

Jack Spalding-Jamieson (Jack S-J)
jacksj@uwaterloo.ca

Independent

Scalable k -Means Clustering for Large k via Seeded Approximate Nearest-Neighbor Search

Joint work with Eliot Robson and Da Wei Zheng

k -Means/Sum of Squares Clustering: Quick Review

Input



- ❖ k : a parameter (e.g., $k = 3$)
- ❖ n d -dimensional vectors

k -Means/Sum of Squares Clustering: Quick Review

Input



- ❖ k : a parameter (e.g., $k = 3$)
- ❖ n d -dimensional vectors

Output



- ❖ Clusters $\{C_i\}_{i \in [k]}$ with centroids $\{\mu_i\}_{i \in [k]}$
- ❖ Objective:
$$\min \sum_{i=1}^k \sum_{x, y \in C_i} \|x - y\|^2$$

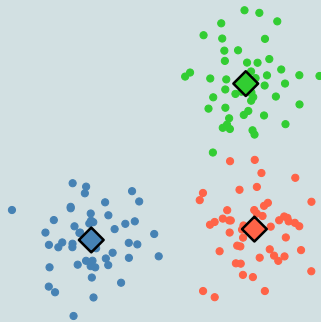
k -Means/Sum of Squares Clustering: Quick Review

Input



- ❖ k : a parameter (e.g., $k = 3$)
- ❖ n d -dimensional vectors

Output

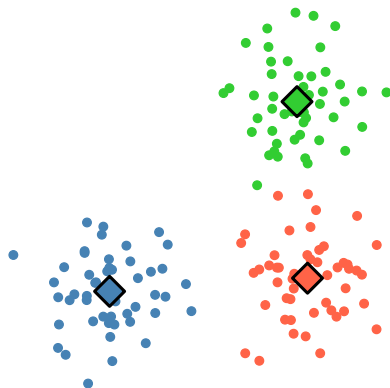


- ❖ Clusters $\{C_i\}_{i \in [k]}$ with centroids $\{\mu_i\}_{i \in [k]}$
- ❖ **Objective:**
$$\min \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} \|\mathbf{x} - \mu_i\|^2$$

(Equivalent by Huygens' theorem)

k -Means/Sum of Squares Clustering: Complexity and Approaches

k -means is NP-hard, even for $k = 2$.¹



¹Aloise et al., *NP-hardness of Euclidean sum-of-squares clustering*

²Arthur & Vassilvitskii, *k-means++: The Advantages of Careful Seeding*

³Bahmani et al., *Scalable k-means++*

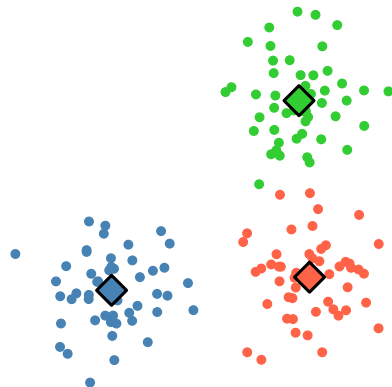
⁴Lloyd, *Least square quantization in PCM*

⁵Bahmani's slides on k -means | |

k -Means/Sum of Squares Clustering: Complexity and Approaches

k -means is NP-hard, even for $k = 2$.¹

Two broad approaches for good (practical) solutions:



¹Aloise et al., *NP-hardness of Euclidean sum-of-squares clustering*

²Arthur & Vassilvitskii, *k-means++: The Advantages of Careful Seeding*

³Bahmani et al., *Scalable k-means++*

⁴Lloyd, *Least square quantization in PCM*

⁵Bahmani's slides on k -means | |

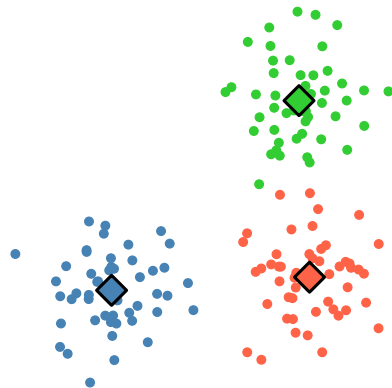
k -Means/Sum of Squares Clustering: Complexity and Approaches

k -means is NP-hard, even for $k = 2$.¹

Two broad approaches for good (practical) solutions:

- ▣ **Local search: Lloyd's algorithm**⁴

- ▣ Time Complexity per iteration: $O(nkd)$



¹Aloise et al., *NP-hardness of Euclidean sum-of-squares clustering*

²Arthur & Vassilvitskii, *k-means++: The Advantages of Careful Seeding*

³Bahmani et al., *Scalable k-means++*

⁴Lloyd, *Least square quantization in PCM*

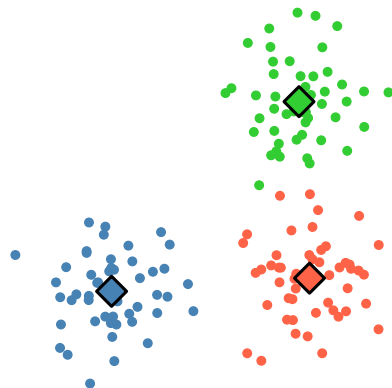
⁵Bahmani's slides on k -means |

k -Means/Sum of Squares Clustering: Complexity and Approaches

k -means is NP-hard, even for $k = 2$.¹

Two broad approaches for good (practical) solutions:

- ❖ **Local search:** Lloyd's algorithm⁴
 - ❖ Time Complexity per iteration: $O(nkd)$
- ❖ **Approximation algorithms:**



¹Aloise et al., *NP-hardness of Euclidean sum-of-squares clustering*

²Arthur & Vassilvitskii, *k-means++: The Advantages of Careful Seeding*

³Bahmani et al., *Scalable k-means++*

⁴Lloyd, *Least square quantization in PCM*

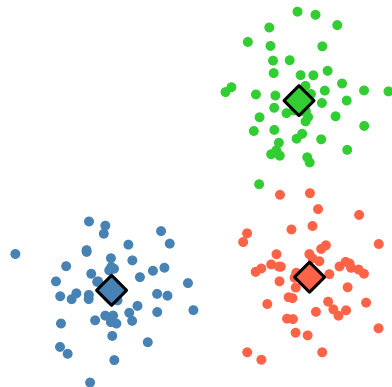
⁵Bahmani's slides on k -means |

k -Means/Sum of Squares Clustering: Complexity and Approaches

k -means is NP-hard, even for $k = 2$.¹

Two broad approaches for good (practical) solutions:

- ❖ **Local search:** Lloyd's algorithm⁴
 - ❖ Time Complexity per iteration: $O(nkd)$
- ❖ **Approximation algorithms:**
 - ❖ **k-means++**²
 - ▶ $O(nkd)$ time
 - ▶ $8(\ln k + 2)$ expected approximation



¹Aloise et al., *NP-hardness of Euclidean sum-of-squares clustering*

²Arthur & Vassilvitskii, *k-means++: The Advantages of Careful Seeding*

³Bahmani et al., *Scalable k-means++*

⁴Lloyd, *Least square quantization in PCM*

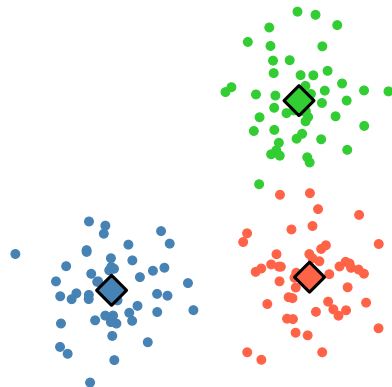
⁵Bahmani's slides on k -means |

k -Means/Sum of Squares Clustering: Complexity and Approaches

k -means is NP-hard, even for $k = 2$.¹

Two broad approaches for good (practical) solutions:

- ❖ **Local search:** Lloyd's algorithm⁴
 - ❖ Time Complexity per iteration: $O(nkd)$
- ❖ **Approximation algorithms:**
 - ❖ k -means++²
 - ▶ $O(nkd)$ time
 - ▶ $8(\ln k + 2)$ expected approximation
 - ❖ k -means | ^{3,5}
 - ▶ $O(c_1 \cdot nd + O(c_2 \cdot k^2 d))$ time, c_1, c_2 small in practice
 - ▶ $O(\log k)$ expected approximation



¹Aloise et al., *NP-hardness of Euclidean sum-of-squares clustering*

²Arthur & Vassilvitskii, *k-means++: The Advantages of Careful Seeding*

³Bahmani et al., *Scalable k-means++*

⁴Lloyd, *Least square quantization in PCM*

⁵Bahmani's slides on k -means |

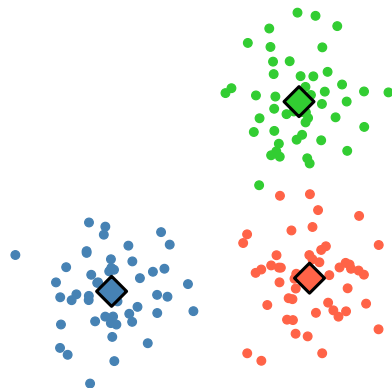
k -Means/Sum of Squares Clustering: Complexity and Approaches

k -means is NP-hard, even for $k = 2$.¹

Two broad approaches for good (practical) solutions:

- ❖ **Local search:** Lloyd's algorithm⁴
 - ❖ Time Complexity per iteration: $O(nkd)$
- ❖ **Approximation algorithms:**
 - ❖ k -means++²
 - ▶ $O(nkd)$ time
 - ▶ $8(\ln k + 2)$ expected approximation
 - ❖ k -means | ^{3,5}
 - ▶ $O(c_1 \cdot nd + O(c_2 \cdot k^2 d))$ time, c_1, c_2 small in practice
 - ▶ $O(\log k)$ expected approximation

Better approximations known, not used in practice.



¹Aloise et al., *NP-hardness of Euclidean sum-of-squares clustering*

²Arthur & Vassilvitskii, *k-means++: The Advantages of Careful Seeding*

³Bahmani et al., *Scalable k-means++*

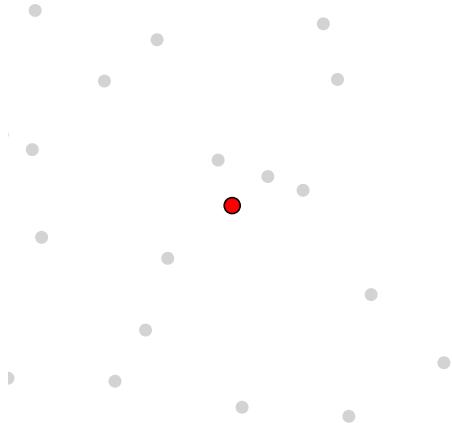
⁴Lloyd, *Least square quantization in PCM*

⁵Bahmani's slides on k -means |

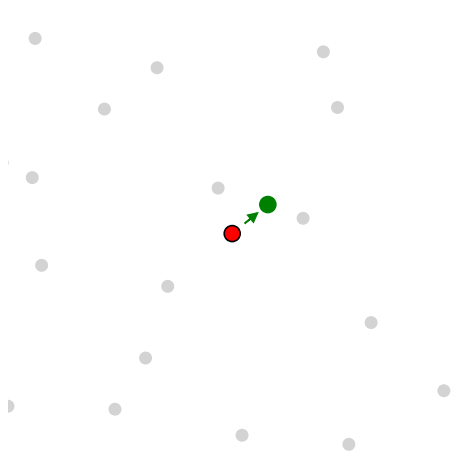
Nearest-Neighbor Search



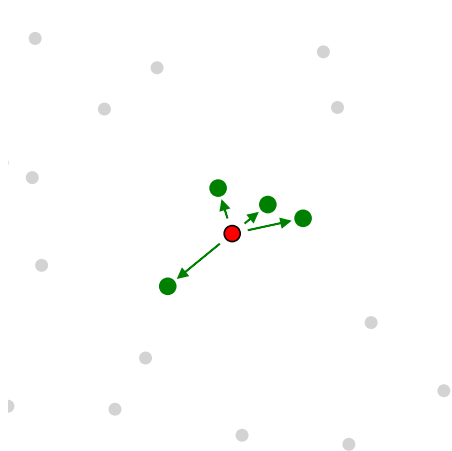
Nearest-Neighbor Search



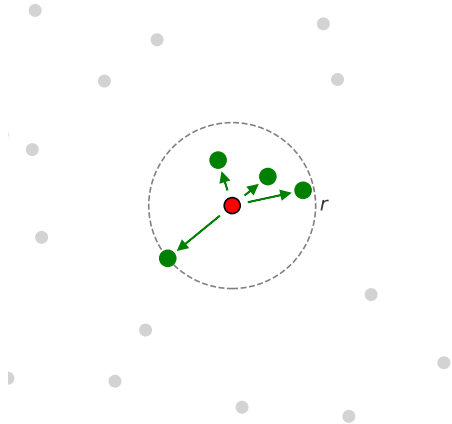
Nearest-Neighbor Search



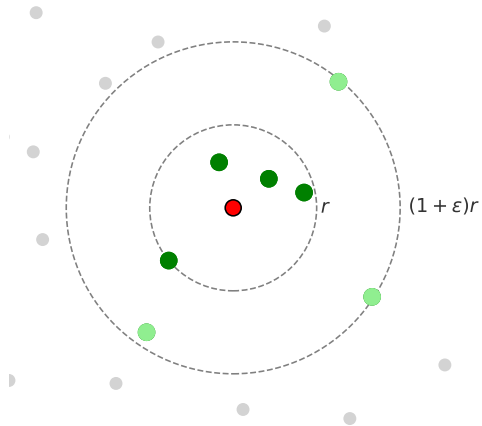
k -Nearest-Neighbor Search



k -Nearest-Neighbor Search

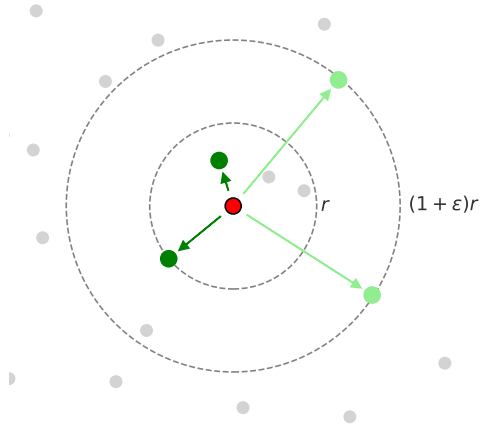


Approximate k -Nearest-Neighbor Search



Approximation factor: $(1 + \epsilon)$

Approximate k -Nearest-Neighbor Search



Approximation factor: $(1 + \varepsilon)$

Recall: $\frac{2}{4} = 0.5$

Dataset sizes vs Existing Approaches

ANNS has three broad families of practical approaches:

- ❑ Quantization
- ❑ Space-partitioning/Clustering
- ❑ Search-graphs

Quantization is normally used alongside the other two.

Dataset sizes vs Existing Approaches

ANNS has three broad families of practical approaches:

- ❑ Quantization
- ❑ Space-partitioning/Clustering
- ❑ Search-graphs

Quantization is normally used alongside the other two.

Important fact: Search-graph approaches are (almost universally) best,

Dataset sizes vs Existing Approaches

ANNS has three broad families of practical approaches:

- ❑ Quantization
- ❑ Space-partitioning/Clustering
- ❑ Search-graphs

Quantization is normally used alongside the other two.

Important fact: Search-graph approaches are (almost universally) best, but **only work if your data or quantized data fits in RAM** (prohibitively slow otherwise).

Dataset sizes vs Existing Approaches

ANNS has three broad families of practical approaches:

- ❑ Quantization
- ❑ Space-partitioning/Clustering
- ❑ Search-graphs

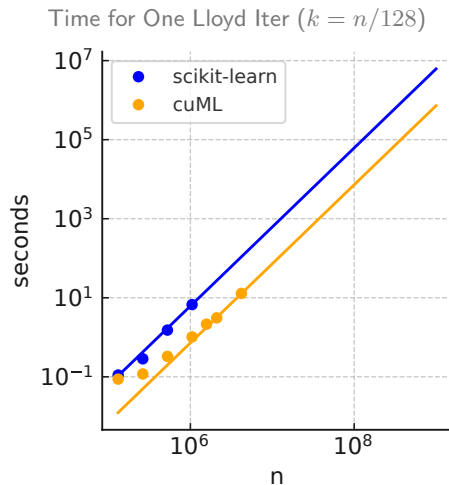
Quantization is normally used alongside the other two.

Important fact: Search-graph approaches are (almost universally) best, but **only work if your data or quantized data fits in RAM** (prohibitively slow otherwise).

e.g. 200-dimensional f32 dataset with 20,000,000 points requires 16GB (plus data structure size).

Some Motivation

“Vector Similarity Search” is very popular in machine learning recently, with a lot of active development.



¹Sivic & Zisserman, *Video Google: A text retrieval approach to object matching in videos*

²Jégou et al., *Product Quantization for Nearest Neighbor Search*

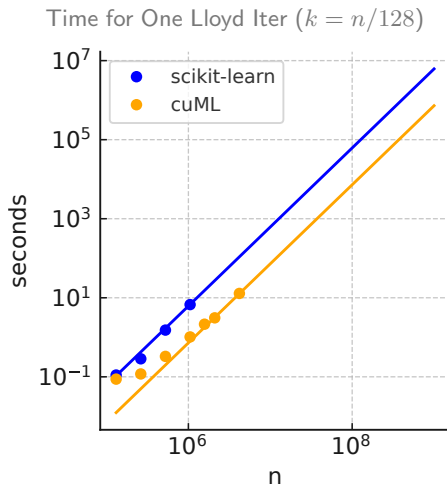
³Chen et al., *SPANN: Highly-Efficient Billion-scale Approximate Nearest Neighbor Search*

⁴Sun et al., *SOAR: Improved Indexing for Approximate Nearest Neighbor Search*

Some Motivation

“Vector Similarity Search” is very popular in machine learning recently, with a lot of active development.

- ❖ **Many methods used in industry for massive datasets^{1,2,3,4} solve k -means as a sub-routine, and would benefit from very large k .**



¹Sivic & Zisserman, *Video Google: A text retrieval approach to object matching in videos*

²Jégou et al., *Product Quantization for Nearest Neighbor Search*

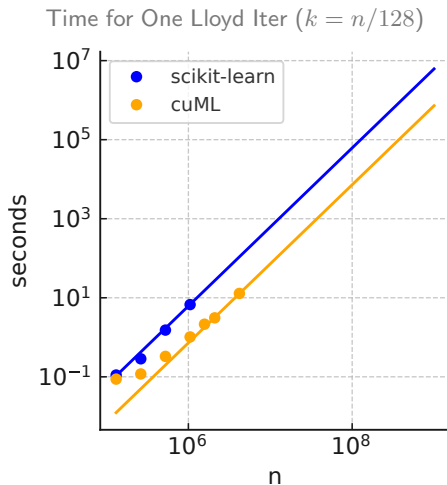
³Chen et al., *SPANN: Highly-Efficient Billion-scale Approximate Nearest Neighbor Search*

⁴Sun et al., *SOAR: Improved Indexing for Approximate Nearest Neighbor Search*

Some Motivation

“Vector Similarity Search” is very popular in machine learning recently, with a lot of active development.

- ❖ Many methods used in industry for massive datasets^{1,2,3,4} solve k -means as a sub-routine, and would benefit from very large k .
- ❖ **Many would benefit from large k ($k \approx n/c$ for small c), but current methods are too slow.**



¹Sivic & Zisserman, *Video Google: A text retrieval approach to object matching in videos*

²Jégou et al., *Product Quantization for Nearest Neighbor Search*

³Chen et al., *SPANN: Highly-Efficient Billion-scale Approximate Nearest Neighbor Search*

⁴Sun et al., *SOAR: Improved Indexing for Approximate Nearest Neighbor Search*

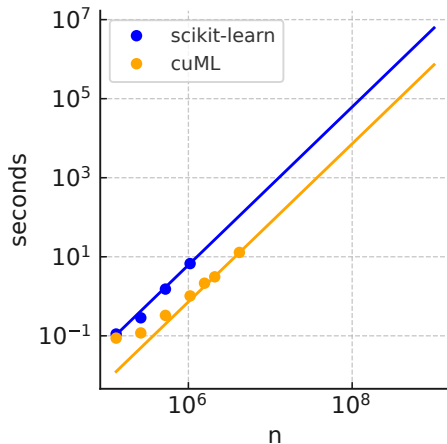
Some Motivation

“Vector Similarity Search” is very popular in machine learning recently, with a lot of active development.

- ❖ Many methods used in industry for massive datasets^{1,2,3,4} solve k -means as a sub-routine, and would benefit from very large k .
- ❖ Many would benefit from large k ($k \approx n/c$ for small c), but current methods are too slow.

All practical methods take $\Omega(k^2)$ time.

Time for One Lloyd Iter ($k = n/128$)



¹Sivic & Zisserman, *Video Google: A text retrieval approach to object matching in videos*

²Jégou et al., *Product Quantization for Nearest Neighbor Search*

³Chen et al., *SPANN: Highly-Efficient Billion-scale Approximate Nearest Neighbor Search*

⁴Sun et al., *SOAR: Improved Indexing for Approximate Nearest Neighbor Search*

Some Motivation

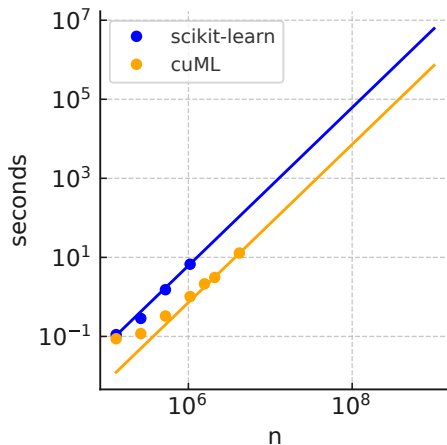
“Vector Similarity Search” is very popular in machine learning recently, with a lot of active development.

- ❖ Many methods used in industry for massive datasets^{1,2,3,4} solve k -means as a sub-routine, and would benefit from very large k .
- ❖ Many would benefit from large k ($k \approx n/c$ for small c), but current methods are too slow.

All practical methods take $\Omega(k^2)$ time.

Focus on large $n \in [10^6, 10^9]$, $k \approx n/c$, and $d \geq 100$.

Time for One Lloyd Iter ($k = n/128$)



¹Sivic & Zisserman, *Video Google: A text retrieval approach to object matching in videos*

²Jégou et al., *Product Quantization for Nearest Neighbor Search*

³Chen et al., *SPANN: Highly-Efficient Billion-scale Approximate Nearest Neighbor Search*

⁴Sun et al., *SOAR: Improved Indexing for Approximate Nearest Neighbor Search*

Some Motivation

“Vector Similarity Search” is very popular in machine learning recently, with a lot of active development.

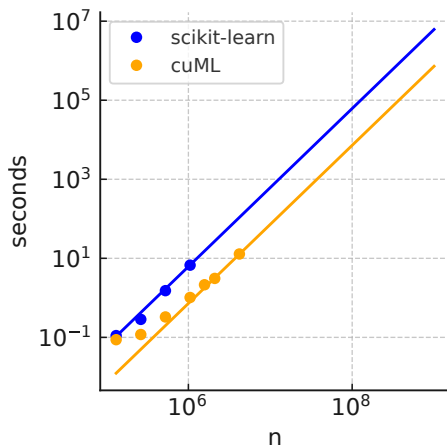
- ❖ Many methods used in industry for massive datasets^{1,2,3,4} solve k -means as a sub-routine, and would benefit from very large k .
- ❖ Many would benefit from large k ($k \approx n/c$ for small c), but current methods are too slow.

All practical methods take $\Omega(k^2)$ time.

Focus on large $n \in [10^6, 10^9]$, $k \approx n/c$, and $d \geq 100$.

- ❖ **Days to weeks with current methods!** →

Time for One Lloyd Iter ($k = n/128$)



¹Sivic & Zisserman, *Video Google: A text retrieval approach to object matching in videos*

²Jégou et al., *Product Quantization for Nearest Neighbor Search*

³Chen et al., *SPANN: Highly-Efficient Billion-scale Approximate Nearest Neighbor Search*

⁴Sun et al., *SOAR: Improved Indexing for Approximate Nearest Neighbor Search*

Some Motivation

“Vector Similarity Search” is very popular in machine learning recently, with a lot of active development.

- ❖ Many methods used in industry for massive datasets^{1,2,3,4} solve k -means as a sub-routine, and would benefit from very large k .
- ❖ Many would benefit from large k ($k \approx n/c$ for small c), but current methods are too slow.

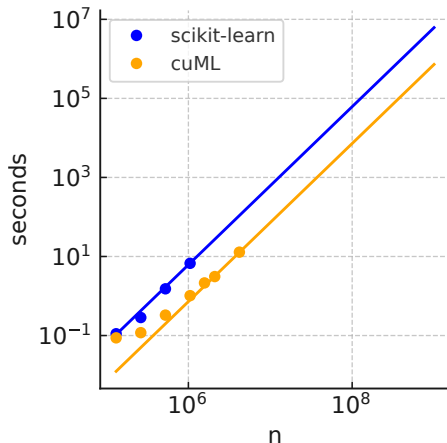
All practical methods take $\Omega(k^2)$ time.

Focus on large $n \in [10^6, 10^9]$, $k \approx n/c$, and $d \geq 100$.

- ❖ Days to weeks with current methods! →

Focus on real performance, not asymptotics.

Time for One Lloyd Iter ($k = n/128$)



¹Sivic & Zisserman, *Video Google: A text retrieval approach to object matching in videos*

²Jégou et al., *Product Quantization for Nearest Neighbor Search*

³Chen et al., *SPANN: Highly-Efficient Billion-scale Approximate Nearest Neighbor Search*

⁴Sun et al., *SOAR: Improved Indexing for Approximate Nearest Neighbor Search*

Typical approach:

Typical approach:

- **Init centroids with one of:**

- `k-means++`
- `k-means ||`
- **uniform random sampling**

Typical approach:

- ❖ Init centroids with one of:
 - ❖ k-means++
 - ❖ k-means ||
 - ❖ uniform random sampling
- ❖ **Lloyd's algorithm for local search**

Typical approach:

- ❖ Init centroids with one of:
 - ❖ k-means++
 - ❖ k-means ||
 - ❖ uniform random sampling
- ❖ Lloyd's algorithm for local search
- ❖ **Iterate until timeout**

Typical approach:

- ❖ Init centroids with one of:
 - ❖ k-means++
 - ❖ k-means ||
 - ❖ uniform random sampling
- ❖ Lloyd's algorithm for local search
- ❖ Iterate until timeout
- ❖ **Accelerate each step on GPU**

Bottleneck Testing

Typical approach:

- ❖ Init centroids with one of:
 - ❖ k-means++
 - ❖ k-means ||
 - ❖ uniform random sampling
- ❖ Lloyd's algorithm for local search
- ❖ Iterate until timeout
- ❖ Accelerate each step on GPU

Small k : k-means++ and || are very good

Bottleneck Testing

Typical approach:

- ❖ Init centroids with one of:
 - ❖ `k-means++`
 - ❖ `k-means ||`
 - ❖ uniform random sampling
- ❖ Lloyd's algorithm for local search
- ❖ Iterate until timeout
- ❖ Accelerate each step on GPU

Small k : `k-means++` and `||` are very good

We have large n and k , e.g. $n = 5e6$ and $k = 1e4$.

Bottleneck Testing

Typical approach:

- ❖ Init centroids with one of:
 - ❖ k-means++
 - ❖ k-means||
 - ❖ uniform random sampling
- ❖ Lloyd's algorithm for local search
- ❖ Iterate until timeout
- ❖ Accelerate each step on GPU

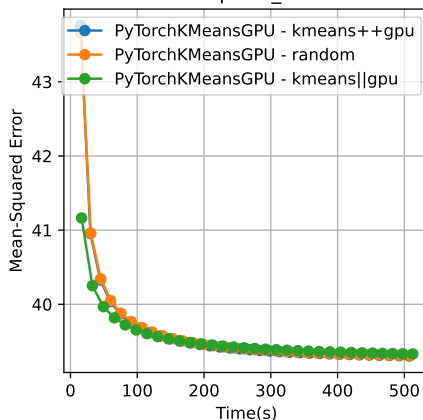
Small k : k-means++ and || are very good

We have large n and k , e.g. $n = 5e6$ and $k = 1e4$.

Large n/k : Only Lloyd's algorithm matters (right)

Conclusion: Want to accelerate Lloyd's algorithm

Score over time in dpr5m_base with $k=10000$



Out-of-Core Similarity Search

(used by your favourite semantic search engine)

Basically out-of-core ANNS with n base points



k -Means Clustering



Exact NNS

k base points

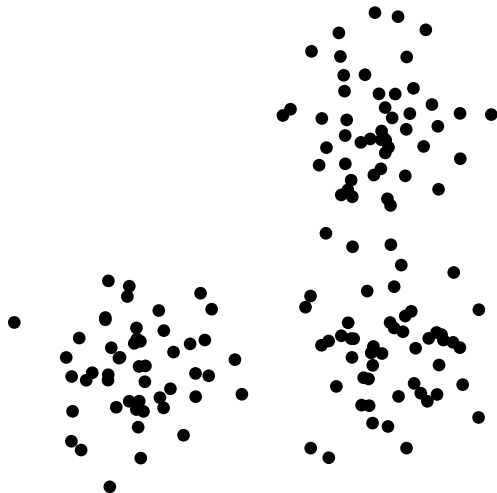


GPU Acceleration

Lloyd's k -Means Method

1. **Initialization:** Sample k centroids uniformly from the dataset.
2. Iterate (local search):
 - ❖ **Assignment:** Assign each point to the nearest centroid.
 - ❖ **Mean Computation:** Update each centroid to be the average of points assigned to it.

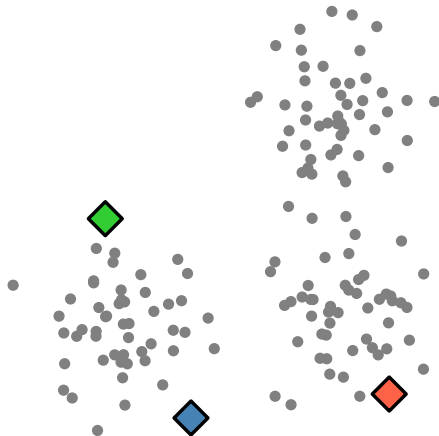
Bottleneck is **Assignment** step.



1. **Initialization: Sample k centroids uniformly from the dataset.**
2. Iterate (local search):
 - **Assignment:** Assign each point to the nearest centroid.
 - **Mean Computation:** Update each centroid to be the average of points assigned to it.

Bottleneck is **Assignment** step.

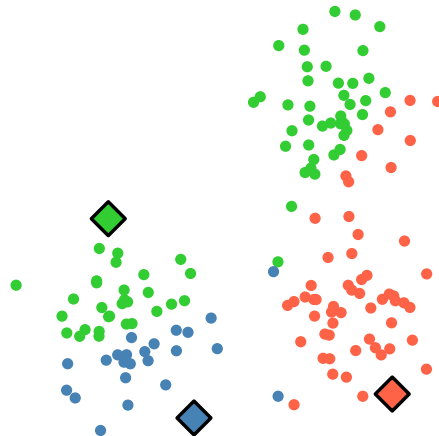
Initialize Centroids



1. **Initialization:** Sample k centroids uniformly from the dataset.
2. Iterate (local search):
 - **Assignment:** Assign each point to the nearest centroid.
 - **Mean Computation:** Update each centroid to be the average of points assigned to it.

Bottleneck is **Assignment** step.

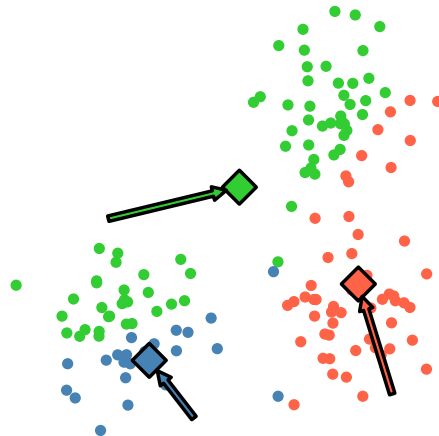
Iter 1a: Assign Labels



Iter 1b: Compute Means

1. **Initialization:** Sample k centroids uniformly from the dataset.
2. Iterate (local search):
 - **Assignment:** Assign each point to the nearest centroid.
 - **Mean Computation:** Update each centroid to be the average of points assigned to it.

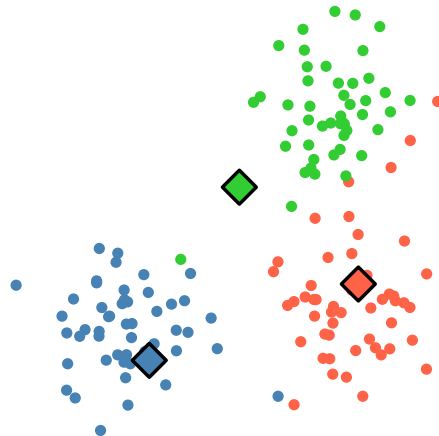
Bottleneck is **Assignment** step.



1. **Initialization:** Sample k centroids uniformly from the dataset.
2. Iterate (local search):
 - **Assignment:** Assign each point to the nearest centroid.
 - **Mean Computation:** Update each centroid to be the average of points assigned to it.

Bottleneck is **Assignment** step.

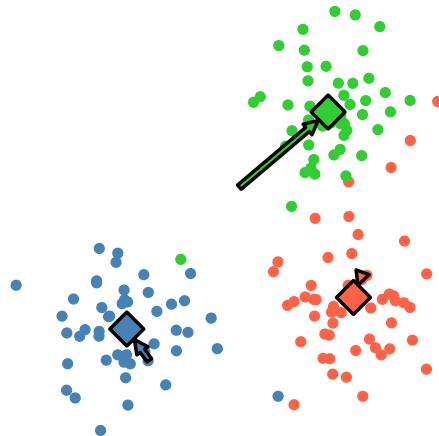
Iter 2a: Assign Labels



Iter 2b: Compute Means

1. **Initialization:** Sample k centroids uniformly from the dataset.
2. Iterate (local search):
 - **Assignment:** Assign each point to the nearest centroid.
 - **Mean Computation:** Update each centroid to be the average of points assigned to it.

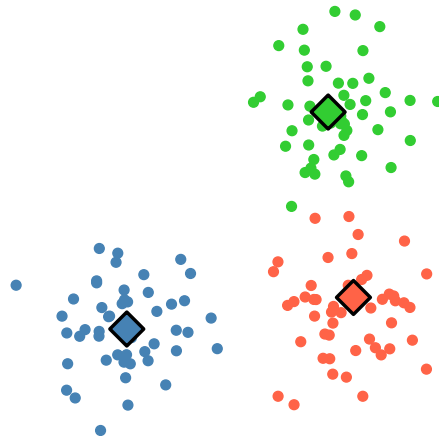
Bottleneck is **Assignment** step.



1. **Initialization:** Sample k centroids uniformly from the dataset.
2. Iterate (local search):
 - **Assignment:** Assign each point to the nearest centroid.
 - **Mean Computation:** Update each centroid to be the average of points assigned to it.

Bottleneck is **Assignment** step.

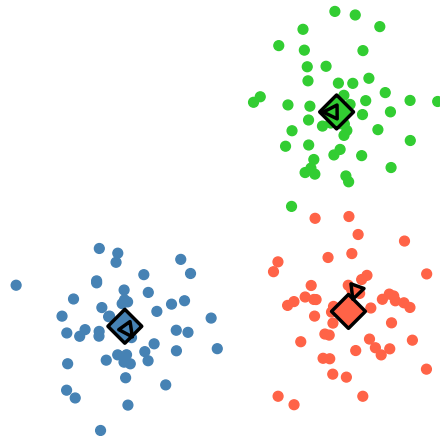
Iter 3a: Assign Labels



Iter 3b: Compute Means

1. **Initialization:** Sample k centroids uniformly from the dataset.
2. Iterate (local search):
 - **Assignment:** Assign each point to the nearest centroid.
 - **Mean Computation:** Update each centroid to be the average of points assigned to it.

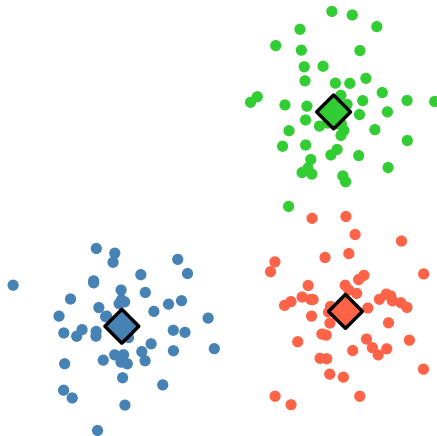
Bottleneck is **Assignment** step.



1. **Initialization:** Sample k centroids uniformly from the dataset.
2. Iterate (local search):
 - **Assignment:** Assign each point to the nearest centroid.
 - **Mean Computation:** Update each centroid to be the average of points assigned to it.

Bottleneck is **Assignment** step.

Iter 4a: Assign Labels

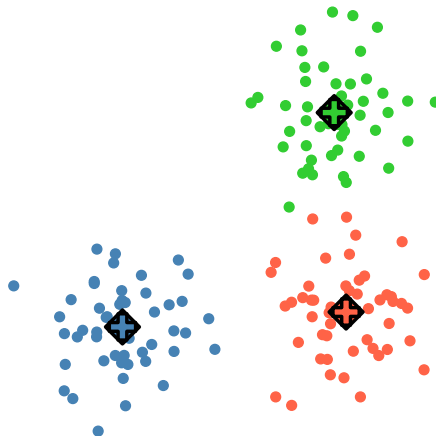


Iter 4b: Compute Means

1. **Initialization:** Sample k centroids uniformly from the dataset.
2. Iterate (local search):
 - **Assignment:** Assign each point to the nearest centroid.
 - **Mean Computation:** Update each centroid to be the average of points assigned to it.

Bottleneck is **Assignment** step.

Key observation: Assignment step is a nearest-neighbour problem.



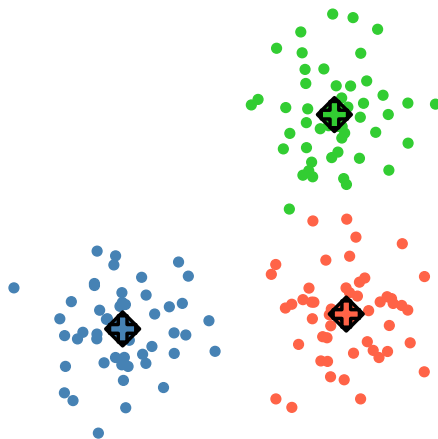
Iter 4b: Compute Means

1. **Initialization:** Sample k centroids uniformly from the dataset.
2. Iterate (local search):
 - **Assignment:** Assign each point to the nearest centroid.
 - **Mean Computation:** Update each centroid to be the average of points assigned to it.

Bottleneck is **Assignment** step.

Key observation: Assignment step is a nearest-neighbour problem.

**Lloyd's algorithm is very limited in theory.
But it's very good in practice.**



Lloyd Iterations with a Black-box ANNS Data Structure

Alternative iteration approach:

- ❖ **Build:** Construct an **approximate nearest neighbor search (ANNS) data structure** on the centroids.
- ❖ **Assignment:** **Use the data structure** to assign each data point to its nearest centroid approximately.
- ❖ **Mean Computation:** Unchanged

¹Borodin et al., *Lower Bounds for High Dimensional Nearest Neighbor Search*

²Liu, *A strong lower bound for approximate nearest neighbor searching*

³Malkov & Yashunin, *Efficient and Robust [ANNS] Using Hierarchical Navigable Small World Graphs*

⁴Raschka et al., *Machine Learning in Python: Main developments and technology trends [...]*

Lloyd Iterations with a Black-box ANNS Data Structure

Alternative iteration approach:

- ❖ **Build:** Construct an approximate nearest neighbor search (ANNS) data structure on the centroids.
- ❖ **Assignment:** Use the data structure to assign each data point to its nearest centroid approximately.
 - ❖ **Make sure not to regress (compare assignments)**
- ❖ **Mean Computation:** Unchanged

¹Borodin et al., *Lower Bounds for High Dimensional Nearest Neighbor Search*

²Liu, *A strong lower bound for approximate nearest neighbor searching*

³Malkov & Yashunin, *Efficient and Robust [ANNS] Using Hierarchical Navigable Small World Graphs*

⁴Raschka et al., *Machine Learning in Python: Main developments and technology trends [...]*

Lloyd Iterations with a Black-box ANNS Data Structure

Alternative iteration approach:

- ❖ **Build:** Construct an approximate nearest neighbor search (ANNS) data structure on the centroids.
- ❖ **Assignment:** Use the data structure to assign each data point to its nearest centroid approximately.
 - ❖ Make sure not to regress (compare assignments)
- ❖ **Mean Computation:** Unchanged

Why ANNS? Exact-NN has a linear lower bound in high-dimensions¹.

ANNS has strong lower bounds too², but good widely-used heuristics exist.

¹Borodin et al., *Lower Bounds for High Dimensional Nearest Neighbor Search*

²Liu, *A strong lower bound for approximate nearest neighbor searching*

³Malkov & Yashunin, *Efficient and Robust [ANNS] Using Hierarchical Navigable Small World Graphs*

⁴Raschka et al., *Machine Learning in Python: Main developments and technology trends [...]*

Lloyd Iterations with a Black-box ANNS Data Structure

Alternative iteration approach:

- ❖ **Build:** Construct an approximate nearest neighbor search (ANNS) data structure on the centroids.
- ❖ **Assignment:** Use the data structure to assign each data point to its nearest centroid approximately.
 - ❖ Make sure not to regress (compare assignments)
- ❖ **Mean Computation:** Unchanged

Experiment:

- ❖ Baseline: Popular Lloyd implementations:
 - ❖ CPU: scikit
 - ❖ GPU: cuML, simple pytorch impl
- ❖ Suite of (CPU) ANNS data structures:
 - ❖ PQ, SQ, IVF IVFPQ, HNSW
 - ❖ FAISS implementations

¹Borodin et al., *Lower Bounds for High Dimensional Nearest Neighbor Search*

²Liu, *A strong lower bound for approximate nearest neighbor searching*

³Malkov & Yashunin, *Efficient and Robust [ANNS] Using Hierarchical Navigable Small World Graphs*

⁴Raschka et al., *Machine Learning in Python: Main developments and technology trends [...]*

Lloyd Iterations with a Black-box ANNS Data Structure

Alternative iteration approach:

- ❖ **Build:** Construct an approximate nearest neighbor search (ANNS) data structure on the centroids.
- ❖ **Assignment:** Use the data structure to assign each data point to its nearest centroid approximately.
 - ❖ Make sure not to regress (compare assignments)
- ❖ **Mean Computation:** Unchanged

Experiment:

- ❖ **Baseline: Popular Lloyd implementations:**
 - ❖ CPU: scikit
 - ❖ GPU: cuML, simple pytorch impl
- ❖ Suite of (CPU) ANNS data structures:
 - ❖ PQ, SQ, IVF IVFPQ, HNSW
 - ❖ FAISS implementations

¹Borodin et al., *Lower Bounds for High Dimensional Nearest Neighbor Search*

²Liu, *A strong lower bound for approximate nearest neighbor searching*

³Malkov & Yashunin, *Efficient and Robust [ANNS] Using Hierarchical Navigable Small World Graphs*

⁴Raschka et al., *Machine Learning in Python: Main developments and technology trends [...]*

Lloyd Iterations with a Black-box ANNS Data Structure

Alternative iteration approach:

- ❖ **Build:** Construct an approximate nearest neighbor search (ANNS) data structure on the centroids.
- ❖ **Assignment:** Use the data structure to assign each data point to its nearest centroid approximately.
 - ❖ Make sure not to regress (compare assignments)
- ❖ **Mean Computation:** Unchanged

Experiment:

- ❖ **Baseline:** Popular Lloyd implementations:
 - ❖ CPU: scikit
 - ❖ GPU: cuML, simple pytorch impl
- ❖ **Suite of (CPU) ANNS data structures:**
 - ❖ PQ, SQ, IVF IVFPQ, HNSW
 - ❖ FAISS implementations

¹Borodin et al., *Lower Bounds for High Dimensional Nearest Neighbor Search*

²Liu, *A strong lower bound for approximate nearest neighbor searching*

³Malkov & Yashunin, *Efficient and Robust [ANNS] Using Hierarchical Navigable Small World Graphs*

⁴Raschka et al., *Machine Learning in Python: Main developments and technology trends [...]*

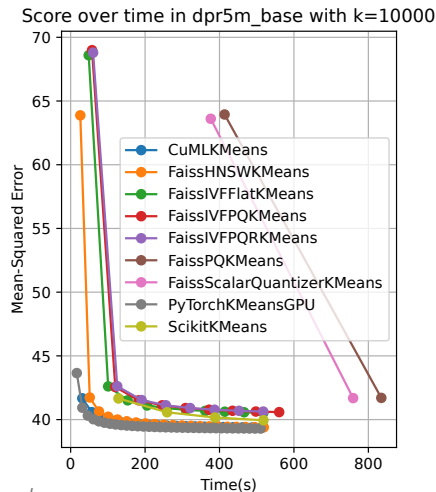
Lloyd Iterations with a Black-box ANNS Data Structure

Alternative iteration approach:

- ❖ **Build:** Construct an approximate nearest neighbor search (ANNS) data structure on the centroids.
- ❖ **Assignment:** Use the data structure to assign each data point to its nearest centroid approximately.
 - ❖ Make sure not to regress (compare assignments)
- ❖ **Mean Computation:** Unchanged

Experiment:

- ❖ **Baseline:** Popular Lloyd implementations:
 - ❖ CPU: scikit
 - ❖ GPU: cuML, simple pytorch impl
- ❖ **Suite of (CPU) ANNS data structures:**
 - ❖ PQ, SQ, IVF IVFPQ, HNSW
 - ❖ FAISS implementations



¹Borodin et al., *Lower Bounds for High Dimensional Nearest Neighbor Search*

²Liu, *A strong lower bound for approximate nearest neighbor searching*

³Malkov & Yashunin, *Efficient and Robust [ANNS] Using Hierarchical Navigable Small World Graphs*

⁴Raschka et al., *Machine Learning in Python: Main developments and technology trends [...]*

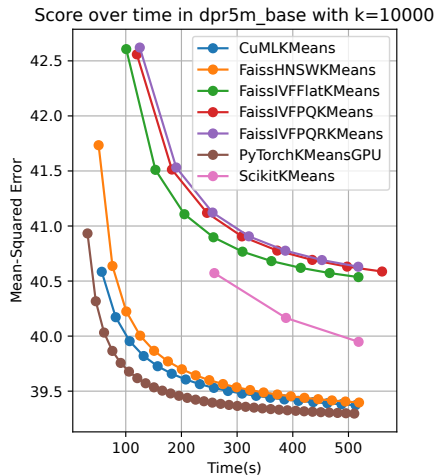
Lloyd Iterations with a Black-box ANNS Data Structure

Alternative iteration approach:

- ❖ **Build:** Construct an approximate nearest neighbor search (ANNS) data structure on the centroids.
- ❖ **Assignment:** Use the data structure to assign each data point to its nearest centroid approximately.
 - ❖ Make sure not to regress (compare assignments)
- ❖ **Mean Computation:** Unchanged

Experiment:

- ❖ **Baseline:** Popular Lloyd implementations:
 - ❖ CPU: scikit
 - ❖ GPU: cuML, simple pytorch impl
- ❖ **Suite of (CPU) ANNS data structures:**
 - ❖ PQ, SQ, IVF IVFPQ, HNSW
 - ❖ FAISS implementations



¹Borodin et al., *Lower Bounds for High Dimensional Nearest Neighbor Search*

²Liu, *A strong lower bound for approximate nearest neighbor searching*

³Malkov & Yashunin, *Efficient and Robust [ANNS] Using Hierarchical Navigable Small World Graphs*

⁴Raschka et al., *Machine Learning in Python: Main developments and technology trends [...]*

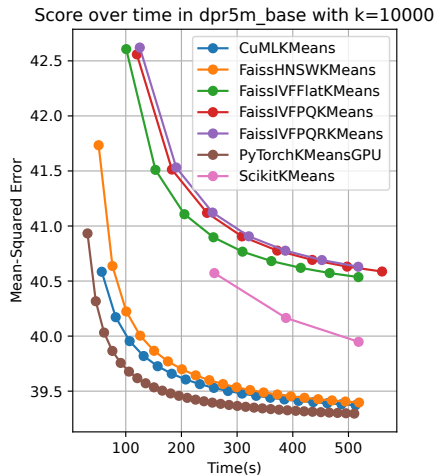
Lloyd Iterations with a Black-box ANNS Data Structure

Alternative iteration approach:

- ❖ **Build:** Construct an approximate nearest neighbor search (ANNS) data structure on the centroids.
- ❖ **Assignment:** Use the data structure to assign each data point to its nearest centroid approximately.
 - ❖ Make sure not to regress (compare assignments)
- ❖ **Mean Computation:** Unchanged

Conclusions:

- ❖ Dimension reduction (quantization) generally bad



¹Borodin et al., *Lower Bounds for High Dimensional Nearest Neighbor Search*

²Liu, *A strong lower bound for approximate nearest neighbor searching*

³Malkov & Yashunin, *Efficient and Robust [ANNS] Using Hierarchical Navigable Small World Graphs*

⁴Raschka et al., *Machine Learning in Python: Main developments and technology trends [...]*

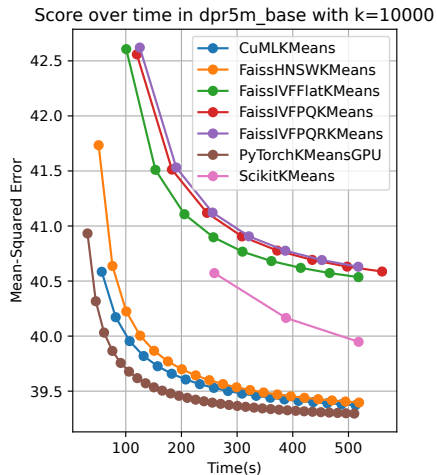
Lloyd Iterations with a Black-box ANNS Data Structure

Alternative iteration approach:

- ❖ **Build:** Construct an approximate nearest neighbor search (ANNS) data structure on the centroids.
- ❖ **Assignment:** Use the data structure to assign each data point to its nearest centroid approximately.
 - ❖ Make sure not to regress (compare assignments)
- ❖ **Mean Computation:** Unchanged

Conclusions:

- ❖ Dimension reduction (quantization) generally bad
- ❖ **HNSW³ is really good!**



¹Borodin et al., *Lower Bounds for High Dimensional Nearest Neighbor Search*

²Liu, *A strong lower bound for approximate nearest neighbor searching*

³Malkov & Yashunin, *Efficient and Robust [ANNS] Using Hierarchical Navigable Small World Graphs*

⁴Raschka et al., *Machine Learning in Python: Main developments and technology trends [...]*

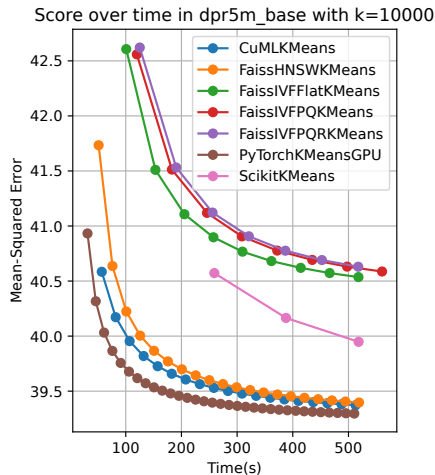
Lloyd Iterations with a Black-box ANNS Data Structure

Alternative iteration approach:

- ❖ **Build:** Construct an approximate nearest neighbor search (ANNS) data structure on the centroids.
- ❖ **Assignment:** Use the data structure to assign each data point to its nearest centroid approximately.
 - ❖ Make sure not to regress (compare assignments)
- ❖ **Mean Computation:** Unchanged

Conclusions:

- ❖ Dimension reduction (quantization) generally bad
- ❖ HNSW³ is *really* good!
 - ❖ Almost as good as Nvidia's own GPU implementation.⁴



¹Borodin et al., *Lower Bounds for High Dimensional Nearest Neighbor Search*

²Liu, *A strong lower bound for approximate nearest neighbor searching*

³Malkov & Yashunin, *Efficient and Robust [ANNS] Using Hierarchical Navigable Small World Graphs*

⁴Raschka et al., *Machine Learning in Python: Main developments and technology trends [...]*

Out-of-Core Similarity Search

(used by your favourite semantic search engine)

Basically out-of-core ANNS with n base points



k -Means Clustering



Exact NNS

k base points

In-Memory ANNS

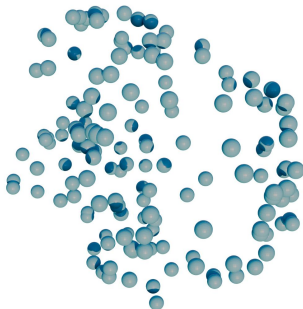
k base points



GPU Acceleration



HNSW

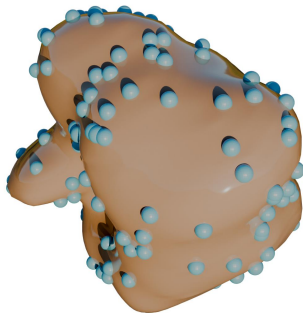


Structure

- ❖ Start with (approximate) NN graph
- ❖ Prune edges with a heuristic
- ❖ Randomly subsample points to get higher layers (similar to skip list)
- ❖ Build/insertions also similar to skip list

Search

- ❖ Start at arbitrary point on top level
- ❖ Greedy local search
- ❖ “Seed” lower layers with result

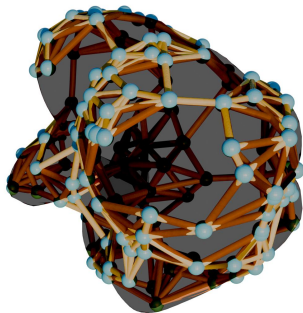


Structure

- ❖ Start with (approximate) NN graph
- ❖ Prune edges with a heuristic
- ❖ Randomly subsample points to get higher layers (similar to skip list)
- ❖ Build/insertions also similar to skip list

Search

- ❖ Start at arbitrary point on top level
- ❖ Greedy local search
- ❖ “Seed” lower layers with result

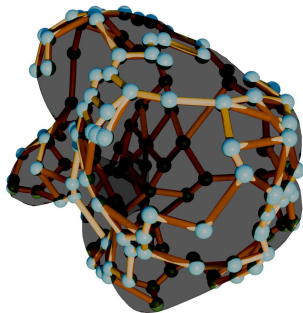


Structure

- ❖ **Start with (approximate) NN graph**
- ❖ Prune edges with a heuristic
- ❖ Randomly subsample points to get higher layers (similar to skip list)
- ❖ Build/insertions also similar to skip list

Search

- ❖ Start at arbitrary point on top level
- ❖ Greedy local search
- ❖ “Seed” lower layers with result

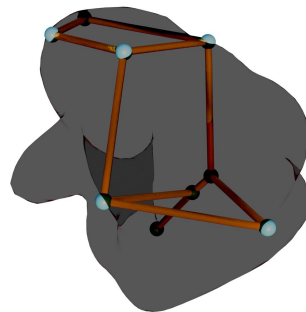
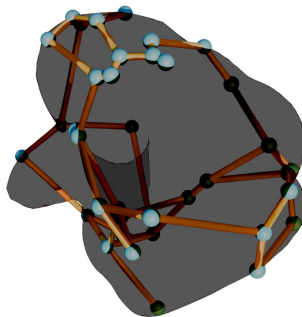
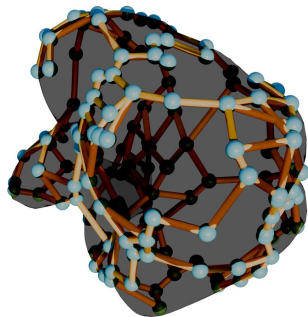


Structure

- ❖ Start with (approximate) NN graph
- ❖ **Prune edges with a heuristic**
- ❖ Randomly subsample points to get higher layers (similar to skip list)
- ❖ Build/insertions also similar to skip list

Search

- ❖ Start at arbitrary point on top level
- ❖ Greedy local search
- ❖ “Seed” lower layers with result

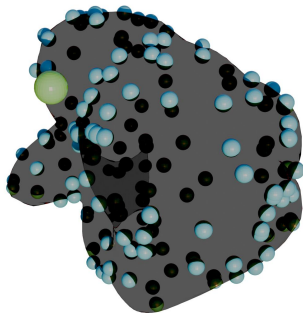


Structure

- ❖ Start with (approximate) NN graph
- ❖ Prune edges with a heuristic
- ❖ **Randomly subsample points to get higher layers (similar to skip list)**
- ❖ Build/insertions also similar to skip list

Search

- ❖ Start at arbitrary point on top level
- ❖ Greedy local search
- ❖ “Seed” lower layers with result

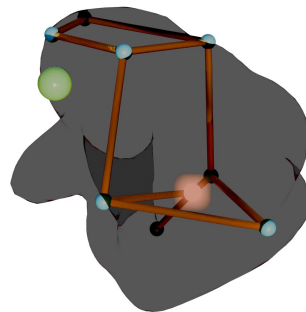
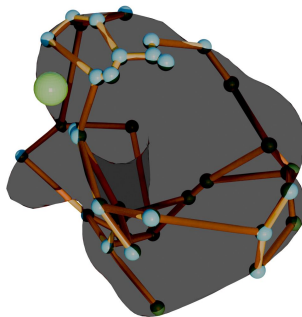
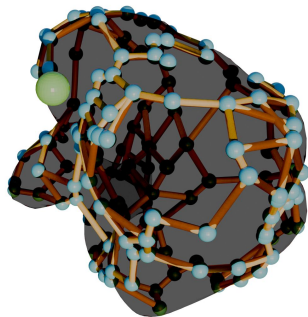


Structure

- ❖ Start with (approximate) NN graph
- ❖ Prune edges with a heuristic
- ❖ Randomly subsample points to get higher layers (similar to skip list)
- ❖ Build/insertions also similar to skip list

Search

- ❖ Start at arbitrary point on top level
- ❖ Greedy local search
- ❖ “Seed” lower layers with result

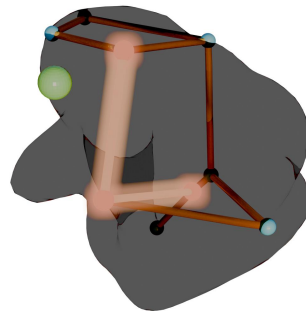
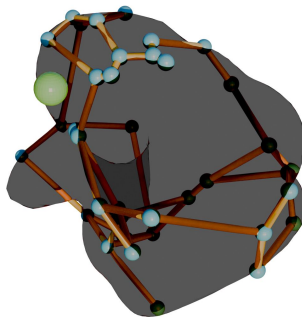
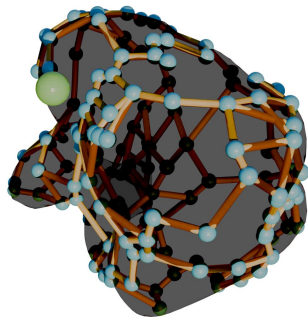


Structure

- ❖ Start with (approximate) NN graph
- ❖ Prune edges with a heuristic
- ❖ Randomly subsample points to get higher layers (similar to skip list)
- ❖ Build/insertions also similar to skip list

Search

- ❖ **Start at arbitrary point on top level**
- ❖ Greedy local search
- ❖ “Seed” lower layers with result

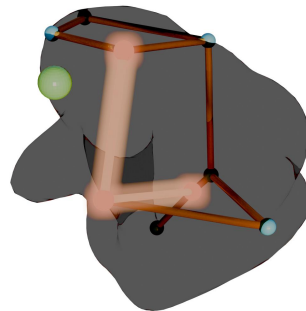
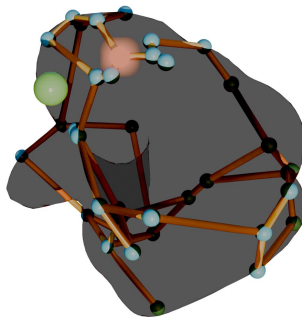
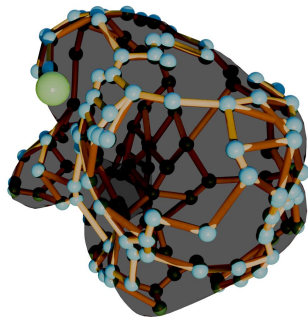


Structure

- ❖ Start with (approximate) NN graph
- ❖ Prune edges with a heuristic
- ❖ Randomly subsample points to get higher layers (similar to skip list)
- ❖ Build/insertions also similar to skip list

Search

- ❖ Start at arbitrary point on top level
- ❖ **Greedy local search**
- ❖ “Seed” lower layers with result

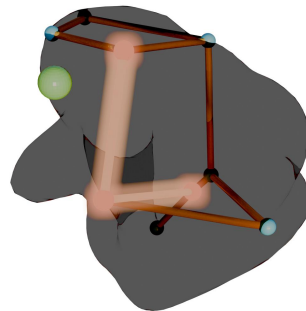
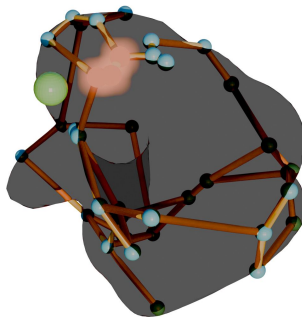
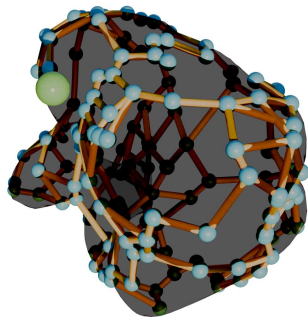


Structure

- ❖ Start with (approximate) NN graph
- ❖ Prune edges with a heuristic
- ❖ Randomly subsample points to get higher layers (similar to skip list)
- ❖ Build/insertions also similar to skip list

Search

- ❖ Start at arbitrary point on top level
- ❖ Greedy local search
- ❖ **“Seed” lower layers with result**

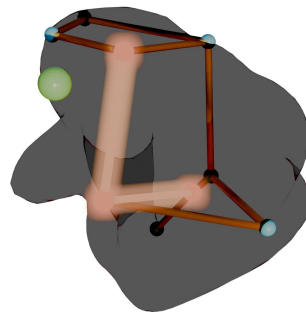
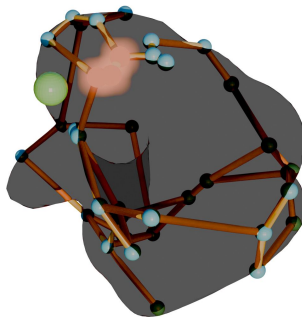
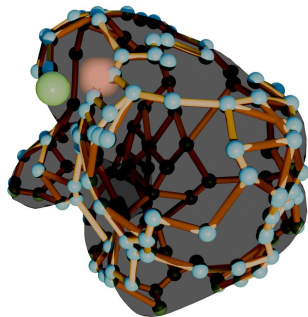


Structure

- ❖ Start with (approximate) NN graph
- ❖ Prune edges with a heuristic
- ❖ Randomly subsample points to get higher layers (similar to skip list)
- ❖ Build/insertions also similar to skip list

Search

- ❖ Start at arbitrary point on top level
- ❖ **Greedy local search**
- ❖ “Seed” lower layers with result

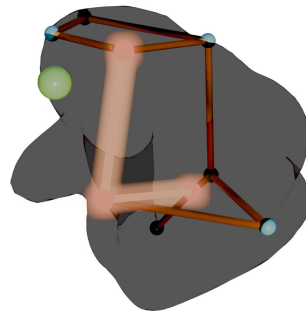
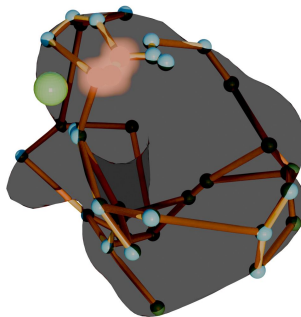
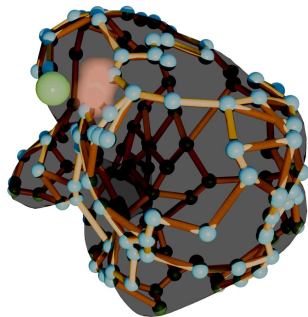


Structure

- ❖ Start with (approximate) NN graph
- ❖ Prune edges with a heuristic
- ❖ Randomly subsample points to get higher layers (similar to skip list)
- ❖ Build/insertions also similar to skip list

Search

- ❖ Start at arbitrary point on top level
- ❖ Greedy local search
- ❖ **“Seed” lower layers with result**



Structure

- ❖ Start with (approximate) NN graph
- ❖ Prune edges with a heuristic
- ❖ Randomly subsample points to get higher layers (similar to skip list)
- ❖ Build/insertions also similar to skip list

Search

- ❖ Start at arbitrary point on top level
- ❖ **Greedy local search**
- ❖ “Seed” lower layers with result

Normal ANNS queries: given a query point, find a good close neighbor.

Normal ANNS queries: given a query point, find a good close neighbor.

Seeded ANNS queries: also given **seed points** (candidate close neighbors).
Very simple learning-augmented form of ANNS.

Normal ANNS queries: given a query point, find a good close neighbor.

Seeded ANNS queries: also given **seed points** (candidate close neighbors).

Very simple learning-augmented form of ANNS.

Want:

- ❑ robustness — good solutions even if seed points bad
- ❑ consistency — better solutions if seed points good

Normal ANNS queries: given a query point, find a good close neighbor.

Seeded ANNS queries: also given **seed points** (candidate close neighbors).

Very simple learning-augmented form of ANNS.

Want:

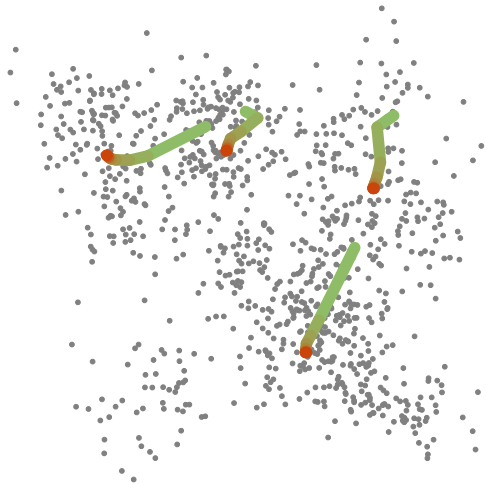
- ❑ robustness — good solutions even if seed points bad
- ❑ consistency — better solutions if seed points good

In HNSW: Add seed points right at start of search on last level. Similar for other search-graph methods (result: **seeded search-graphs**).

Main Improvements: Beating Black-Box HNSW

Centroids “slow down” over time:

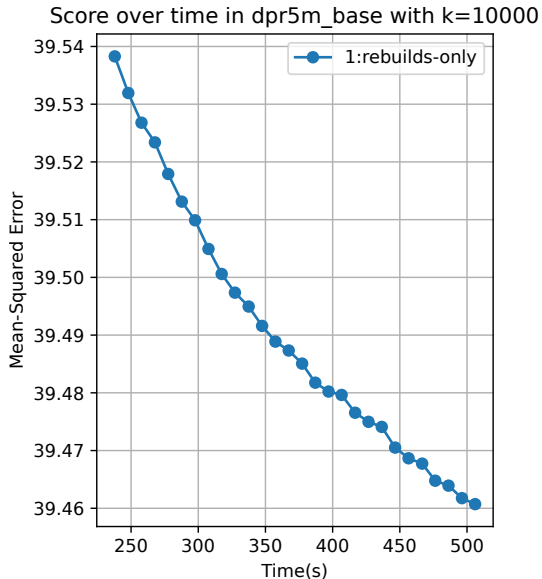
- ✦ Rebuilds (HNSW as a kind of kinetic data structure, omitting details)
- ✦ Extra “seed points” from prev assignment:
Seeded ANNS



Main Improvements: Beating Black-Box HNSW

Centroids “slow down” over time:

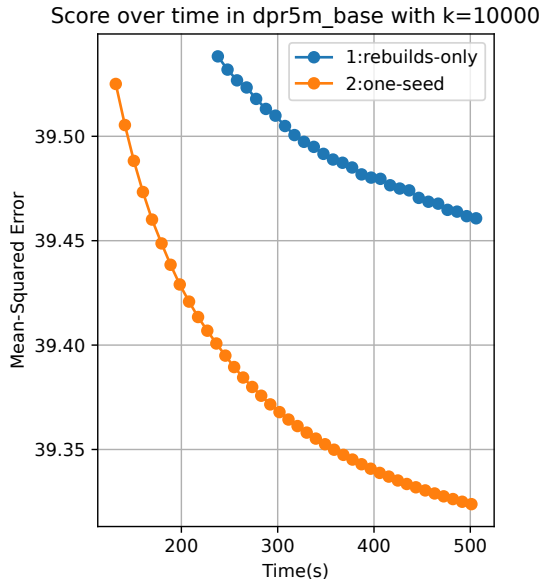
- ✚ **Rebuilds** (HNSW as a kind of kinetic data structure, omitting details)
- ✚ Extra “seed points” from prev assignment: Seeded ANNS



Main Improvements: Beating Black-Box HNSW

Centroids “slow down” over time:

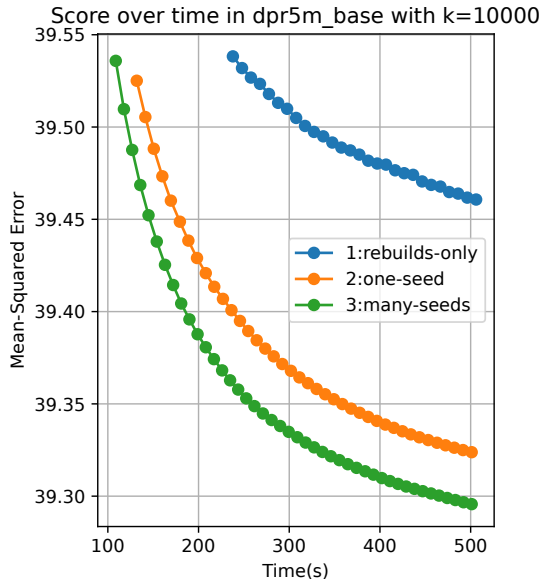
- ❖ Rebuilds (HNSW as a kind of kinetic data structure, omitting details)
- ❖ **Extra “seed points” from prev assignment:**
Seeded ANNS



Main Improvements: Beating Black-Box HNSW

Centroids “slow down” over time:

- ❖ Rebuilds (HNSW as a kind of kinetic data structure, omitting details)
- ❖ **Multiple** extra “seed points” from prev assignment: Seeded ANNS



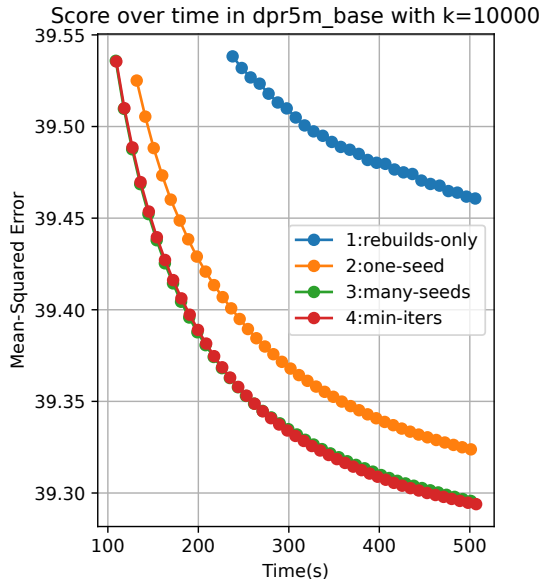
Main Improvements: Beating Black-Box HNSW

Centroids “slow down” over time:

- ❖ Rebuilds (HNSW as a kind of kinetic data structure, omitting details)
- ❖ Multiple extra “seed points” from prev assignment: Seeded ANNS

Often also improve:

- ❖ **Min iteration threshold**
- ❖ Bulk queries for more seed points



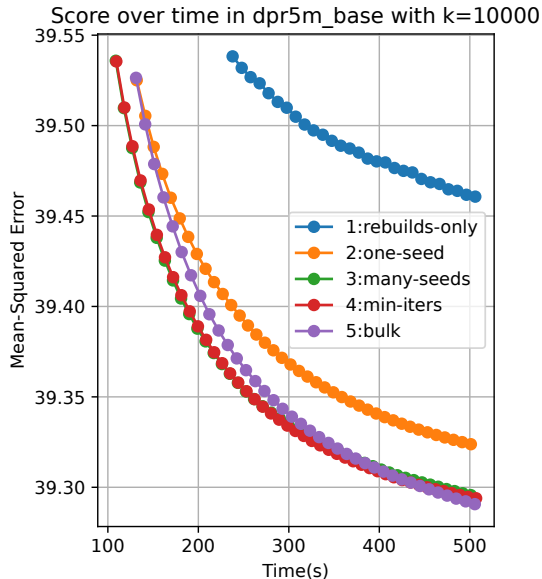
Main Improvements: Beating Black-Box HNSW

Centroids “slow down” over time:

- ❑ Rebuilds (HNSW as a kind of kinetic data structure, omitting details)
- ❑ Multiple extra “seed points” from prev assignment: Seeded ANNS

Often also improve:

- ❑ Min iteration threshold
- ❑ **Bulk queries for more seed points**



Main Improvements: Beating Black-Box HNSW

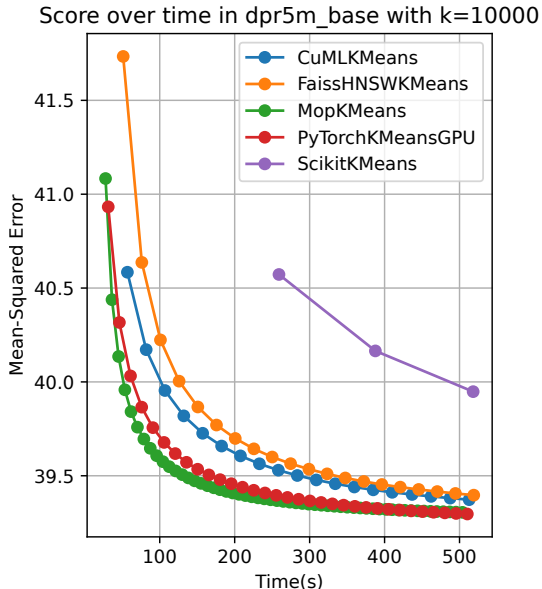
Centroids “slow down” over time:

- ❖ Rebuilds (HNSW as a kind of kinetic data structure, omitting details)
- ❖ Multiple extra “seed points” from prev assignment: Seeded ANNS

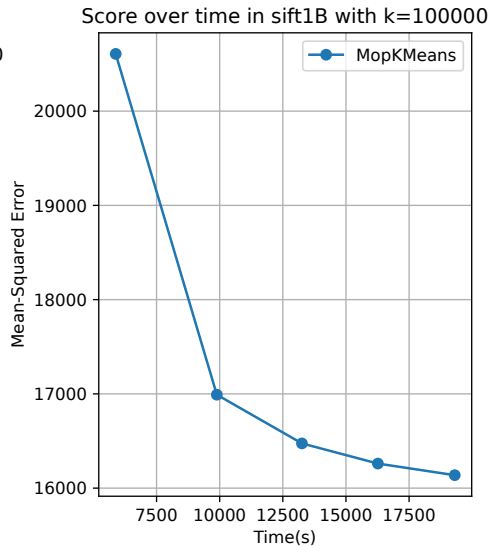
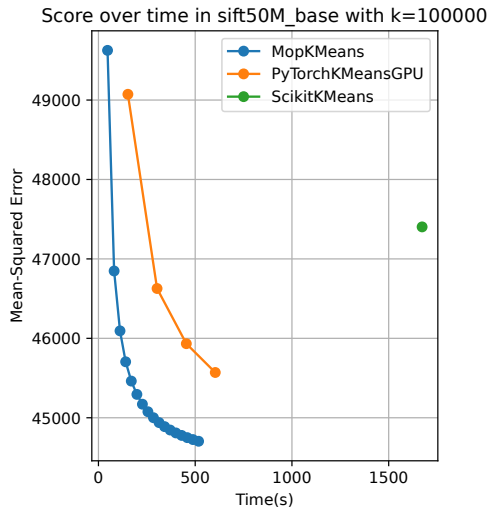
Often also improve:

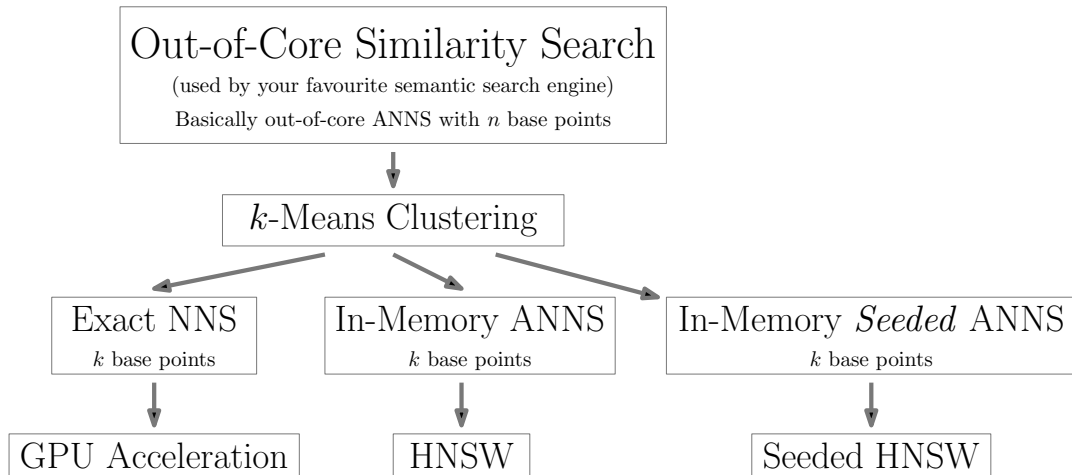
- ❖ Min iteration threshold
- ❖ Bulk queries for more seed points

Now (mostly) **beating GPU implementations with CPU**.



More Results





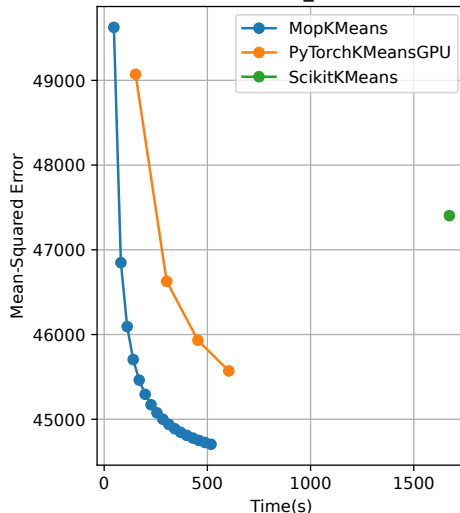
Open Problems

- ❖ GPU-acceleration of SANNS? No current methods for ANNS on GPU extend well.
- ❖ Perfect consistency guarantees for SANNS in fixed doubling dimension? Search-graph methods don't seem to work.
- ❖ General open ANNS problem: Better theoretical understanding of why search-graphs work well? Best known is (tunable) additive approximation known for one specific graph algorithm, with fixed doubling-dimension¹.



arxiv.org/abs/2502.06163

Score over time in sift50M_base with k=100000



¹Indyk & Xu, *Worst-case performance of popular approximate nearest neighbor search implementations* [...]

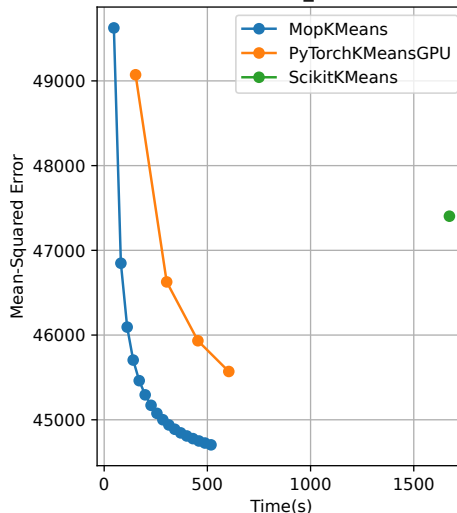
Open Problems

- ❖ GPU-acceleration of SANNS? No current methods for ANNS on GPU extend well.
- ❖ Perfect consistency guarantees for SANNS in fixed doubling dimension? Search-graph methods don't seem to work.
- ❖ General open ANNS problem: Better theoretical understanding of why search-graphs work well? Best known is (tunable) additive approximation known for one specific graph algorithm, with fixed doubling-dimension¹.



arxiv.org/abs/2502.06163

Score over time in sift50M_base with k=100000



Fin.

¹Indyk & Xu, *Worst-case performance of popular approximate nearest neighbor search implementations [...]*