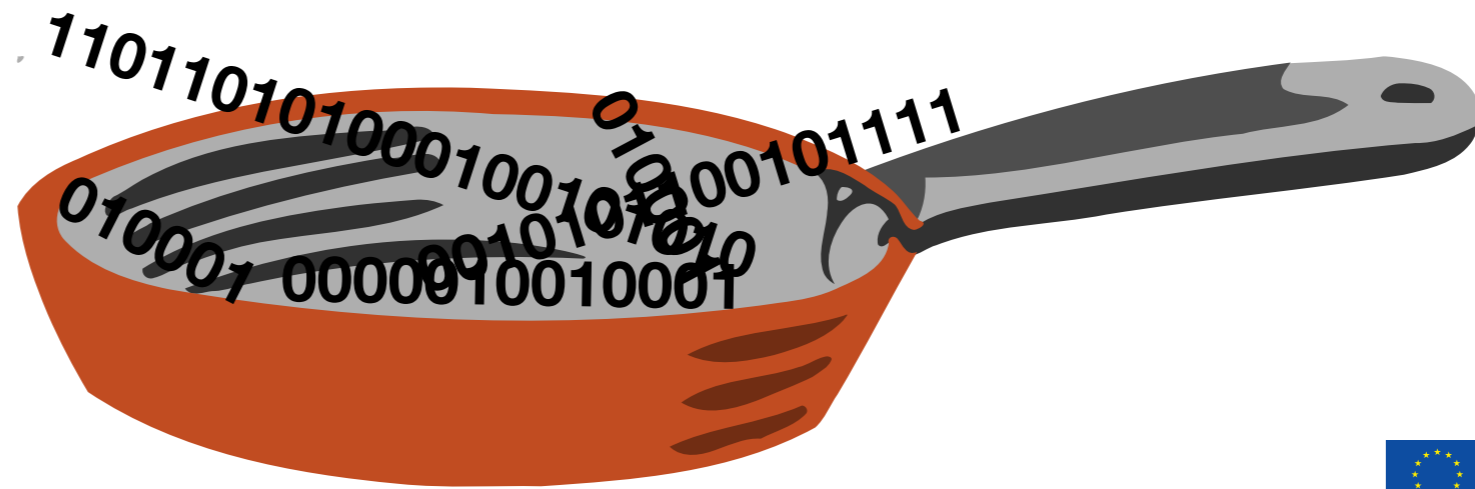


Efficient Reductions for k-Nearest Neighbor Search

Joint work with Tobias Christiani and Mikkel Thorup



Rasmus Pagh

BARC and IT University of Copenhagen

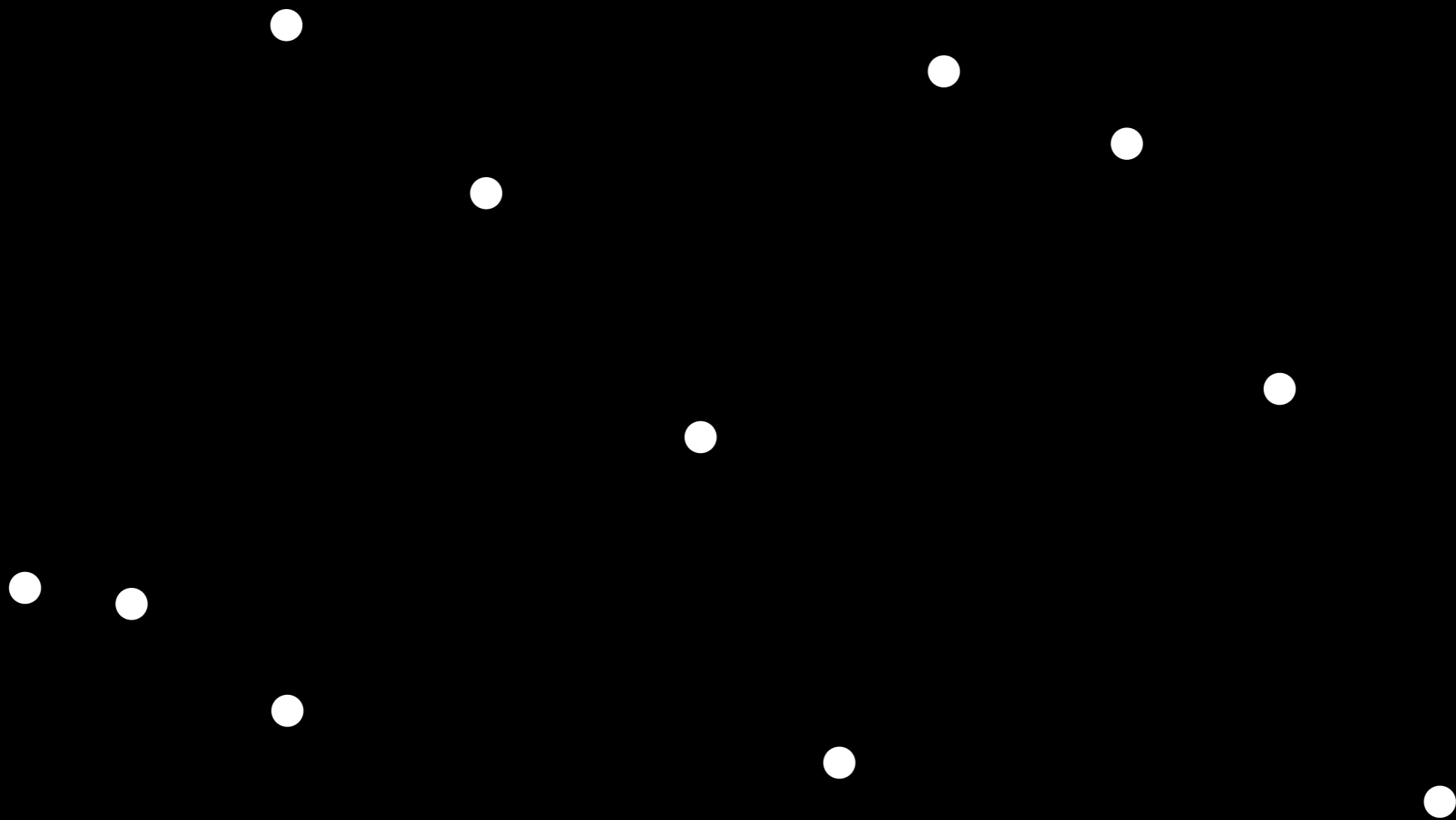
Workshop on Sublinear Algorithms
and Nearest Neighbor Search

Simons Institute for Theory of Computing
November 29, 2018



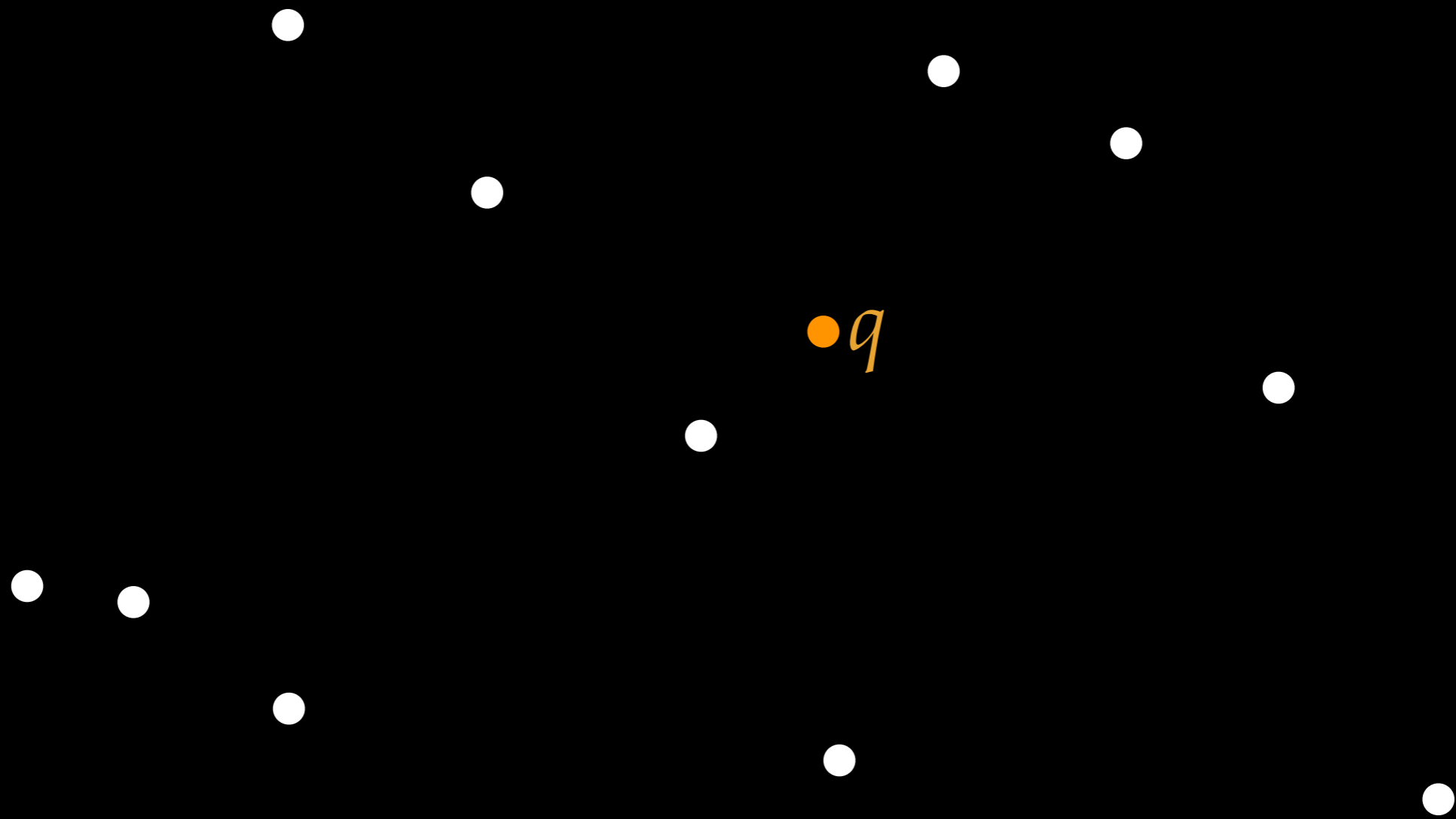
**SCALABLE
SIMILARITY
SEARCH**

Nearest neighbor search



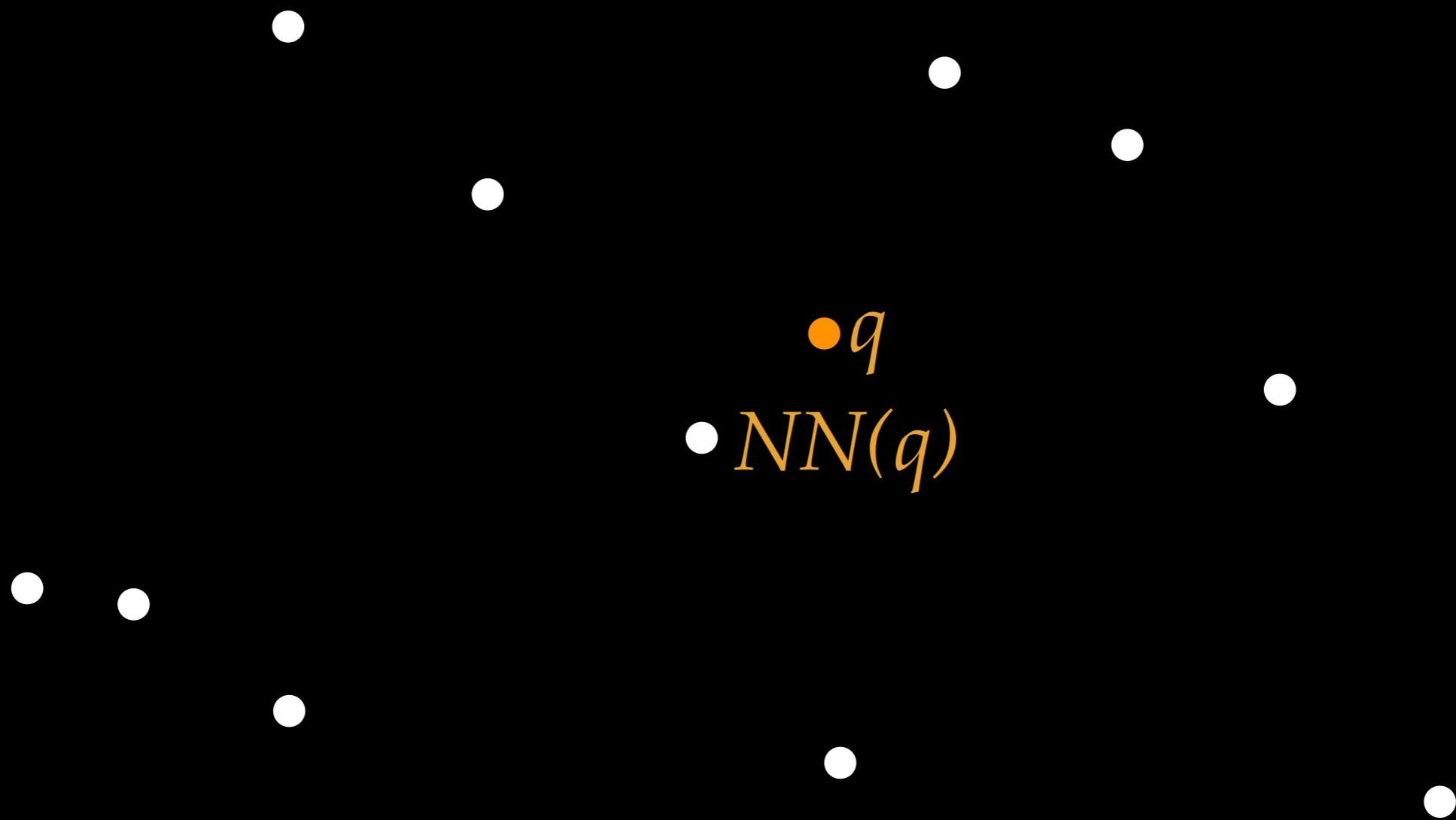
$$P = \{x_1, \dots, x_n\}$$

Nearest neighbor search



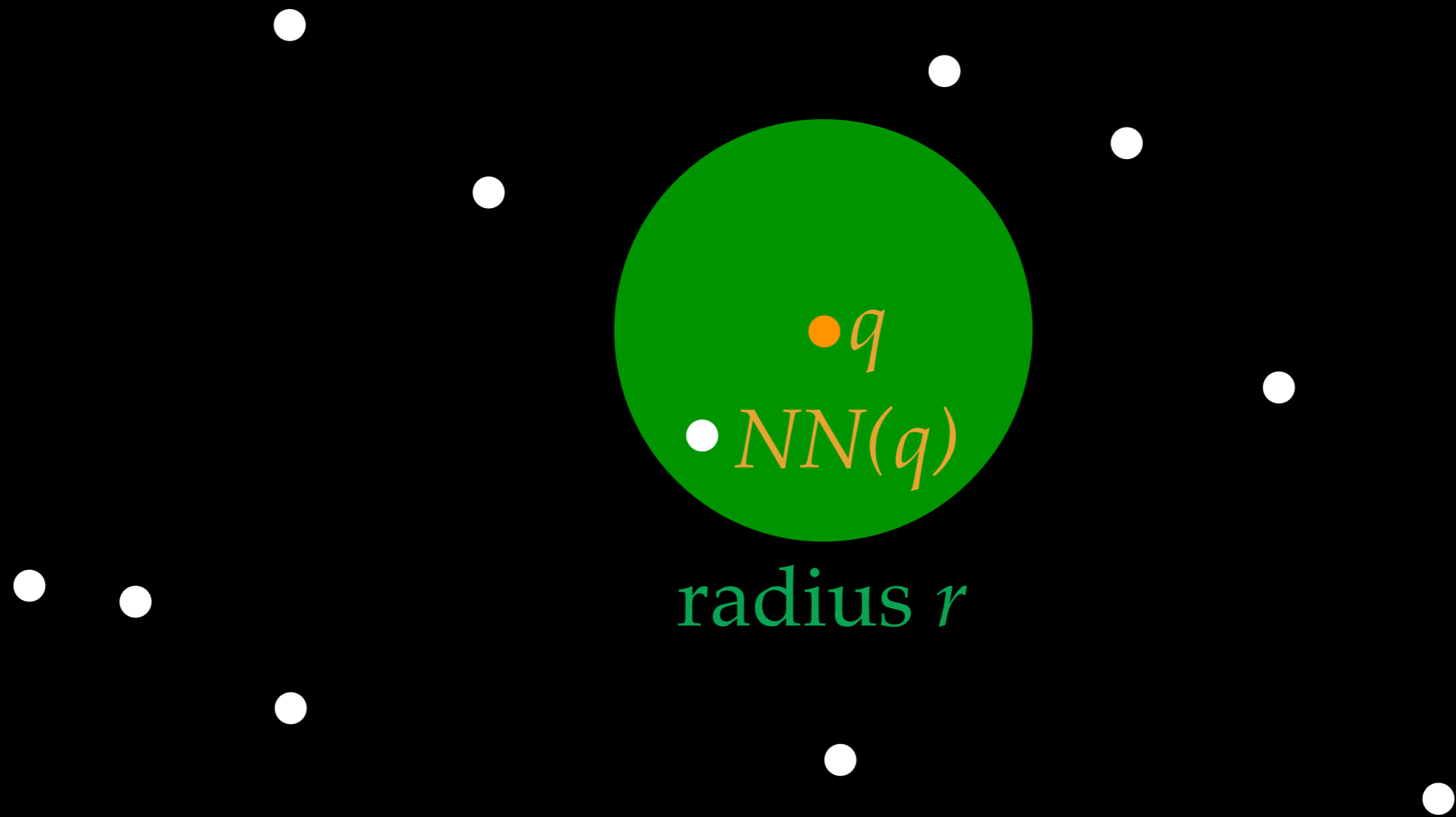
$$P = \{x_1, \dots, x_n\}$$

Nearest neighbor search



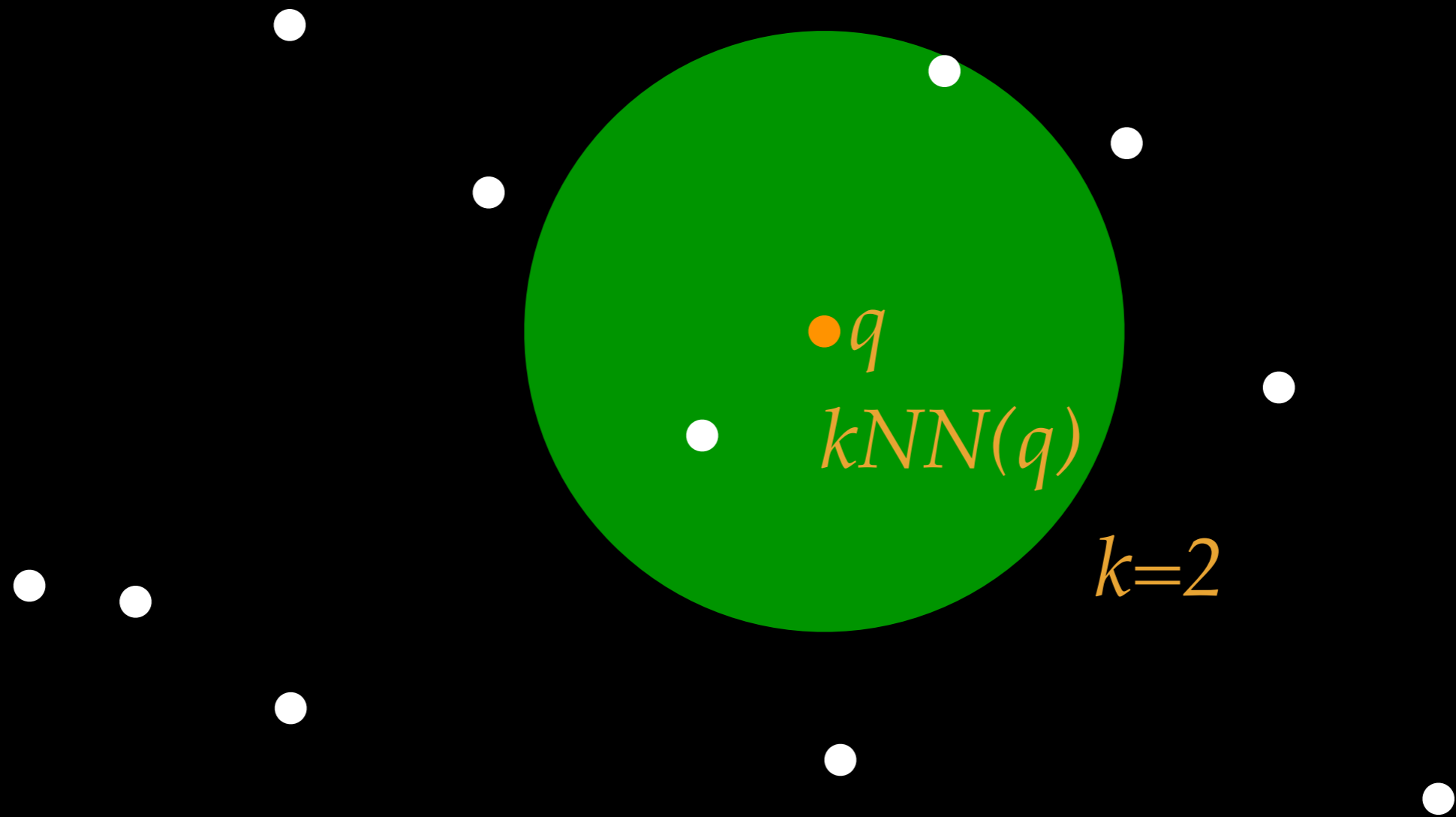
$$P = \{x_1, \dots, x_n\}$$

Nearest ~~neighbor~~ search



$$P = \{x_1, \dots, x_n\}$$

k -nearest neighbor search



$$P = \{x_1, \dots, x_n\}$$

Some application areas

Some application areas

- Association rule mining
- Automation
- Bio-chemistry (finding motifs)
- Bio-informatics (homology search)
- Clustering
- Computer vision and pattern recognition
- Databases
- Data cleaning
- Data stream computation
- Data privacy
- First story detection (with application to Twitter)
- Identifying trends in time series
- Linear algebra
- Motion planning for robots
- Near-duplicate detection
- News personalization (collaborative filtering)
- Privacy preserving data mining
- Search engines for 3D models
- Sensor networks
- ...

Hardness of NN search

- [Williams '04], [Alman & Williams '15]:
NN search on $P \subseteq \{0,1\}^d$ in time $n^{0.99} 2^{o(d)}$ with
preprocessing time $\text{poly}(n) 2^{o(d)} \implies$
 k -SAT w. n variables can be solved in time c^n , $c < 2$

Hardness of NN search

- [Williams '04], [Alman & Williams '15]:
NN search on $P \subseteq \{0,1\}^d$ in time $n^{0.99} 2^{o(d)}$ with
preprocessing time $\text{poly}(n) 2^{o(d)} \implies$
 k -SAT w. n variables can be solved in time c^n , $c < 2$

Under strong exponential time hypothesis, this is not possible!

Hardness of NN search

- [Williams '04], [Alman & Williams '15]:
NN search on $P \subseteq \{0,1\}^d$ in time $n^{0.99} 2^{o(d)}$ with
preprocessing time $\text{poly}(n) 2^{o(d)} \implies$
 k -SAT w. n variables can be solved in time c^n , $c < 2$

Under strong exponential time hypothesis, this is not possible!

In practice: "Curse of dimensionality" makes NN search slow in high dimension.

Approximate nearest neighbor

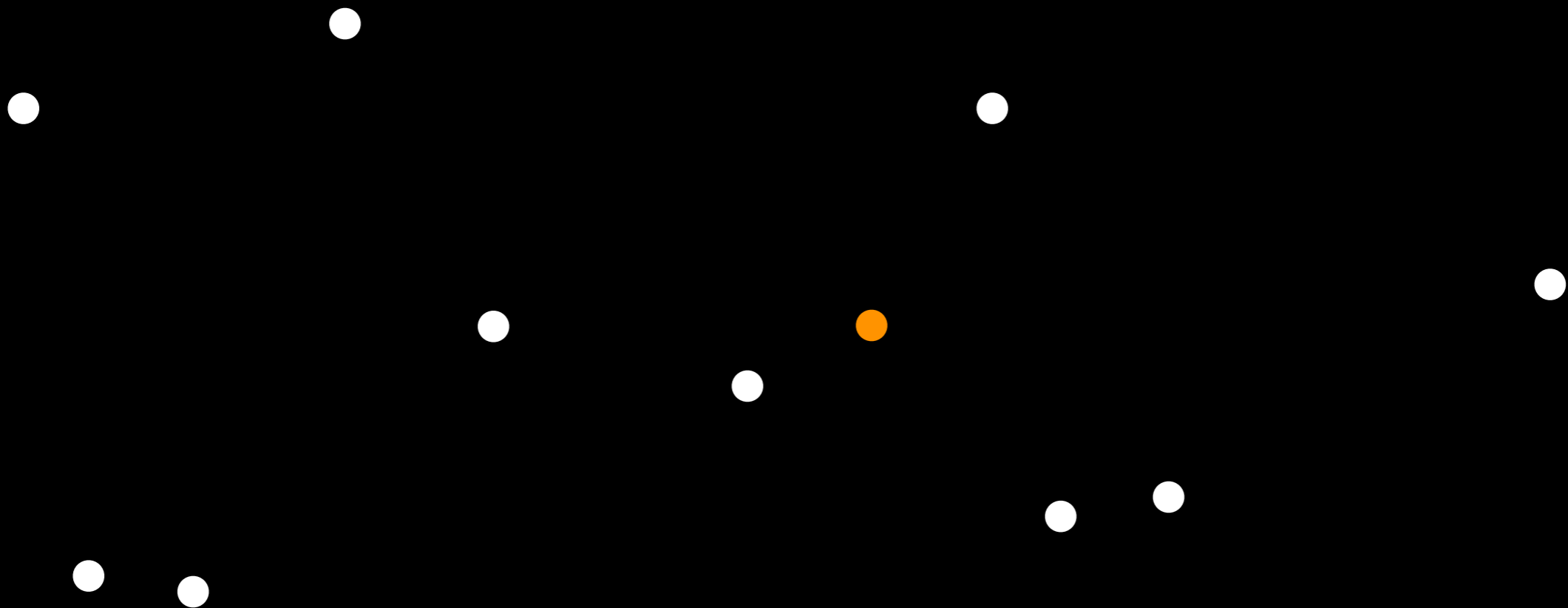
Approximate nearest neighbors: towards removing the curse of dimensionality

3794

1998

P Indyk, R Motwani

Proceedings of the thirtieth annual ACM symposium on Theory of computing ...



Approximate nearest neighbor

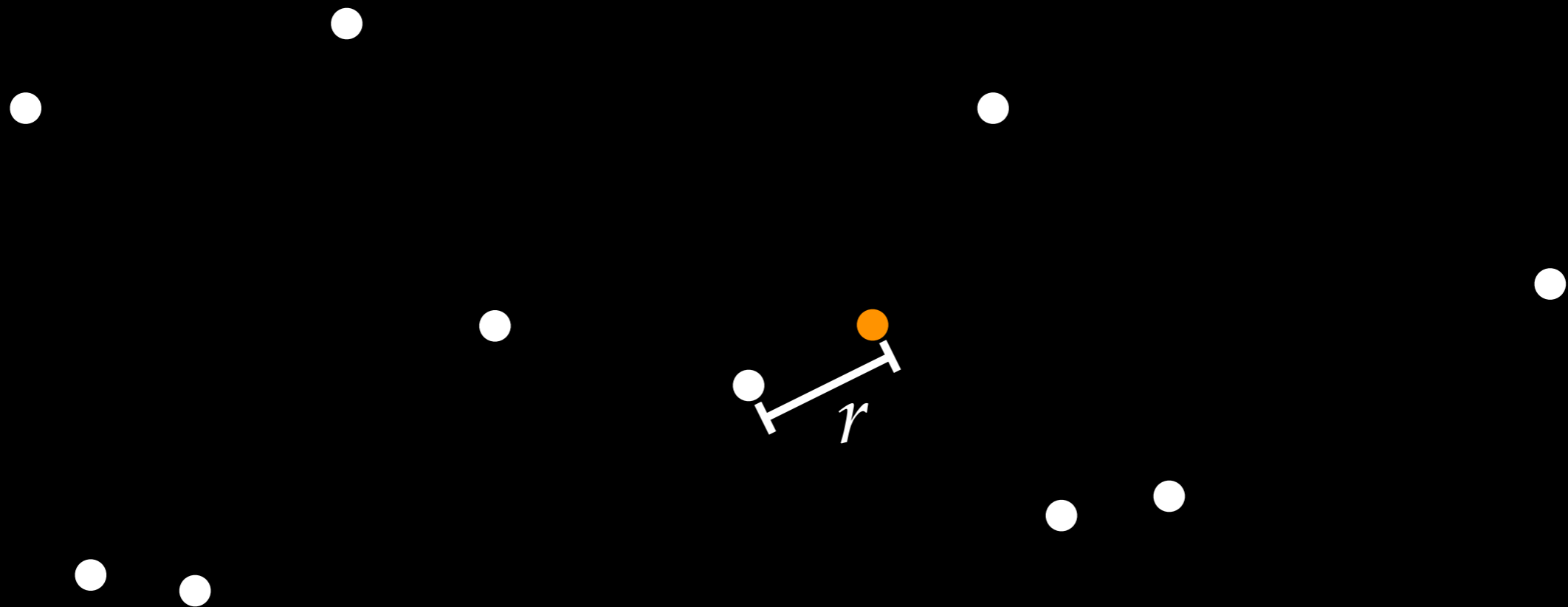
Approximate nearest neighbors: towards removing the curse of dimensionality

3794

1998

P Indyk, R Motwani

Proceedings of the thirtieth annual ACM symposium on Theory of computing ...



Approximate nearest neighbor

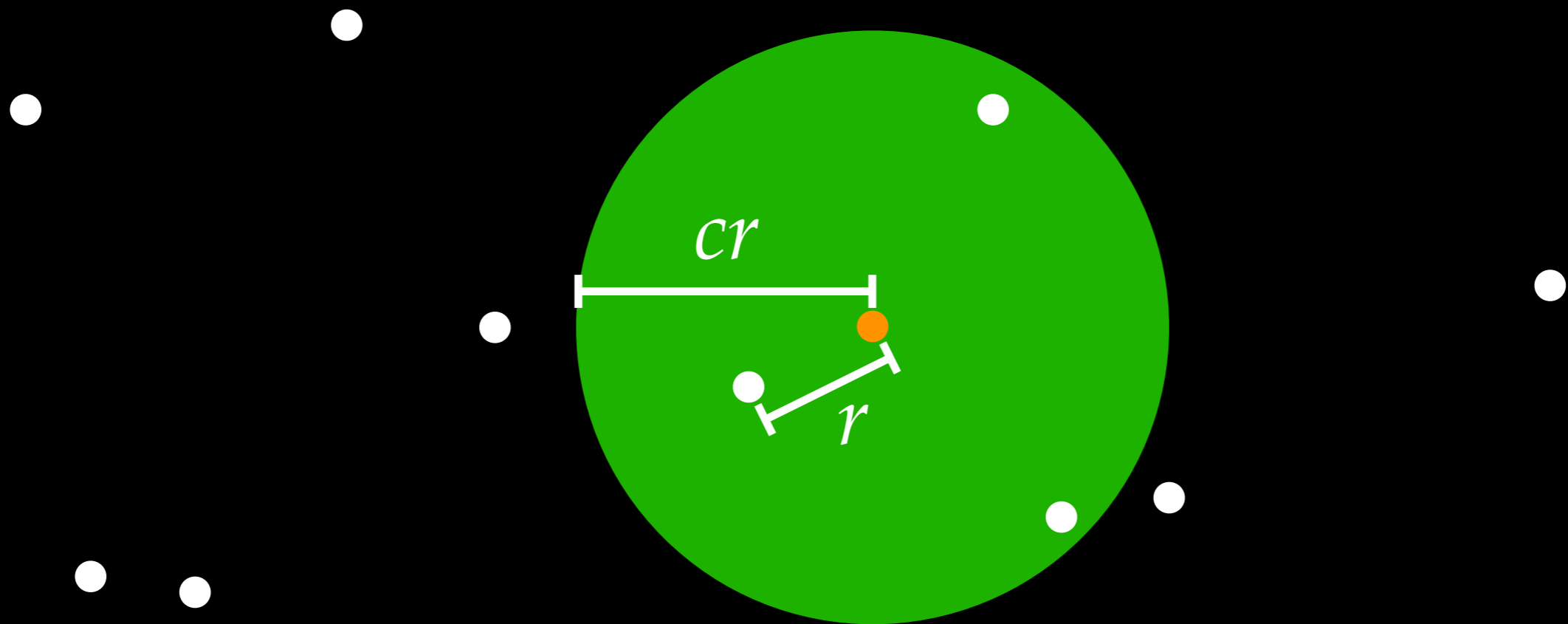
Approximate nearest neighbors: towards removing the curse of dimensionality

3794

1998

P Indyk, R Motwani

Proceedings of the thirtieth annual ACM symposium on Theory of computing ...



Approximate nearest neighbor

Approximate nearest neighbors: towards removing the curse of dimensionality

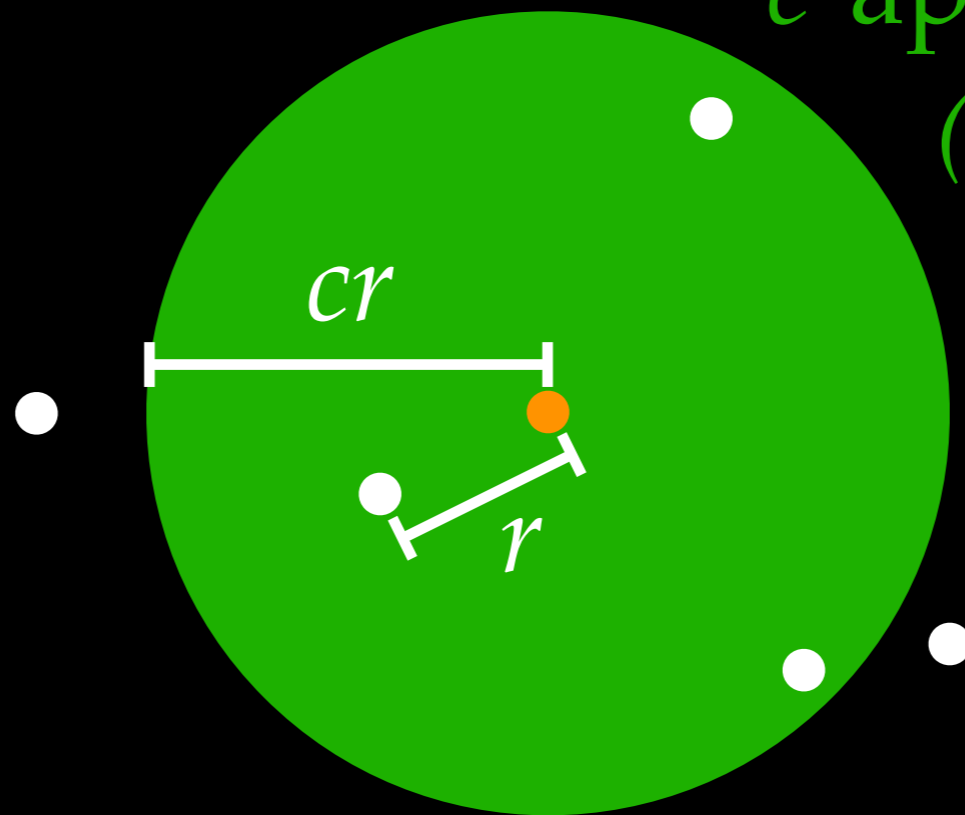
3794

1998

P Indyk, R Motwani

Proceedings of the thirtieth annual ACM symposium on Theory of computing ...

c-approximate NN
(return any one)



Approximate nearest neighbor

Approximate nearest neighbors: towards removing the curse of dimensionality

3794

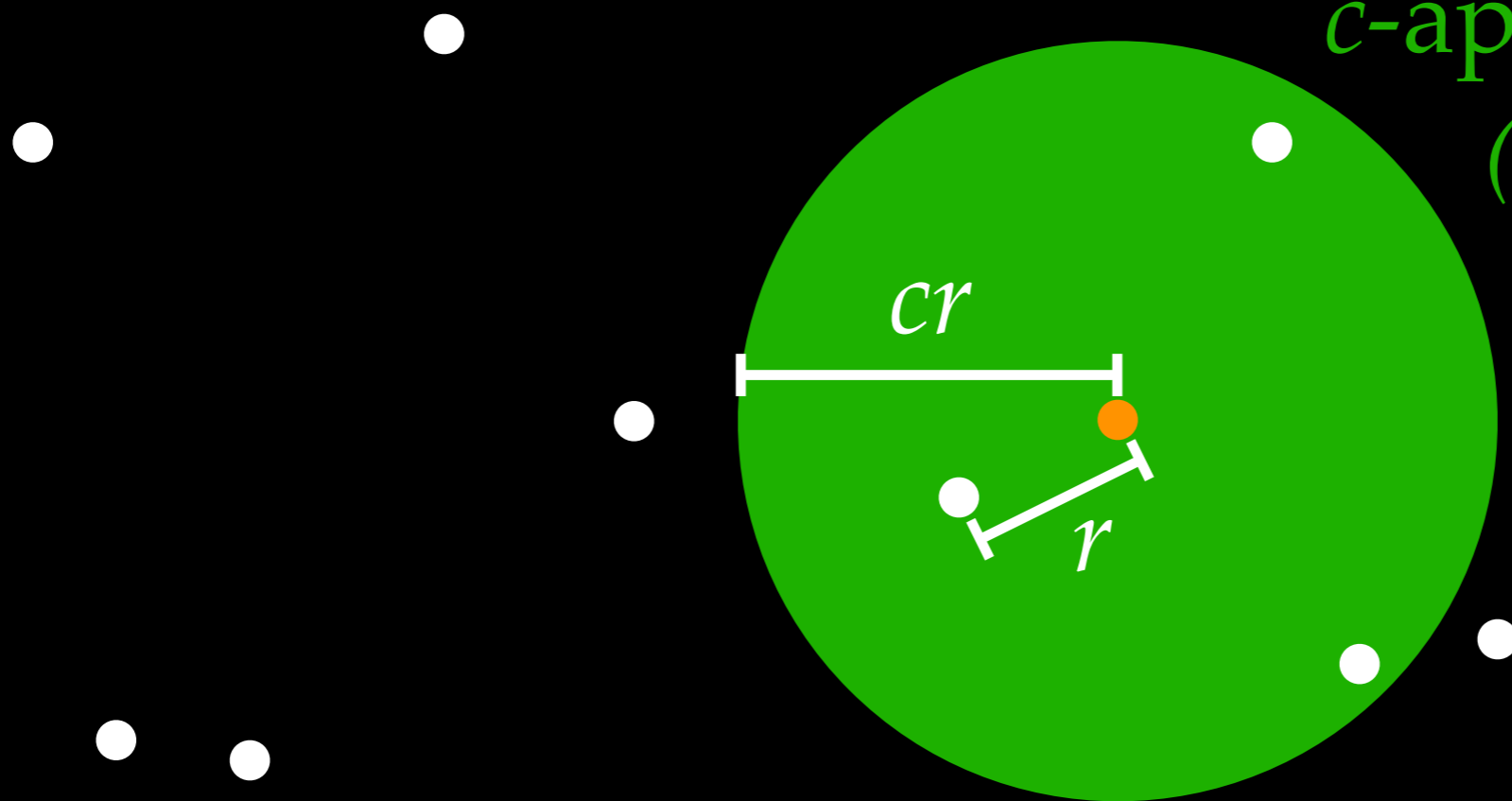
1998

P Indyk, R Motwani

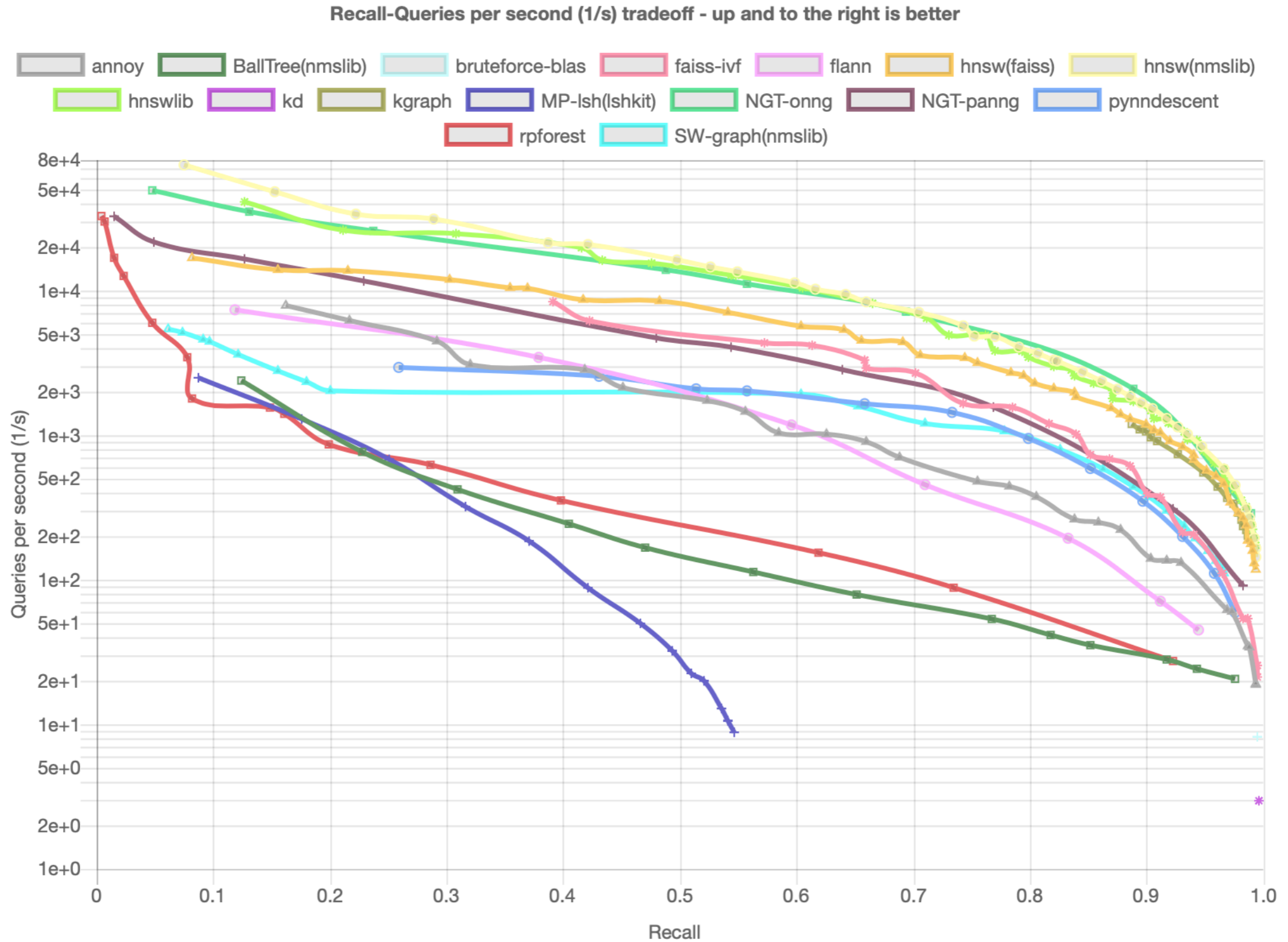
Proceedings of the thirtieth annual ACM symposium on Theory of computing ...

Time $n^{\rho(c)}$ **and space** $n^{1+\rho(c)}$, $\rho(c) < 1$ **for** $c > 1$

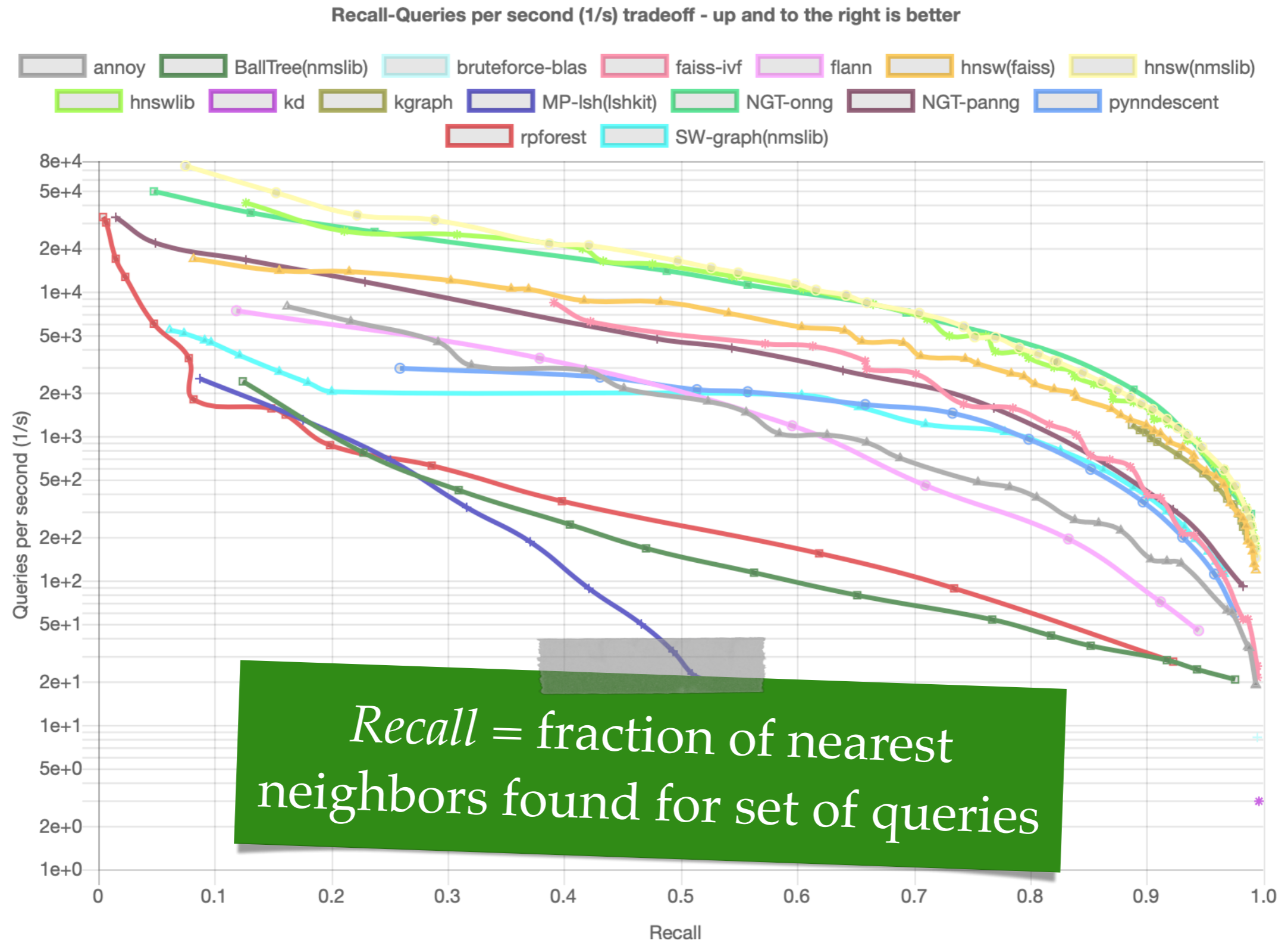
c-approximate NN
(return any one)



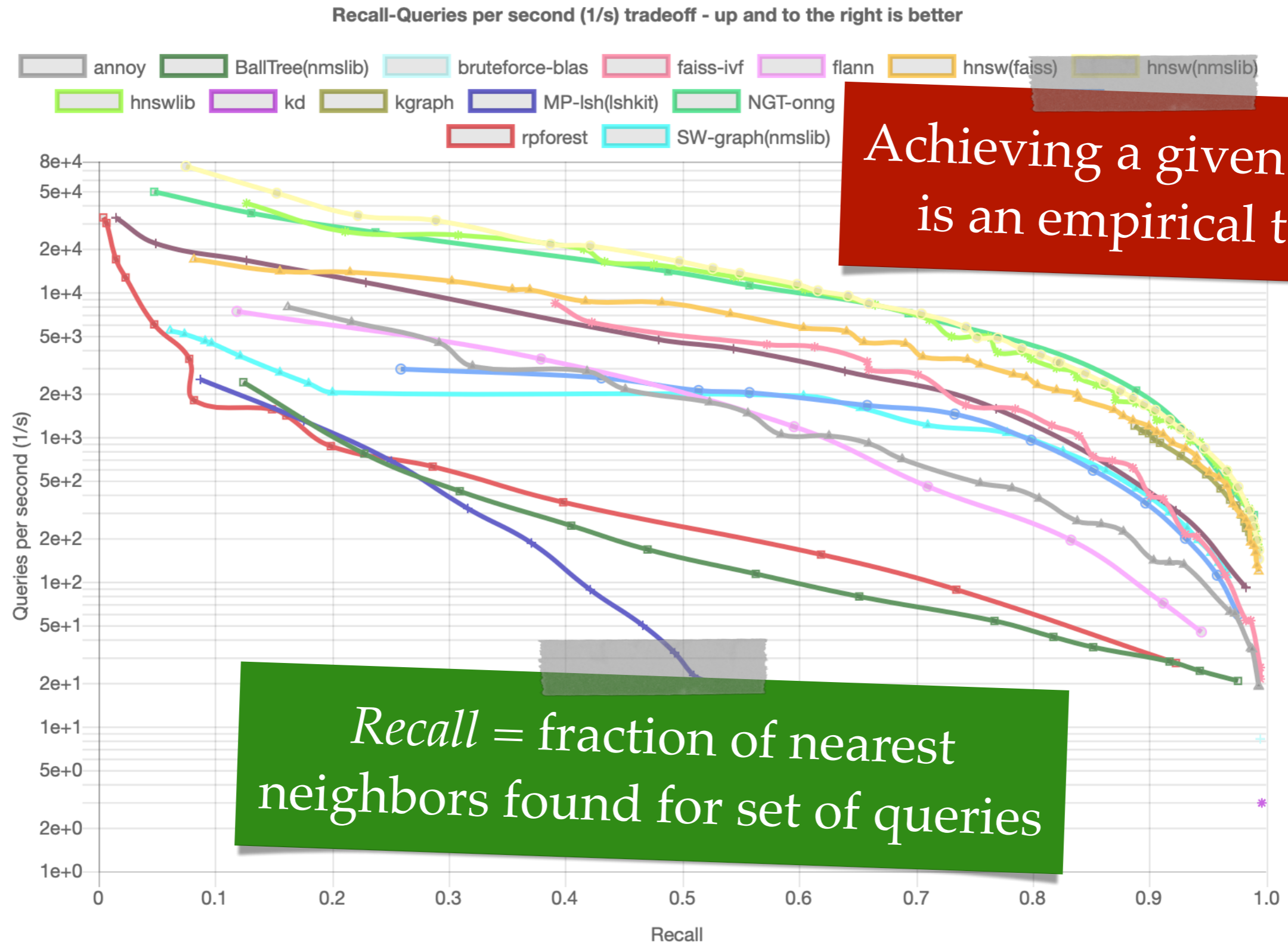
Approximate NN in practice



Approximate NN in practice



Approximate NN in practice

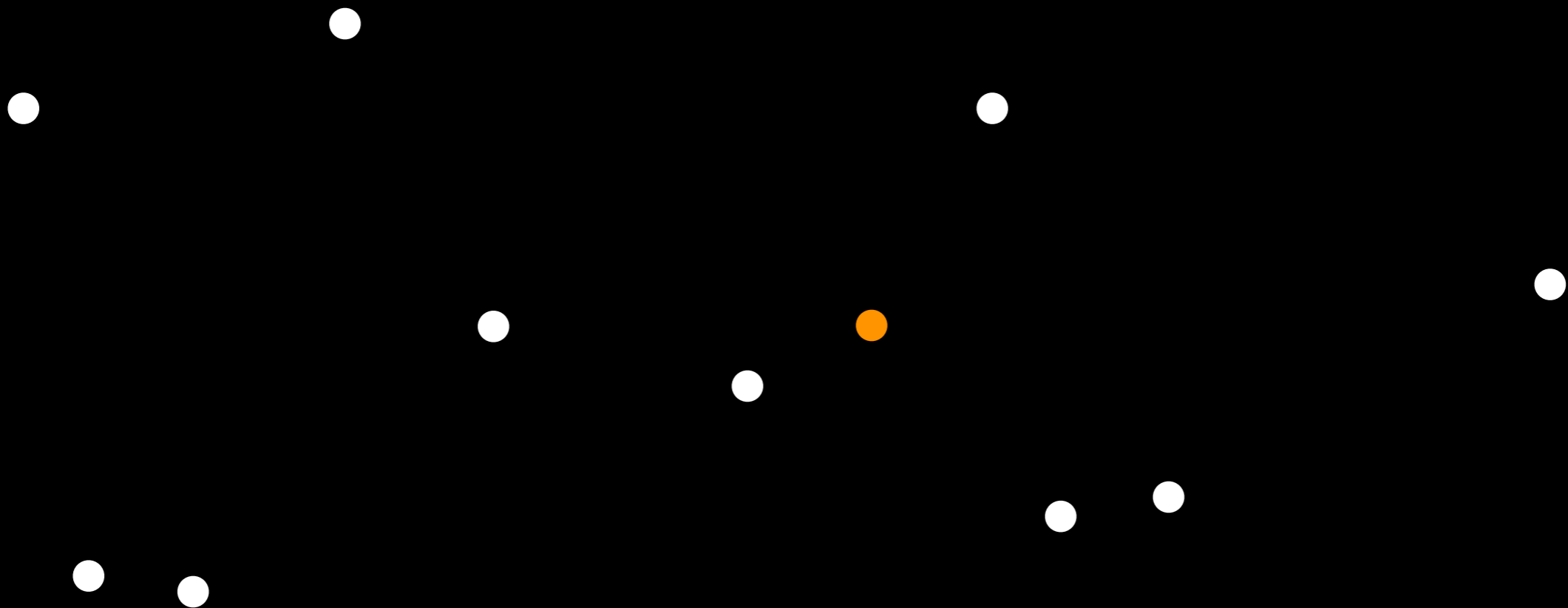


Achieving a given recall is an empirical task

Recall = fraction of nearest neighbors found for set of queries

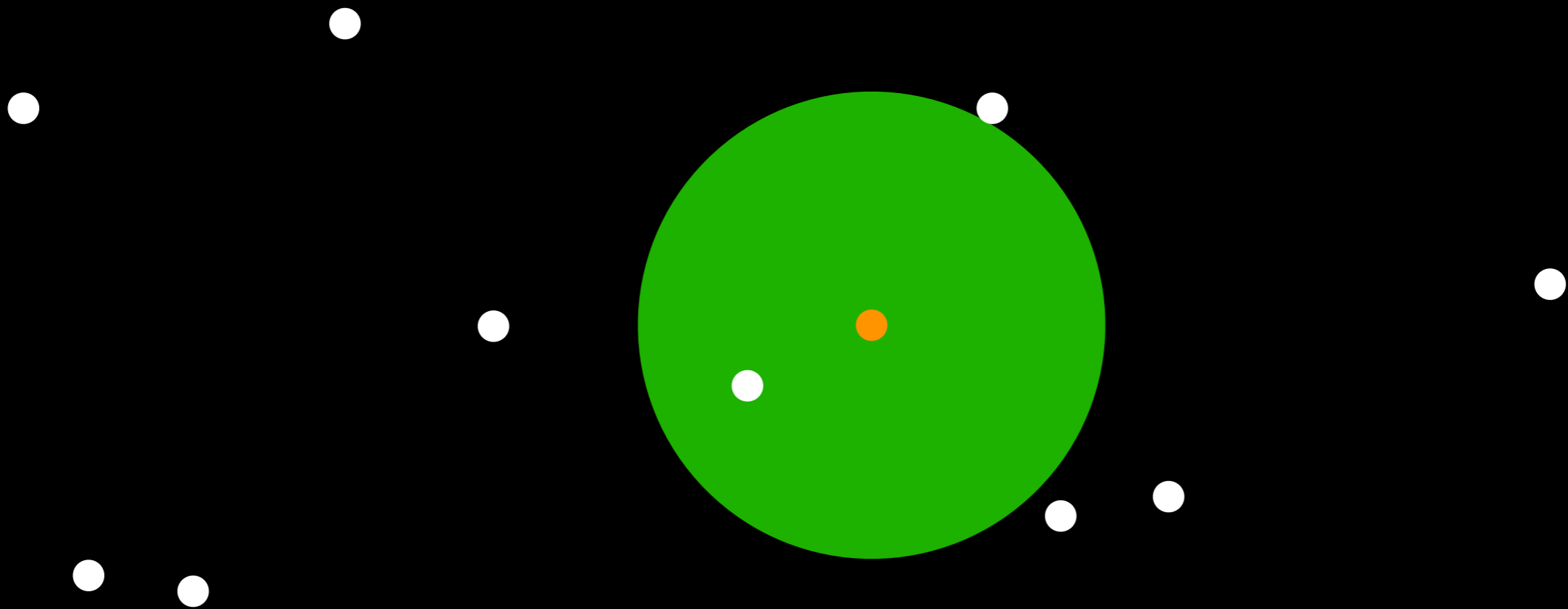
NN *recall* guarantee?

Black-box reduction: Choose c small enough to distinguish nearest neighbor from other points.



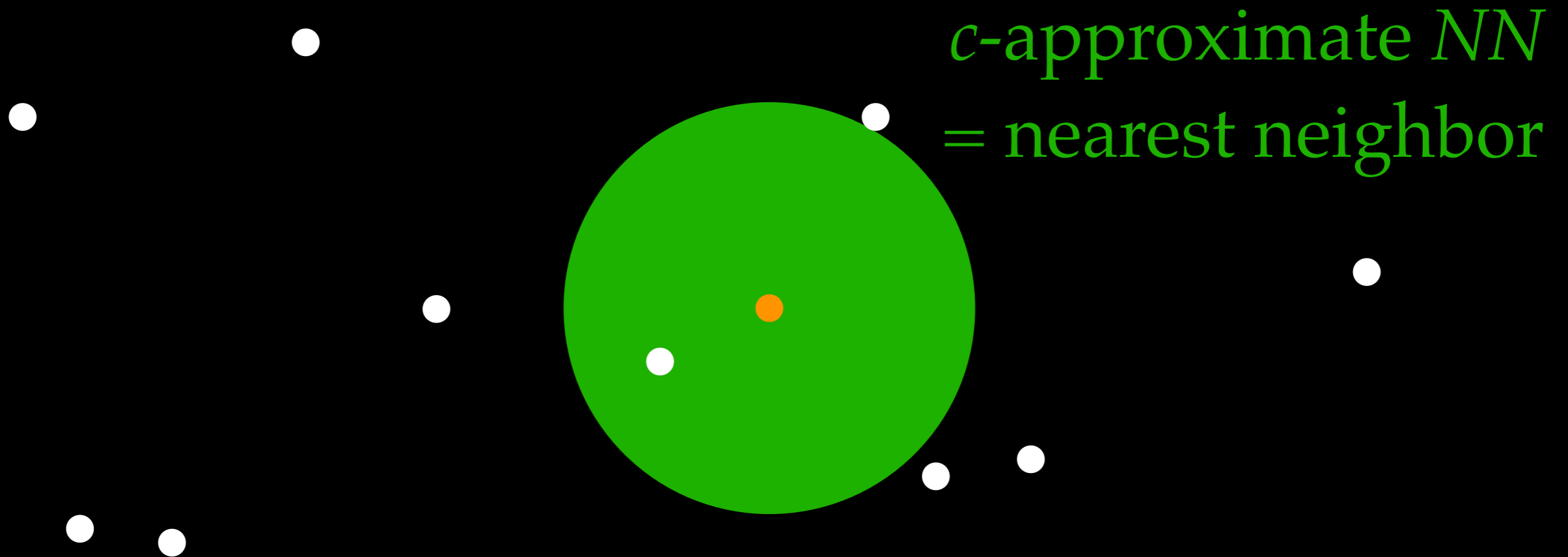
NN *recall* guarantee?

Black-box reduction: Choose c small enough to distinguish nearest neighbor from other points.



NN *recall* guarantee?

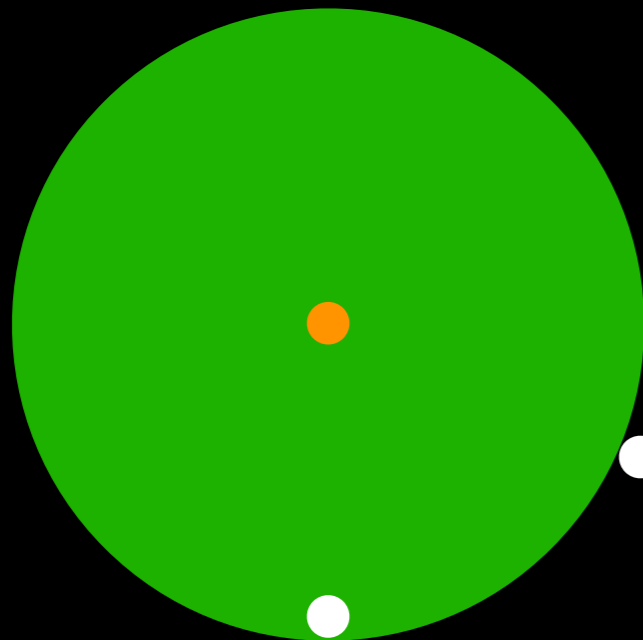
Black-box reduction: Choose c small enough to distinguish nearest neighbor from other points.



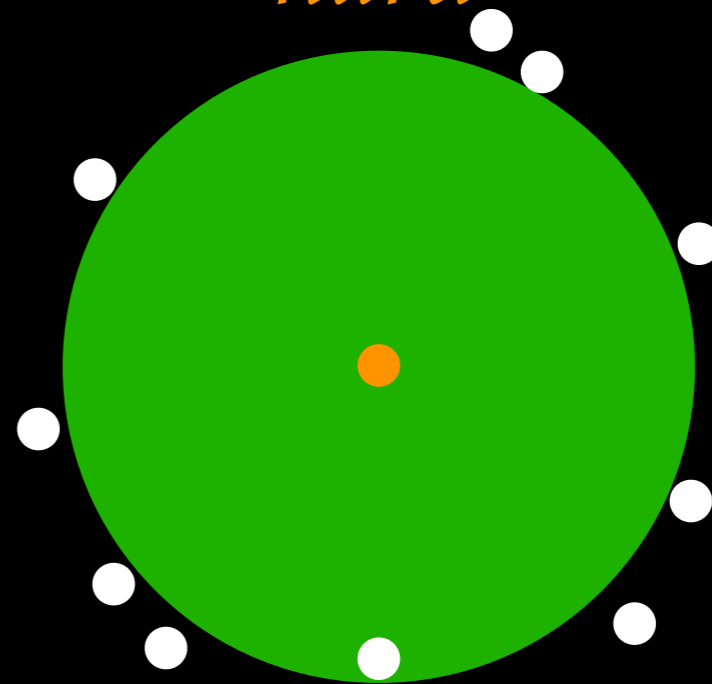
NN *recall* guarantee?

Black-box reduction does not distinguish easy and hard cases

easy



hard



LSH with multi-probing

- Building block: Linear space data structure (a hash table) that allows us to retrieve x_i with probability p_i .
 - Let $p_1 =$ retrieval probability of nearest neighbor

LSH with multi-probing

- Building block: Linear space data structure (a hash table) that allows us to retrieve x_i with probability p_i .

- Let $p_1 =$ retrieval probability of nearest neighbor

- Expected cost $O\left(K + T + \sum_i p_i\right)$

Parameters K, T

LSH with multi-probing

- Building block: Linear space data structure (a hash table) that allows us to retrieve x_i with probability p_i .
 - Let $p_1 =$ retrieval probability of nearest neighbor
- Expected cost $O\left(K + T + \sum_i p_i\right)$ Parameters K, T
- Repeat independently L times; fails to find NN with probability $(1 - p_1)^L$

LSH with multi-probing

- Building block: Linear space data structure (a hash table) that allows us to retrieve x_i with probability p_i .
 - Let $p_1 =$ retrieval probability of nearest neighbor
- Expected cost $O\left(K + T + \sum_i p_i\right)$
- Repeat independently L times; fails to find NN with probability

Parameters K, T

$$(1 - p_1)^L$$

Need $L \approx \ln(1 / \delta) / p_1$ for expected recall $1 - \delta$

From FALCONN documentation



ilyaraz

152 commits **165,070 ++** **63,535 --**

#1



ludwigschmidt

90 commits **278,244 ++** **16,349 --**

#2

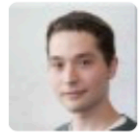
Choosing parameters

All in all, we have three parameters:

- K , the number of hash functions (space partitions) per hash table.
- L , the number of hash tables (each of which has K hash functions).
- T , the number of probes (total number of buckets probed across all hash tables).

Usually, it is a good idea to choose L first based on the available memory. Then, we have a trade-off between K and T : the larger K is, the more probes we need to achieve a given probability of success, and vice versa. The best way to choose K and T is usually the following parameter search: Try increasing values of K , and for each value of K , find the right number of probes T so that we get the desired accuracy on a set of sample queries. Varying the parameter T does not require rebuilding the hash table (as opposed to K and L). Moreover, we can search over T using a binary search. Usually, this means that we can find the optimal parameter setting fairly quickly.

From FALCONN documentation



ilyaraz

152 commits **165,070 ++** **63,535 --**

#1



ludwigschmidt

90 commits **278,244 ++** **16,349 --**

#2

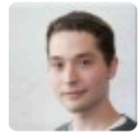
Choosing parameters

All in all, we have three parameters:

- K , the number of hash functions (space partitions) per hash table.
- L , the number of hash tables (each of which has K hash functions).
- T , the number of probes (total number of buckets probed across all hash tables).

Usually, it is a good idea to choose L first based on the available memory. Then, we have a trade-off between K and T : the larger K is, the more probes we need to achieve a given probability of success, and vice versa. The best way to choose K and T is usually the following parameter search: Try increasing values of K , and for each value of K , find the right number of probes T so that we get the desired accuracy on a set of sample queries. Varying the parameter T does not require rebuilding the hash table (as opposed to K and L). Moreover, we can search over T using a binary search. Usually, this means that we can find the optimal parameter setting fairly quickly.

From FALCONN documentation



ilyaraz

152 commits 165,070 ++ 63,535 --

#1



ludwigschmidt

90 commits 278,244 ++ 16,349 --

#2

Choosing parameters

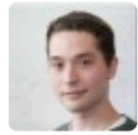
All in all, we have three parameters:

- K , the number of hash functions (space partitions) per hash table.
- L , the number of hash tables (each of which has K hash functions).
- T , the number of probes (total number of buckets probed across all hash tables).

Usually, it is a good idea to choose L first based on the available memory. Then, we have a trade-off between K and T : the larger K is, the more probes we need to achieve a given probability of success, and vice versa. The best way to choose K and T is usually the following parameter search: Try increasing values of K , and for each value of K , find the right number of probes T so that we get the desired accuracy on a set of sample queries. Varying the parameter T does not require rebuilding the hash table (as opposed to K and L). Moreover, we can search over T using a binary search. Usually, this means that we can find the optimal parameter setting fairly quickly.

Achieving a given recall using LSH methods is an empirical task

From FALCONN documentation



ilyaraz

152 commits 165,070 ++ 63,535 --

#1



ludwigschmidt

90 commits 278,244 ++ 16,349 --

#2

Choosing parameters

All in all, we have three parameters:

- K , the number of hash functions (space partitions) per hash table.
- L , the number of hash tables (each of which has K hash functions).
- T , the number of probes (total number of buckets probed across all hash tables).

Usually, it is a good idea to choose L first based on the available memory. Then, we have a trade-off between K and T : the larger K is, the more probes we need to achieve a given probability of success, and vice versa. Try increasing values of K , the desired accuracy on a rebuilding the hash table search. Usually, this means that we can find the optimal parameter setting fairly quickly.

Achieving a given recall using LSH methods is an empirical task

In theory: "Only polylog n different parameter choices"

Adaptive stopping

Idea:

- Suppose that after searching t hash tables the nearest neighbor retrieved so far is x^* .
- Let p^* denote the probability that x^* is retrieved in a hash table.
- If $t > \ln(1 / \delta) / p^*$ then stop and return x^* .

Adaptive stopping

Idea:

- Suppose that after searching t hash tables the nearest neighbor retrieved so far is x^* .
- Let p^* denote the probability that x^* is retrieved in a hash table.
- If $t > \ln(1/\delta)/p^*$ then stop and return x^* .

Yields expected recall $1-\delta$ if the nearest neighbor has the largest retrieval probability

Adaptive stopping

Idea:

- Suppose that after searching t hash tables the nearest neighbor retrieved so far is x^* .
- Let p^* denote the probability that x^* is retrieved in a hash table.
- If $t > \ln(1/\delta)/p^*$ then stop and return x^* .

Yields expected recall $1-\delta$ if the nearest neighbor has the largest retrieval probability

Issue: Need way of computing p^* .

Abstract view

• q

Abstract view

A wide, bright blue diagonal band runs from the top-left towards the bottom-right of the slide. In the center of this band, there is a small orange dot followed by the lowercase letter 'q' in a matching orange color.

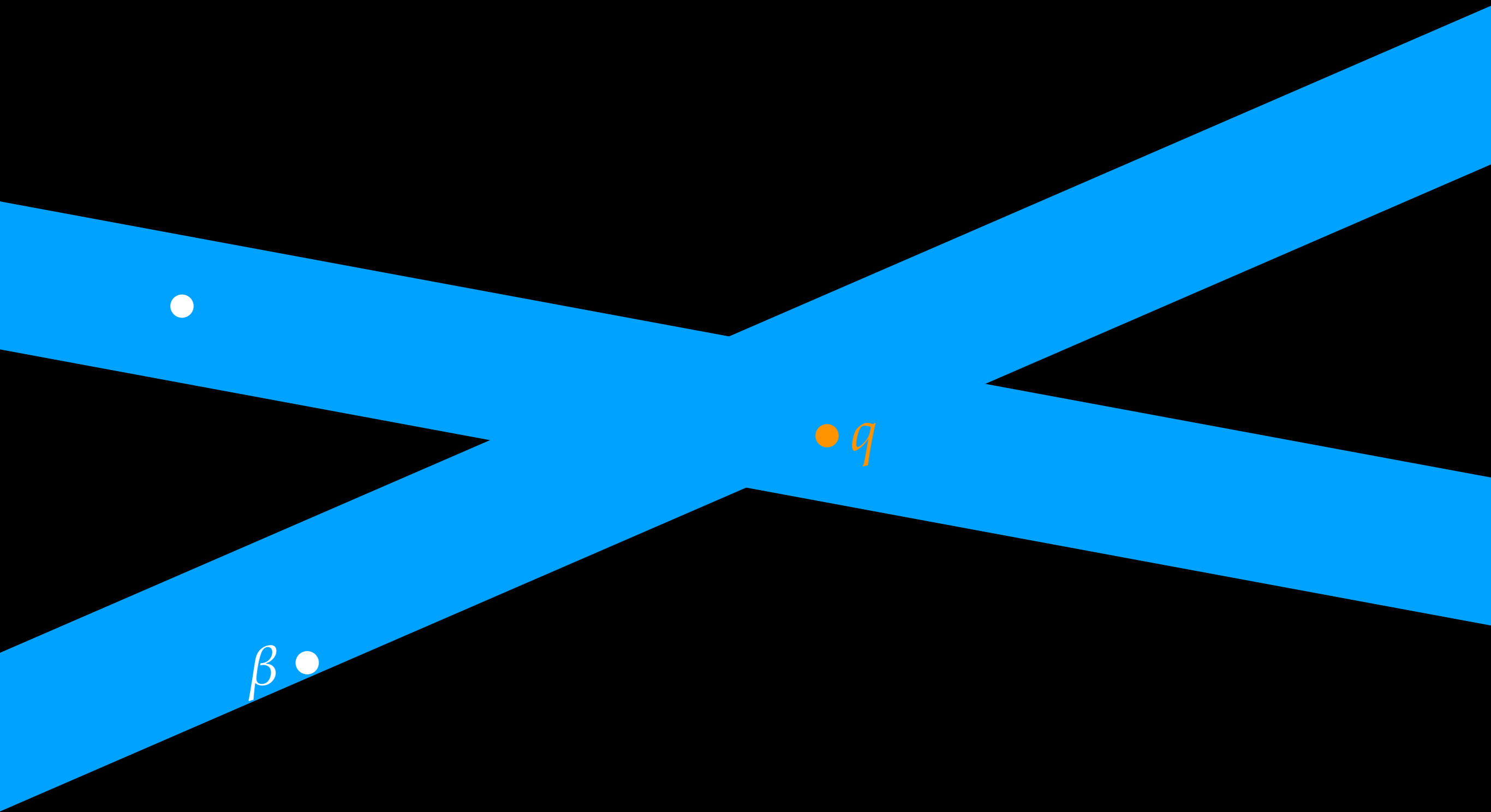
• *q*

Abstract view

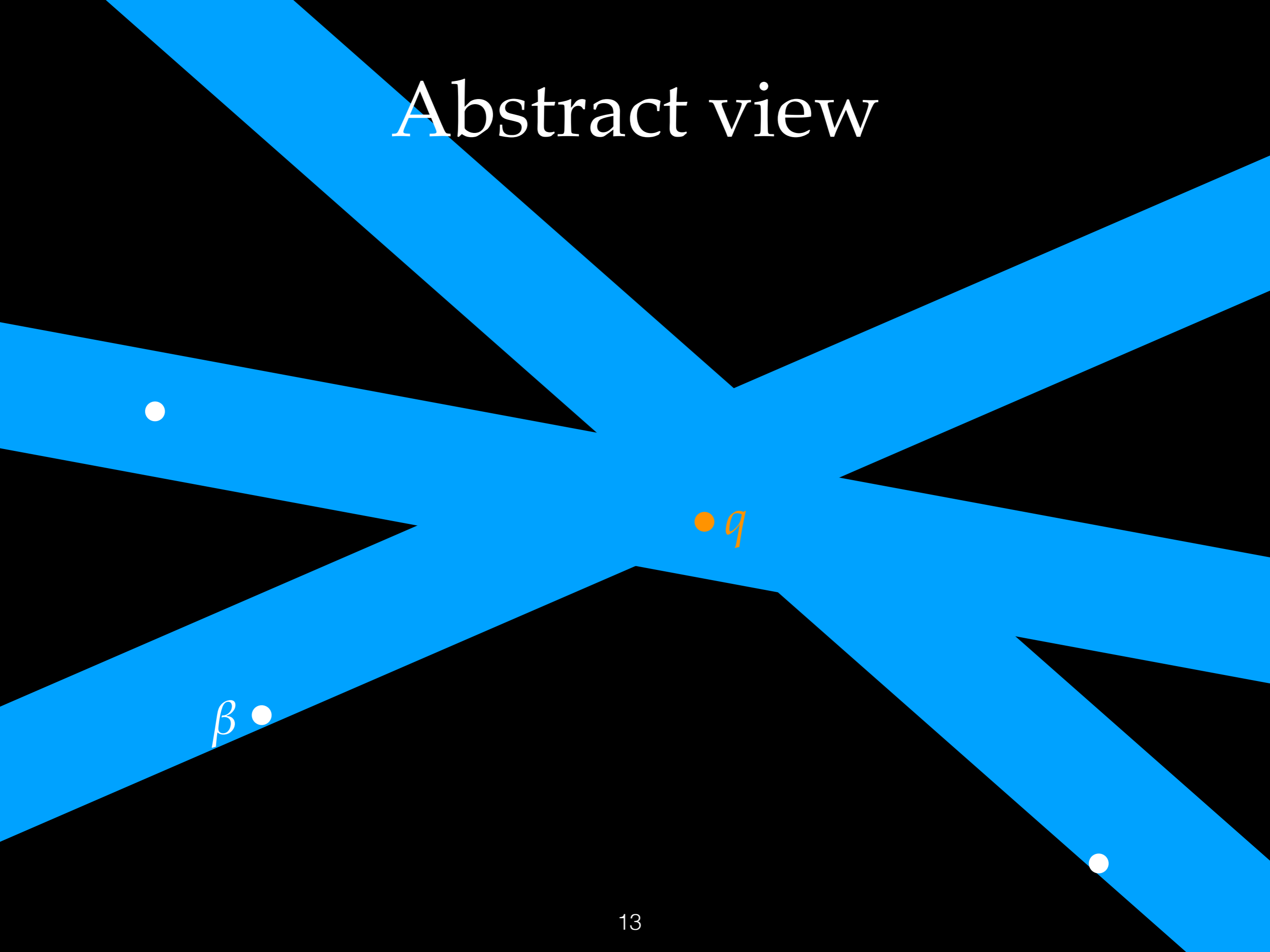
$\beta \bullet$

$\bullet q$

Abstract view



Abstract view

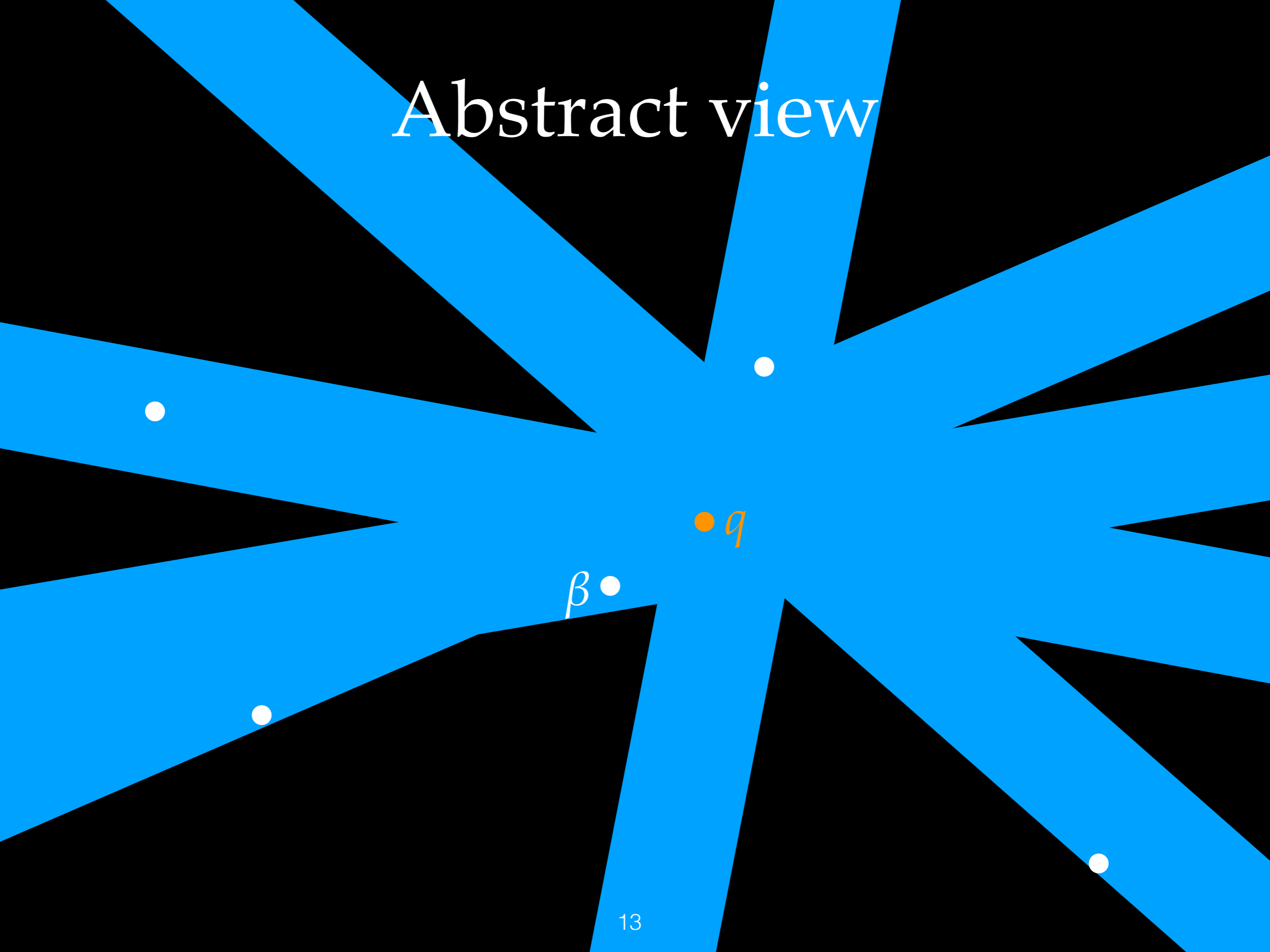


Abstract view

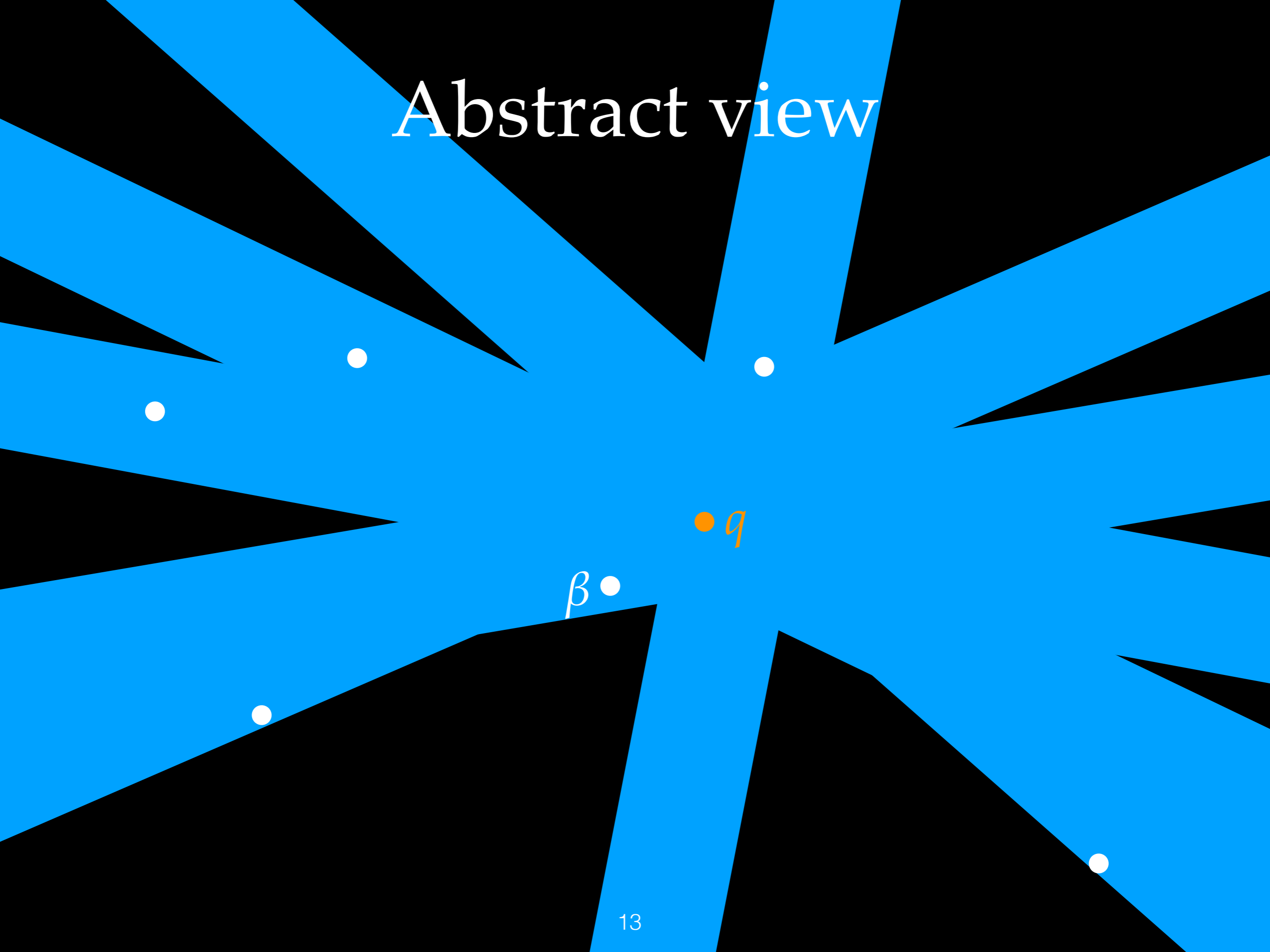
$\beta \bullet$

$\bullet q$

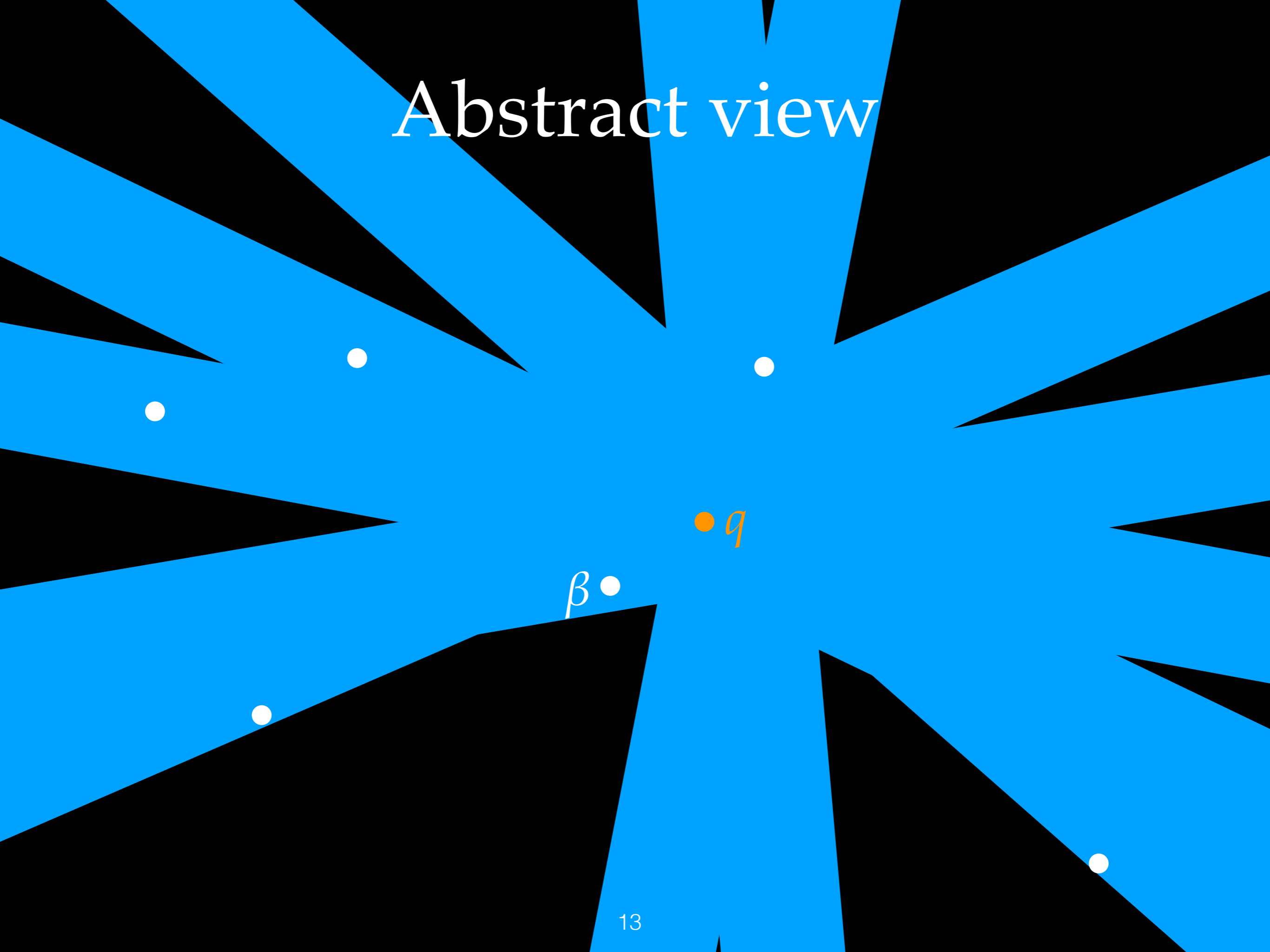
Abstract view



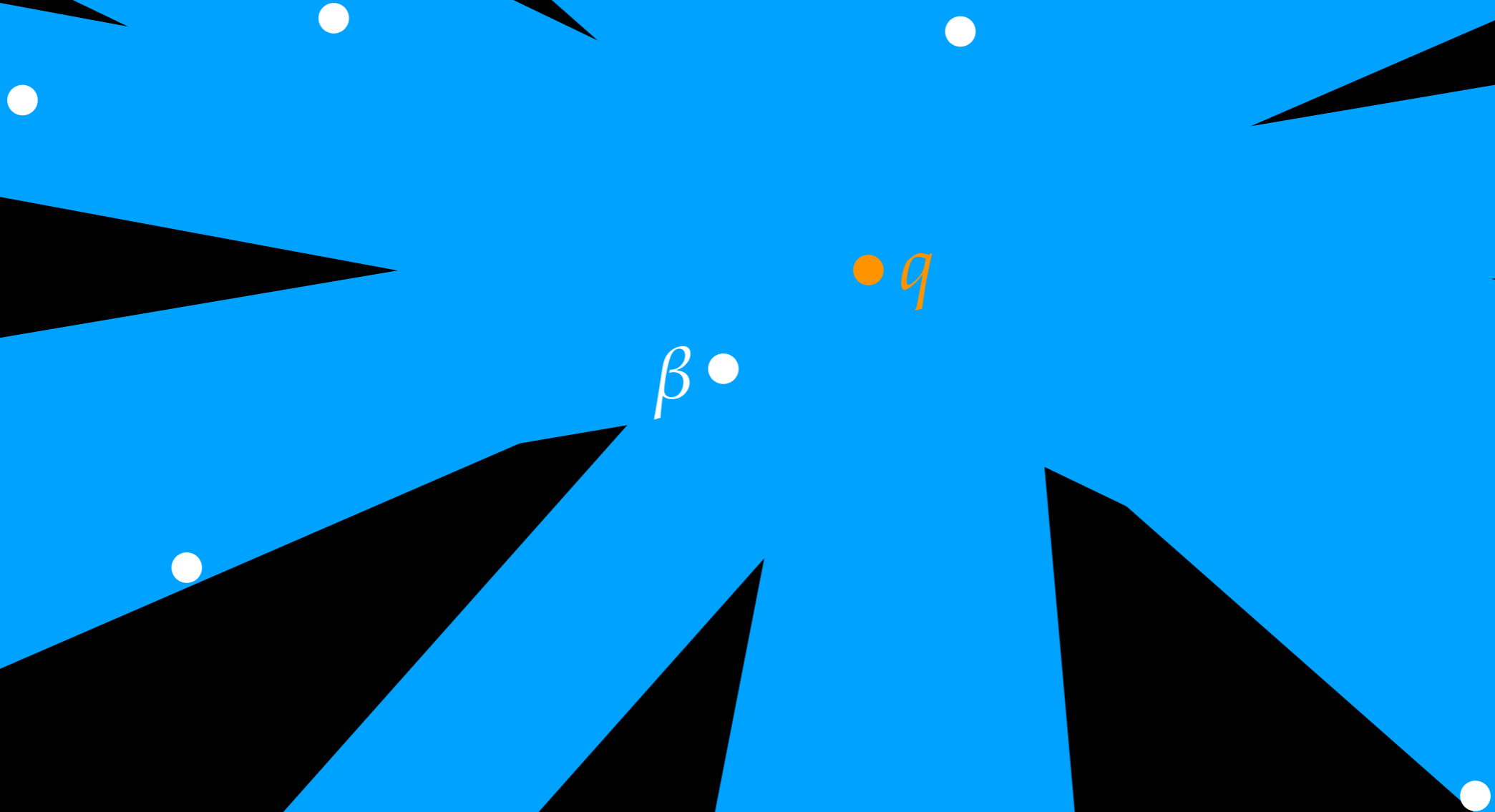
Abstract view



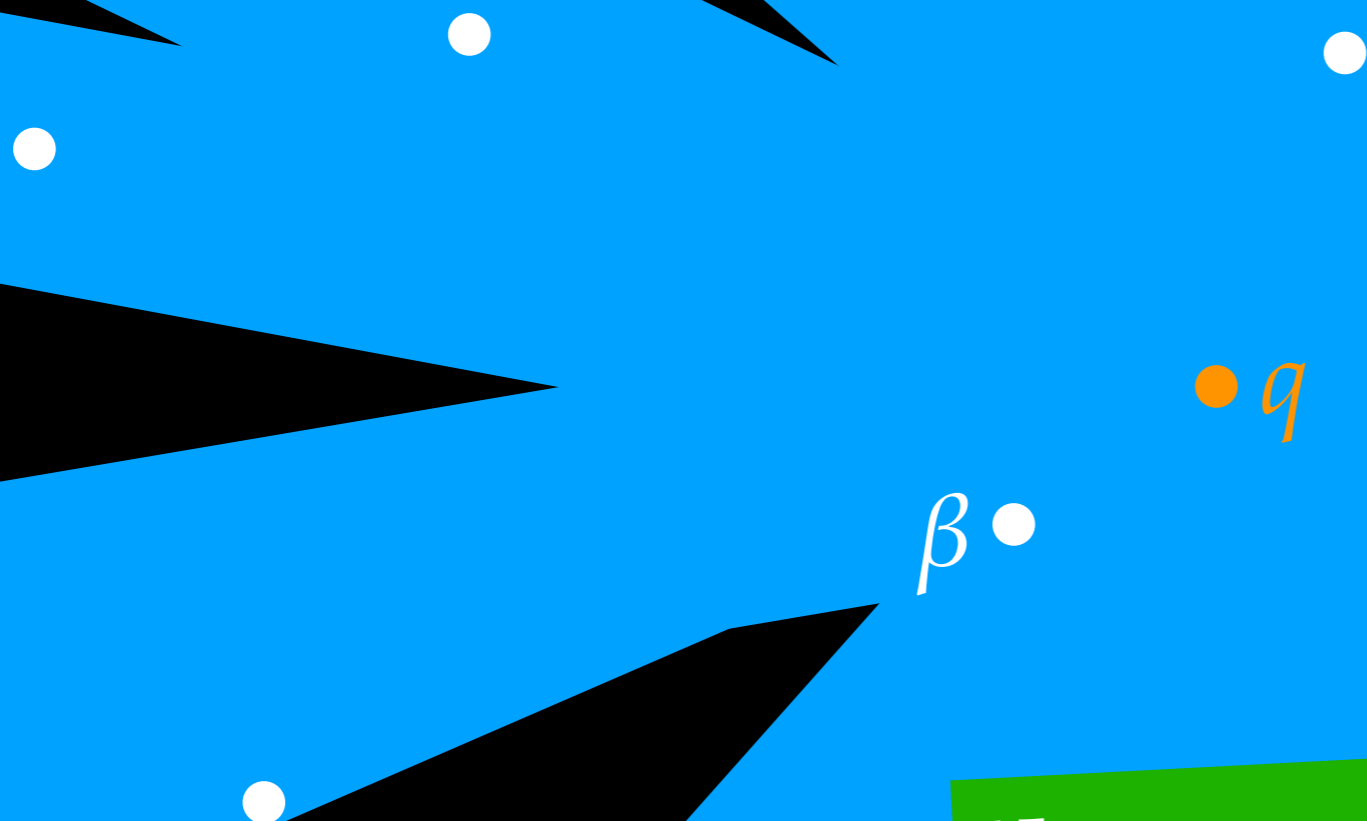
Abstract view



Abstract view

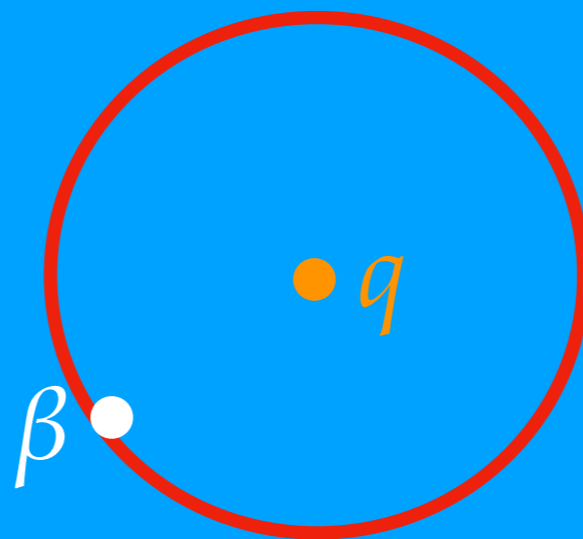


Abstract view



How can we know if β is likely to be the nearest neighbor?

Abstract view



How can we know if β is likely to be the nearest neighbor?

Confirmation sampling

- Notation: Repetition i produces *one* nearest neighbor candidate y_i , independently for each i .
- Let β_i be the nearest neighbor of q in $Y_i = \{y_1, \dots, y_i\}$.

Confirmation sampling

- Notation: Repetition i produces *one* nearest neighbor candidate y_i , independently for each i .
- Let β_i be the nearest neighbor of q in $Y_i = \{y_1, \dots, y_i\}$.
- Stop after repetition i if β_i appears $t+1$ times in Y_i .

Confirmation sampling

- Notation: Repetition i produces *one* nearest neighbor candidate y_i , independently for each i .
- Let β_i be the nearest neighbor of q in $Y_i = \{y_1, \dots, y_i\}$.
- Stop after repetition i if β_i appears $t+1$ times in Y_i .

NN candidate
is "confirmed"
 t times

Confirmation sampling

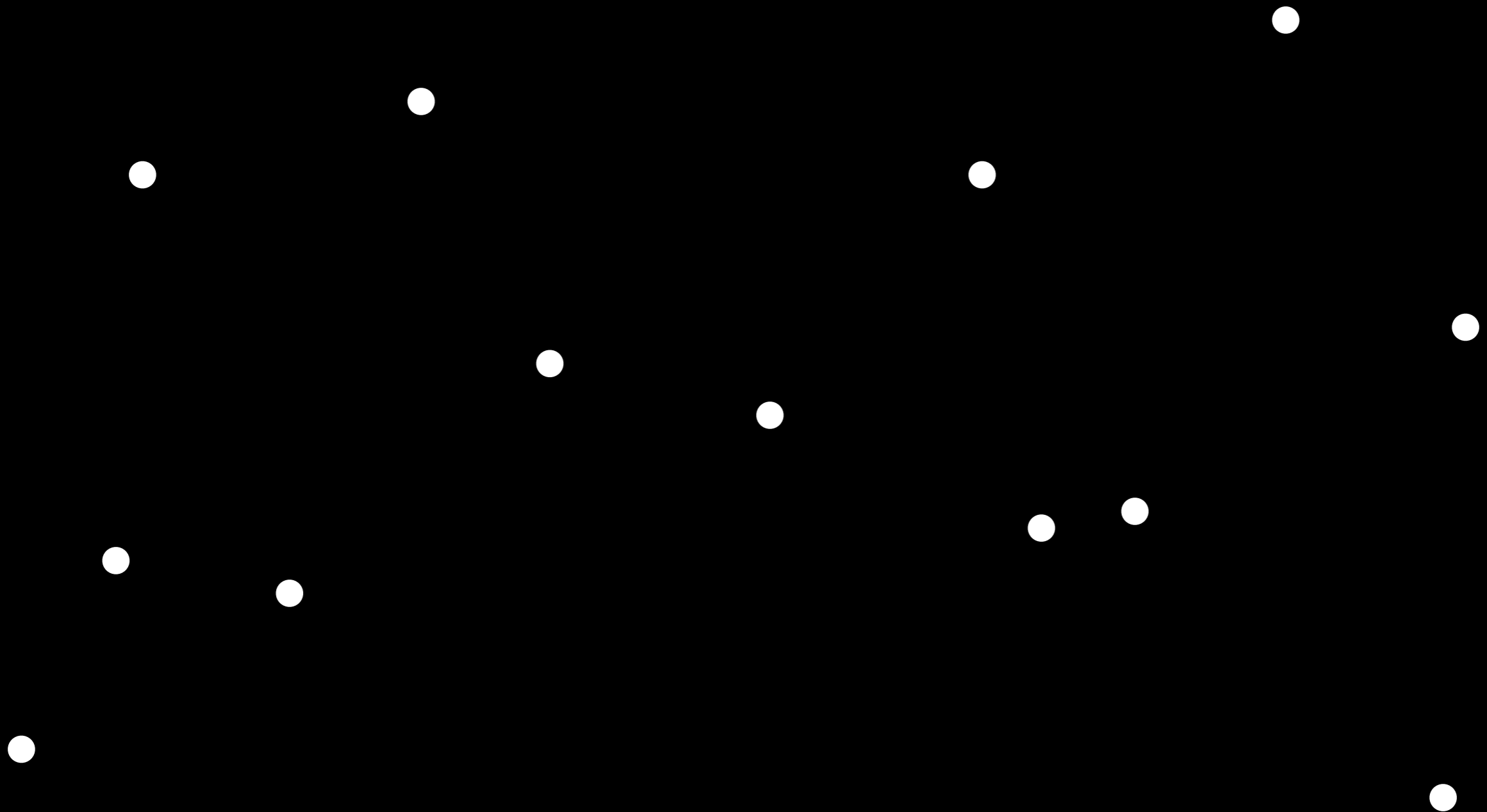
- Notation: Repetition i produces *one* nearest neighbor candidate y_i , independently for each i .
- Let β_i be the nearest neighbor of q in $Y_i = \{y_1, \dots, y_i\}$.
- Stop after repetition i if β_i appears $t+1$ times in Y_i .

Claim: If nearest neighbor is most likely to be a candidate, the probability of returning a different point is at most 2^{-t} .

NN candidate is "confirmed" t times

Confirmation sampling

• $t = 1$



Confirmation sampling

• $t = 1$

• q

Confirmation sampling

• $t = 1$

• y

• q

Confirmation sampling

• $t = 1$

$\beta \bullet \gamma$

• q

Confirmation sampling

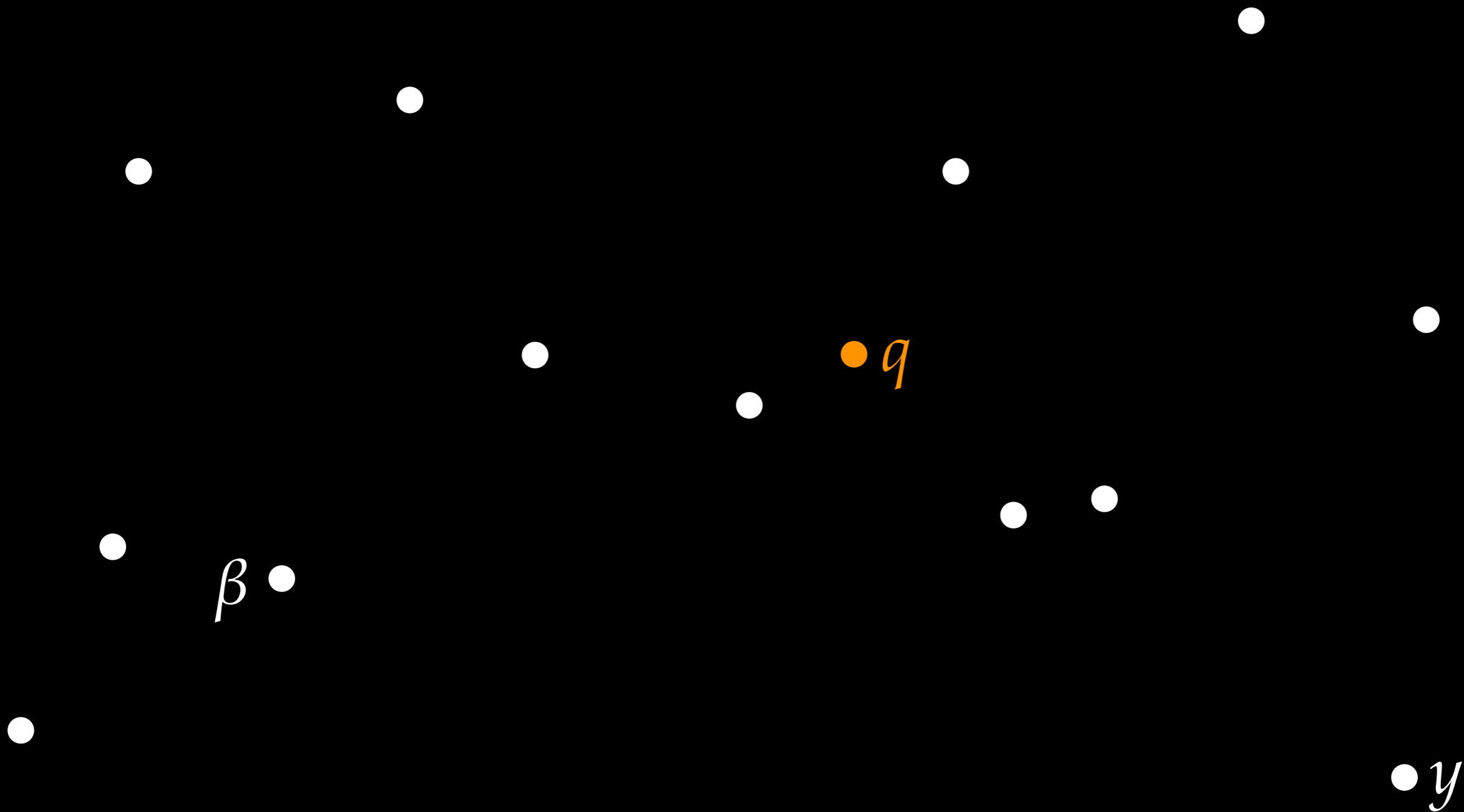
• $t = 1$

• q

• $\beta \cdot \gamma$

Confirmation sampling

• $t = 1$



Confirmation sampling

• $t = 1$

• $\beta \cdot \gamma$

• q

Confirmation sampling

• $t = 1$

• β

• q

• y

Confirmation sampling

• $t = 1$

• $\beta \cdot y$ • q

Confirmation sampling

• $t = 1$

• y

• β

• q

Confirmation sampling

• $t = 1$

• $\beta \cdot y$ • q

Confirmation:
Return β

Why confirmation sampling works

- Suppose x_1 is the nearest neighbor.
- For $i > 1$, if $y_i = \beta_{i-1}$ but $y_i \neq x_1$ we say step i is a *false confirmation*.

Why confirmation sampling works

- Suppose x_1 is the nearest neighbor.
- For $i > 1$, if $y_i = \beta_{i-1}$ but $y_i \neq x_1$ we say step i is a *false confirmation*.
- Consider the first t steps where we either sample x_1 or produce a false confirmation:
 - If sampling x_1 is more likely than sampling $\beta_{i-1} \neq x_1$, the probability of false confirmation is at most $1/2$ in each step.
 - Probability of t false confirmations is at most 2^{-t} .

Abstract confirmation sampling

Algorithm 1: CONFIRMATION SAMPLING(\mathcal{Q}, t, \prec)

```
1  $\beta \leftarrow \infty, \text{count} \leftarrow 0$ 
2 while  $\text{count} < t$  do
3   sample  $X \sim \mathcal{Q}$ 
4   if  $X = \beta$  then
5     |  $\text{count} \leftarrow \text{count} + 1$ 
6   else if  $X \prec \beta$  then
7     |  $\beta \leftarrow X$ 
8     |  $\text{count} \leftarrow 0$ 
9 return  $\beta$ 
```

Algorithm 1: CONFIRMATIONSAMPLING(\mathcal{Q}, t, \prec)

```
1  $\beta \leftarrow \infty$ , count  $\leftarrow 0$ 
2 while count  $< t$  do
3   sample  $X \sim \mathcal{Q}$ 
4   if  $X = \beta$  then
5     count  $\leftarrow$  count + 1
6   else if  $X \prec \beta$  then
7      $\beta \leftarrow X$ 
8     count  $\leftarrow 0$ 
9 return  $\beta$ 
```

Theorem 3. Let \mathcal{Q} denote a probability distribution with finite support S . For $x_1 = \min(S)$ and $X \sim \mathcal{Q}$ let $p_1 = \Pr[X = x_1]$ and let $p_2 = \max\{\Pr[X = x] \mid x \in S \setminus \{x_1\}\}$ be the largest sampling probability among elements of S other than x_1 . Then:

$$\Pr[\text{CONFIRMATIONSAMPLING}(\mathcal{Q}, t) \neq x_1] \leq (1 - p_1) \left(\frac{p_2}{p_1 + p_2} \right)^t$$

The expected number of samples made by CONFIRMATIONSAMPLING is bounded by $(t + 1)/p_1$.

Application to nearest neighbor

Theorem 1. *Suppose there is a sequence of independent, randomized data structures $\mathcal{D}_1, \mathcal{D}_2, \dots$, such that on query q , \mathcal{D}_i returns the nearest neighbor of q in P with probability at least p_q and each other point in P with probability at most p_q . Let $\delta > 0$ be given. There is an algorithm that depends on δ but not on p_q that on input q queries data structures $\mathcal{D}_1, \dots, \mathcal{D}_{j_q}$, performs j_q distance computations, where $\mathbb{E}[j_q] = O(\ln(1/\delta)/p_1)$, and returns the nearest neighbor of q with probability at least $1 - \delta$.*

Application to nearest neighbor

Theorem 1. *Suppose there is a sequence of independent, randomized data structures $\mathcal{D}_1, \mathcal{D}_2, \dots$, such that on query q , \mathcal{D}_i returns the nearest neighbor of q in P with probability at least p_q and each other point in P with probability at most p_q . Let $\delta > 0$ be given. There is an algorithm that depends on δ but not on p_q that on input q queries data structures $\mathcal{D}_1, \dots, \mathcal{D}_{j_q}$, performs j_q distance computations, where $\mathbb{E}[j_q] = O(\ln(1/\delta)/p_1)$, and returns the nearest neighbor of q with probability at least $1 - \delta$.*

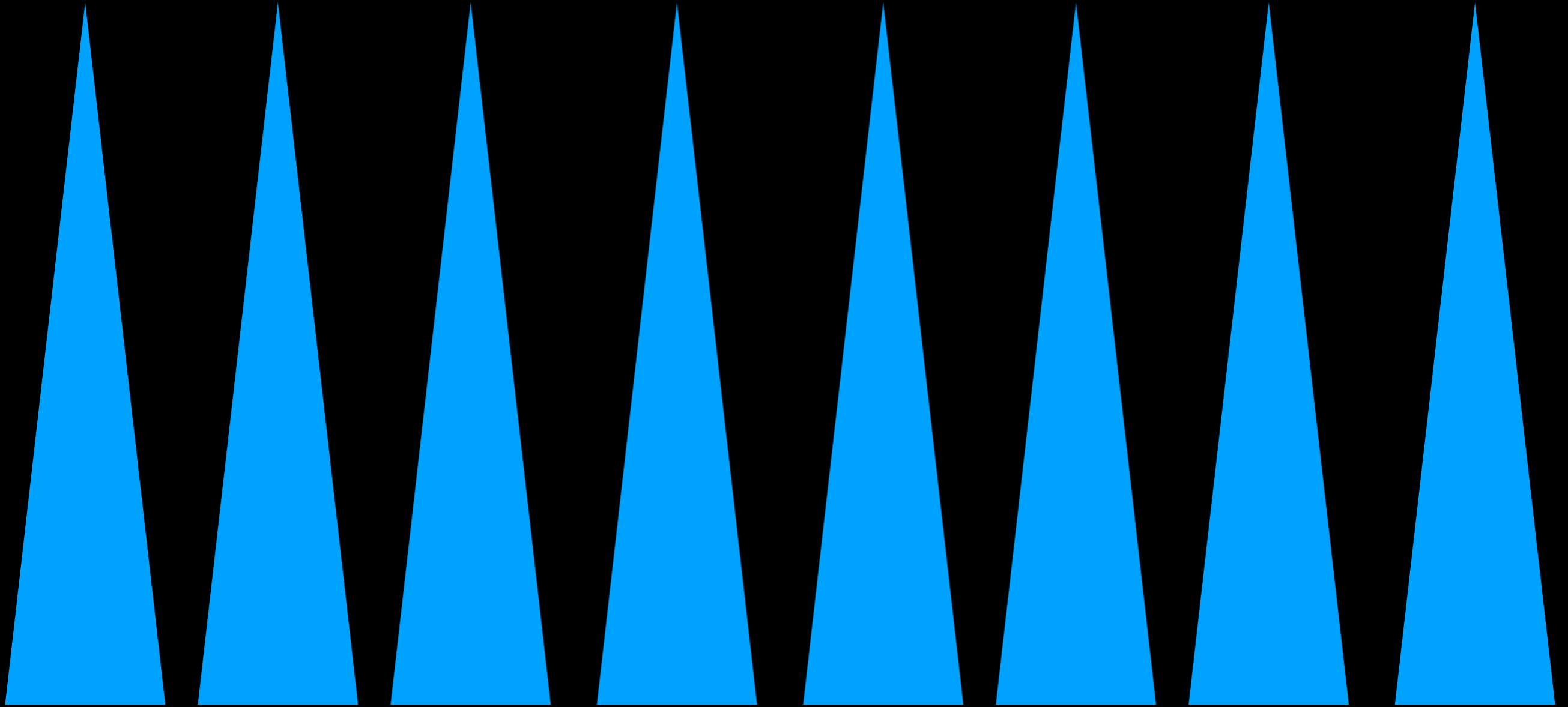
- Finding the nearest neighbor quickly boils down to minimizing the product of the expected time for querying \mathcal{D}_i and $1/p_1$.

Application to nearest neighbor

Theorem 1. *Suppose there is a sequence of independent, randomized data structures $\mathcal{D}_1, \mathcal{D}_2, \dots$, such that on query q , \mathcal{D}_i returns the nearest neighbor of q in P with probability at least p_q and each other point in P with probability at most p_q . Let $\delta > 0$ be given. There is an algorithm that depends on δ but not on p_q that on input q queries data structures $\mathcal{D}_1, \dots, \mathcal{D}_{j_q}$, performs j_q distance computations, where $\mathbb{E}[j_q] = O(\ln(1/\delta)/p_1)$, and returns the nearest neighbor of q with probability at least $1 - \delta$.*

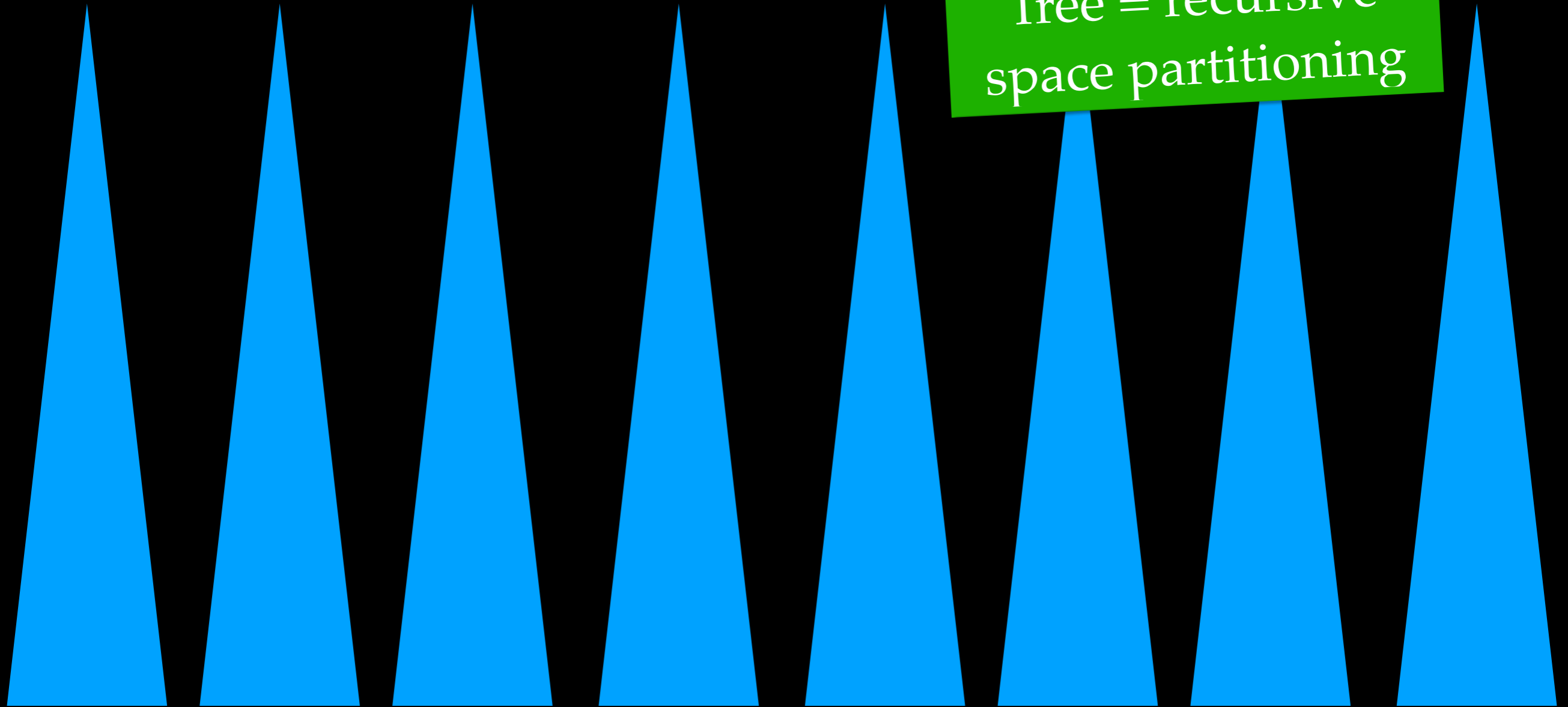
- Finding the nearest neighbor quickly boils down to minimizing the product of the expected time for querying \mathcal{D}_i and $1/p_1$.
- **Question:** Can this be done optimally without knowledge of the distance to the nearest neighbor?

Partial answer for LSH forest



Partial answer for LSH forest

Tree = recursive
space partitioning



Partial answer for LSH forest

Tree = recursive
space partitioning

Standard query algorithm: Search level where number of points in q 's space partition is $O(1)$ on average.

Partial answer for LSH forest

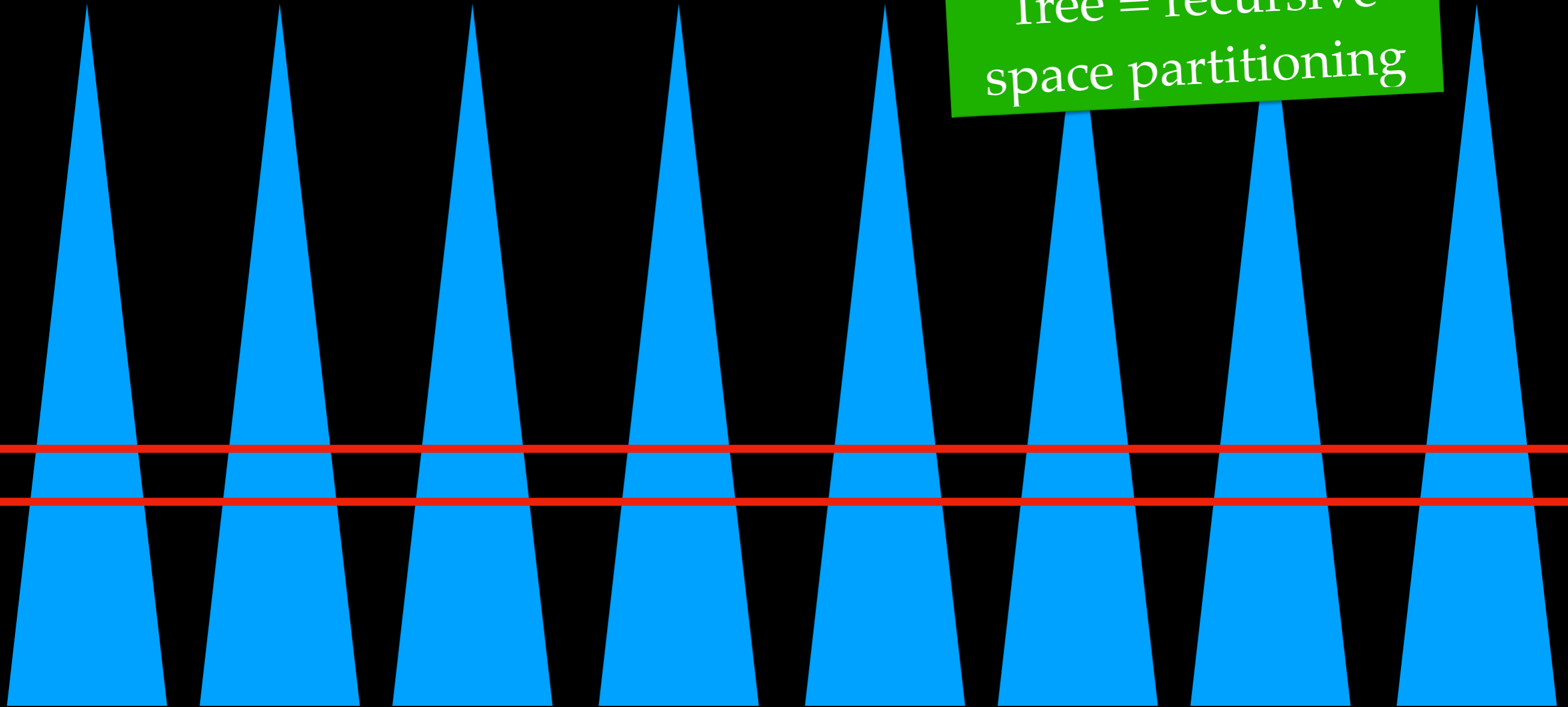
Tree = recursive
space partitioning

Standard query algorithm: Search level where number of points in q 's space partition is $O(1)$ on average.

Works if number of trees is high enough,
depending on distance to nearest neighbor.

Partial answer for LSH forest

Tree = recursive
space partitioning



Partial answer for LSH forest

Tree = recursive
space partitioning

Modified query algorithm: Search higher levels until
confirmation search says we found the nearest neighbor.

Partial answer for LSH forest

Tree = recursive
space partitioning

Modified query algorithm: Search higher levels until
confirmation search says we found the nearest neighbor.

Partial answer for LSH forest

Tree = recursive
space partitioning

Modified query algorithm: Search higher levels until
confirmation search says we found the nearest neighbor.

Partial answer for LSH forest

Tree = recursive
space partitioning

Modified query algorithm: Search higher levels until
confirmation search says we found the nearest neighbor.

Also consider partially searching a level.

Partial answer for LSH forest

Tree = recursive
space partitioning

Competitive with best way of
searching LSH forest for given query

Modified query algorithm: Search higher levels until
confirmation search says we found the nearest neighbor.

Also consider partially searching a level.

Open question

- Is it possible to achieve space and time that is $O(1)$ -competitive with the best LSH scheme, adapted to the *query and to the data distribution*, for a given expected recall?