIMD Register
➤ Process eight floats at once

➤ 128bit SIMD Register

$\|x - y\|_2^2$ by SIMD

Rename variables for the sake of explanation

```

Ref.
def l2sqr(x, y):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (x[d] - y[d])**2
    return diff

```

```

float fvec_L2sqr (const float * x,
                 const float * y,
                 size_t d)

```

```

{
    __m256 msum1 = _mm256_setzero_ps();

```

```

while (d >= 8) {
    __m256 mx = _mm256_loadu_ps (x); x += 8;
    __m256 my = _mm256_loadu_ps (y); y += 8;
    const __m256 a_m_b1 = mx - my;
    msum1 += a_m_b1 * a_m_b1;
    d -= 8;
}

```

```

__m128 msum2 = _mm256_extractf128_ps(msum1, 1);
msum2 += _mm256_extractf128_ps(msum1, 0);

```

```

if (d >= 4) {
    __m128 mx = _mm_loadu_ps (x); x += 4;
    __m128 my = _mm_loadu_ps (y); y += 4;
    const __m128 a_m_b1 = mx - my;
    msum2 += a_m_b1 * a_m_b1;
    d -= 4;
}

```

```

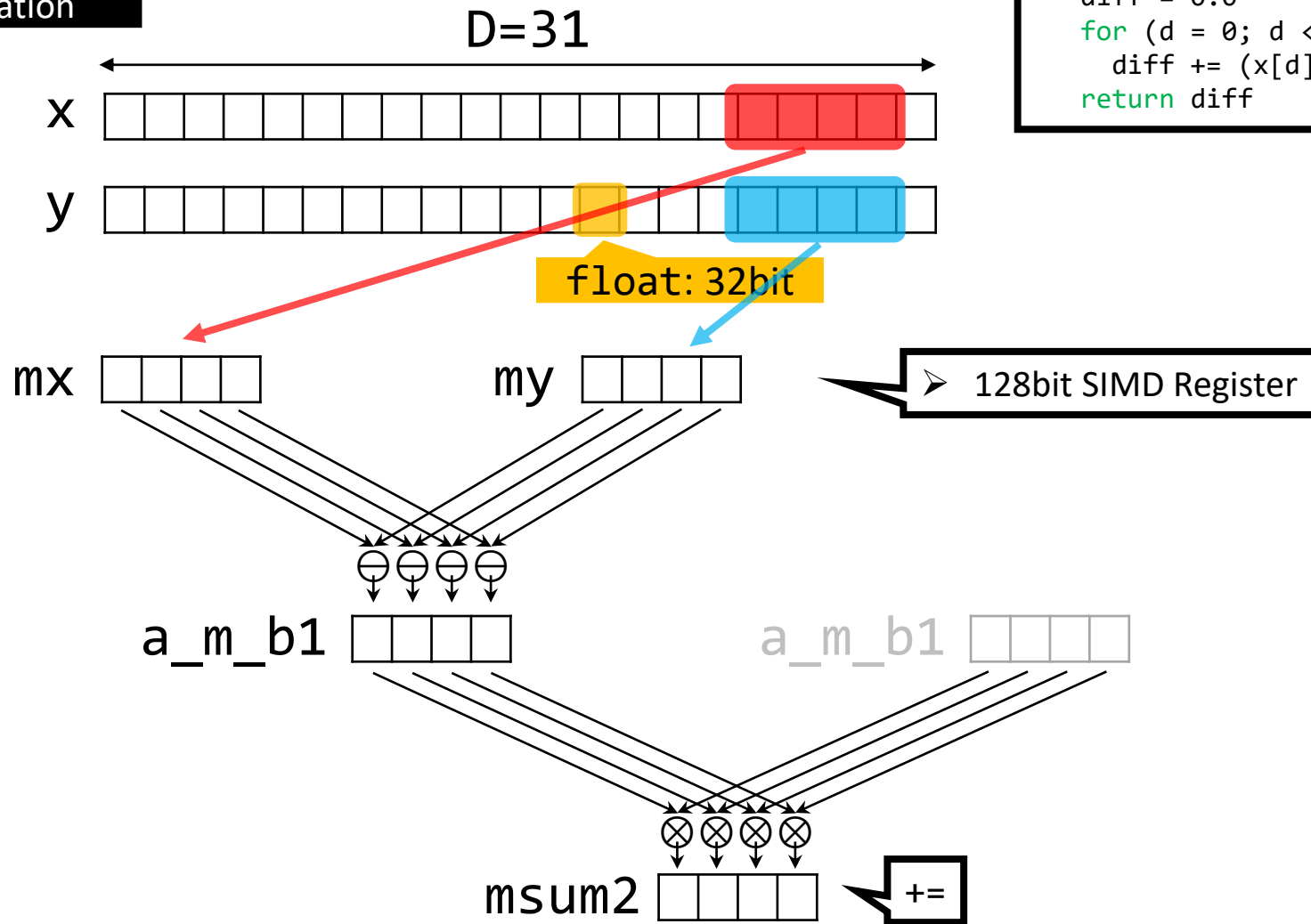
if (d > 0) {
    __m128 mx = masked_read (d, x);
    __m128 my = masked_read (d, y);
    __m128 a_m_b1 = mx - my;
    msum2 += a_m_b1 * a_m_b1;
}

```

```

msum2 = _mm_hadd_ps (msum2, msum2);
msum2 = _mm_hadd_ps (msum2, msum2);
return _mm_cvtss_f32 (msum2);
}

```



$\|x - y\|_2^2$ by SIMD

Rename variables for the sake of explanation

```
float fvec_L2sqr (const float * x,
                 const float * y,
                 size_t d)
{
    __m256 msum1 = _mm256_setzero_ps();

    while (d >= 8) {
        __m256 mx = _mm256_loadu_ps (x); x += 8;
        __m256 my = _mm256_loadu_ps (y); y += 8;
        const __m256 a_m_b1 = mx - my;
        msum1 += a_m_b1 * a_m_b1;
        d -= 8;
    }

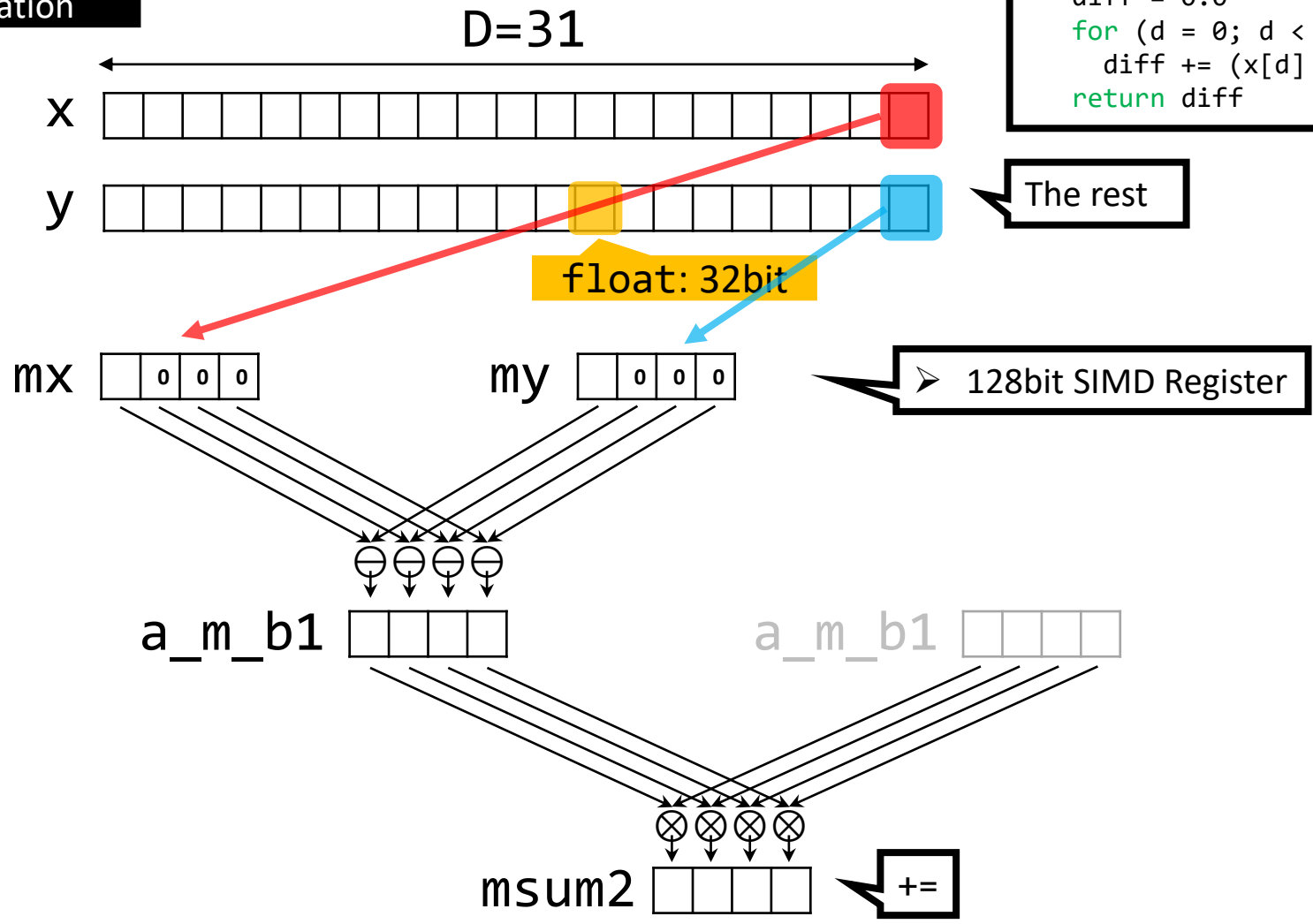
    __m128 msum2 = _mm256_extractf128_ps(msum1, 1);
    msum2 += _mm256_extractf128_ps(msum1, 0);

    if (d >= 4) {
        __m128 mx = _mm_loadu_ps (x); x += 4;
        __m128 my = _mm_loadu_ps (y); y += 4;
        const __m128 a_m_b1 = mx - my;
        msum2 += a_m_b1 * a_m_b1;
        d -= 4;
    }

    if (d > 0) {
        __m128 mx = masked_read (d, x);
        __m128 my = masked_read (d, y);
        __m128 a_m_b1 = mx - my;
        msum2 += a_m_b1 * a_m_b1;
    }

    msum2 = _mm_hadd_ps (msum2, msum2);
    msum2 = _mm_hadd_ps (msum2, msum2);
    return _mm_cvtss_f32 (msum2);
}
```

```
Ref.
def l2sqr(x, y):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (x[d] - y[d])**2
    return diff
```



$\|x - y\|_2^2$ by SIMD

Rename variables for the sake of explanation

```
float fvec_L2sqr (const float * x,
                 const float * y,
                 size_t d)
{
    __m256 msum1 = _mm256_setzero_ps();

    while (d >= 8) {
        __m256 mx = _mm256_loadu_ps (x); x += 8;
        __m256 my = _mm256_loadu_ps (y); y += 8;
        const __m256 a_m_b1 = mx - my;
        msum1 += a_m_b1 * a_m_b1;
        d -= 8;
    }

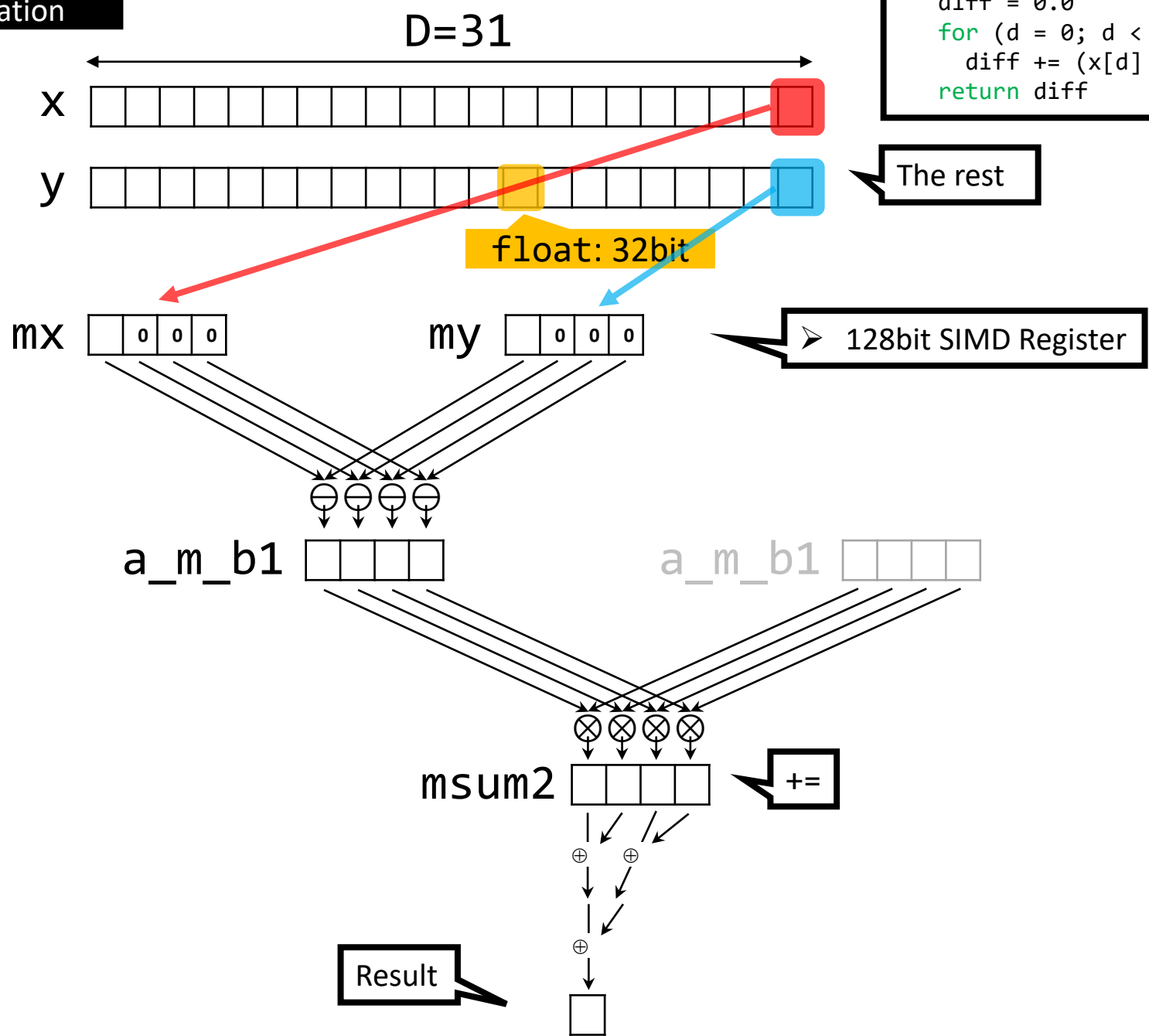
    __m128 msum2 = _mm256_extractf128_ps(msum1, 1);
    msum2 += _mm256_extractf128_ps(msum1, 0);

    if (d >= 4) {
        __m128 mx = _mm_loadu_ps (x); x += 4;
        __m128 my = _mm_loadu_ps (y); y += 4;
        const __m128 a_m_b1 = mx - my;
        msum2 += a_m_b1 * a_m_b1;
        d -= 4;
    }

    if (d > 0) {
        __m128 mx = masked_read (d, x);
        __m128 my = masked_read (d, y);
        __m128 a_m_b1 = mx - my;
        msum2 += a_m_b1 * a_m_b1;
    }

    msum2 = _mm_hadd_ps (msum2, msum2);
    msum2 = _mm_hadd_ps (msum2, msum2);
    return _mm_cvtss_f32 (msum2);
}
```

```
Ref.
def l2sqr(x, y):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (x[d] - y[d])**2
    return diff
```



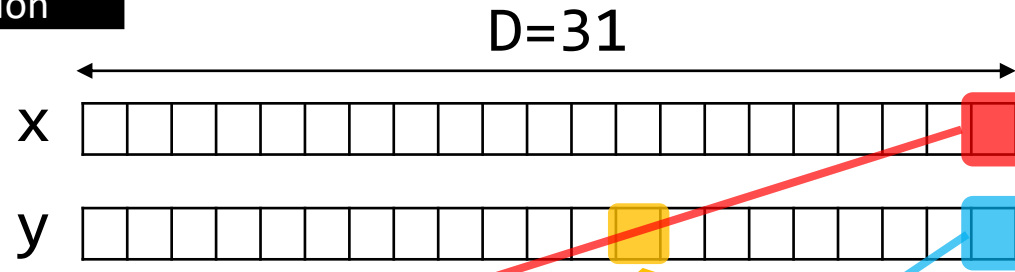
$\|x - y\|_2^2$ by SIMD

Rename variables for the sake of explanation

Ref.

```
def l2sqr(x, y):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (x[d] - y[d])**2
    return diff
```

```
float fvec_L2sqr (const float * x,
                 const float * y,
                 size_t d)
{
    __m256 msum1 = _mm256_setzero_ps();
```

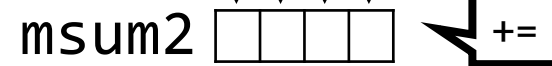


- SIMD codes of faiss are simple and easy to read
- Being able to read SIMD codes comes in handy sometimes; *why this impl is super fast*
- Another example of SIMD L2sqr from HNSW:

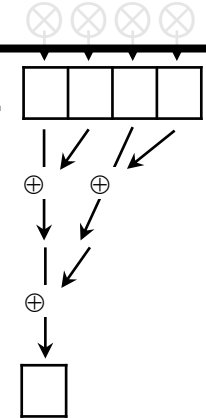
https://github.com/nmslib/hnswlib/blob/master/hnswlib/space_l2.h

```
__m128 mx = masked_read (d, x);
__m128 my = masked_read (d, y);
__m128 a_m_b1 = mx - my;
msum2 += a_m_b1 * a_m_b1;
```

```
msum2 = _mm_hadd_ps (msum2, msum2);
msum2 = _mm_hadd_ps (msum2, msum2);
return _mm_cvtss_f32 (msum2);
```



Result



M D -dim query vectors

$$Q = \{q_1, q_2, \dots, q_M\}$$

N D -dim database vectors

$$X = \{x_1, x_2, \dots, x_N\}$$

Task : Given $q \in Q$ and $x \in X$, compute $\|q - x\|_2^2$

Naïve impl.

```
def l2sqr(q, x):
    diff = 0.0
    for (d = 0; d < D; ++d):
        diff += (q[d] - x[d])**2
    return diff
```

```
parfor q in Q:
    for x in X:
        l2sqr(q, x)
```

Parallelize query-side

Select min by heap, but omit it now

faiss impl.

if $M < 20$:

compute $\|q - x\|_2^2$ by SIMD

else :

compute $\|q - x\|_2^2 = \|q\|_2^2 - 2q^T x + \|x\|_2^2$ by BLAS

Compute $\|\mathbf{q} - \mathbf{x}\|_2^2 = \|\mathbf{q}\|_2^2 - 2\mathbf{q}^\top \mathbf{x} + \|\mathbf{x}\|_2^2$ with BLAS

Stack M D -dim query vectors to a $D \times M$ matrix: $Q = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_M] \in \mathbb{R}^{D \times M}$

Stack N D -dim database vectors to a $D \times N$ matrix: $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$

Compute tables

SIMD-accelerated function

```
q_norms = norms(Q) #  $\|\mathbf{q}_1\|_2^2, \|\mathbf{q}_2\|_2^2, \dots, \|\mathbf{q}_M\|_2^2$   
x_norms = norms(X) #  $\|\mathbf{x}_1\|_2^2, \|\mathbf{x}_2\|_2^2, \dots, \|\mathbf{x}_N\|_2^2$   
ip = sgemm_(Q, X, ...) #  $Q^\top X$ 
```

Scan and sum

```
parfor (m = 0; m < M; ++m):  
  for (n = 0; n < N; ++n):
```

```
    dist = q_norms[m] + x_norms[n] - ip[m][n]
```

- Matrix multiplication by BLAS
- Dominant if Q and X are large
- The difference of the background matters:
 - ✓ Intel MKL is 30% faster than OpenBLAS

$$\|\mathbf{q}_m - \mathbf{x}_n\|_2^2$$

$$\|\mathbf{q}_m\|_2^2$$

$$\|\mathbf{x}_n\|_2^2$$

$$(Q^\top X)_{mn}$$

NN in GPU (faiss-gpu) is 10x faster than NN in CPU (faiss-cpu)

Benchmark: <https://github.com/facebookresearch/faiss/wiki/Low-level-benchmarks>

➤ NN-GPU always compute $\|\mathbf{q}\|_2^2 - 2\mathbf{q}^\top \mathbf{x} + \|\mathbf{x}\|_2^2$

➤ k-means for 1M vectors (D=256, K=20000)

✓ 11 min on CPU

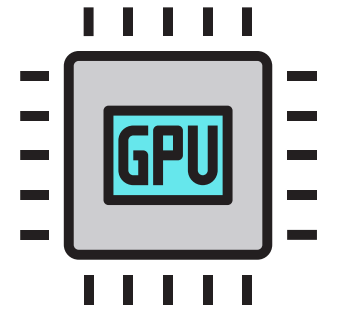
✓ 55 sec on 1 Pascal-class P100 GPU (float32 math)

✓ 34 sec on 1 Pascal-class P100 GPU (float16 math)

✓ 21 sec on 4 Pascal-class P100 GPUs (float32 math)

✓ 16 sec on 4 Pascal-class P100 GPUs (float16 math)

x10 faster



➤ If GPU is available and its memory is enough, try GPU-NN

➤ The behavior is little bit different (e.g., a restriction for top-k)

Reference

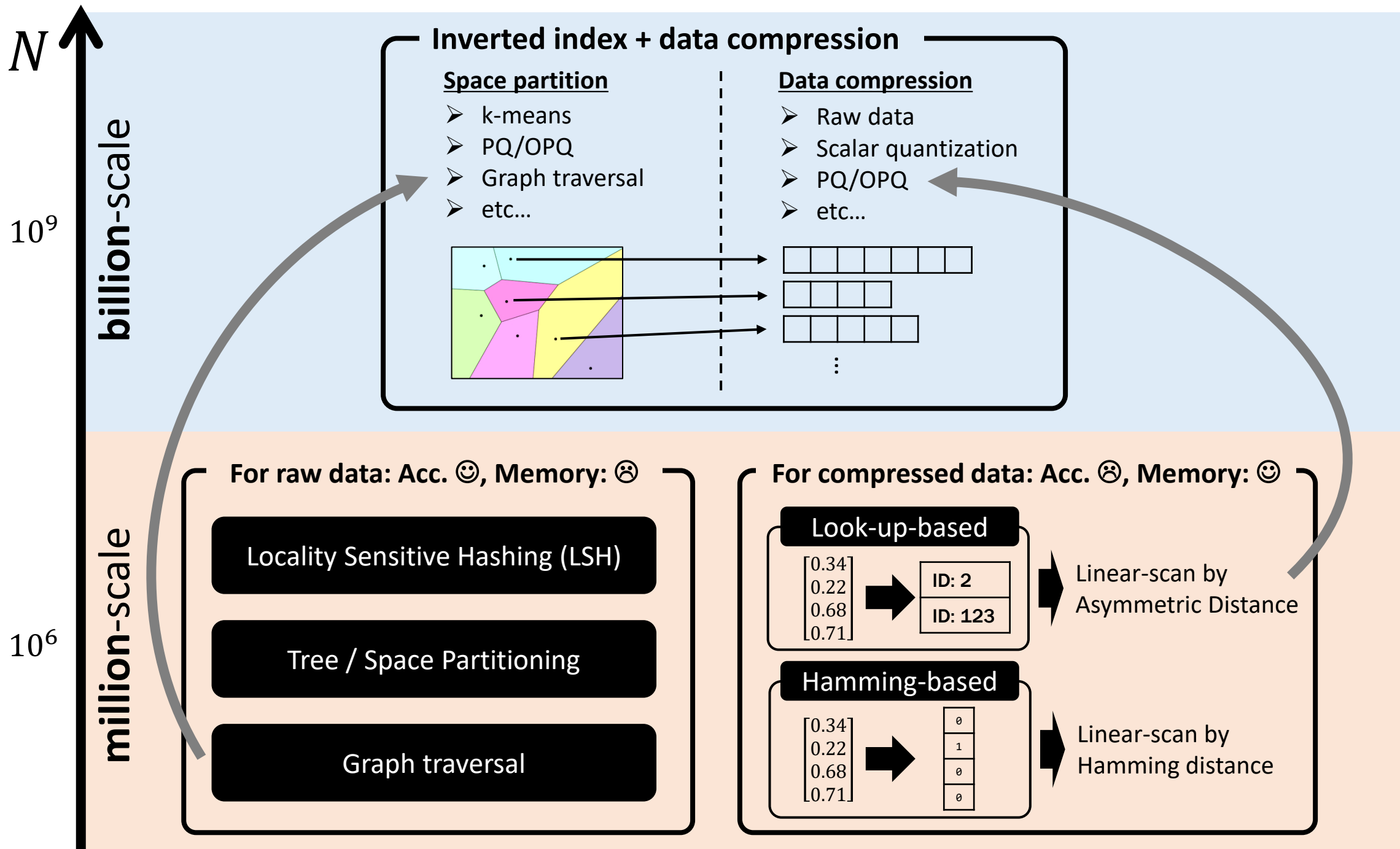
- Switch implementation of L2sqr in faiss:
[<https://github.com/facebookresearch/faiss/wiki/Implementation-notes#matrix-multiplication-to-do-many-l2-distance-computations>]
- Introduction to SIMD: a lecture by Markus Püschel (ETH) [[How to Write Fast Numerical Code - Spring 2019](#)], especially [[SIMD vector instructions](#)]
 - ✓ <https://acl.inf.ethz.ch/teaching/fastcode/2019/>
 - ✓ <https://acl.inf.ethz.ch/teaching/fastcode/2019/slides/07-simd.pdf>
- SIMD codes for faiss [https://github.com/facebookresearch/faiss/blob/master/utils/distances_simd.cpp]
- L2sqr benchmark including AVX512 for faiss-L2sqr
[<https://gist.github.com/matsui528/583925f88fcb08240319030202588c74>]

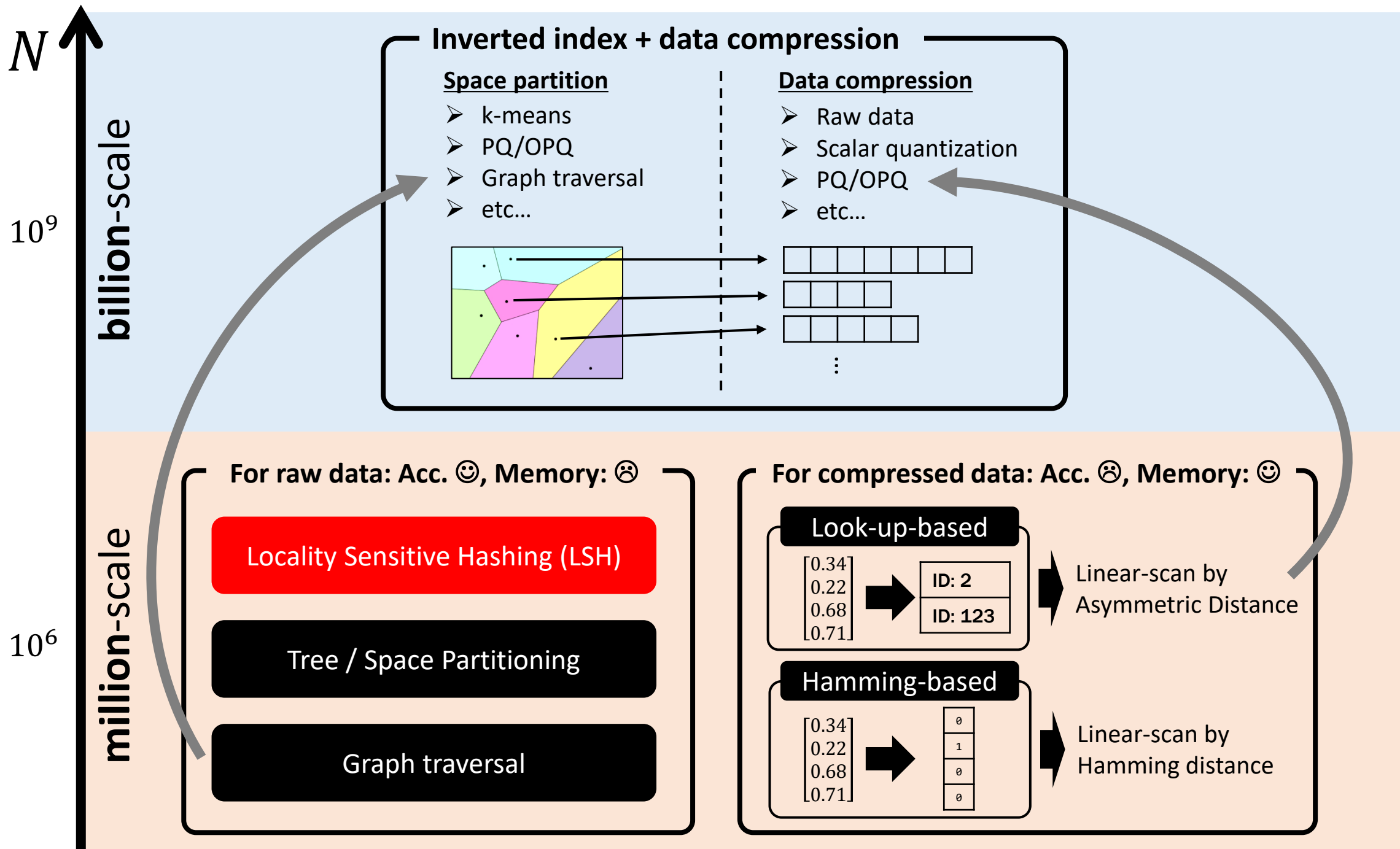
Part 1:

Nearest Neighbor Search

Part 2:

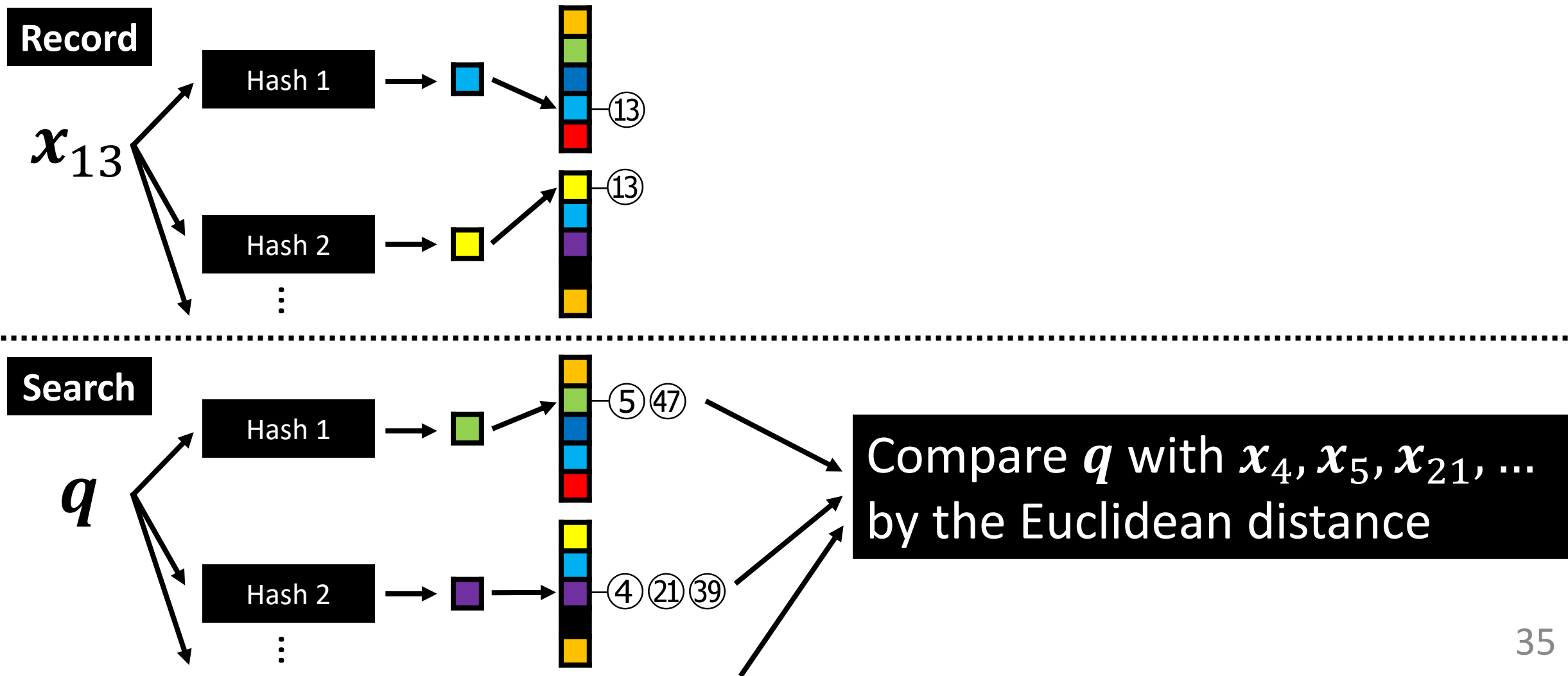
Approximate Nearest Neighbor Search





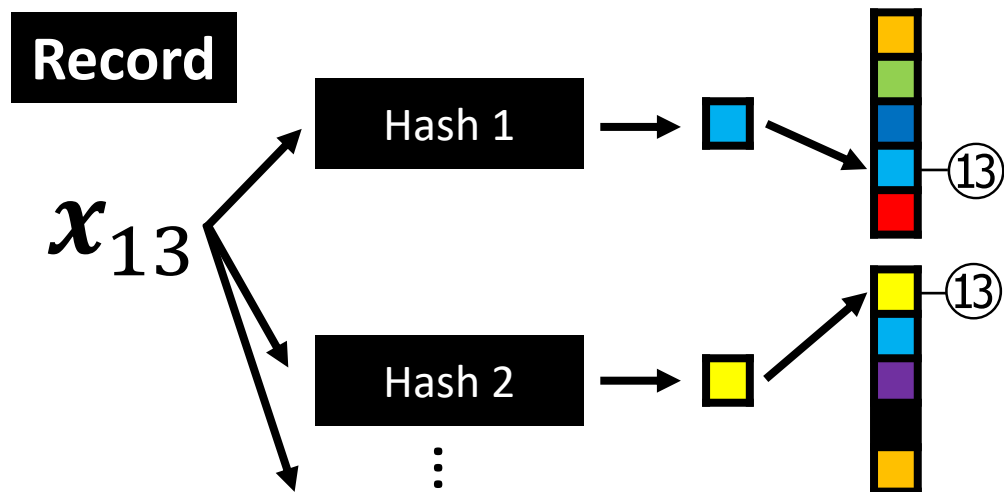
Locality Sensitive Hashing (LSH)

- LSH = Hash functions + Hash tables
- Map similar items to the same symbol with a high probability



Locality Sensitive Hashing (LSH)

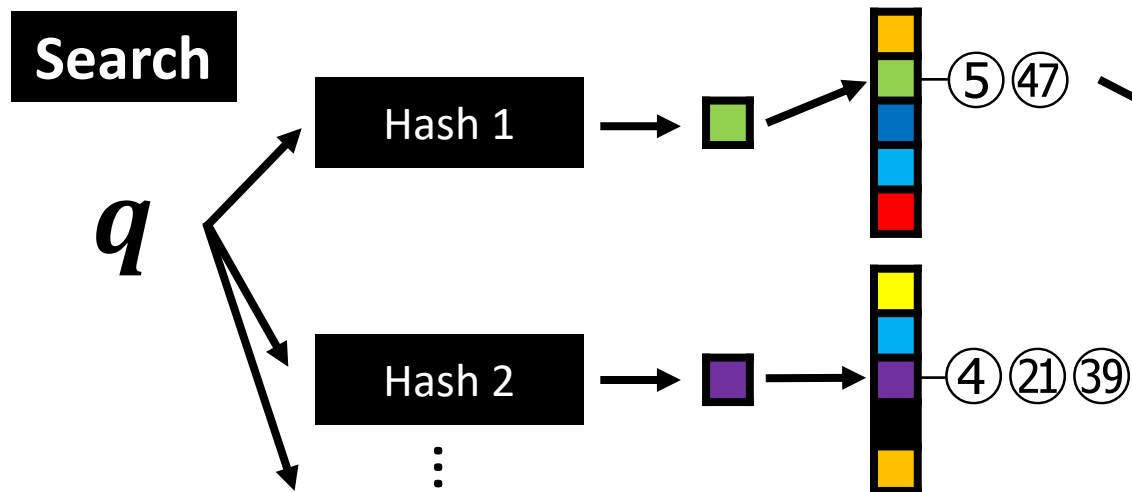
- LSH = Hash functions + Hash tables
- Map similar items to the same symbol with a high probability



E.g., random projection [Datar+, SCG 04]

$$H(\mathbf{x}) = [h_1(\mathbf{x}), \dots, h_M(\mathbf{x})]^T$$

$$h_m(\mathbf{x}) = \left\lfloor \frac{\mathbf{a}^T \mathbf{x} + b}{W} \right\rfloor$$



Compare q with x_4, x_5, x_{21}, \dots
by the Euclidean distance

Locality Sensitive Hashing (LSH)

➤ LSH = Hash functions + Hash tables

😊: Map similar items to the same symbol with a high probability

➤ Math-friendly

➤ Popular in the theory area (FOCS, STOC, ...)

☹: Large memory cost

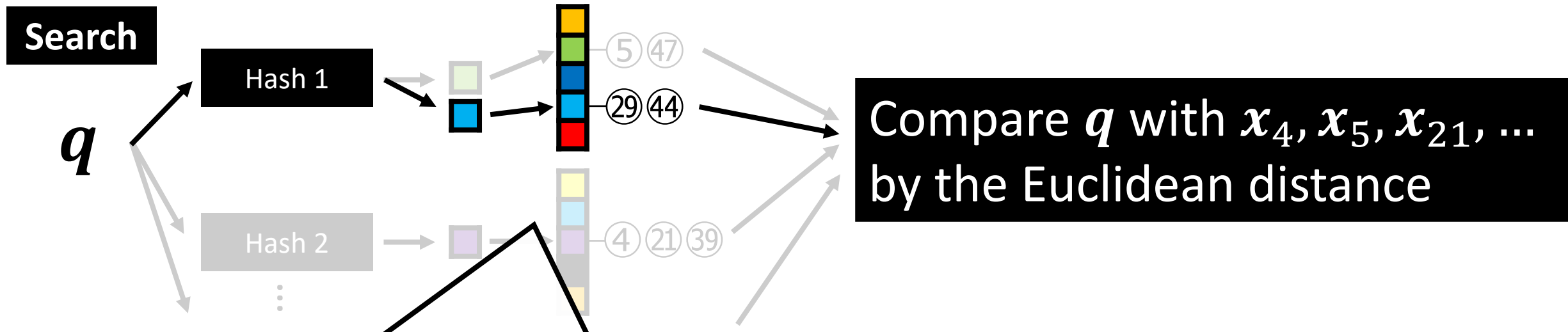
- ✓ Need several tables to boost the accuracy

- ✓ Need to store the original data, $\{x_n\}_{n=1}^N$, on memory

➤ Data-dependent methods such as PQ are better for real-world data

➤ Thus, in recent CV papers, LSH has been treated as a classic-method ☹☹☹

$$H(\mathbf{x}) = [h_1(\mathbf{x}), \dots, h_M(\mathbf{x})]^T$$
$$h_m(\mathbf{x}) = \left\lfloor \frac{\mathbf{a}^T \mathbf{x} + b}{W} \right\rfloor$$



In fact:

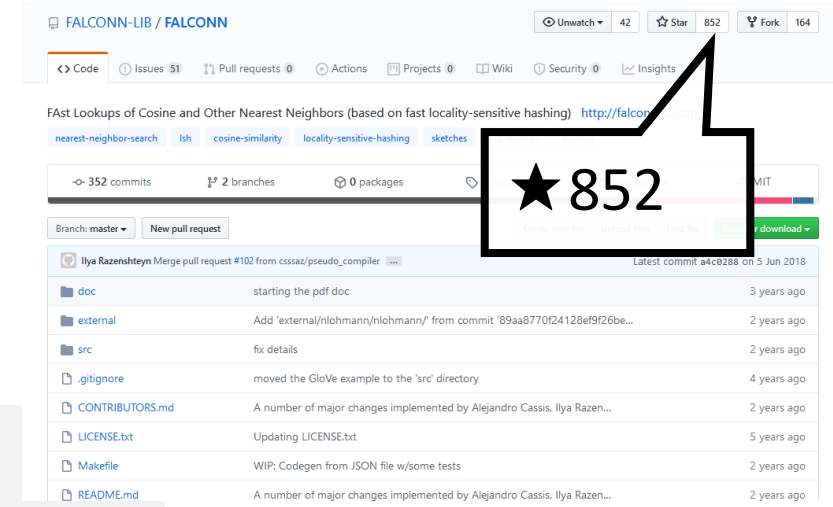
- Consider a next candidate → practical memory consumption (Multi-Probe [Lv+, VLDB 07])
- A library based on the idea: **FALCONN**

Falconn

<https://github.com/falconn-lib/falconn>

```
$> pip install FALCONN
```

```
table = falconn.LSHIndex(params_cp)
table.setup(X-center)
query_object = table.construct_query_object()
# query parameter config here
query_object.find_nearest_neighbor(Q-center, topk)
```



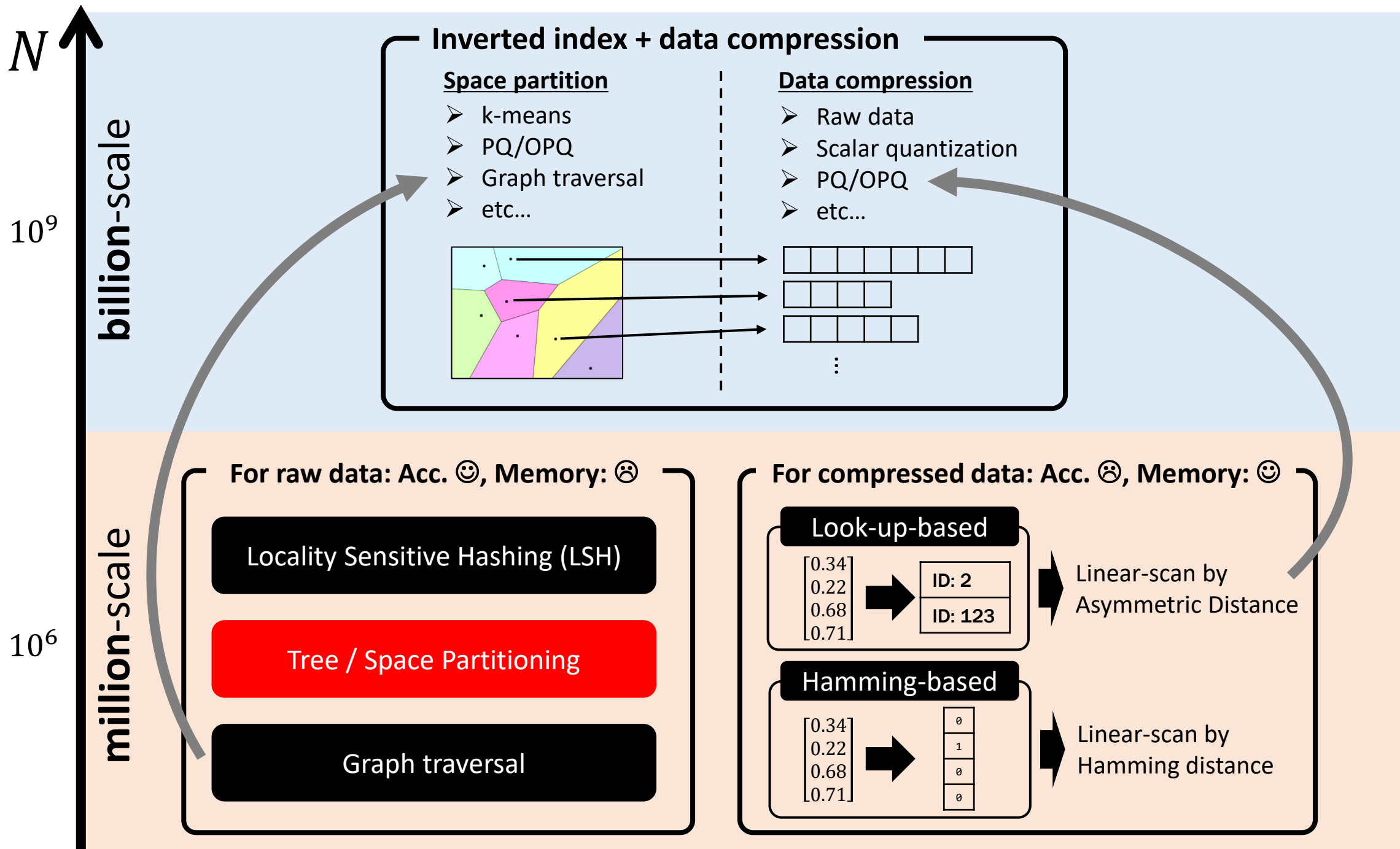
😊 Faster data addition (than annoy, nmslib, ivfpq)

😊 Useful for on-the-fly addition

😞 Parameter configuration seems a bit non-intuitive

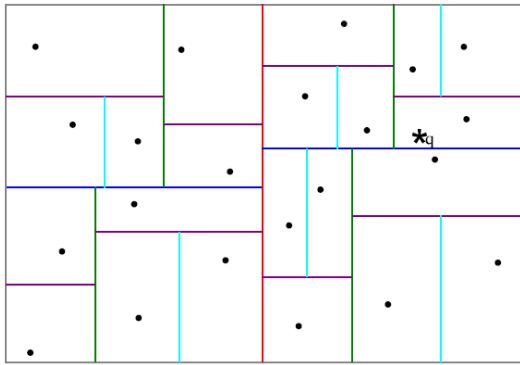
Reference

- Good summaries on this field: CVPR 2014 Tutorial on Large-Scale Visual Recognition, Part I: Efficient matching, H. Jégou
[<https://sites.google.com/site/lsvrtutorialcvpr14/home/efficient-matching>]
- Practical Q&A: FAQ in Wiki of FALCONN [<https://github.com/FALCONN-LIB/FALCONN/wiki/FAQ>]
- Hash functions: M. Datar et al., “Locality-sensitive hashing scheme based on p-stable distributions,” SCG 2004.
- Multi-Probe: Q. Lv et al., “Multi-Probe LSH: Efficient Indexing for High-Dimensional Similarity Search”, VLDB 2007
- Survey: A. Andoni and P. Indyk, “Near-Optimal Hashing Algorithms for Approximate Nearest Neighbor in High Dimensions,” Comm. ACM 2008

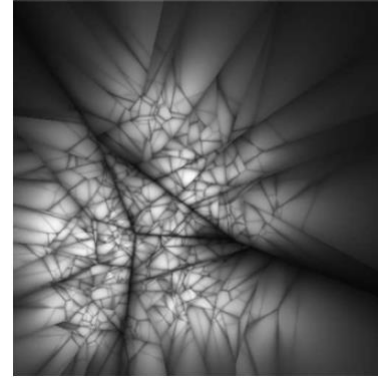
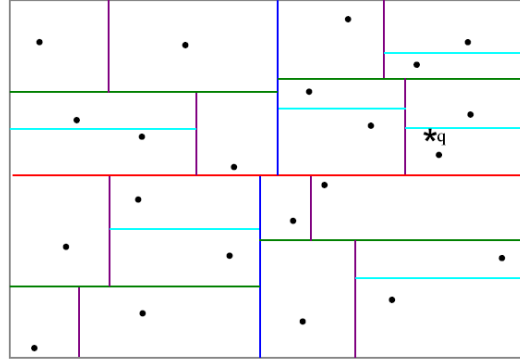


FLANN: Fast Library for Approximate Nearest Neighbors

Images are from [Muja and Lowe, TPAMI 2014]



Randomized KD Tree



k-means Tree

➤ Automatically select “Randomized KD Tree” or “k-means Tree”

<https://github.com/mariusmuja/flann>

☺ Good code base. Implemented in OpenCV and PCL

☺ Very popular in the late 00's and early 10's

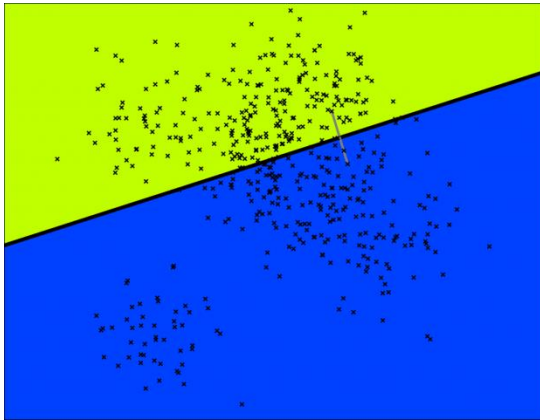
☹ Large memory consumption. The original data need to be stored

☹ Not actively maintained now

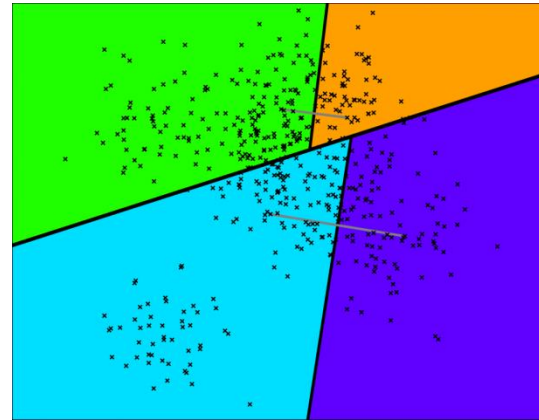
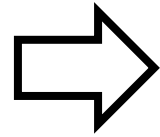
Annoy

“2-means tree”+ “multiple-trees” + “shared priority queue”

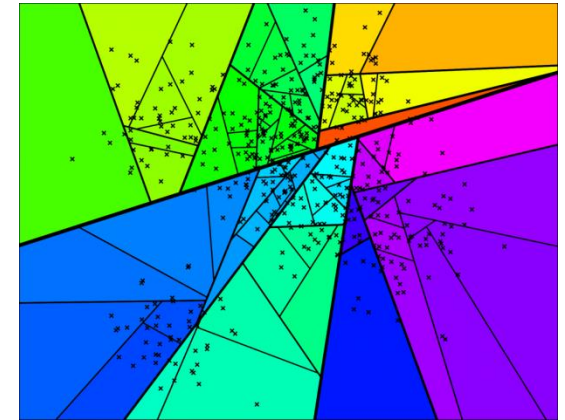
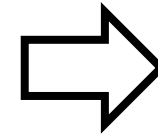
Record



Select two points randomly

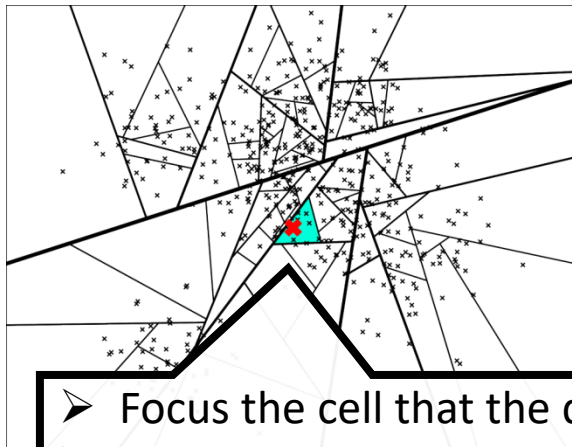


Divide up the space

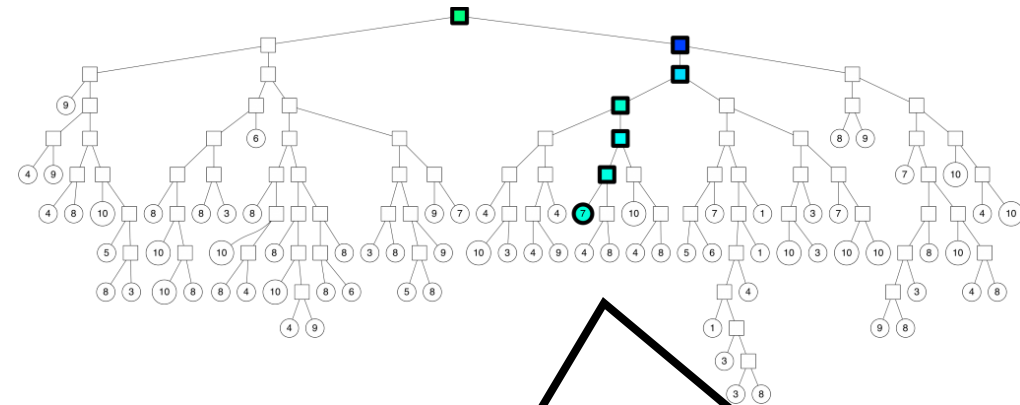


Repeat hierarchically

Search



- Focus the cell that the query lives
- Compare the distances



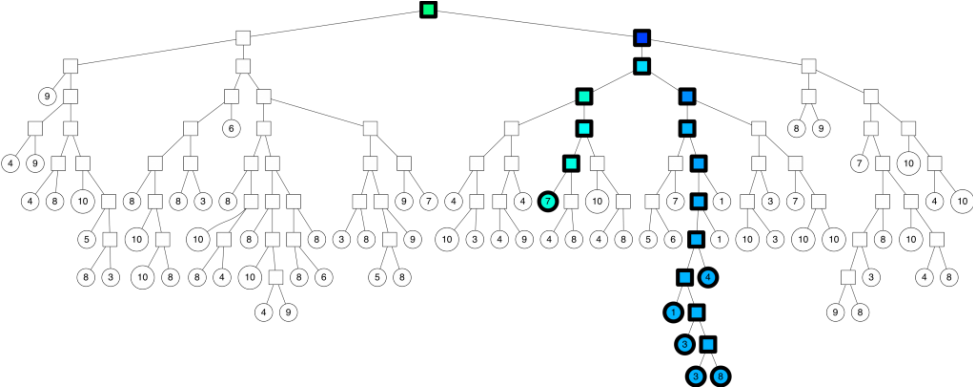
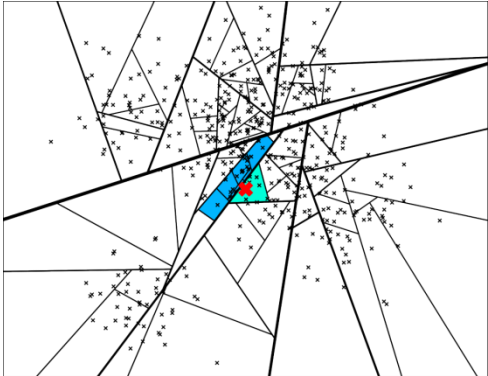
Can traverse the tree by log-times comparisons

Annoy

“2-means tree”+ “multiple-trees” + “shared priority queue”

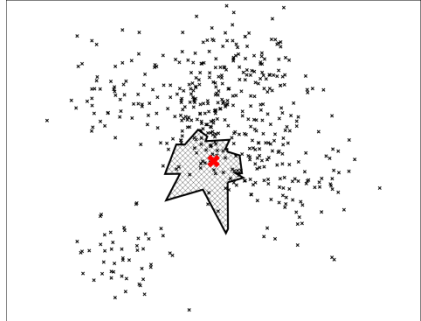
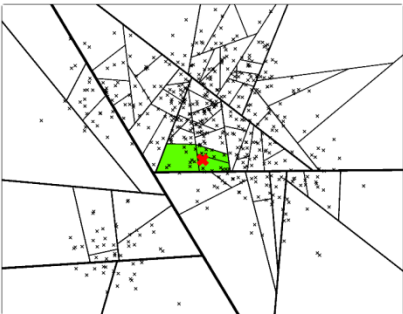
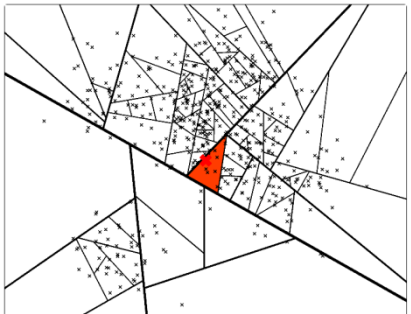
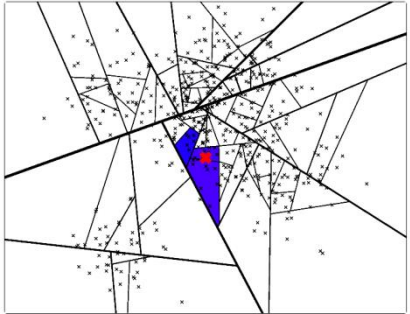
Feature 1

If we need more data points, use a priority queue



Feature 2

Boost the accuracy by multi-tree with a shared priority queue



Annoy

<https://github.com/erikbern/annoy>

```
$> pip install annoy
```

```
t = AnnoyIndex(D)
for n, x in enumerate(X):
    t.add_item(n, x)
t.build(n_trees)

t.get_nns_by_vector(q, topk)
```

spotify / annoy

Used by 1k Unwatch 332 Unstar 7.1k Fork 778

<> Code Issues 21 Pull requests 3 Actions Projects 0 Wiki Security 0 Insights

Approximate Nearest Neighbors in C++/Python optimized for memory usage and load

c-plus-plus python nearest-neighbor-search locality-sensitive-hashing approximate-nearest-neighbors

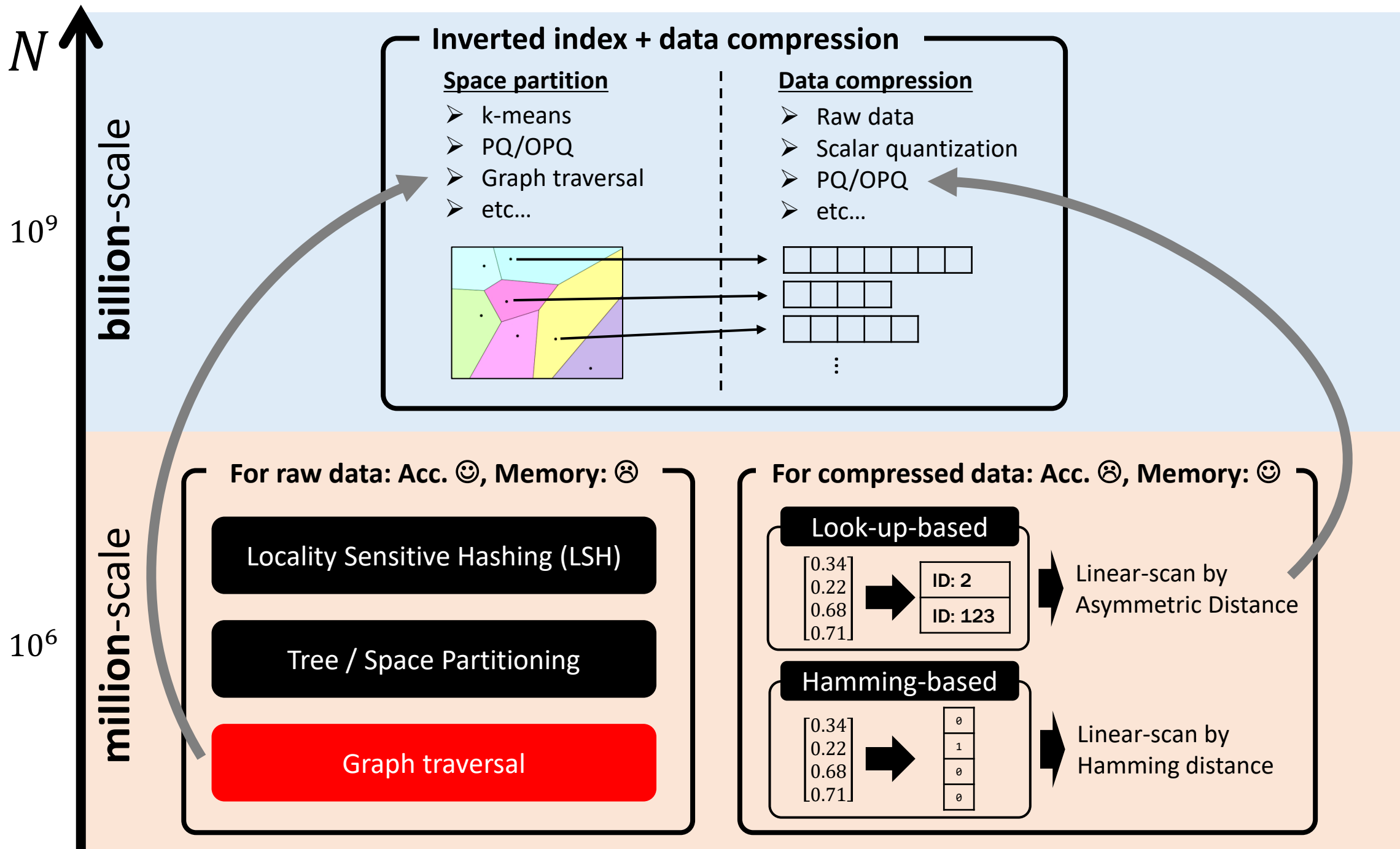
745 commits 19 branches 0 packages 23 releases 48 contributors Apache-2.0

Branch: master New pull request Create new file Upload files Find file Clone or download

erikbern Update README.rst Latest commit 8b6a825 on 12 May

annoy	remove subclass	3 years ago
debian	removed boost from debian/control and .travis.yml	5 years ago
examples	fix another futurewarning	11 months ago
src	A more informative error for #423	last month
test	Fix misc minor compilation warnings	3 months ago
.gitignore	Improve .gitignore coverage of files created by tests	8 months ago
.travis.yml	unrelated os x failure, try bumping python versoin	6 months ago
LICENSE	added Apache license	7 years ago

- 😊 Developed at Spotify. Well-maintained. Stable
- 😊 Simple interface with only a few parameters
- 😊 Baseline for million-scale data
- 😊 Support mmap, i.e., can be accessed from several processes
- 😞 Large memory consumption
- 😞 Runtime itself is slower than HNSW



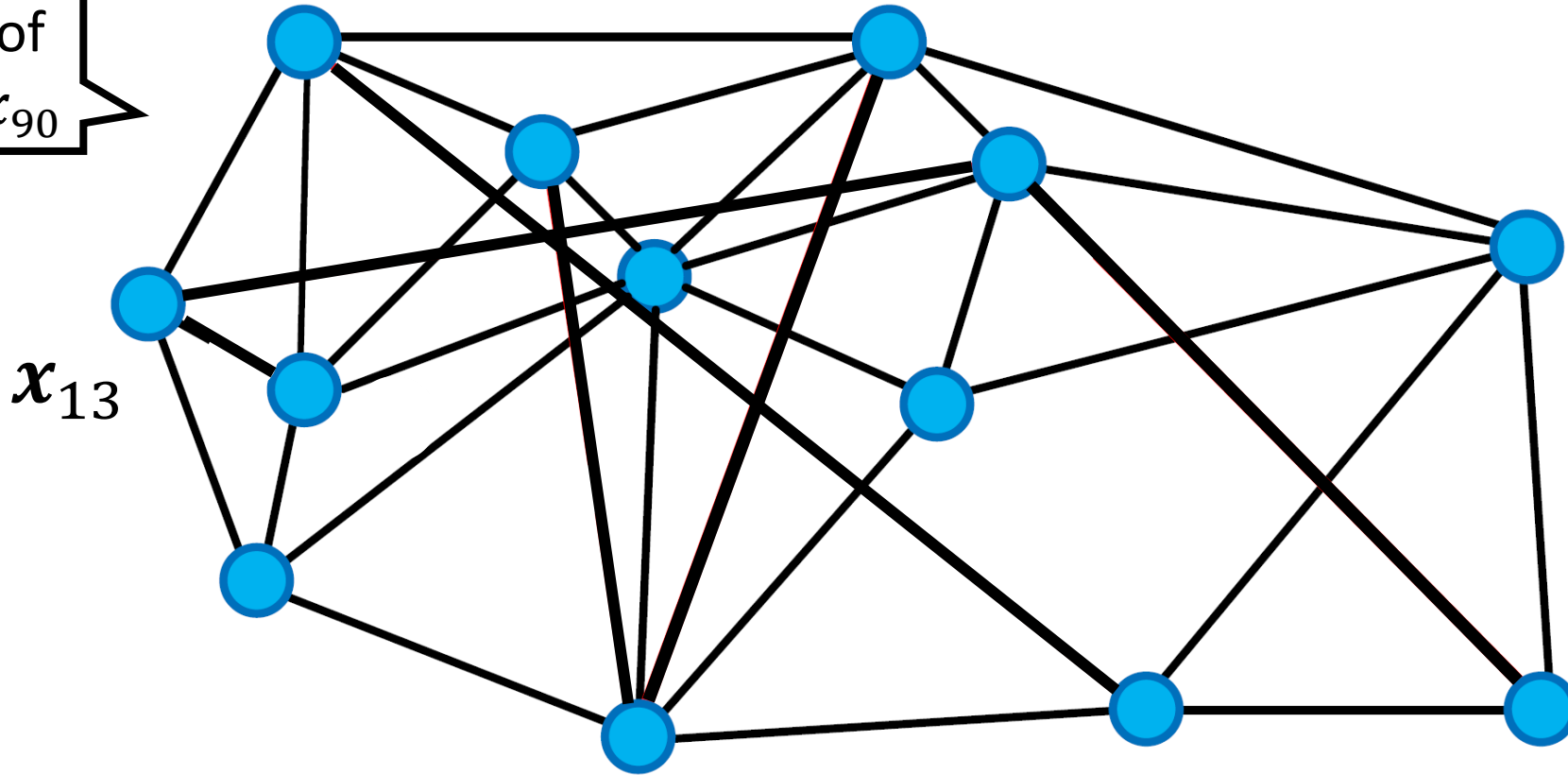
Graph traversal

- **Very popular** in recent years
- Around 2017, it turned out that the graph-traversal-based methods work well for million-scale data
- Pioneer:
 - ✓ Navigable Small World Graphs (NSW)
 - ✓ Hierarchical NSW (**HNSW**)
- Implementation: nmslib, hnsw, faiss

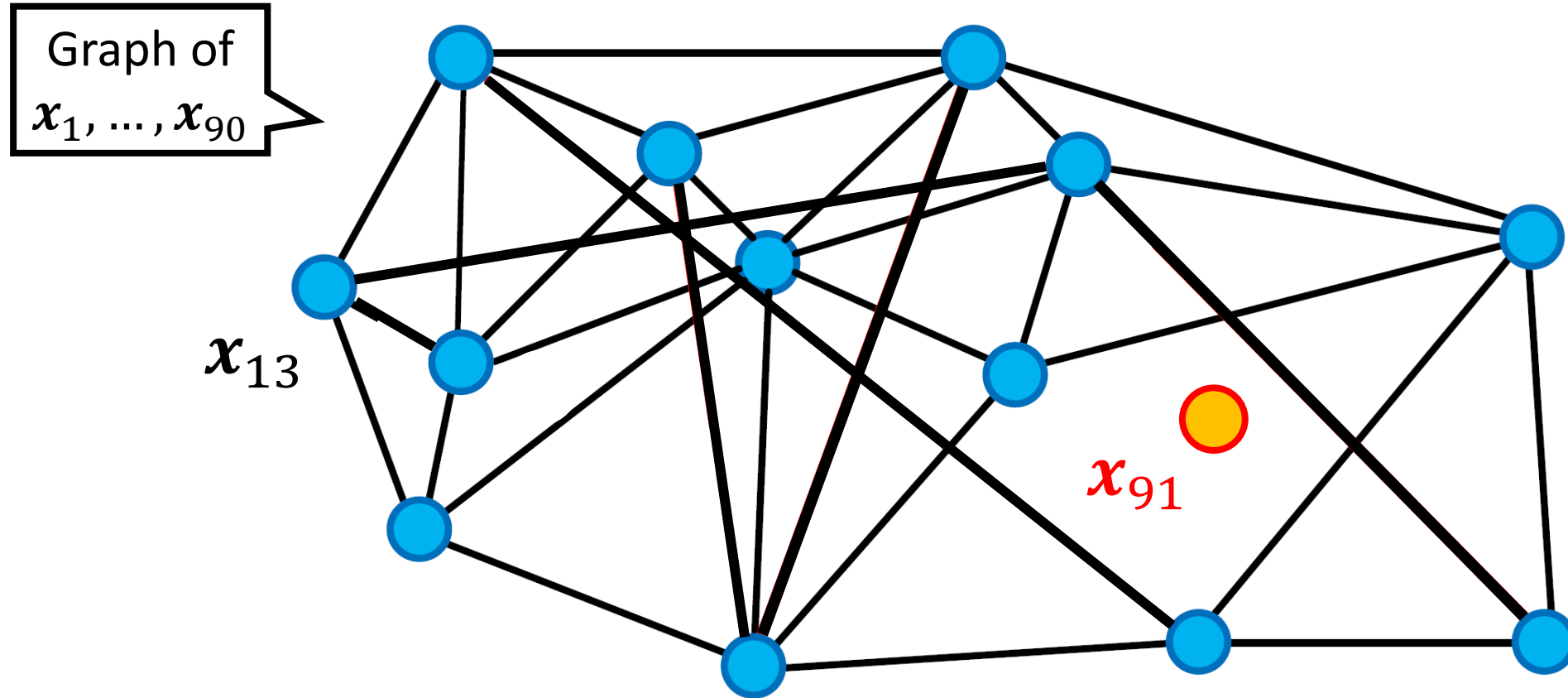
Record

Images are from [Malkov+, Information Systems, 2013]

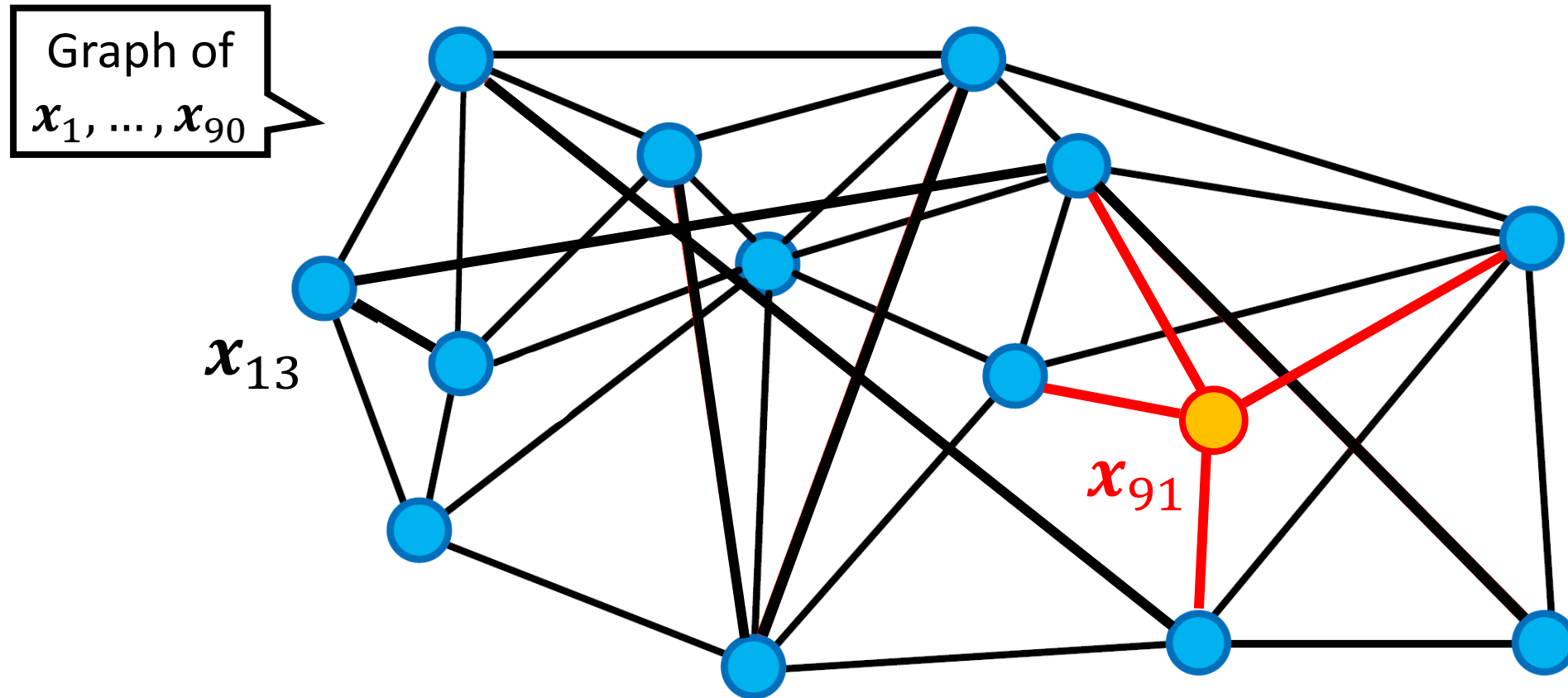
Graph of x_1, \dots, x_{90}



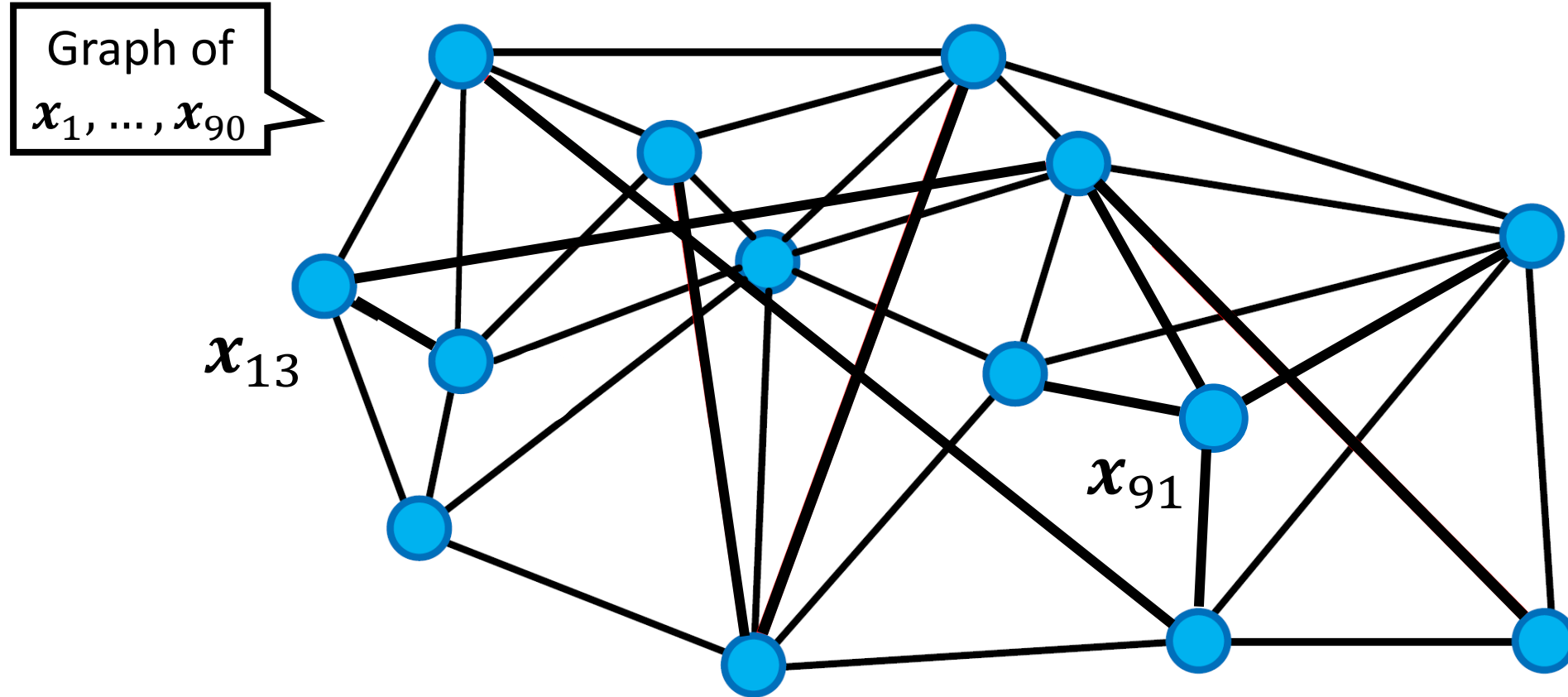
➤ Each node is a database vector



- Each node is a database vector
- Given a new database vector, create new edges to neighbors



- Each node is a database vector
- Given a new database vector, create new edges to neighbors

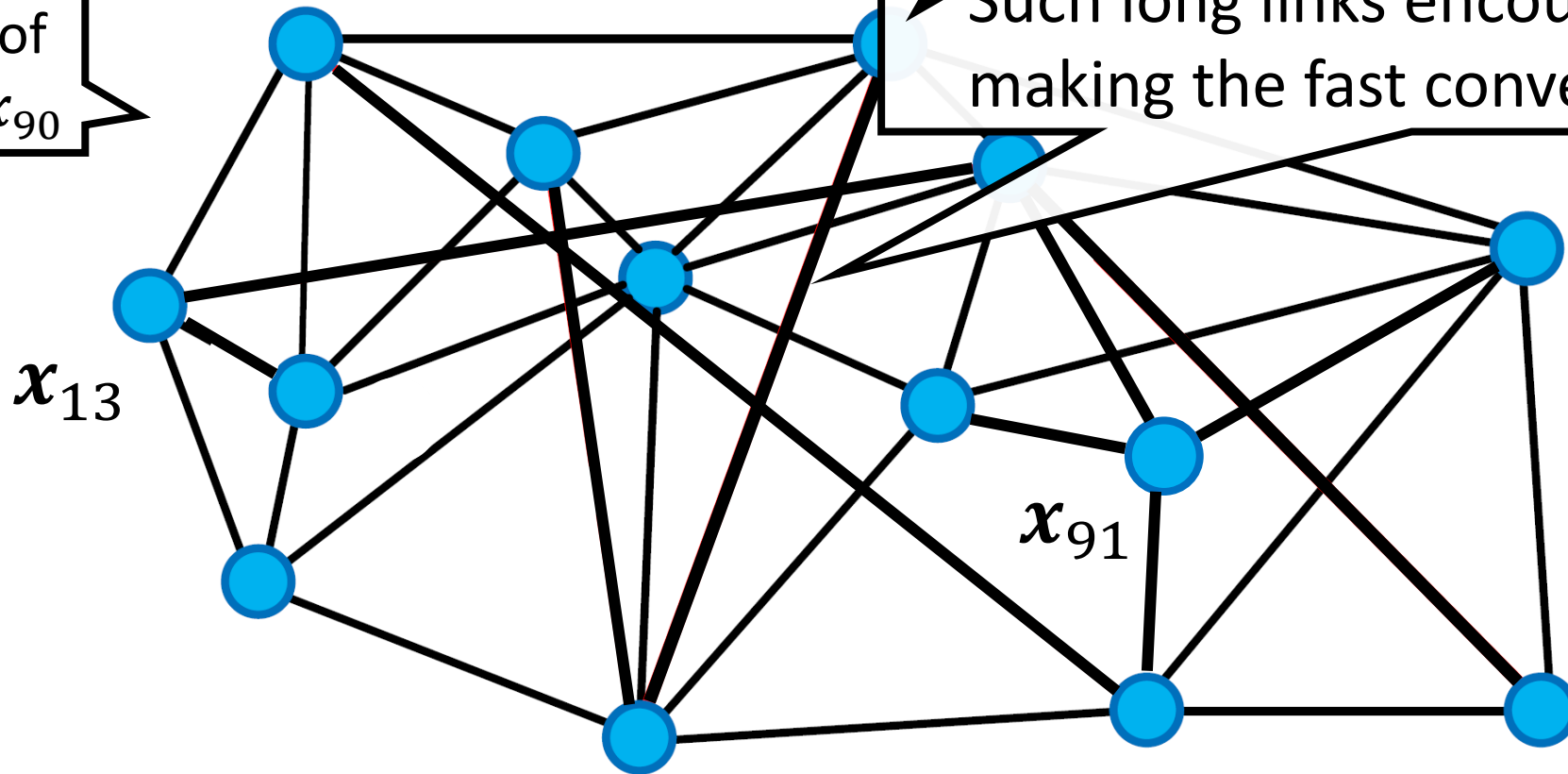


- Each node is a database vector
- Given a new database vector, create new edges to neighbors

Record

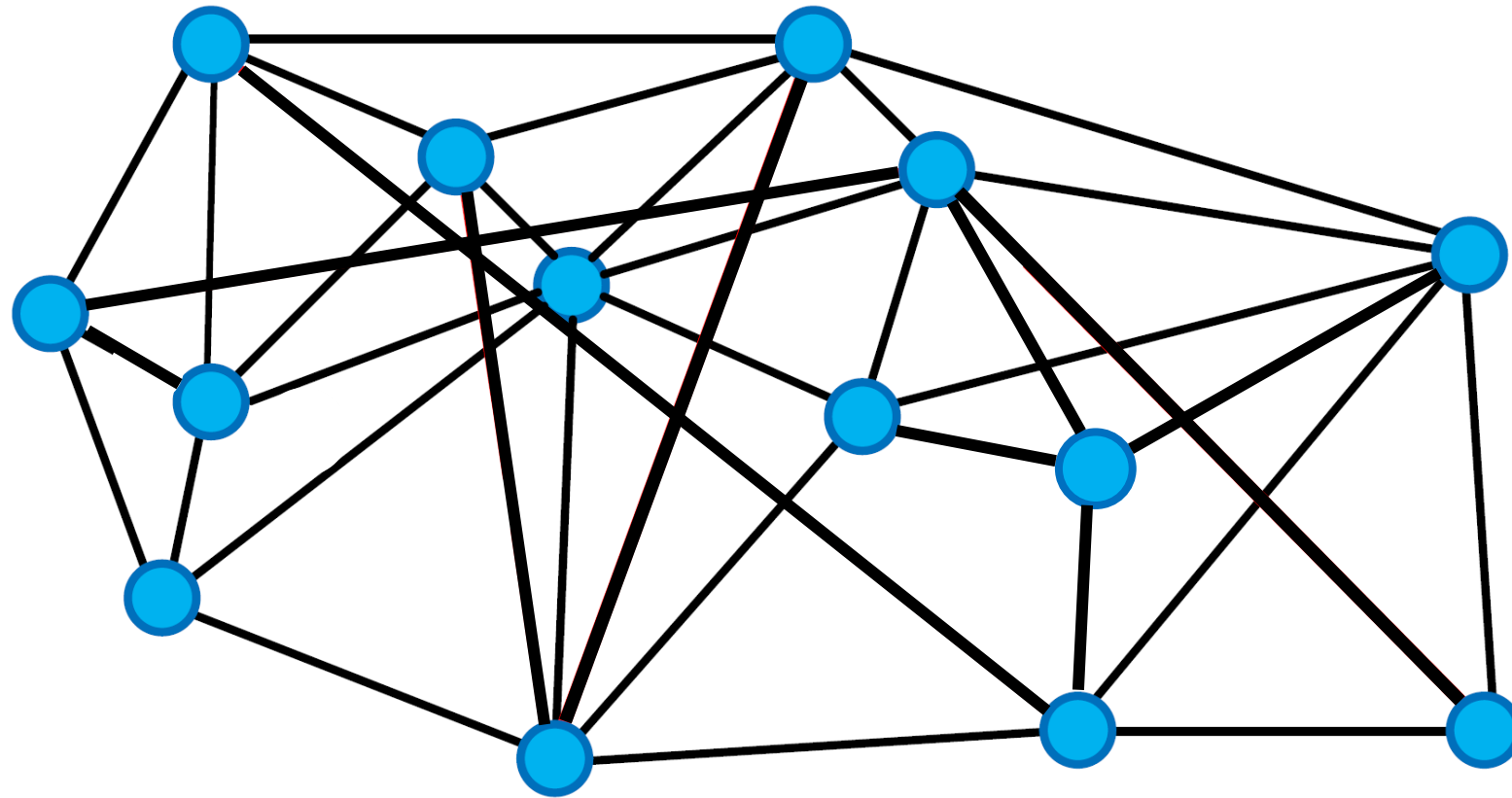
Images are from [Malkov+, Information Systems, 2013]

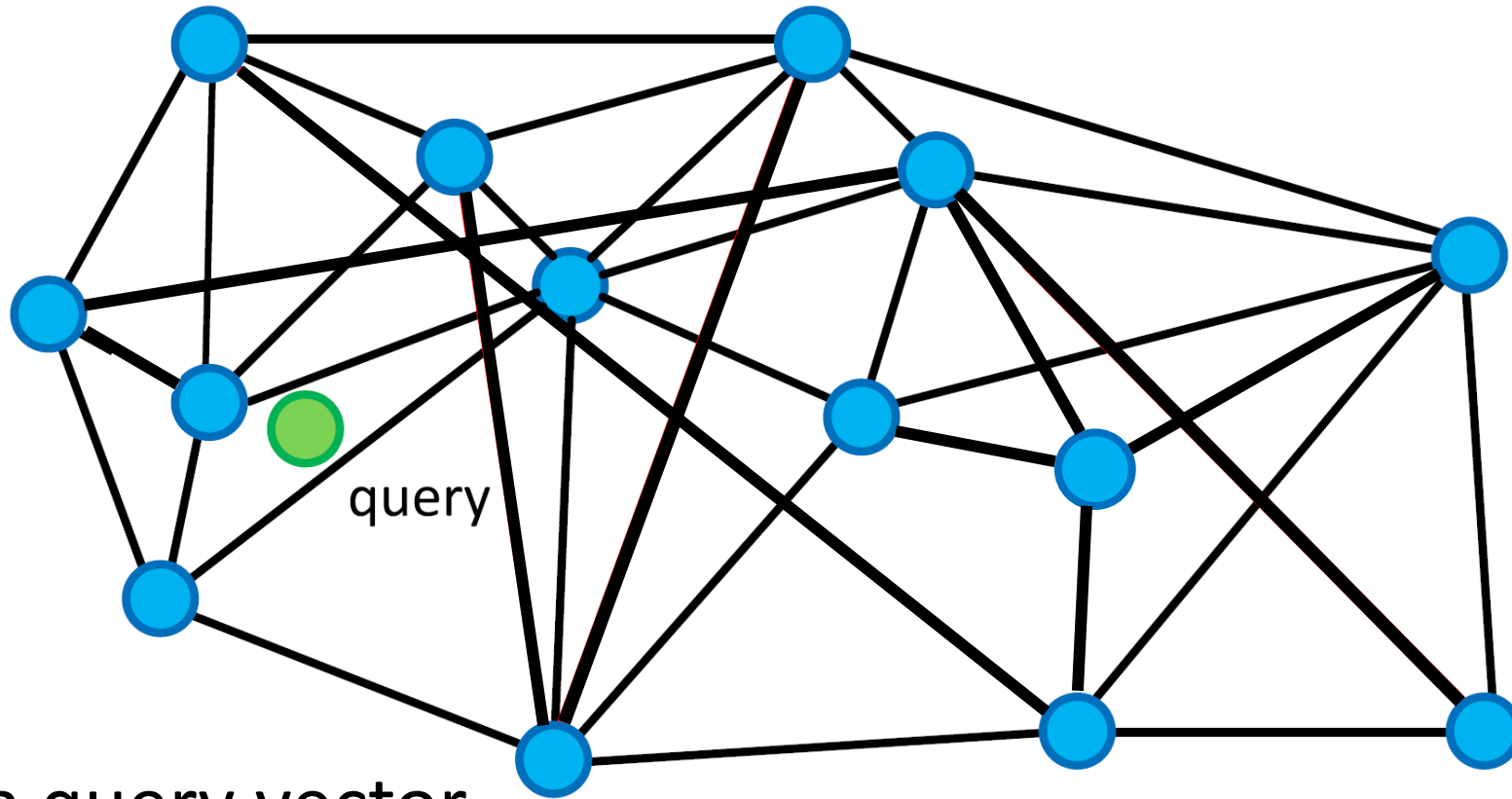
Graph of
 x_1, \dots, x_{90}



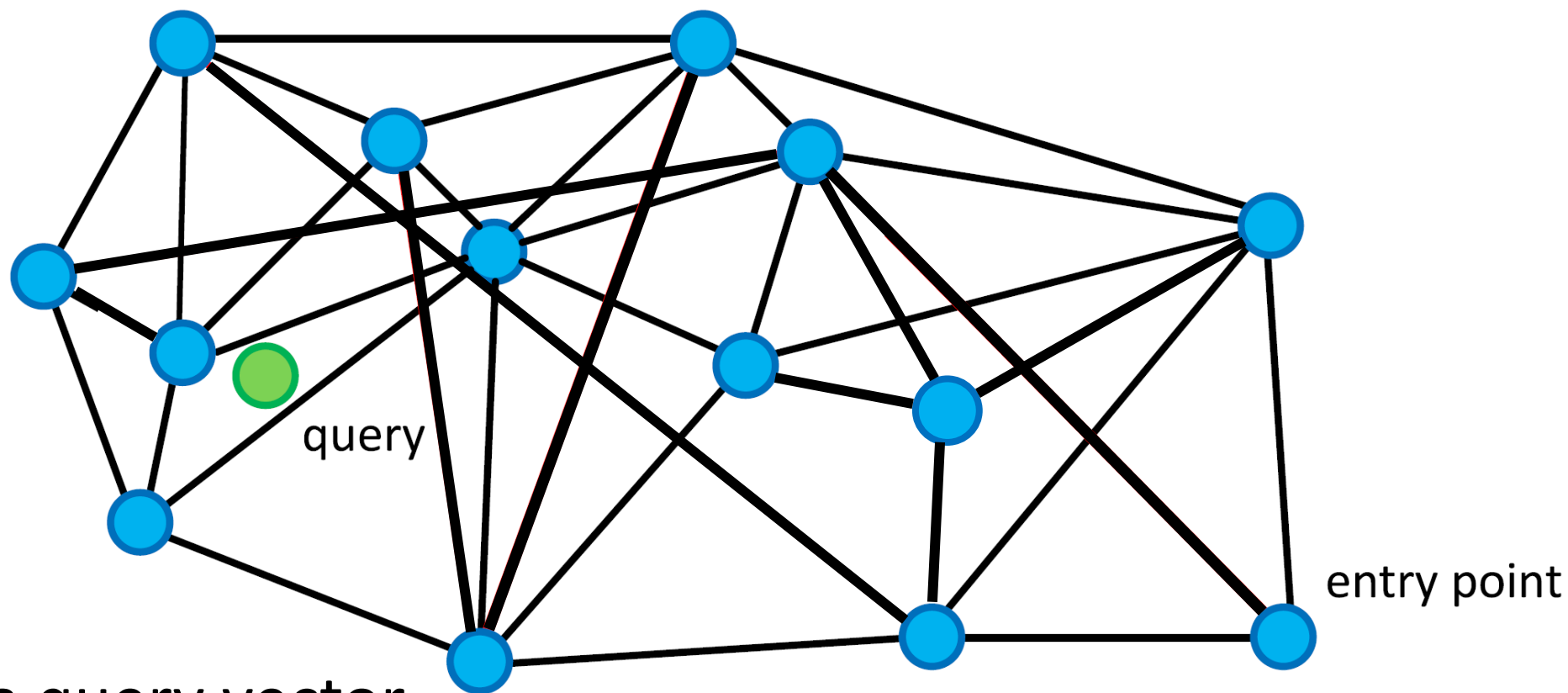
- Early links can be long
- Such long links encourage a large hop, making the fast convergence for search

- Each node is a database vector
- Given a new database vector, create new edges to neighbors

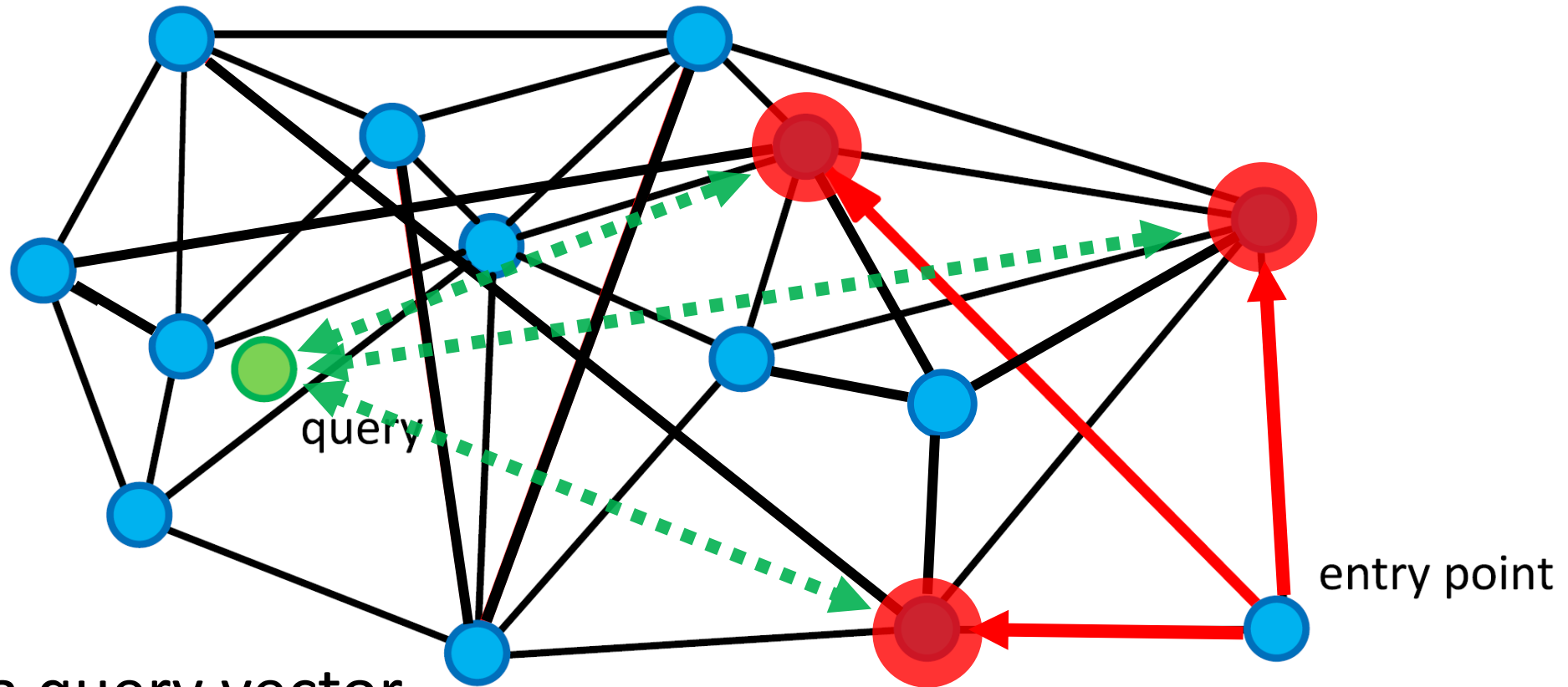




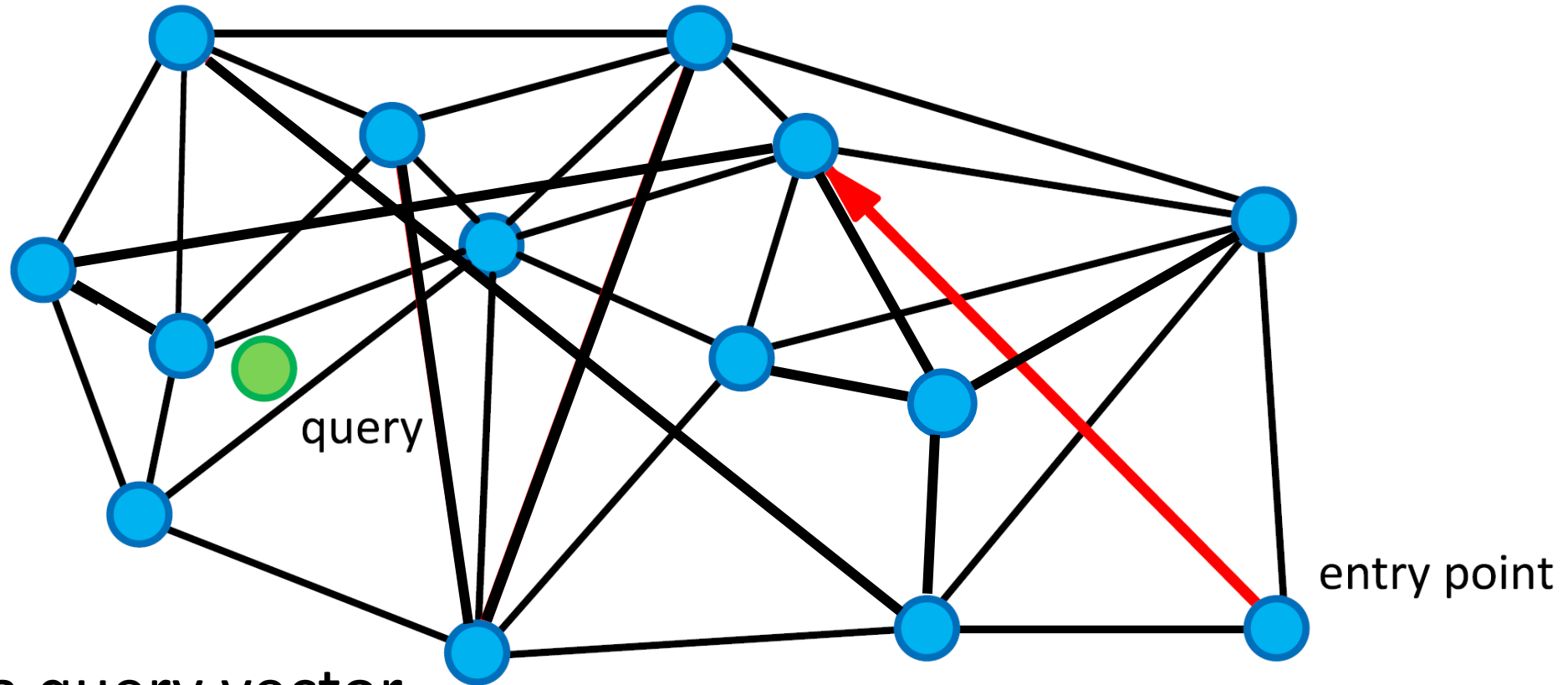
➤ Given a query vector



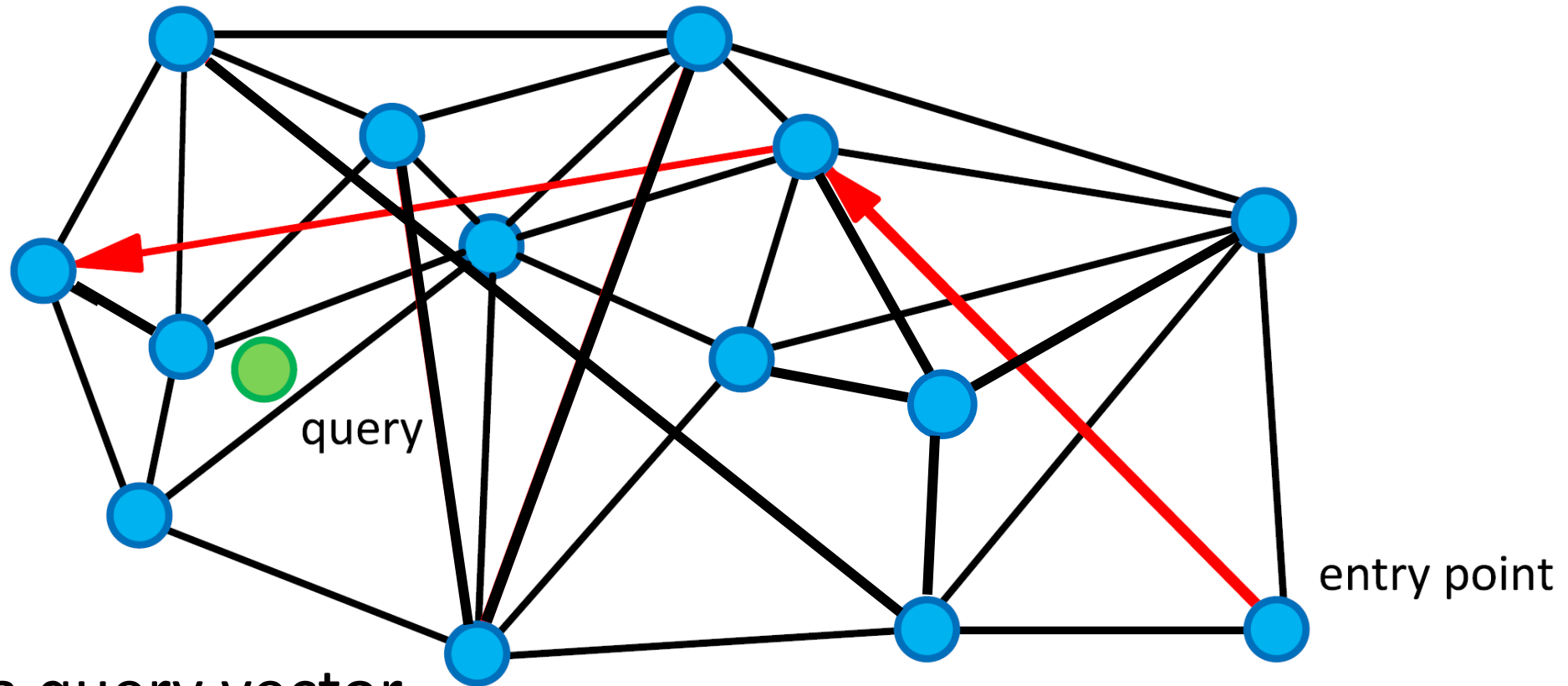
- Given a query vector
- Start from a random point



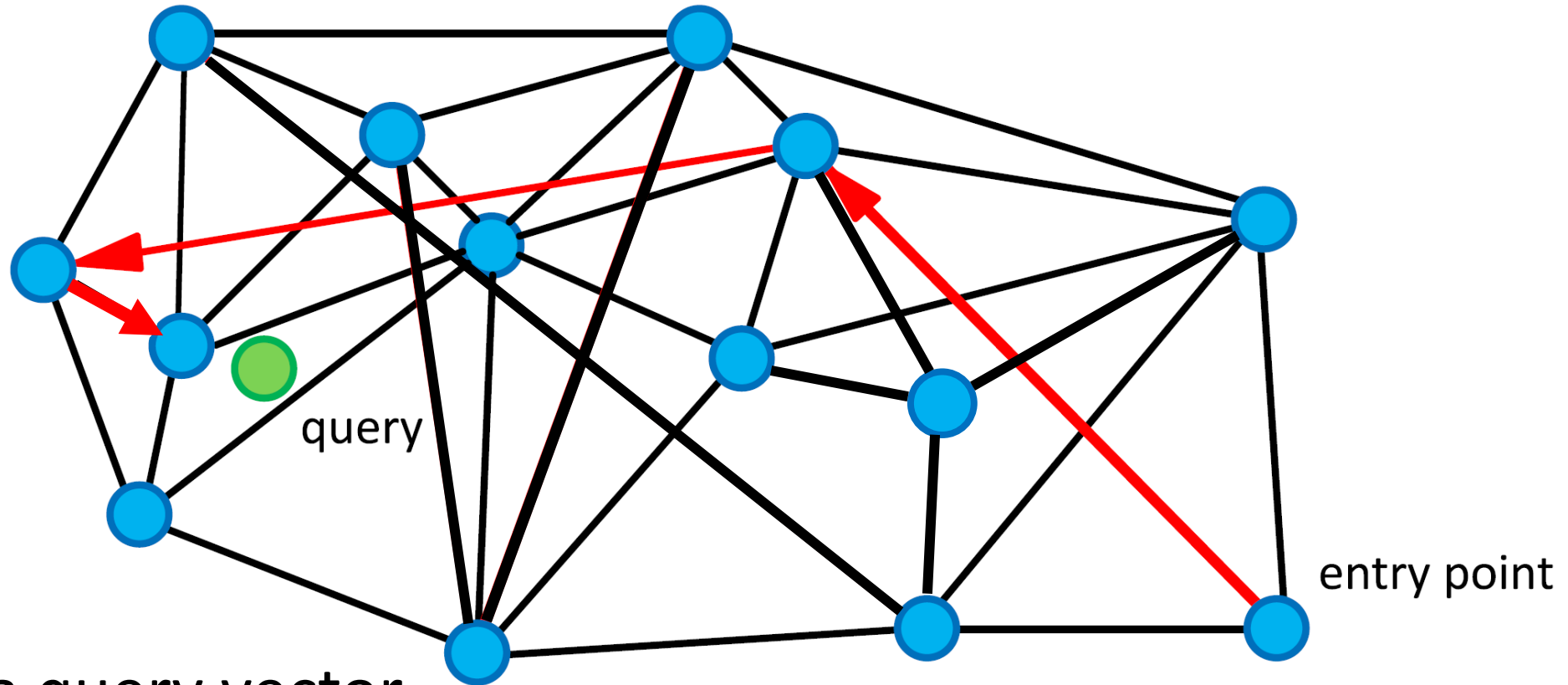
- Given a query vector
- Start from a random point
- From the connected nodes, find the closest one to the query



- Given a query vector
- Start from a random point
- From the connected nodes, find the closest one to the query



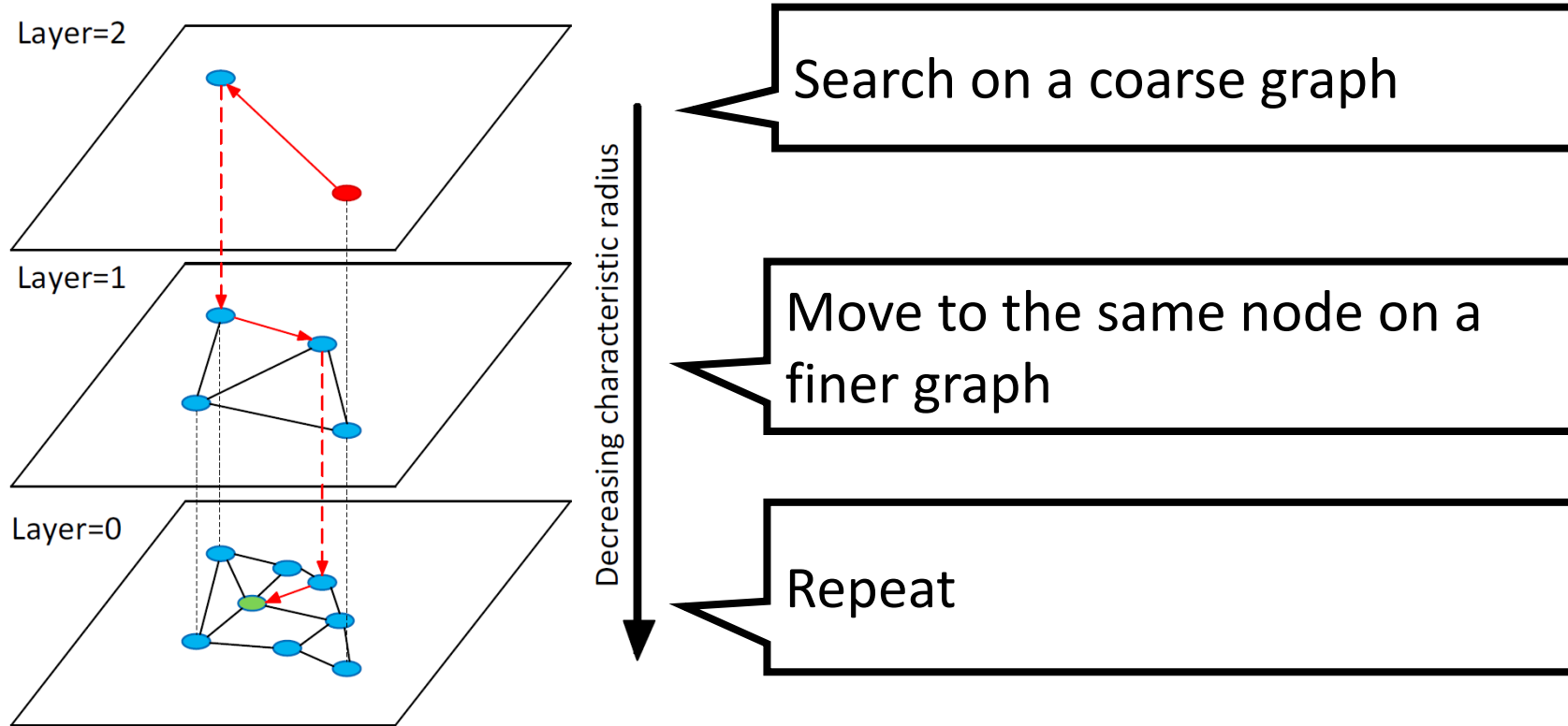
- Given a query vector
- Start from a random point
- From the connected nodes, find the closest one to the query
- Traverse in a greedy manner



- Given a query vector
- Start from a random point
- From the connected nodes, find the closest one to the query
- Traverse in a greedy manner

Extension: Hierarchical NSW; HNSW

- Construct the graph hierarchically [Malkov and Yashunin, TPAMI, 2019]
- This structure works pretty well for real-world data



[Malkov and Yashunin, TPAMI, 2019]

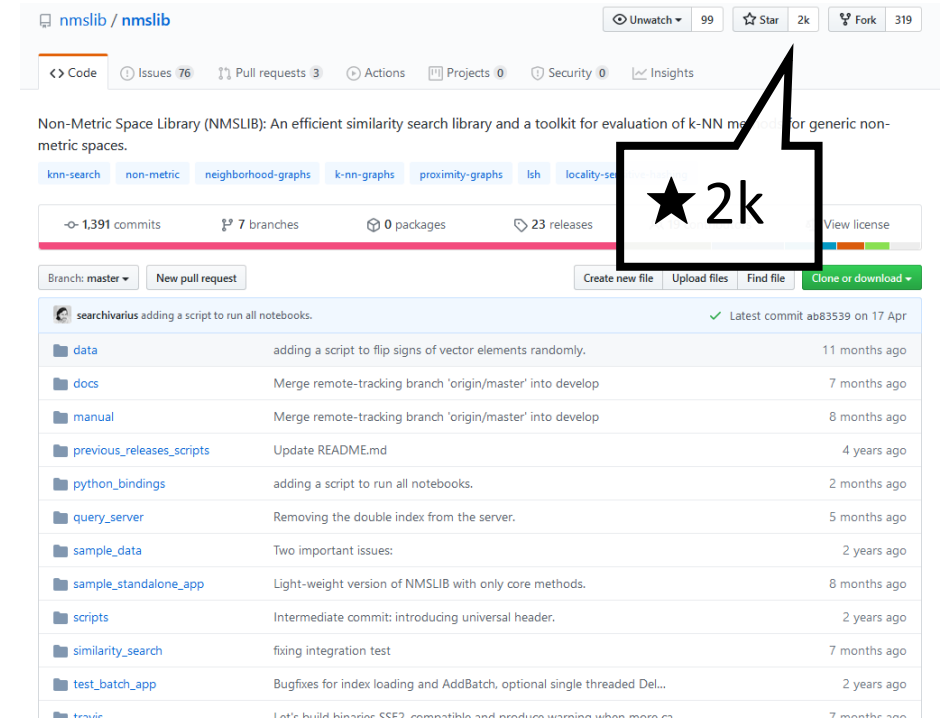
NMSLIB (Non-Metric Space Library)

<https://github.com/nmslib/nmslib>

```
$> pip install nmslib
```

```
index = nmslib.init(method='hnsw')
index.addDataPointBatch(X)
index.createIndex(params1)

index.setQueryTimeParams(params2)
index.knnQuery(q, topk)
```



nmslib / nmslib

Unwatch 99 Star 2k Fork 319

Code Issues 76 Pull requests 3 Actions Projects 0 Security 0 Insights

Non-Metric Space Library (NMSLIB): An efficient similarity search library and a toolkit for evaluation of k-NN methods for generic non-metric spaces.

kn-search non-metric neighborhood-graphs k-nn-graphs proximity-graphs lsh locality-sensitive hashing

1,391 commits 7 branches 0 packages 23 releases

Branch: master New pull request

Create new file Upload files Find file Clone or download

searchivarius adding a script to run all notebooks. Latest commit ab83539 on 17 Apr

data	adding a script to flip signs of vector elements randomly.	11 months ago
docs	Merge remote-tracking branch 'origin/master' into develop	7 months ago
manual	Merge remote-tracking branch 'origin/master' into develop	8 months ago
previous_releases_scripts	Update README.md	4 years ago
python_bindings	adding a script to run all notebooks.	2 months ago
query_server	Removing the double index from the server.	5 months ago
sample_data	Two important issues:	2 years ago
sample_standalone_app	Light-weight version of NMSLIB with only core methods.	8 months ago
scripts	Intermediate commit: introducing universal header.	2 years ago
similarity_search	fixing integration test	7 months ago
test_batch_app	Bugfixes for index loading and AddBatch, optional single threaded Del...	2 years ago
travis	Let's build binary SSE3 compatible and produce warning when more...	7 months ago

- 😊 The “hnsw” is the best method as of 2020 for million-scale data
- 😊 Simple interface
- 😊 If memory consumption is not the problem, try this
- 😞 Large memory consumption
- 😞 Data addition is not fast

Other implementations of HNSW

Hnswlib: <https://github.com/nmslib/hnswlib>

- Spin-off library from nmslib
- Include only hns w
- Simpler; may be useful if you want to extend hns w

Faiss: <https://github.com/facebookresearch/faiss>

- Libraries for PQ-based methods. Will introduce later
- This lib also includes hns w

Other graph-based approaches

➤ From Alibaba:

C. Fu et al., “Fast Approximate Nearest Neighbor Search with the Navigating Spreading-out Graph”, VLDB19

<https://github.com/ZJULearning/nsg>

➤ From Microsoft Research Asia. Used inside Bing:

J. Wang and S. Lin, “Query-Driven Iterated Neighborhood Graph Search for Large Scale Indexing”, ACMMM12 (This seems the backbone paper)

<https://github.com/microsoft/SPTAG>

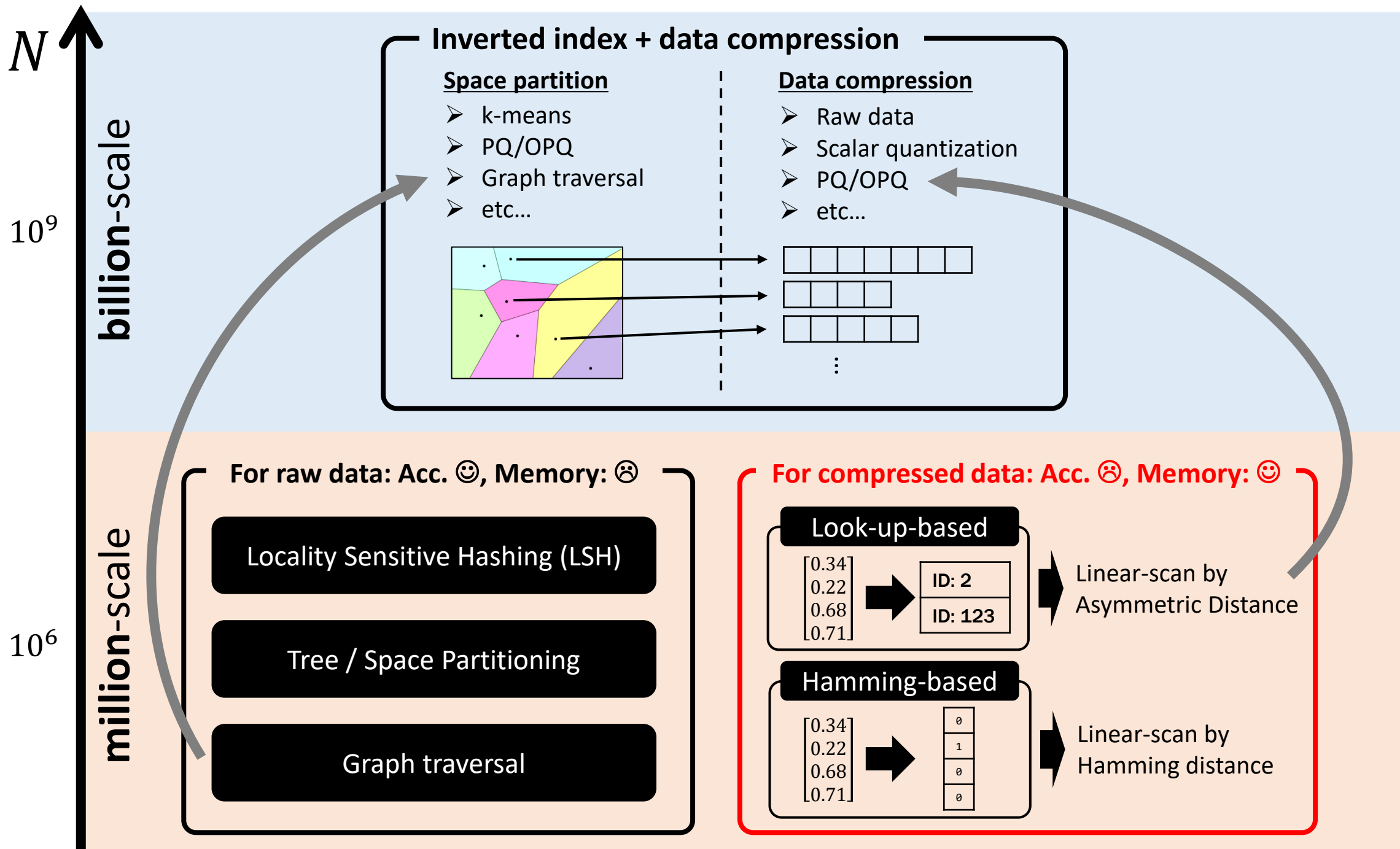
➤ From Yahoo Japan. Competing with NMSLIB for the 1st place of benchmark:

M. Iwasaki and D. Miyazaki, “Optimization of Indexing Based on k-Nearest Neighbor Graph for Proximity Search in High-dimensional Data”, arXiv18

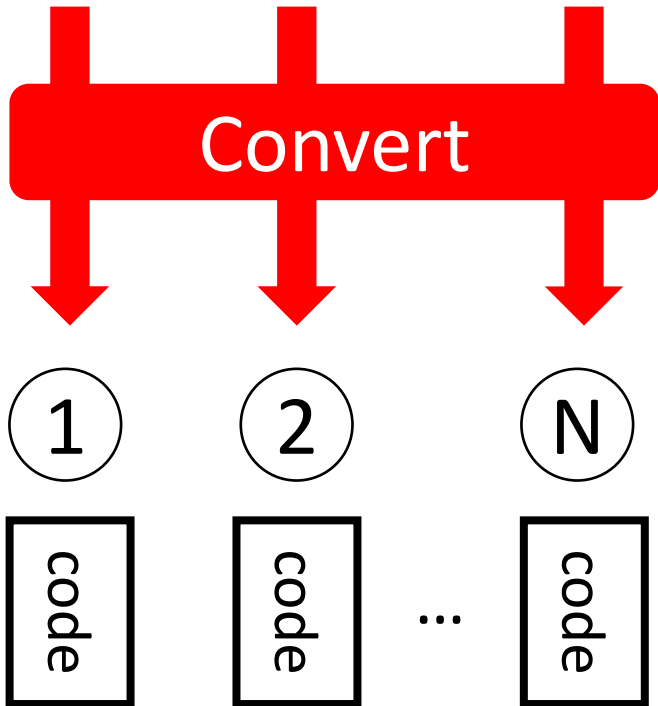
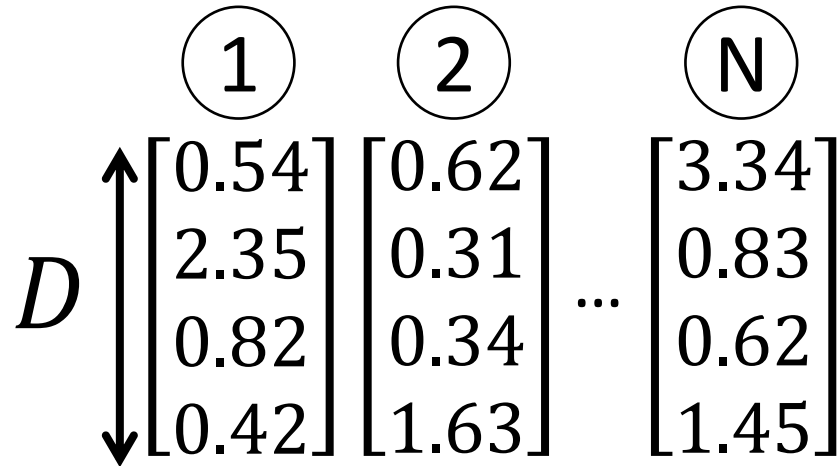
<https://github.com/yahoojapan/NGT>

Reference

- The original paper of Navigable Small World Graph: Y. Malkov et al., “Approximate Nearest Neighbor Algorithm based on Navigable Small World Graphs,” Information Systems 2013
- The original paper of Hierarchical Navigable Small World Graph: Y. Malkov and D. Yashunin, “Efficient and Robust Approximate Nearest Neighbor search using Hierarchical Navigable Small World Graphs,” IEEE TPAMI 2019

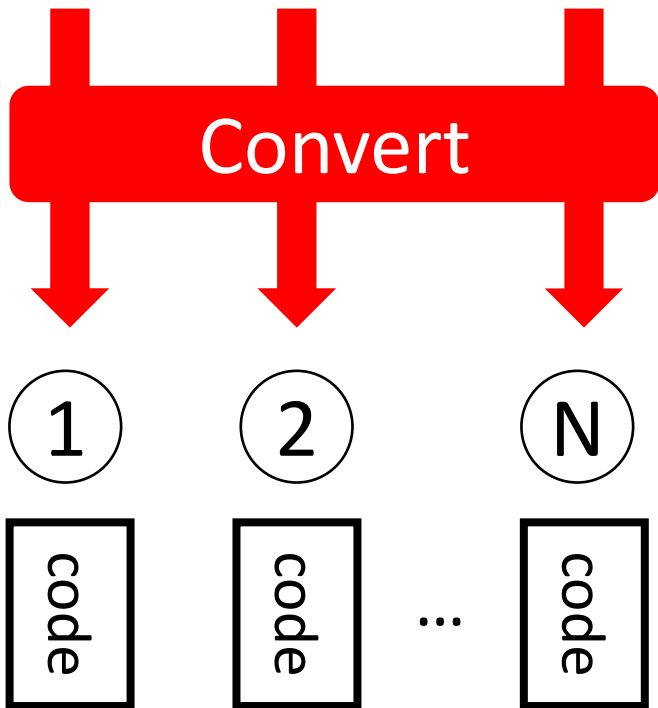
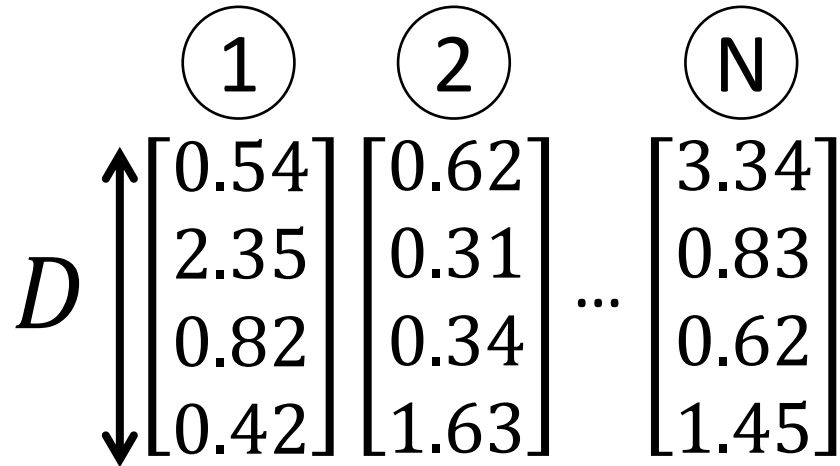


Basic idea



- Need $4ND$ byte to represent N real-valued vectors using floats
- If N or D is too large, we cannot read the data on memory
 - ✓ E.g., 512 GB for $D = 128, N = 10^9$
- Convert each vector to a **short-code**
- Short-code is designed as memory-efficient
 - ✓ E.g., 4 GB for the above example, with 32-bit code
- Run search for short-codes

Basic idea



➤ Need $4ND$ byte to represent N real-valued vectors

What kind of conversion is preferred?

➤ If N or D is too large, we cannot read the data on memory
E.g. 512 GB for $D = 128, N = 10^9$

1. The “distance” between two codes can be calculated (e.g., Hamming-distance)

➤ Convert each vector to a short-code

2. The distance can be computed quickly

➤ Short-code is designed as memory-efficient

3. That distance approximates the distance between the original vectors (e.g., L_2)

➤ Run search for short-codes

4. Sufficiently small length of codes can achieve the above three criteria

- Convert x to a B -bit binary vector:

$$f(x) = b \in \{0, 1\}^B$$

- Hamming distance

$$d_H(b_1, b_2) = |b_1 \oplus b_2| \sim d(x_1, x_2)$$

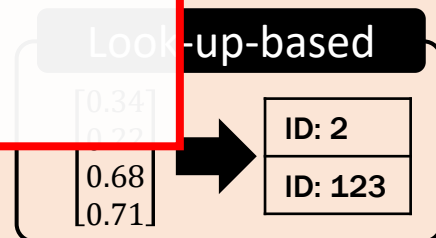
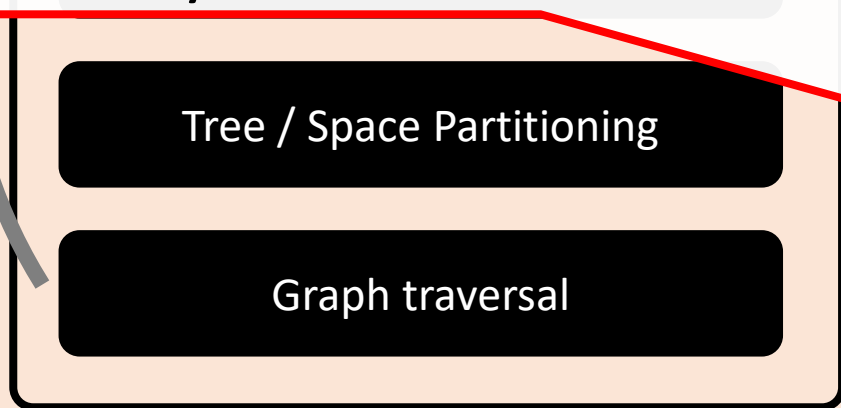
- A lot of methods:

- ✓ J. Wang et al., "Learning to Hash for Indexing Big Data - A Survey", Proc. IEEE 2015
- ✓ J. Wang et al., "A Survey on Learning to Hash", TPAMI 2018

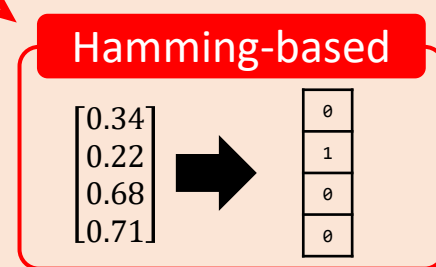
- Not the main scope of this tutorial;
PQ is usually more accurate

10^6

million-scale

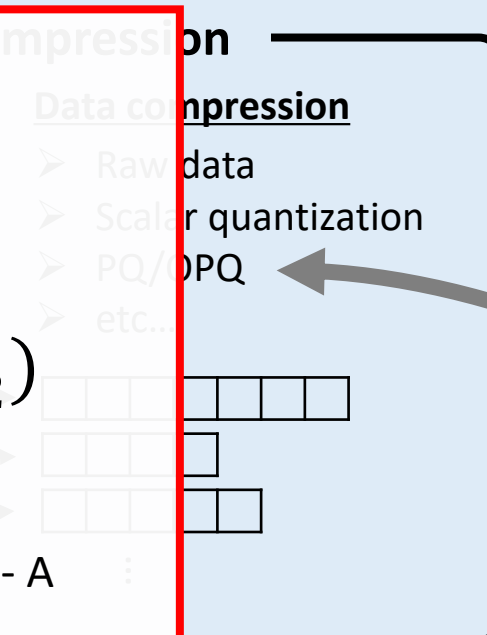


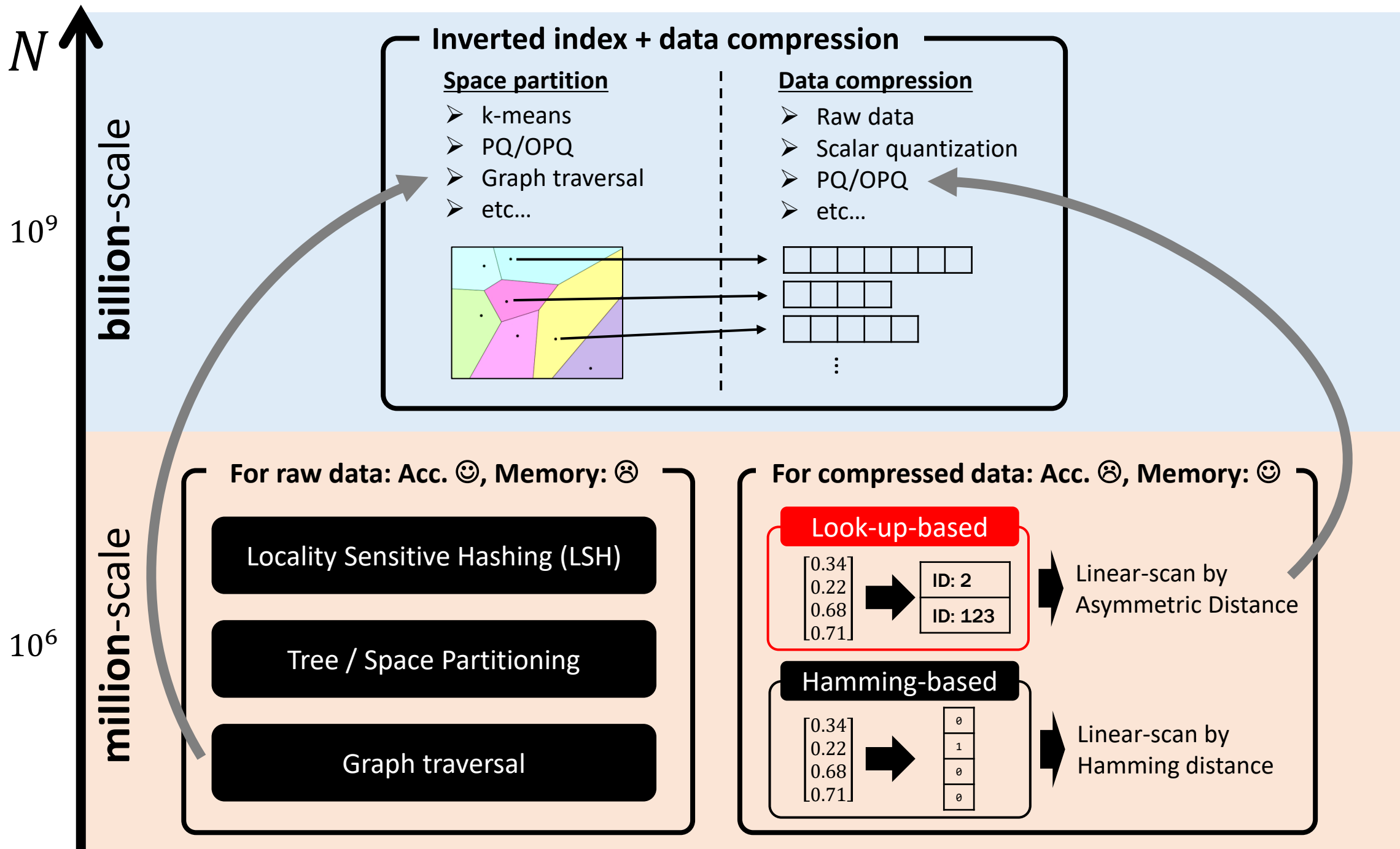
Linear-scan by Asymmetric Distance



Linear-scan by Hamming distance

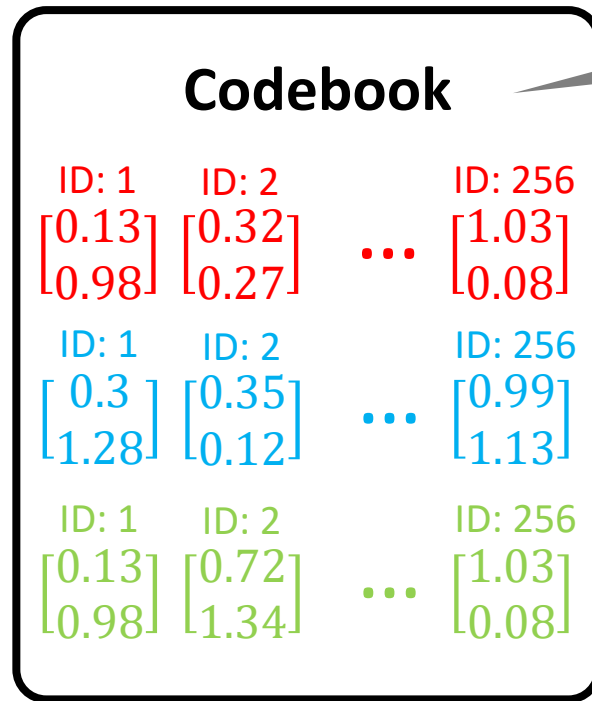
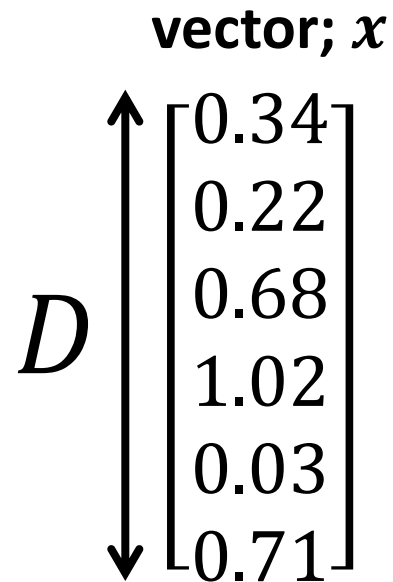
For compressed data: Acc. ☹️, Memory: 😊



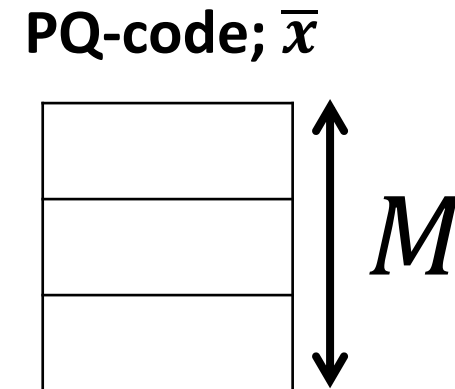


Product Quantization; PQ [Jégou, TPAMI 2011]

- Split a vector into sub-vectors, and quantize each sub-vector

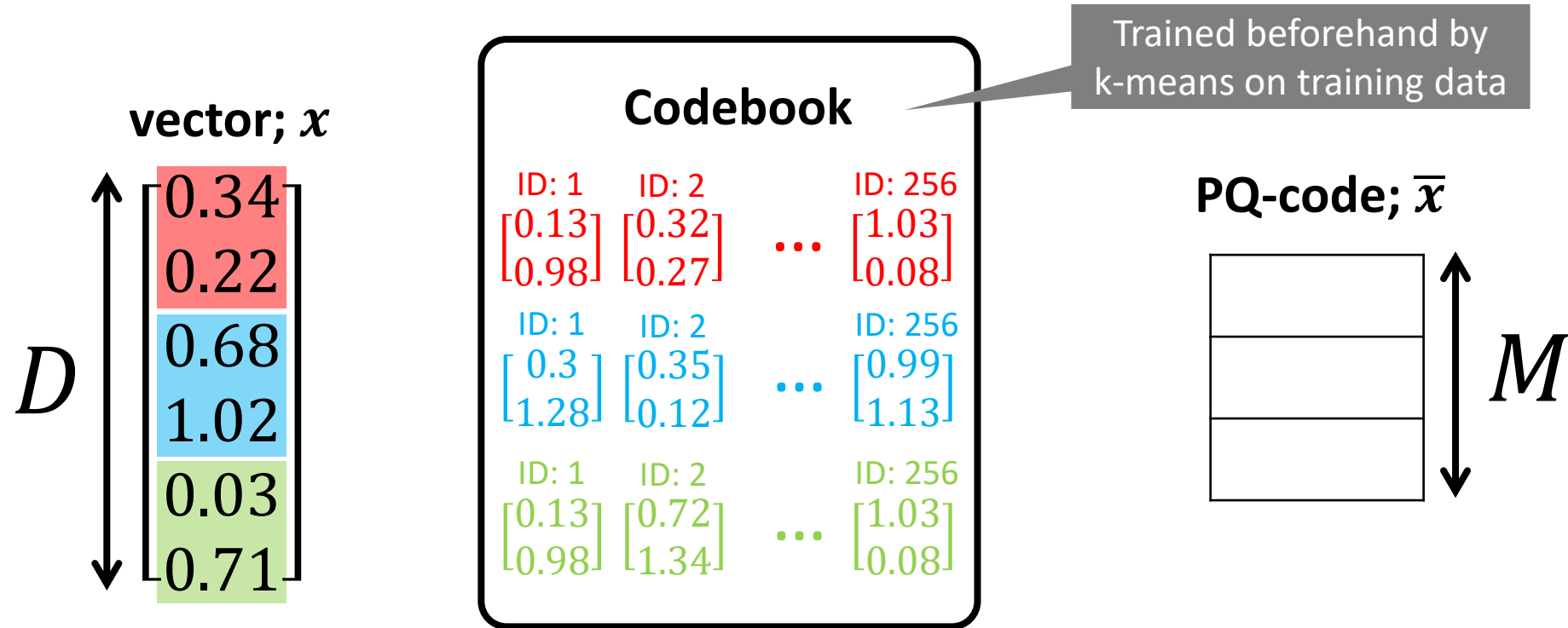


Trained beforehand by k-means on training data



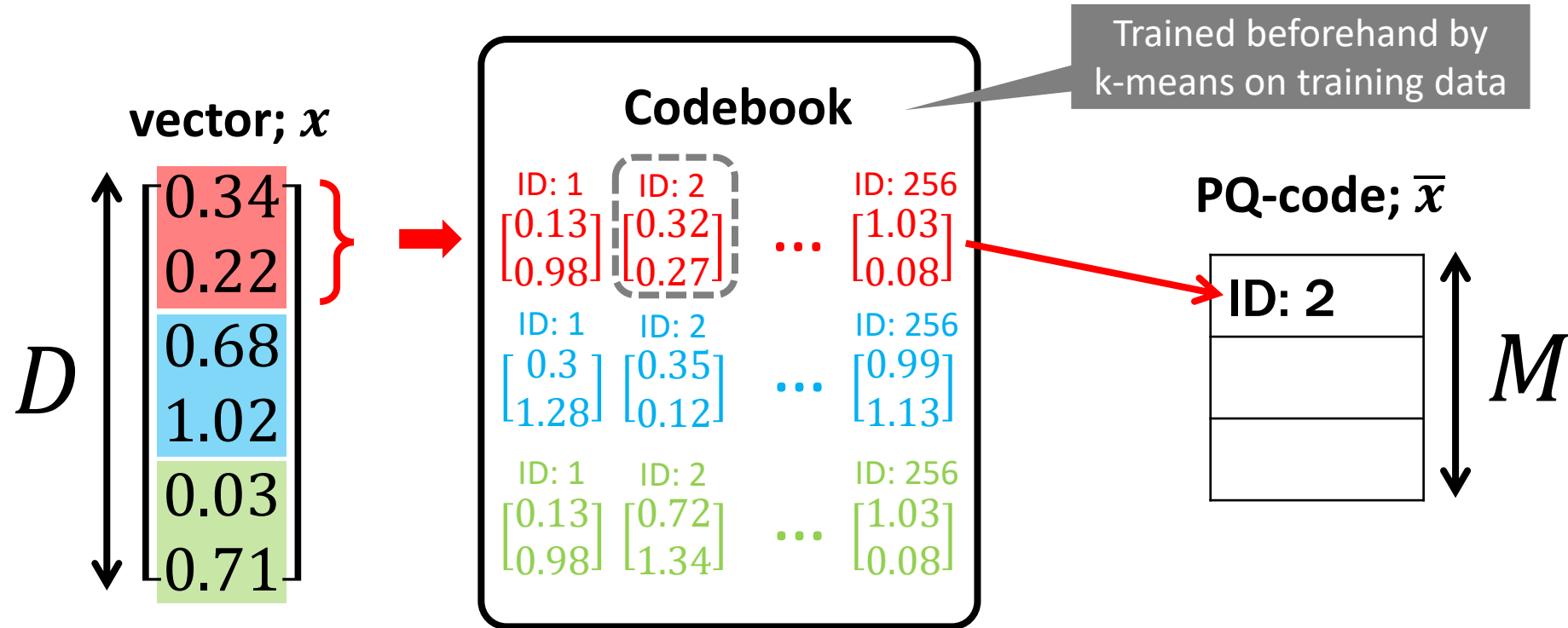
Product Quantization; PQ [Jégou, TPAMI 2011]

- Split a vector into sub-vectors, and quantize each sub-vector



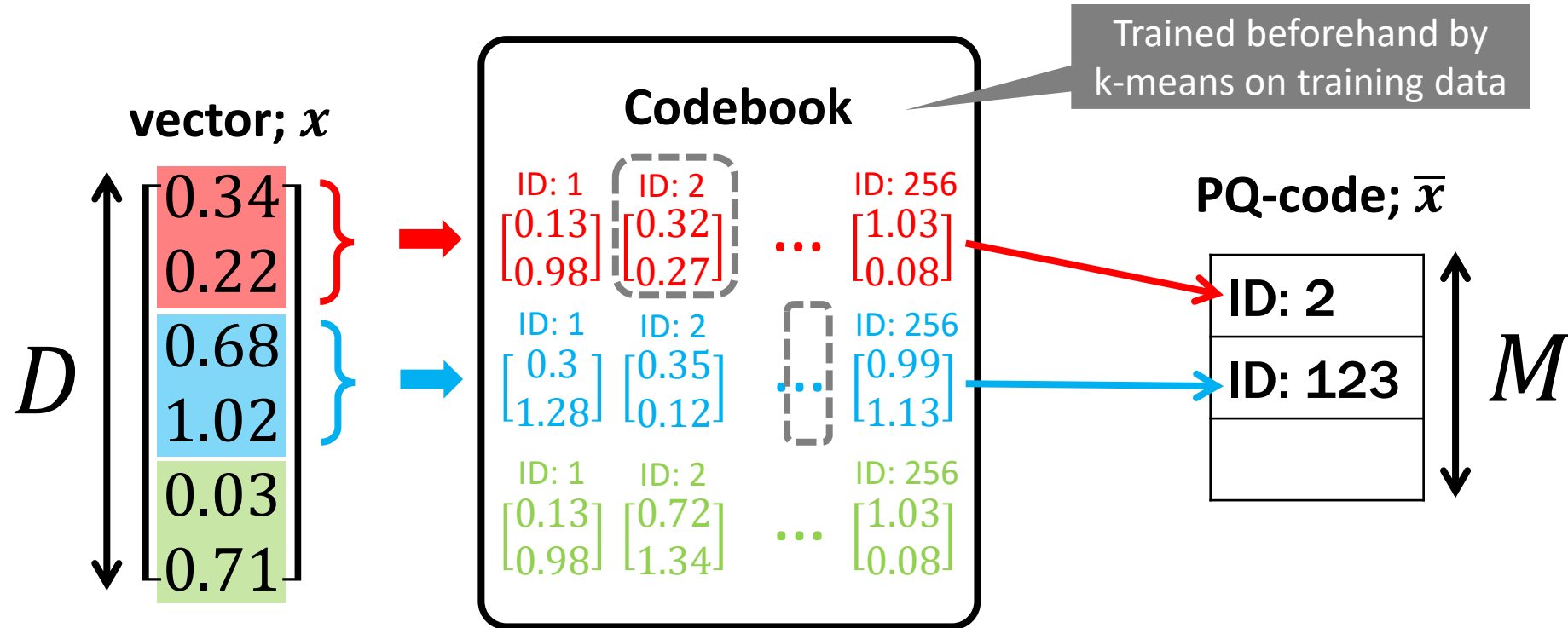
Product Quantization; PQ [Jégou, TPAMI 2011]

- Split a vector into sub-vectors, and quantize each sub-vector



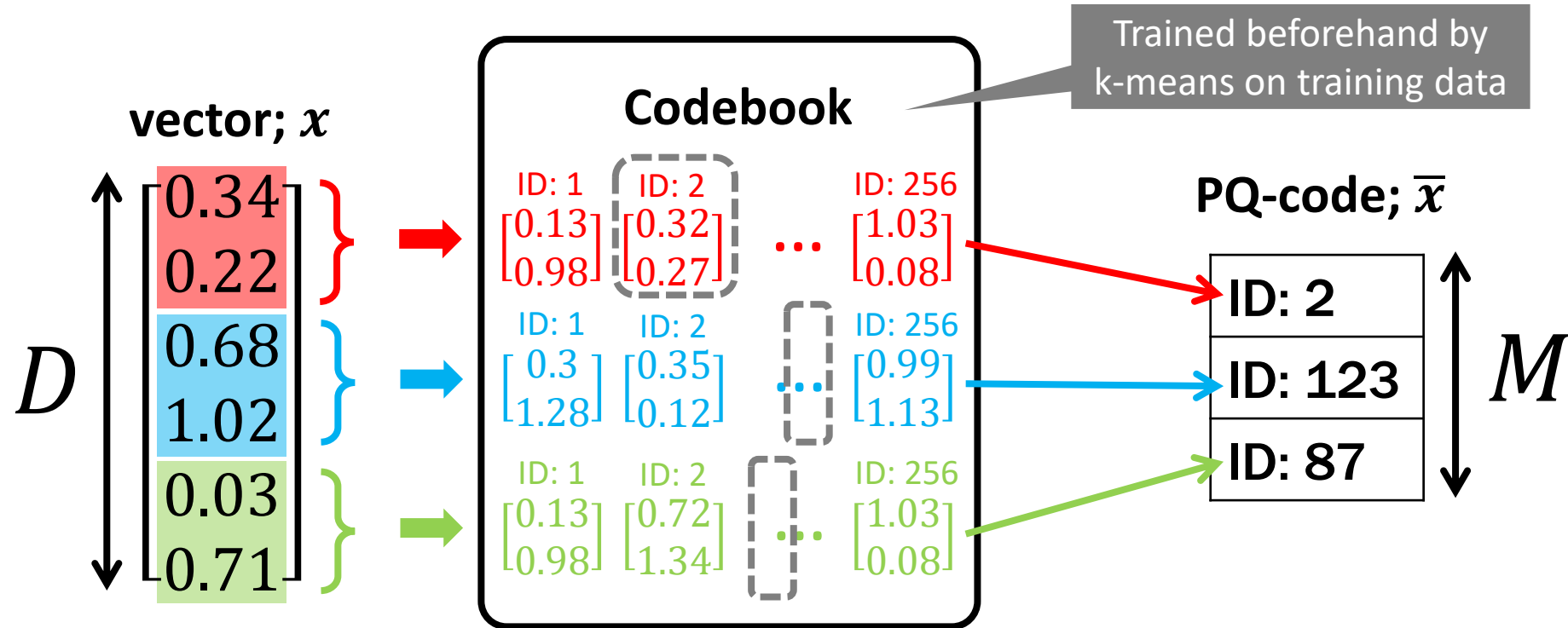
Product Quantization; PQ [Jégou, TPAMI 2011]

- Split a vector into sub-vectors, and quantize each sub-vector



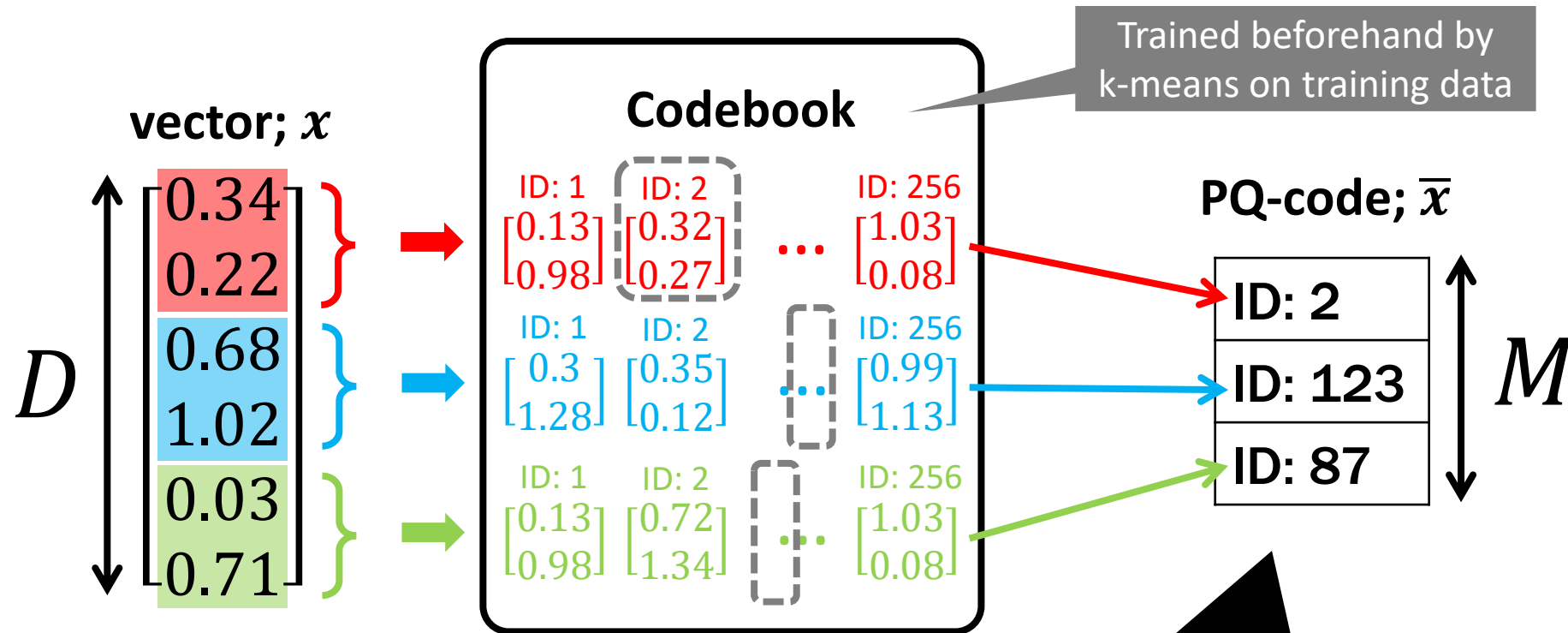
Product Quantization; PQ [Jégou, TPAMI 2011]

- Split a vector into sub-vectors, and quantize each sub-vector



Product Quantization; PQ [Jégou, TPAMI 2011]

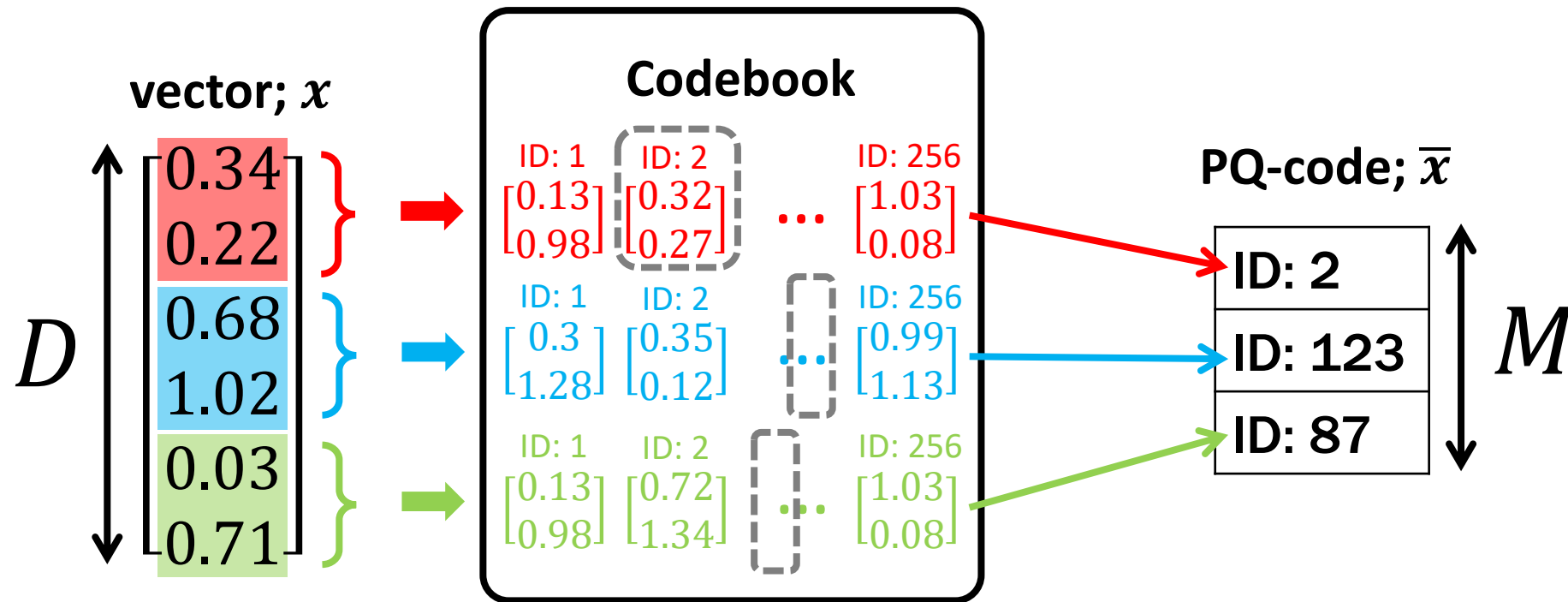
- Split a vector into sub-vectors, and quantize each sub-vector



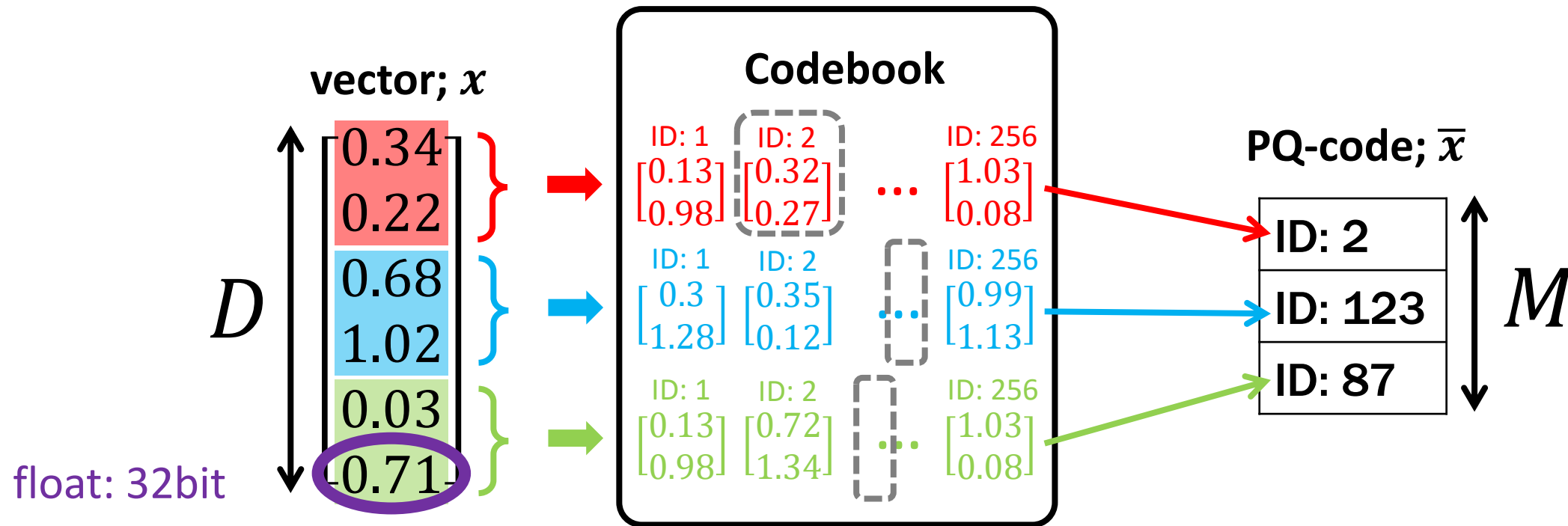
- Simple
- Memory efficient
- Distance can be estimated

Bar notation for PQ-code in this tutorial:
 $x \in \mathbb{R}^D \mapsto \bar{x} \in \{1, \dots, 256\}^M$

Product Quantization: **Memory efficient**



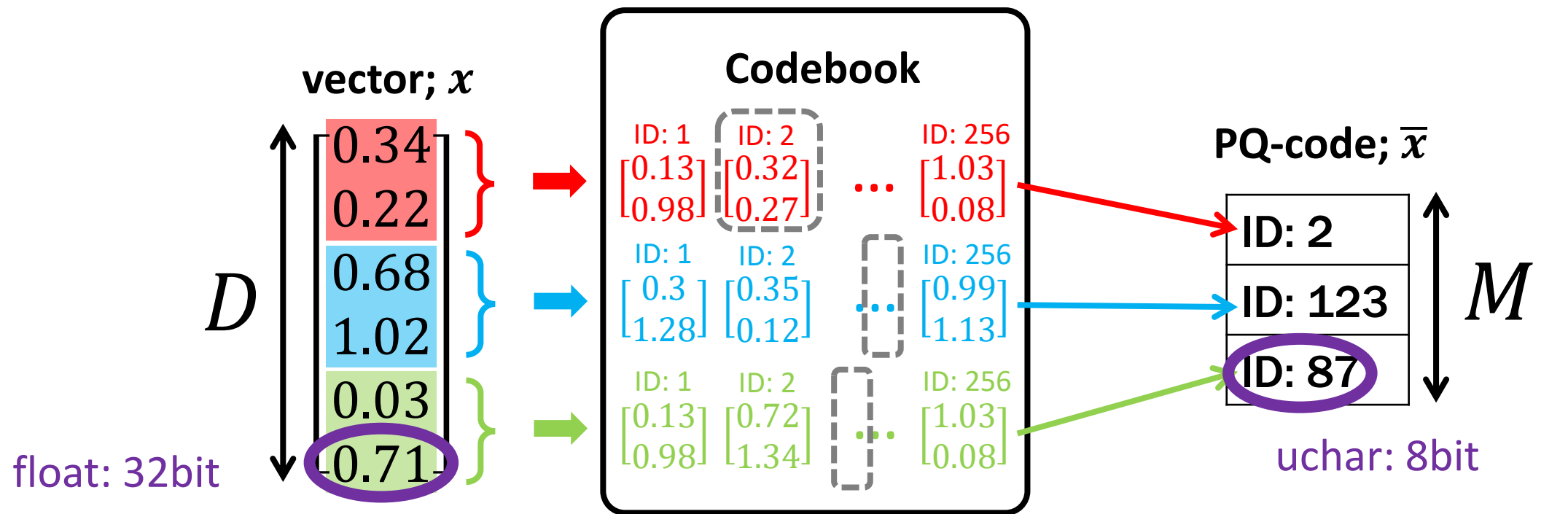
Product Quantization: **Memory efficient**



e.g., $D = 128$

$$128 \times 32 = 4096 \text{ [bit]}$$

Product Quantization: **Memory efficient**



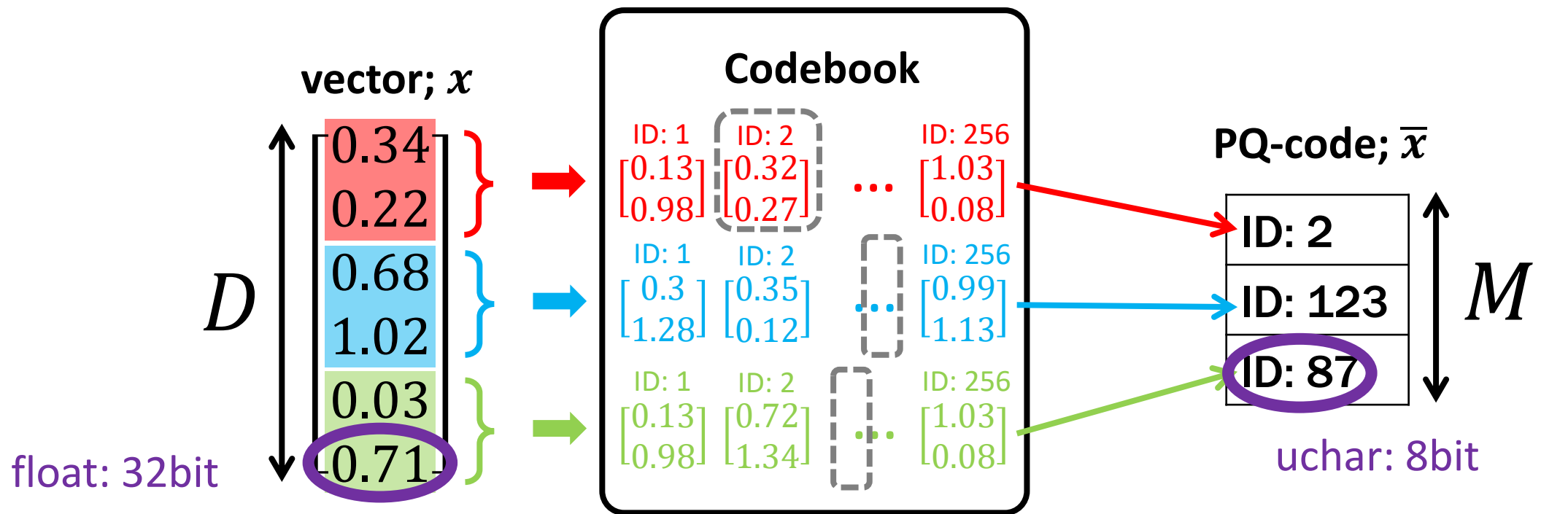
e.g., $D = 128$

$$128 \times 32 = 4096 \text{ [bit]}$$

e.g., $M = 8$

$$8 \times 8 = 64 \text{ [bit]}$$

Product Quantization: **Memory efficient**



e.g., $D = 128$

$$128 \times 32 = 4096 \text{ [bit]}$$

e.g., $M = 8$

$$8 \times 8 = 64 \text{ [bit]}$$

$1/64$

Product Quantization: Distance estimation

Query; $q \in \mathbb{R}^D$

$$\begin{bmatrix} 0.34 \\ 0.22 \\ 0.68 \\ 1.02 \\ 0.03 \\ 0.71 \end{bmatrix}$$

Database vectors

x_1

$$\begin{bmatrix} 0.54 \\ 2.35 \\ 0.82 \\ 0.42 \\ 0.14 \\ 0.32 \end{bmatrix}$$

x_2

$$\begin{bmatrix} 0.62 \\ 0.31 \\ 0.34 \\ 1.63 \\ 1.43 \\ 0.74 \end{bmatrix}$$

...

x_N

$$\begin{bmatrix} 3.34 \\ 0.83 \\ 0.62 \\ 1.45 \\ 0.12 \\ 2.32 \end{bmatrix}$$

Product Quantization: Distance estimation

Query; $q \in \mathbb{R}^D$

$\begin{bmatrix} 0.34 \\ 0.22 \\ 0.68 \\ 1.02 \\ 0.03 \\ 0.71 \end{bmatrix}$

Database vectors

$\begin{matrix} x_1 & x_2 & & x_N \\ \begin{bmatrix} 0.54 \\ 2.35 \\ 0.82 \\ 0.42 \\ 0.14 \\ 0.32 \end{bmatrix} & \begin{bmatrix} 0.62 \\ 0.31 \\ 0.34 \\ 1.63 \\ 1.43 \\ 0.74 \end{bmatrix} & \dots & \begin{bmatrix} 3.34 \\ 0.83 \\ 0.62 \\ 1.45 \\ 0.12 \\ 2.32 \end{bmatrix} \end{matrix}$

Product
quantization

Product Quantization: Distance estimation

Query; $q \in \mathbb{R}^D$

$\begin{bmatrix} 0.34 \\ 0.22 \\ 0.68 \\ 1.02 \\ 0.03 \\ 0.71 \end{bmatrix}$

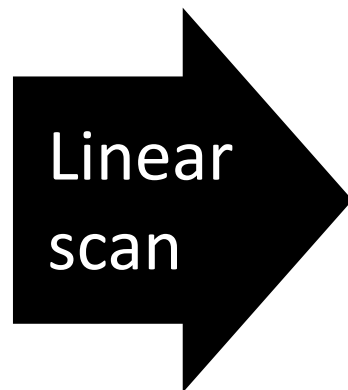
$\bar{x}_1 \in \{1, \dots, 256\}^M$

\bar{x}_1	\bar{x}_2		\bar{x}_N
ID: 42	ID: 221		ID: 99
ID: 67	ID: 143	...	ID: 234
ID: 92	ID: 34		ID: 3

Product Quantization: Distance estimation

Query; $\mathbf{q} \in \mathbb{R}^D$

$\begin{bmatrix} 0.34 \\ 0.22 \\ 0.68 \\ 1.02 \\ 0.03 \\ 0.71 \end{bmatrix}$



$\bar{\mathbf{x}}_1 \in \{1, \dots, 256\}^M$

$\bar{\mathbf{x}}_1$	$\bar{\mathbf{x}}_2$...	$\bar{\mathbf{x}}_N$
ID: 42	ID: 221		ID: 99
ID: 67	ID: 143		ID: 234
ID: 92	ID: 34		ID: 3

Asymmetric distance

- $d(\mathbf{q}, \mathbf{x})^2$ can be efficiently approximated by $d_A(\mathbf{q}, \bar{\mathbf{x}})^2$
- Lookup-trick: Looking up pre-computed distance-tables
- Linear-scan by d_A

Not pseudo codes

```
import numpy as np
from scipy.cluster.vq import vq, kmeans2
from scipy.spatial.distance import cdist

def train(vec, M):
    Ds = int(vec.shape[1] / M) #  $D_s = D / M$ 
    #  $\text{codeword}[m][k] = \mathbf{c}_k^m$ 
    codeword = np.empty((M, 256, Ds), np.float32)

    for m in range(M):
        vec_sub = vec[:, m * Ds : (m + 1) * Ds]
        codeword[m], label = kmeans2(vec_sub, 256)

    return codeword

def encode(codeword, vec): #  $\text{vec} = \{\mathbf{x}_n\}_{n=1}^N$ 
    M, _K, Ds = codeword.shape
    #  $\text{pqcode}[n] = \mathbf{i}(\mathbf{x}_n)$ ,  $\text{pqcode}[n][m] = i^m(\mathbf{x}_n)$ 
    pqcode = np.empty((vec.shape[0], M), np.uint8)

    for m in range(M): # Eq. (2) and Eq. (3)
        vec_sub = vec[:, m * Ds : (m + 1) * Ds]
        pqcode[:, m], dist = vq(vec_sub, codeword[m])

    return pqcode
```

```
def search(codeword, pqcode, query):
    M, _K, Ds = codeword.shape
    #  $\text{dist\_table} = D(m, k)$ 
    dist_table = np.empty((M, 256), np.float32)

    for m in range(M):
        query_sub = query[m * Ds : (m + 1) * Ds]
        dist_table[m, :] = cdist([query_sub],
            ↪ codeword[m], 'sqeuclidean')[0] # Eq. (5)

    # Eq. (6)
    dist = np.sum(dist_table[range(M), pqcode], axis=1)

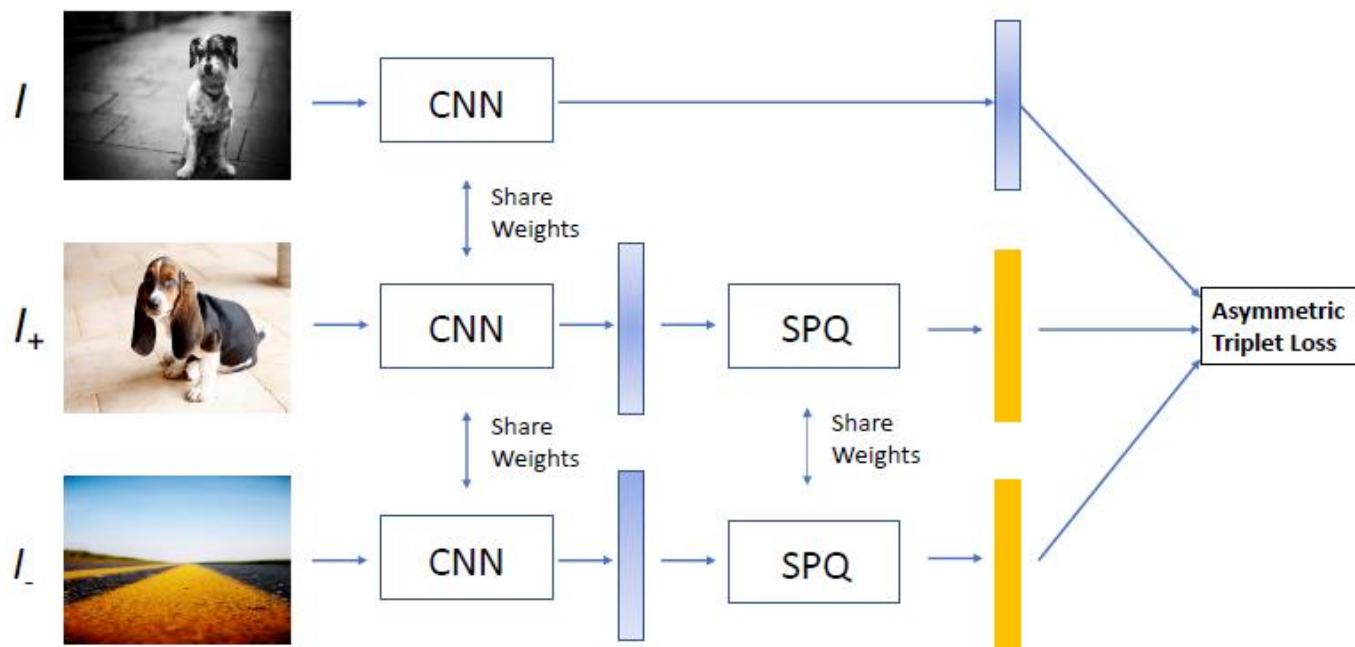
    return dist

if __name__ == "__main__":
    # Read  $\text{vec\_train}$ ,  $\text{vec}$  ( $\{\mathbf{x}_n\}_{n=1}^N$ ), and  $\text{query}$  ( $\mathbf{y}$ )
    M = 4
    codeword = train(vec_train, M)
    pqcode = encode(codeword, vec)
    dist = search(codeword, pqcode, query)
    print(dist)
```

- Only tens of lines in Python
- Pure Python library: nanopq <https://github.com/matsui528/nanopq>
- `pip install nanopq`

Deep PQ

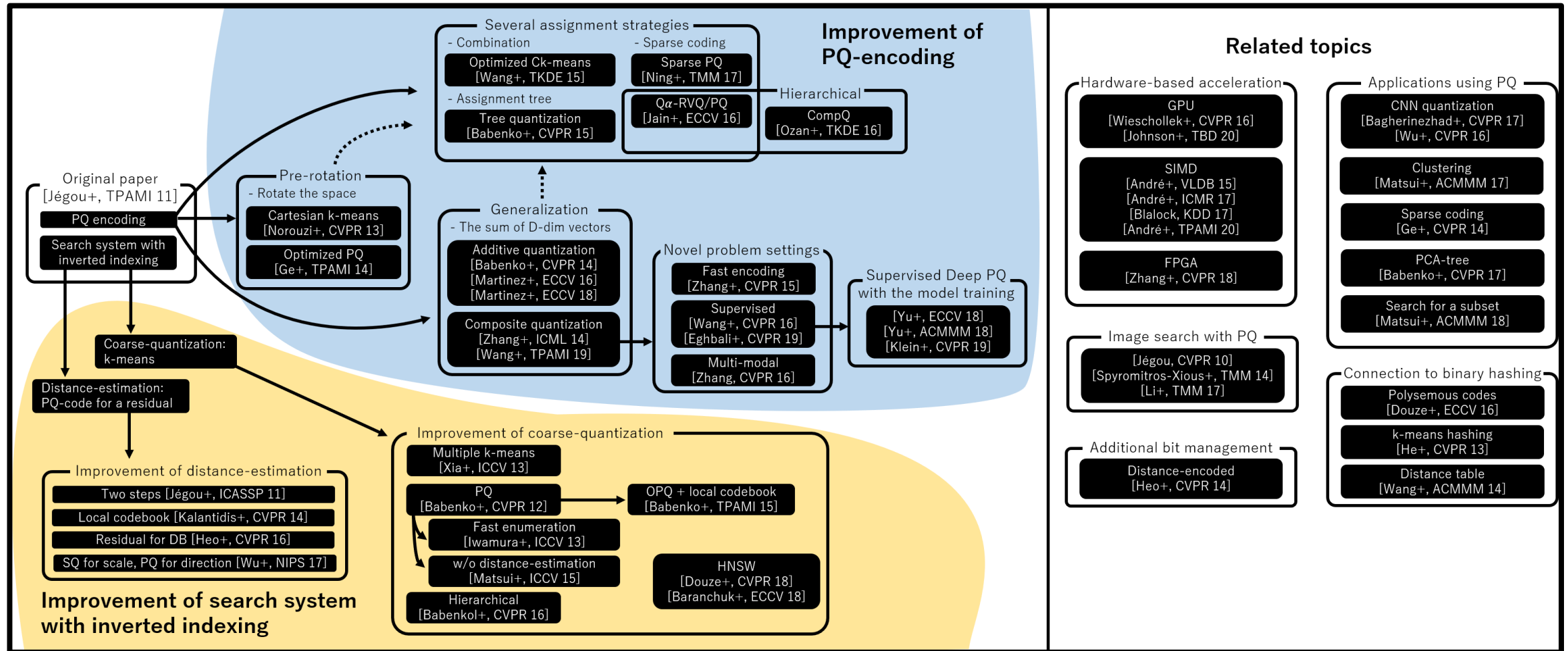
- Supervised search (unlike the original PQ)
- Base-CNN + PQ-like-layer + Some-loss
- Need class information



From T. Yu et al., "Product Quantization Network for Fast Image Retrieval", ECCV 18

- T. Yu et al., "Product Quantization Network for Fast Image Retrieval", ECCV 18, IJCV20
- L. Yu et al., "Generative Adversarial Product Quantisation", ACM MM 18
- B. Klein et al., "End-to-End Supervised Product Quantization for Image Search and Retrieval", CVPR 19

More extensive survey for PQ



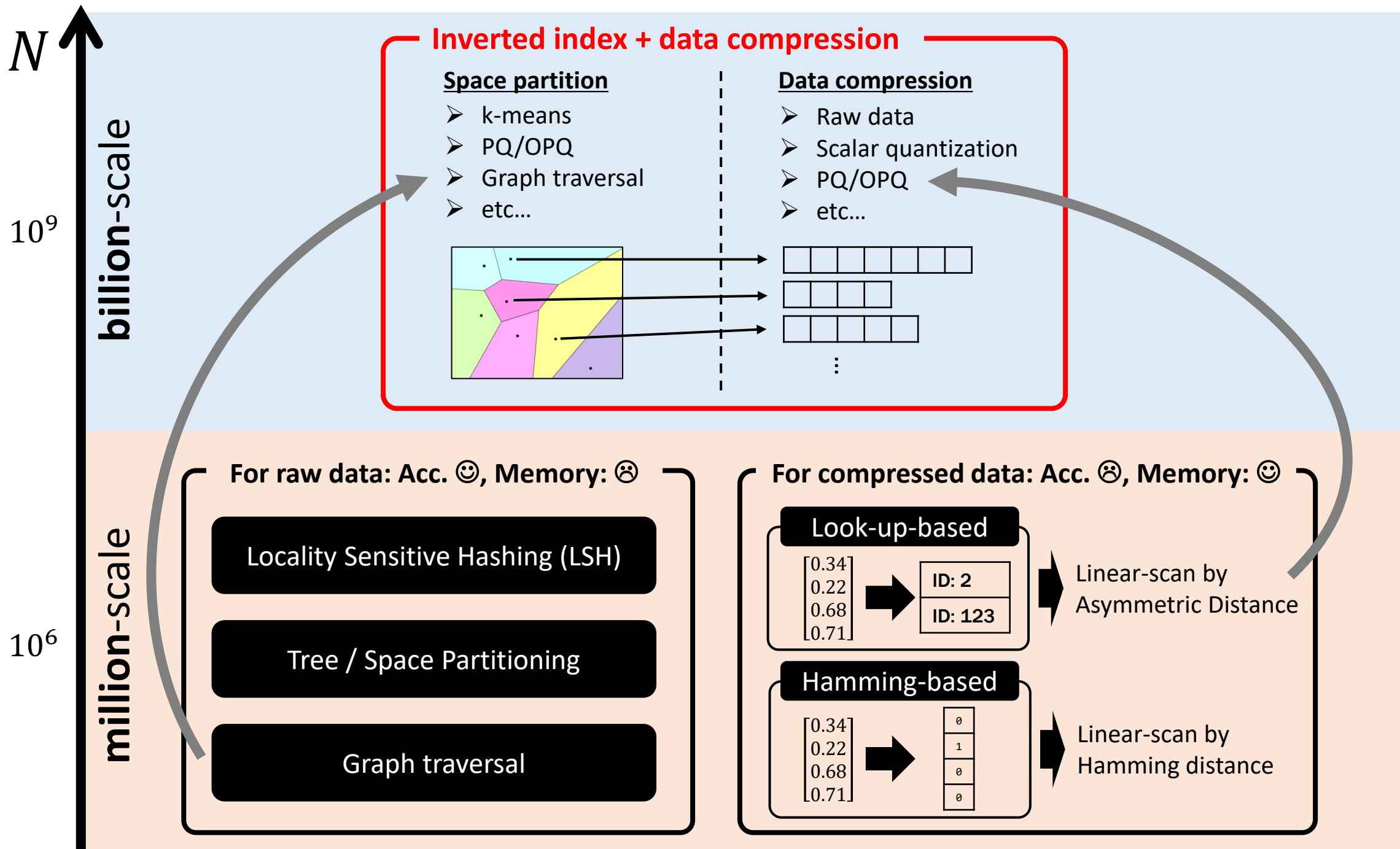
➤ <https://github.com/facebookresearch/faiss/wiki#research-foundations-of-faiss>

➤ http://yusukematsui.me/project/survey_pq/survey_pq_jp.html

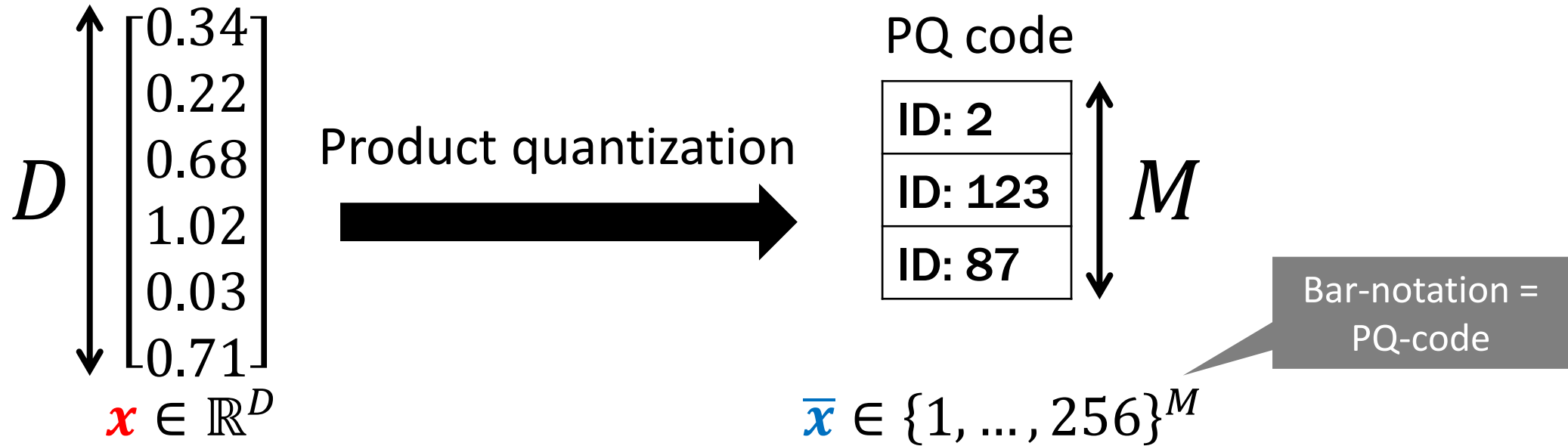
➤ Y. Matsui, Y. Uchida, H. Jégou, S. Satoh “A Survey of Product Quantization”, ITE 2018.

Hamming-based vs Look-up-based

	<div style="background-color: black; color: white; padding: 5px; border-radius: 10px; display: inline-block;">Hamming-based</div>	<div style="background-color: black; color: white; padding: 5px; border-radius: 10px; display: inline-block;">Look-up-based</div>
Representation	Binary code : $\{0, 1\}^B$	PQ code : $\{1, \dots, 256\}^M$
Distance	Hamming distance	Asymmetric distance
Approximation	😊	😊😊
Runtime	😊😊	😊
Pros	No auxiliary structure	Can reconstruct the original vector
Cons	Cannot reconstruct the original vector	Require an auxiliary structure (codebook)



Inverted index + PQ: Recap the notation

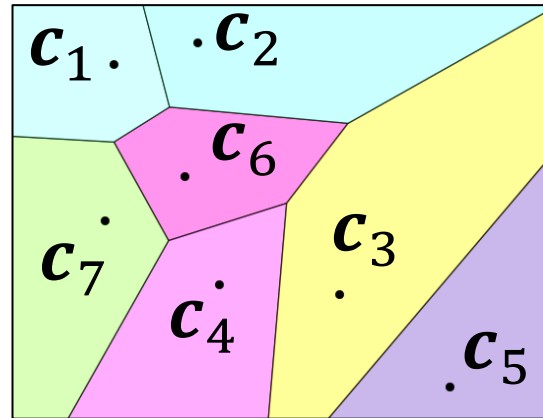


- Suppose $\mathbf{q}, \mathbf{x} \in \mathbb{R}^D$, where \mathbf{x} is quantized to $\bar{\mathbf{x}}$
- $d(\mathbf{q}, \mathbf{x})^2$ can be efficiently approximated by $\bar{\mathbf{x}}$:

$$d(\mathbf{q}, \mathbf{x})^2 \sim d_A(\mathbf{q}, \bar{\mathbf{x}})^2$$

Just by a PQ-code.
Not the original vector \mathbf{x}

Inverted index + PQ: Record



Coarse quantizer

$$k = 1$$

$$k = 2$$

⋮

$$k = K$$

Prepare a **coarse quantizer**

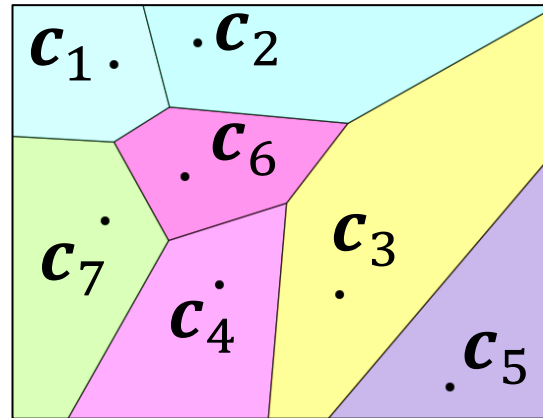
✓ Split the space into K sub-spaces

✓ $\{\mathbf{c}_k\}_{k=1}^K$ are created by running k-means on training data

Inverted index + PQ: Record

Record x_1

$\begin{bmatrix} 1.02 \\ 0.73 \\ 0.56 \\ 1.37 \\ 1.37 \\ 0.72 \end{bmatrix}$
 x_1



Coarse quantizer

$$k = 1$$

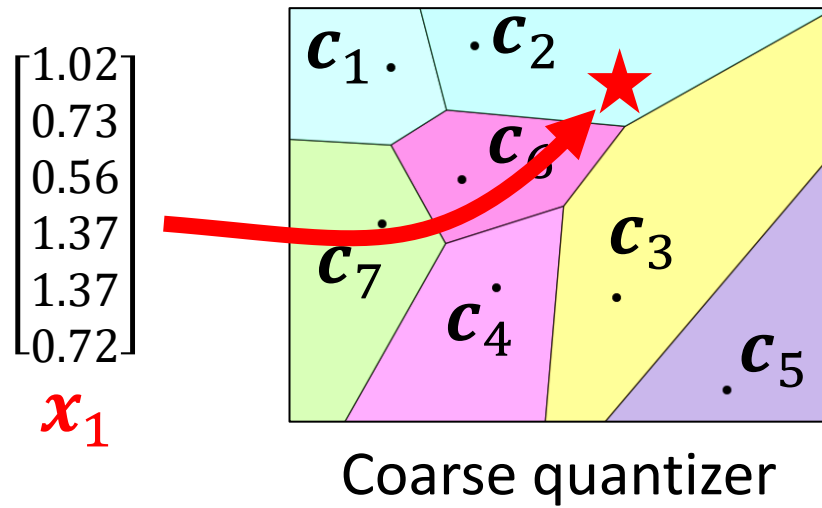
$$k = 2$$

\vdots

$$k = K$$

Inverted index + PQ: Record

Record x_1



$$k = 1$$

$$k = 2$$

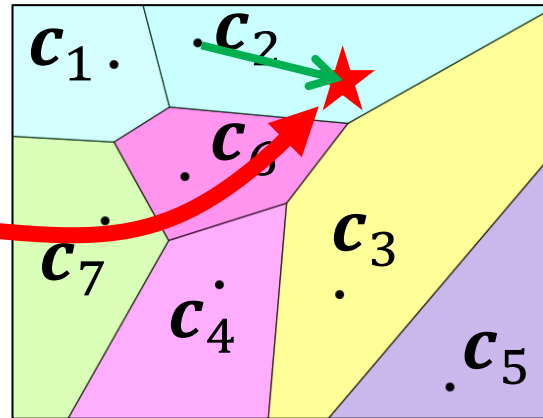
\vdots

$$k = K$$

Inverted index + PQ: Record

Record x_1

$\begin{bmatrix} 1.02 \\ 0.73 \\ 0.56 \\ 1.37 \\ 1.37 \\ 0.72 \end{bmatrix}$
 x_1



Coarse quantizer

$k = 1$

$k = 2$

\vdots

$k = K$

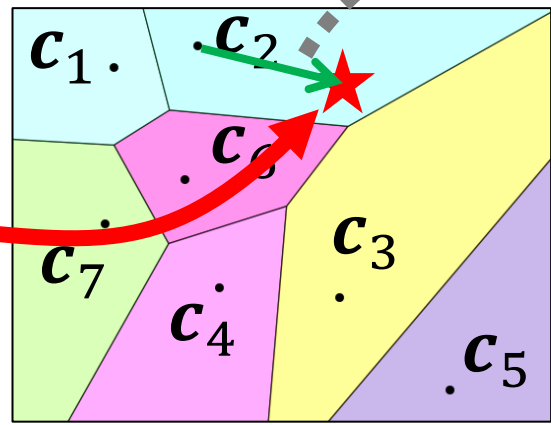
- c_2 is closest to x_1
- Compute a **residual** r_1 between x_1 and c_2 :

$$r_1 = x_1 - c_2 \quad (\rightarrow)$$

Inverted index + PQ: Record

Record x_1

$\begin{bmatrix} 1.02 \\ 0.73 \\ 0.56 \\ 1.37 \\ 1.37 \\ 0.72 \end{bmatrix}$
 x_1



Coarse quantizer

$k = 1$
 $k = 2$
 \vdots
 $k = K$

1
ID: 42
ID: 37
ID: 9

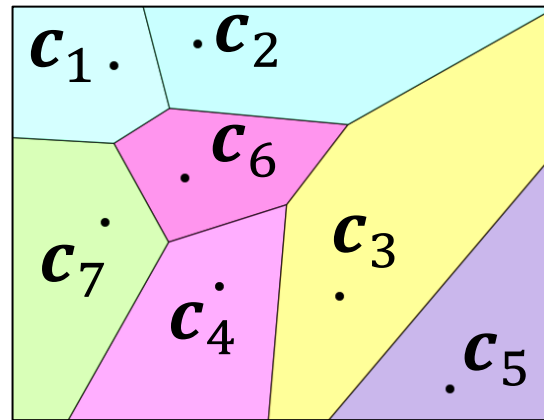
i
 \bar{r}_i

- c_2 is closest to x_1
- Compute a residual r_1 between x_1 and c_2 :
 $r_1 = x_1 - c_2$ (\rightarrow)

- Quantize r_1 to \bar{r}_1 by PQ
- Record it with the ID, "1"
- i.e., record (i, \bar{r}_i)

Inverted index + PQ: Record

➤ For all database vectors, record [ID + PQ(res)] as pointing lists



coarse quantizer

$k = 1$

245	12	1932	1721
ID: 42	ID: 25	ID: 38	ID: 16
ID: 37	ID: 47	ID: 49	ID: 72
ID: 9	ID: 32	ID: 72	ID: 95

$k = 2$

1	3721
ID: 42	ID: 18
ID: 37	ID: 4
ID: 9	ID: 96

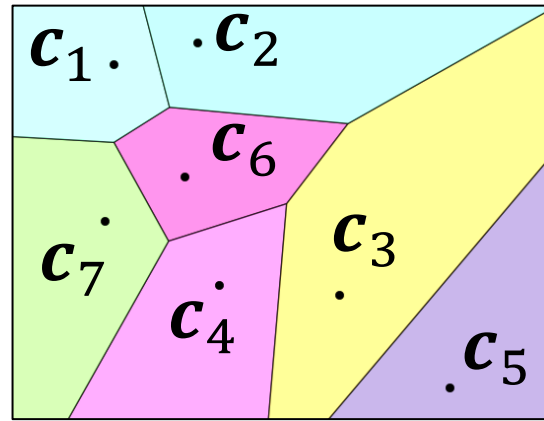
⋮

⋮

$k = K$

8621	145	324
ID: 24	ID: 77	ID: 32
ID: 54	ID: 21	ID: 11
ID: 23	ID: 5	ID: 85

Inverted index + PQ: Search



coarse quantizer

$k = 1$

245	12	1932	1721
ID: 42	ID: 25	ID: 38	ID: 16
ID: 37	ID: 47	ID: 49	ID: 72
ID: 9	ID: 32	ID: 72	ID: 95

$k = 2$

1	3721
ID: 42	ID: 18
ID: 37	ID: 4
ID: 9	ID: 96

⋮

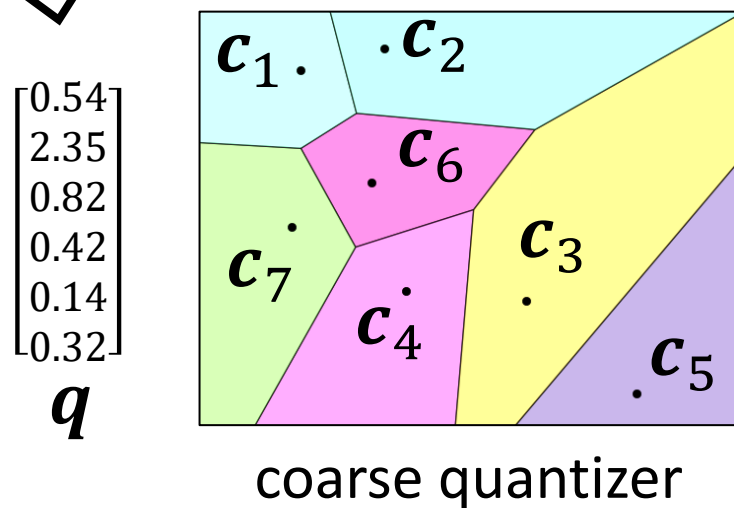
⋮

$k = K$

8621	145	324
ID: 24	ID: 77	ID: 32
ID: 54	ID: 21	ID: 11
ID: 23	ID: 5	ID: 85

Inverted index + PQ: Search

Find the nearest vector to q



$k = 1$

245	12	1932	1721
ID: 42	ID: 25	ID: 38	ID: 16
ID: 37	ID: 47	ID: 49	ID: 72
ID: 9	ID: 32	ID: 72	ID: 95

$k = 2$

1	3721
ID: 42	ID: 18
ID: 37	ID: 4
ID: 9	ID: 96

⋮

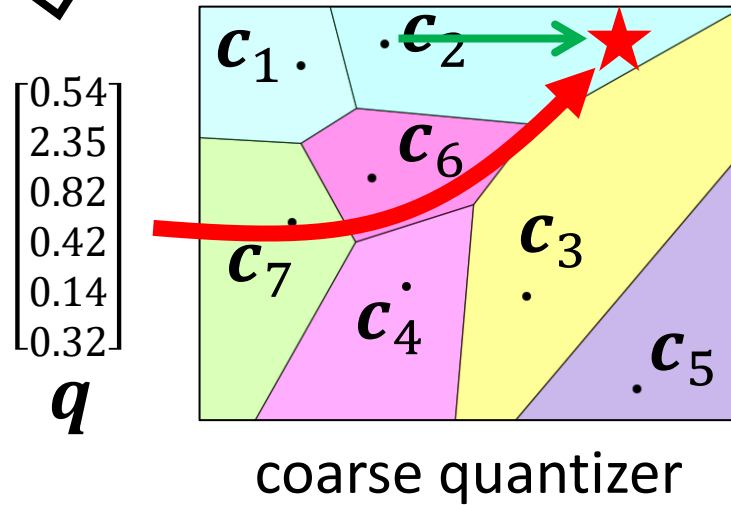
⋮

$k = K$

8621	145	324
ID: 24	ID: 77	ID: 32
ID: 54	ID: 21	ID: 11
ID: 23	ID: 5	ID: 85

Inverted index + PQ: Search

Find the nearest vector to q



$k = 1$

245	12	1932	1721
ID: 42	ID: 25	ID: 38	ID: 16
ID: 37	ID: 47	ID: 49	ID: 72
ID: 9	ID: 32	ID: 72	ID: 95

$k = 2$

1	3721
ID: 42	ID: 18
ID: 37	ID: 4
ID: 9	ID: 96

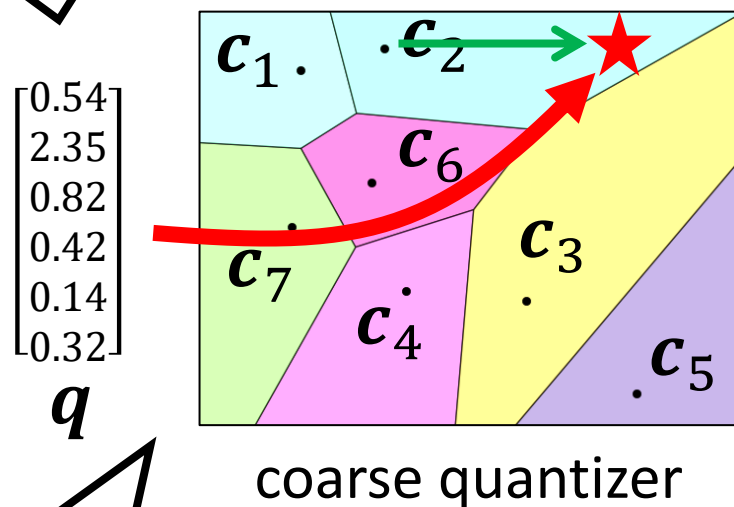
⋮

$k = K$

8621	145	324
ID: 24	ID: 77	ID: 32
ID: 54	ID: 21	ID: 11
ID: 23	ID: 5	ID: 85

Inverted index + PQ: Search

Find the nearest vector to q



- c_2 is the closest to q
- Compute the residual: $r_q = q - c_2$

$k = 1$

245	12	1932	1721
ID: 42	ID: 25	ID: 38	ID: 16
ID: 37	ID: 47	ID: 49	ID: 72
ID: 9	ID: 32	ID: 72	ID: 95

$k = 2$

1	3721
ID: 42	ID: 18
ID: 37	ID: 4
ID: 9	ID: 96

⋮

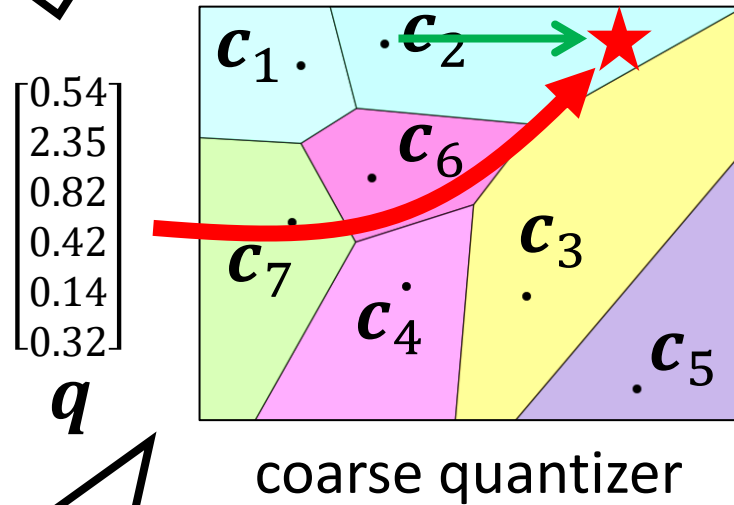
⋮

$k = K$

8621	145	324
ID: 24	ID: 77	ID: 32
ID: 54	ID: 21	ID: 11
ID: 23	ID: 5	ID: 85

Inverted index + PQ: Search

Find the nearest vector to q



- c_2 is the closest to q
- Compute the residual: $r_q = q - c_2$

$k = 1$

245	12	1932	1721
ID: 42	ID: 25	ID: 38	ID: 16
ID: 37	ID: 47	ID: 49	ID: 72
ID: 9	ID: 32	ID: 72	ID: 95

$k = 2$

1	3721
ID: 42	ID: 18
ID: 37	ID: 4
ID: 9	ID: 96

i

\bar{r}_i

- For all (i, \bar{r}_i) in $k = 2$, compare \bar{r}_i with r_q :

$$d(q, x_i)^2 = d(q - c_2, x_i - c_2)^2$$

$$= d(r_q, r_i)^2 \sim d_A(r_q, \bar{r}_i)^2$$
- Find the smallest one (several strategies)

Faiss

<https://github.com/facebookresearch/faiss>

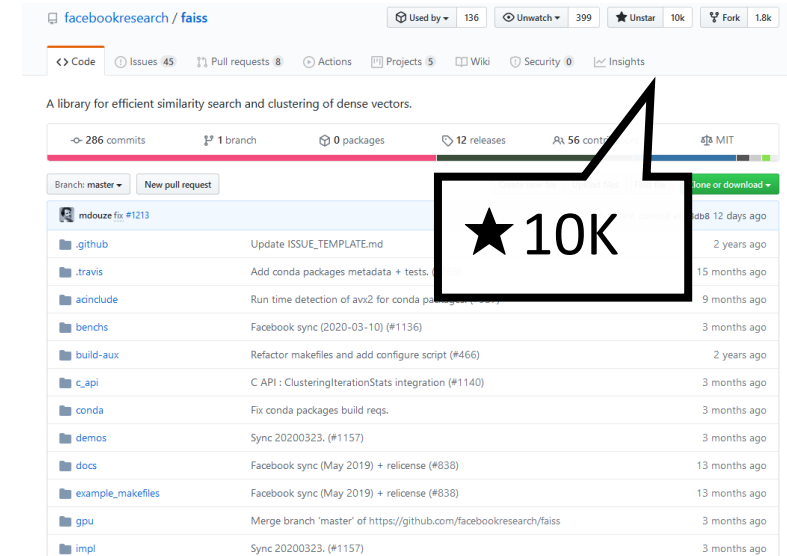
```
$> conda install faiss-cpu -c pytorch
```

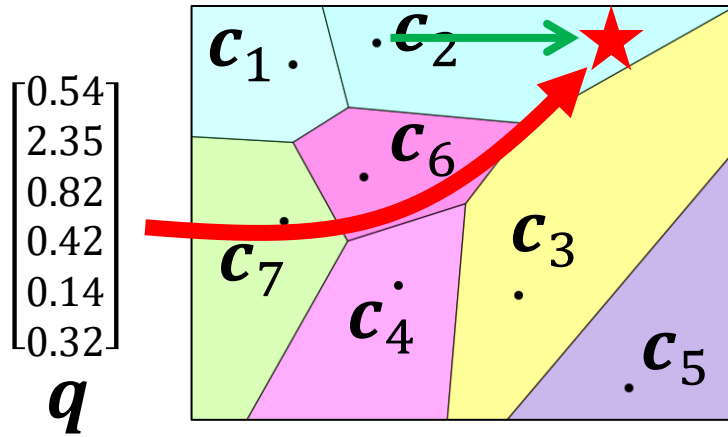
```
$> conda install faiss-gpu -c pytorch
```

- From the original authors of the PQ and a GPU expert, FAIR
- CPU-version: all PQ-based methods
- GPU-version: some PQ-based methods
- Bonus:
 - NN (not ANN) is also implemented, and quite fast
 - k-means (CPU/GPU). Fast.

Benchmark of k-means:

https://github.com/DwangoMediaVillage/pqkmeans/blob/master/tutorial/4_comparison_to_faiss.ipynb





coarse quantizer

$k = 1$

245	12	1932	1721
ID: 42	ID: 25	ID: 38	ID: 16
ID: 37	ID: 47	ID: 49	ID: 72
ID: 9	ID: 32	ID: 72	ID: 95

⋮

$k = K$

8621	145	324
ID: 24	ID: 77	ID: 32
ID: 54	ID: 21	ID: 11
ID: 23	ID: 5	ID: 85

M

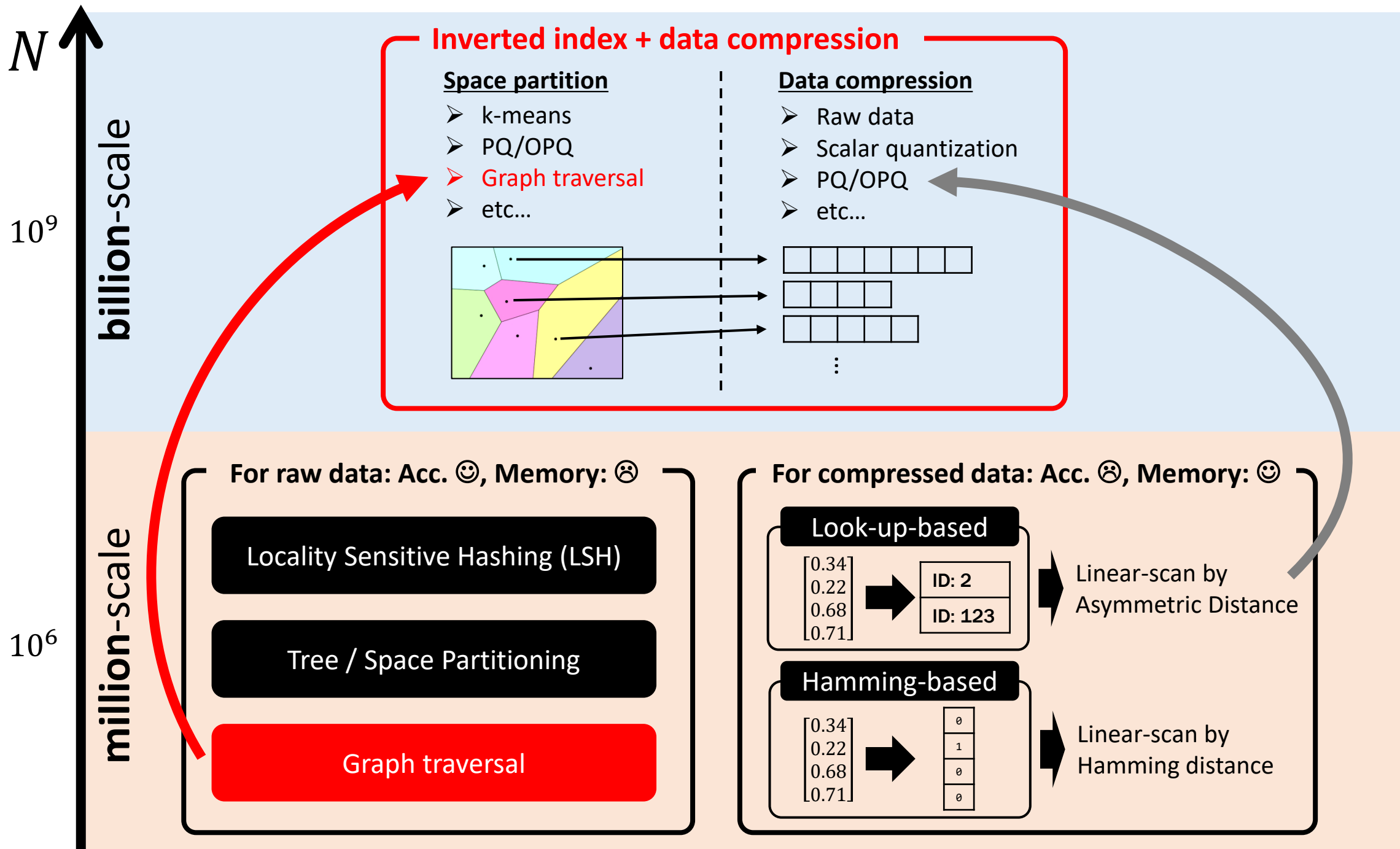
Simple linear scan

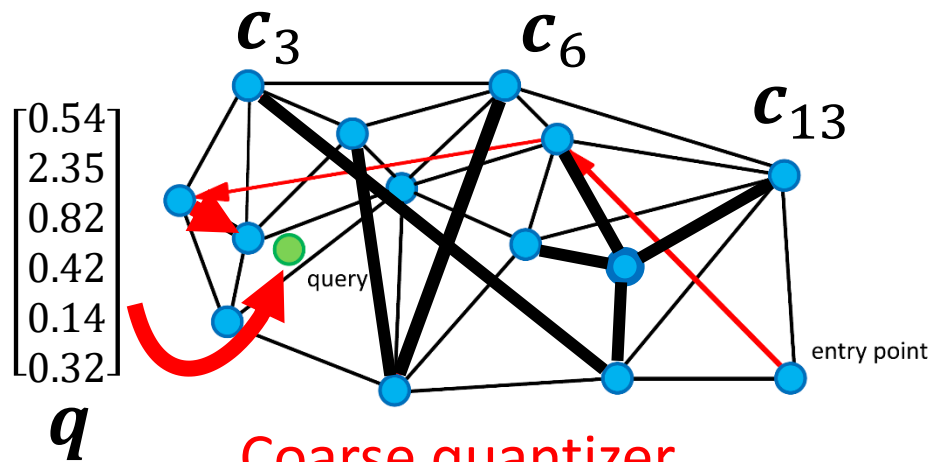
```
quantizer = faiss.IndexFlatL2(D)
index = faiss.IndexIVFPQ(quantizer, D, nlist, M, nbits)
```

Select a coarse quantizer

Usually, 8 bit

```
index.train(Xt) # Train
index.add(X) # Record data
index.nprobe = nprobe # Search parameter
dist, ids = index.search(Q, topk) # Search
```





$k = 1$

245	12	1932	1721
ID: 42	ID: 25	ID: 38	ID: 16
ID: 37	ID: 47	ID: 49	ID: 72
ID: 9	ID: 32	ID: 72	ID: 95

⋮

$k = K$

8621	145	324
ID: 24	ID: 77	ID: 32
ID: 54	ID: 21	ID: 11
ID: 23	ID: 5	ID: 85

M

```
quantizer = faiss.IndexHNSWFlat(D, hnsw_m)
index = faiss.IndexIVFPQ(quantizer, D, nlist, M, nbits)
```

Select a coarse quantizer

Usually, 8 bit

- Switch a coarse quantizer from linear-scan to HNSW
- The best approach for billion-scale data as of 2020
- The backbone of [Douze+, CVPR 2018] [Baranchuk+, ECCV 2018]

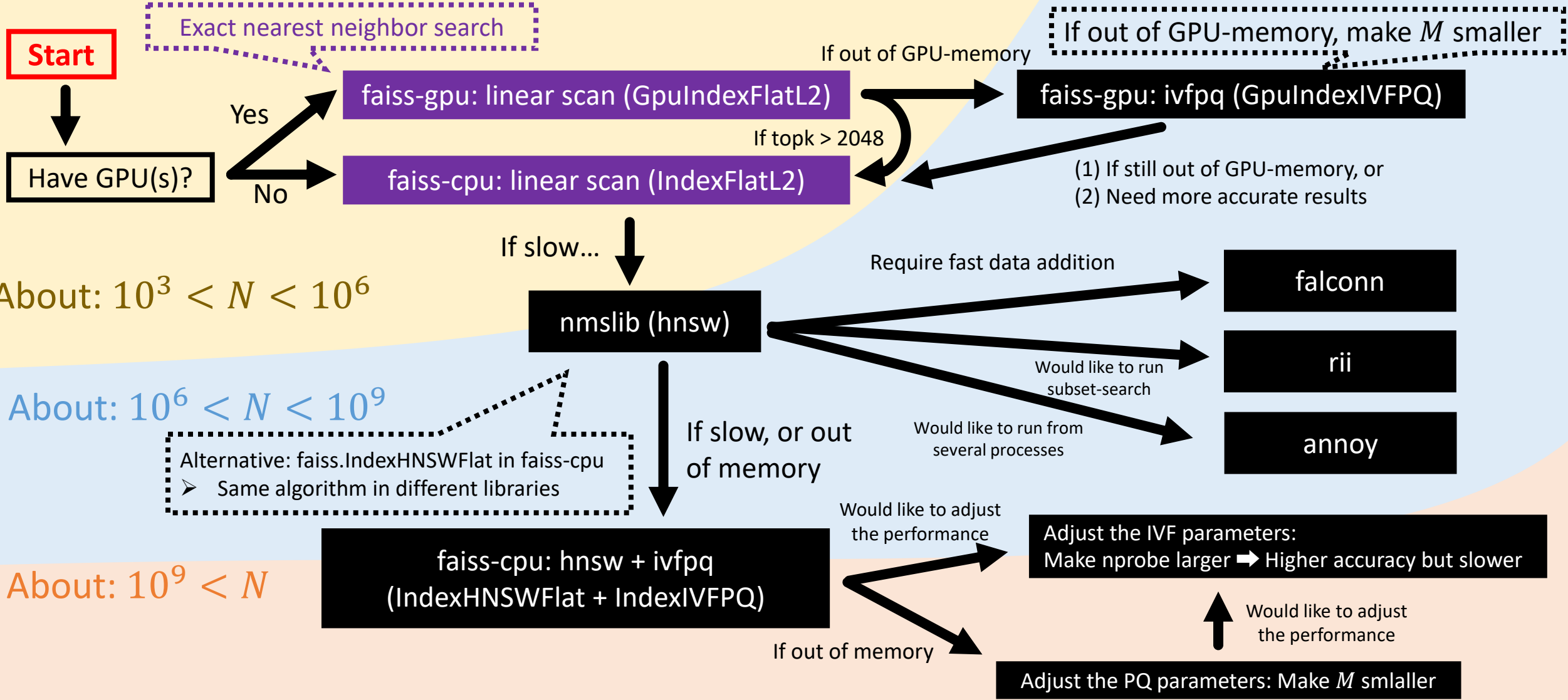
- 😊 From the original authors of PQ. Extremely efficient (theory & impl)
- 😊 Used in a real-world product (Mercari, etc)
- 😊 For billion-scale data, Faiss is the best option
- 😊 Especially, large-batch-search is fast; #query is large

- 😞 Lack of documentation (especially, python binding)
- 😞 Hard for a novice user to select a suitable algorithm
- 😞 As of 2020, anaconda is required. Pip is not supported officially

Reference

- Faiss wiki: [<https://github.com/facebookresearch/faiss/wiki>]
- Faiss tips: [https://github.com/matsui528/faiss_tips]
- Julia implementation of lookup-based methods [<https://github.com/una-dinosauria/Rayuela.jl>]
- PQ paper: H. Jégou et al., “Product quantization for nearest neighbor search,” TPAMI 2011
- IVFADC + HNSW (1): M. Douze et al., “Link and code: Fast indexing with graphs and compact regression codes,” CVPR 2018
- IVFADC + NHSW (2): D. Baranchuk et al., “Revisiting the Inverted Indices for Billion-Scale Approximate Nearest Neighbors,” ECCV 2018

cheat-sheet for ANN in Python (as of 2020. Can be installed by conda or pip)



About: $10^3 < N < 10^6$

About: $10^6 < N < 10^9$

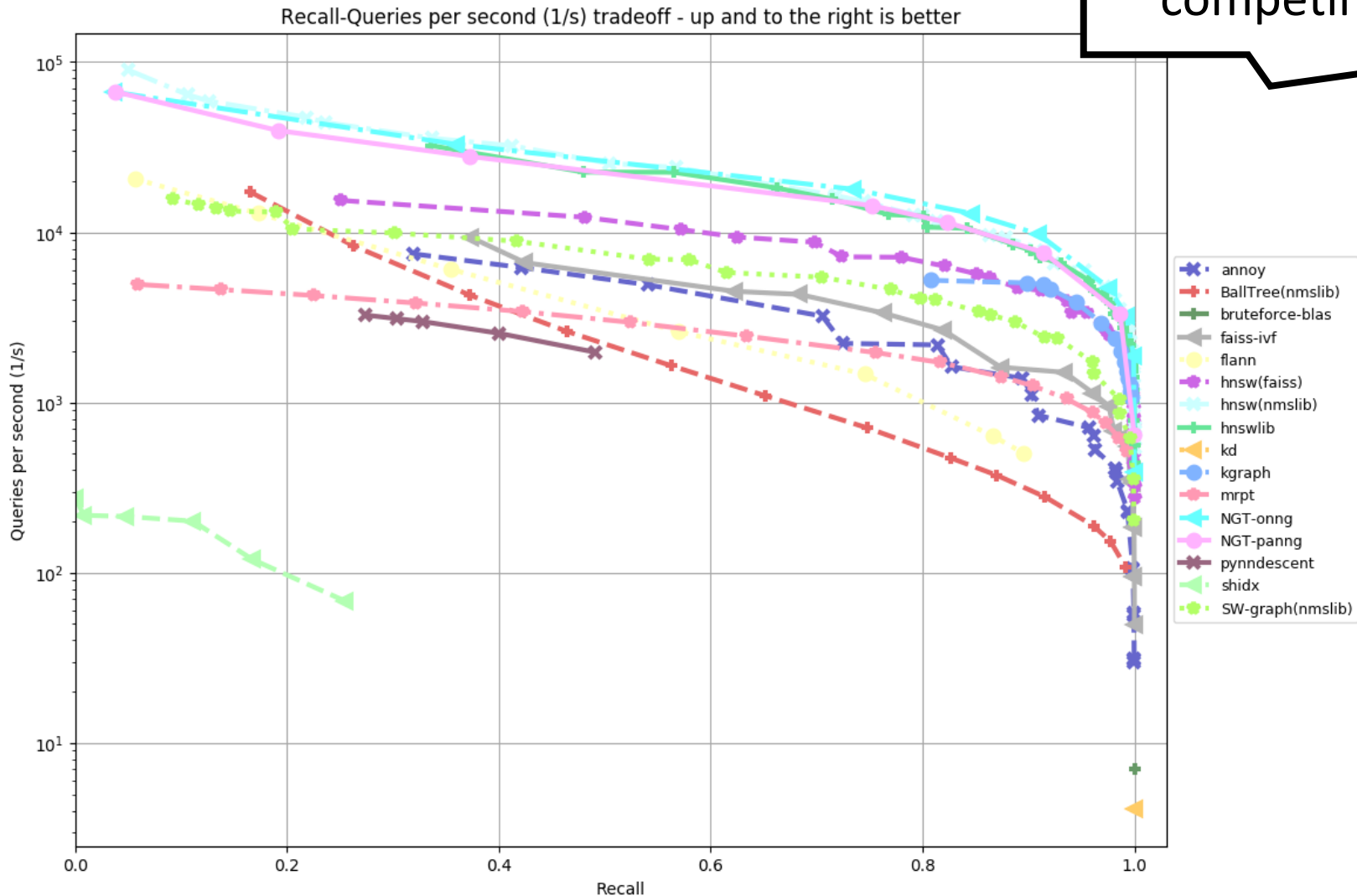
About: $10^9 < N$

Note: Assuming $D \cong 100$. The size of the problem is determined by DN . If $100 \ll D$, run PCA to reduce D to 100

Benchmark 1: ann-benchmarks

- <https://github.com/erikbern/ann-benchmarks>
- Comprehensive and thorough benchmarks for various libraries. Docker-based

- Top right, the better
- As of June, 2020, NMSLIB and NGT are competing each other for the first place



Benchmark 2: annbench

- <https://github.com/matsui528/annbench>
- Lightweight, easy-to-use

Install libraries

```
pip install -r requirements.txt
```

Download dataset on ./dataset

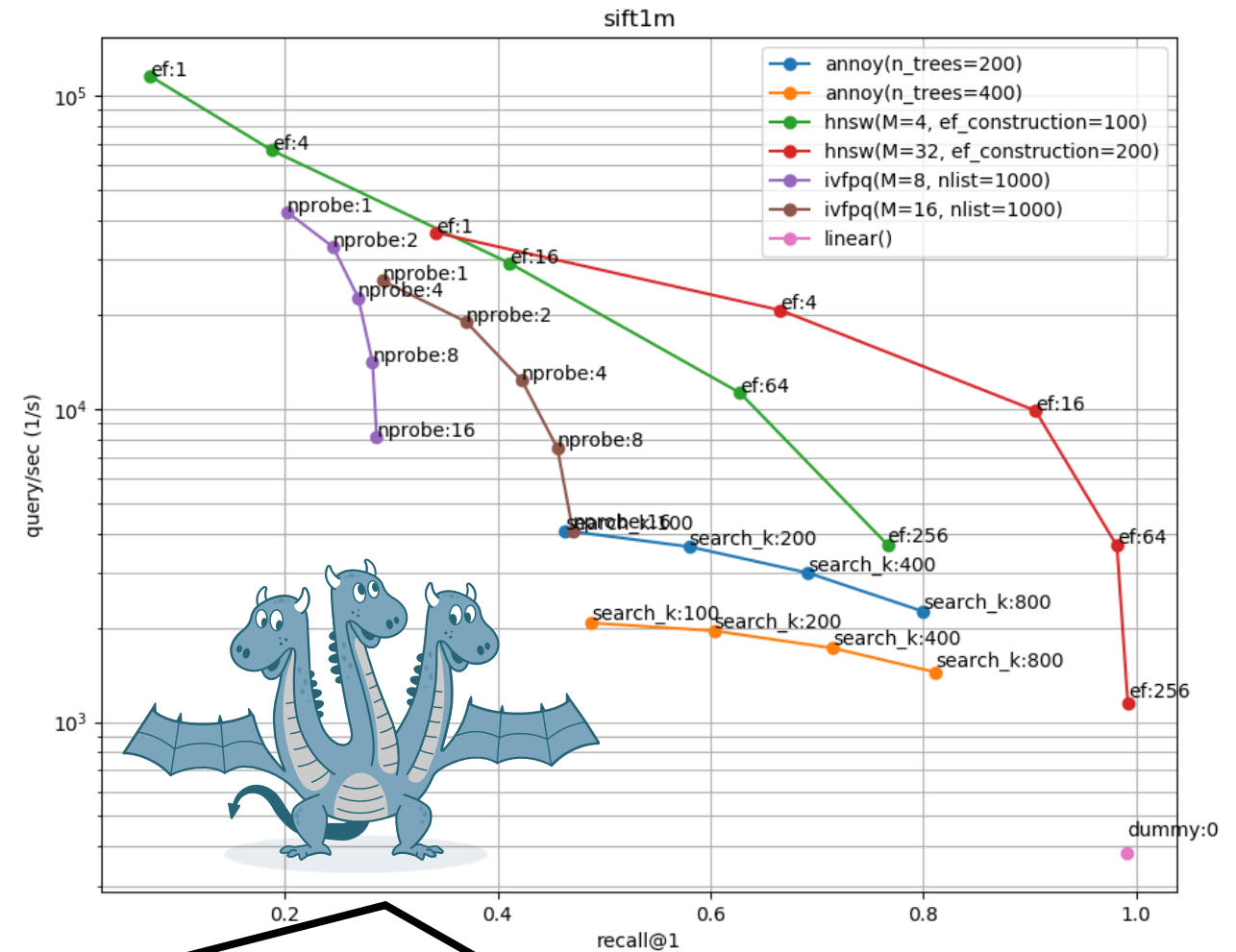
```
python download.py dataset=siftsmall
```

Evaluate algos. Results are on ./output

```
python run.py dataset=siftsmall algo=annoy
```

Visualize





```
python plot.py
```



Multi-run by Hydra

```
python run.py --multirun dataset=siftsmall,sift1m algo=linear,annoy,ivfpq,hnsw
```

Search for a “subset”

ID	Img	Tag
1		“cat”
2		“bird”
⋮		
125		“zebra”
126		“elephant”
⋮		

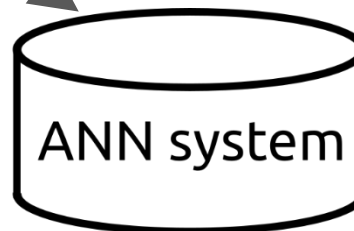
(1) Tag-based search:
tag == “zebra”

Target IDs:
[125, 223, 365, ...]

(2) Image search with a query q

$\{x_n\}_{n=1}^N$

q
Query vector
[125, 223, 365, ..., 881]
Target IDs



Ranked list

dist	ID
0.13	365
0.24	223
⋮	

Subset-search

Trillion-scale search: $N = 10^{12}$ (1T)

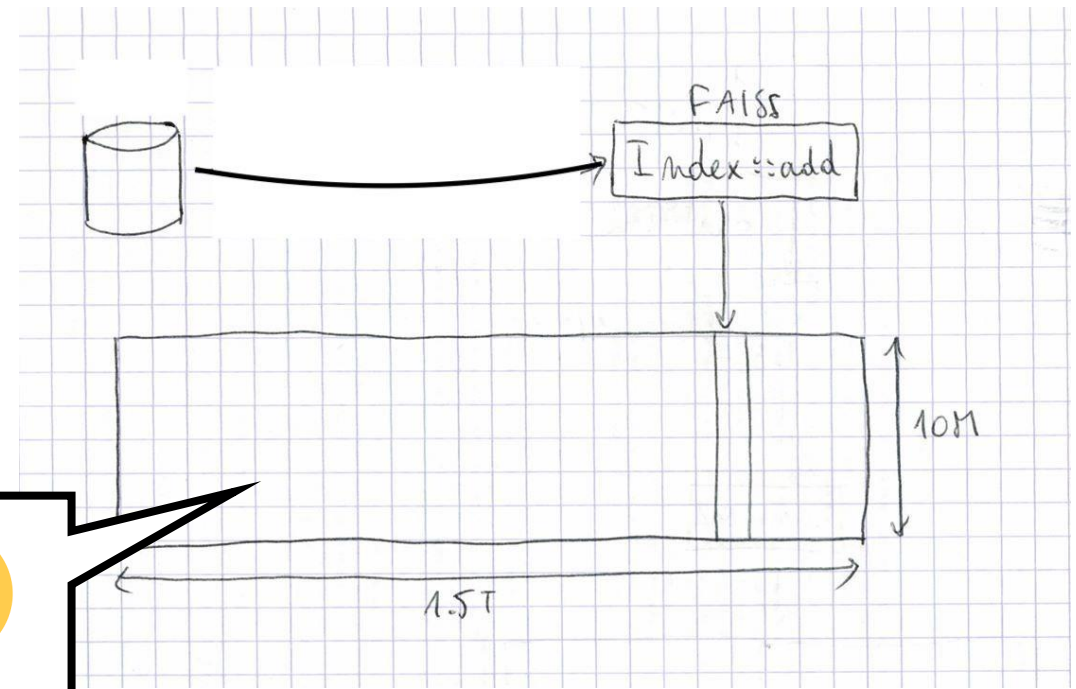
Sense of scale

- $K(= 10^3)$ Just in a second on a local machine
- $M(= 10^6)$ All data can be on memory. Try several approaches
- $G(= 10^9)$ Need to compress data by PQ. Only two datasets are available (SIFT1B, Deep1B)
- $T(= 10^{12})$ Cannot even imagine

<https://github.com/facebookresearch/faiss/wiki/Indexing-1T-vectors>

- Only in Faiss wiki
- Distributed, mmap, etc.

<https://github.com/facebookresearch/faiss/wiki/Indexing-1T-vectors>



A sparse matrix of 15 Exa elements?



Nearest neighbor search engine: something like ANN + SQL

- The algorithm inside is faiss, nmslib, or NGT



<https://github.com/vearch/vearch>



Elasticsearch KNN

<https://github.com/opendistro-for-elasticsearch/k-NN>



<https://github.com/milvus-io/milvus>

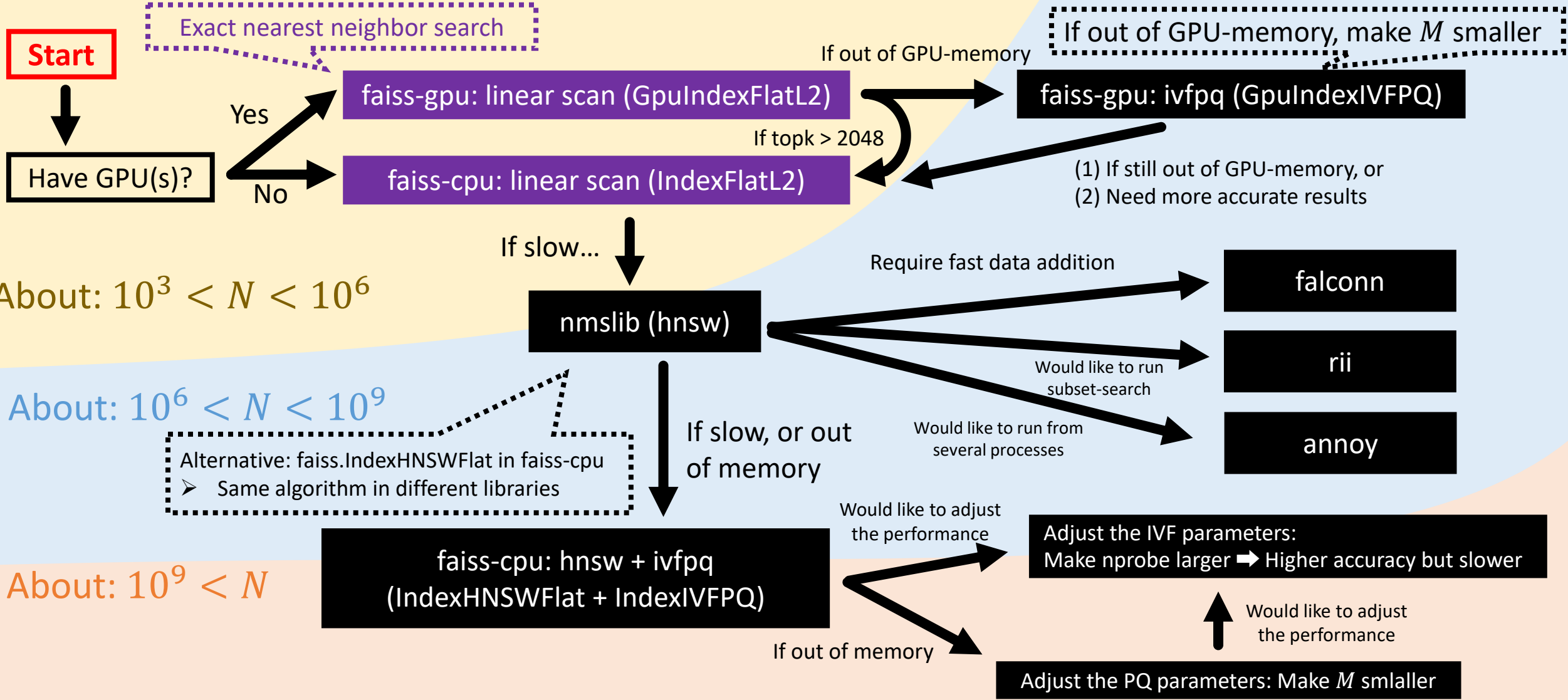


<https://github.com/vdaas/vald>

Problems of ANN

- No mathematical background.
 - ✓ Only actual measurements matter: recall and runtime
 - ✓ The ANN problem was mathematically defined 10+ years ago (LSH), but recently no one cares the definition.
- Thus, when the score is high, it's not clear the reason:
 - ✓ The method is good?
 - ✓ The implementation is good?
 - ✓ Just happens to work well for the target dataset?
 - ✓ E.g.: The difference of math library (OpenBLAS vs Intel MKL) matters.
- If one can explain “why this approach works good for this dataset”, it would be a great contribution to the field.
- Not enough dataset. Currently, only two datasets are available for billion-scale data: SIFT1B and Deep1B

cheat-sheet for ANN in Python (as of 2020. Can be installed by conda or pip)



Note: Assuming $D \cong 100$. The size of the problem is determined by DN . If $100 \ll D$, run PCA to reduce D to 100.