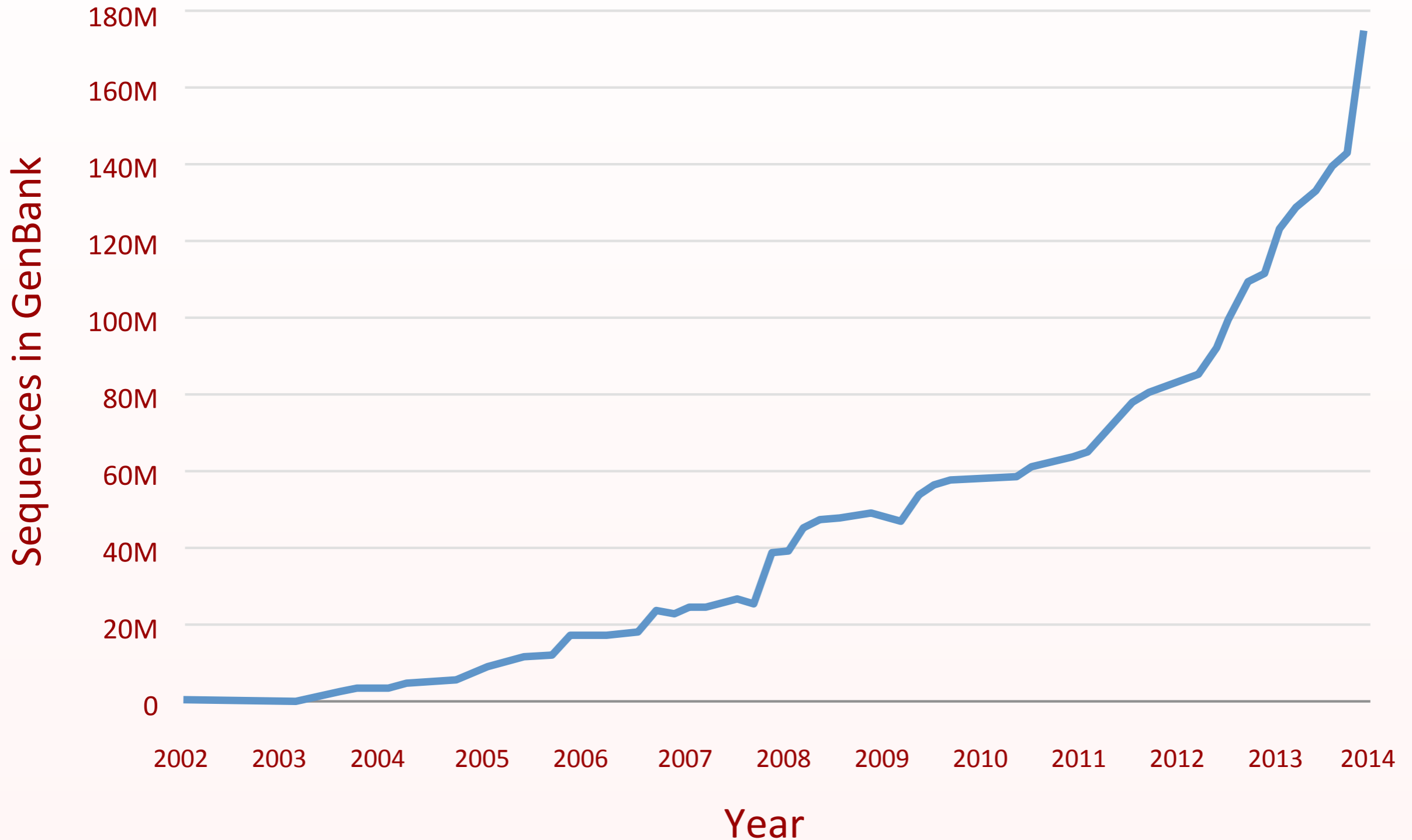

Genomic Compression: Storage, Transmission, and Analytics


Noah M. Daniels
Bonnie Berger

Genomic data are growing exponentially



A bigger cloud?

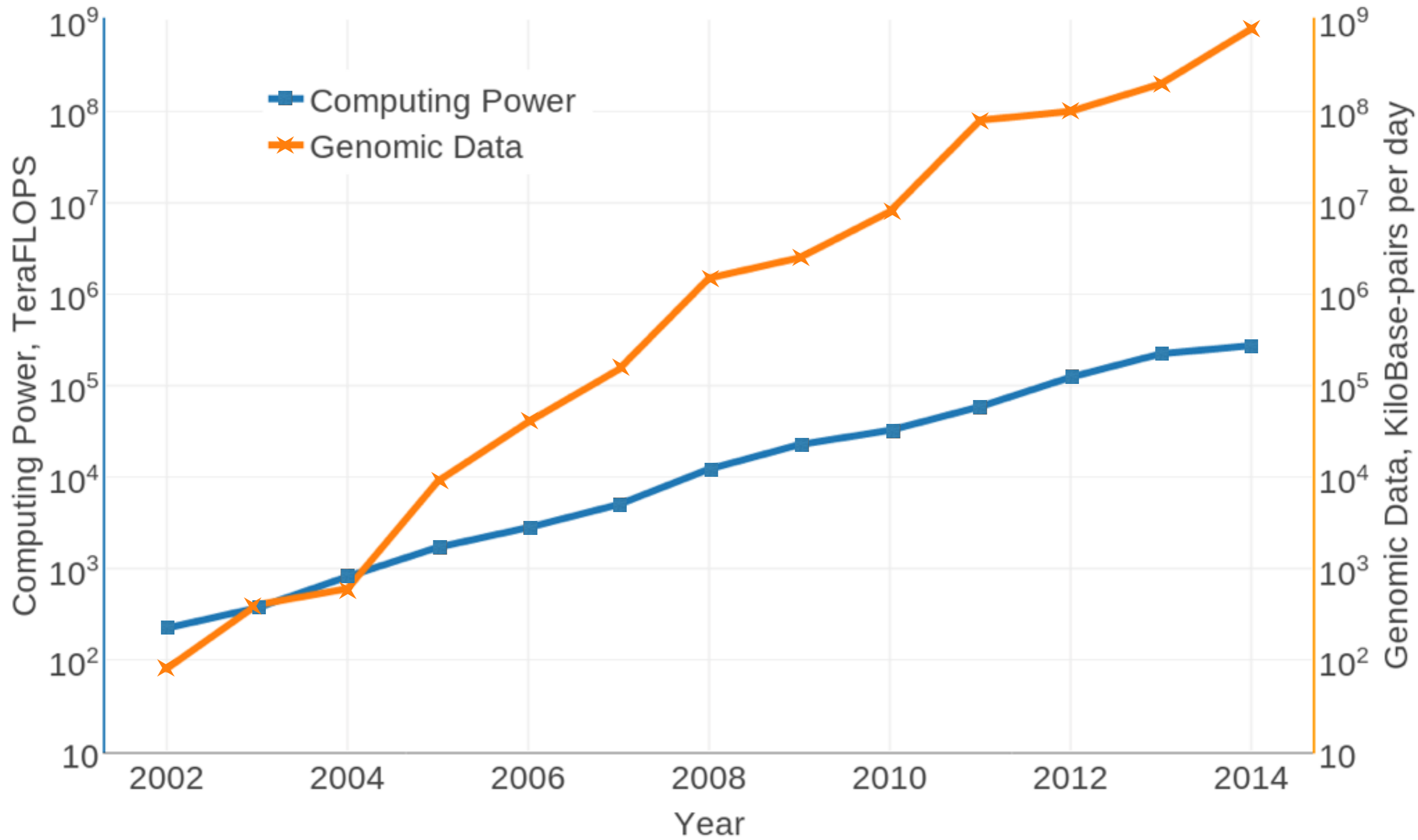




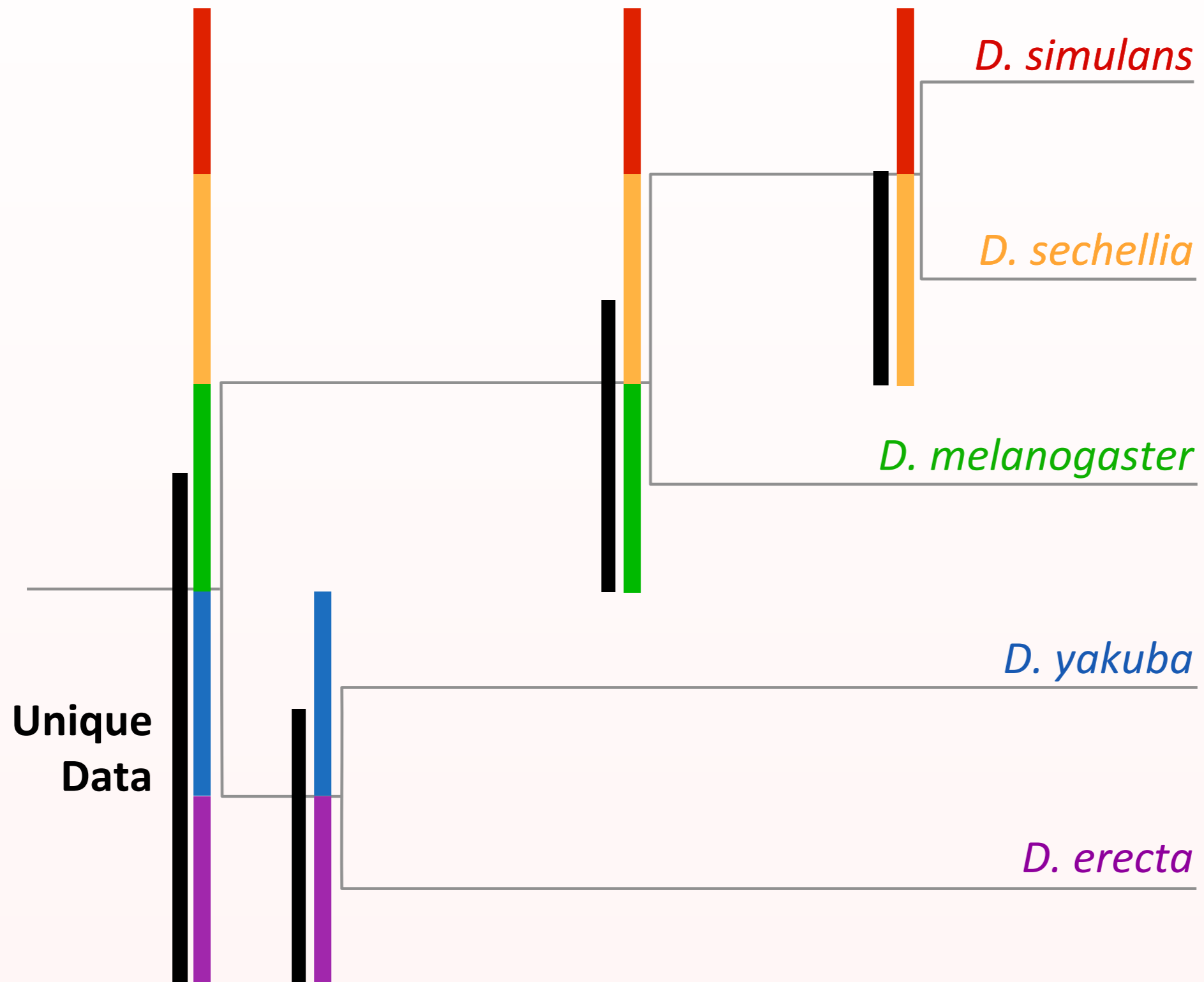
A bigger cloud?

Still constrained by Moore

Genomic data are growing exponentially



Redundancy in data



Outline

Compression background

Storage and transmission

Analysis

Limits to compression

Source coding theorem [Shannon, 1948]

N i.i.d. random variables from source X
each with entropy $H(X)$

cannot be compressed into fewer than $NH(X)$
bits without loss of information

$$H(X) = - \sum_i P(x_i) \log_b P(x_i)$$

What can we do?

Approach Shannon limit as efficiently as possible (fast, lossless compression)

Be willing to throw away information (*lossy* compression)

Approaching Shannon limit

Huffman coding [1952]

Arithmetic coding [Rissanen 1976]

Lempel-Ziv [1977]

Burrows-Wheeler Transform [1994]

Huffman coding

Huffman coding

Variable-length prefix code

Huffman coding

Variable-length prefix code

More frequent symbols take fewer bits

Huffman coding

Variable-length prefix code

More frequent symbols take fewer bits

DNA (A,C,T,G)

Huffman coding

Variable-length prefix code

More frequent symbols take fewer bits

DNA (A,C,T,G)

2 bits/symbol: 00, 01, 10, 11

Huffman coding

Variable-length prefix code

More frequent symbols take fewer bits

DNA (A,C,T,G)

2 bits/symbol: 00, 01, 10, 11

Suppose AT bias

Huffman coding

Variable-length prefix code

More frequent symbols take fewer bits

DNA (A,C,T,G)

2 bits/symbol: 00, 01, 10, 11

Suppose AT bias

P. falciparum 20% GC

A

T

C

G

A 0.4

T 0.4

C 0.1

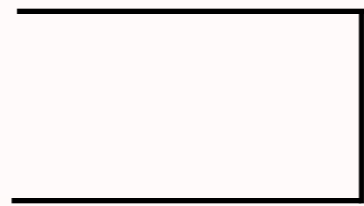
G 0.1

A 0.4

T 0.4

C 0.1

G 0.1



A 0.4

T 0.4

C 0.1

G 0.1

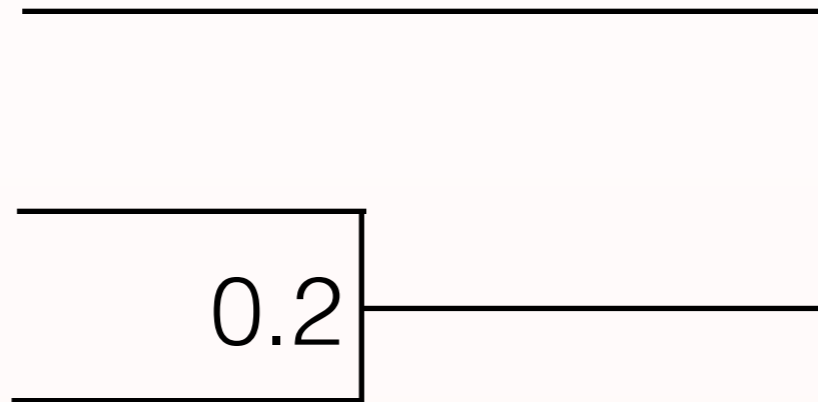


A 0.4

T 0.4

C 0.1

G 0.1



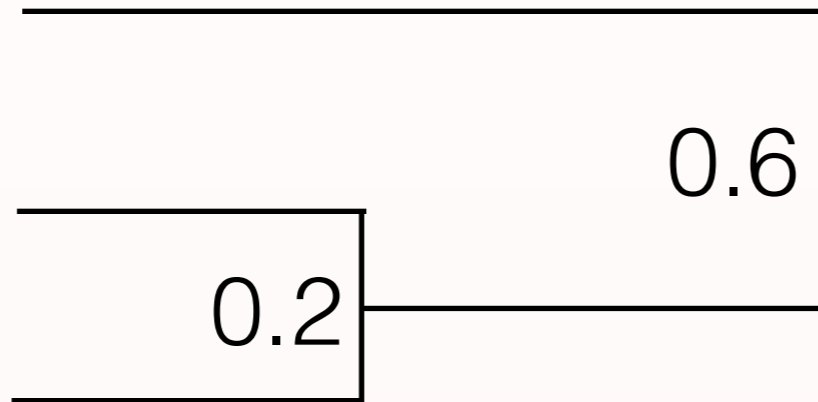
0.2

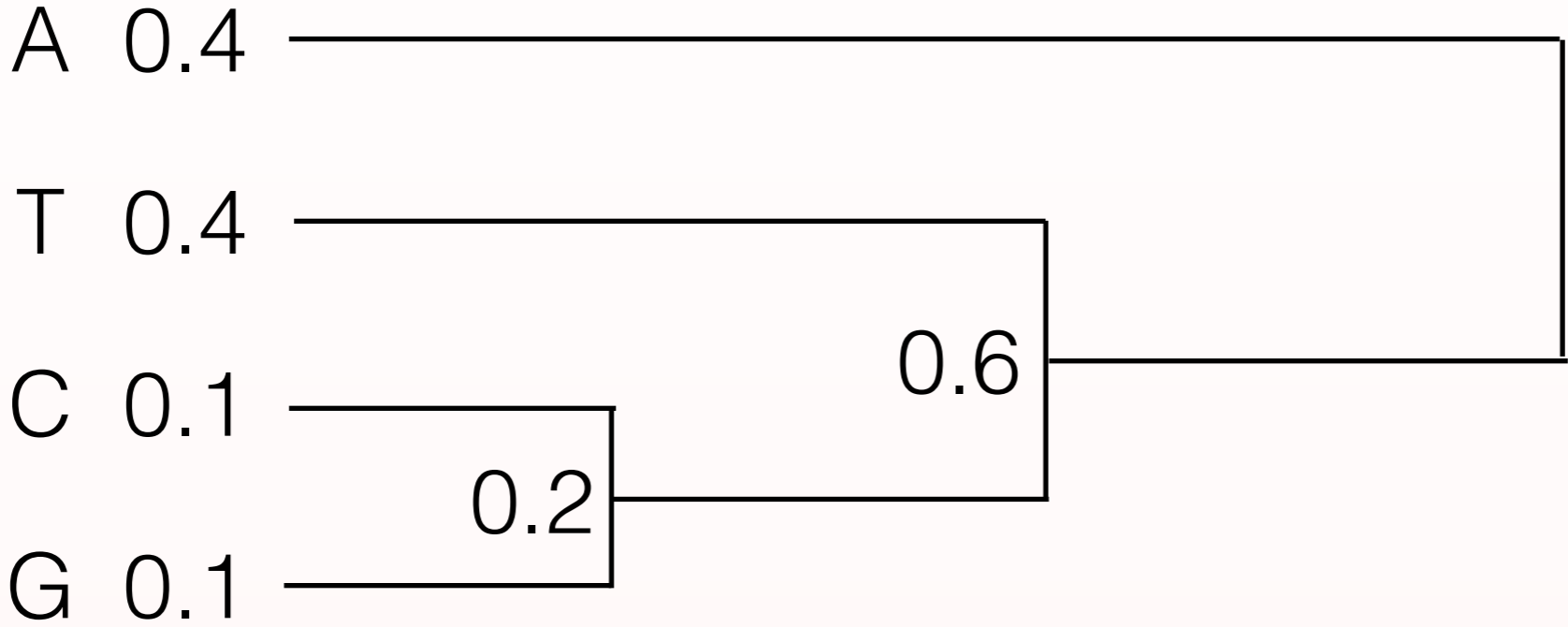
A 0.4

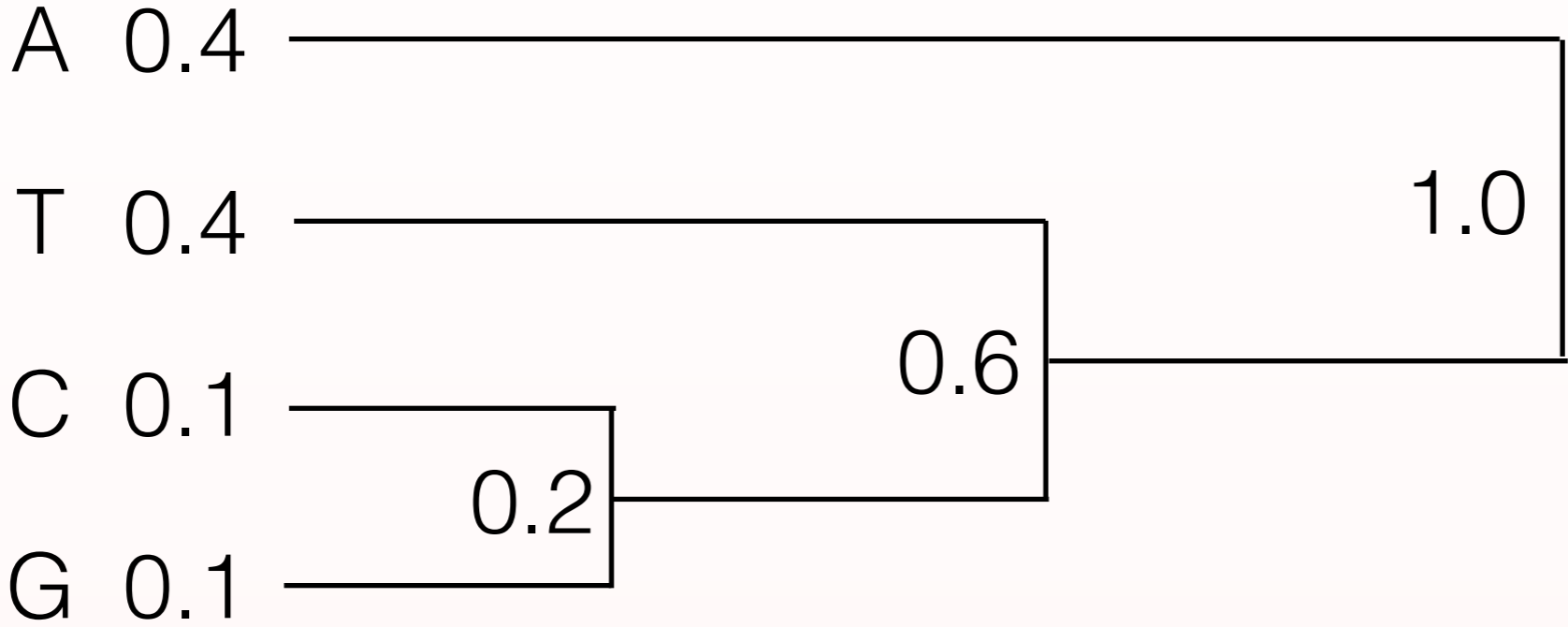
T 0.4

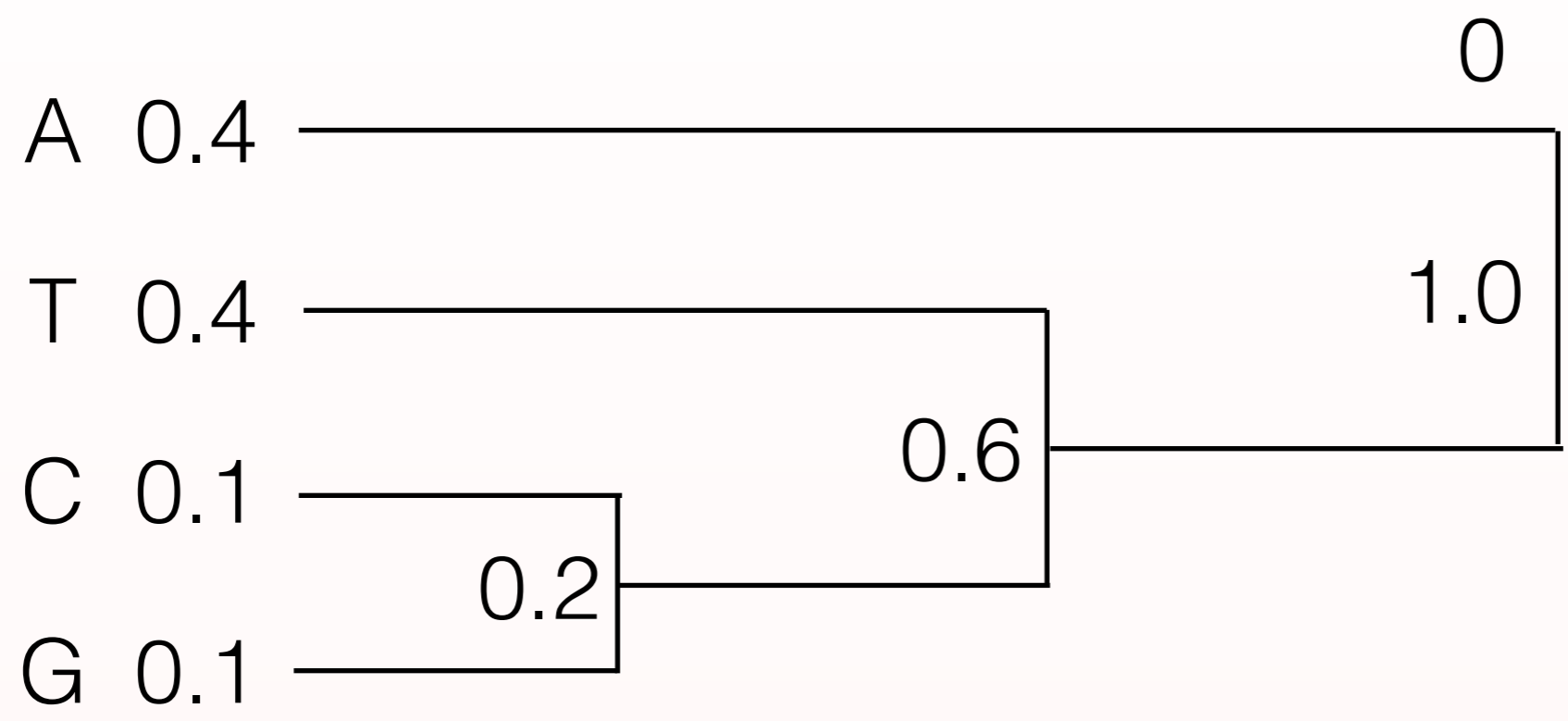
C 0.1

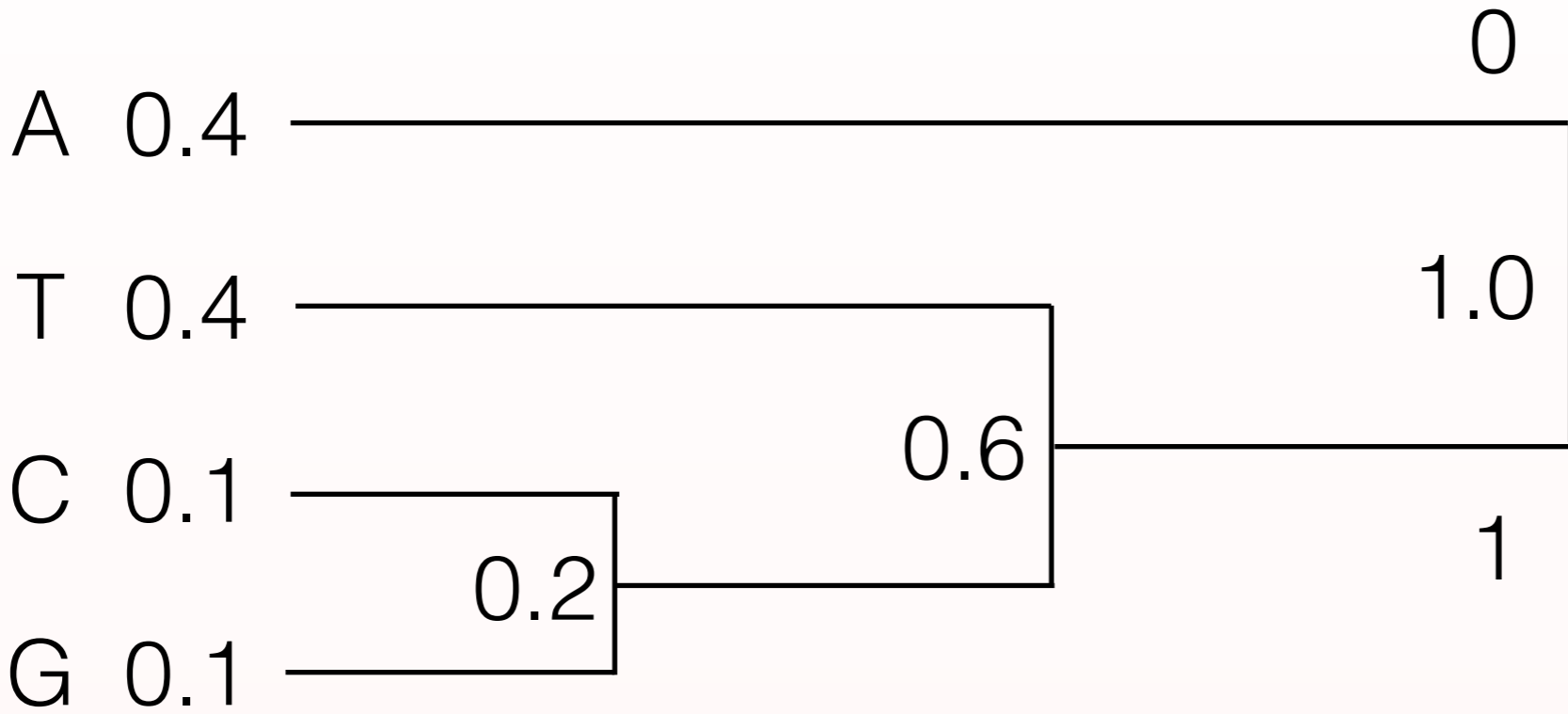
G 0.1

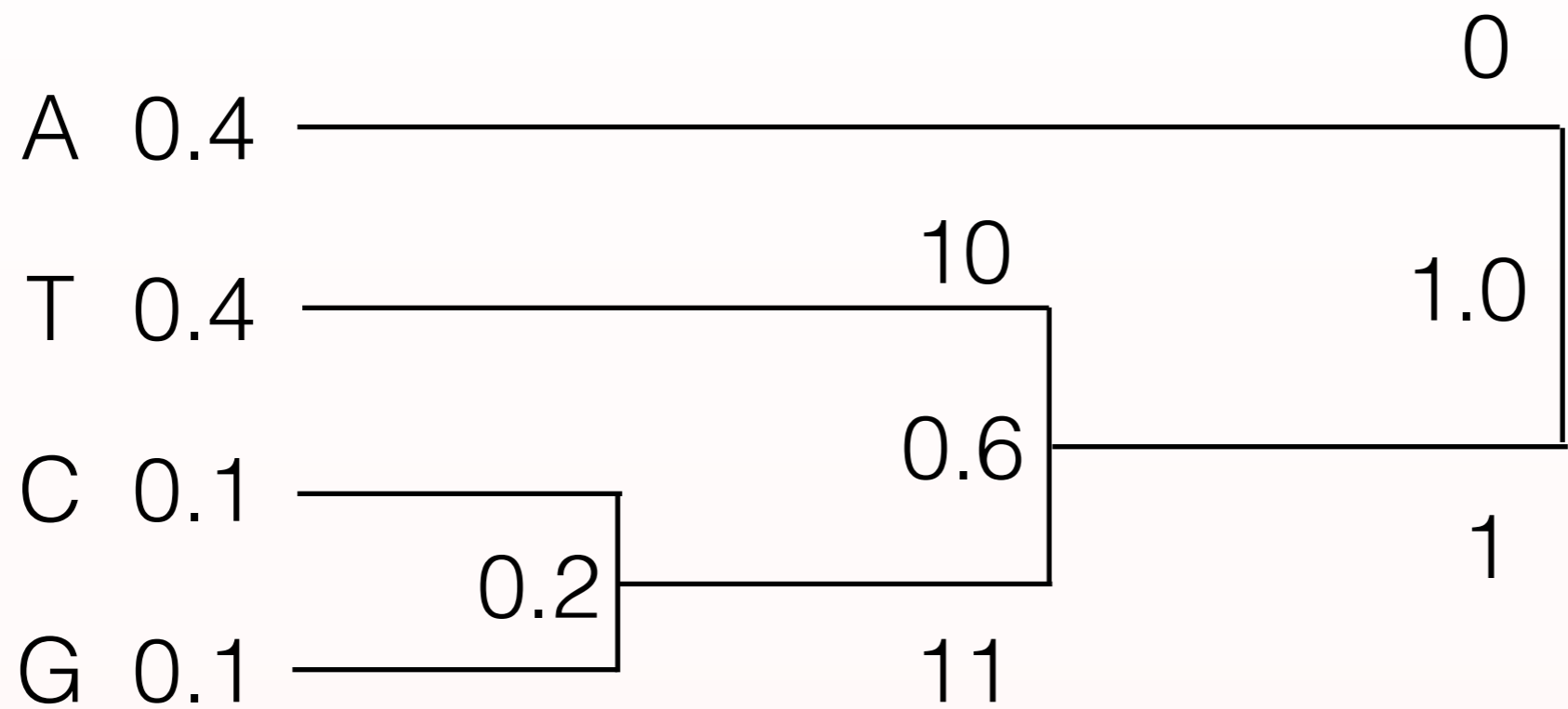


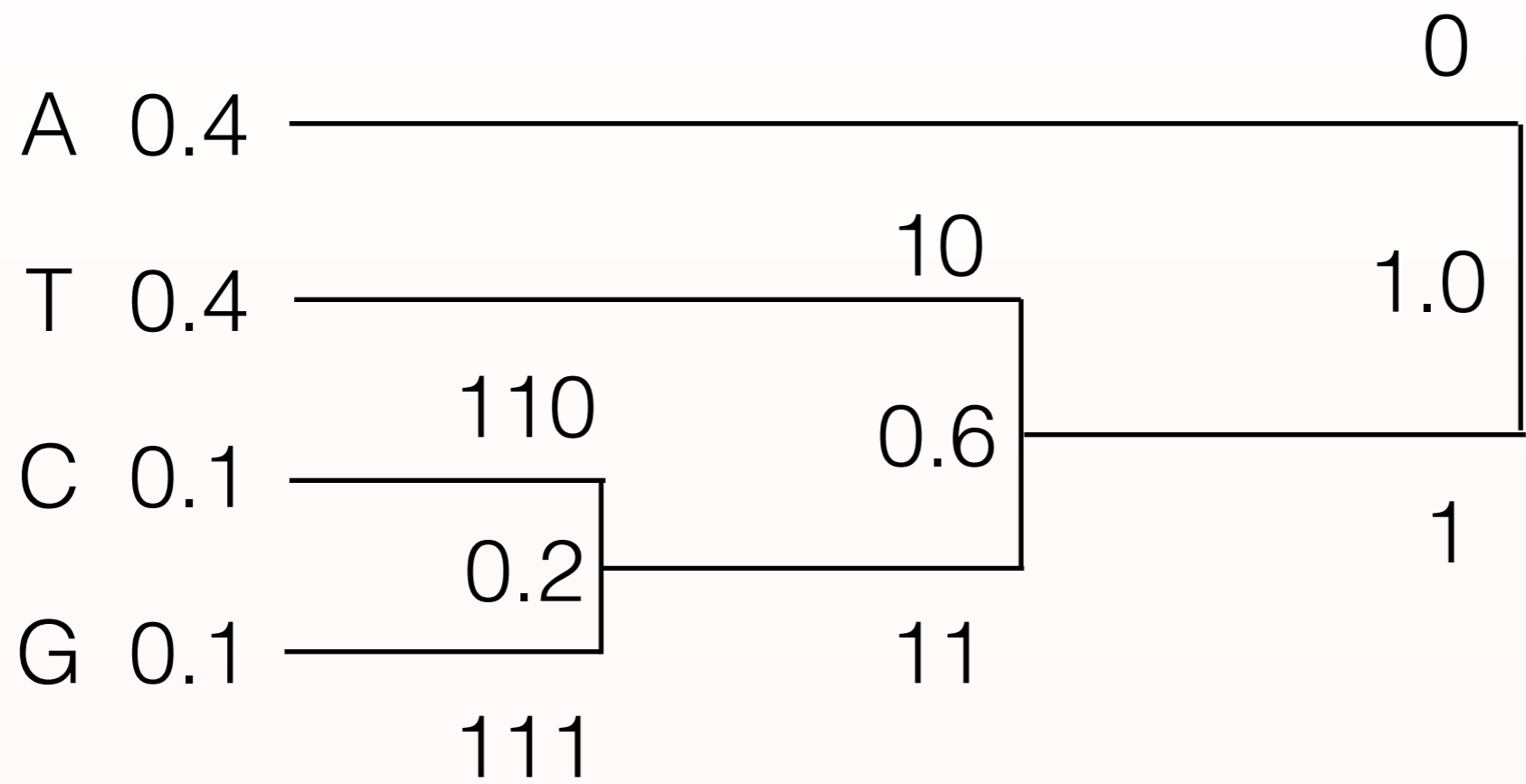












Huffman coding

Naive encoding: 2 bits/base

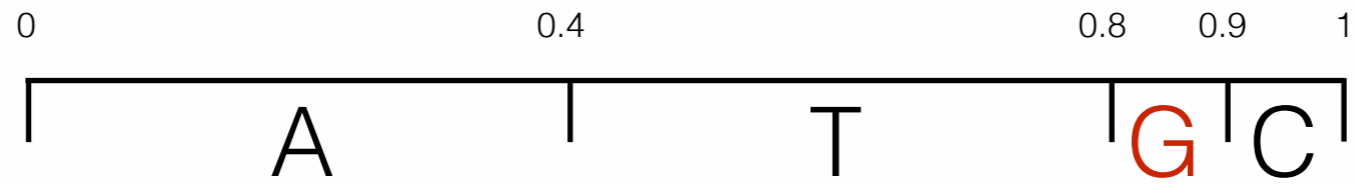
Huffman coding: 1.8 bits/base

	A	T	G	C
F	0.4	0.4	0.1	0.1
Codeword	0	10	110	111
Bits	1	2	3	3
Avg Bits	0.4	0.8	0.3	0.3

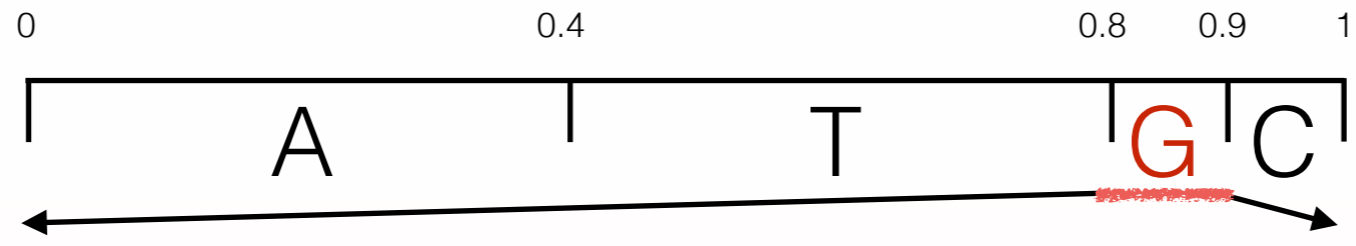
Arithmetic coding

Encode entire message as a number

	A	T	G	C
F	0.4	0.4	0.1	0.1
Interval	$[0,0.4)$	$[0.4,0.8)$	$[0.8,0.9)$	$[0.9,1)$

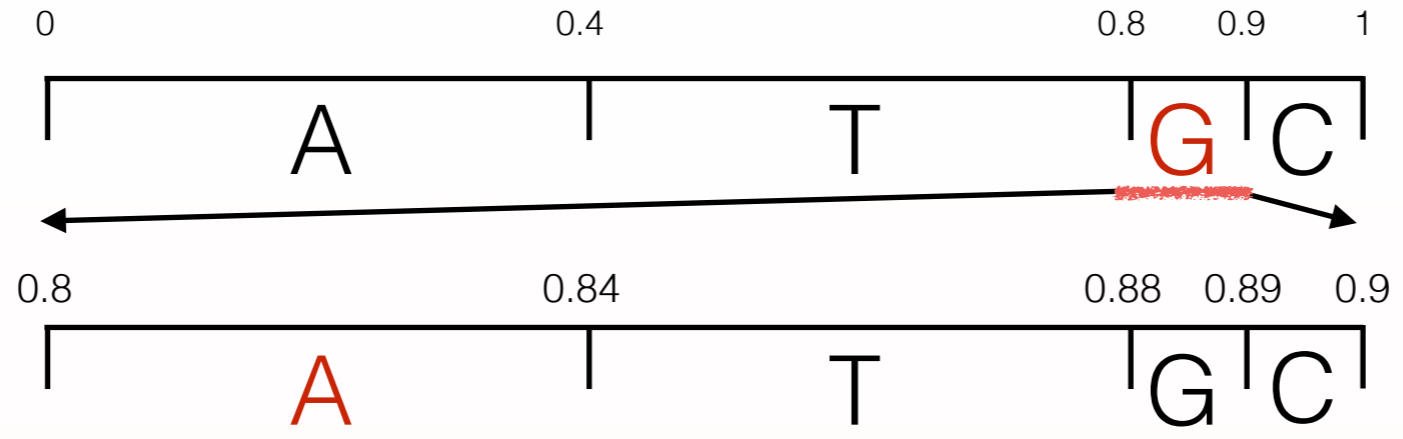


GATTACA

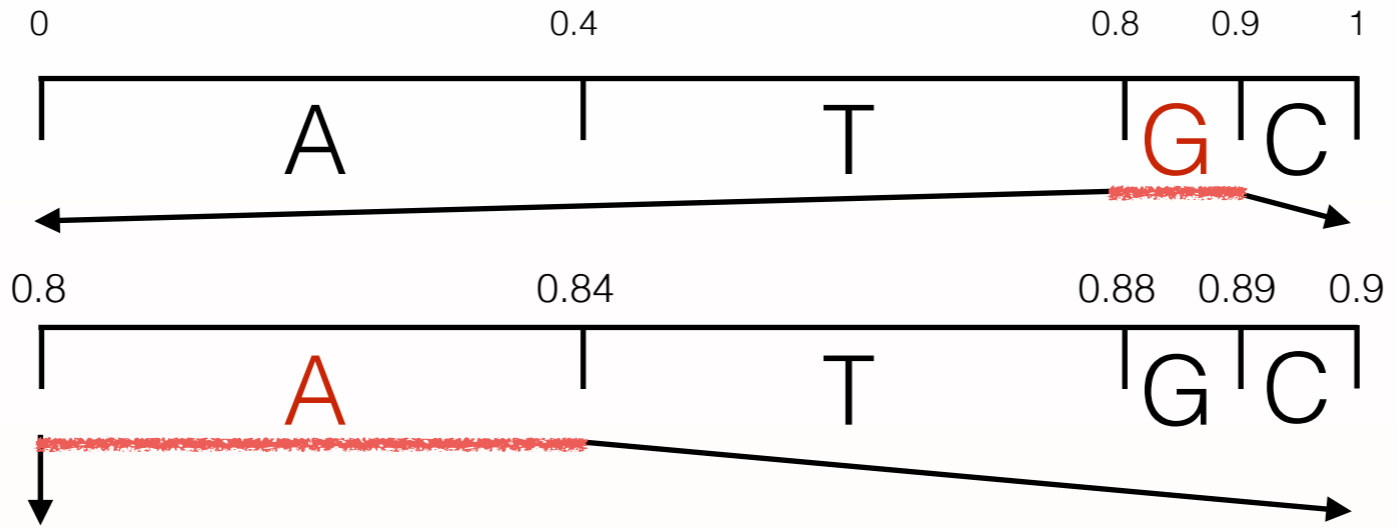


GATTACA

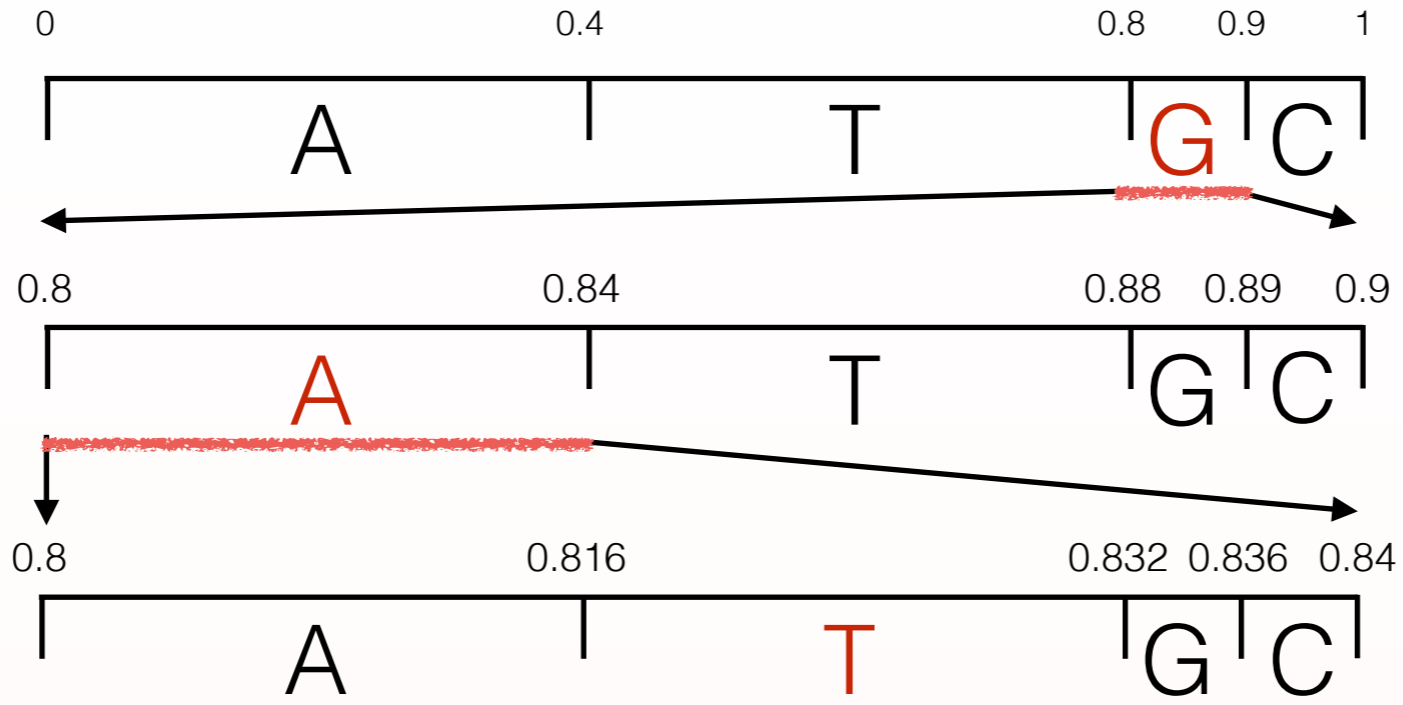
GATTACA



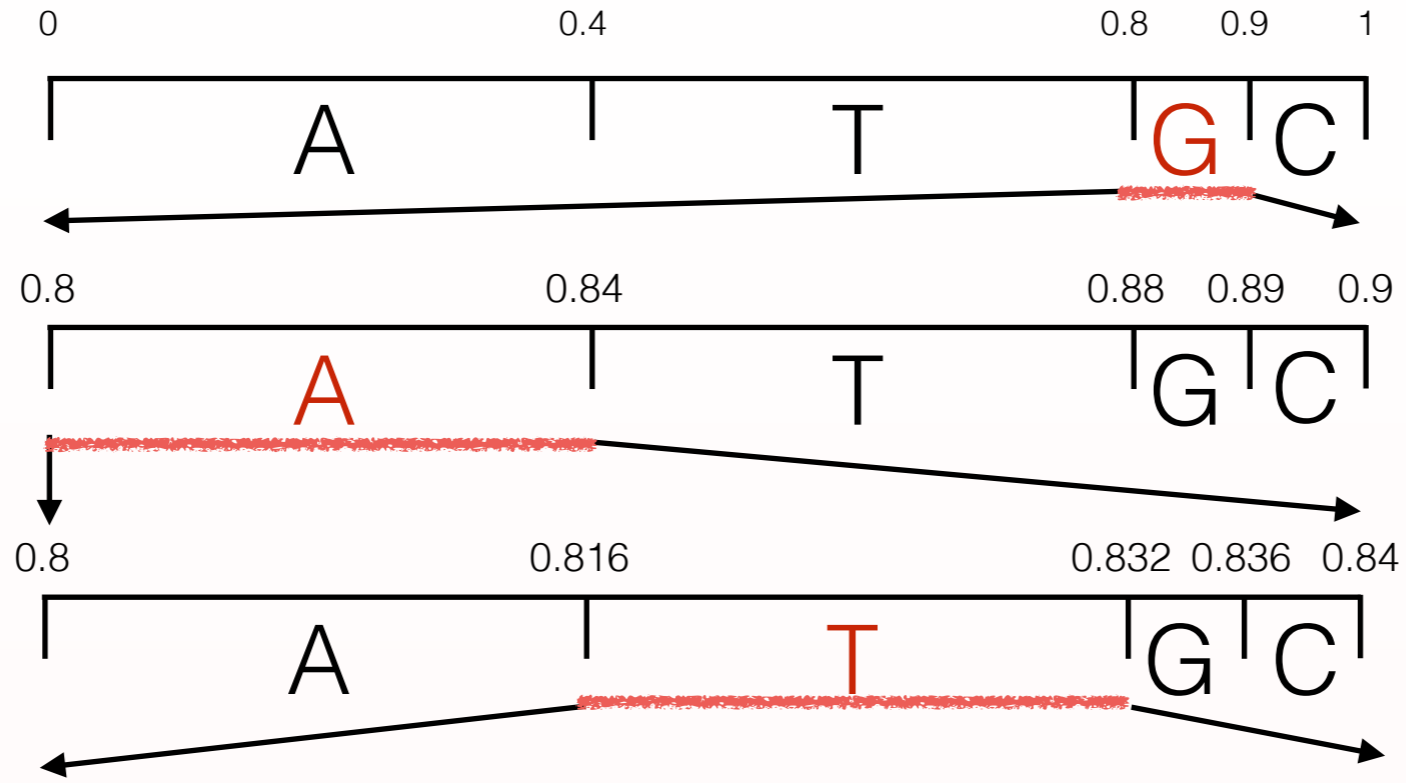
GATTACA



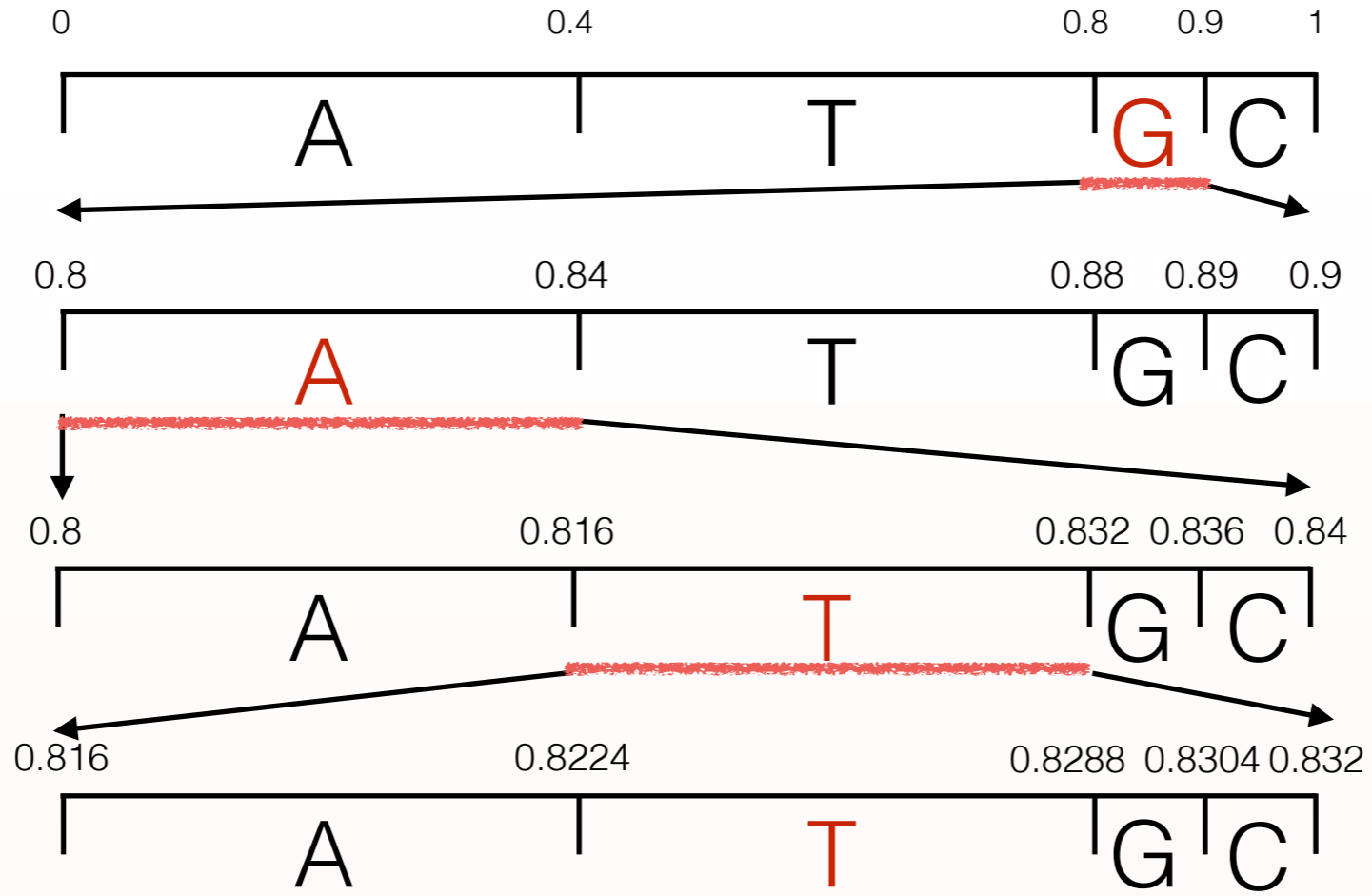
GATTACA



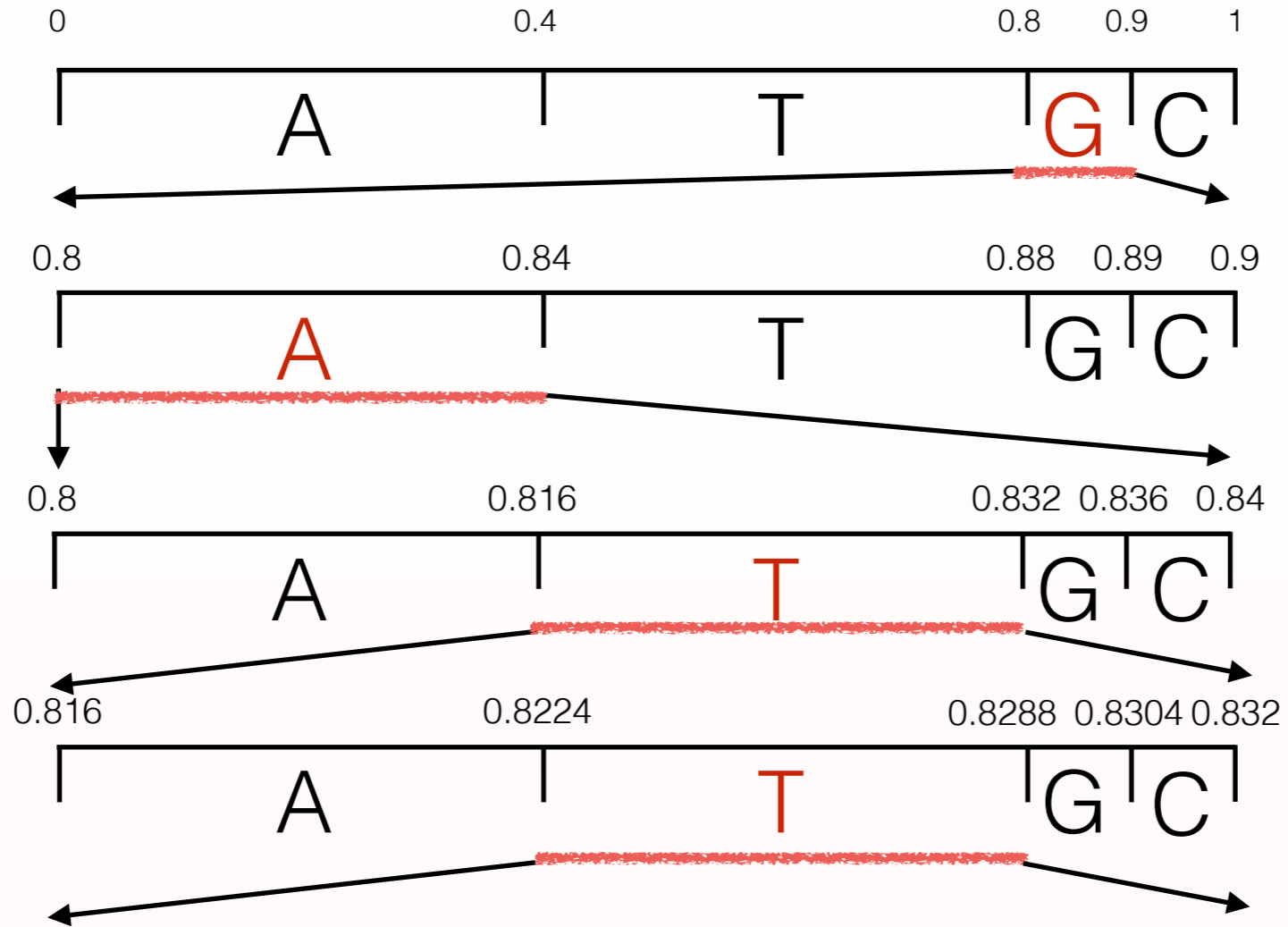
GATTACA



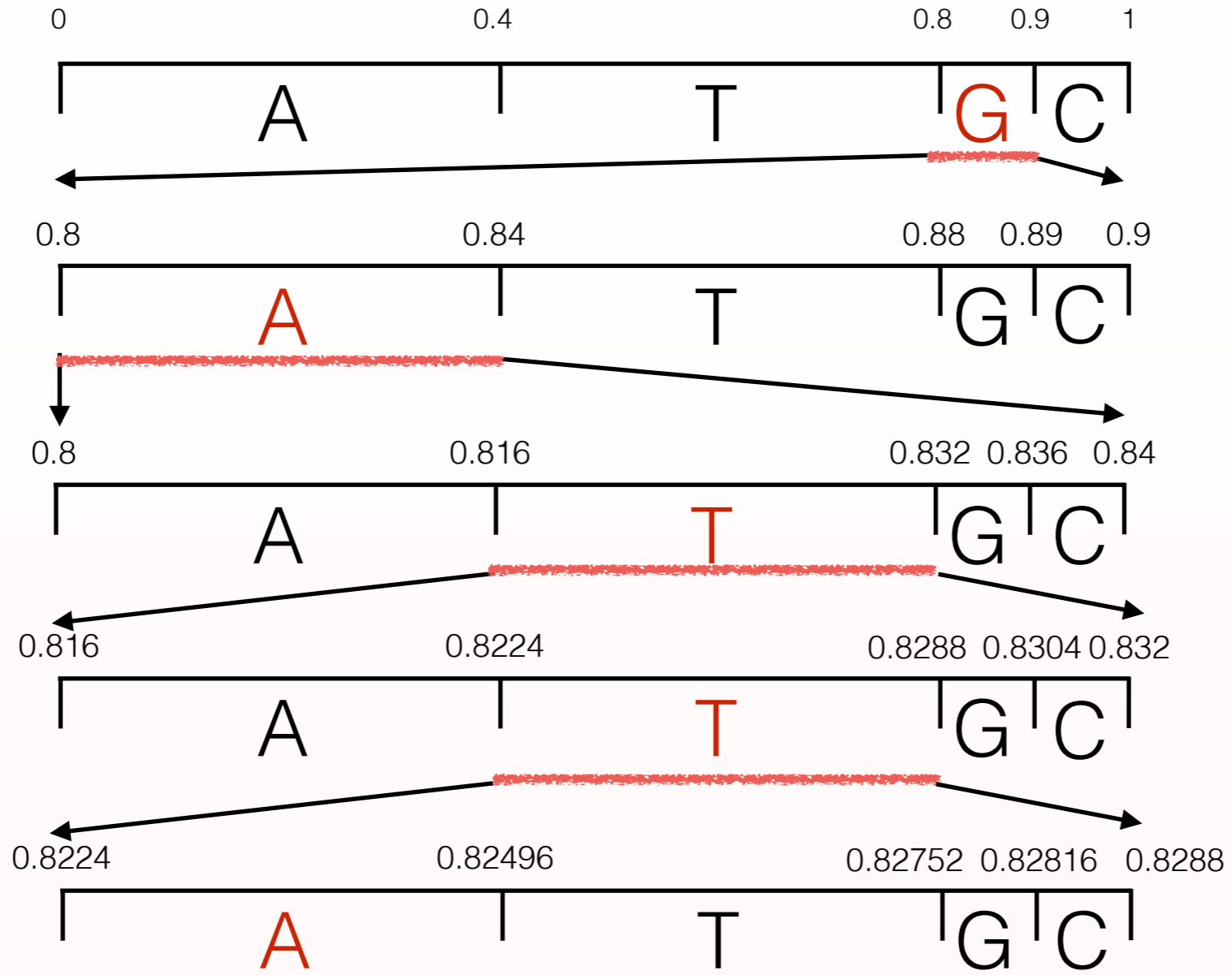
GATTACA



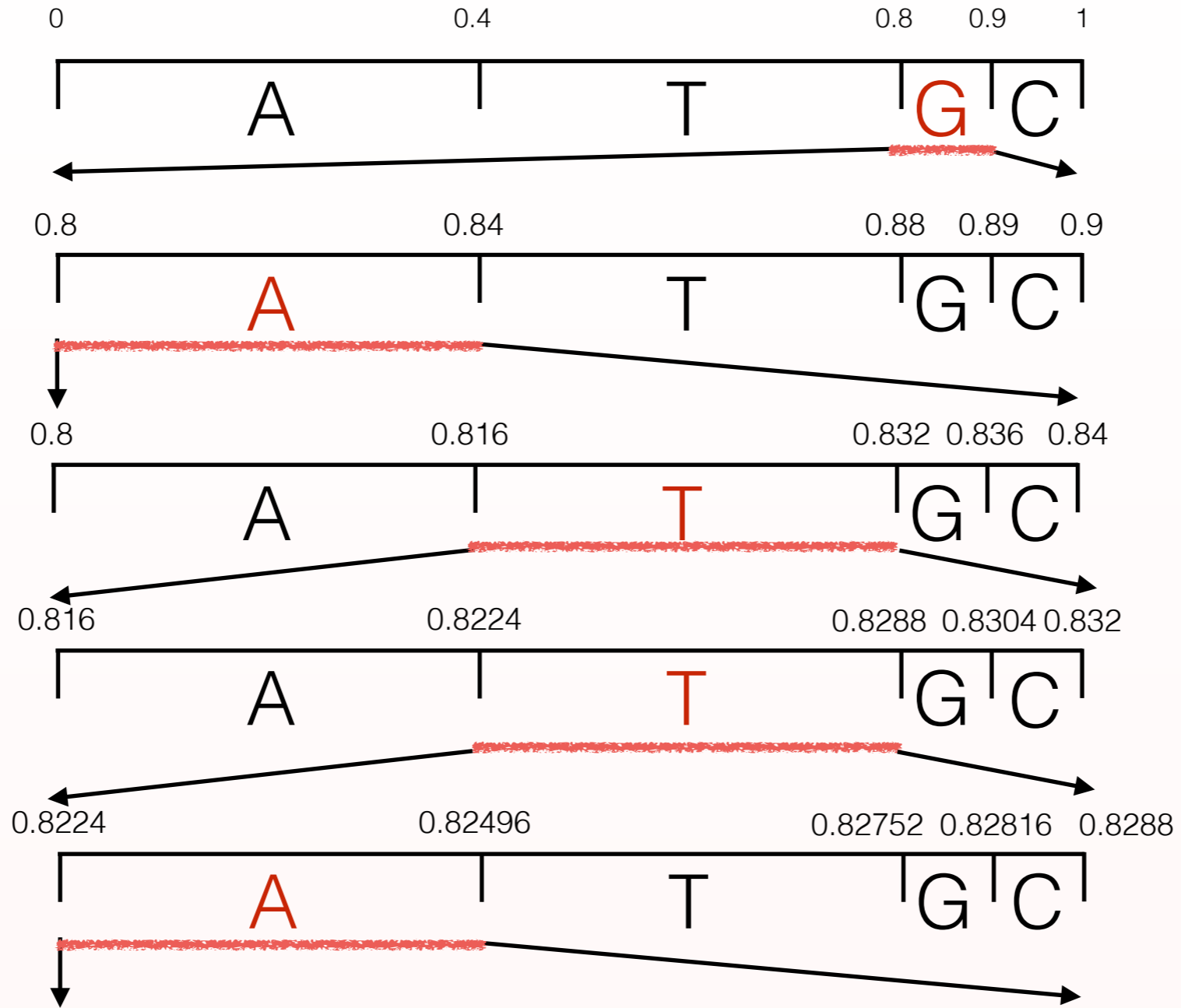
GATTACA



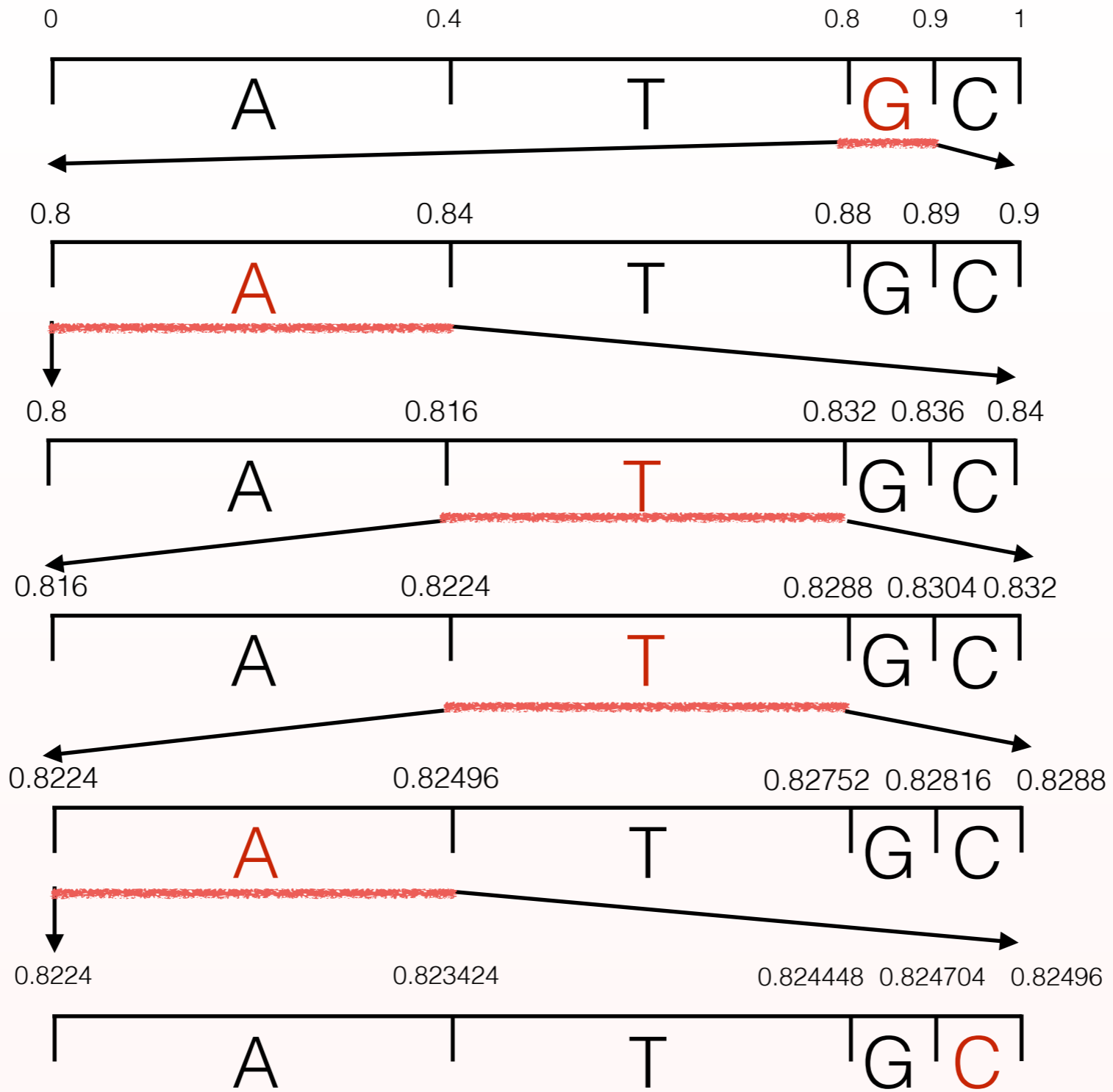
GATTACA



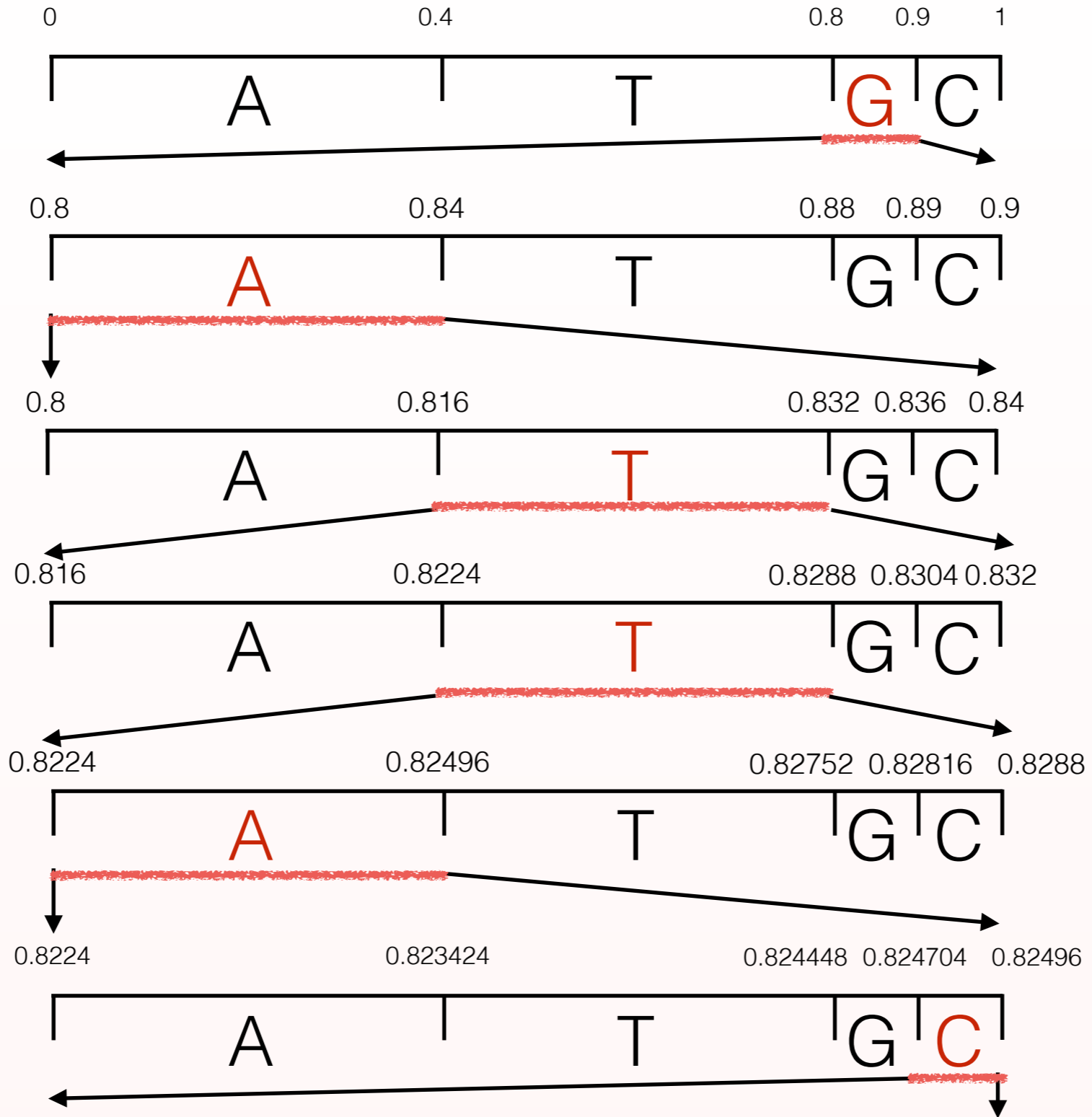
GATTACA



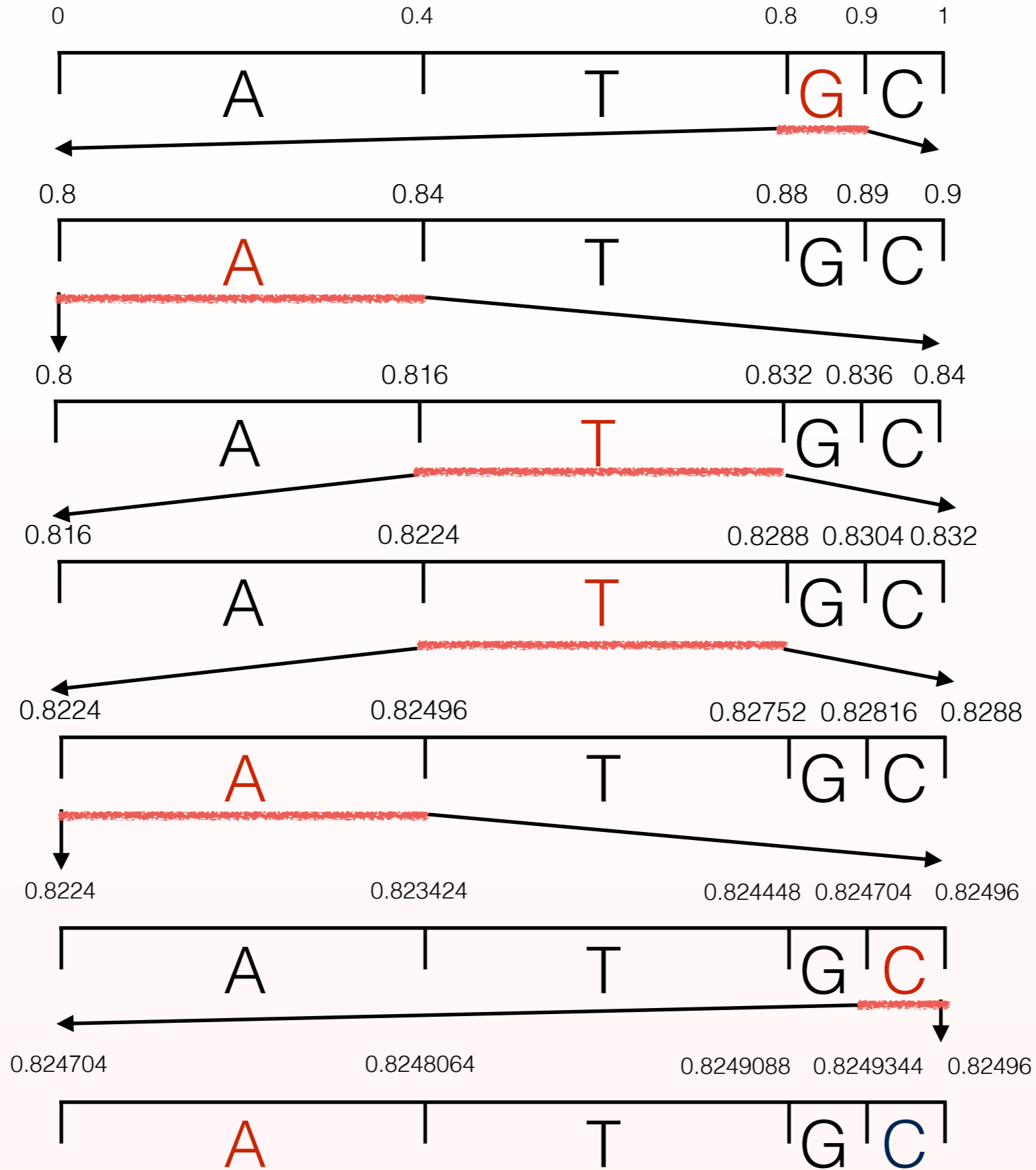
GATTACA



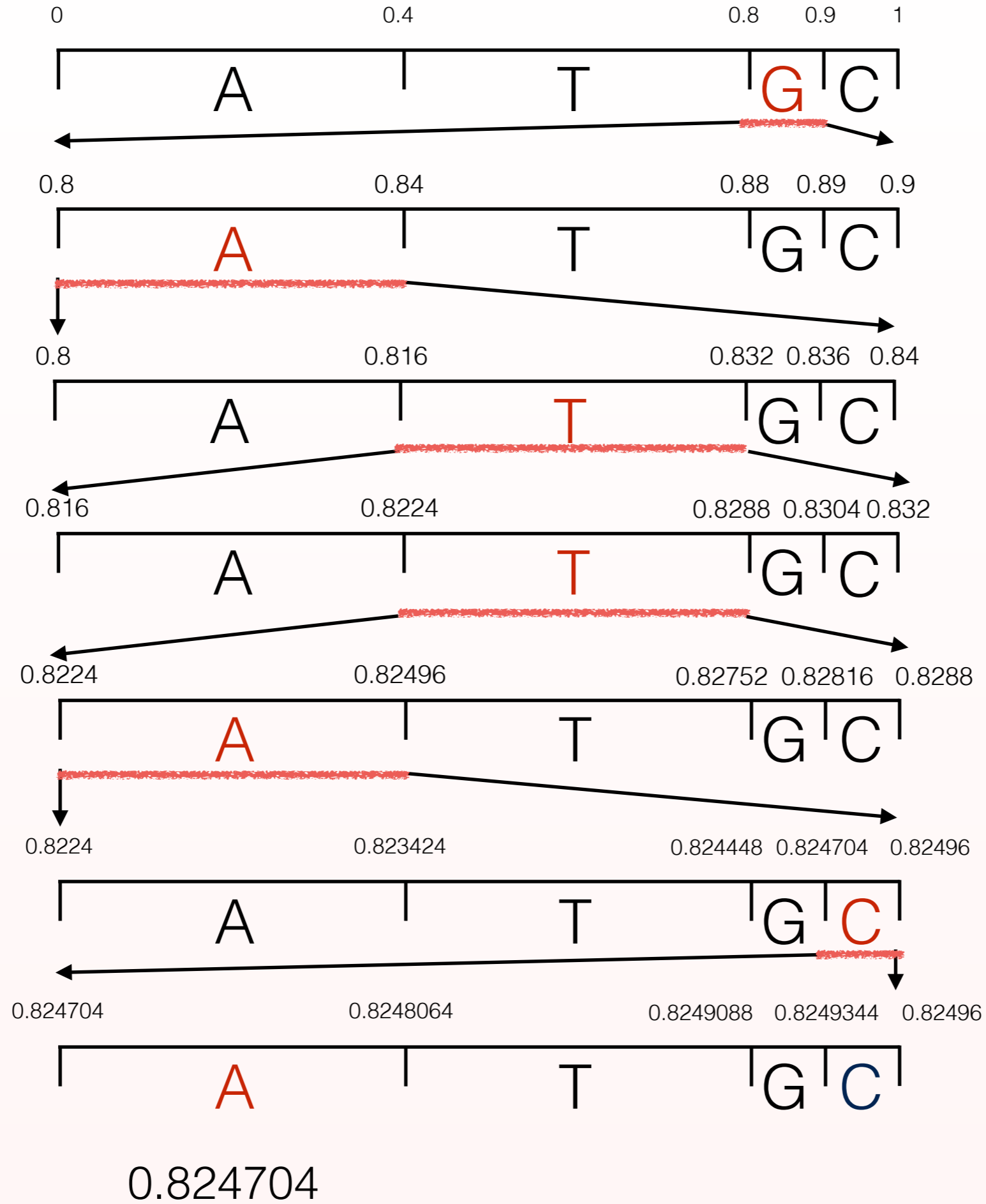
GATTACA



GATTACA



GATTACA



Arithmetic coding

Arithmetic coding

How do we know when the message ends?

Arithmetic coding

How do we know when the message ends?

I lied; we need a stop symbol

Arithmetic coding

How do we know when the message ends?

I lied; we need a stop symbol

Bigger problem: floating point precision

Arithmetic coding

How do we know when the message ends?

I lied; we need a stop symbol

Bigger problem: floating point precision

Integer version

Arithmetic coding

How do we know when the message ends?

I lied; we need a stop symbol

Bigger problem: floating point precision

Integer version

Rodionov & Volkov, 2007 & 2010 (p-adic arithmetic)

Arithmetic coding

How do we know when the message ends?

I lied; we need a stop symbol

Bigger problem: floating point precision

Integer version

Rodionov & Volkov, 2007 & 2010 (p-adic arithmetic)

Huffman coding is a specialized case, less likely to be optimal

Lempel-Ziv

Slide a fixed-length window along sequence

Replace already-seen patterns (of some maximum length)

Used by gzip (& others)

GATTACATTA

Lempel-Ziv

Slide a fixed-length window along sequence

Replace already-seen patterns (of some maximum length)

Used by gzip (& others)

GATTACATTA

Lempel-Ziv

Slide a fixed-length window along sequence

Replace already-seen patterns (of some maximum length)

Used by gzip (& others)

GATTACATTA

Lempel-Ziv

Slide a fixed-length window along sequence

Replace already-seen patterns (of some maximum length)

Used by gzip (& others)

GATTACATTA

Lempel-Ziv

Slide a fixed-length window along sequence

Replace already-seen patterns (of some maximum length)

Used by gzip (& others)

GATTACATTA

Lempel-Ziv

Slide a fixed-length window along sequence

Replace already-seen patterns (of some maximum length)

Used by gzip (& others)

GATTACATTA

Lempel-Ziv

Slide a fixed-length window along sequence

Replace already-seen patterns (of some maximum length)

Used by gzip (& others)

GATTACATTA

Lempel-Ziv

Slide a fixed-length window along sequence

Replace already-seen patterns (of some maximum length)

Used by gzip (& others)

GATTACATTA

GATTAC<1, 4>

Lempel-Ziv

Slide a fixed-length window along sequence

Replace already-seen patterns (of some maximum length)

Used by gzip (& others)

GATTAC**ATTA**

GATTAC<1, 4>

Pointer



Lempel-Ziv

Slide a fixed-length window along sequence

Replace already-seen patterns (of some maximum length)

Used by gzip (& others)

GATTAC**ATTA**

GATTAC<1, 4>

Pointer

Length

Burrows-Wheeler Transform (BWT)

Burrows-Wheeler Transform (BWT)

block-sorting compression

Burrows-Wheeler Transform (BWT)

block-sorting compression

not itself a compressor, but a way to sort
input, reversibly

Burrows-Wheeler Transform (BWT)

block-sorting compression

not itself a compressor, but a way to sort input, reversibly

Lends itself to run-length encoding

Burrows-Wheeler Transform (BWT)

block-sorting compression

not itself a compressor, but a way to sort input, reversibly

Lends itself to run-length encoding

Genomic compressors

Genomic compressors

Why not just gzip it?

Genomic compressors

Why not just gzip it?

sequence

Genomic compressors

Why not just gzip it?

sequence

quality scores (one per base)

Genomic compressors

Why not just gzip it?

sequence

quality scores (one per base)

sequence: 4-letter alphabet

Genomic compressors

Why not just gzip it?

sequence

quality scores (one per base)

sequence: 4-letter alphabet

quality scores: 40-value range

Genomic compressors

Why not just gzip it?

sequence

quality scores (one per base)

sequence: 4-letter alphabet

quality scores: 40-value range

sequence: must be lossless

Genomic compressors

Why not just gzip it?

sequence

quality scores (one per base)

sequence: 4-letter alphabet

quality scores: 40-value range

sequence: must be lossless

quality scores: lossy might be ok

Genomic compressors

Genomic compressors

Take advantage of structure

Genomic compressors

Take advantage of structure

Store differences from a **reference**

Genomic compressors

Take advantage of structure

Store differences from a **reference**

Cluster similar reads (**reference-free**)

Genomic compressors

Take advantage of structure

Store differences from a **reference**

Cluster similar reads (**reference-free**)

don't have references for all species

SCALCE [Hach et al. 2012]

SCALCE [Hach et al. 2012]

Reference-free

SCALCE [Hach et al. 2012]

Reference-free

Locally Consistent Encoding

SCALCE [Hach et al. 2012]

Reference-free

Locally Consistent Encoding

Buckets reads by substrings

SCALCE [Hach et al. 2012]

Reference-free

Locally Consistent Encoding

Buckets reads by substrings

improves LZ77 runtime AND compression

SCALCE [Hach et al. 2012]

SCALCE [Hach et al. 2012]

Identify local maxima

SCALCE [Hach et al. 2012]

Identify local maxima

Identify local minima

SCALCE [Hach et al. 2012]

Identify local maxima
Identify local minima
Partition at each side

SCALCE [Hach et al. 2012]

Identify local maxima

Identify local minima

Partition at each side

Extend L and R to “core blocks”

SCALCE [Hach et al. 2012]

Identify local maxima

Identify local minima

Partition at each side

Extend L and R to “core blocks”

X0X

SCALCE [Hach et al. 2012]

Identify local maxima

Identify local minima

Partition at each side

Extend L and R to “core blocks”

X0X

21312032102021312032102

SCALCE [Hach et al. 2012]

Identify local maxima

Identify local minima

Partition at each side

Extend L and R to “core blocks”

XOX

21312032102021312032102

SCALCE [Hach et al. 2012]

Identify local maxima

Identify local minima

Partition at each side

Extend L and R to “core blocks”

XOX

21312032102021312032102

| 213 | 12 | 03 | 2102 | 02 | 13 | 12 | 03 | 2102 |

SCALCE [Hach et al. 2012]

Identify local maxima

Identify local minima

Partition at each side

Extend L and R to “core blocks”

XOX

21312032102021312032102

| 213 | 12 | 03 | 2102 | 02 | 13 | 12 | 03 | 2102 |

**2131, 3120, 2032, 321020, 2021, 2131,
3120, 2032, 32102**

SCALCE [Hach et al. 2012]

Core blocks between 8 & 20

99% of all HTS reads of length ≥ 50

include at least one core of length ≤ 14

Identify overlapping reads by core blocks

Reorder reads to favor LZ77

DeeZ [Hach, et al. 2014]

DeeZ [Hach, et al. 2014]

Given a mapping (reads to reference)

DeeZ [Hach, et al. 2014]

Given a mapping (reads to reference)

Partition reads into blocks according to locus

DeeZ [Hach, et al. 2014]

Given a mapping (reads to reference)

Partition reads into blocks according to locus

Build contig covering all reads within a block

DeeZ [Hach, et al. 2014]

Given a mapping (reads to reference)

Partition reads into blocks according to locus

Build contig covering all reads within a block

Encode locus for each read within contig

DeeZ [Hach, et al. 2014]

Given a mapping (reads to reference)

Partition reads into blocks according to locus

Build contig covering all reads within a block

Encode locus for each read within contig

Encode rare differences (sequencing errors)

DeeZ [Hach, et al. 2014]

GTCGTCCTACA T
CGTCGTCCTACA
AACCTCGTCTAC
GTCTACA TCTA
ACGTGCTAAACGTCGTCTACAGTCTACAGA

DeeZ [Hach, et al. 2014]

```
      GTCGTC TACA T
      CGTCGTC TACA
AACCTCGTC TAC
      GTCTACA TCTA
ACGTGCTAAACGTCGT TACAG TCTACAGA
ACGTGCTAAACGTCGTCTACA TCTACAGA
```


DeeZ [Hach, et al. 2014]

GTCGTCCTACA T
CGTCGTCCTACA
AACCTCGTCTAC
GTCTACA TCTA
ACGTGCTAAACGTCGTTTACAGTCTACAGA
ACGTGCTAAACGTCGTCCTACA TCTACAGA

DeeZ [Hach, et al. 2014]

DeeZ [Hach, et al. 2014]

Tokenization of read names

DeeZ [Hach, et al. 2014]

Tokenization of read names

LZ77

DeeZ [Hach, et al. 2014]

Tokenization of read names

LZ77

Lossy QS compression from SCALCE

DeeZ [Hach, et al. 2014]

Tokenization of read names

LZ77

Lossy QS compression from SCALCE

Random access

PathEnc [Kingsford & Patro 2015]

Reference-based

–but no aligning

–statistical, generative model of reads

For RNA-seq data (but need not be)

PathEnc [Kingsford & Patro 2015]

GAUU

GAUUAGAUUG

PathEnc [Kingsford & Patro 2015]

GAUU
↓
AUUA

GAUUAGAUUG

PathEnc [Kingsford & Patro 2015]

GAUU
↓
AUUA
↓
UUAG

GAUUAGAUUG

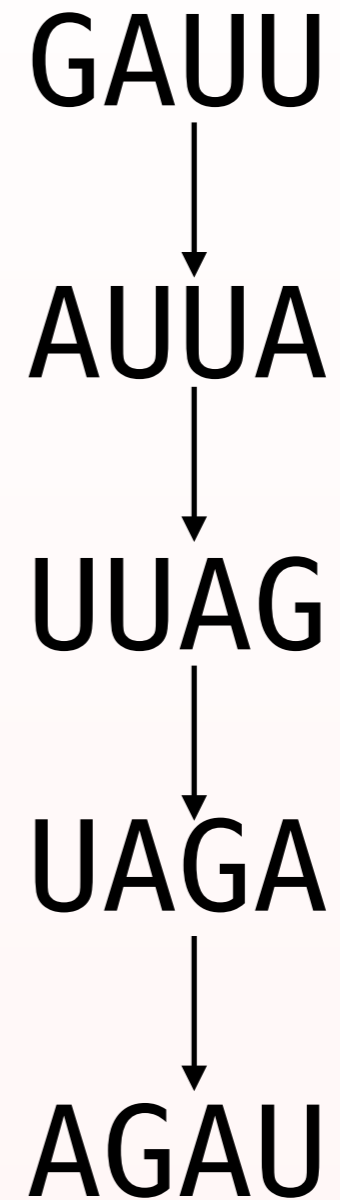
PathEnc [Kingsford & Patro 2015]

GAUU
↓
AUUA
↓
UUAG
↓
UAGA

GAUUAGAUUG

PathEnc [Kingsford & Patro 2015]

GAUUAGAUUG

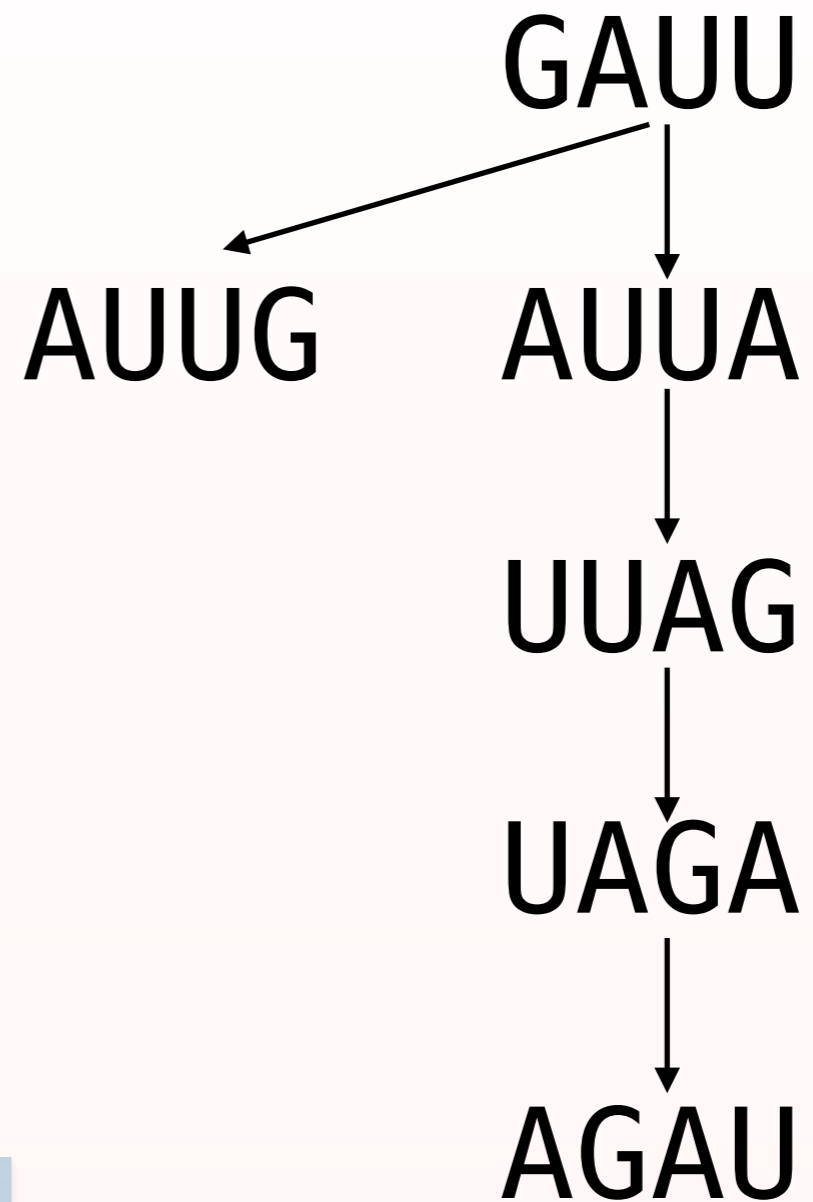


PathEnc [Kingsford & Patro 2015]

GAUUAGAUUG

GAUU
↓
AUUA
↓
UUAG
↓
UAGA
↓
AGAU

PathEnc [Kingsford & Patro 2015]



GAUUAGAUUG

PathEnc [Kingsford & Patro 2015]

PathEnc [Kingsford & Patro 2015]

Path encoding in graph G

PathEnc [Kingsford & Patro 2015]

Path encoding in graph G

1 node per k-mer

PathEnc [Kingsford & Patro 2015]

Path encoding in graph G

1 node per k -mer

edge between k -mers (u, v) if v follows u

PathEnc [Kingsford & Patro 2015]

Path encoding in graph G

1 node per k -mer

edge between k -mers (u, v) if v follows u

Each read encoded as a path in G

PathEnc [Kingsford & Patro 2015]

PathEnc [Kingsford & Patro 2015]

1st node of each read path (read head)

PathEnc [Kingsford & Patro 2015]

1st node of each read path (read head)

4-ary tree of depth k

PathEnc [Kingsford & Patro 2015]

1st node of each read path (read head)

4-ary tree of depth k

edges removed for nonexistent k -mers

PathEnc [Kingsford & Patro 2015]

1st node of each read path (read head)

4-ary tree of depth k

edges removed for nonexistent k -mers

traverse in fixed order, emitting 1 for edge

PathEnc [Kingsford & Patro 2015]

1st node of each read path (read head)

4-ary tree of depth k

edges removed for nonexistent k -mers

traverse in fixed order, emitting 1 for edge

resulting bit string gzipped

PathEnc [Kingsford & Patro 2015]

PathEnc [Kingsford & Patro 2015]

remaining nodes (read tails)

PathEnc [Kingsford & Patro 2015]

remaining nodes (read tails)

arithmetic coding

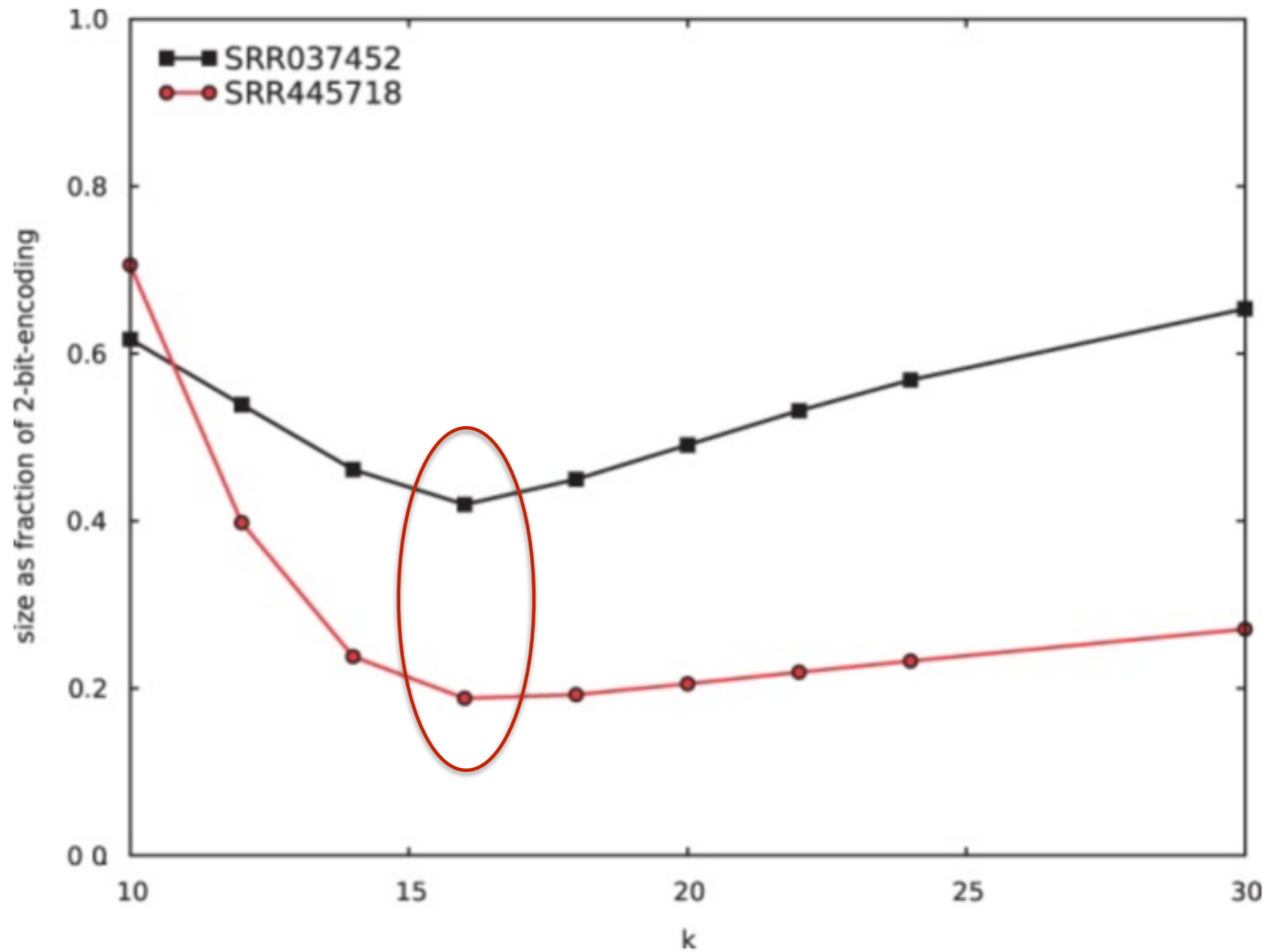
PathEnc [Kingsford & Patro 2015]

remaining nodes (read tails)

arithmetic coding

probability distribution per node in G

PathEnc [Kingsford & Patro 2015]



MINCE [Patro & Kingsford 2015]

MINCE [Patro & Kingsford 2015]

Reference-free

MINCE [Patro & Kingsford 2015]

Reference-free

Buckets reads based on k -mers ($k=15$)

MINCE [Patro & Kingsford 2015]

Reference-free

Buckets reads based on k -mers ($k=15$)

Replace common k -mer with pointer

MINCE [Patro & Kingsford 2015]

Reference-free

Buckets reads based on k -mers ($k=15$)

Replace common k -mer with pointer

Reorder reads within each bucket

MINCE [Patro & Kingsford 2015]

Reference-free

Buckets reads based on k -mers ($k=15$)

Replace common k -mer with pointer

Reorder reads within each bucket

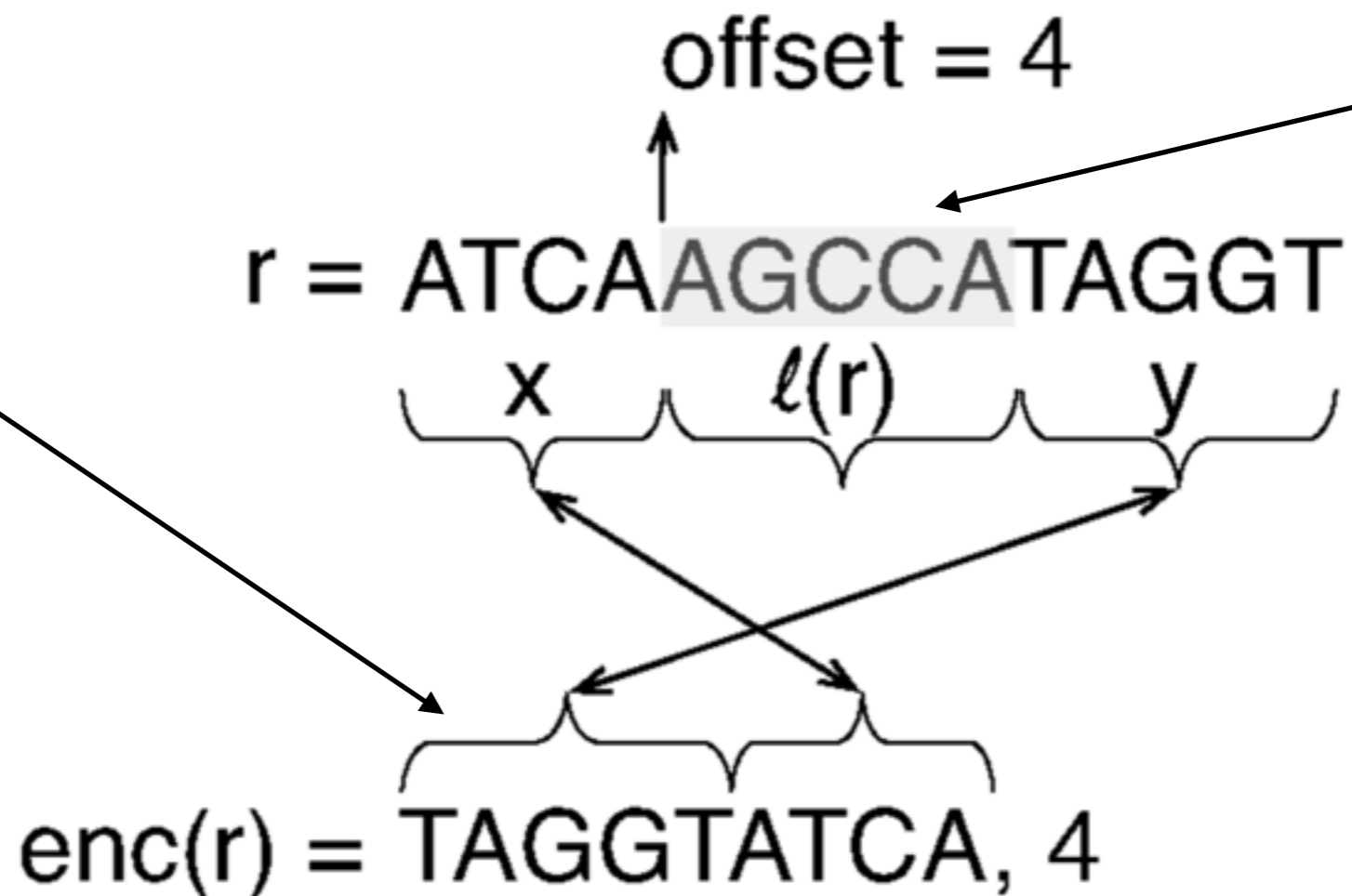
This boosts lzip performance

MINCE [Patro & Kingsford 2015]

Split-swap read transformation

Suffix of
bucket label

common
bucket
label

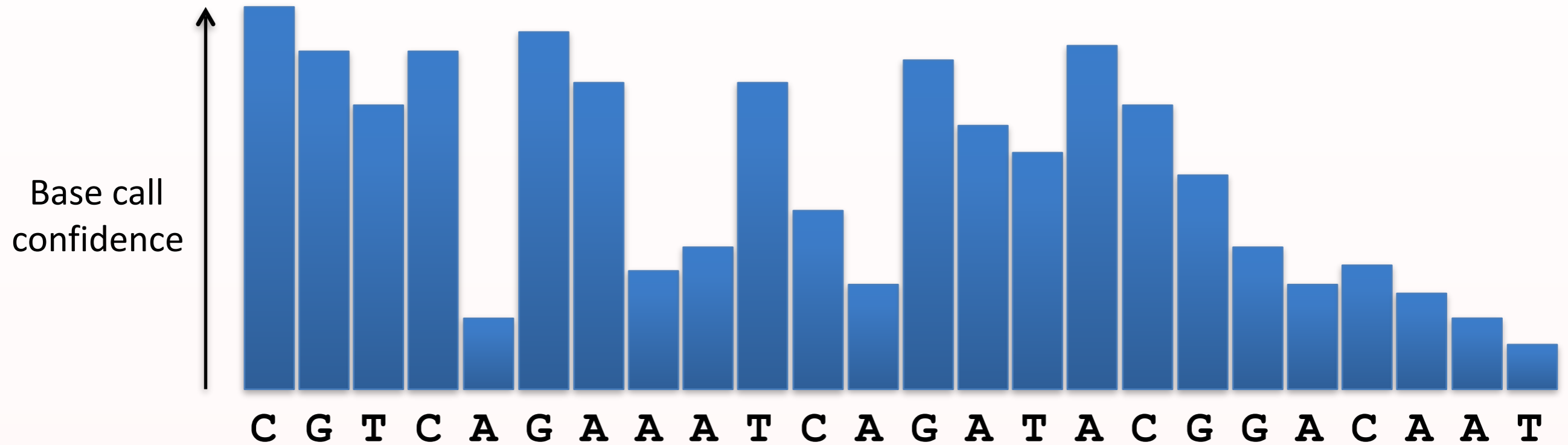


The Quality Score Problem

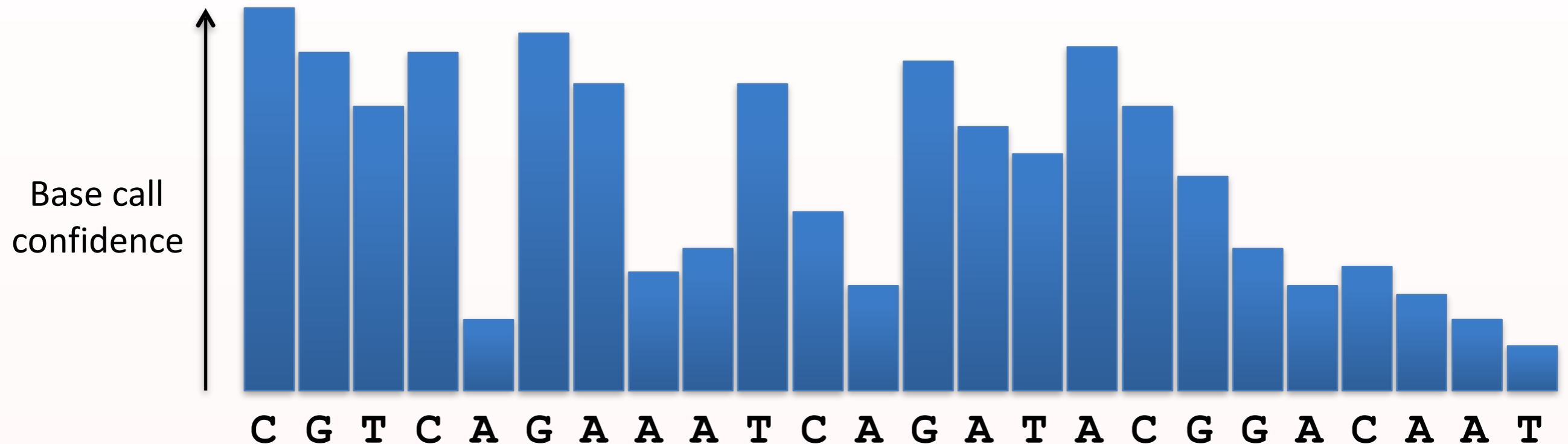
The Quality Score Problem

C G T C A G A A A T C A G A T A C G G A C A A T

The Quality Score Problem



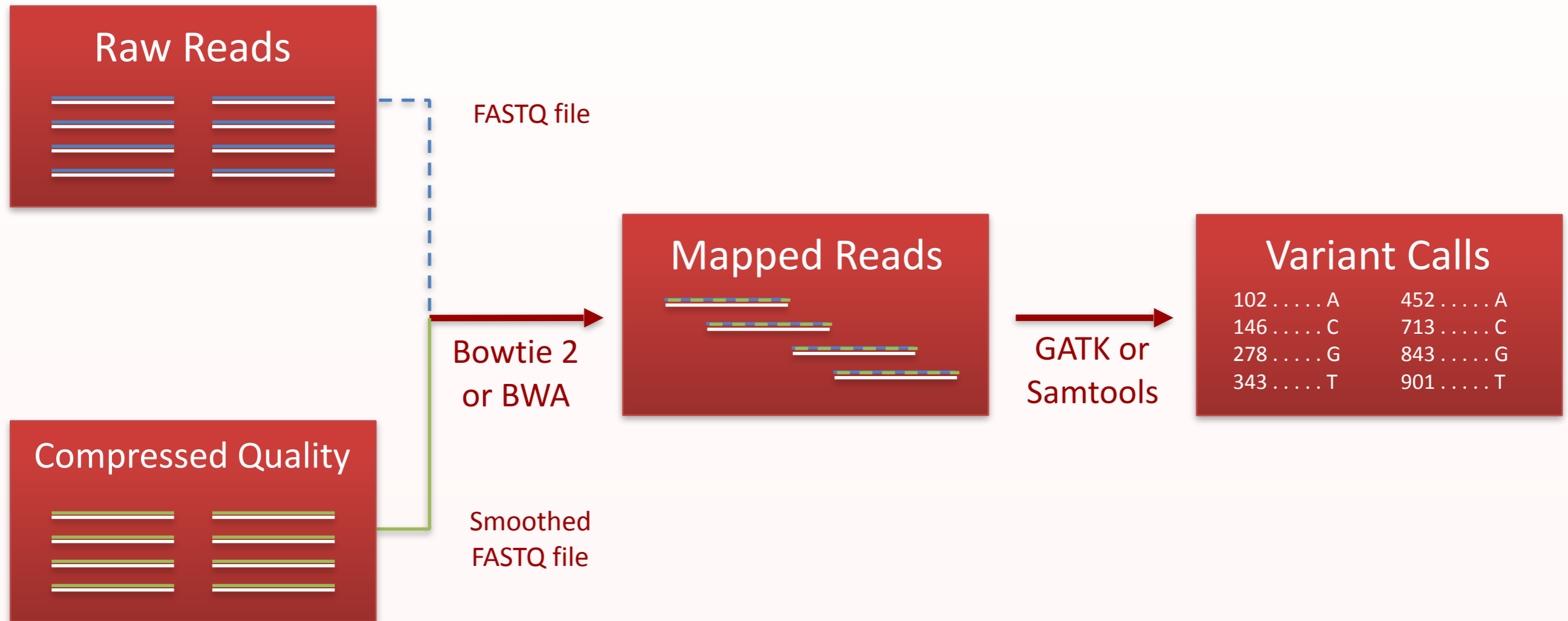
The Quality Score Problem



Given reads from a genome, can we **efficiently** compress quality scores while maintaining or even improving accuracy?

Quality scores in SNP-calling

Quality scores in SNP-calling



Quality score compressors

Quality score compressors

Greater dynamic range than sequence

Quality score compressors

Greater dynamic range than sequence

Can be lossless or lossy

Quality score compressors

Greater dynamic range than sequence

Can be lossless or lossy

Focus here on lossy

8-binning [Illumina]

Universally collapse dynamic range 40 => 8

Small effect on SNP calling error

Available in CRAM

QS bins	New value
N (no call)	N (no call)
2-9	6
10-19	15
20-24	22
25-29	27
30-34	33
35-39	37
≥40	40

SCALCE [Hach et al. 2012]

SCALCE [Hach et al. 2012]

Lossy compression of quality scores

SCALCE [Hach et al. 2012]

Lossy compression of quality scores

Neighboring Qs often similar

SCALCE [Hach et al. 2012]

Lossy compression of quality scores

Neighboring QSs often similar

Reduce dynamic range

SCALCE [Hach et al. 2012]

Lossy compression of quality scores

Neighboring Qs often similar

Reduce dynamic range

Arithmetic coding

SCALCE [Hach et al. 2012]

Lossy compression of quality scores

Neighboring Qs often similar

Reduce dynamic range

Arithmetic coding

Small loss (<0.1%) of SNP calling accuracy

QVZ [Malysa, et al. 2015]

QVZ [Malysa, et al. 2015]

Based on rate-distortion theory

Discard as little information as possible for a desired bit rate

Key insight: neighboring QSs are likely to be correlated

Illumina reads often have lower QSs at end

QVZ [Malysa, et al. 2015]

QVZ [Malysa, et al. 2015]

1. Compute the empirical transition probabilities of a Markov-1 model

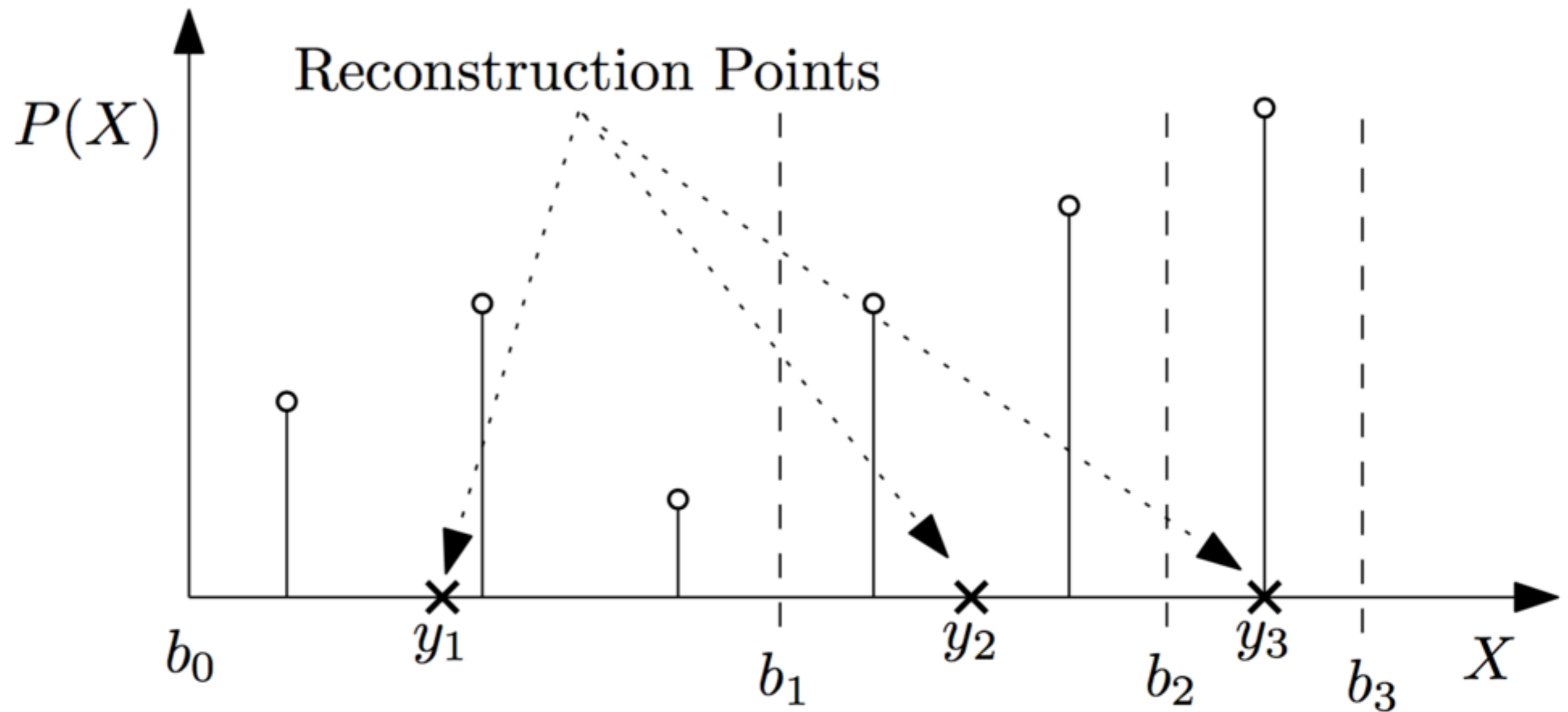
QVZ [Malysa, et al. 2015]

1. Compute the empirical transition probabilities of a Markov-1 model
2. Construct a codebook using the Lloyd-Max algorithm

QVZ [Malysa, et al. 2015]

1. Compute the empirical transition probabilities of a Markov-1 model
2. Construct a codebook using the Lloyd-Max algorithm
3. Quantize the input using the codebook, use arithmetic encoder

QVZ [Malysa, et al. 2015]



Quartz [Yu, et al. 2015]

Preprocessing

Corpus of NGS reads



Dictionary of common k-mers

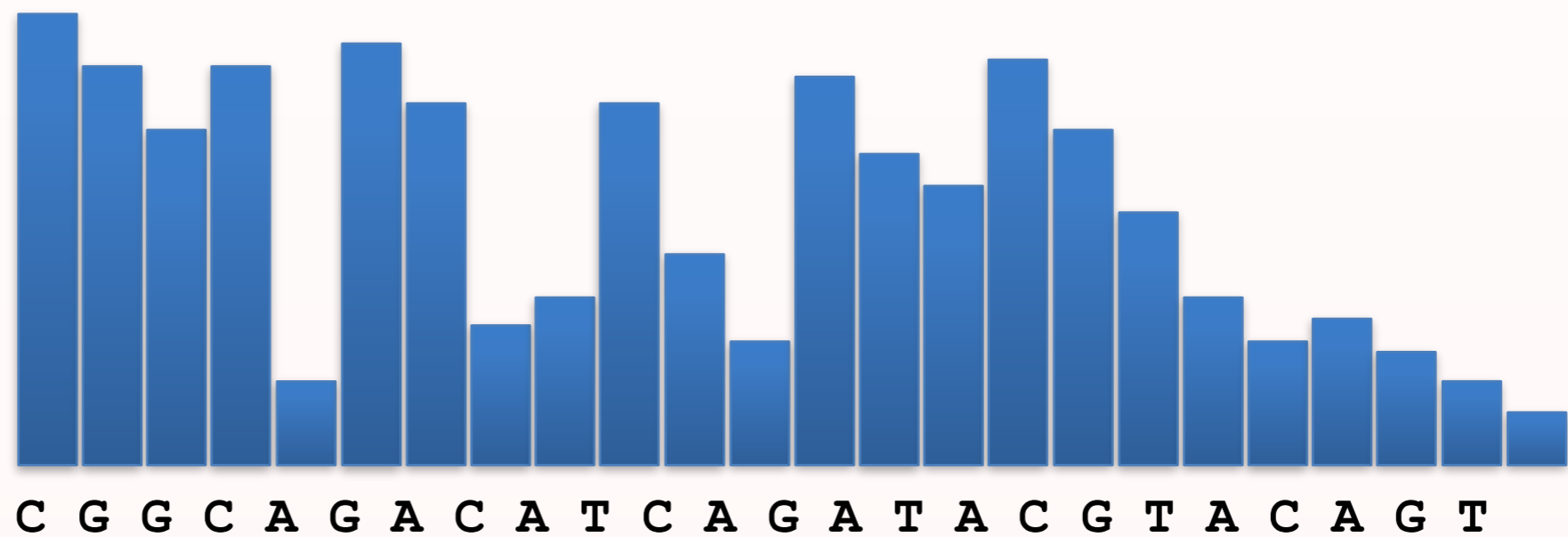


Quartz [Yu, et al. 2015]

Preprocessing



Compression

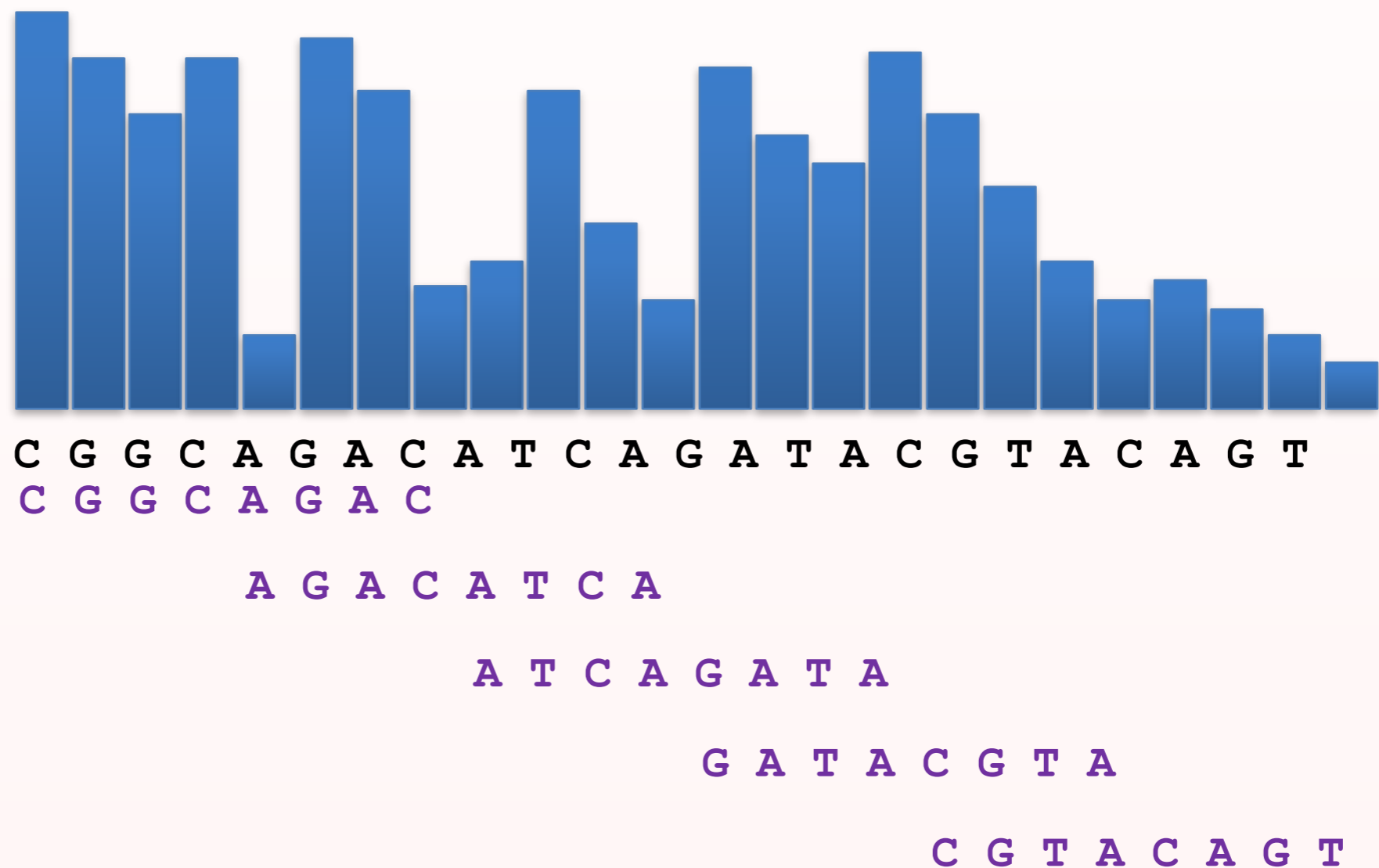


Quartz [Yu, et al. 2015]

Preprocessing

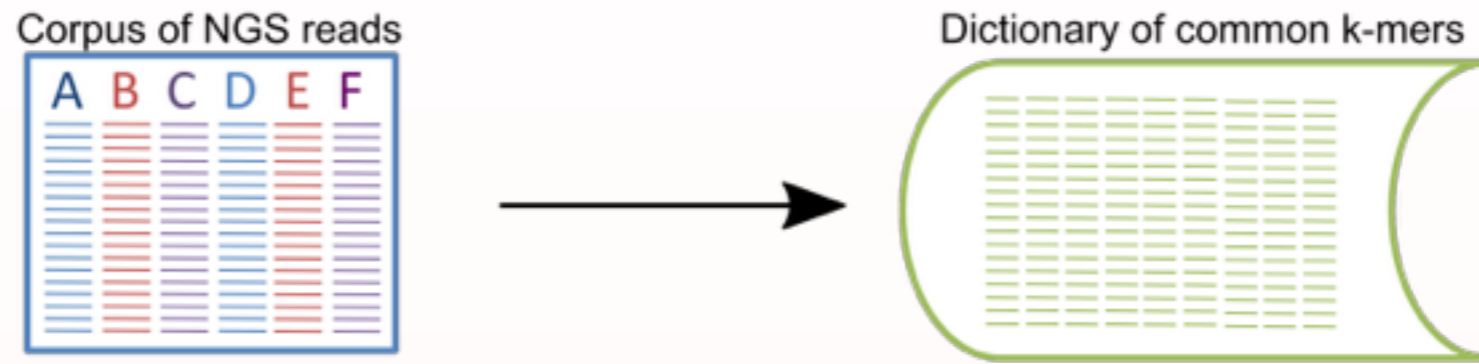


Compression

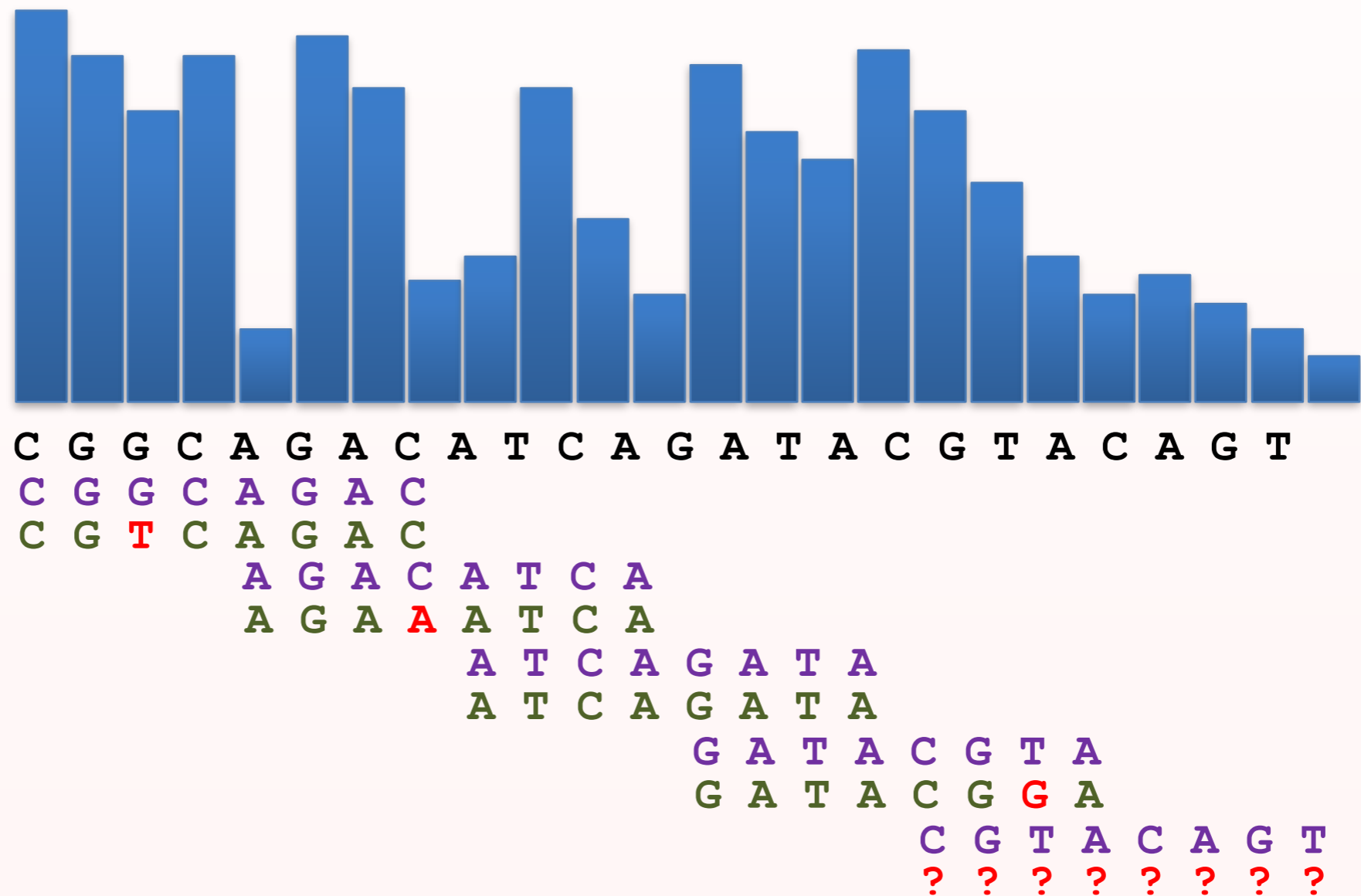


Quartz [Yu, et al. 2015]

Preprocessing

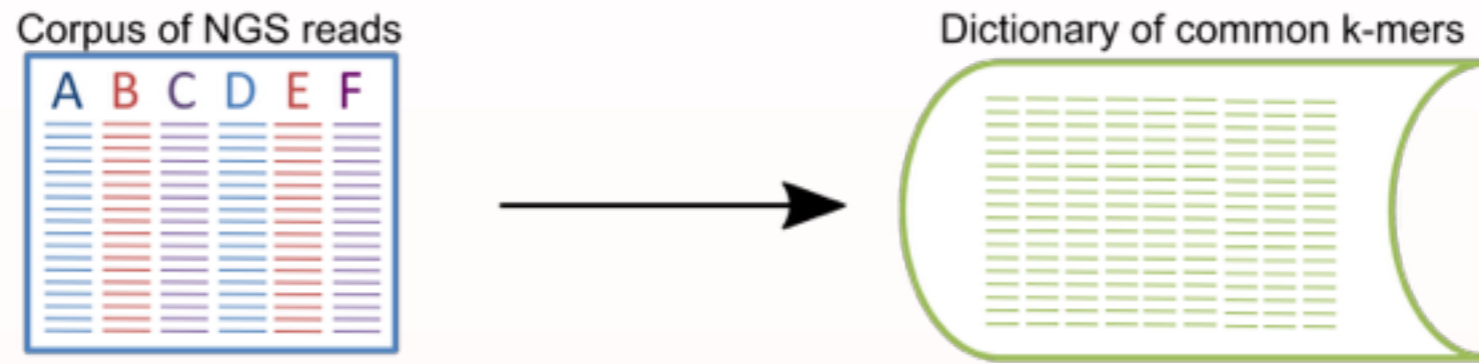


Compression

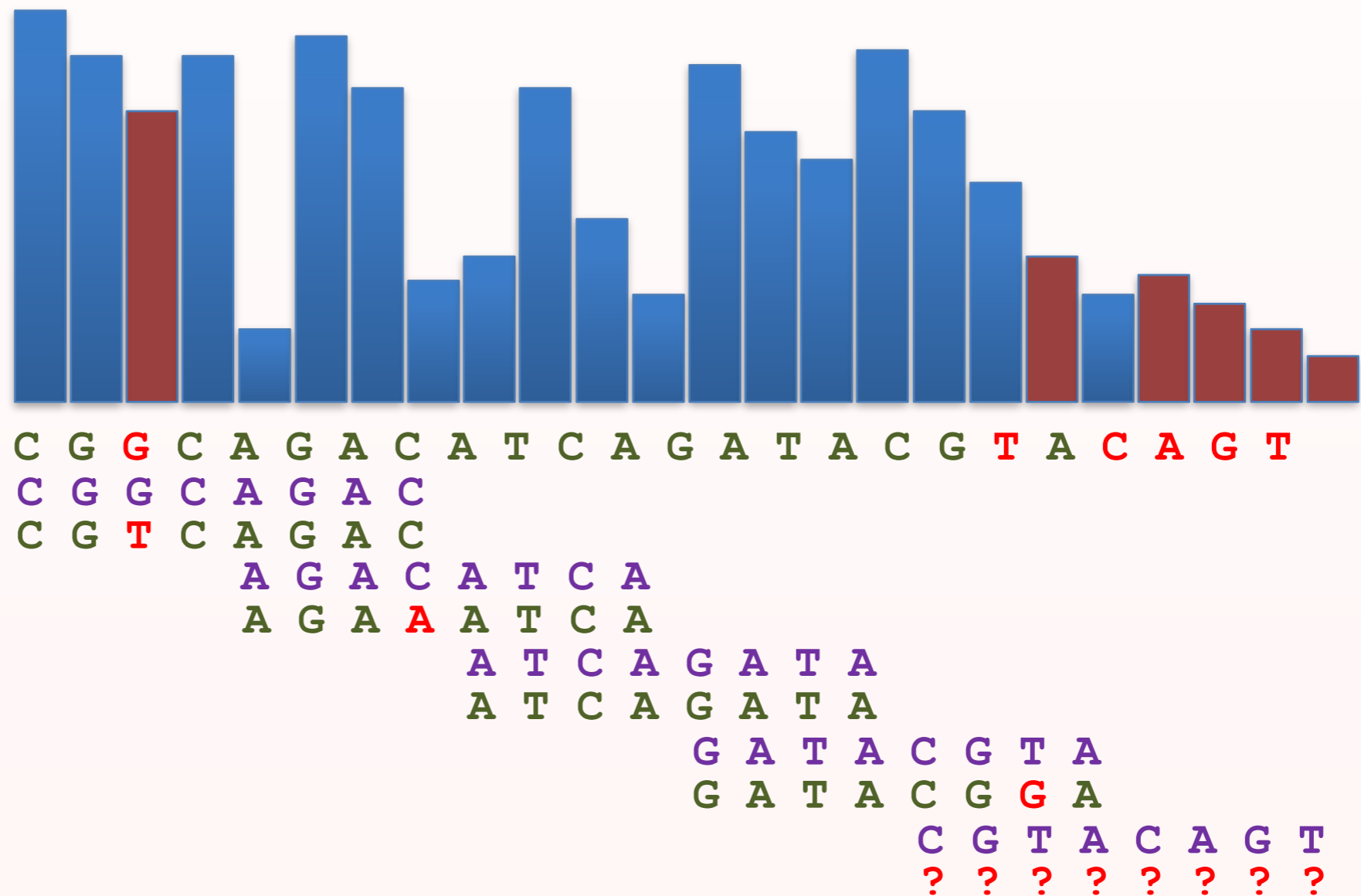


Quartz [Yu, et al. 2015]

Preprocessing



Compression



Approximate k-mer search

Naïve approaches

Approximate k-mer search

Naïve approaches

- Need to quickly find all $3k$ Hamming neighbors of each k-mer from the read in the dictionary to identify mis-matched bases.

Approximate k-mer search

Naïve approaches

- Need to quickly find all $3k$ Hamming neighbors of each k-mer from the read in the dictionary to identify mis-matched bases.
- Naïve approaches
 - Sorted list
 - Memory efficient but binary search is CPU and cache inefficient
 - Hash tables
 - Faster CPU-wise, but memory and cache inefficient

Approximate k-mer search

Locality sensitive hashing

An (R, cR, P_1, P_2) -sensitive LSH family F of hash functions $h : M \rightarrow S$ is defined if

$\forall p, q \in M$, a uniformly random $h \in F$ satisfies:

if $\|p - q\| \leq R$ then $\mathbb{P}(h(p) = h(q)) \geq P_1$

if $\|p - q\| \geq cR$ then $\mathbb{P}(h(p) = h(q)) \leq P_2$

Project k -mers onto random $\frac{k}{2}$ -mers, forming a $(1, c, \frac{1}{2}, 2^{c-1})$ -sensitive family of hash functions under the Hamming metric

Approximate k-mer search

Double hashing for fun and profit

Approximate k-mer search

Double hashing for fun and profit

- Notice that each $h : M \rightarrow S$ comes with an orthogonal projection $h' : M \rightarrow S$

Approximate k-mer search

Double hashing for fun and profit

- Notice that each $h : M \rightarrow S$ comes with an orthogonal projection $h' : M \rightarrow S$
- Also, if $\|p, q\| \leq 1$ then by counting, at least one of $h(p) = h(q)$ or $h'(p) = h'(q)$ must hold.

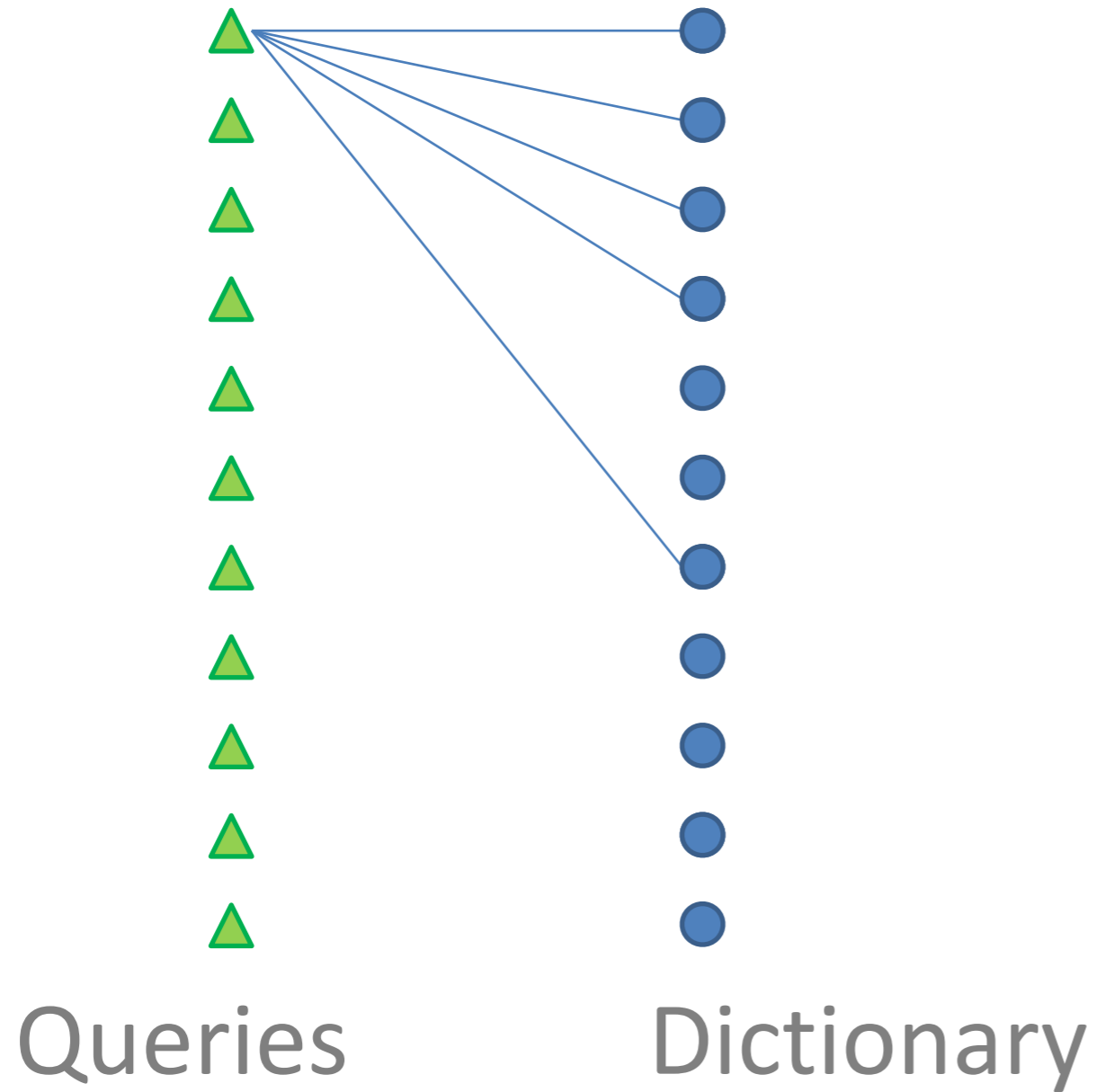
Approximate k-mer search

Double hashing for fun and profit

- Notice that each $h : M \rightarrow S$ comes with an orthogonal projection $h' : M \rightarrow S$
- Also, if $\|p, q\| \leq 1$ then by counting, at least one of $h(p) = h(q)$ or $h'(p) = h'(q)$ must hold.
- Thus, by double hashing, all Hamming neighbors of a k -mer can be found by looking in just two hash buckets.
 - Better cache-efficiency
 - Cheating by carefully choosing the projection and sorting the buckets gives also processor-efficiency.

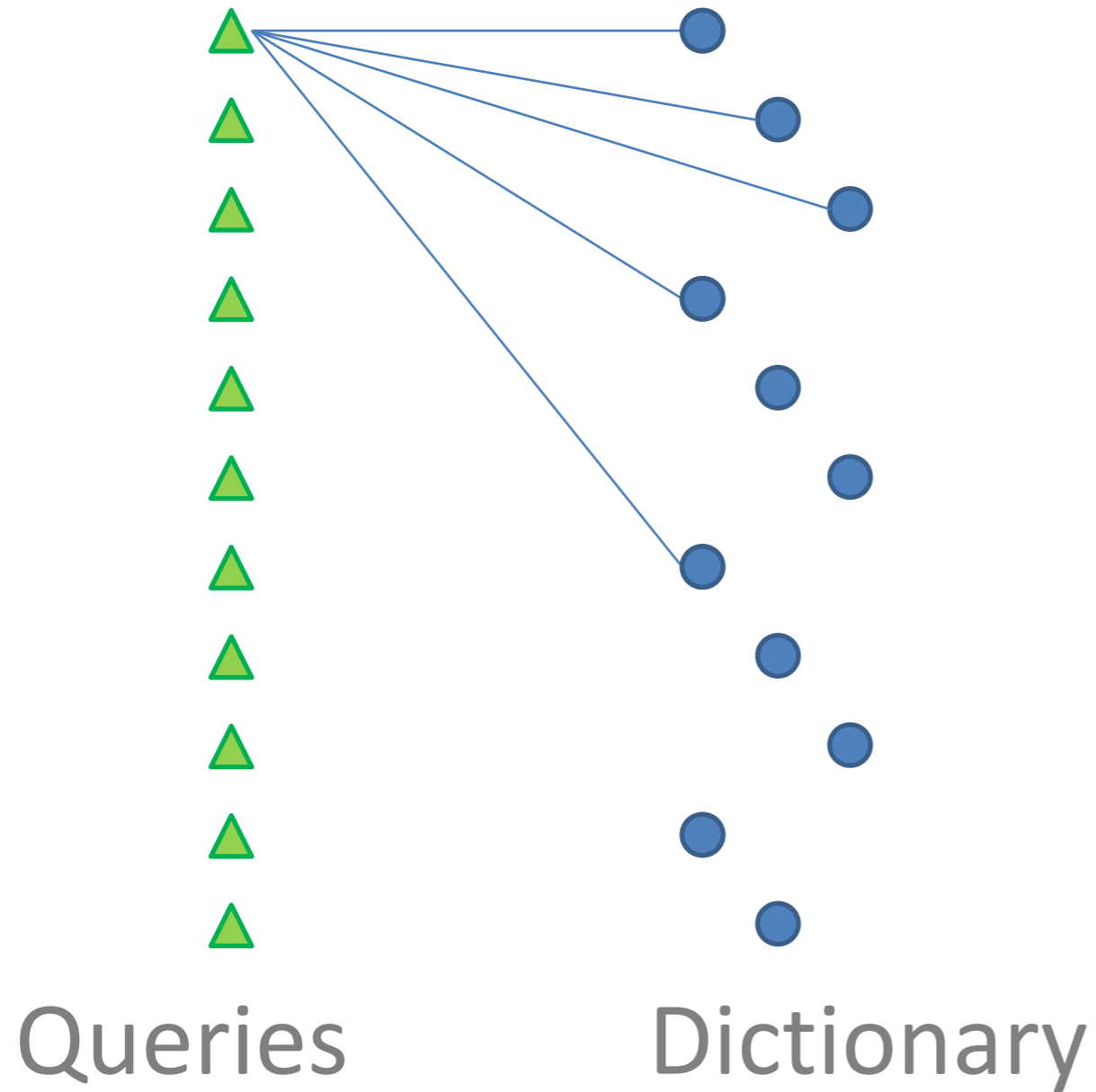
Approximate k-mer search

Fast retrieval of Hamming neighbors



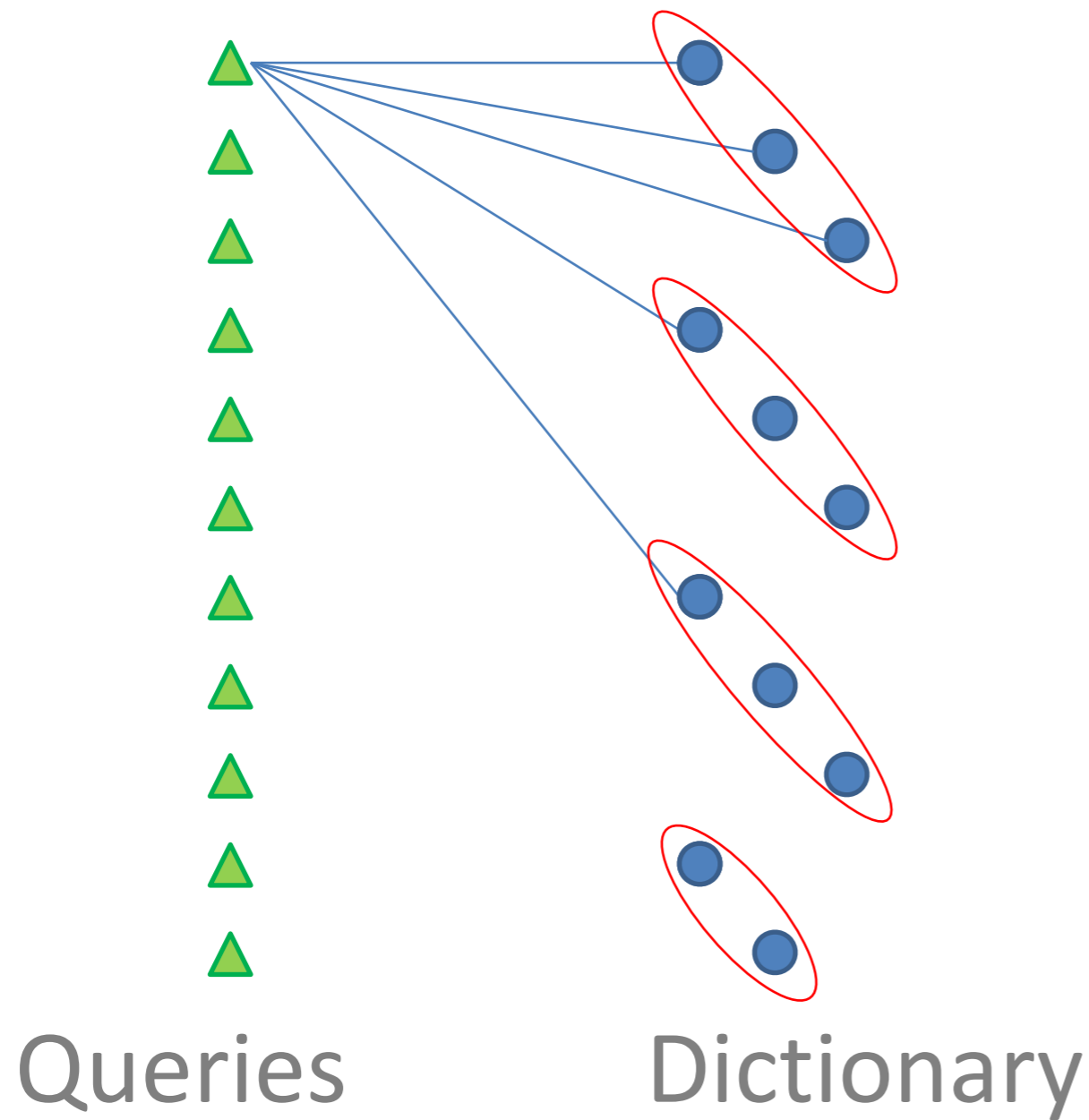
Approximate k-mer search

Fast retrieval of Hamming neighbors



Approximate k-mer search

Fast retrieval of Hamming neighbors

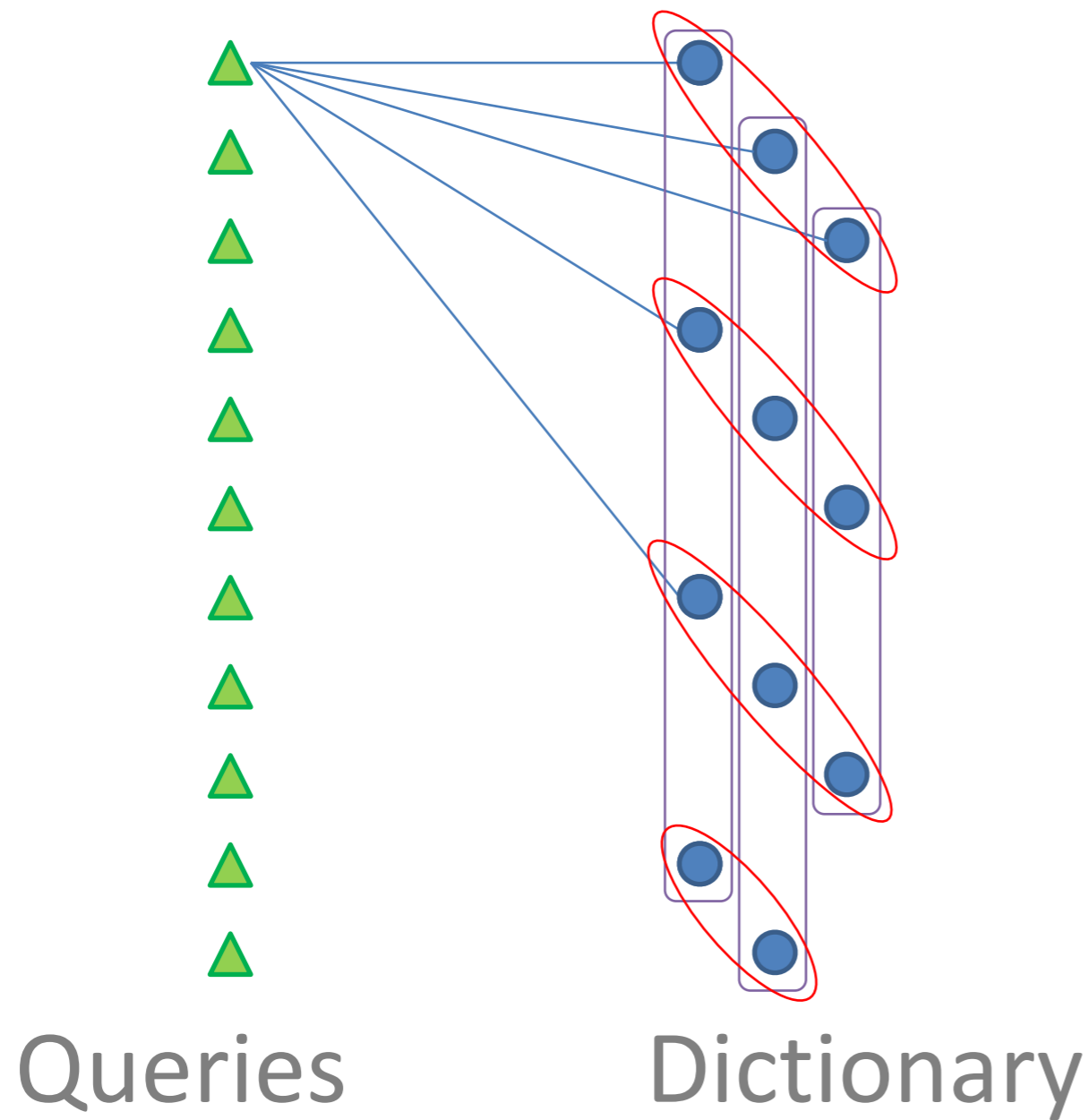


Clustered by front half of k-mer

C	G	G	C	A	G	A	C
C	G	G	C	C	G	A	C
C	G	G	C	A	G	T	C

Approximate k-mer search

Fast retrieval of Hamming neighbors



Clustered by front half of k-mer

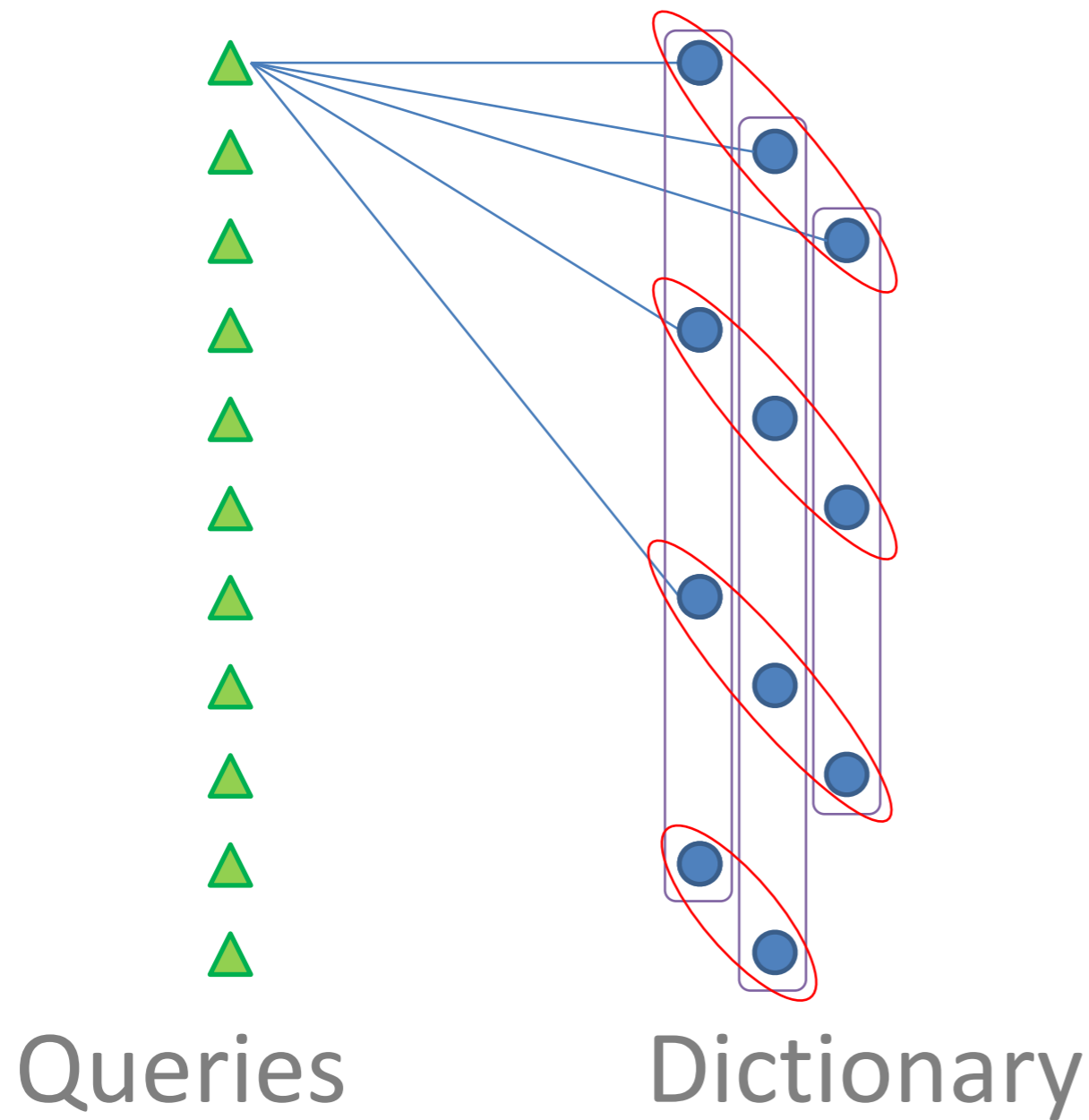
```
C G G C A G A C  
C G G C C G A C  
C G G C A G T C
```

Clustered by back half of k-mer

```
A G G C A G A C  
G G G C A G A C  
C C G C A G A C  
T A T A A G A C
```

Approximate k-mer search

Fast retrieval of Hamming neighbors



Clustered by front half of k-mer

C G G C A G A C
C G G C C G A C
C G G C A G T C

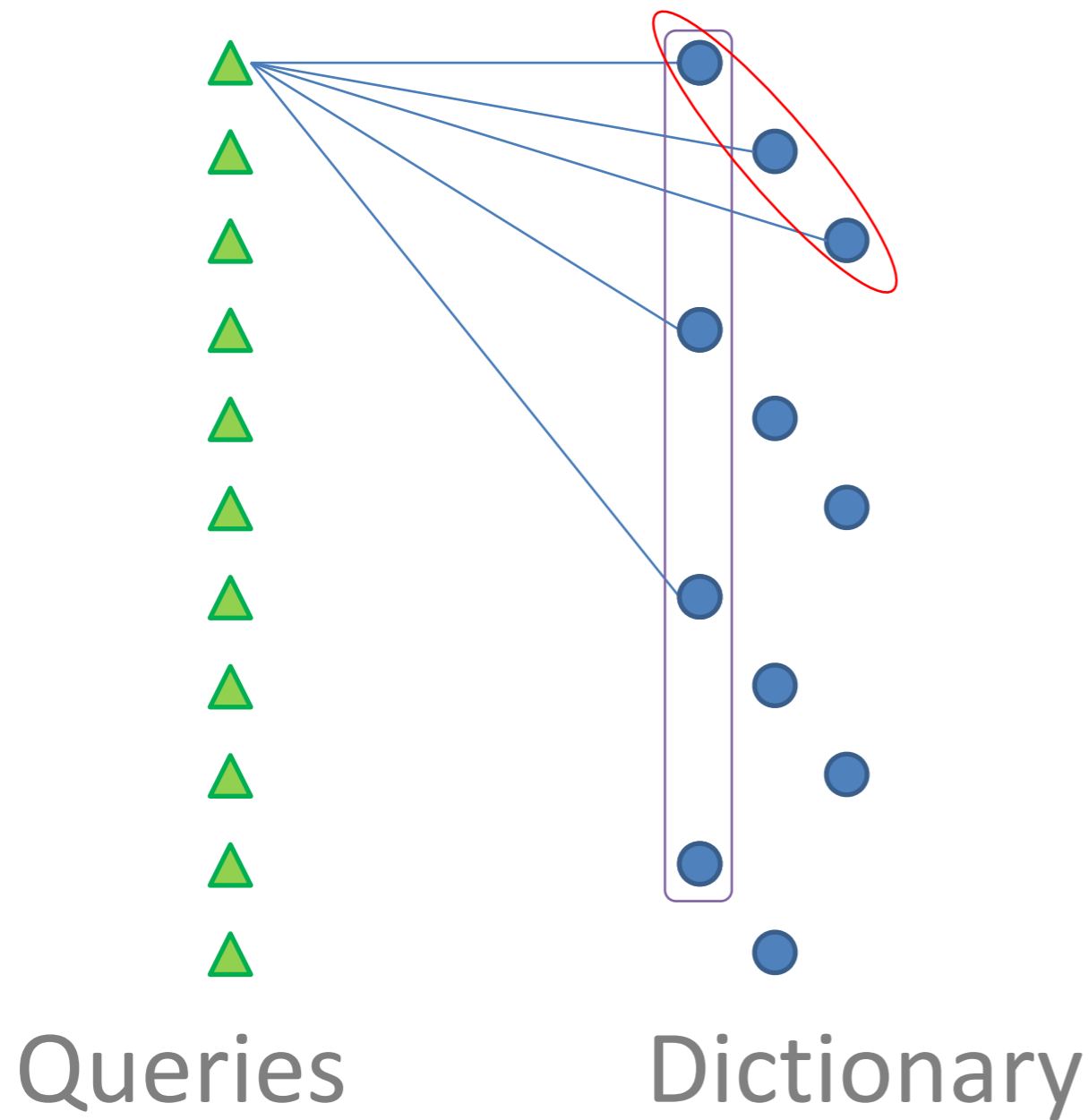
Clustered by back half of k-mer

A G G C A G A C
G G G C A G A C
C C G C A G A C
T A T A A G A C

▲ = C G G C A G A C

Approximate k-mer search

Fast retrieval of Hamming neighbors



Clustered by front half of k-mer

```
C G G C A G A C
C G G C C G A C
C G G C A G T C
```

Clustered by back half of k-mer

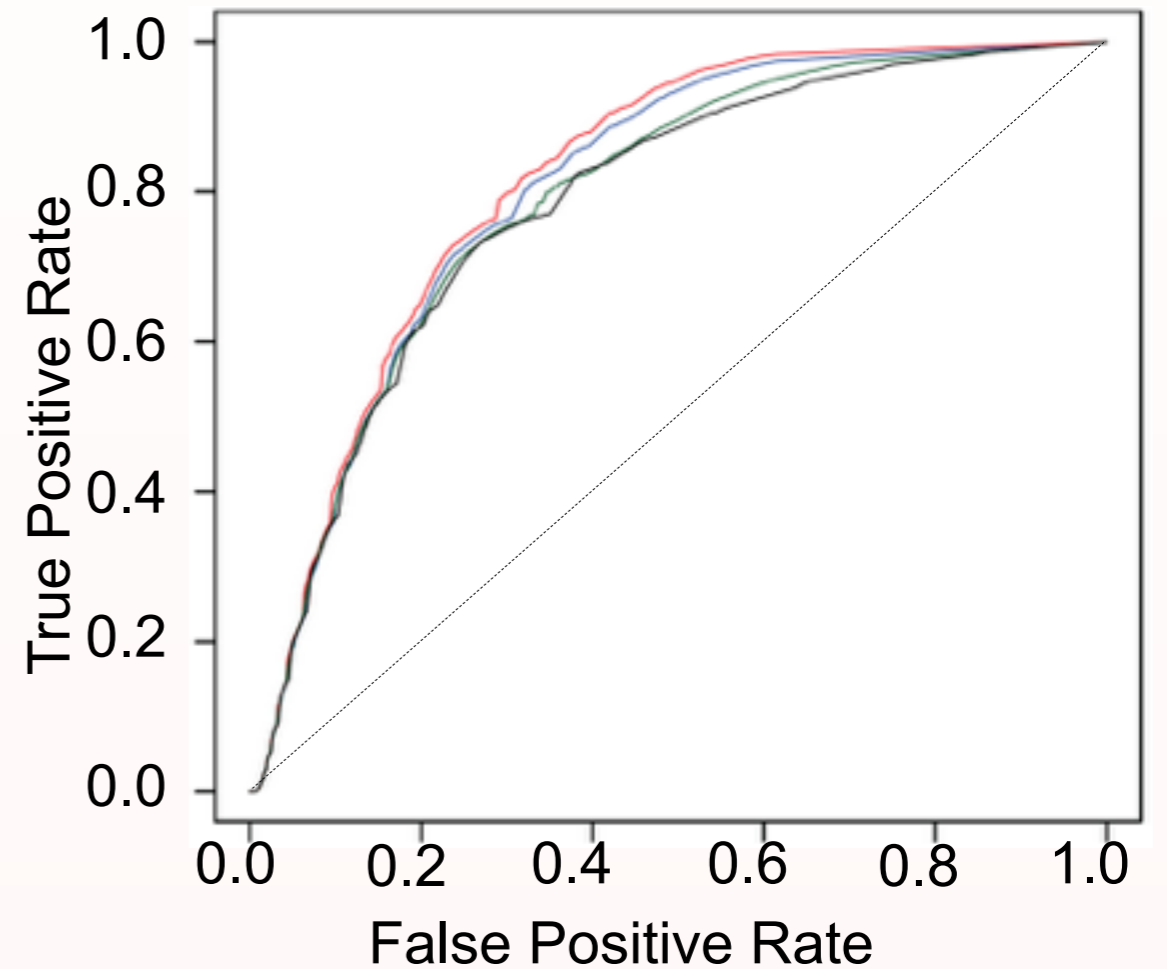
```
A G G C A G A C
G G G C A G A C
C C G C A G A C
T A T A A G A C
```

▲ = C G G C A G A C

Result Highlights

Comparison with other methods

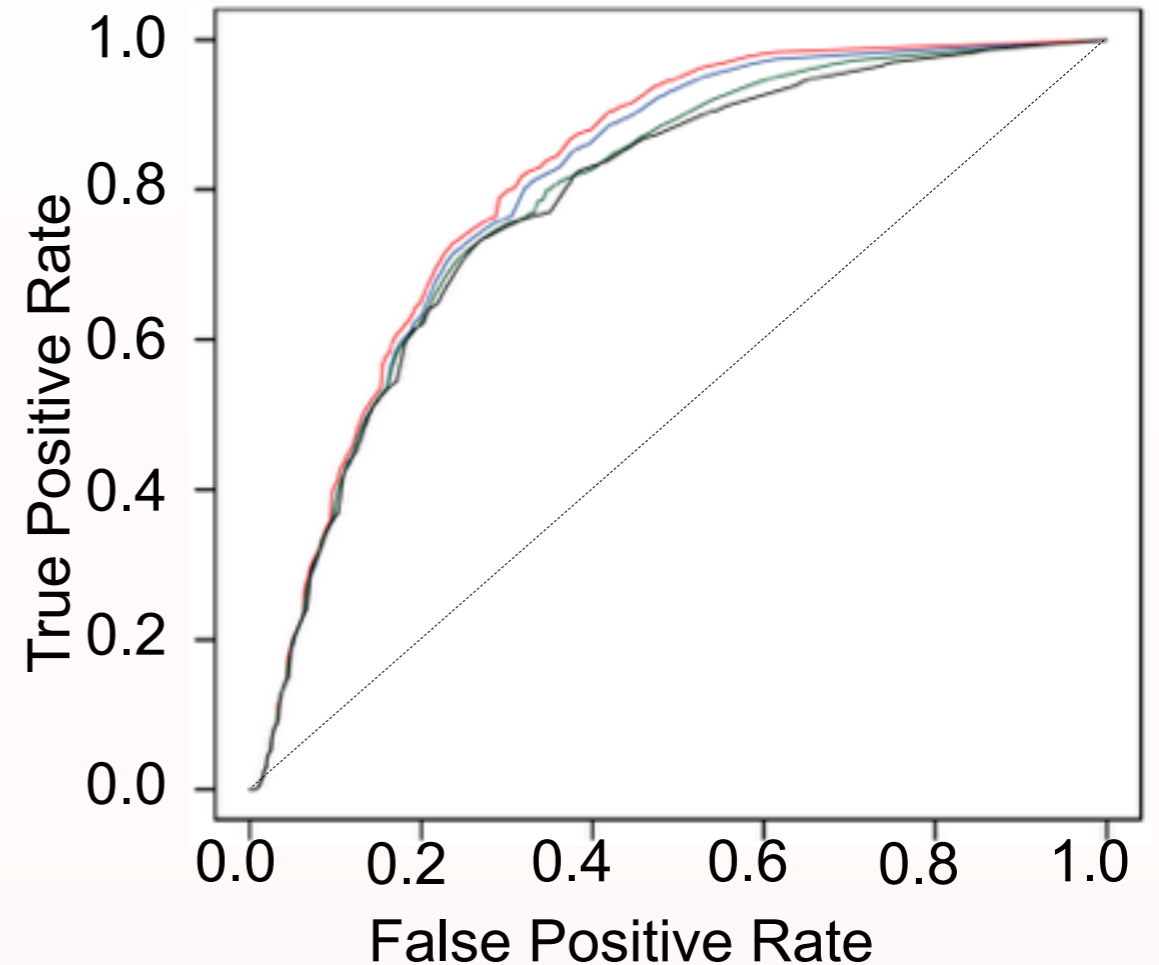
Method	Bits/Q	Time (s)	Area Under ROC Curve
Uncompressed	8	N/A	0.8254
Quartz	0.3564	2,696	0.8288
QualComp	0.5940	33,316	0.8053
Janin et al.	0.5376	164,702	0.8019



Result Highlights

Comparison with other methods

Method	Bits/Q	Time (s)	Area Under ROC Curve
Uncompressed	8	N/A	0.8254
Quartz	0.3564	2,696	0.8288
QualComp	0.5940	33,316	0.8053
Janin et al.	0.5376	164,702	0.8019



Quartz is orders of magnitude faster

A word on k-mers

MINCE: 15-mers optimal (8-mer labels!)

using labels as search heuristic; don't want too many

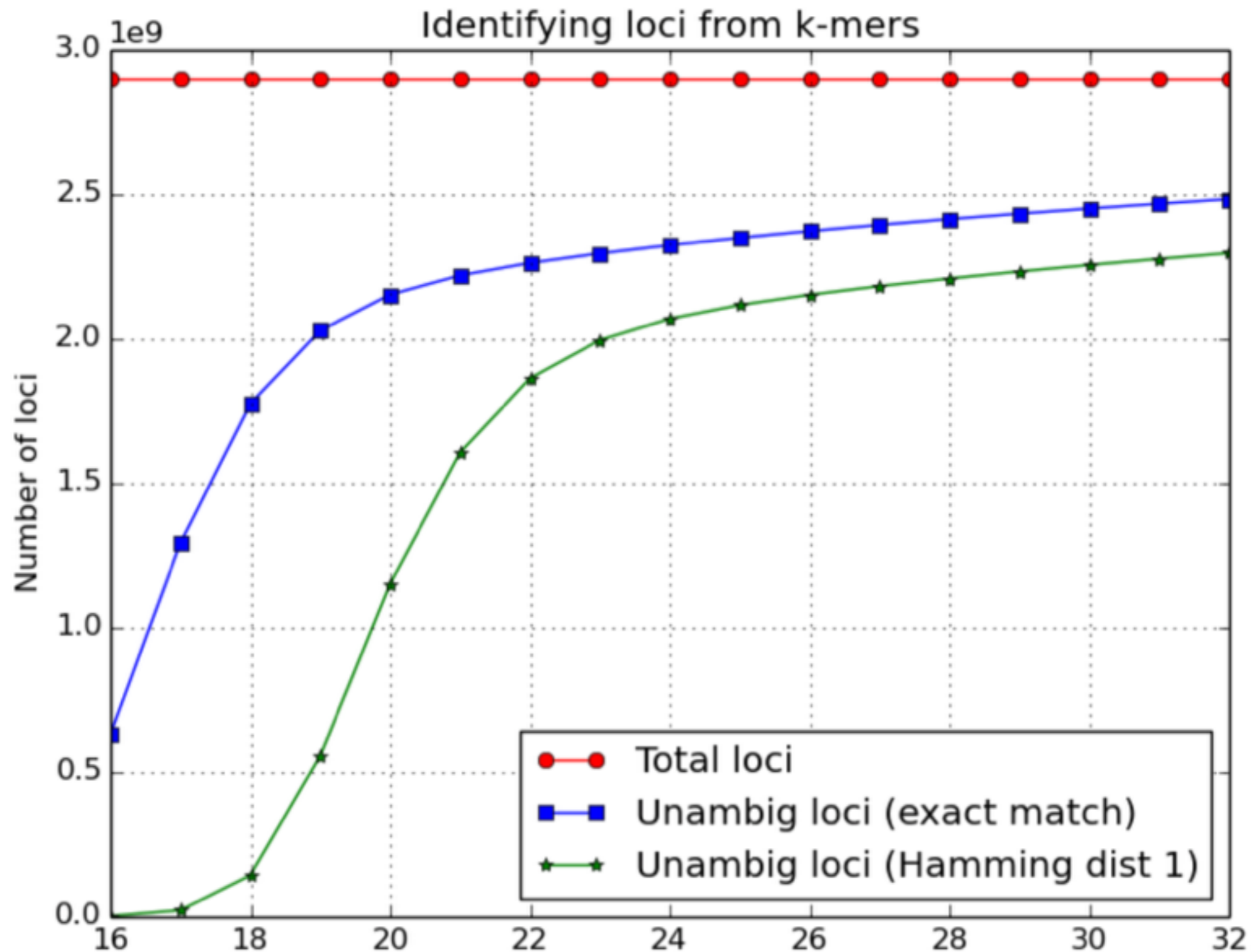
Quartz: 32-mers optimal

in a genome, most 8-mers may exist; most 32-mers will not

and they tend to be unique

	k=16	k=32
unique	21.8%	85.7%
unique at Hamming 1	0.0008%	79.3%

Uniqueness of k-mers (hg19)



Why not longer k-mers?

Want to ensure at most 1 sequencing error per k -mer

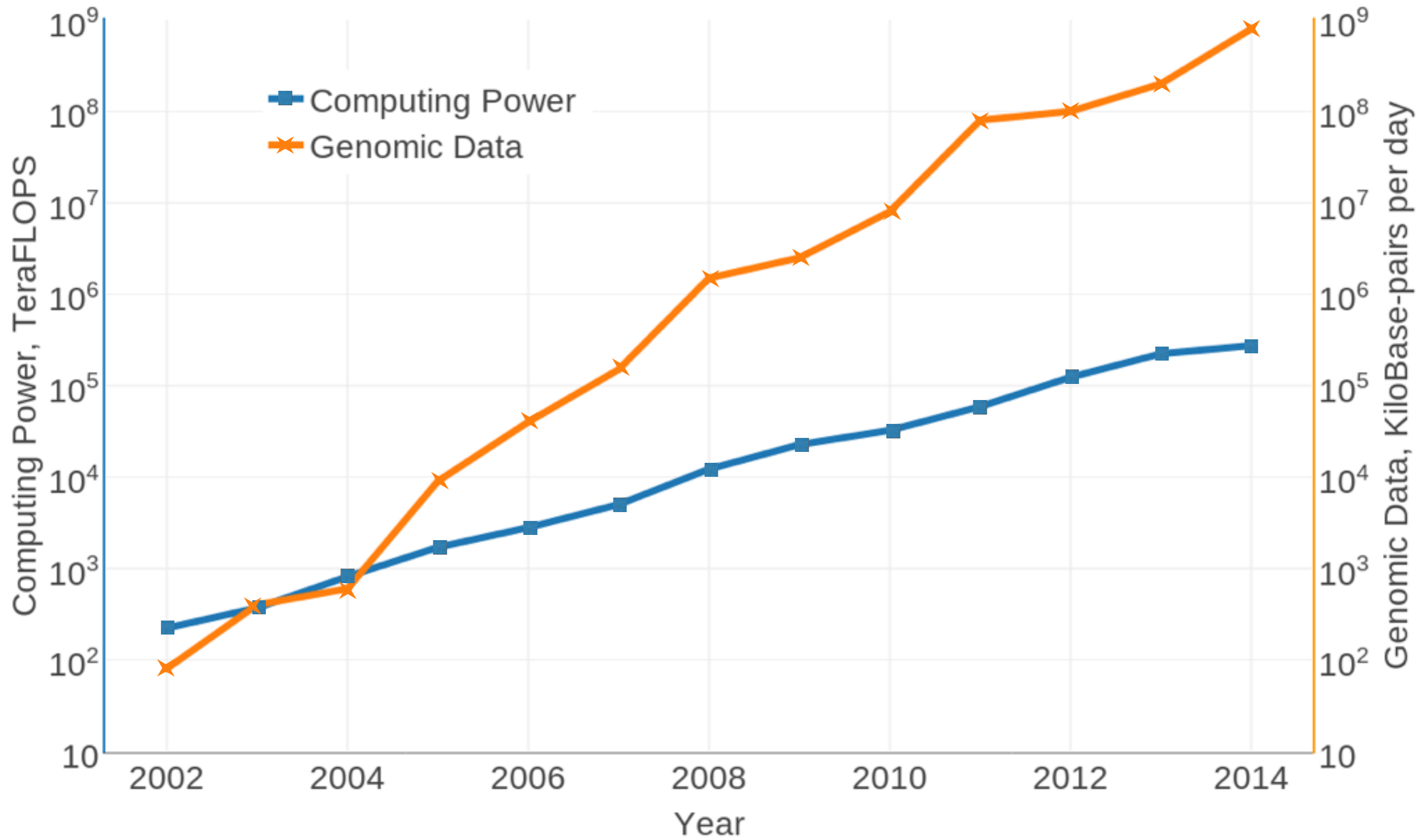
Assume independence of errors
Error rate of p

Likelihood of l errors $\binom{k}{l} (1-p)^{k-l} p^l$

1% error rate
 $k=32 \quad l \geq 2 \quad 2\%$
 $k=64 \quad l \geq 2 \quad 13\%$

2% error rate
 $k=32 \quad l \geq 2 \quad 13\%$
 $k=64 \quad l \geq 2 \quad 36\%$

Compression for speed



Compressive genomics

caBLAST [Loh, et al. 2012]

caBLASTP [Daniels, et al. 2013]

BLAST uses seed-and-extend

but must extend on many fruitless seeds

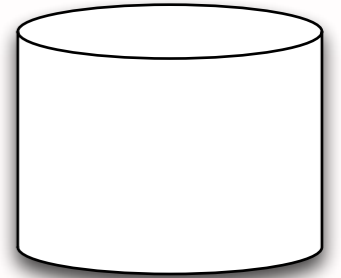
Use compression to reduce the search space

How does this compression work?

How does this compression work?

STAQEPKSAEDSLRARD

Coarse Database

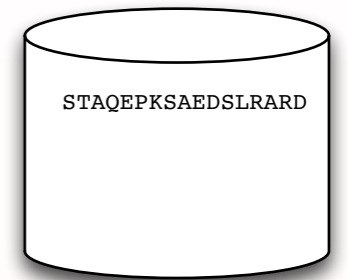


How does this compression work?

STAQEPKSAEDSLRARD



Coarse Database



How does this compression work?

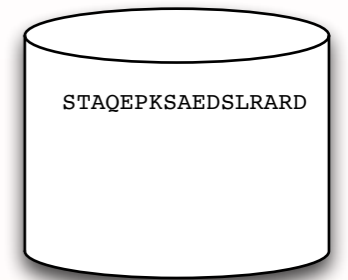
Seed Table



STAQEPKSAEDSLRARD



Coarse Database



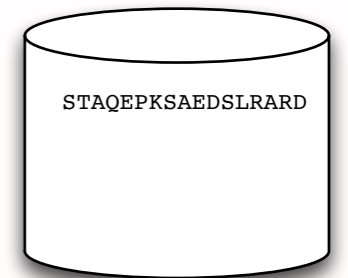
How does this compression work?

Seed Table



STAQEPKSAEDSLRARD

Coarse Database



LQSTAQEPKSAEQRDSVNARDRQRNVIIAQE

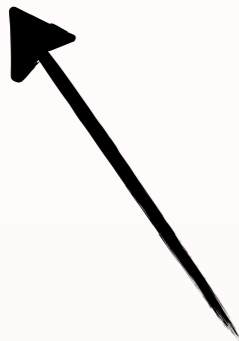
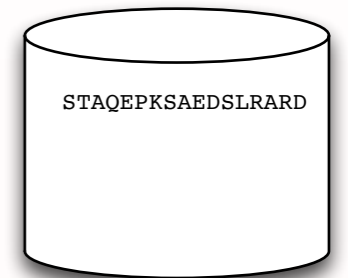
How does this compression work?

Seed Table



STAQEPKSAEDSLRARD

Coarse Database



LQSTAQEPKSAEQRDSVNARDRQRNVIIAQE

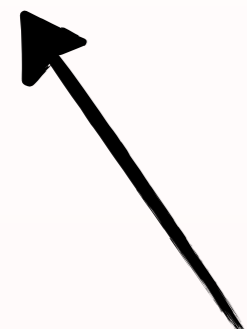
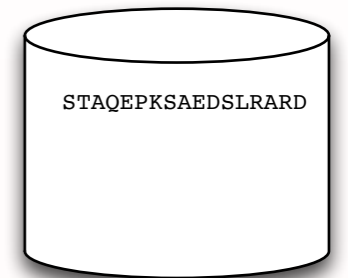
How does this compression work?

Seed Table



STAQEPKSAEDSLRARD

Coarse Database



LQSTAQEPKSAEQRDSVNARDRQRNVIIAQE

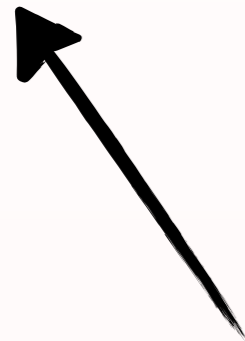
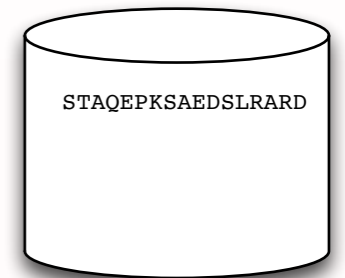
How does this compression work?

Seed Table



STAQEPKSAEDSLRARD

Coarse Database



LQSTAQEPKSAEQRDSVNARDRQRNVIIAQE

How does this compression work?

Seed Table

...
STAP STAQ
STAR STAS
...

STAQEPKSAE DSLRARD

Coarse Database

STAQEPKSAEDSLRARD

LQSTAQEPKSAEQRDSVNARDRQRNVIIAQE

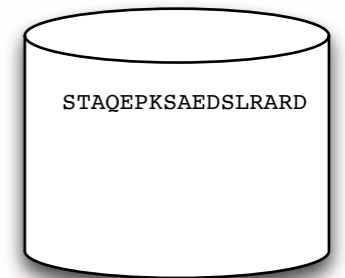
How does this compression work?

Seed Table



STAQEPKSAE DSLRARD

Coarse Database



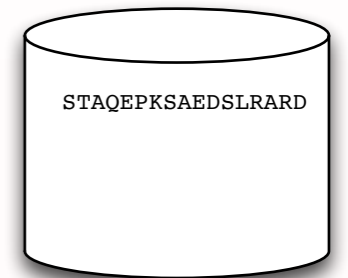
LQSTAQEPKSAEQRDSVNARDRQRNVIIAQE

How does this compression work?

Seed Table



Coarse Database



STAQEPKSAE DSLRARD

-- DSLRARD

QRDSVNARD

LQ STAQEPKSAE QRDSVNARD RQRNVIIAQE

How does this compression work?

Seed Table

...
STAP STAQ
STAR STAS
...

Coarse Database

STAQEPKSAEDSLRARD

STAQEPKSAE DSLRARD

-- DSLRARD

QRDSVNARD

LQ STAQEPKSAE QRDSVNARD RQRNVIIAQE

LQ

QR

VN

How does this compression work?

Seed Table

...
STAP STAQ
STAR STAS
...

STAQEPKSAE DSLRARD

-- DSLRARD

QRDSVNARD

LQ STAQEPKSAE QRDSVNARD RQRNVIIAQE

Coarse Database

STAQEPKSAEDSLRARD

How does this compression work?

Seed Table



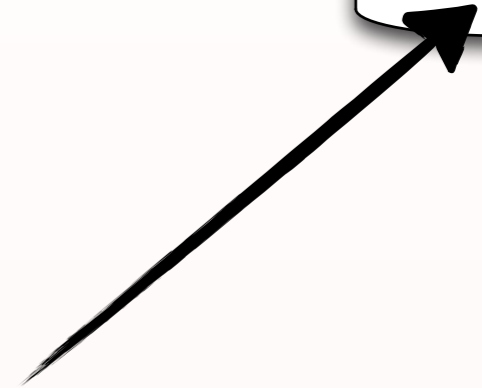
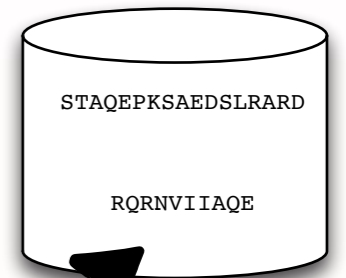
STAQEPKSAE DSLRARD

-- DSLRARD

QRDSVNARD

LQ STAQEPKSAE QRDSVNARD RQRNVIIAQE

Coarse Database



How does this compression work?

Seed Table



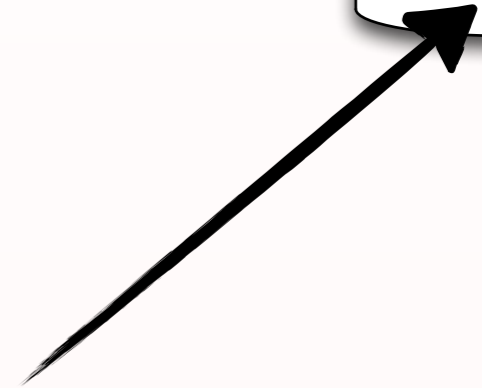
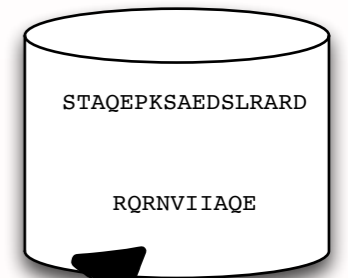
STAQEPKSAE DSLRARD

-- DSLRARD

QRDSVNARD

LQ STAQEPKSAE QRDSVNARD RQRNVIIAQE

Coarse Database



Lossless compression!

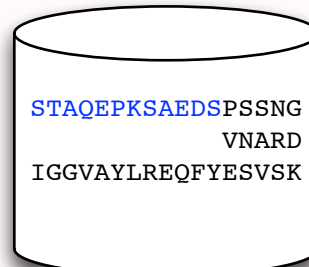
How does compressed search work?

How does compressed search work?

Query Sequence

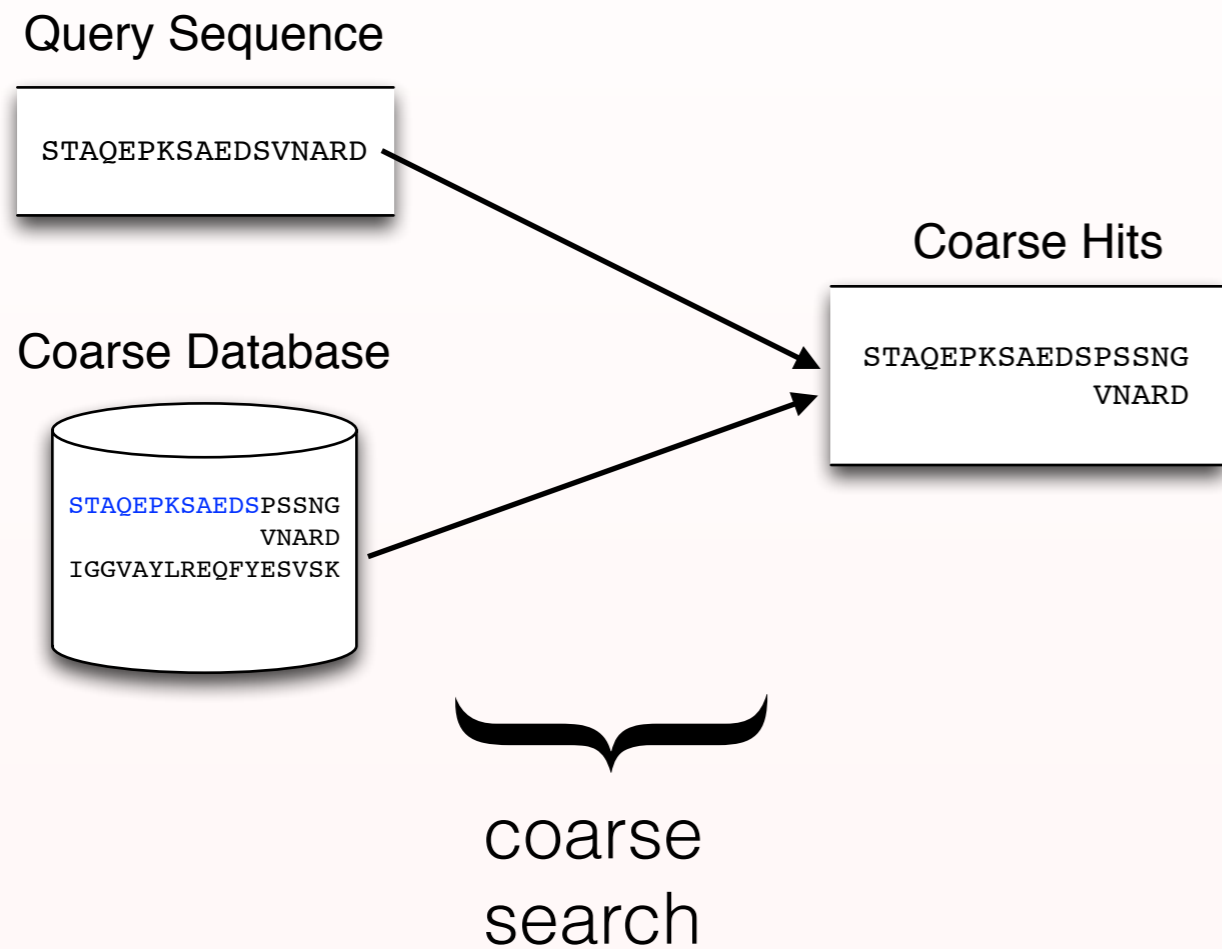
STAQEPKSAEDSVNARD

Coarse Database

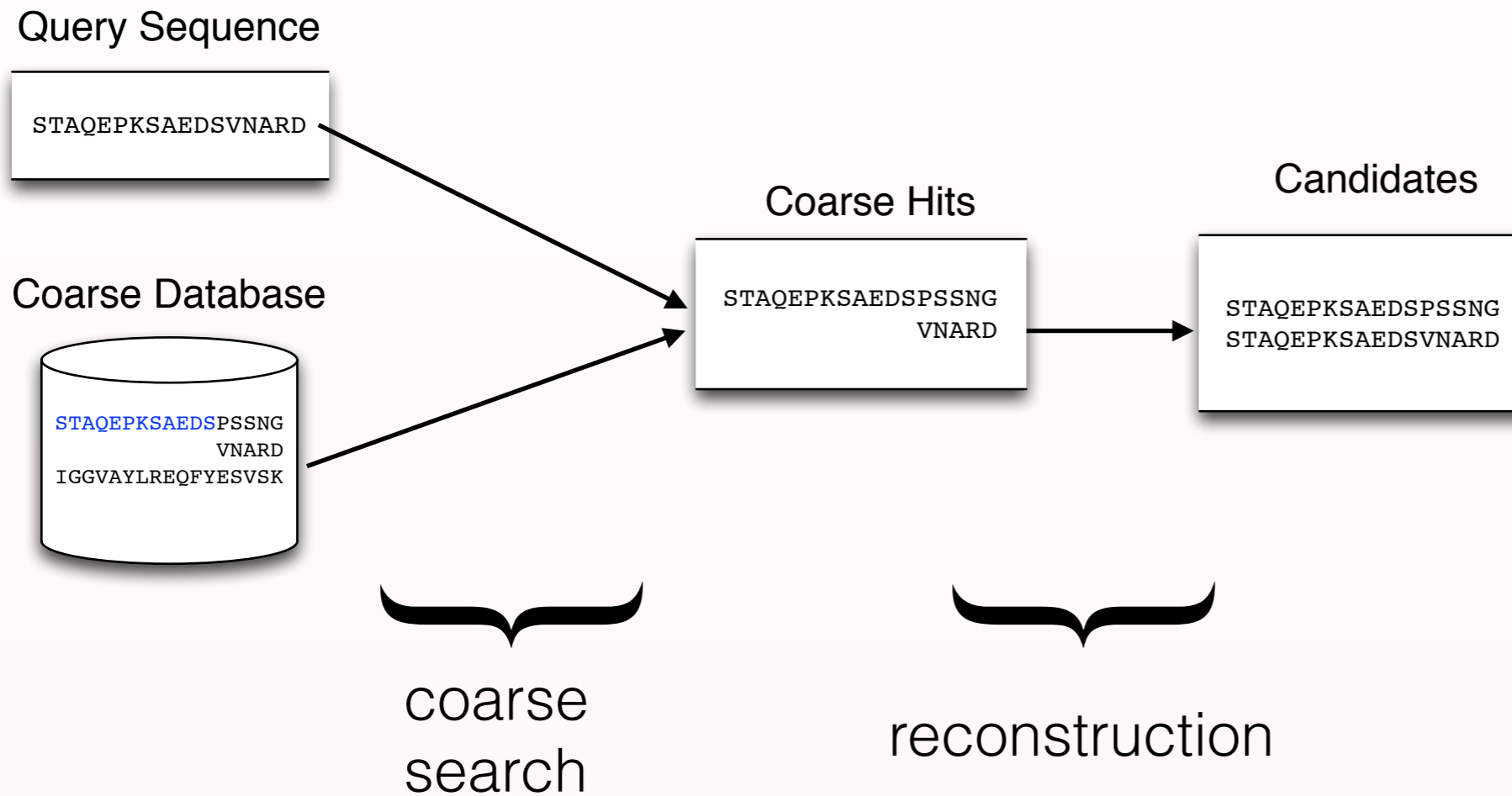


STAQEPKSAEDSPSSNG
VNARD
IGGVAYLREQFYESVSK

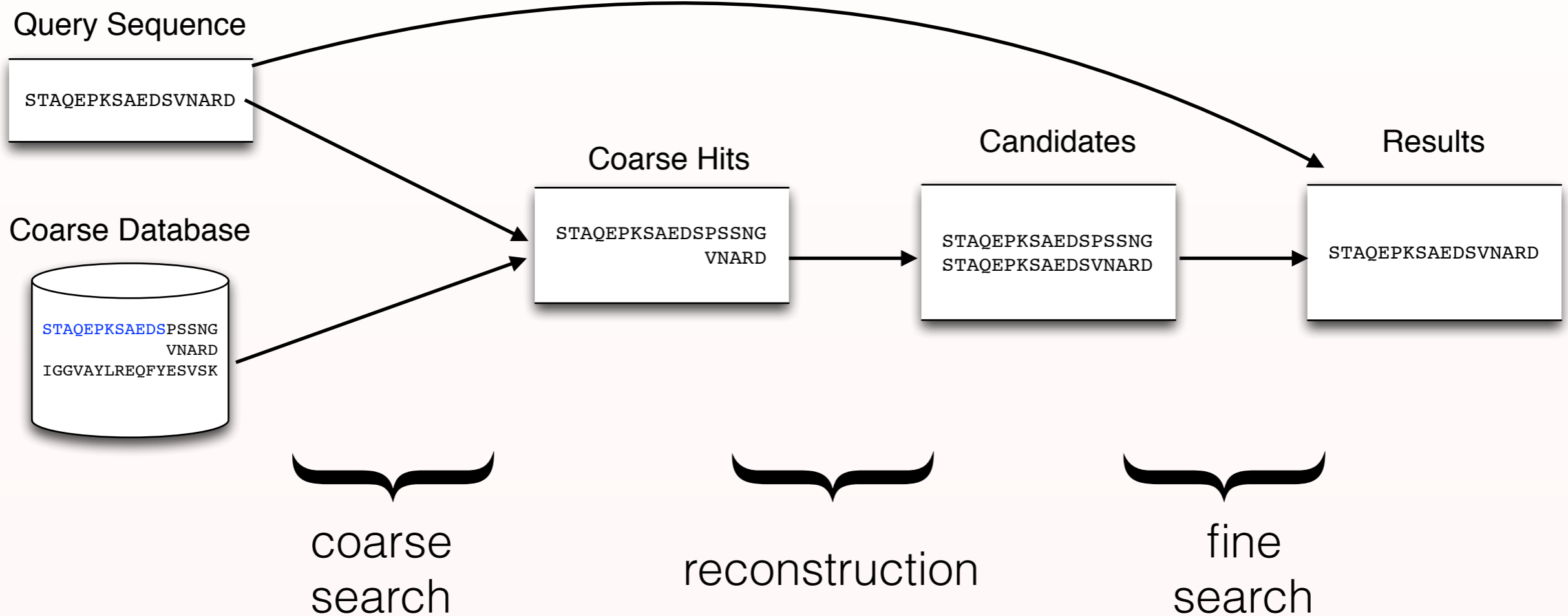
How does compressed search work?



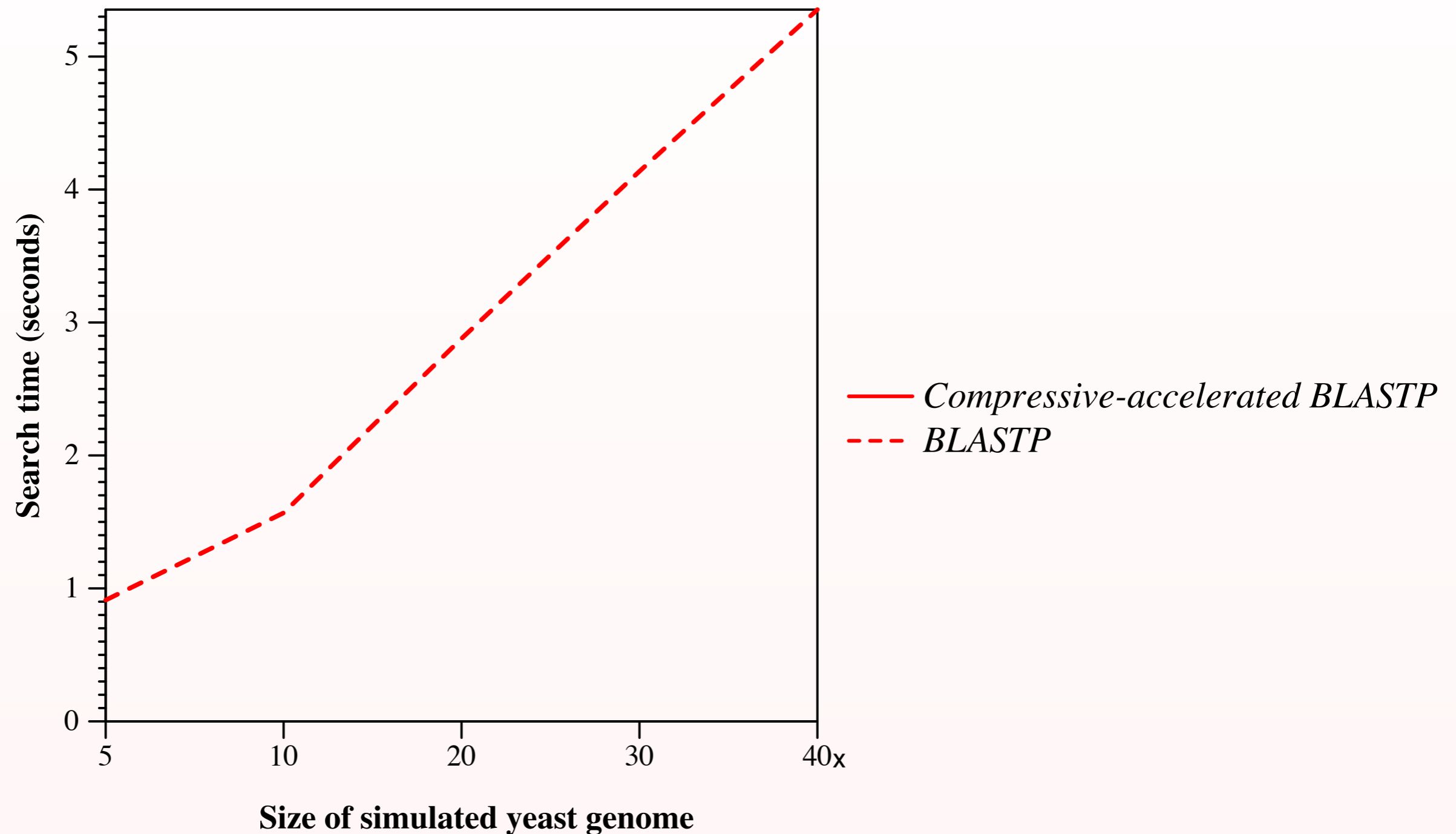
How does compressed search work?



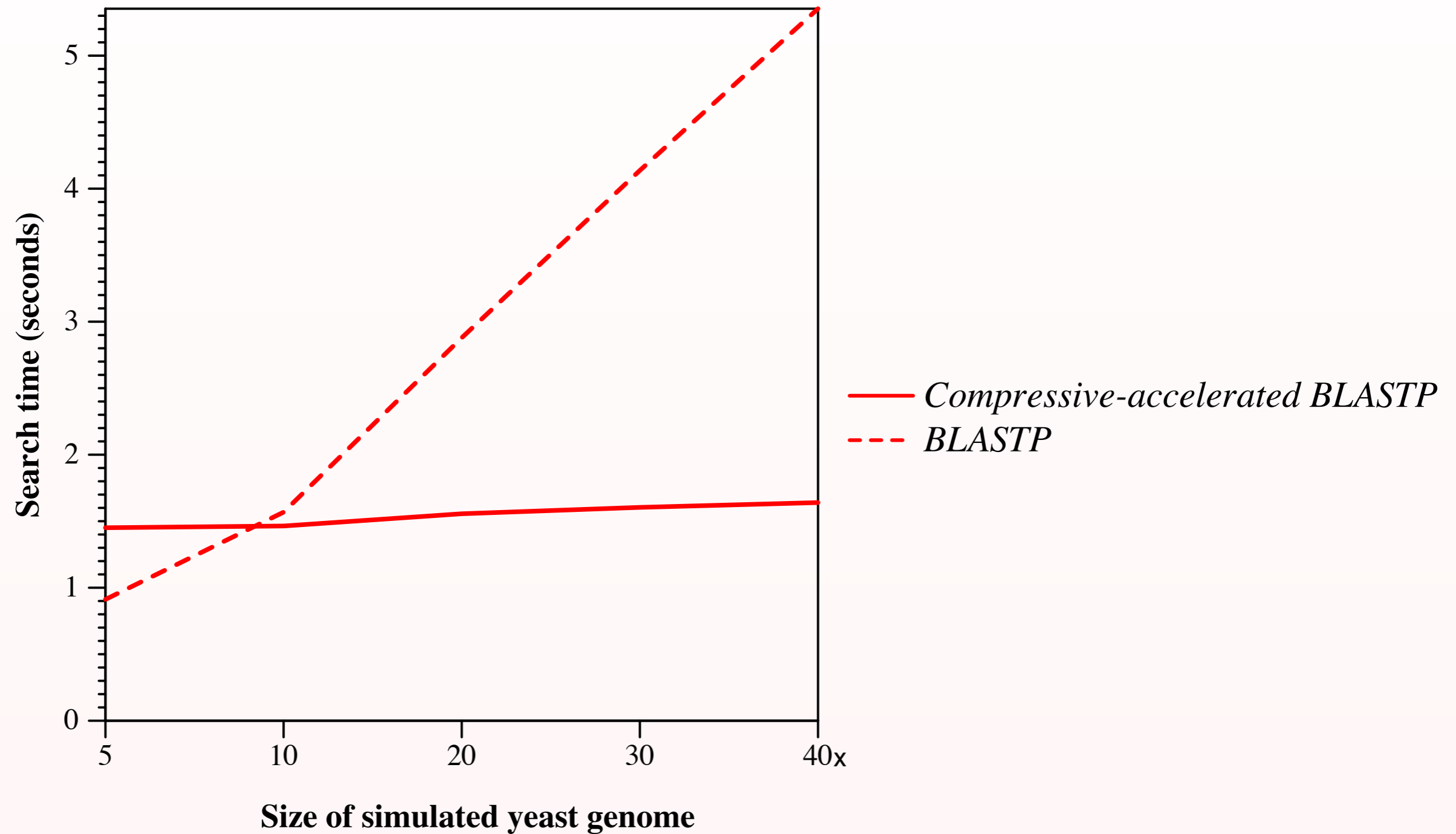
How does compressed search work?



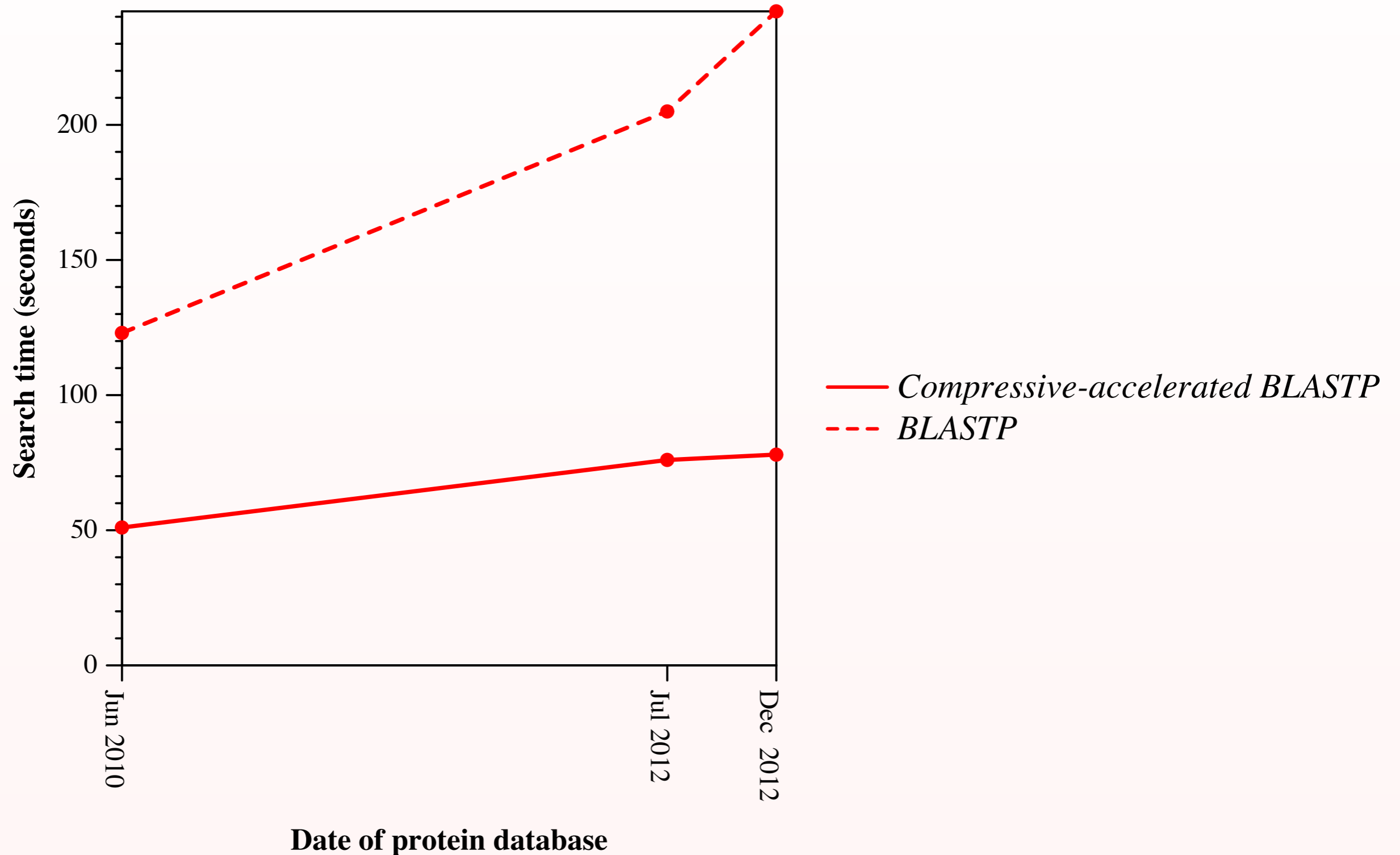
Simulated data growth



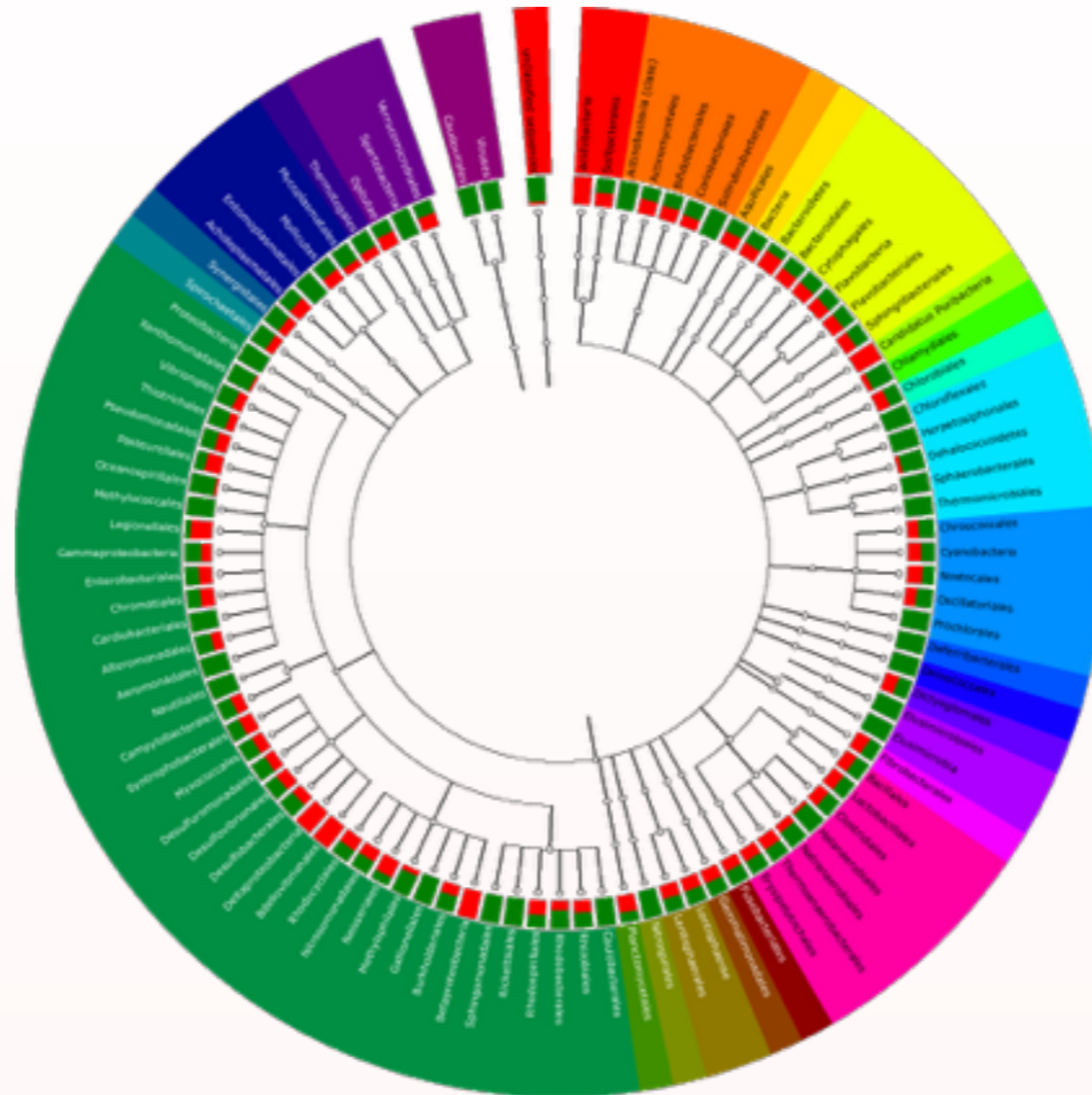
Simulated data growth



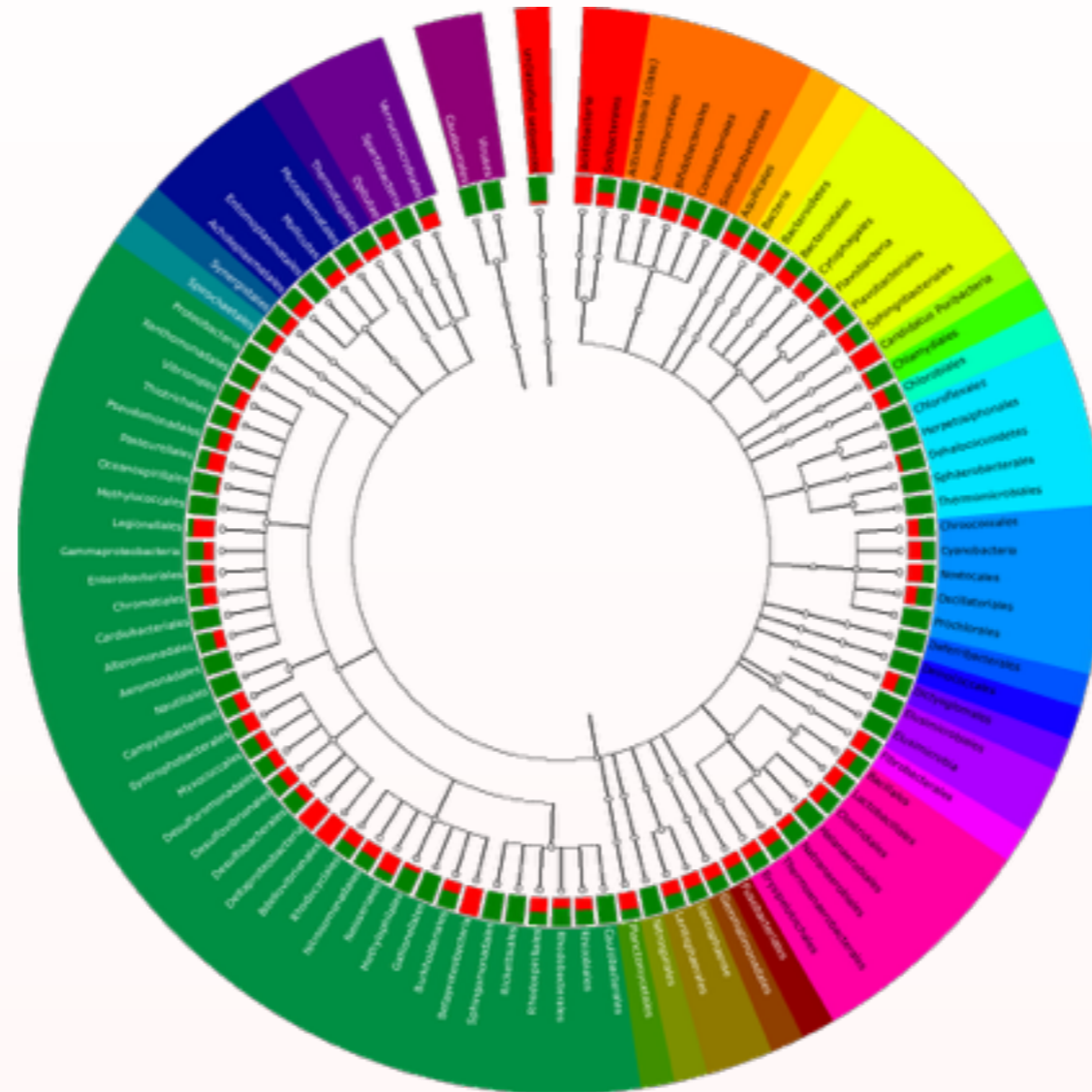
Real data growth on a protein database



Metagenomics



Metagenomics



Problem: Given reads from a microbiome,
match each read with similar proteins.

BLASTX [Altschul, et al. 1990]
RapSearch2 [Zhao, et al. 2012]
DIAMOND [Buchfink, et al. 2015]

Metagenomics

BLASTX [Altschul, et al. 1990]
RapSearch2 [Zhao, et al. 2012]
DIAMOND [Buchfink, et al. 2015]

Metagenomics

Metagenomic Reads

```
>read1
ACGTGGCTATCAACTCGCTAACTAA
>read2
ACGTGGCTATCAACTCGCTAACTAA
>read3
ACGTGGCTATCAACTCGCTAACTAT
...
>readk
TCGTCGAACTACATTACATTTACAG
>readk+1
TCGTCGAACTACATTACAAATACAG
...
>readm
GCTCGGACTATATATAGGCCTAGAA
...
```

BLASTX [Altschul, et al. 1990]
RapSearch2 [Zhao, et al. 2012]
DIAMOND [Buchfink, et al. 2015]

Metagenomics

Metagenomic Reads

```
>read1
ACGTGGCTATCAACTCGCTAACTAA
>read2
ACGTGGCTATCAACTCGCTAACTAA
>read3
ACGTGGCTATCAACTCGCTAACTAT
...
>readk
TCGTCGAACTACATTACATTTACAG
>readk+1
TCGTCGAACTACATTACAAATACAG
...
>readm
GCTCGGACTATATATAGGCCTAGAA
...
```

Translated ORFs

```
>read1-1
TWLSTR
>read1-2
RGYQLAN
>read1-3
VAINSLT
>read1-r1
LVSELIAT
>read1-r2
LAS
>read1-r3
RVDSH
>read2-1
SSNYITFT
>read2-2
RRTTLHLQ
>read2-3
VELHYIY
>read2-r1
CSST
>read2-r2
CKCNVRRR
>read2-r3
VNV
...
```

BLASTX [Altschul, et al. 1990]
RapSearch2 [Zhao, et al. 2012]
DIAMOND [Buchfink, et al. 2015]

Metagenomics

Metagenomic Reads

```
>read1
ACGTGGCTATCAACTCGCTAACTAA
>read2
ACGTGGCTATCAACTCGCTAACTAA
>read3
ACGTGGCTATCAACTCGCTAACTAT
...
>readk
TCGTCGAACTACATTACATTTACAG
>readk+1
TCGTCGAACTACATTACAAATACAG
...
>readm
GCTCGGACTATATATAGGCCTAGAA
...
```

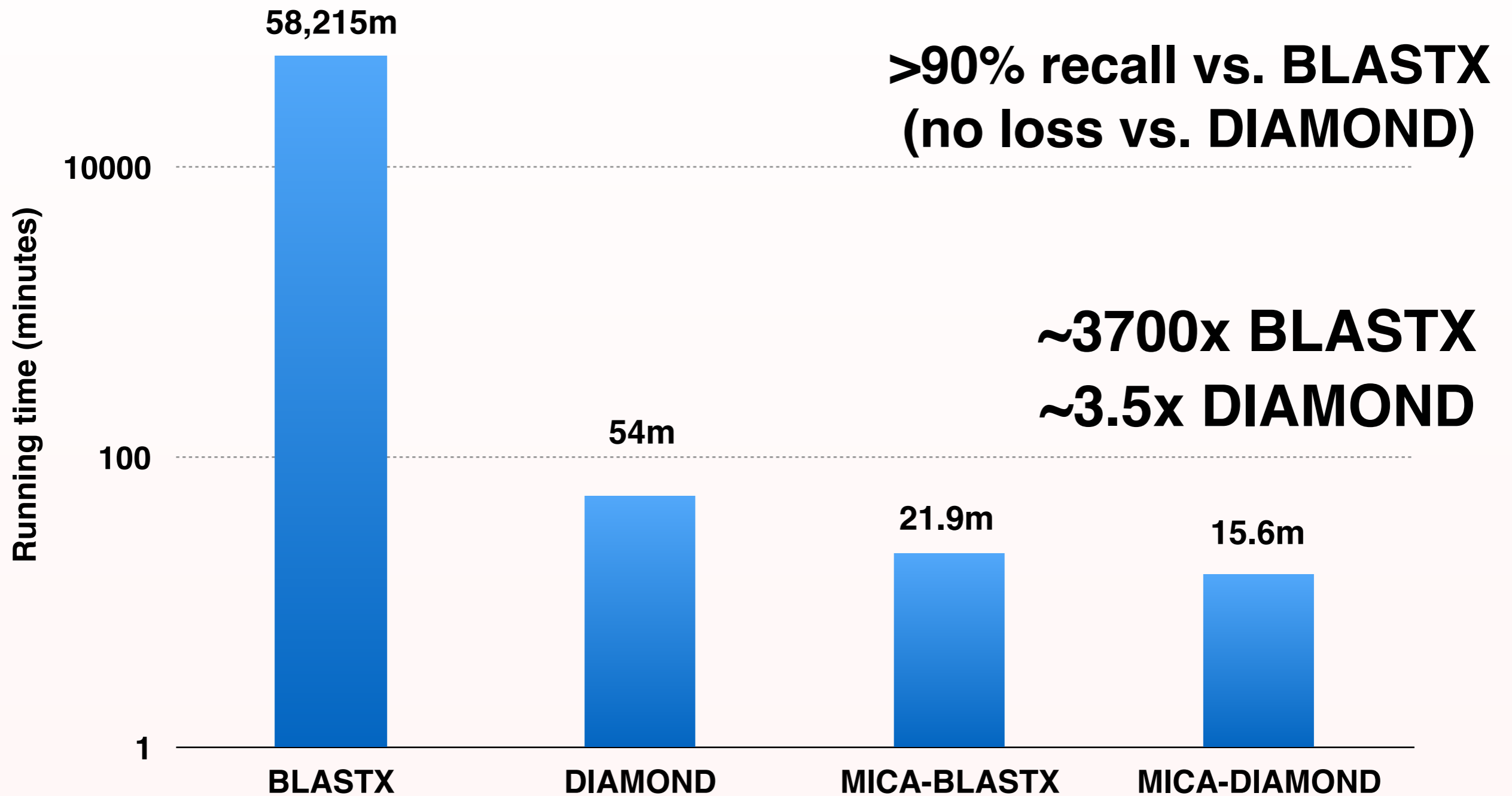
Translated ORFs

```
>read1-1
TWLSTR
>read1-2
RGYQLAN
>read1-3
VAINSLT
>read1-r1
LVSELIAT
>read1-r2
LAS
>read1-r3
RVDSH
>read2-1
SSNYITFT
>read2-2
RRTTLHLQ
>read2-3
VELHYIY
>read2-r1
CSST
>read2-r2
CKCNVRR
>read2-r3
VNV
...
```

Protein Database (NR)

```
>protein 1
MRVLVINSGSSSIKYQLIEM
>protein 2
EILGKLEELKIITCHIGNGASVAAVKY
>protein 3
LKKLLESSGCRLVRYGNILIG
...
```

MICA [Yu, Daniels, et al. 2015]



American gut microbiome project NGS reads,
searching NCBI NR database

Compression = clustering

Compression = clustering

Cluster similar entities

Compression = clustering

Cluster similar entities

or sub-entities

Compression = clustering

Cluster similar entities

or sub-entities

Only store the common parts once

Compression = clustering

Cluster similar entities

or sub-entities

Only store the common parts once

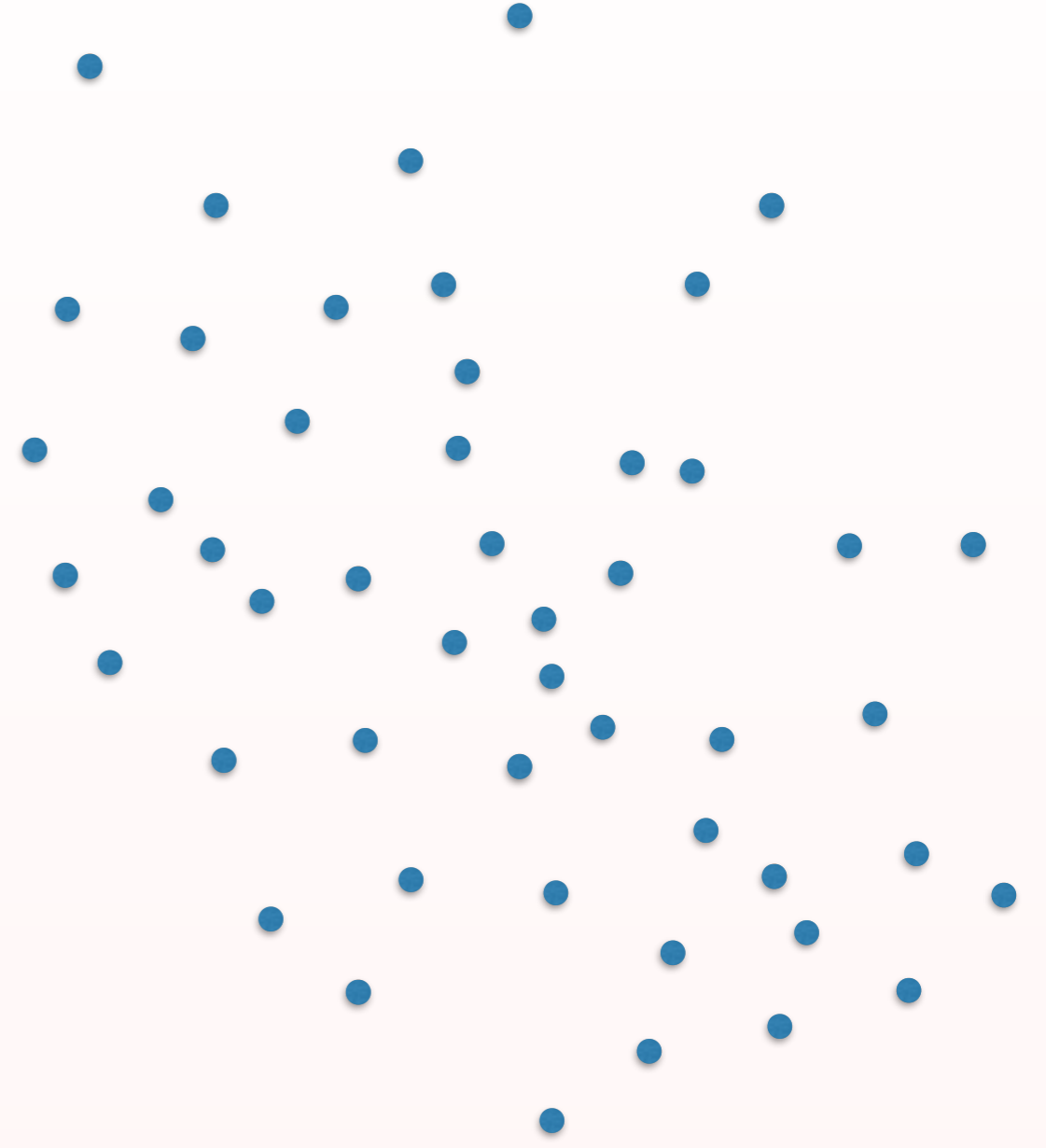
Only analyze the common parts once



Queries



Database

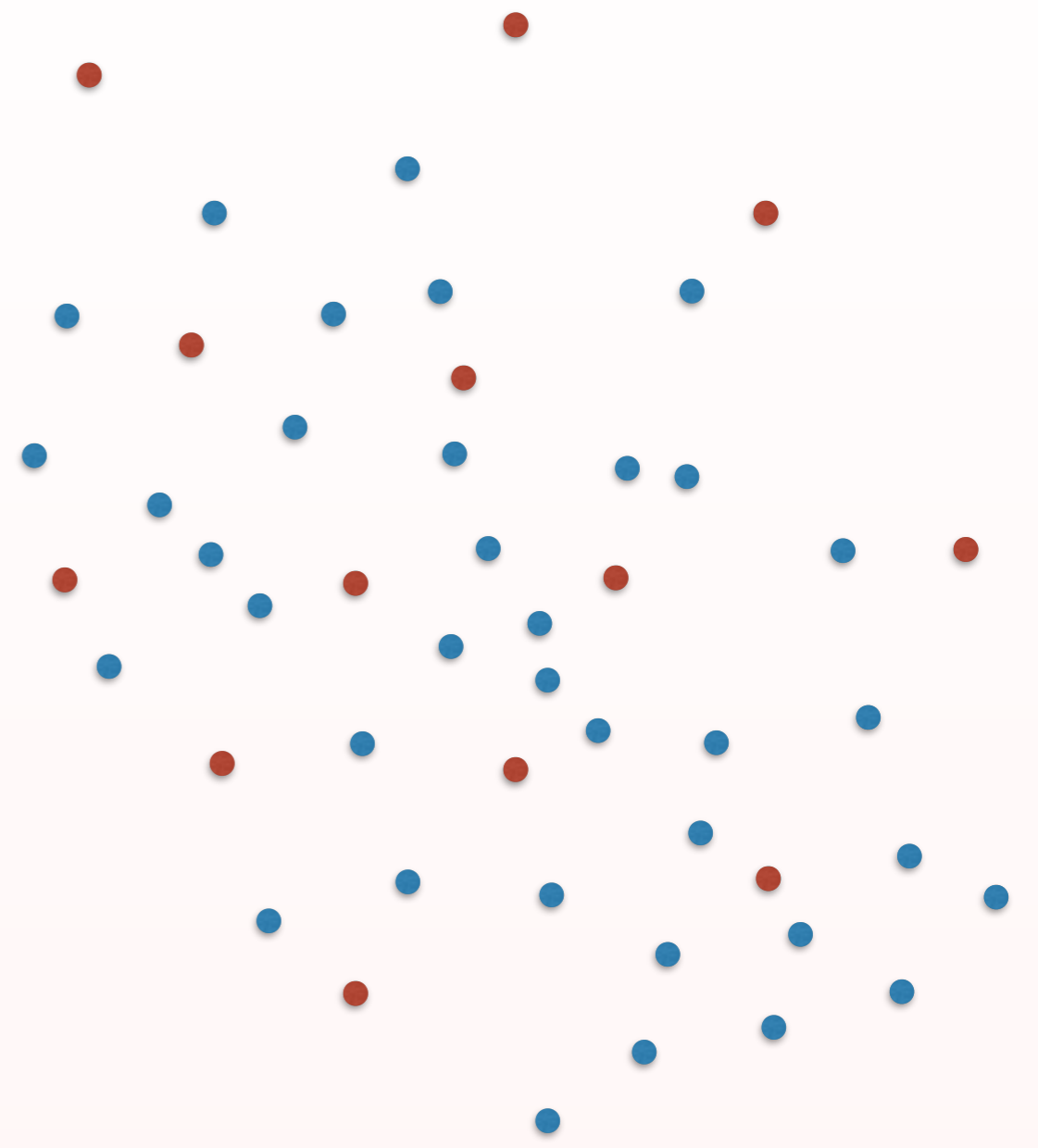




Queries



Database

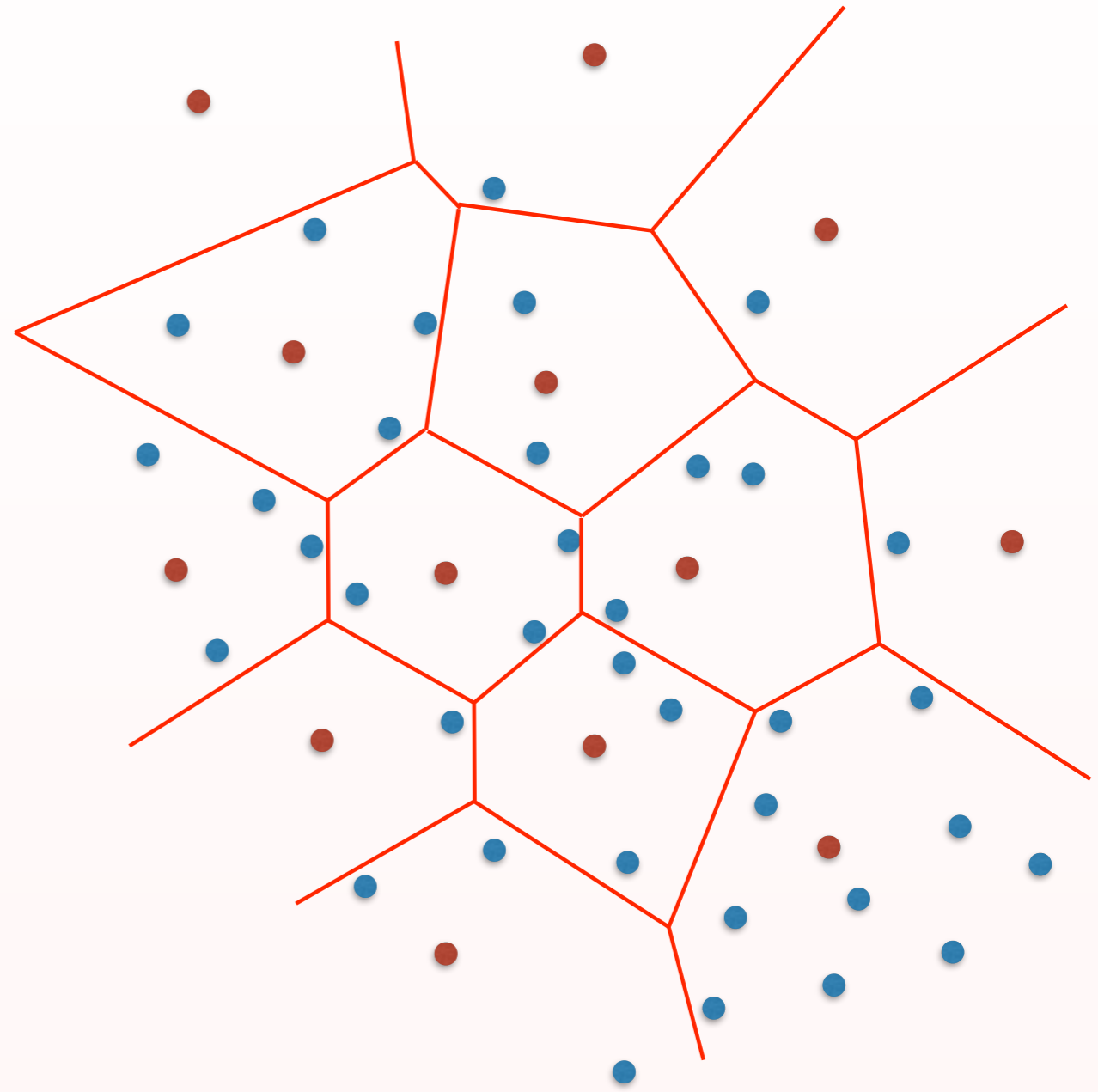




Queries



Database

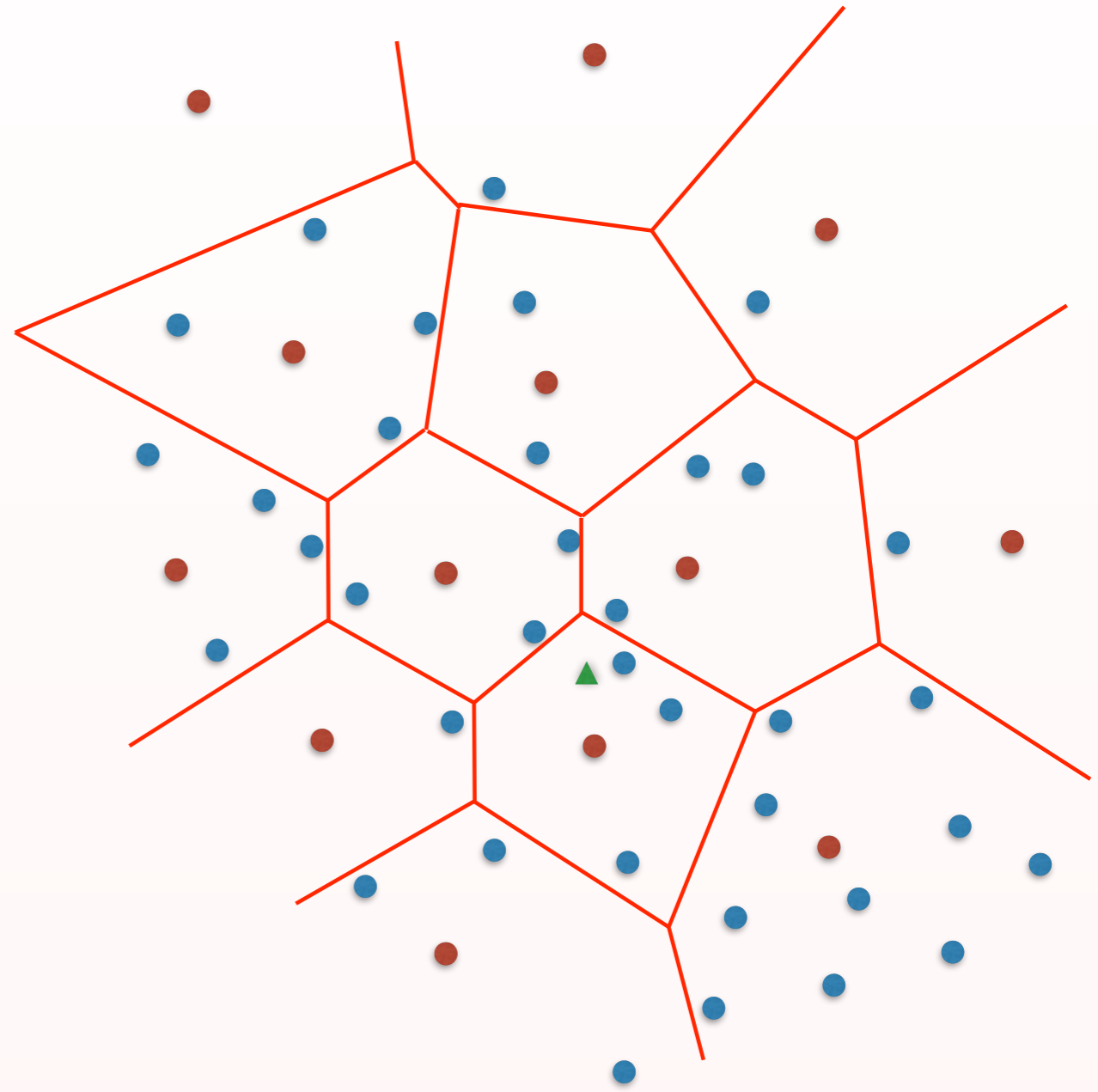


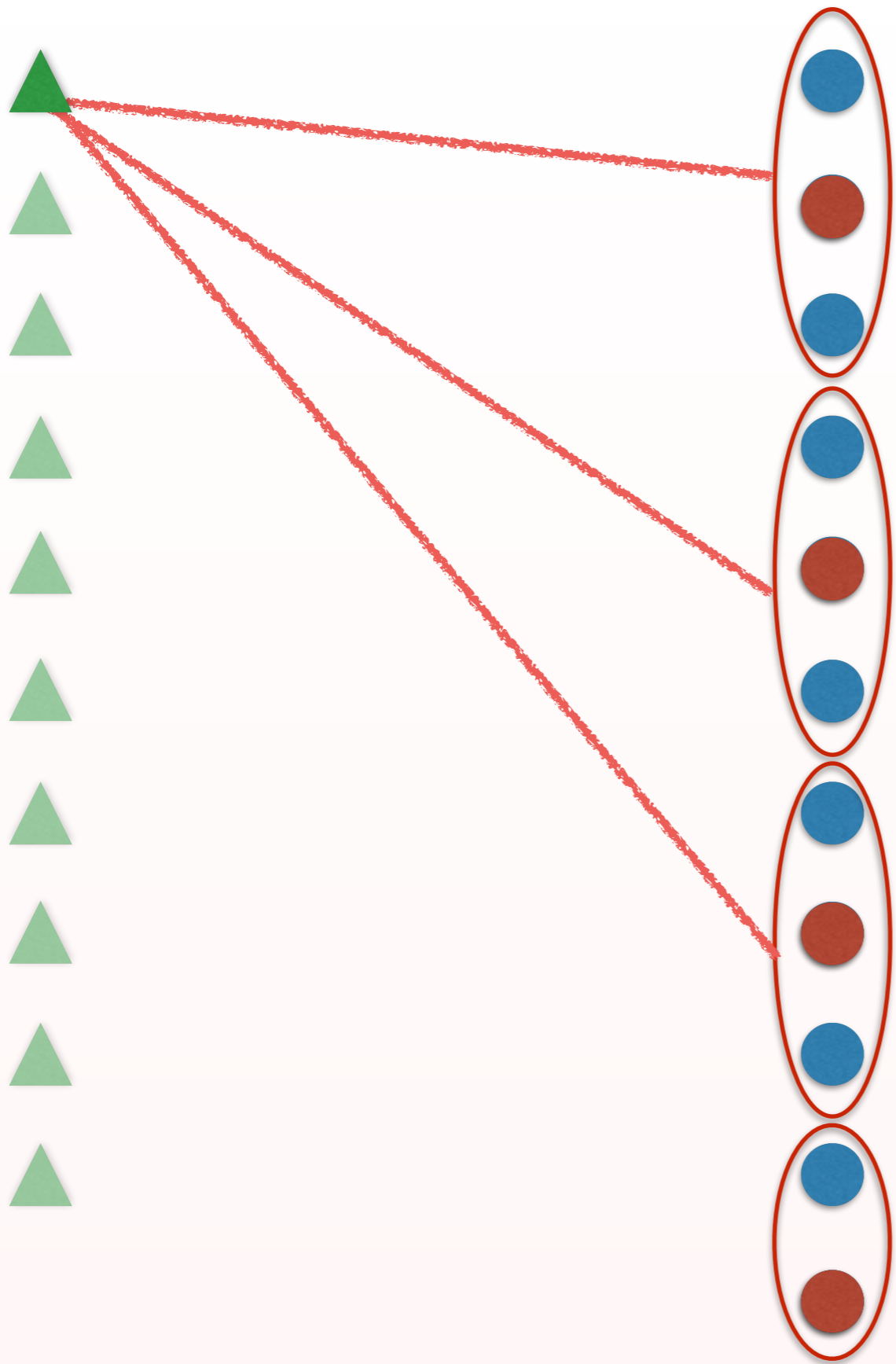


Queries

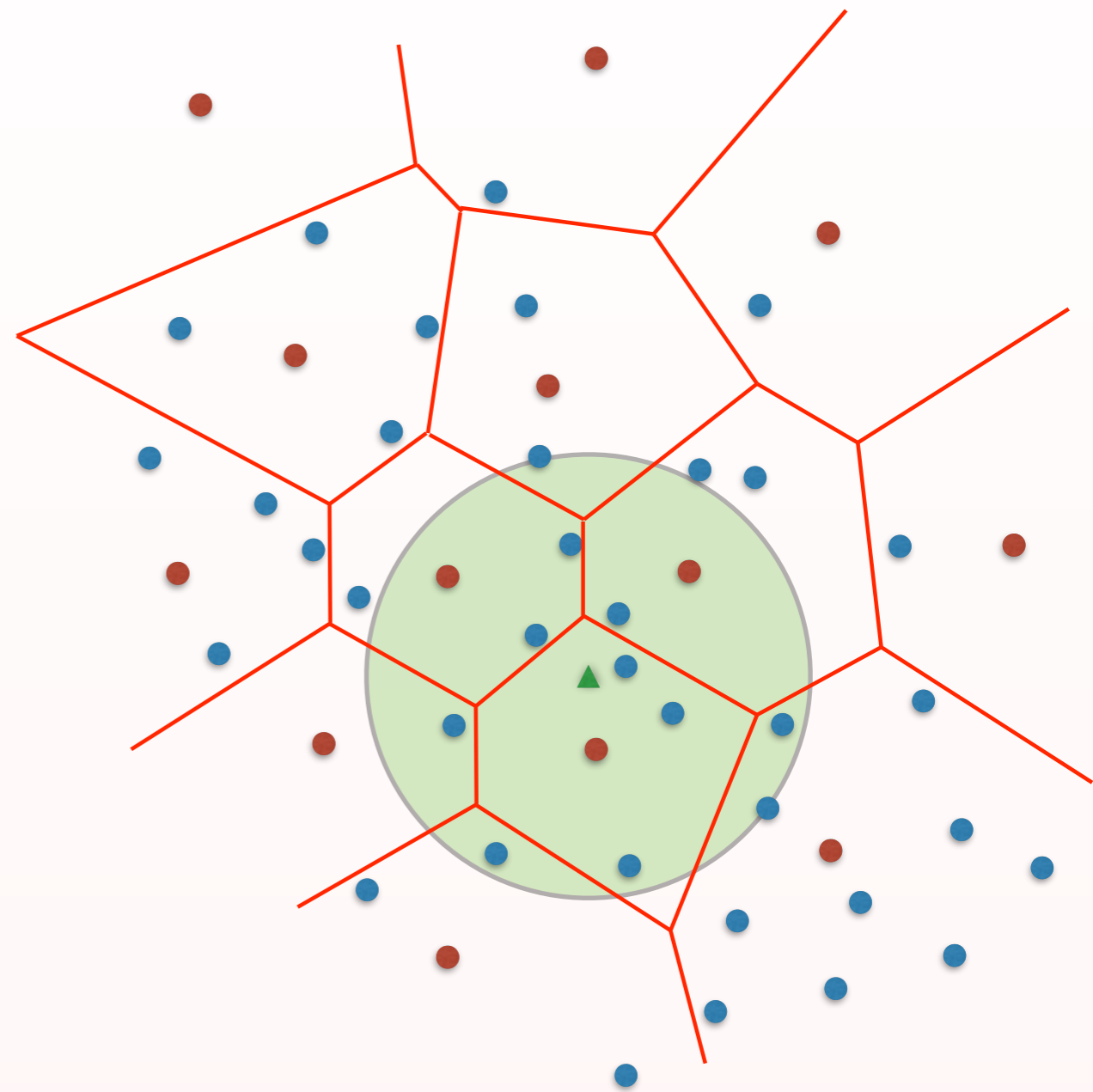


Database

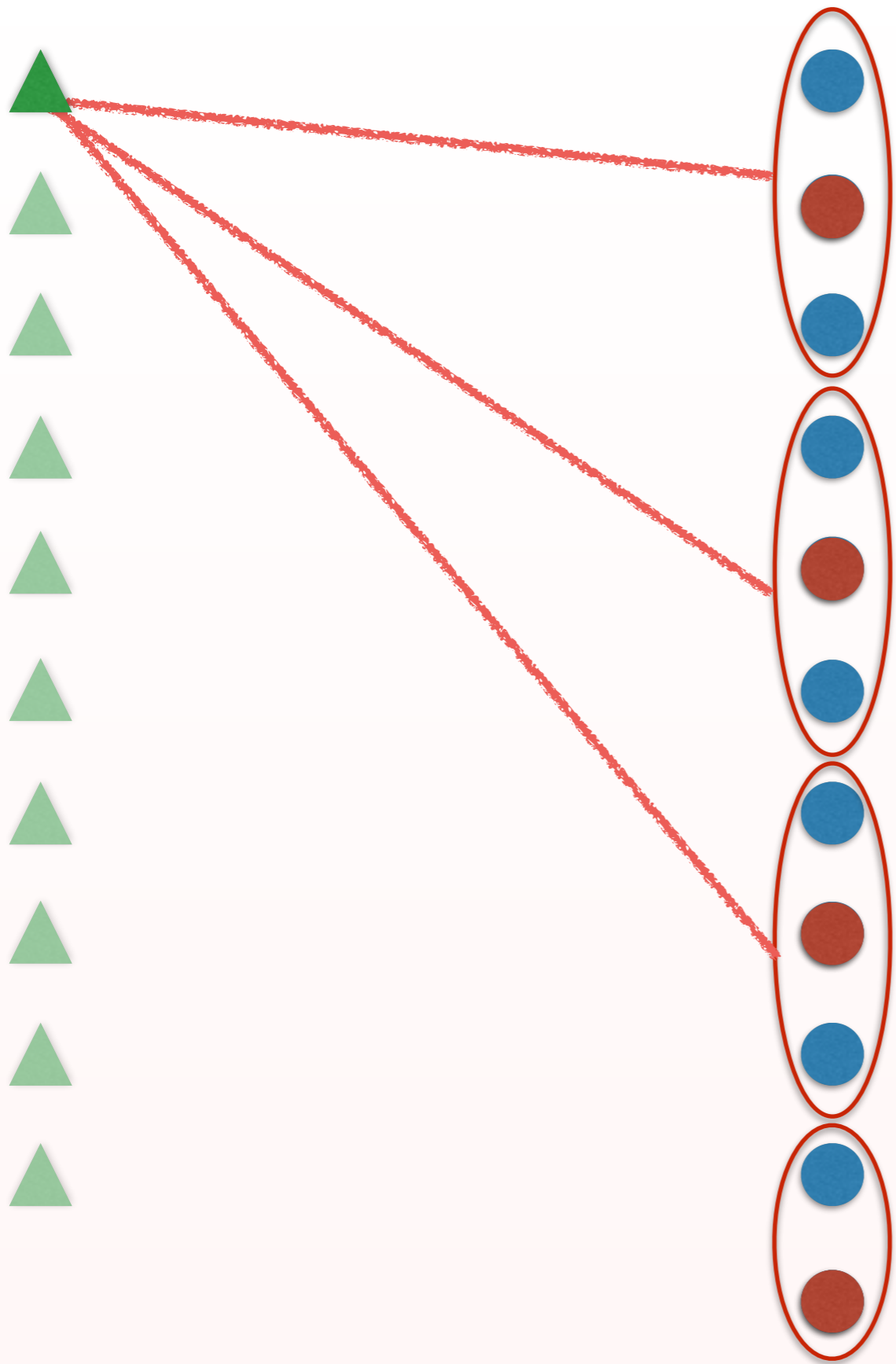




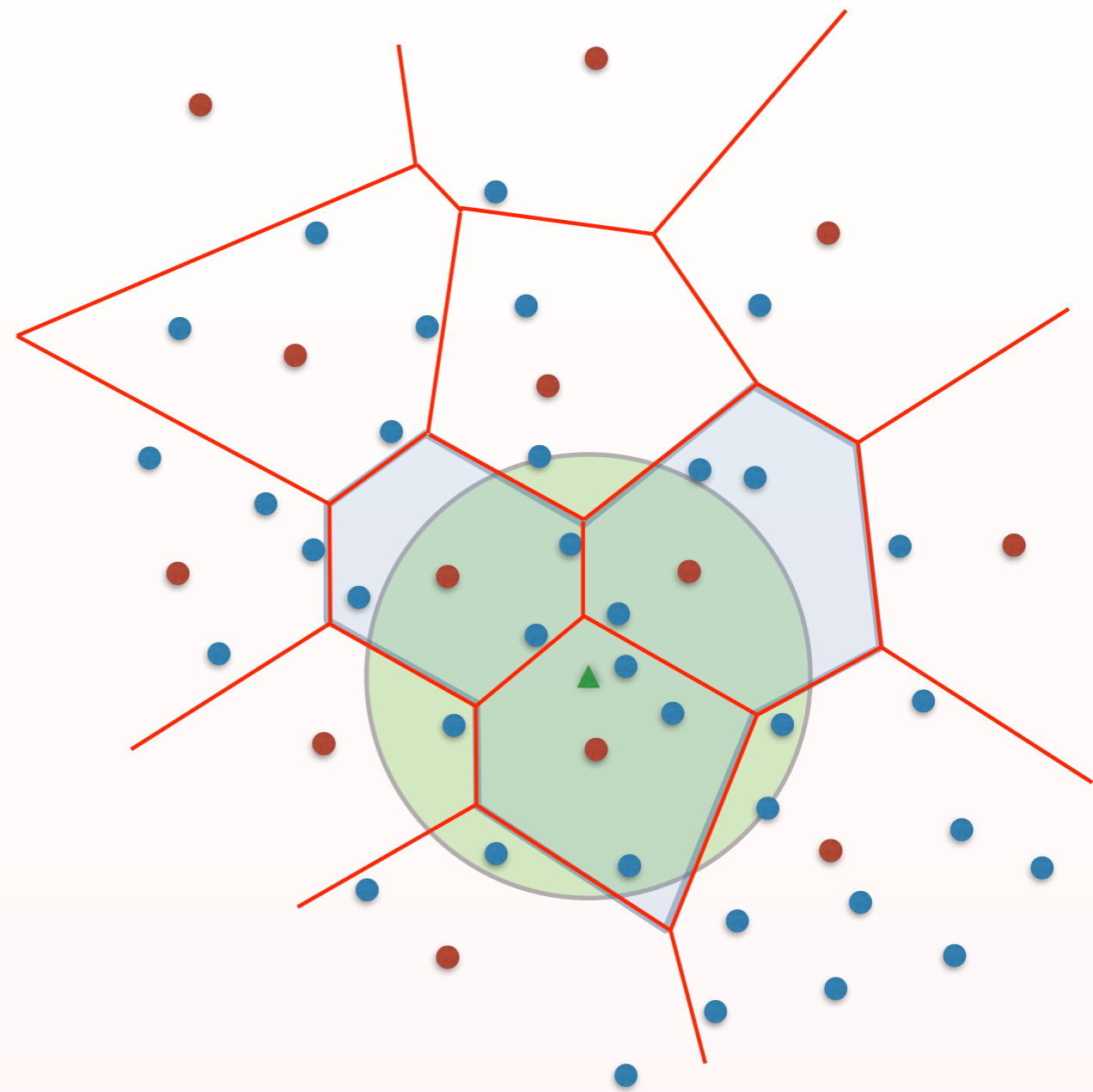
Queries



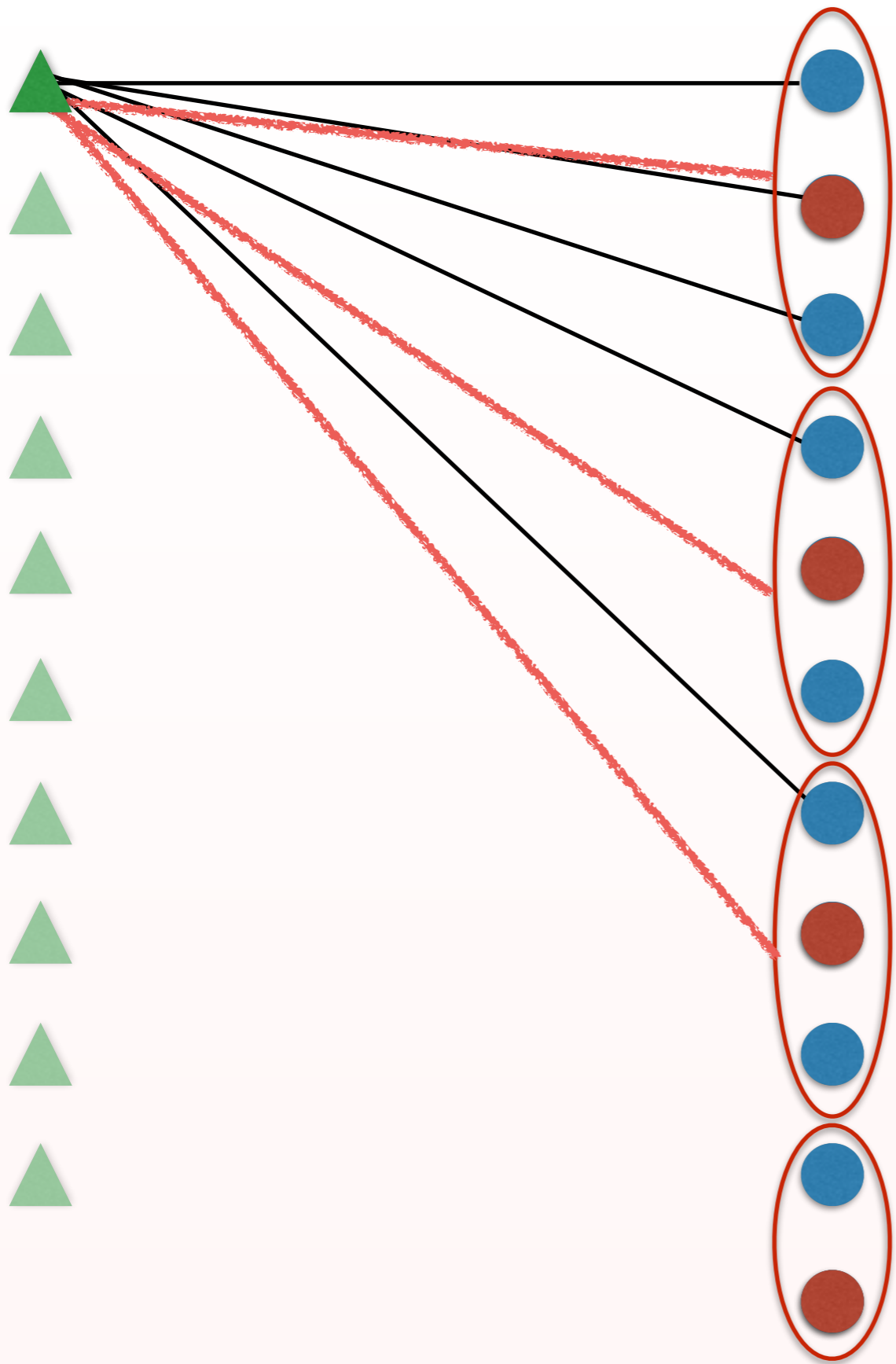
Database



Queries

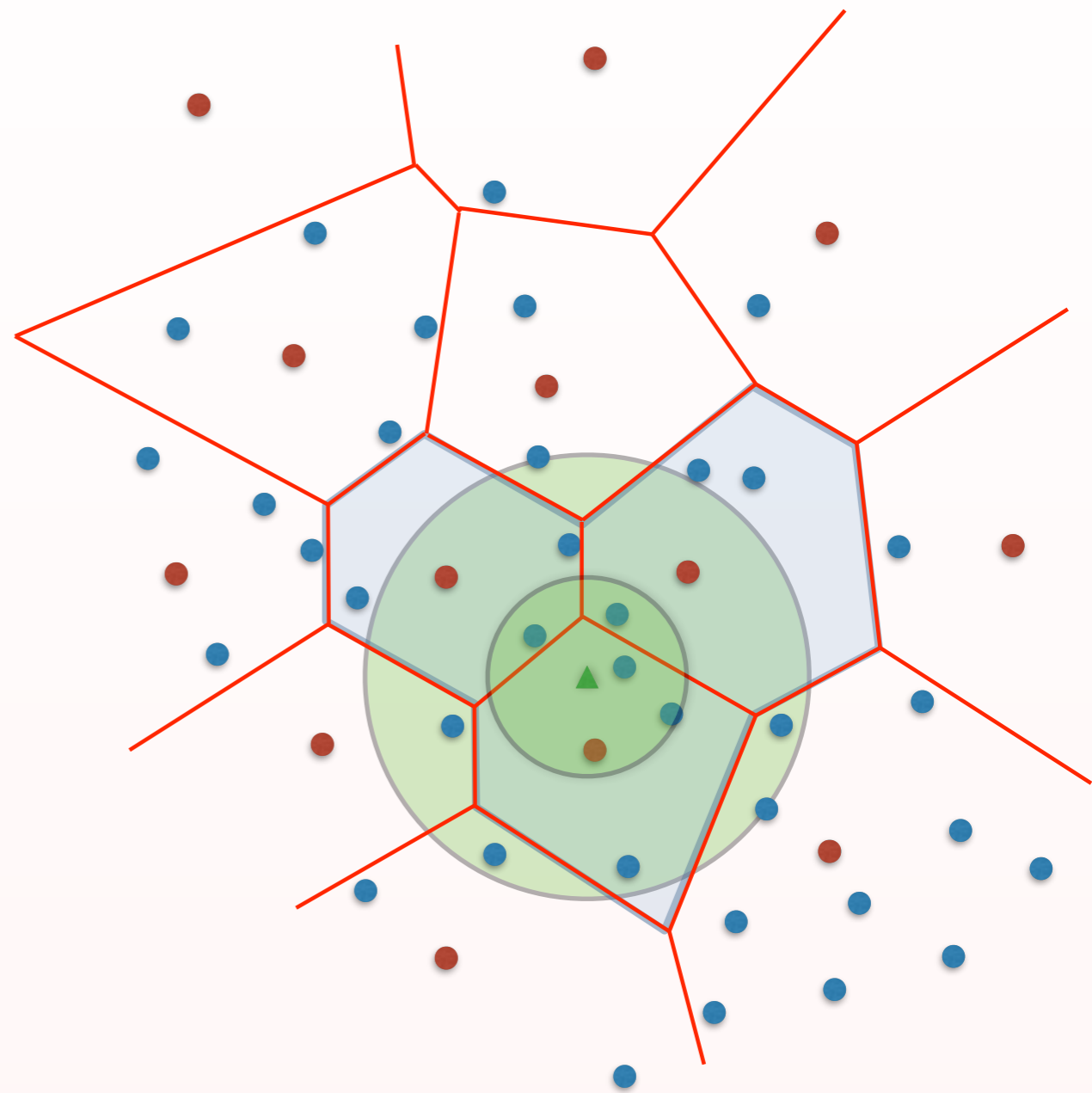


Database

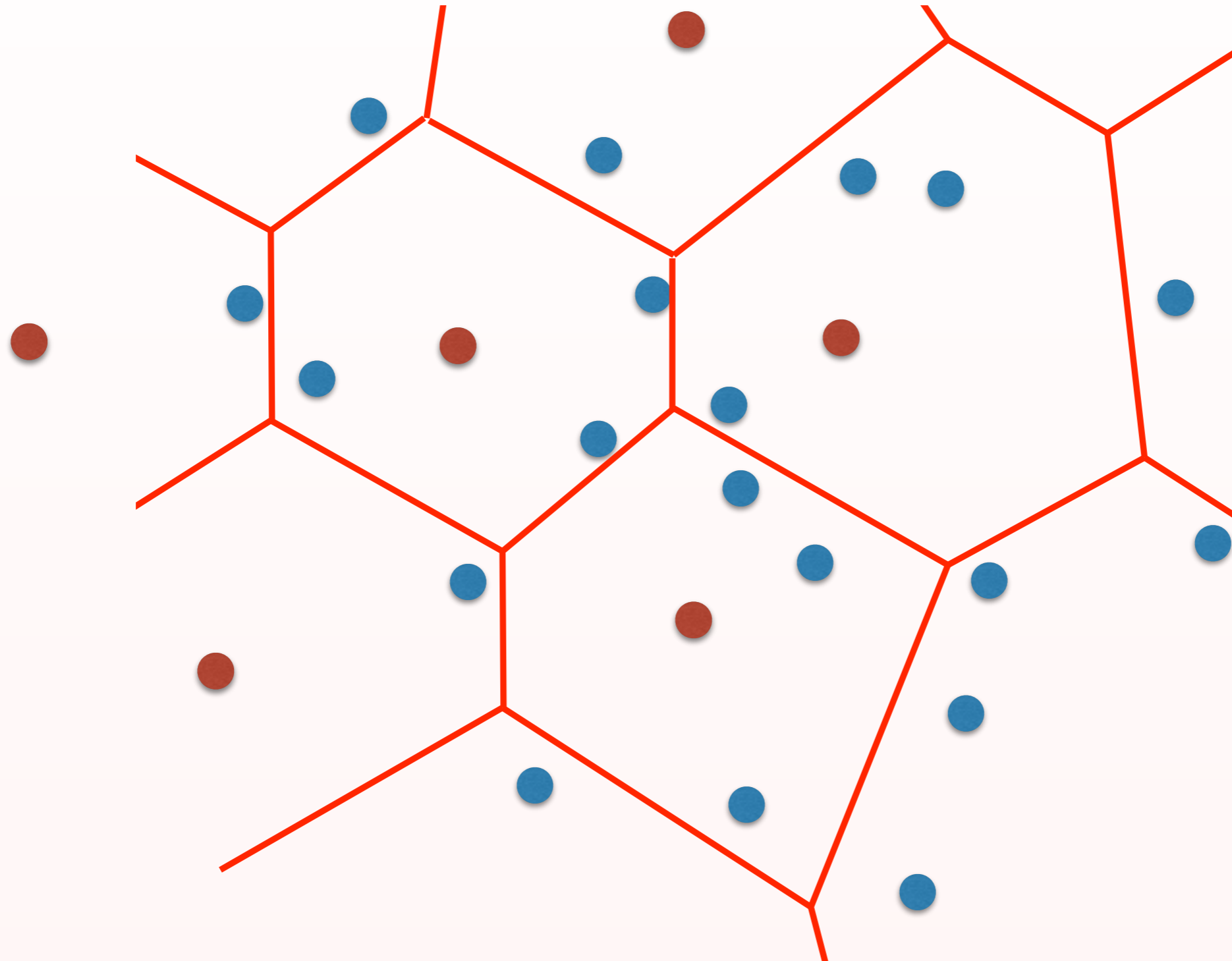


Queries

Database

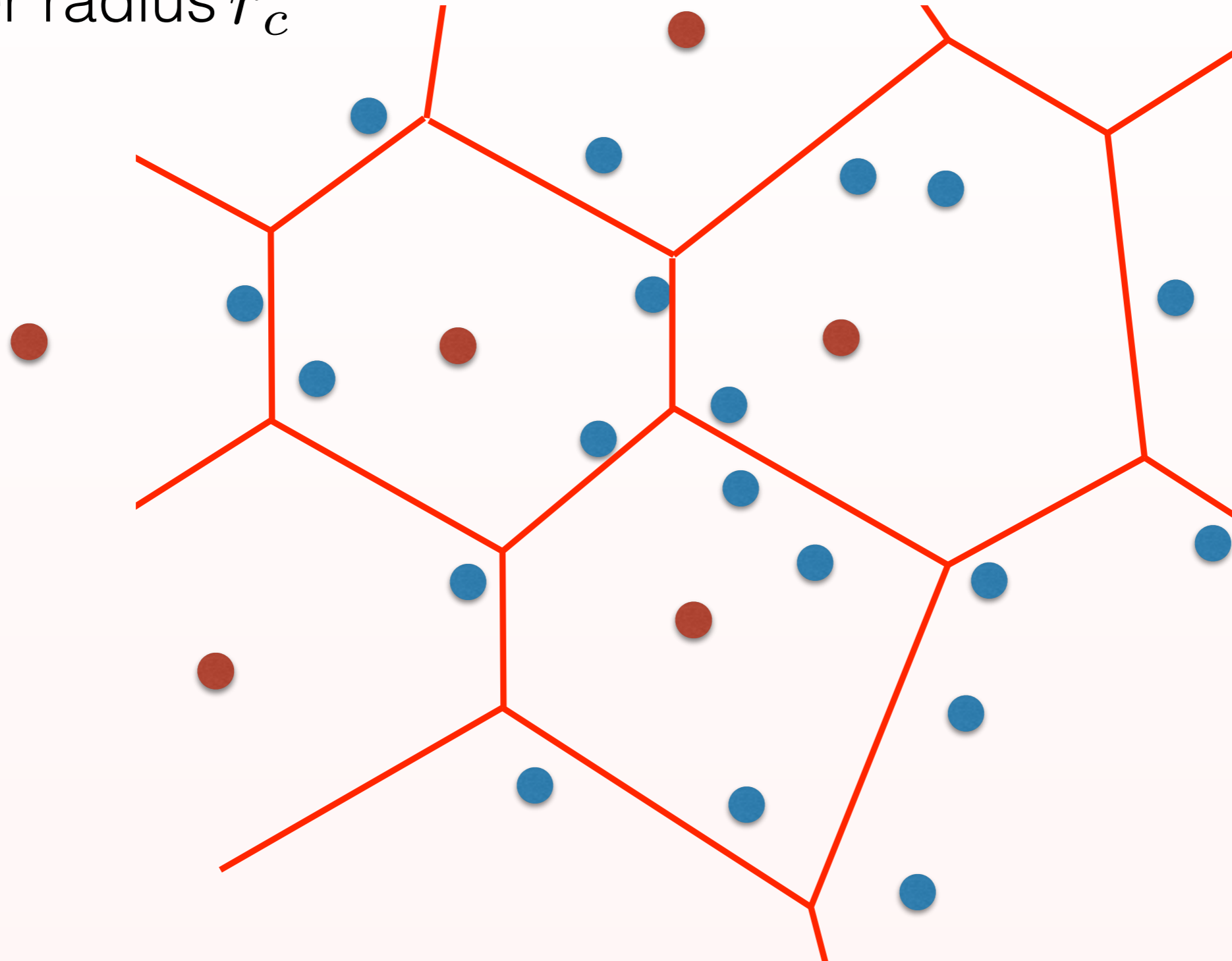


Triangle inequality allows this



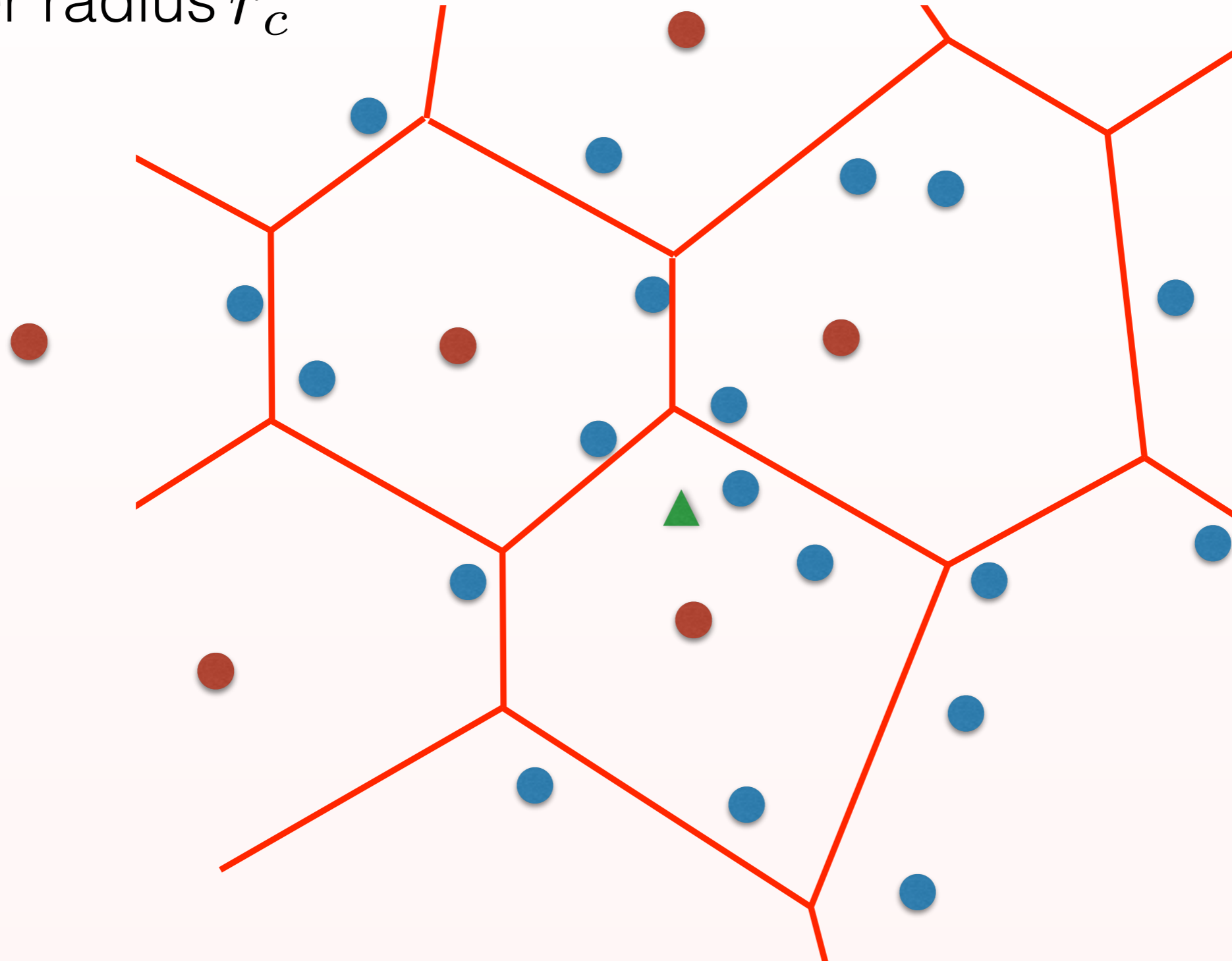
Triangle inequality allows this

cluster radius r_c



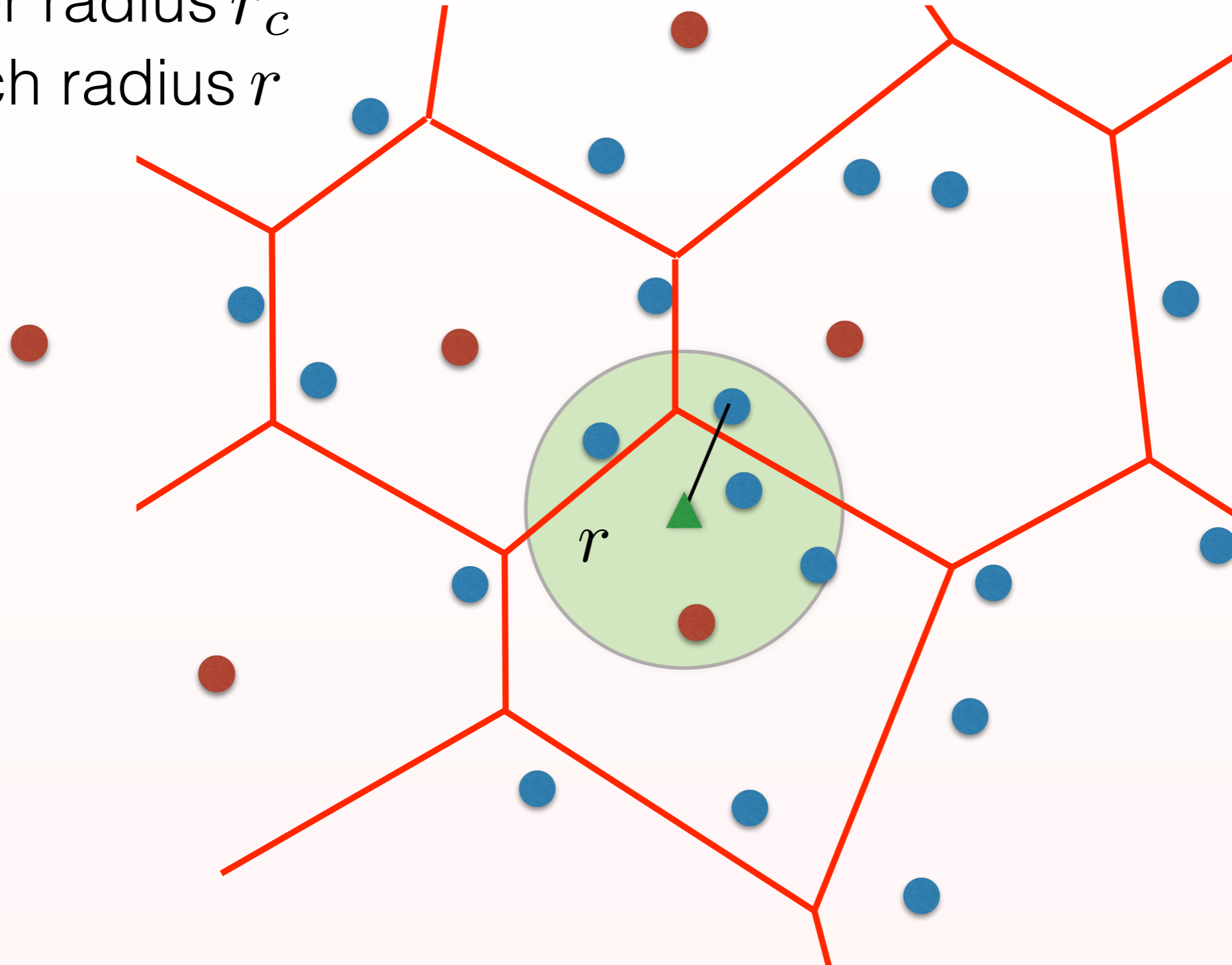
Triangle inequality allows this

cluster radius r_c



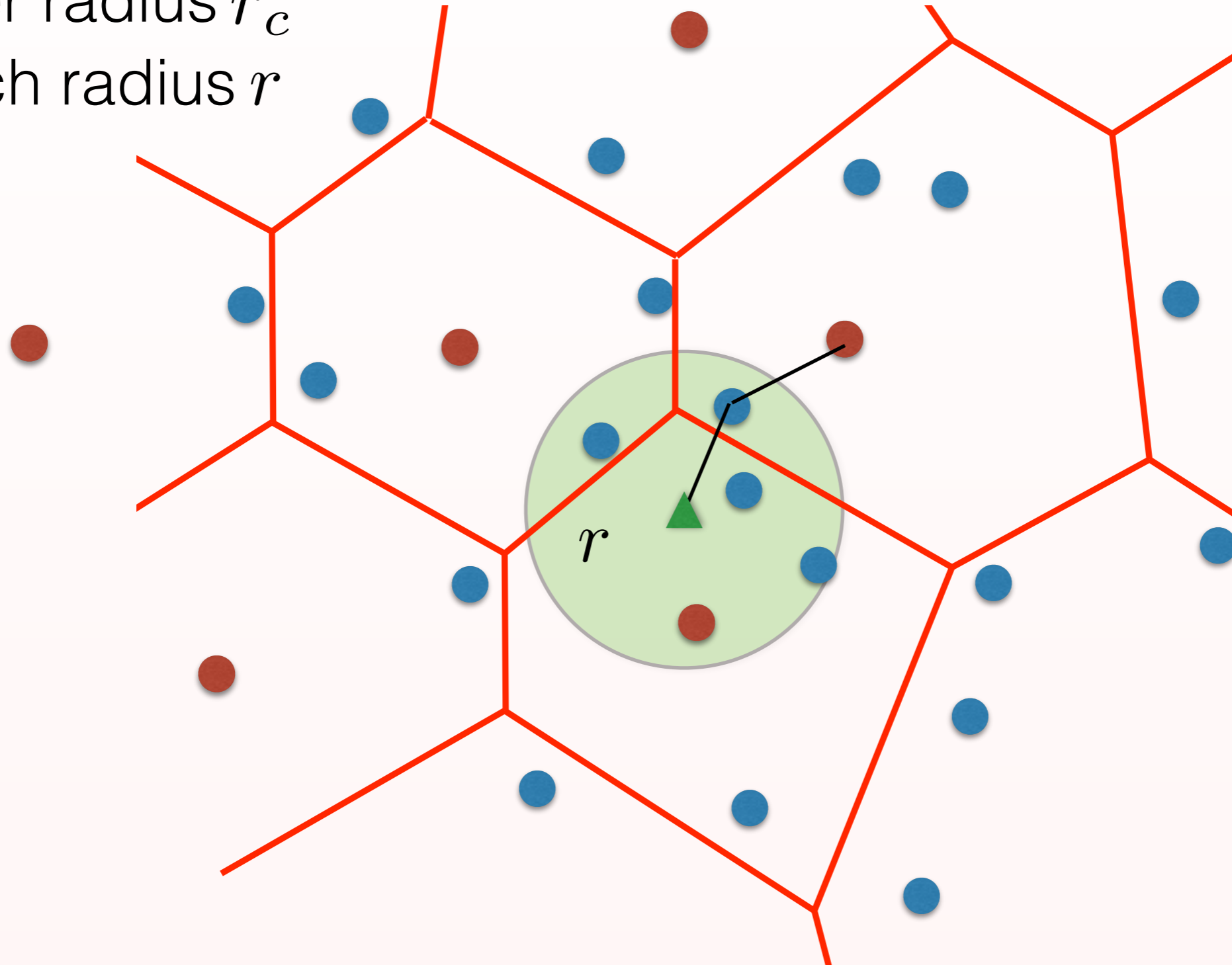
Triangle inequality allows this

cluster radius r_c
search radius r



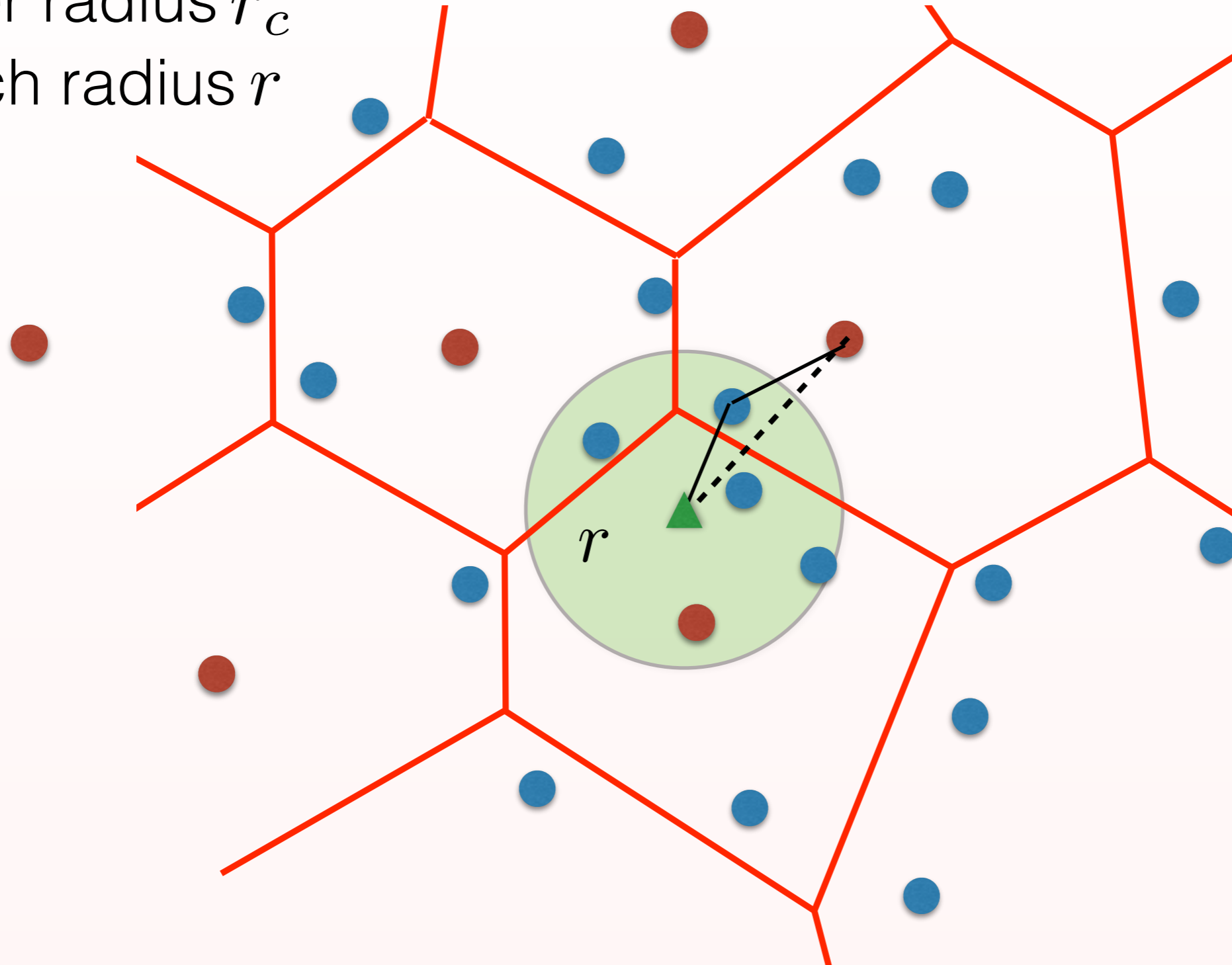
Triangle inequality allows this

cluster radius r_c
search radius r



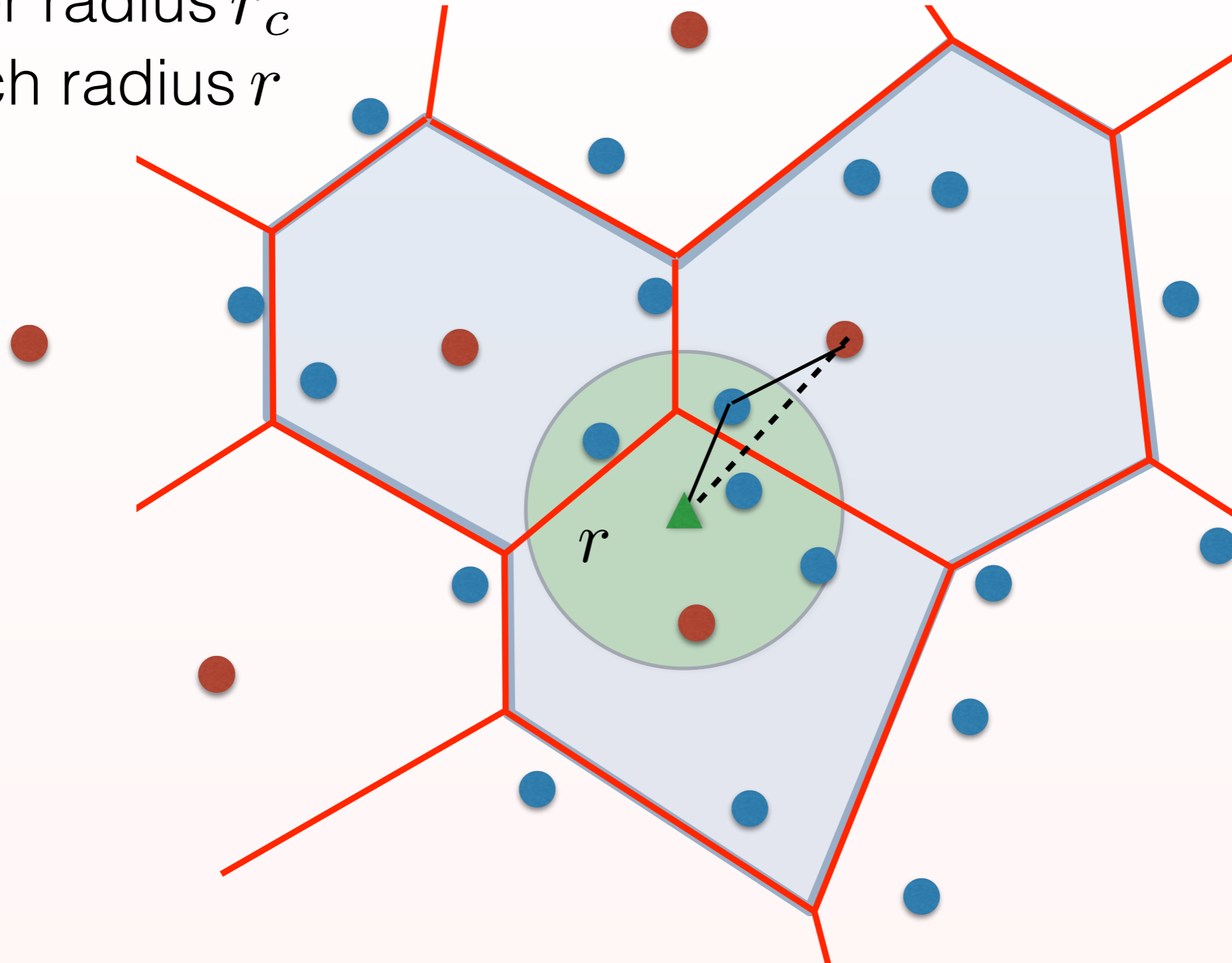
Triangle inequality allows this

cluster radius r_c
search radius r



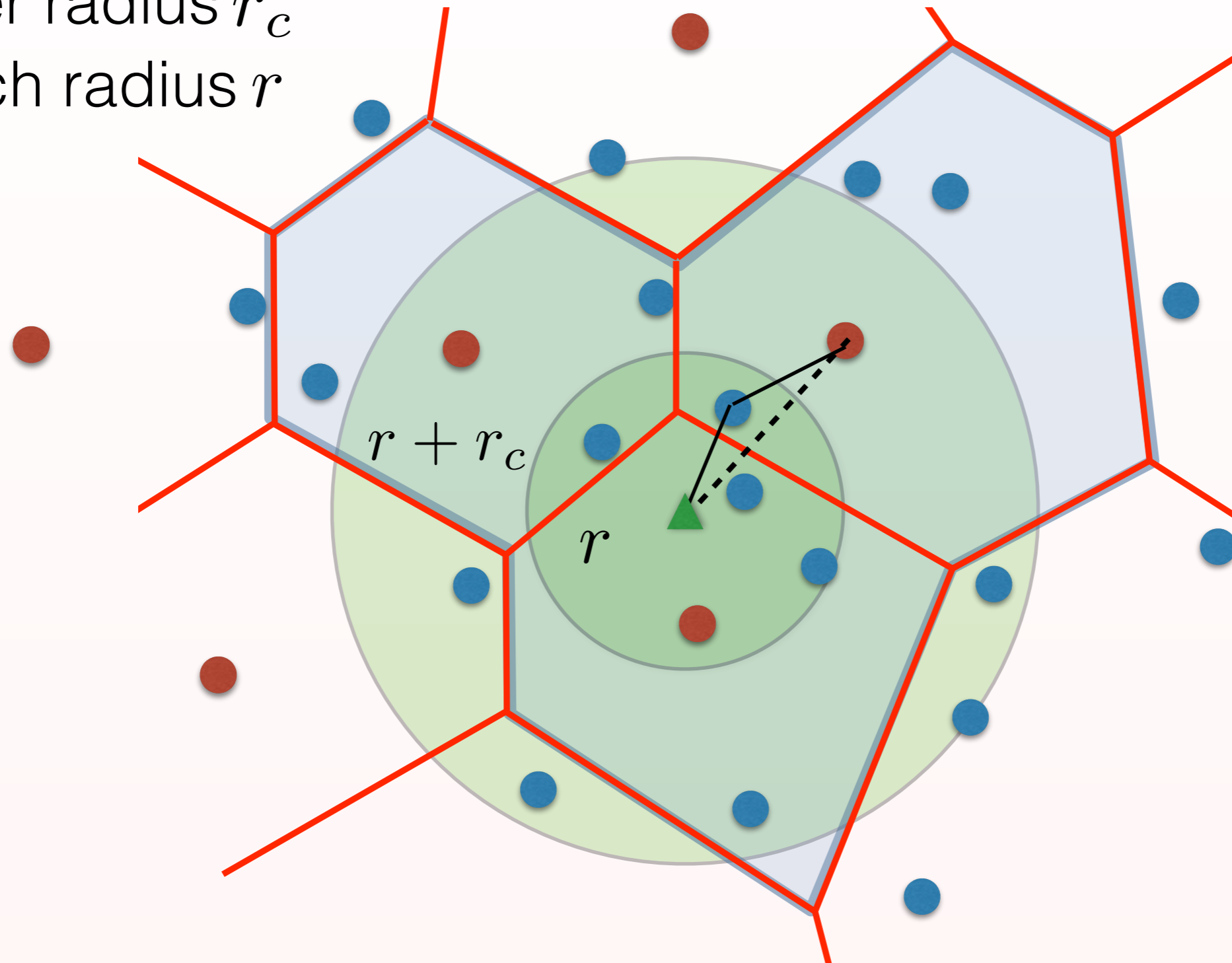
Triangle inequality allows this

cluster radius r_c
search radius r



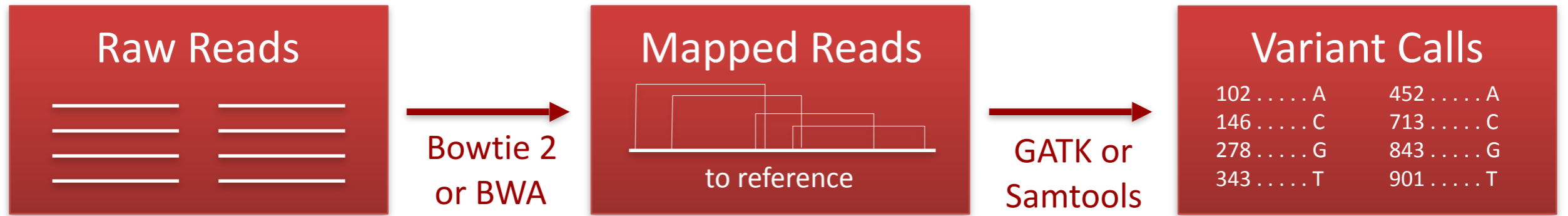
Triangle inequality allows this

cluster radius r_c
search radius r

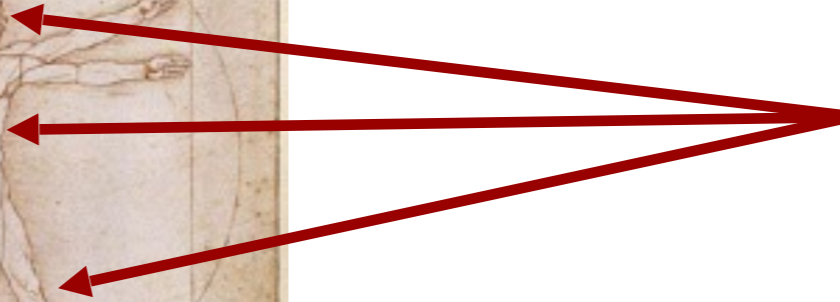
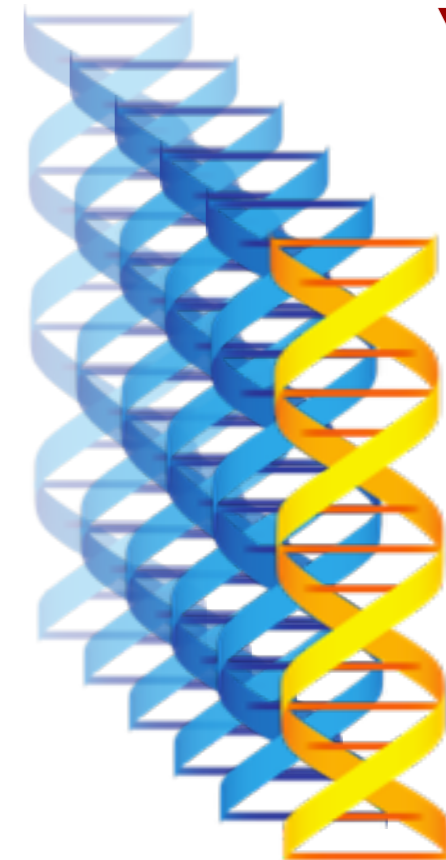
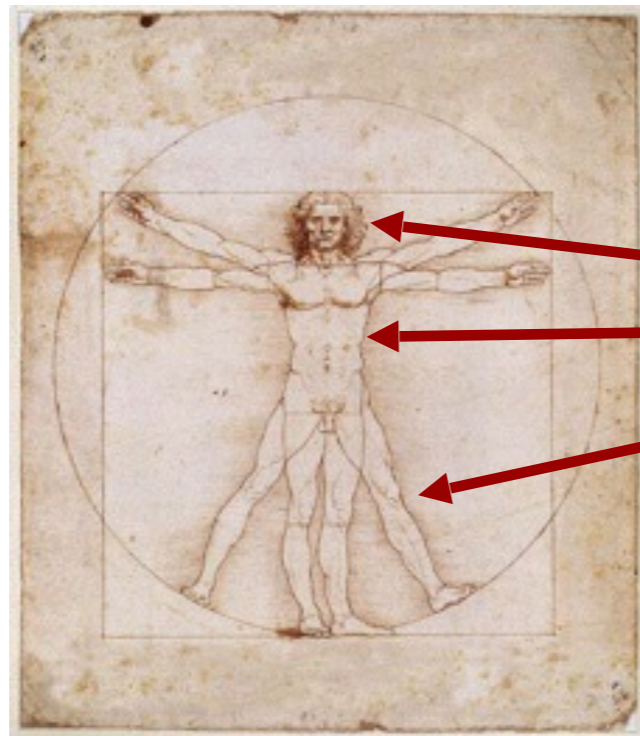
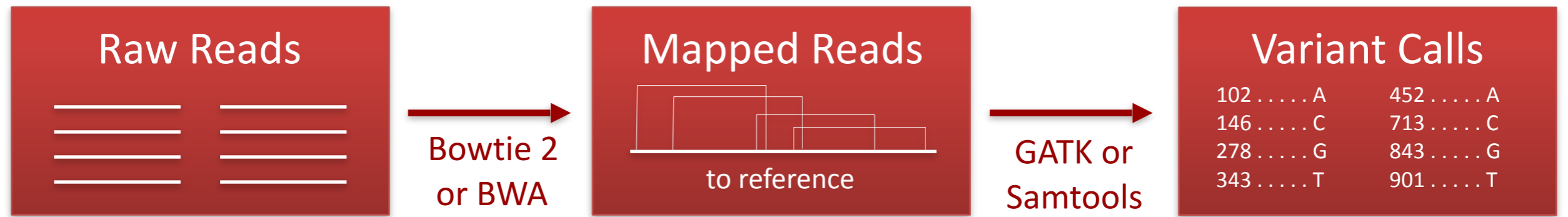


Large-scale NGS Analysis

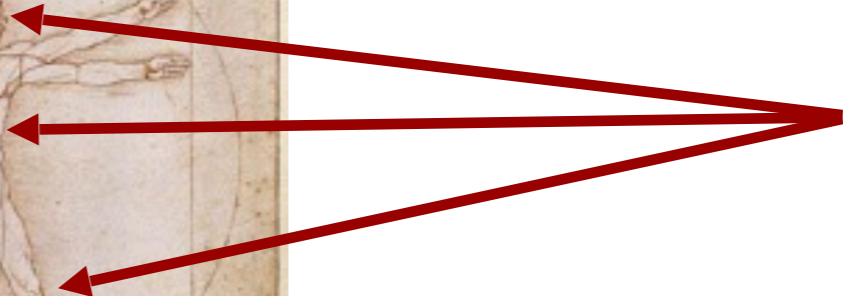
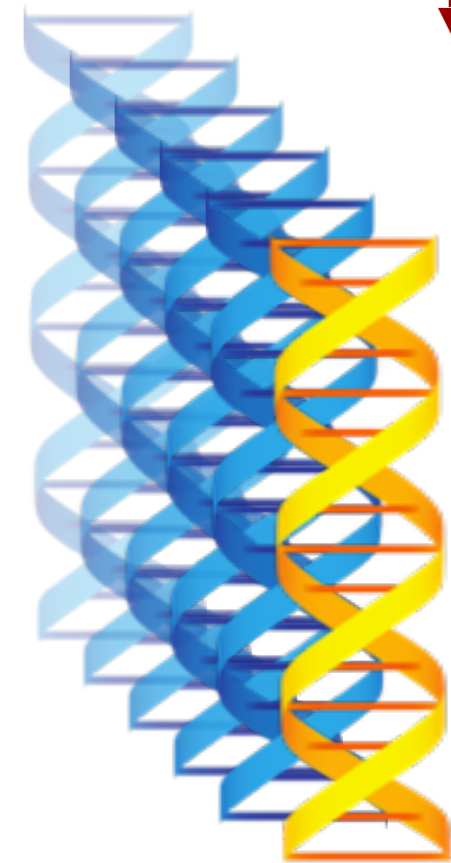
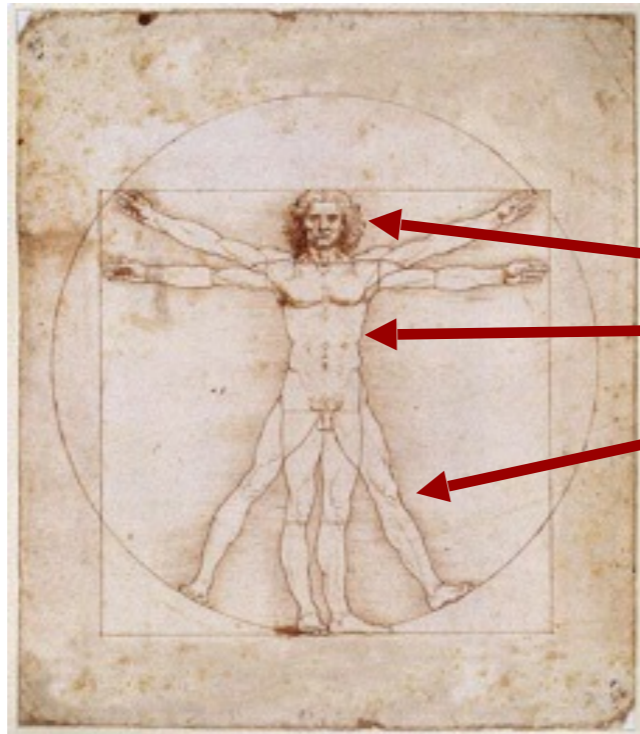
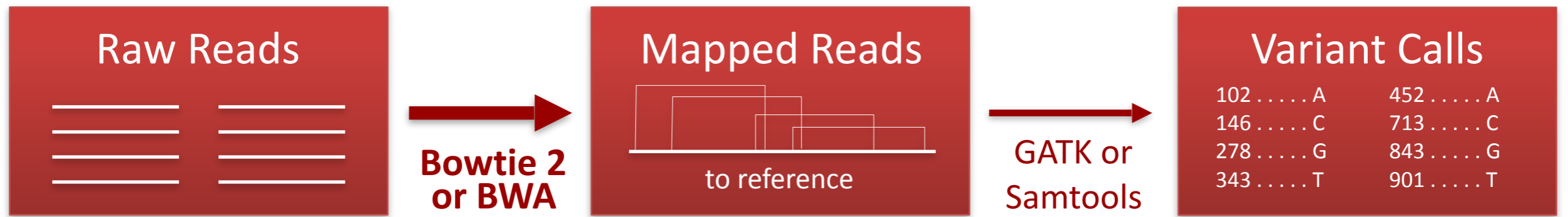
Large-scale NGS Analysis



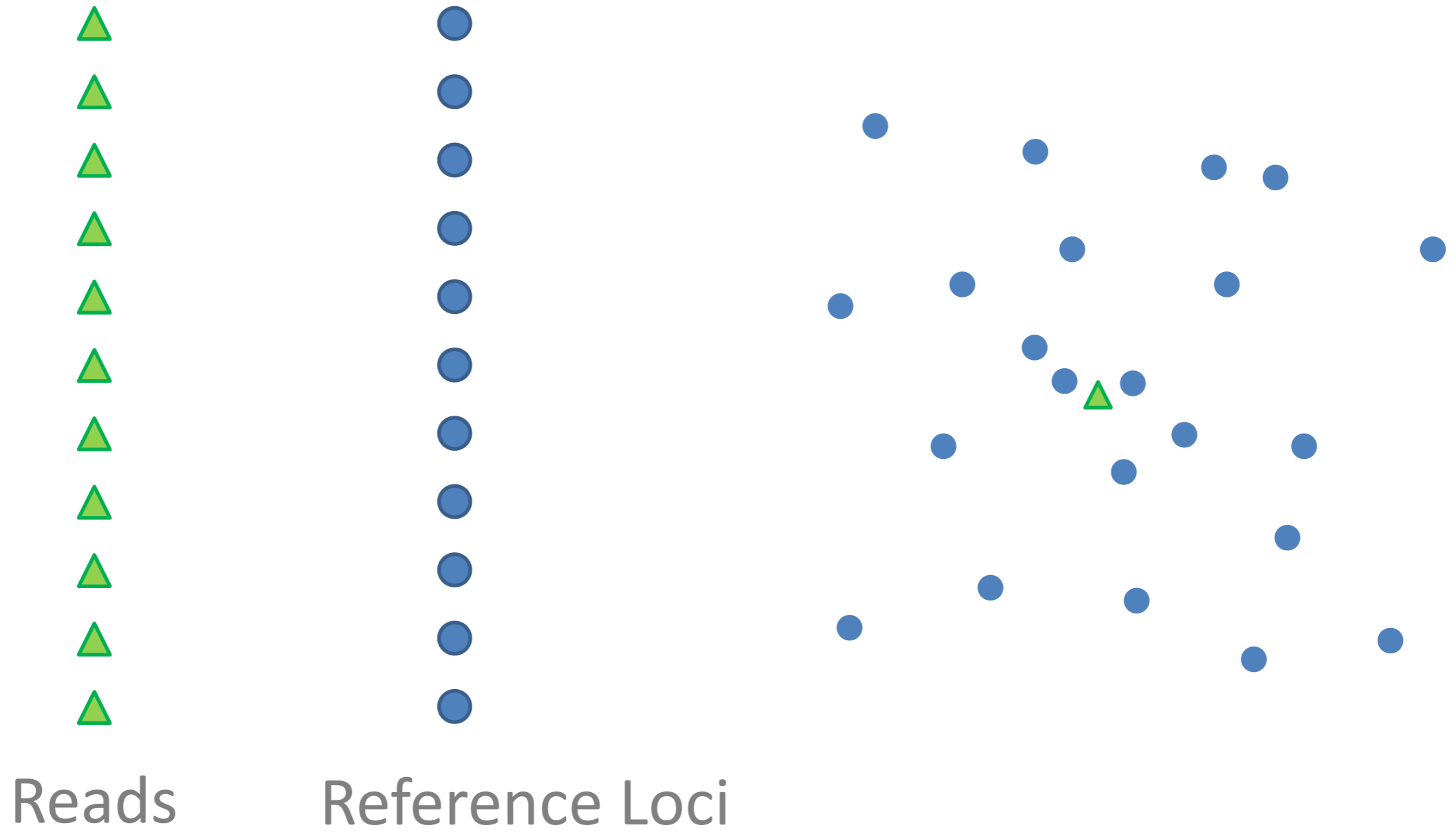
Large-scale NGS Analysis



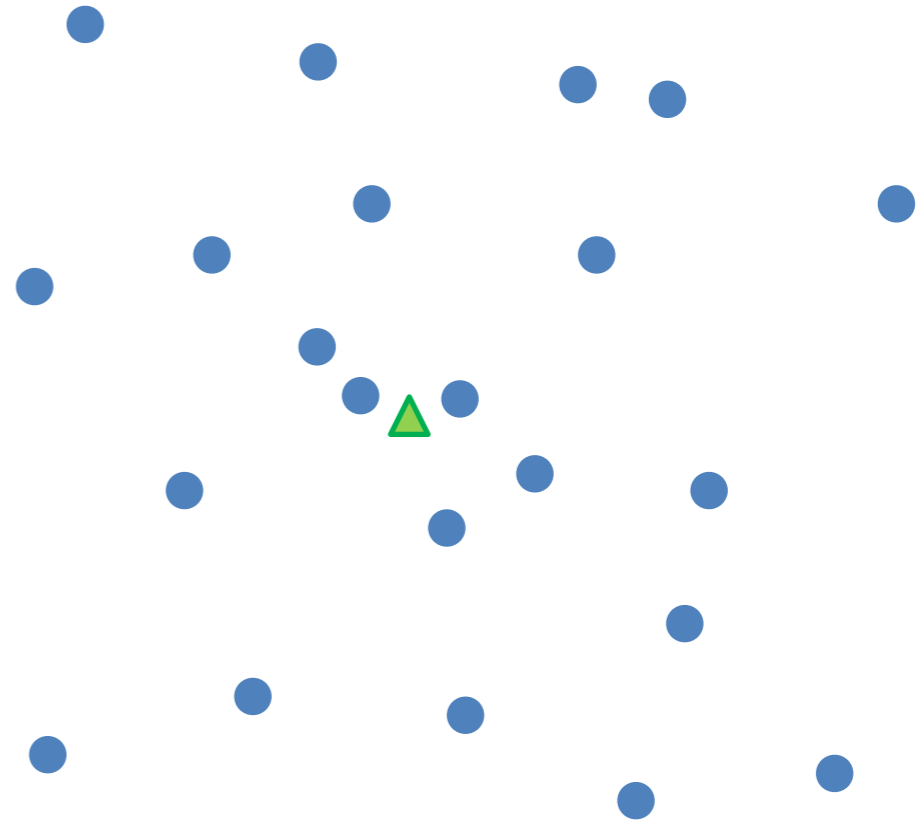
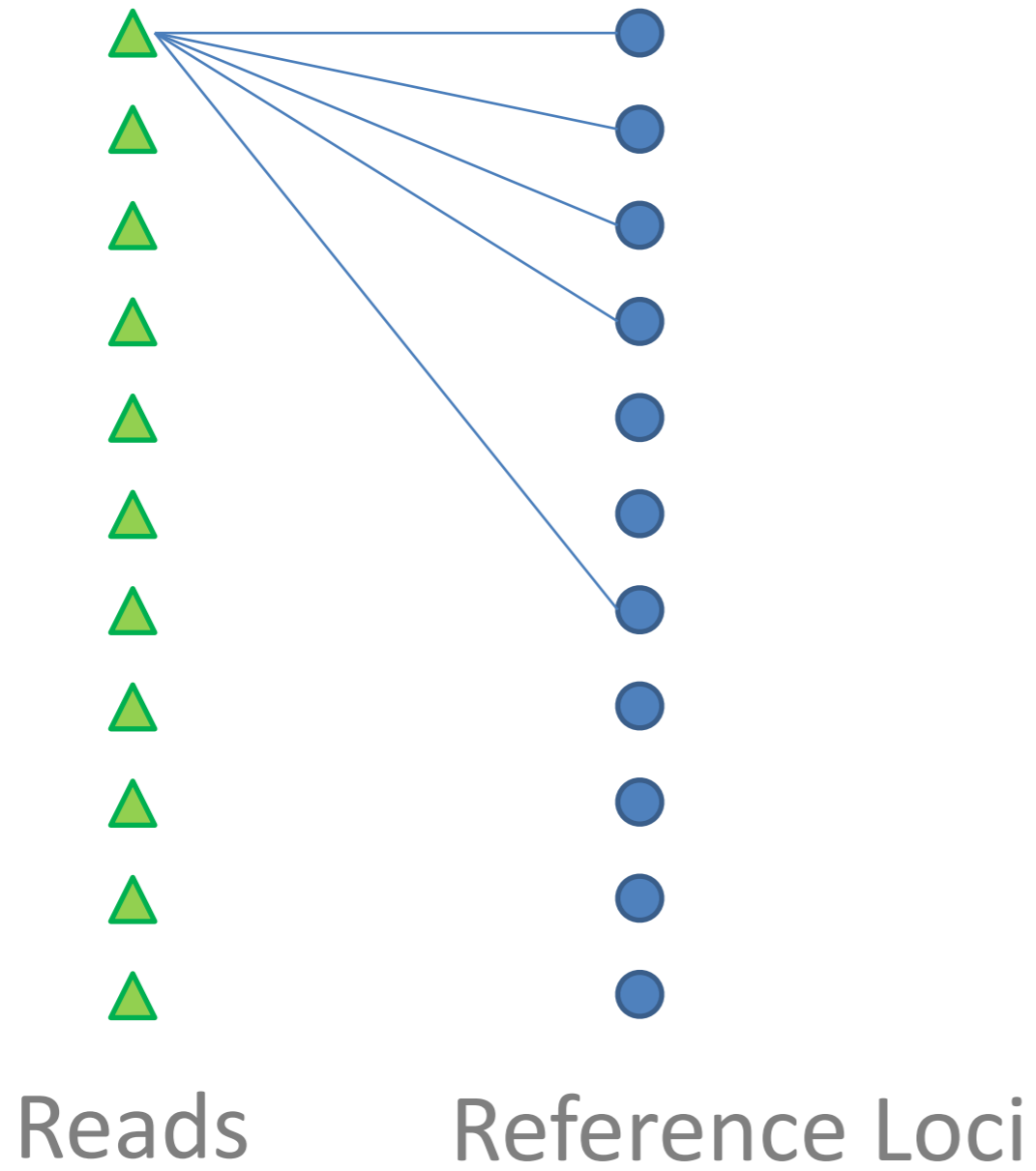
Large-scale NGS Analysis



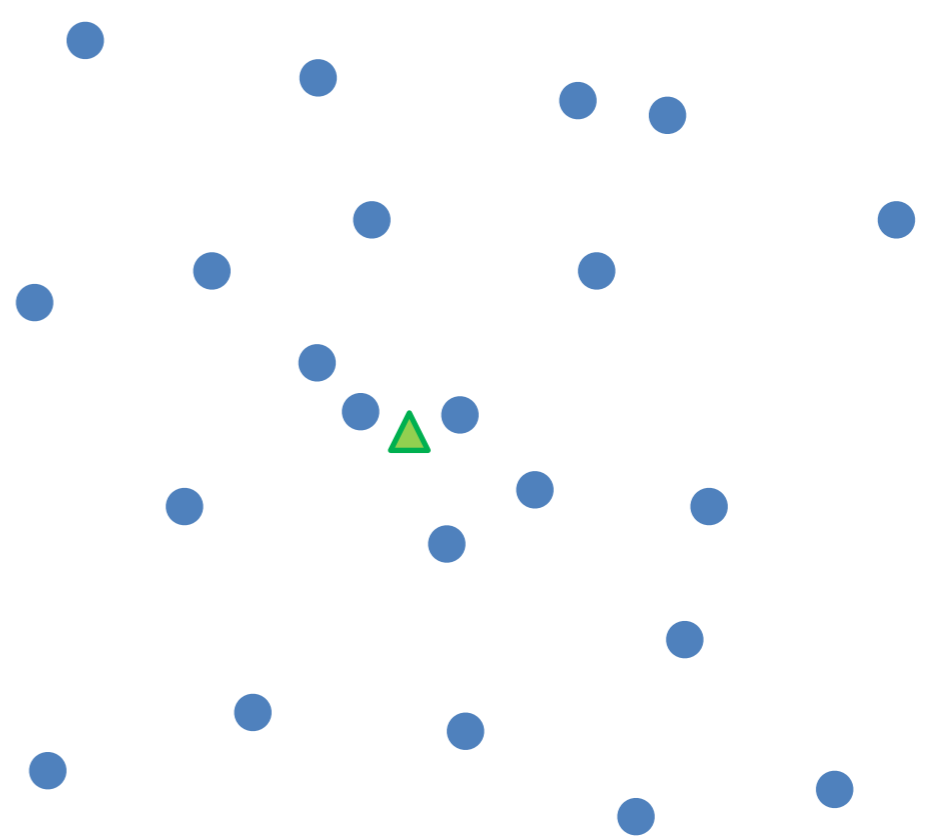
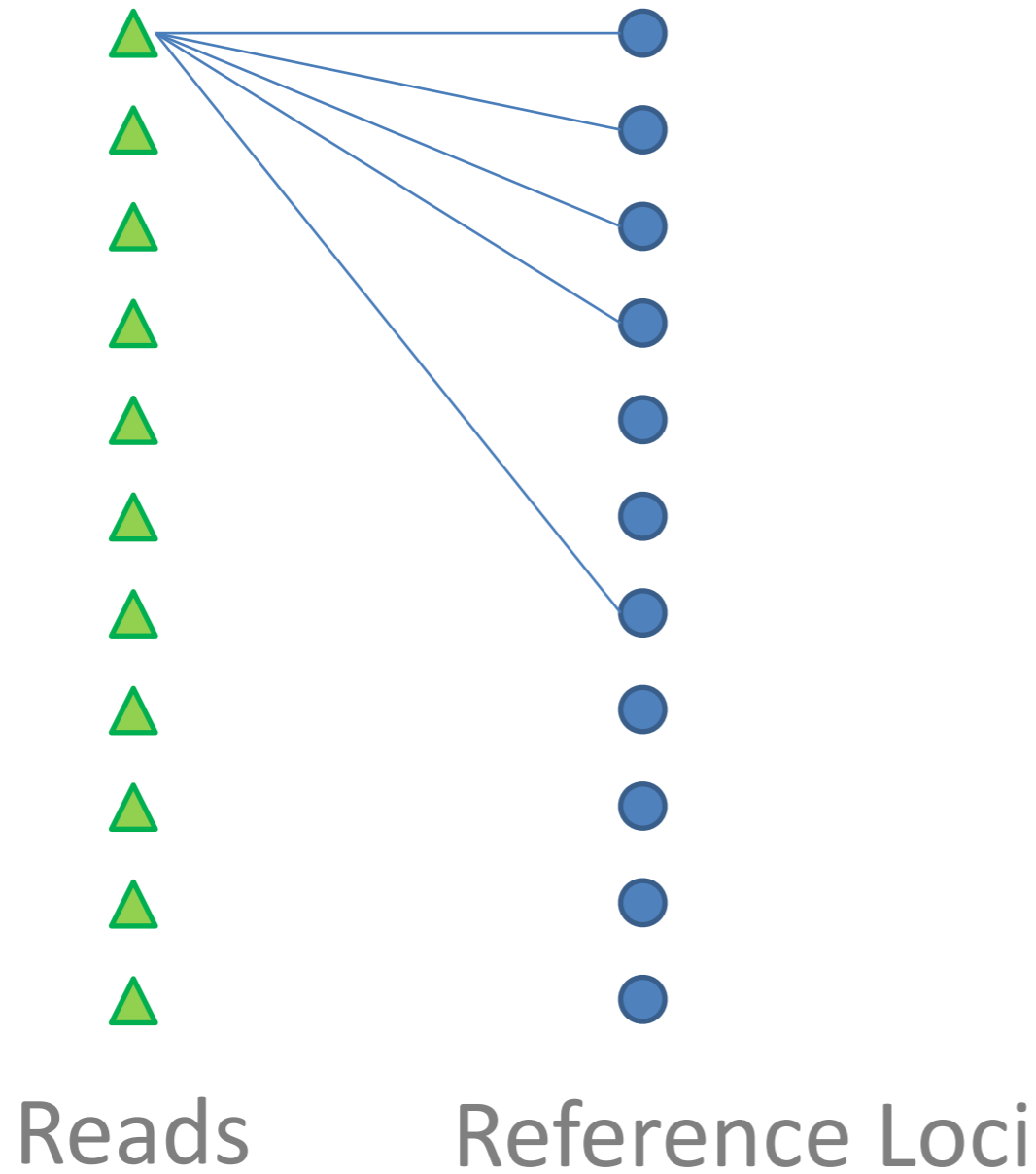
NGS-Mapping in Personal Genomics Era



NGS-Mapping in Personal Genomics Era

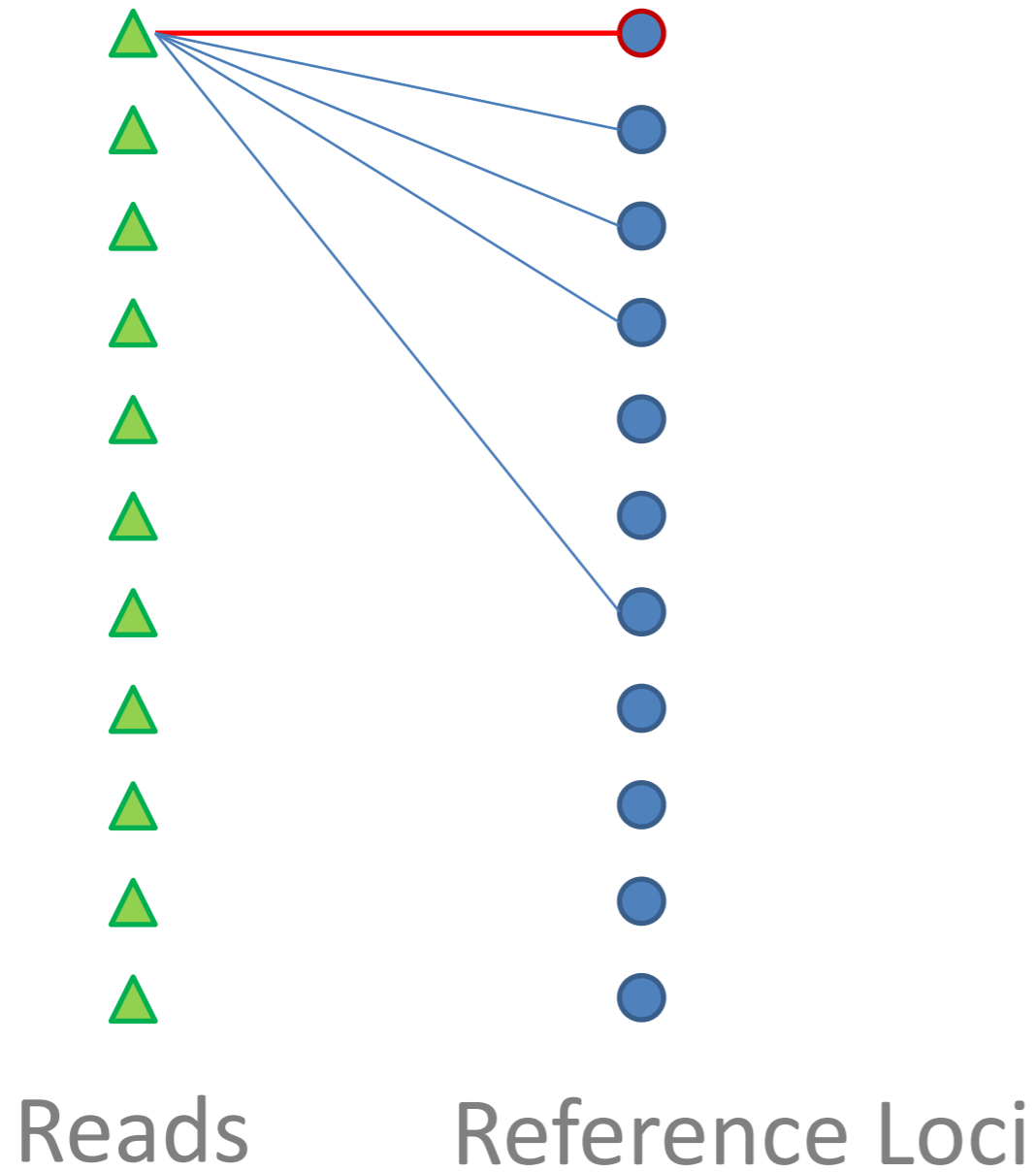


NGS-Mapping in Personal Genomics Era

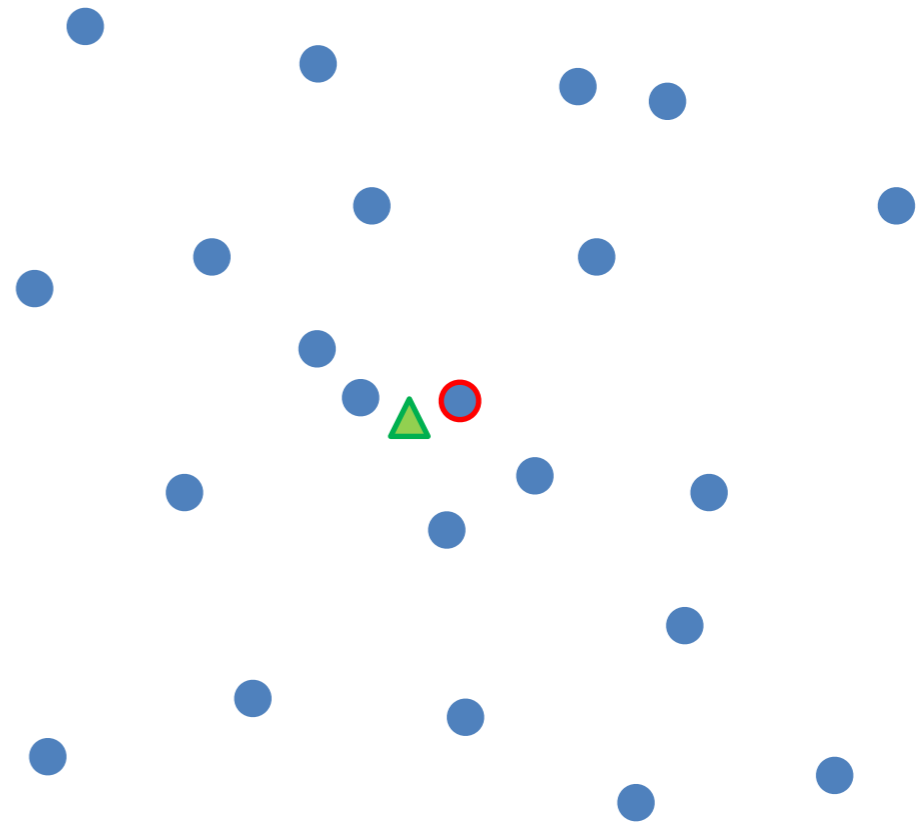


Key Intuition: Reads have lots of redundancy!

NGS-Mapping in Personal Genomics Era



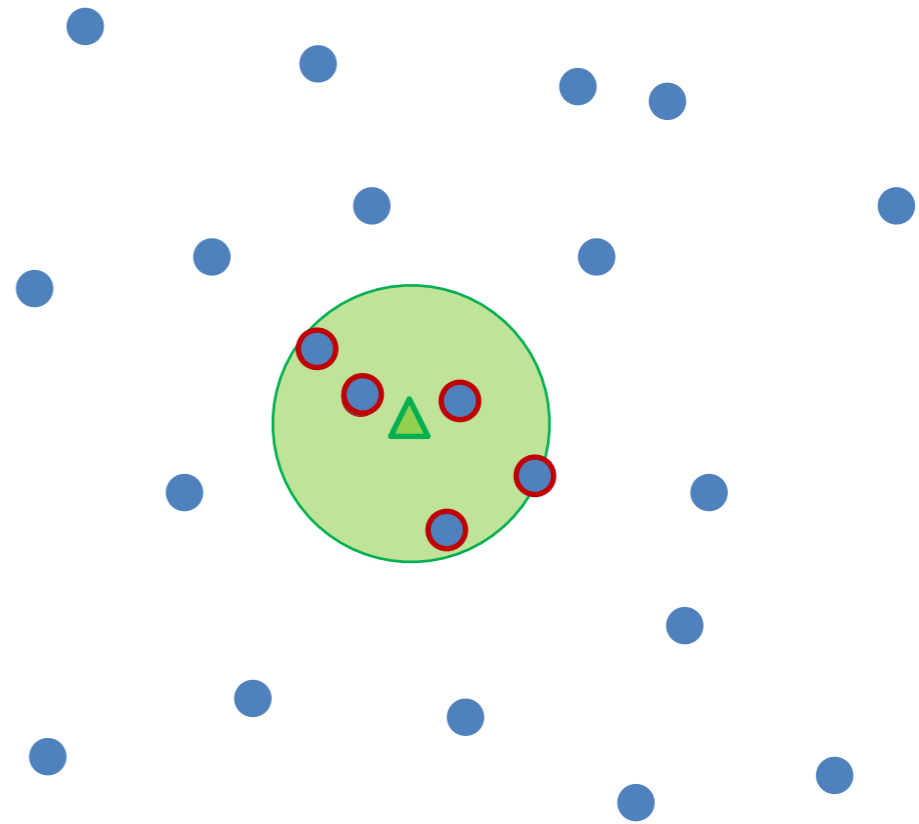
Best Mapping



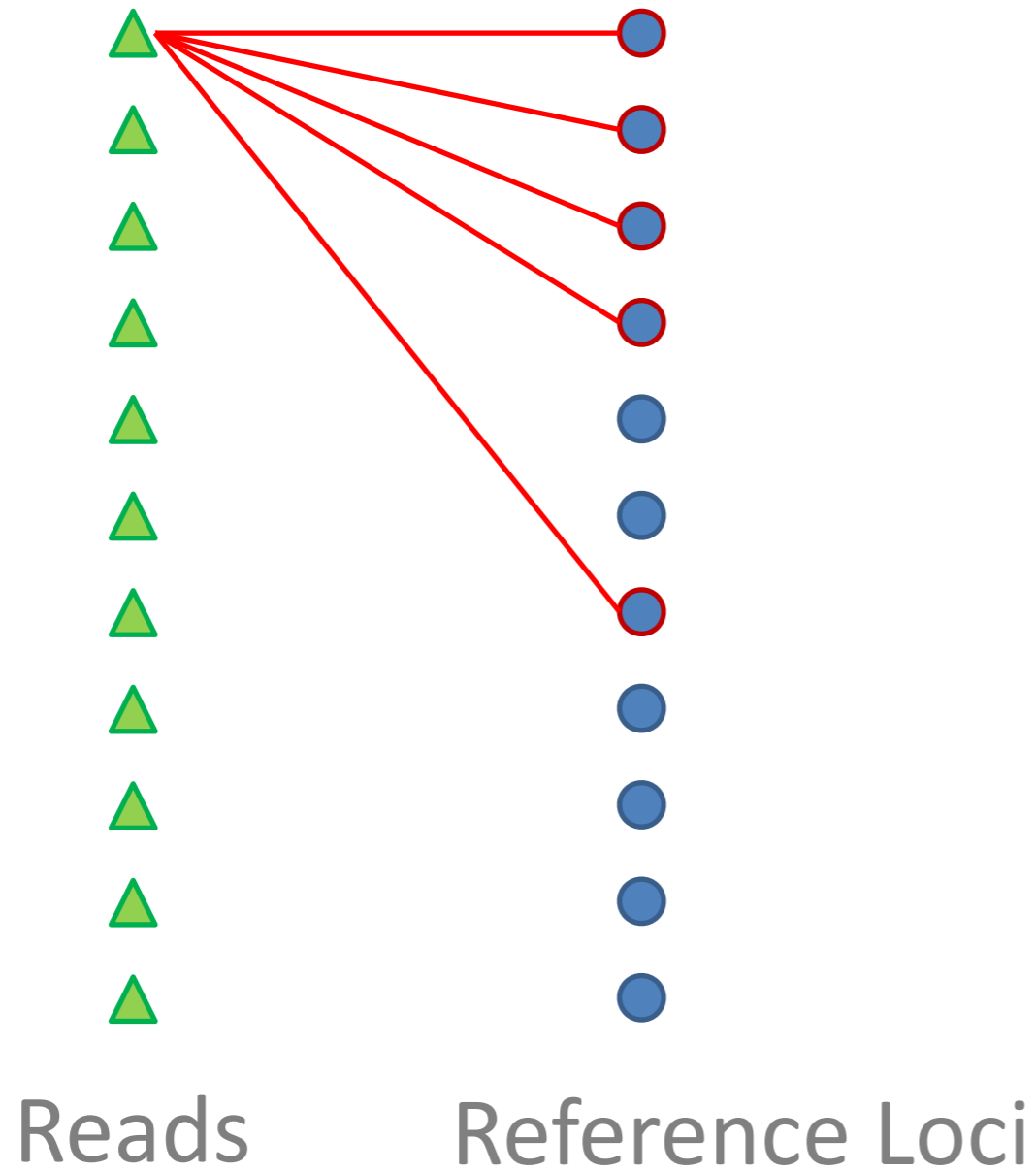
NGS-Mapping in Personal Genomics Era



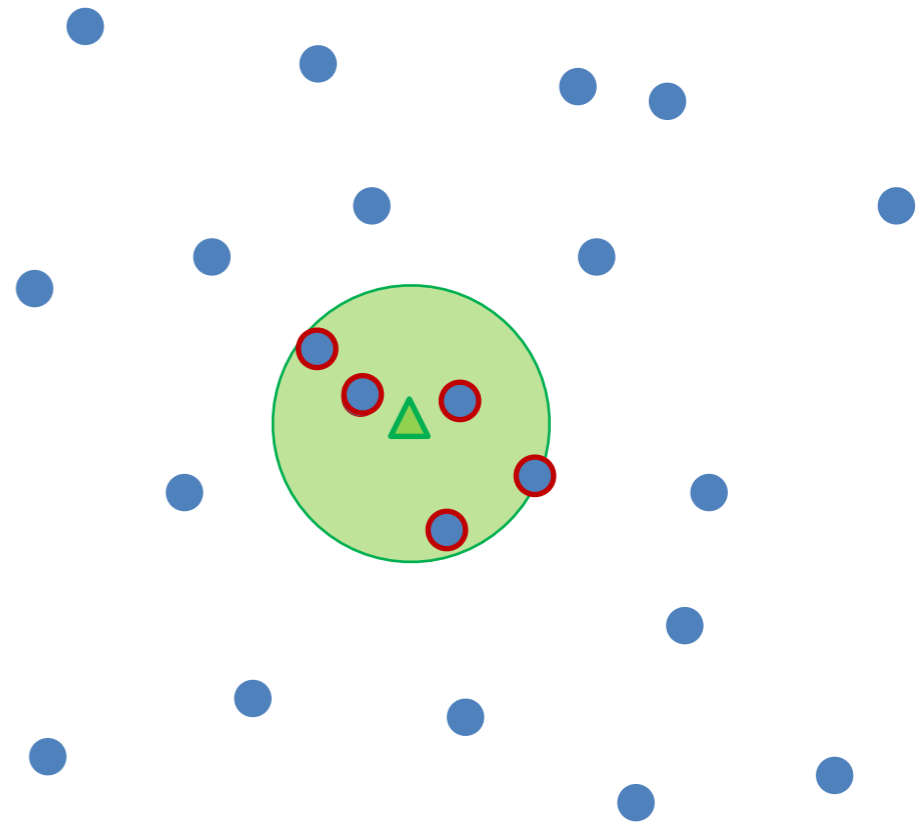
All Mapping



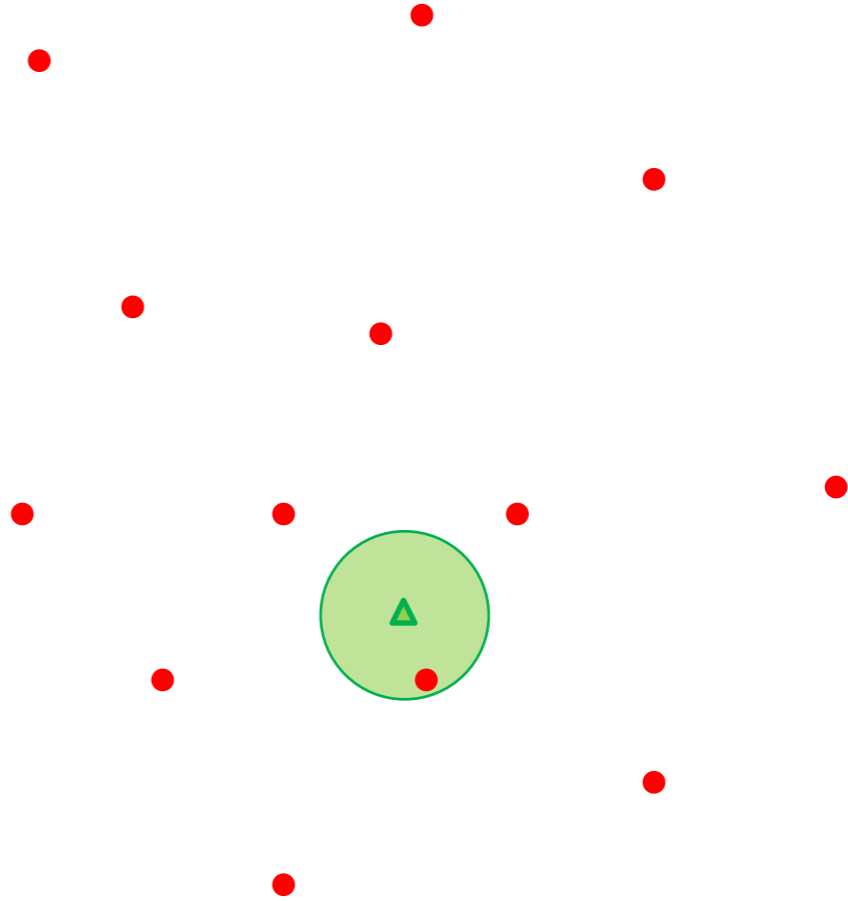
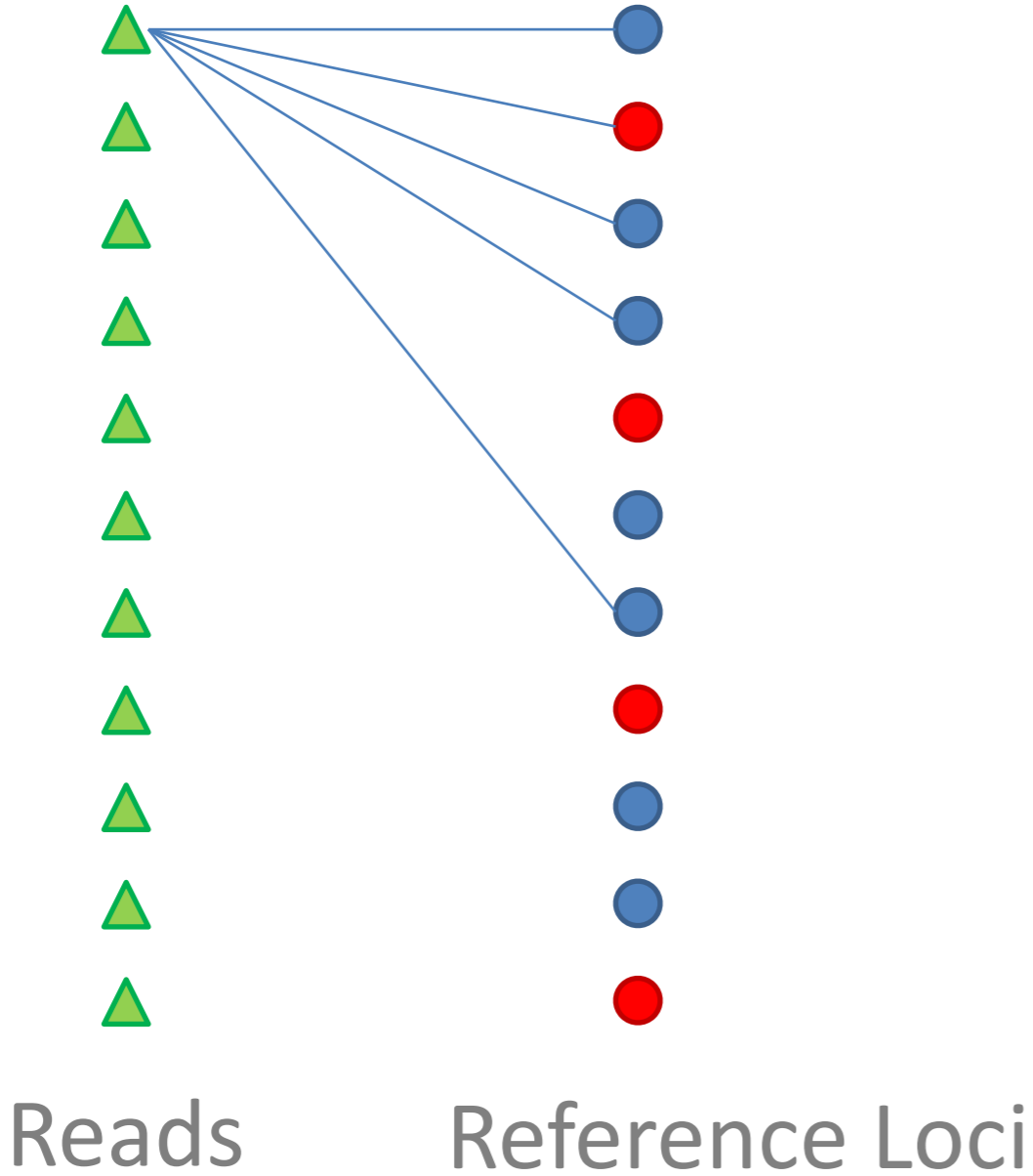
NGS-Mapping in Personal Genomics Era



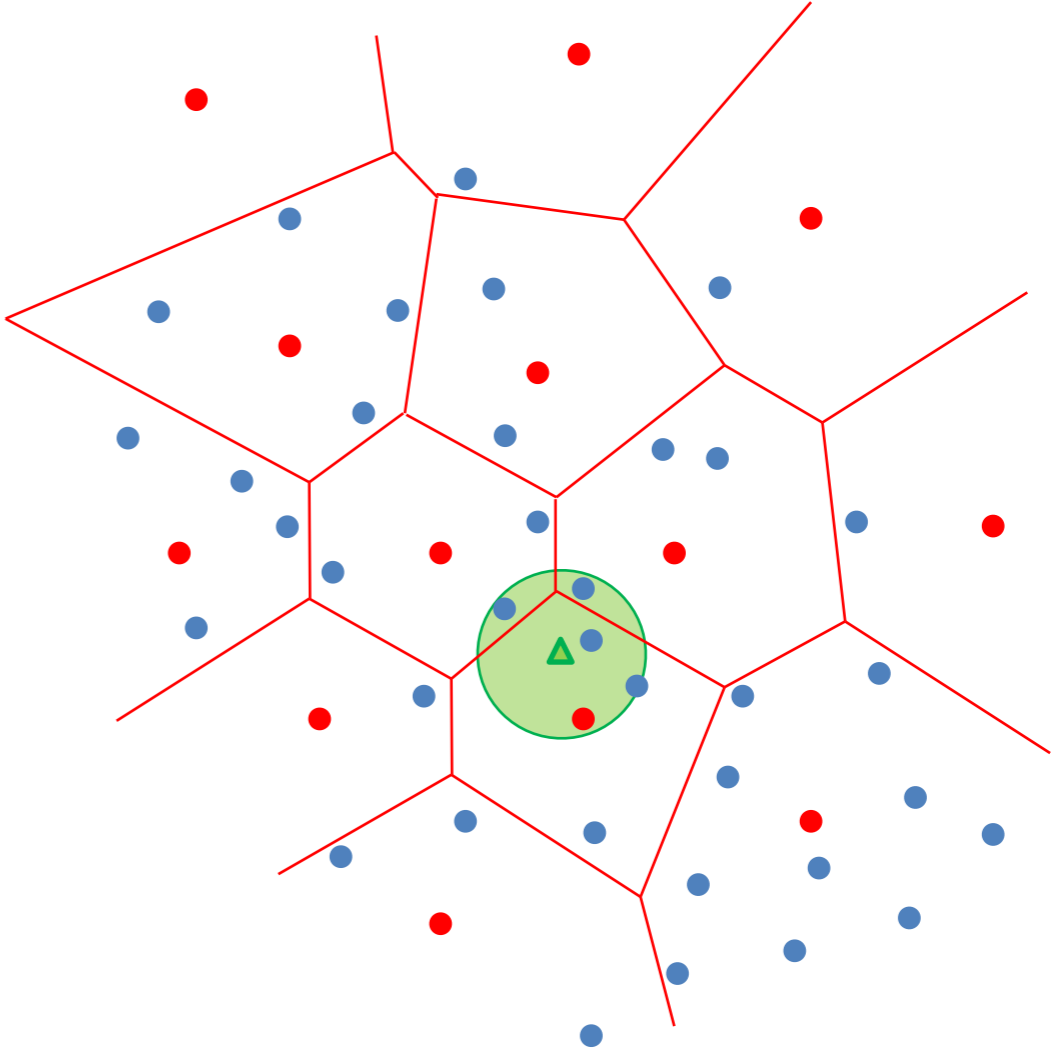
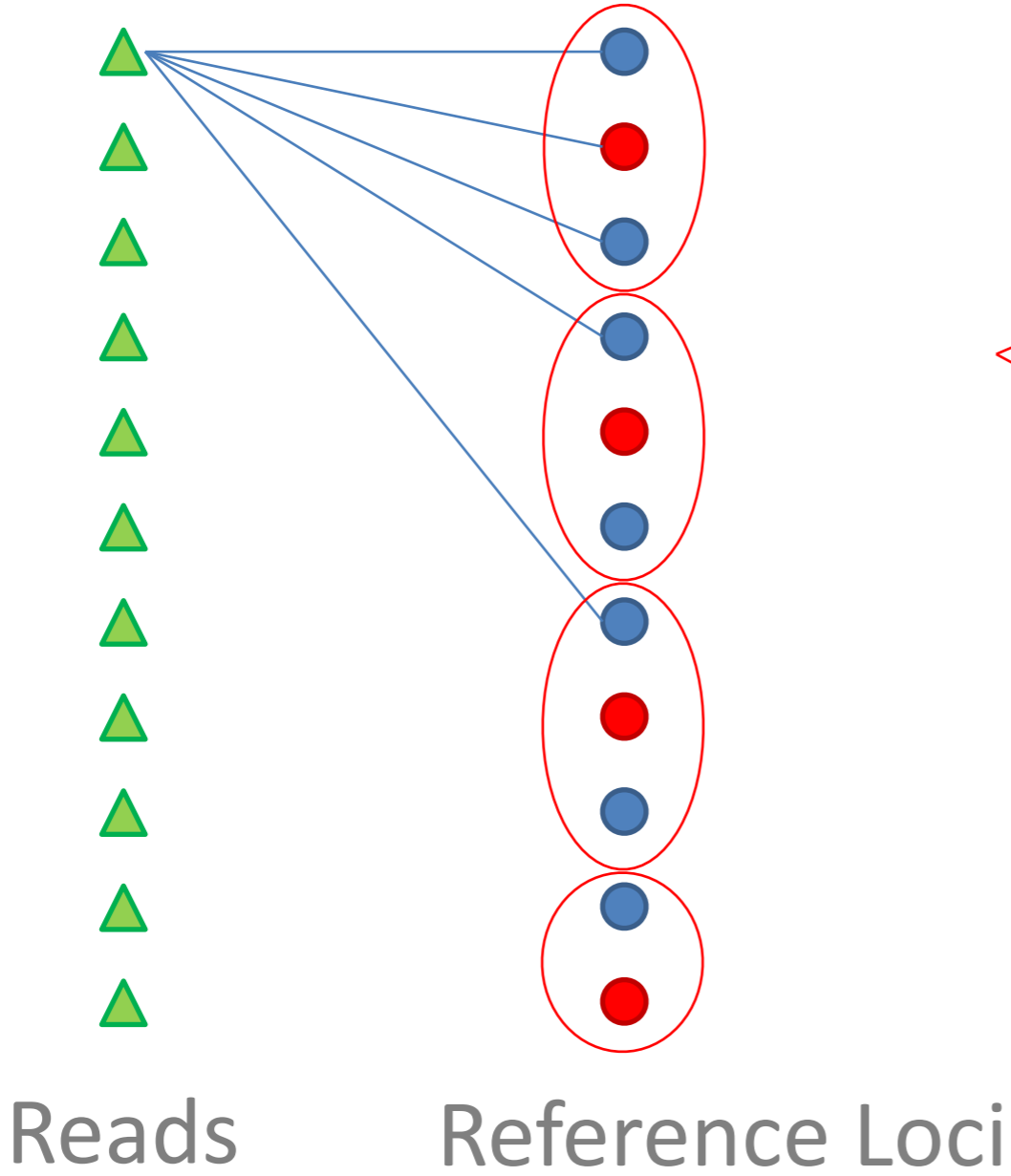
All Mapping



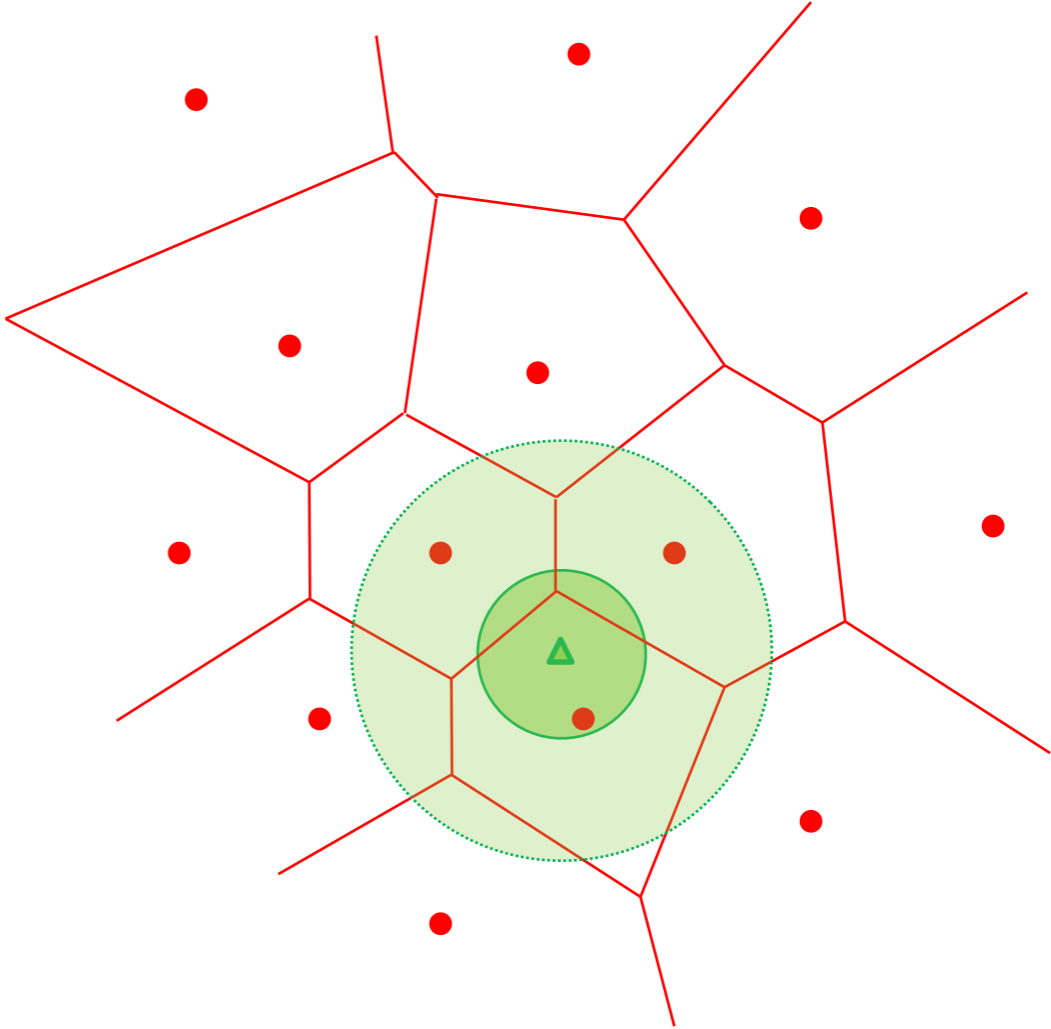
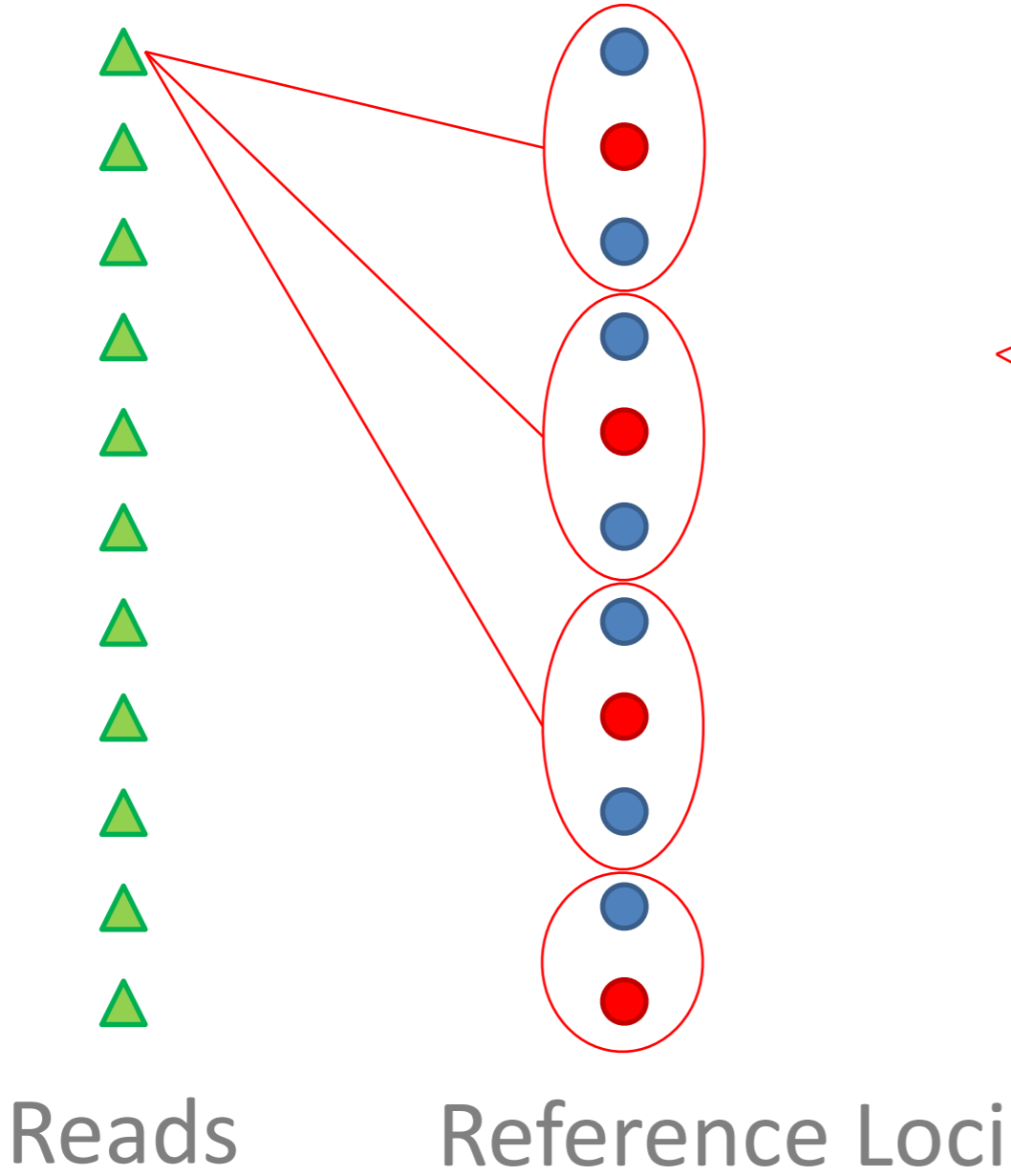
CORA [Yorukoglu, et al. accepted for publication]



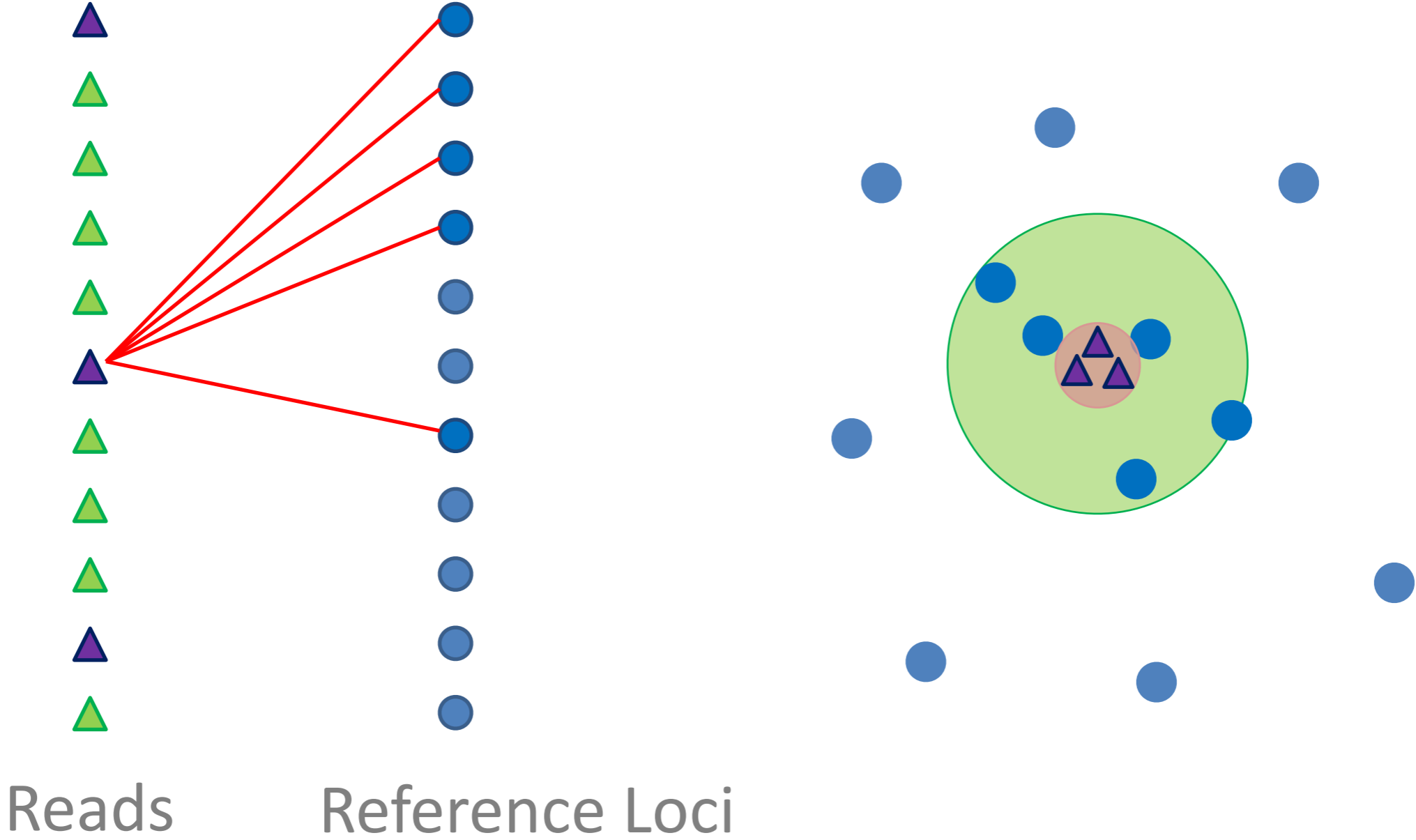
CORA [Yorukoglu, et al. accepted for publication]



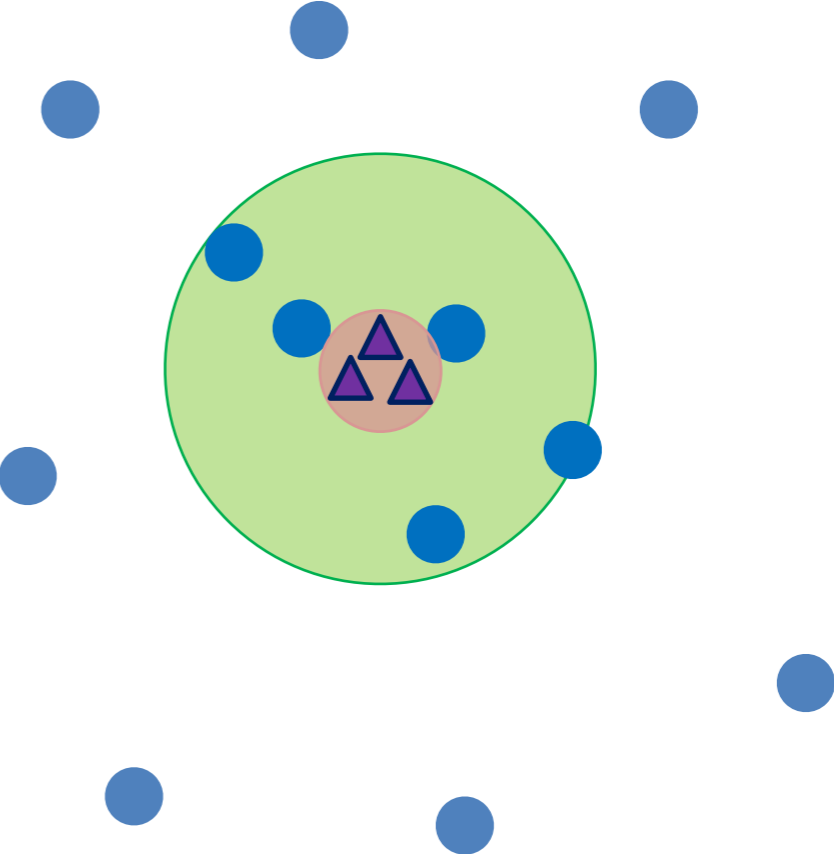
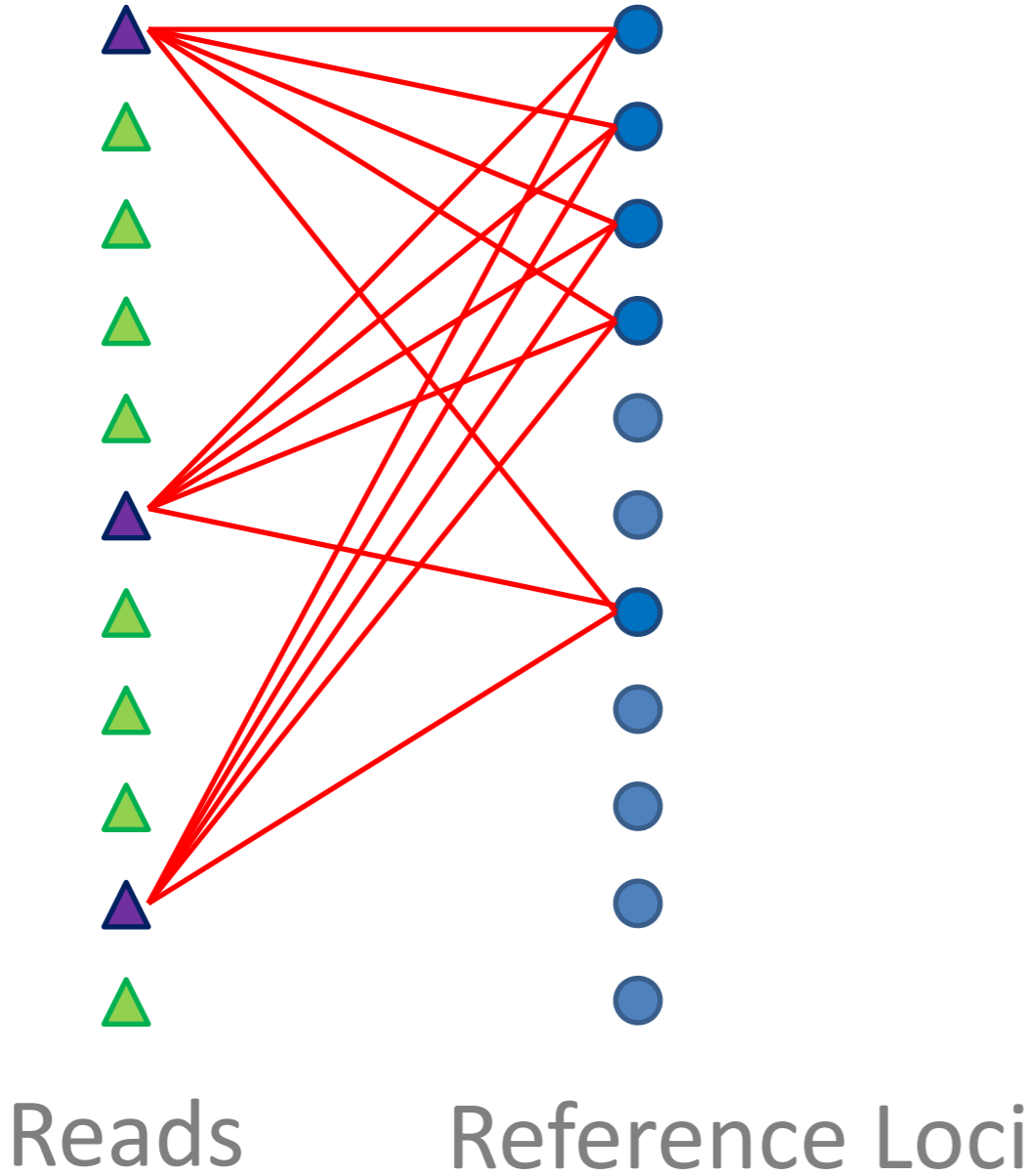
CORA [Yorukoglu, et al. accepted for publication]



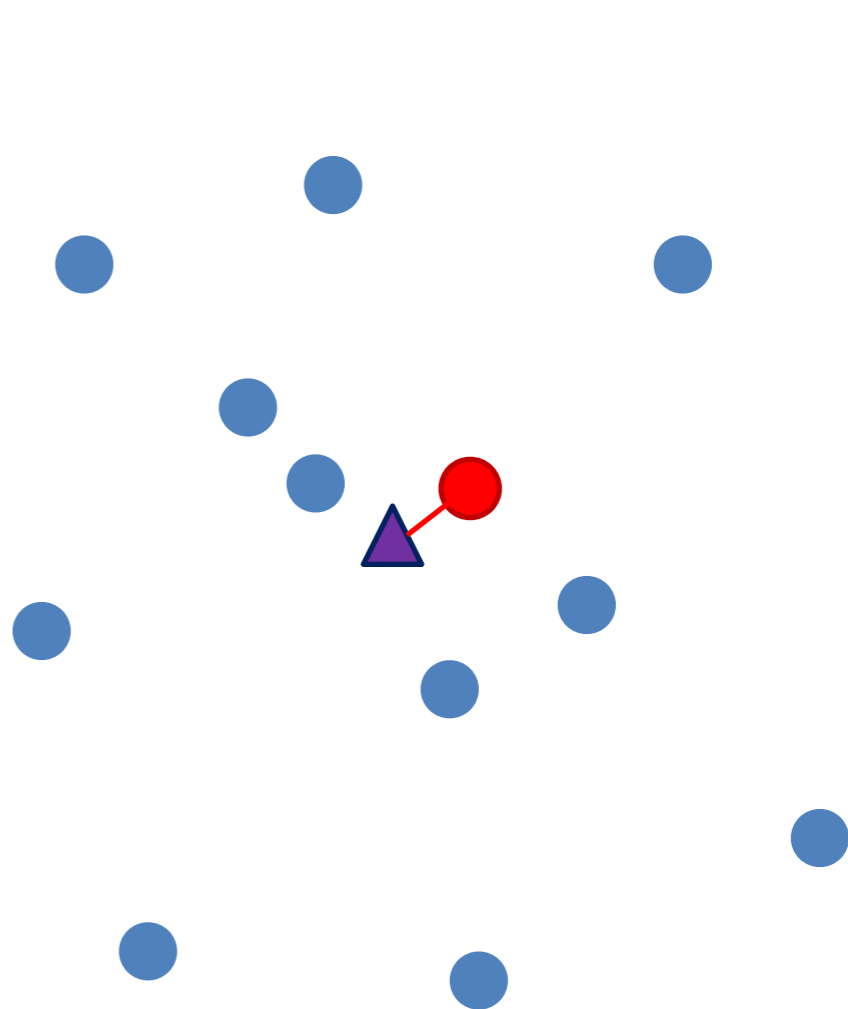
CORA [Yorukoglu, et al. accepted for publication]



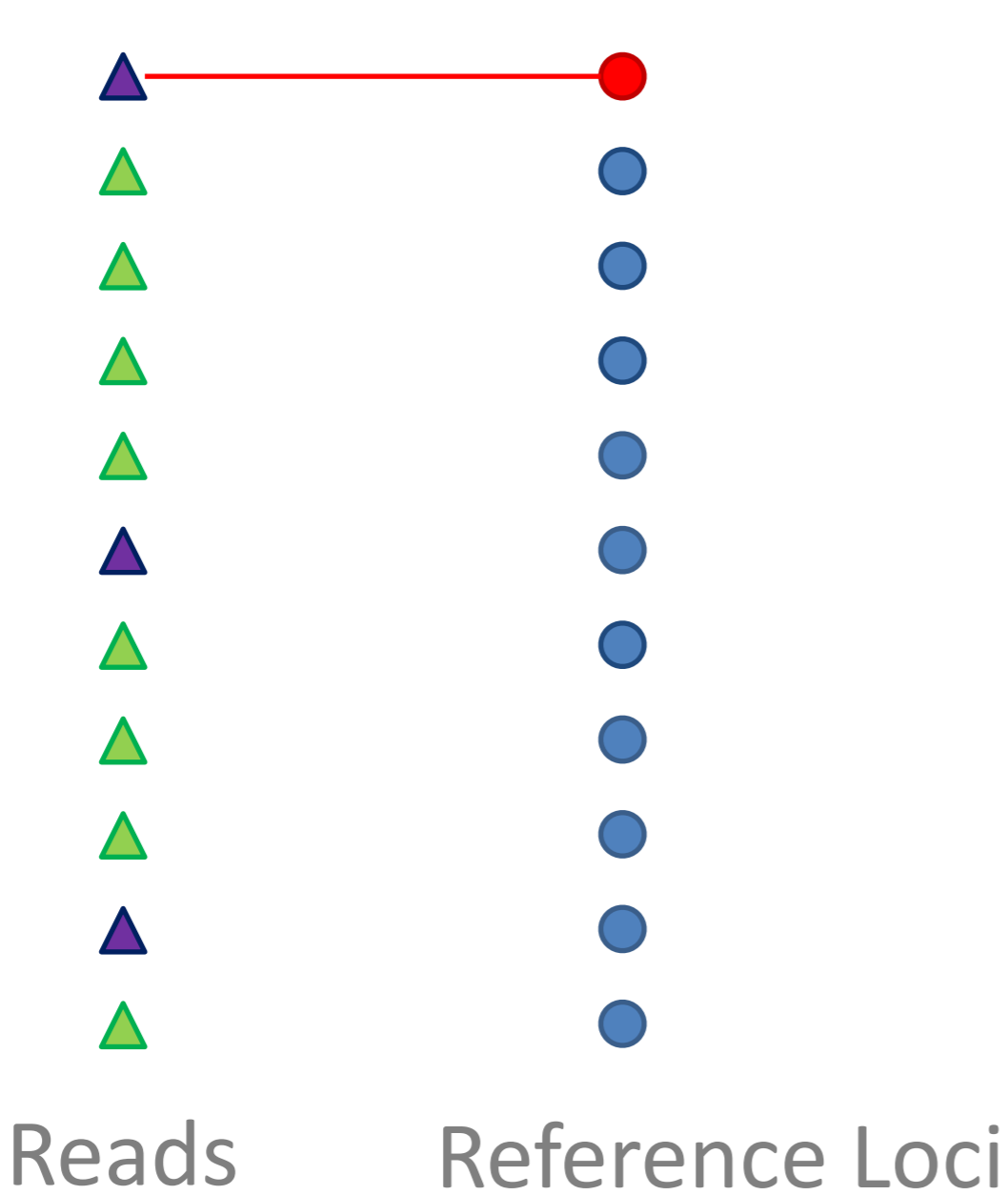
CORA [Yorukoglu, et al. accepted for publication]



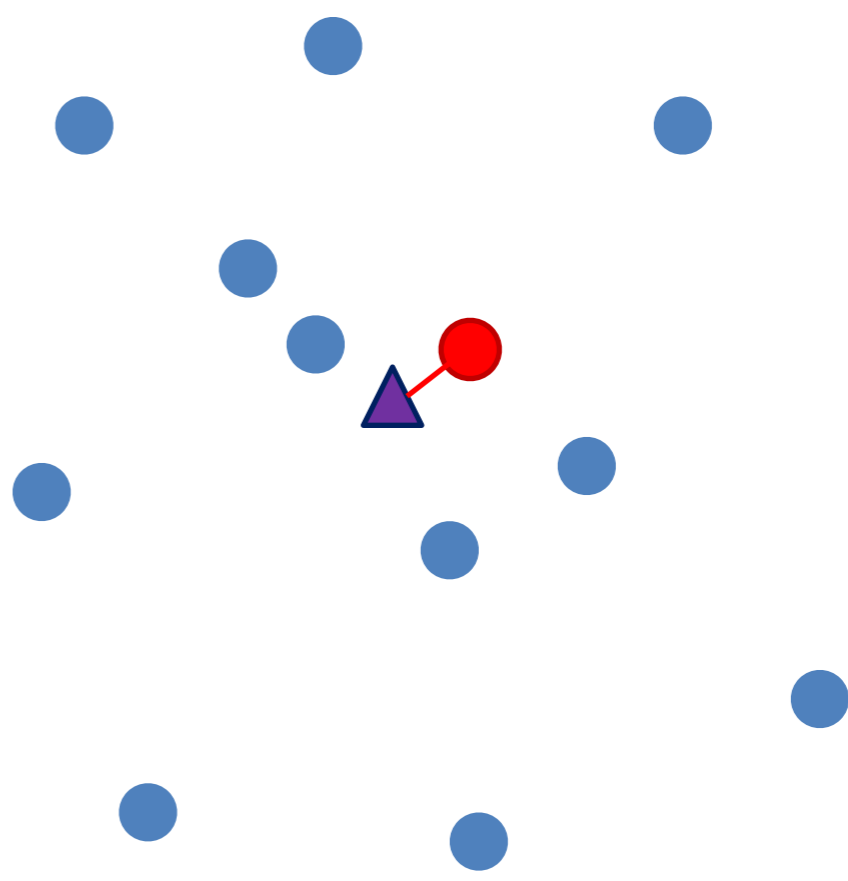
CORA [Yorukoglu, et al. accepted for publication]



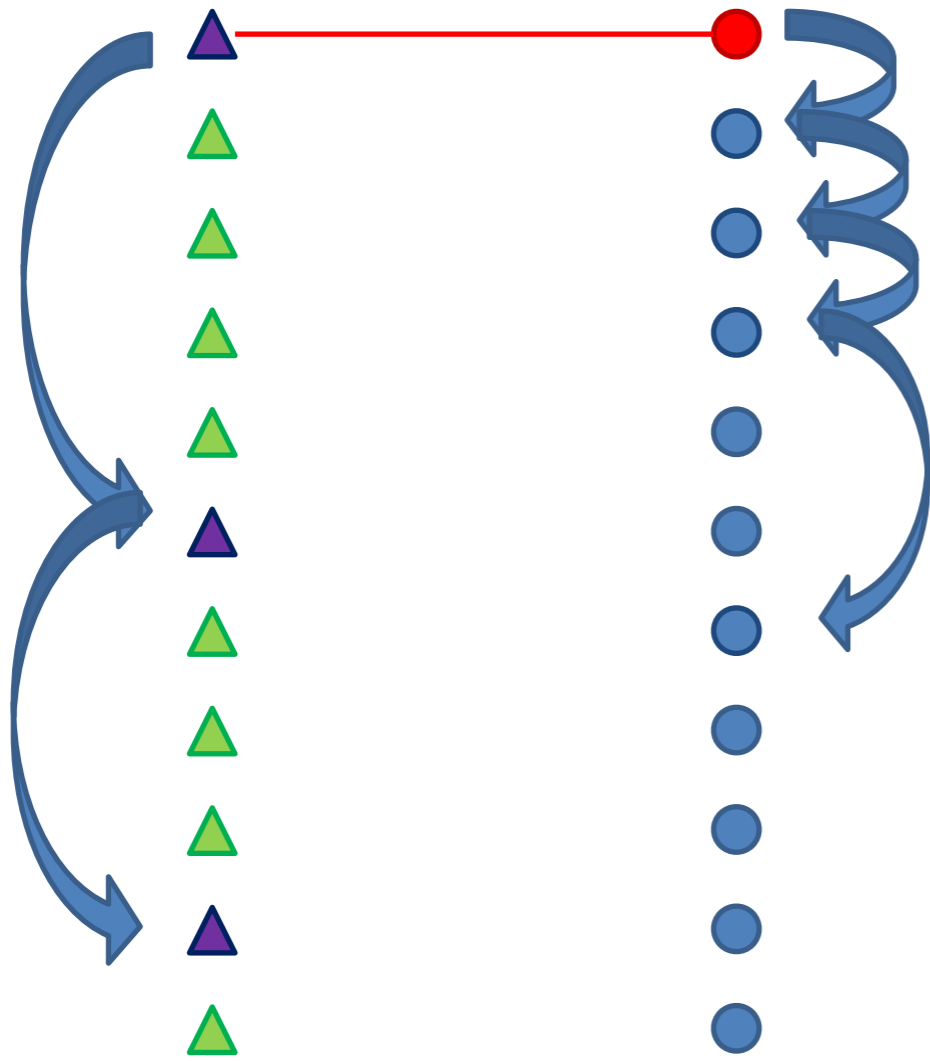
CORA [Yorukoglu, et al. accepted for publication]



Coarse Mapping



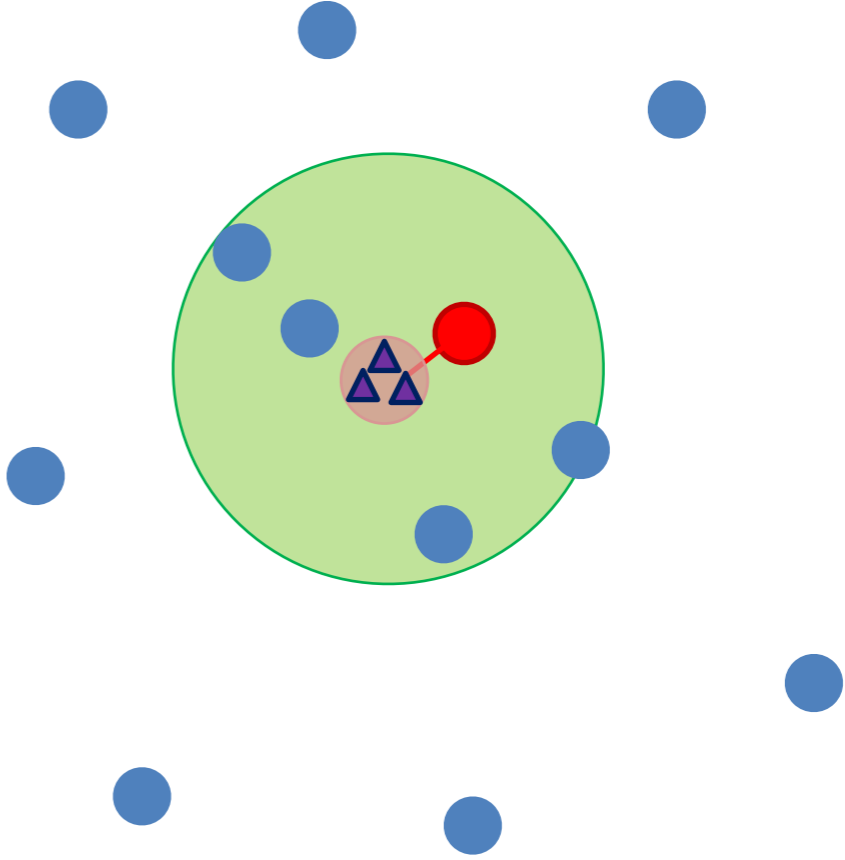
CORA [Yorukoglu, et al. accepted for publication]



Reads

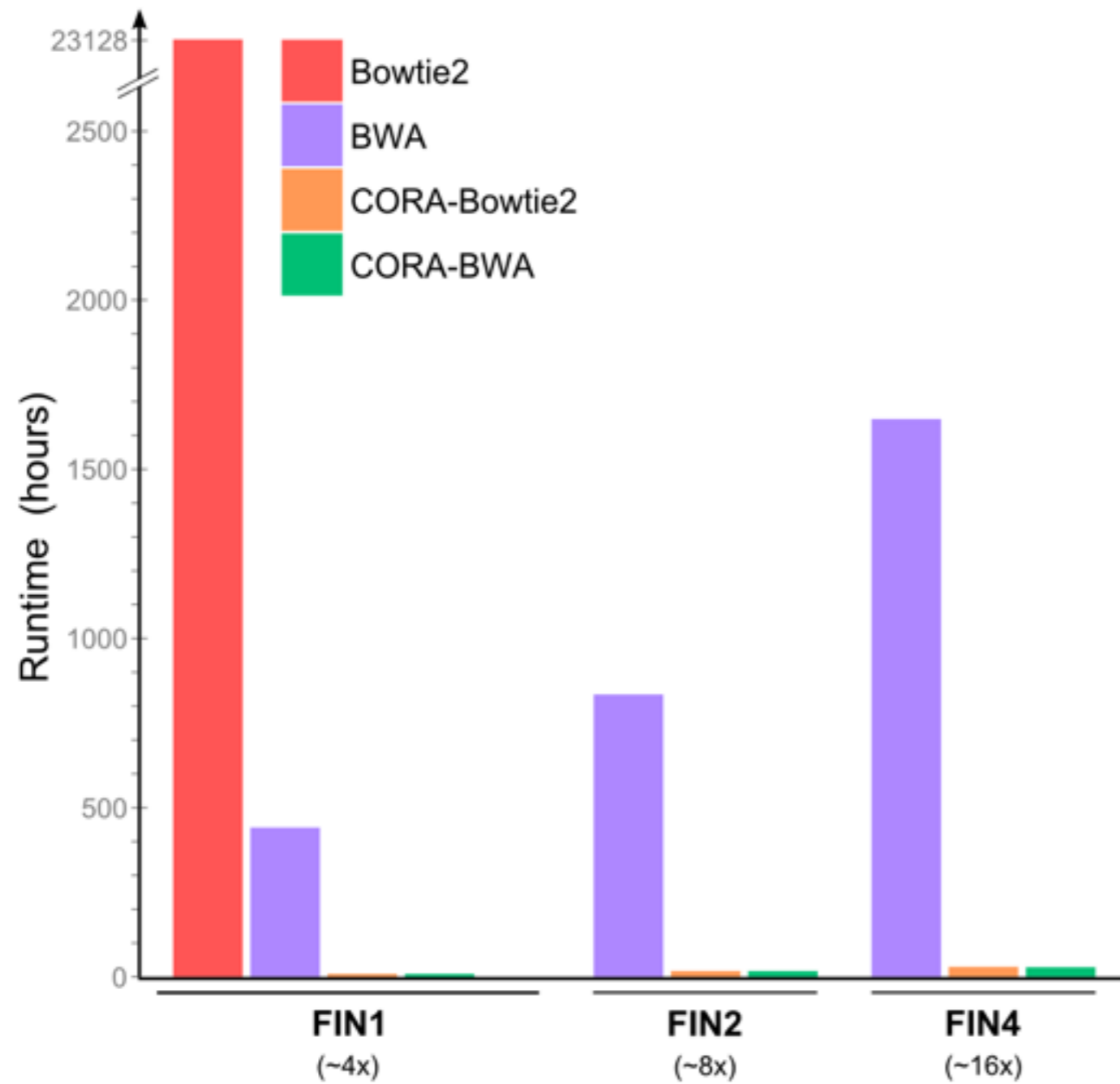
Reference Loci

Fine Mapping



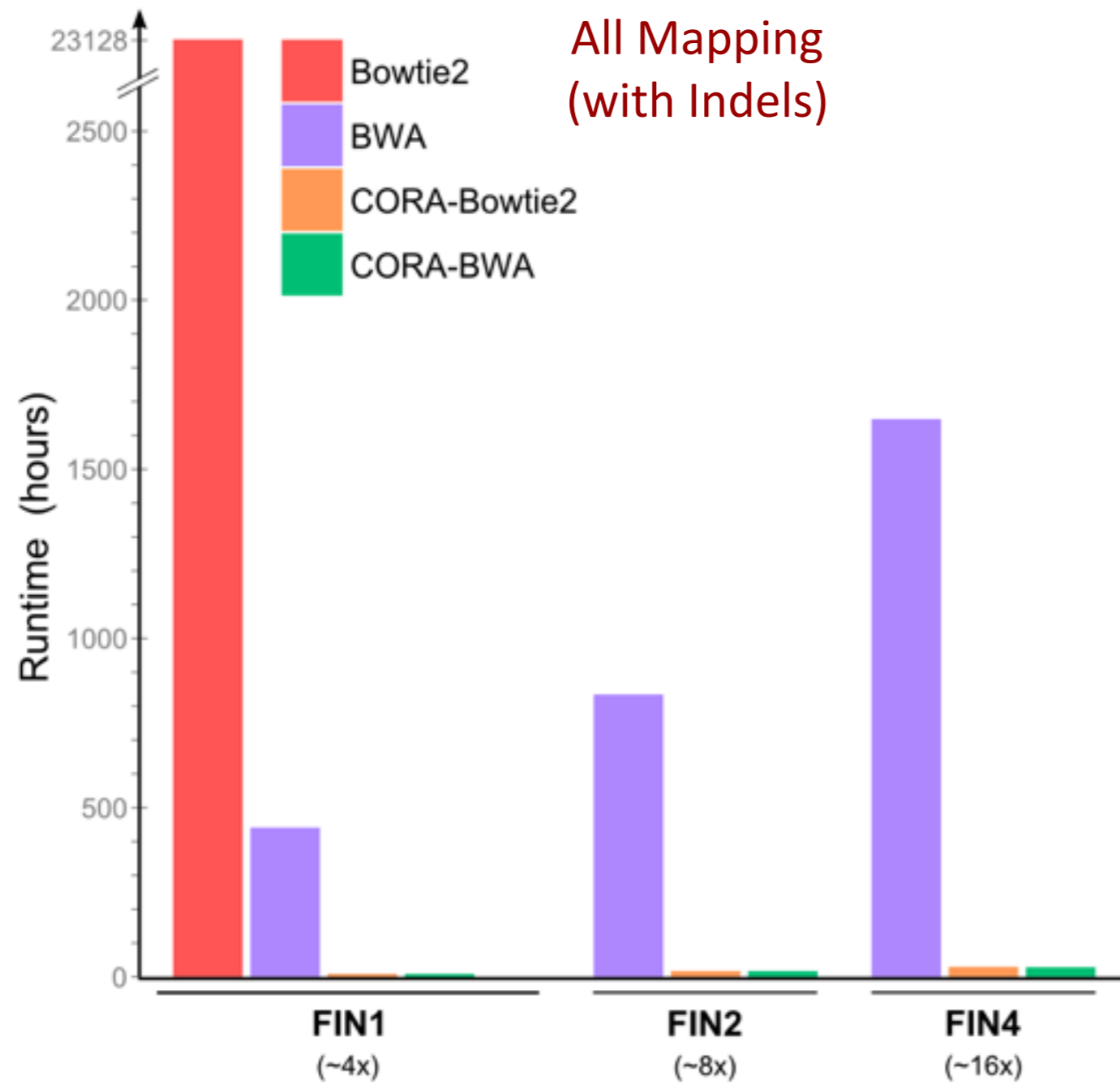
Result Highlights

CORA: 1000 Genomes Project



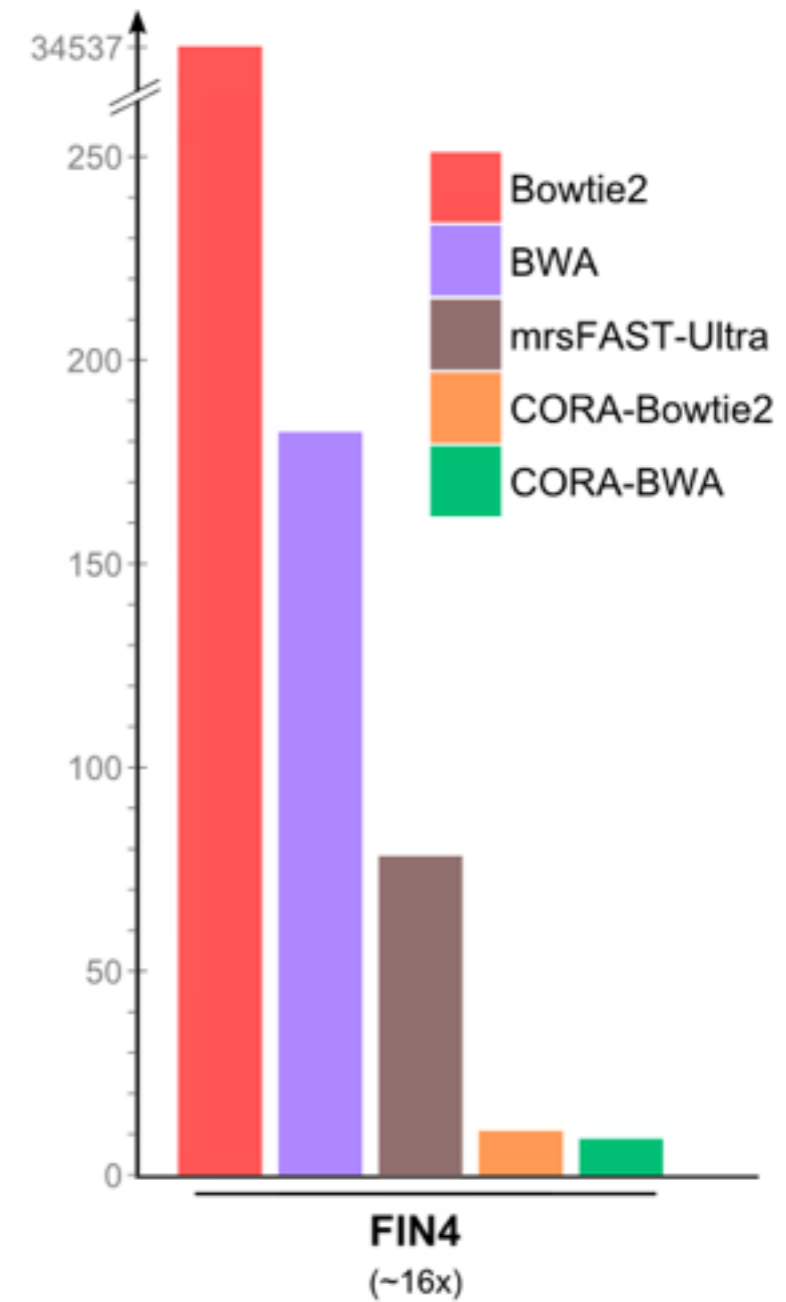
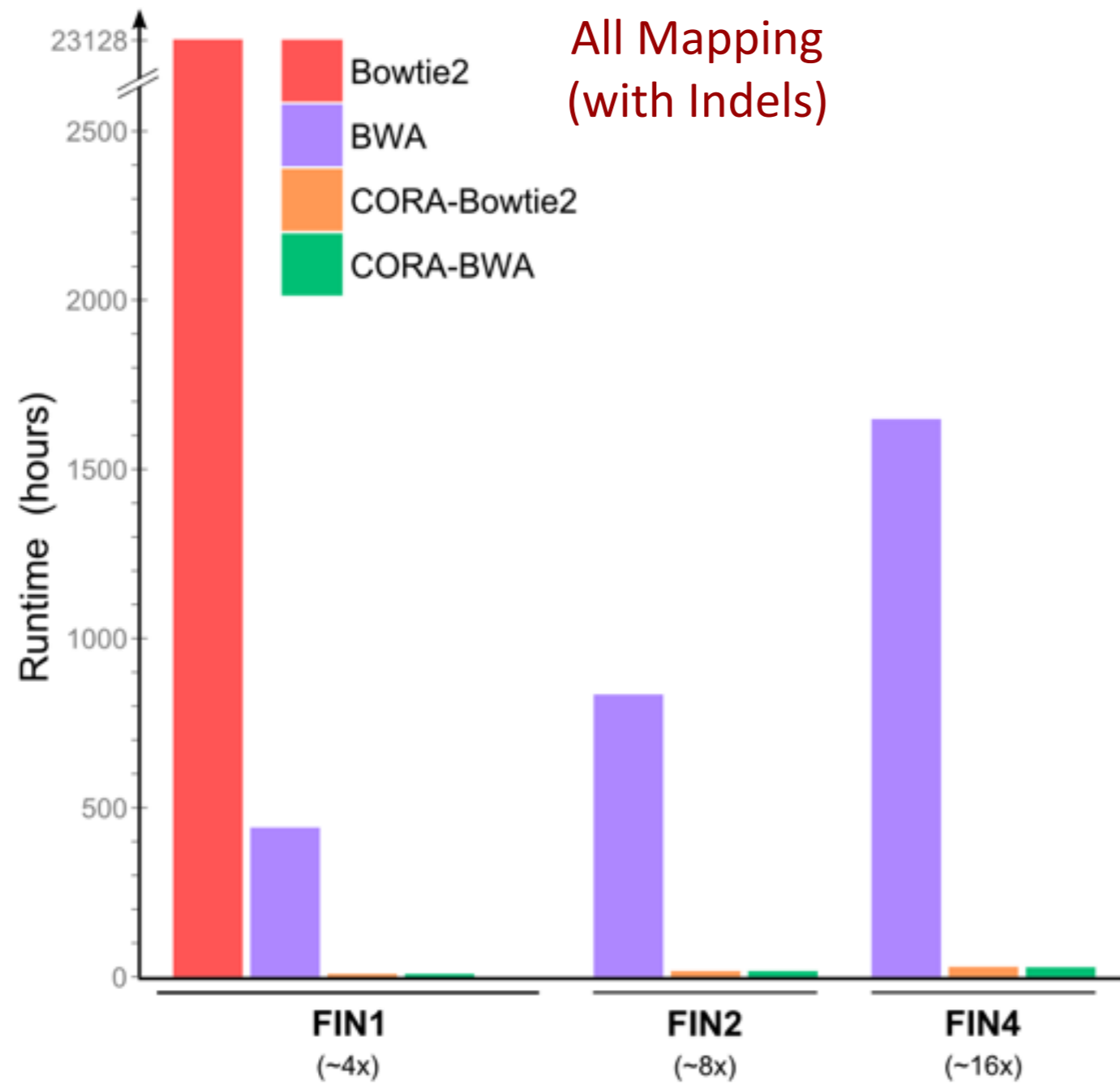
Result Highlights

CORA: 1000 Genomes Project



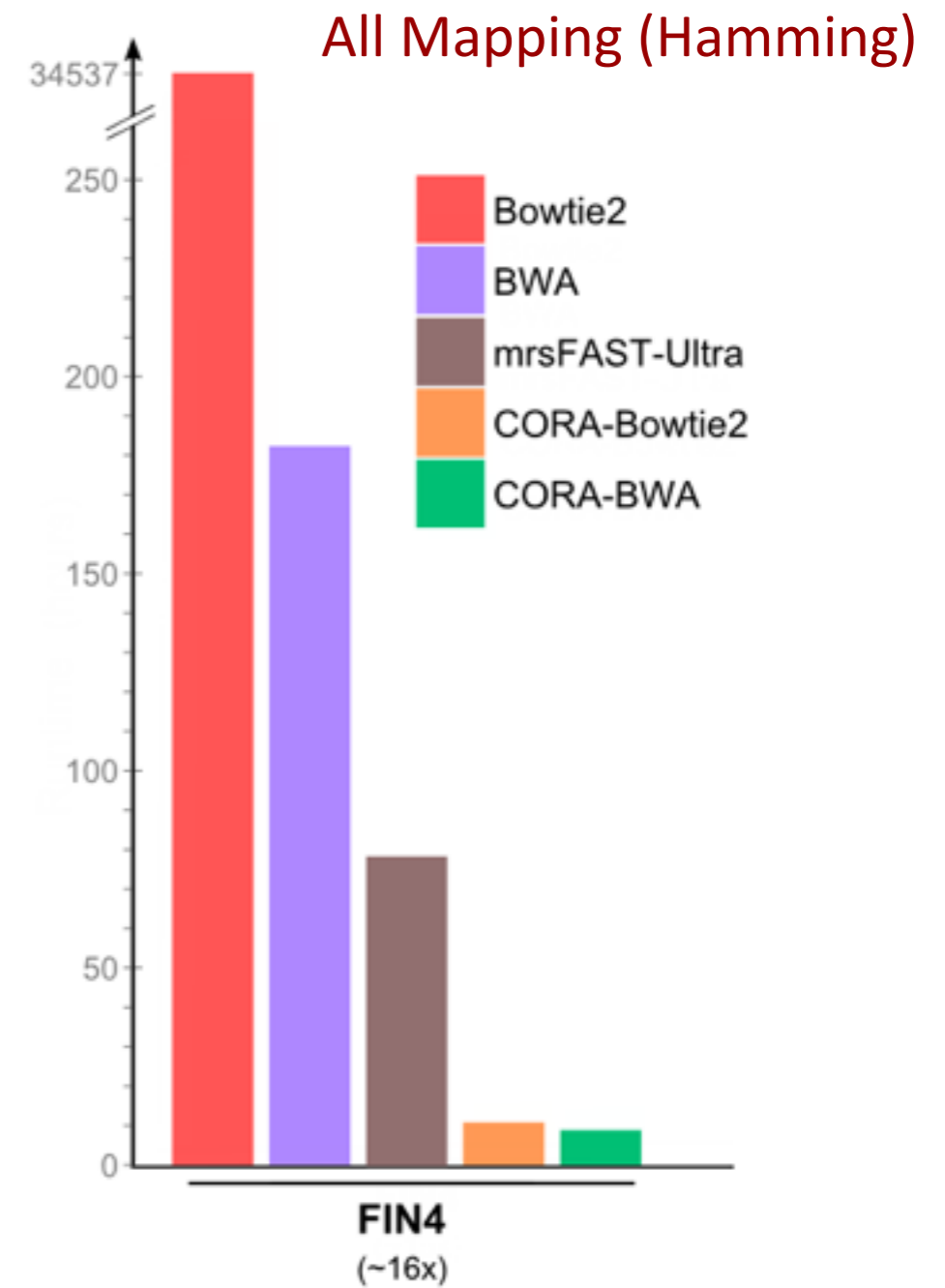
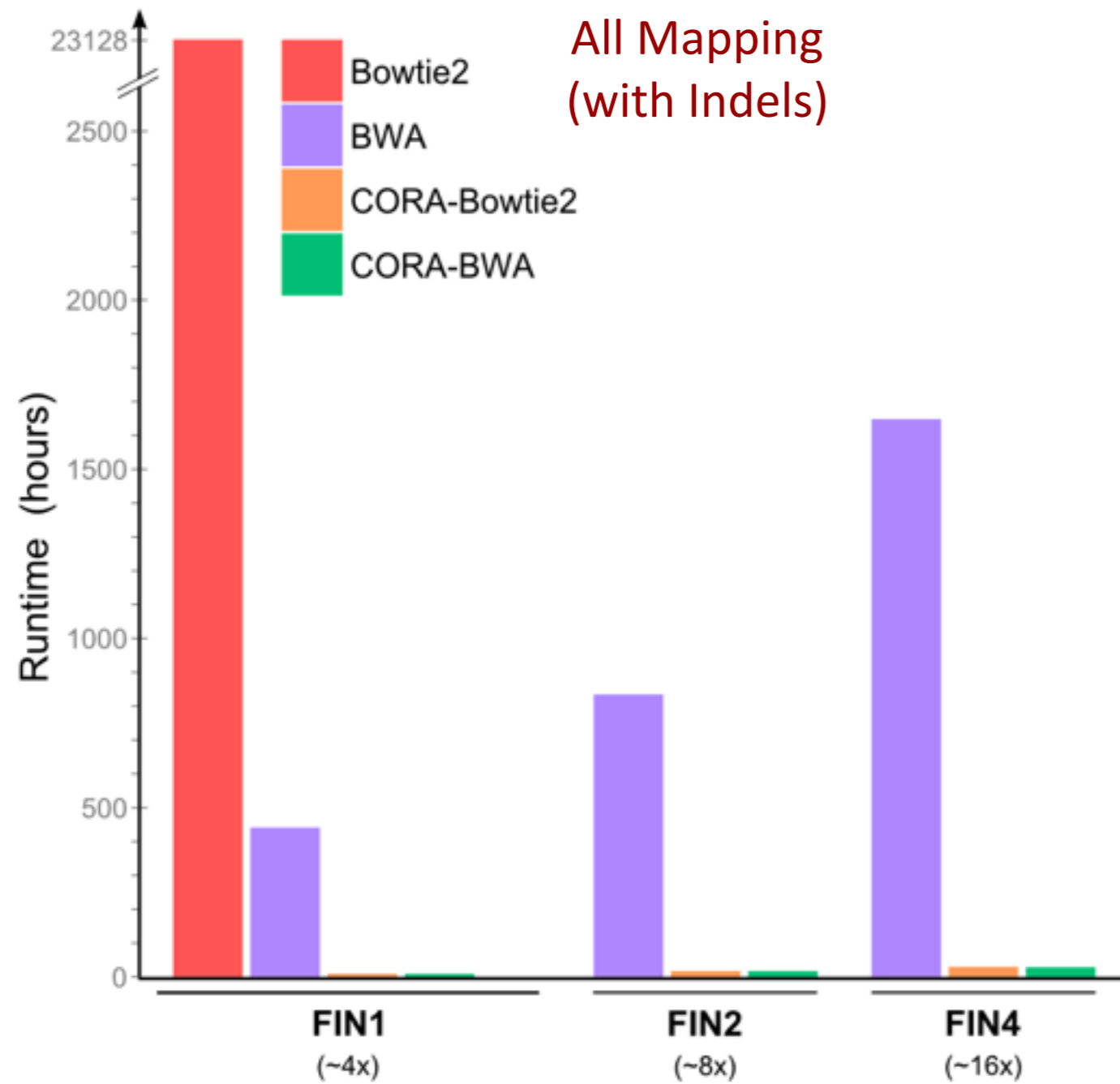
Result Highlights

CORA: 1000 Genomes Project



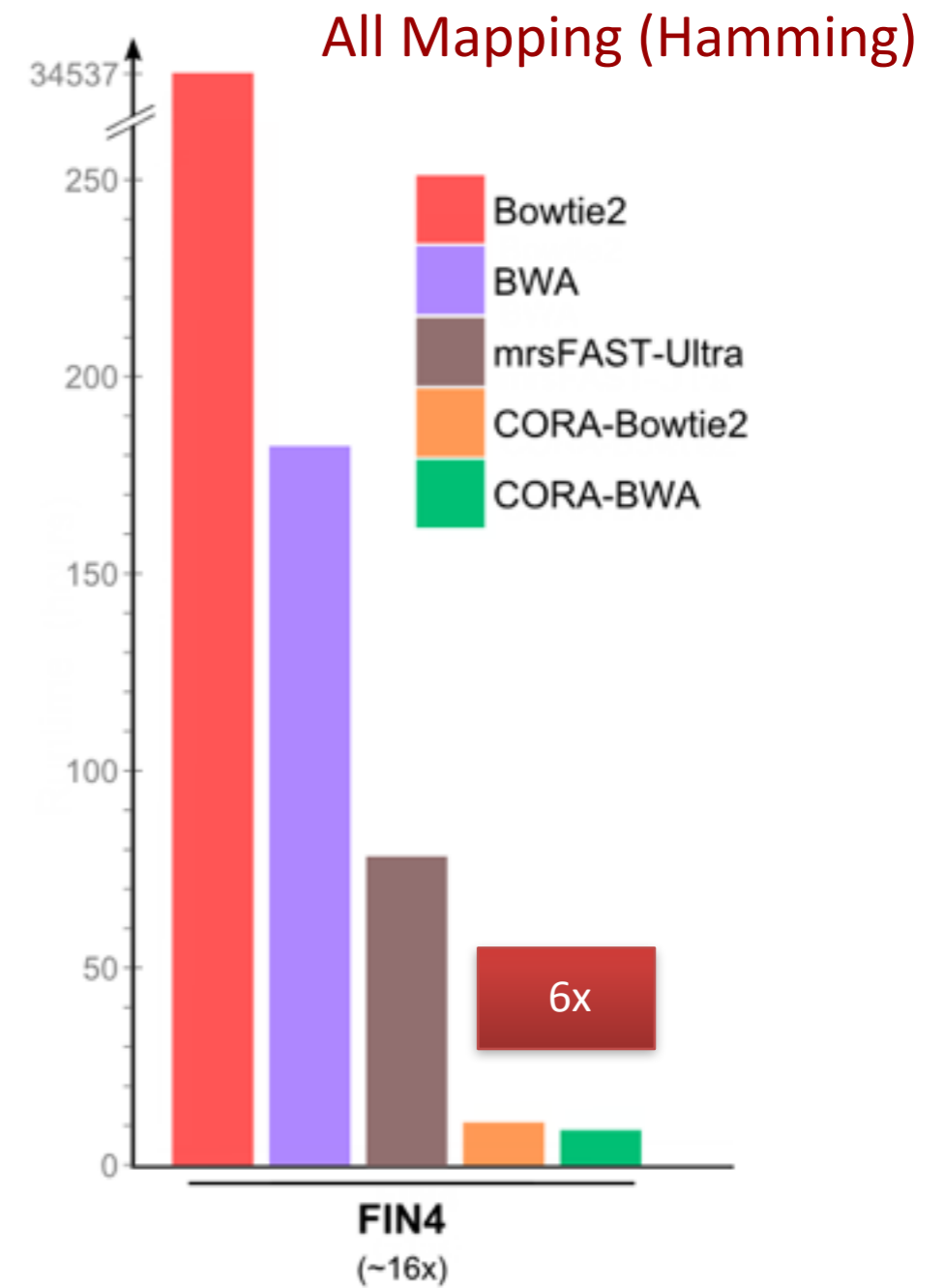
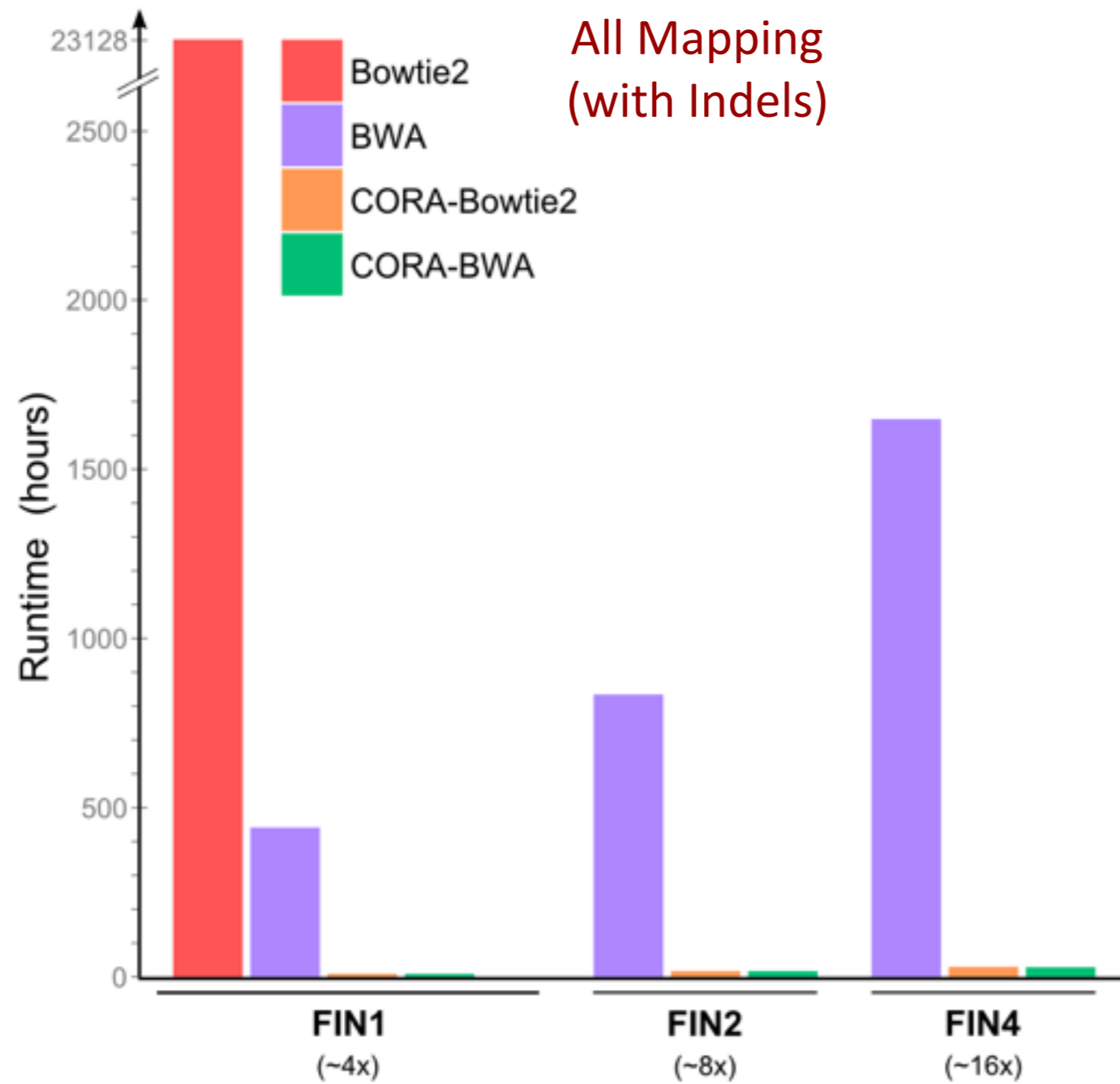
Result Highlights

CORA: 1000 Genomes Project



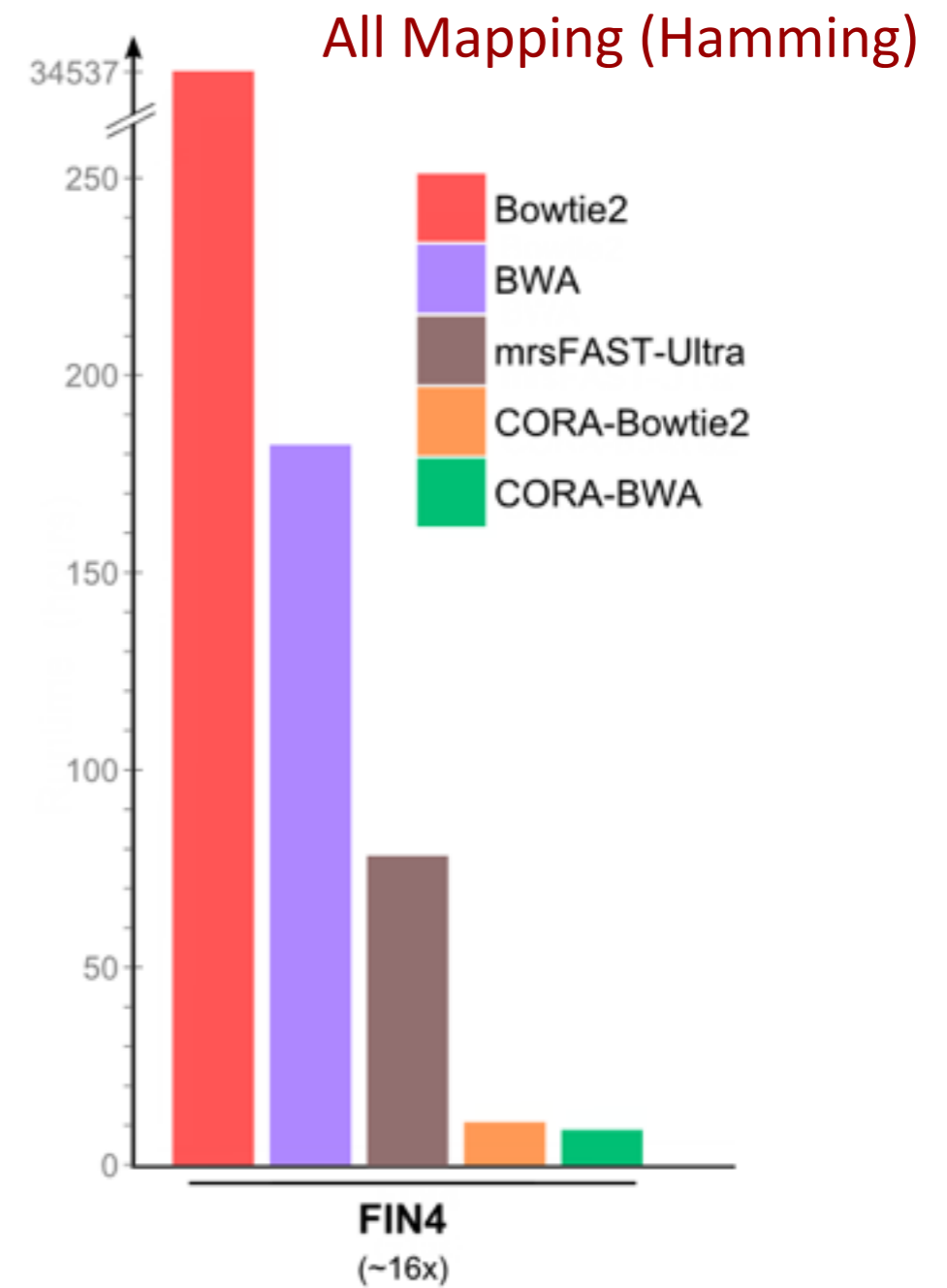
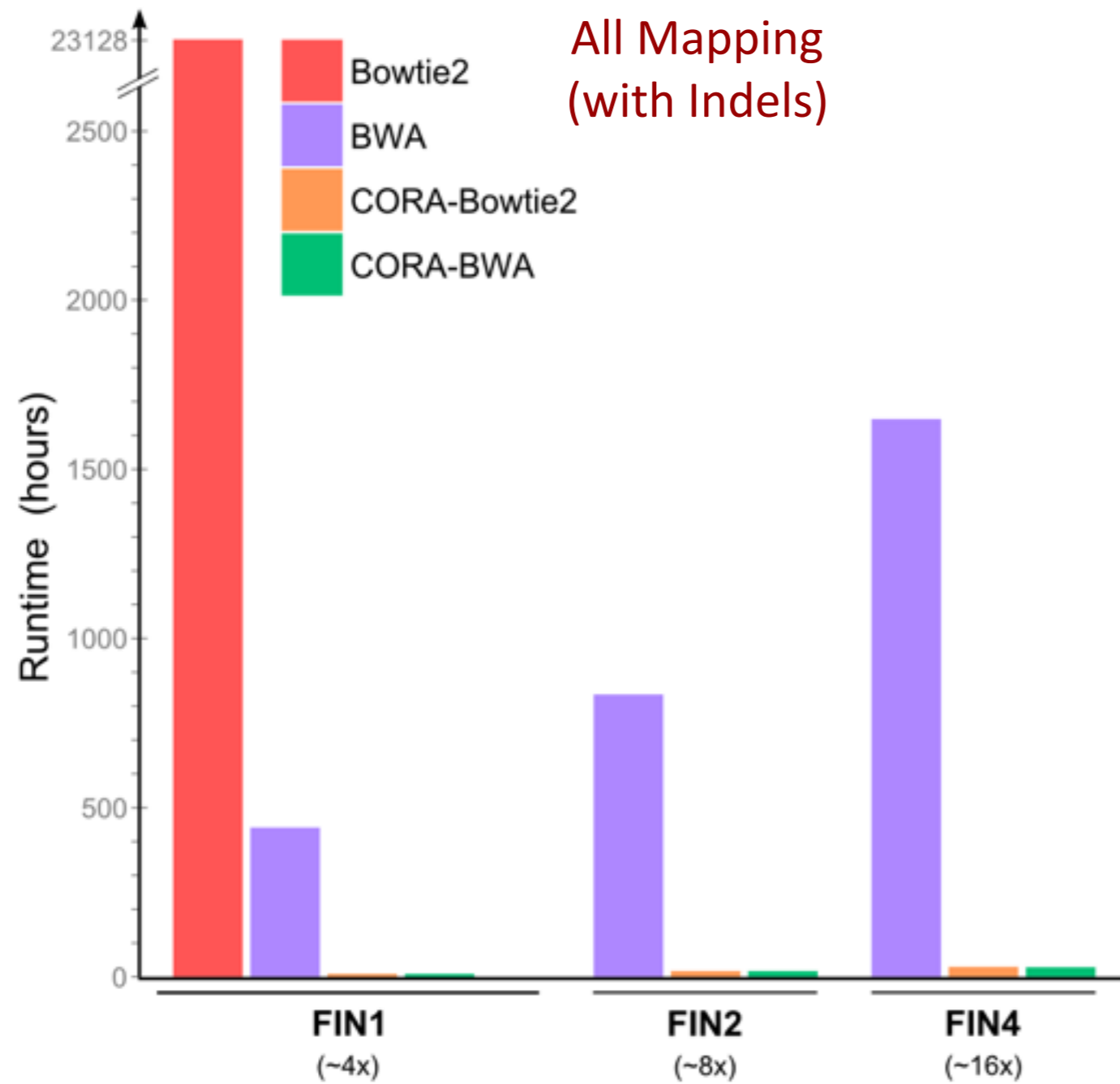
Result Highlights

CORA: 1000 Genomes Project



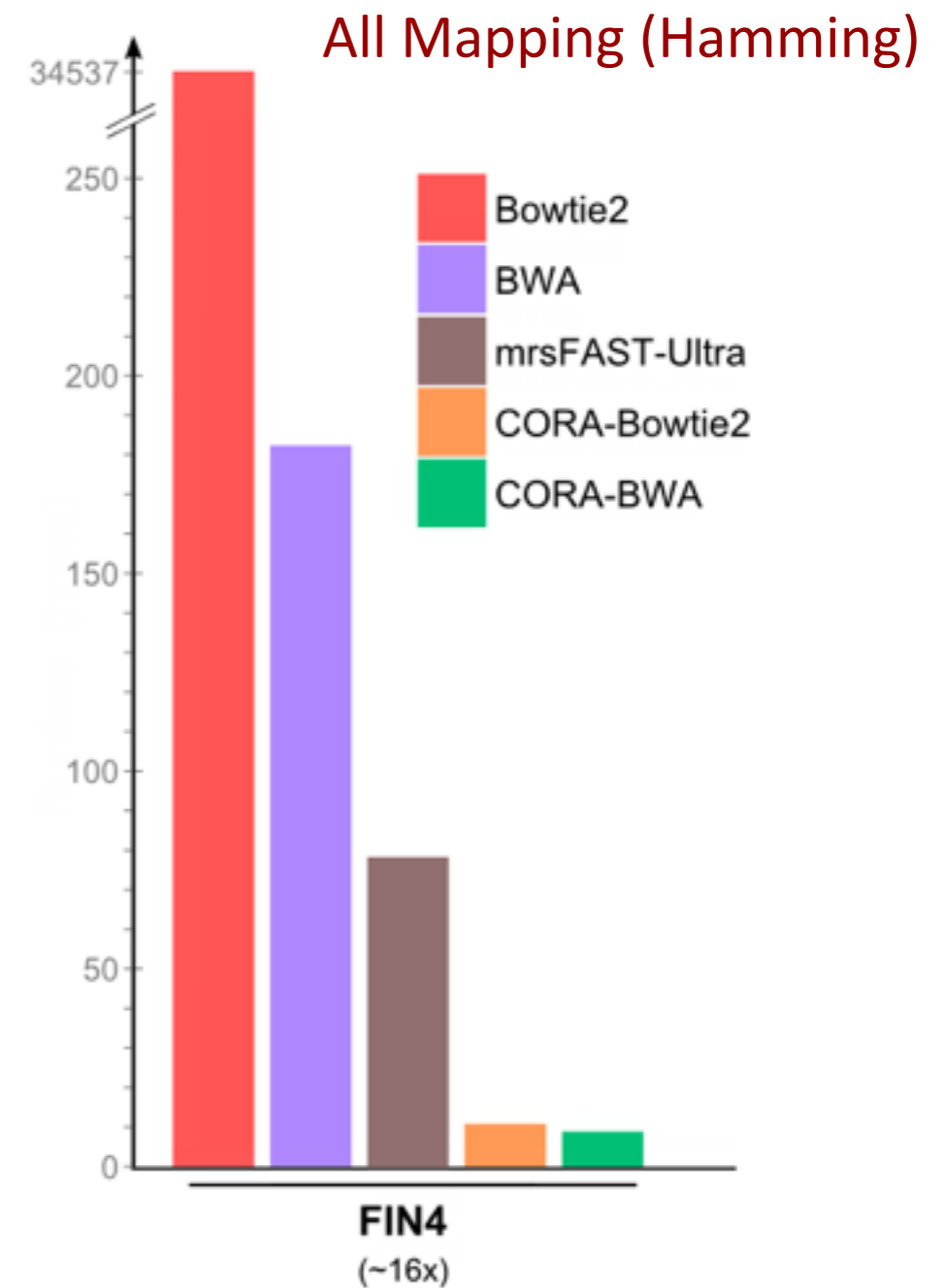
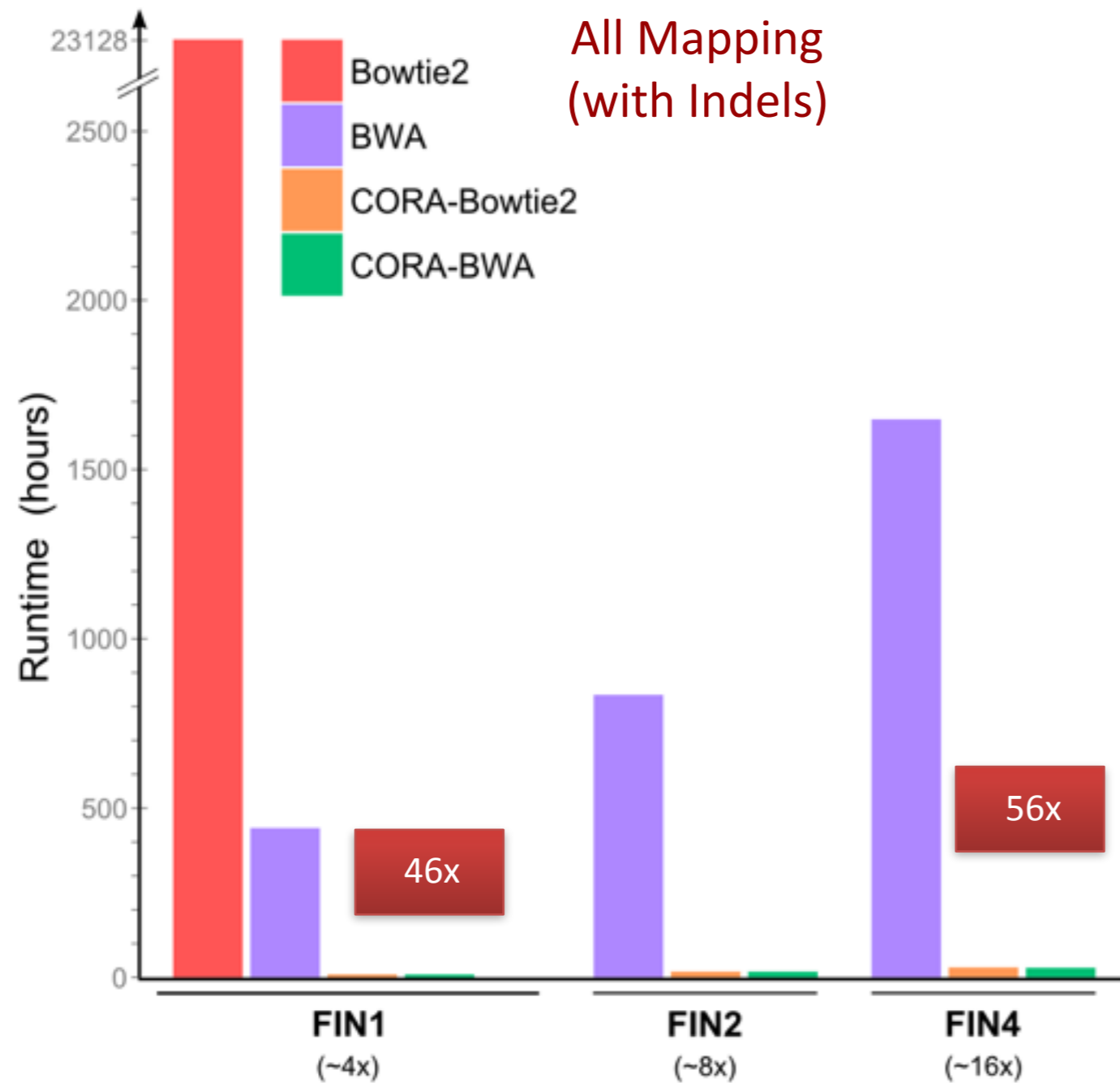
Result Highlights

CORA: 1000 Genomes Project



Result Highlights

CORA: 1000 Genomes Project

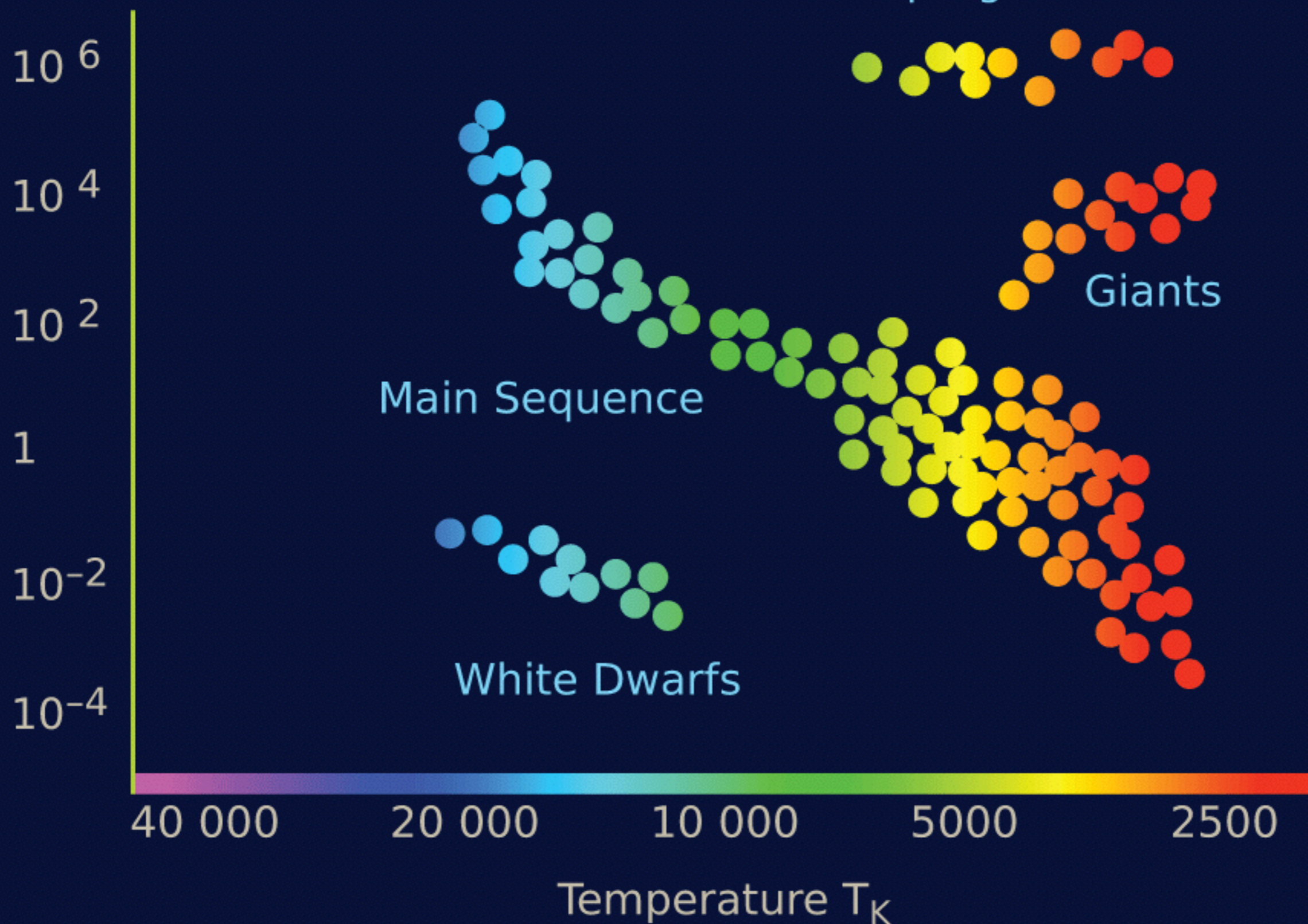


Scales sublinearly with very high sensitivity on real genomic data

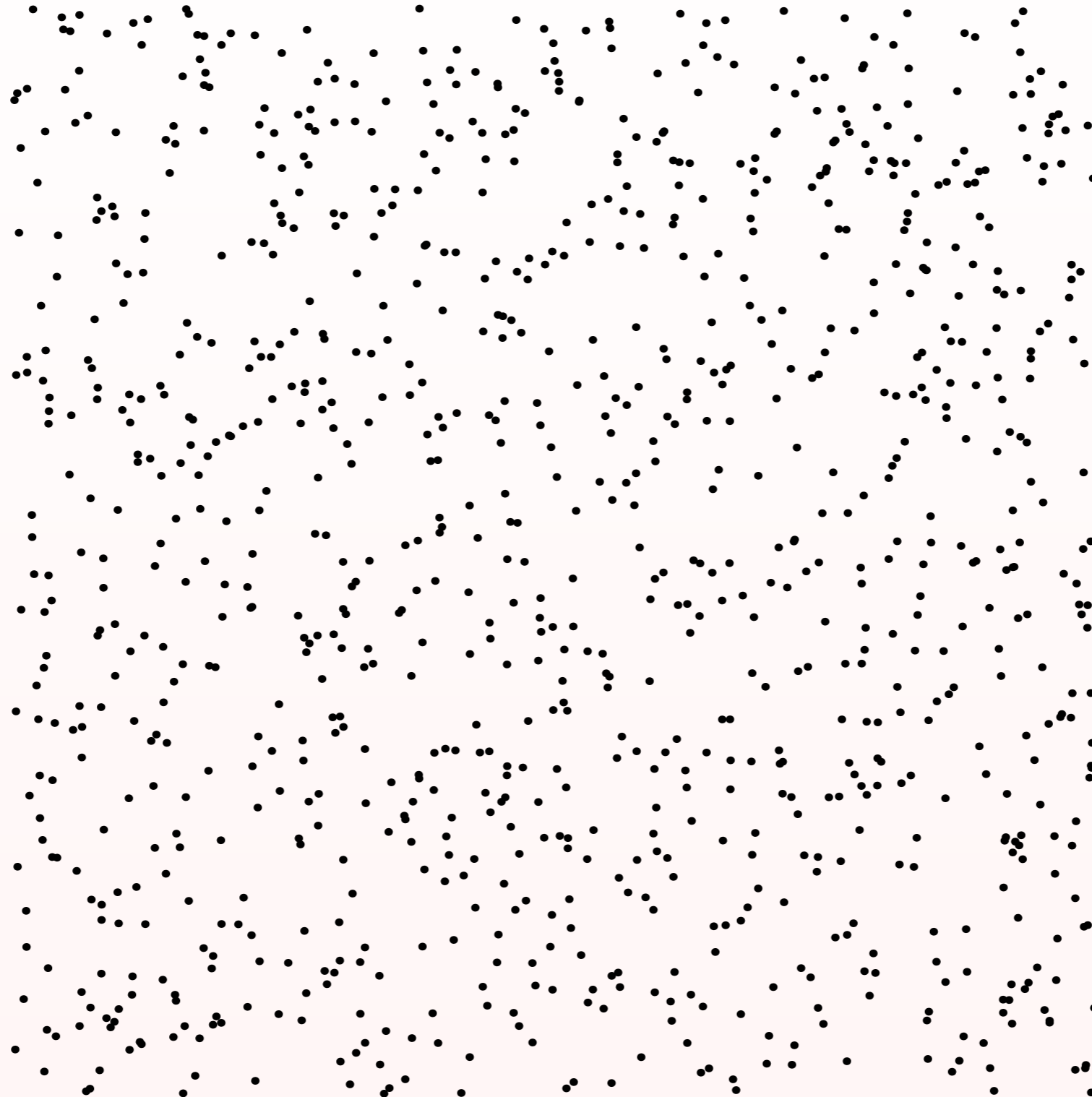
Data constrained by physical process

Hertzprung-Russell Diagram

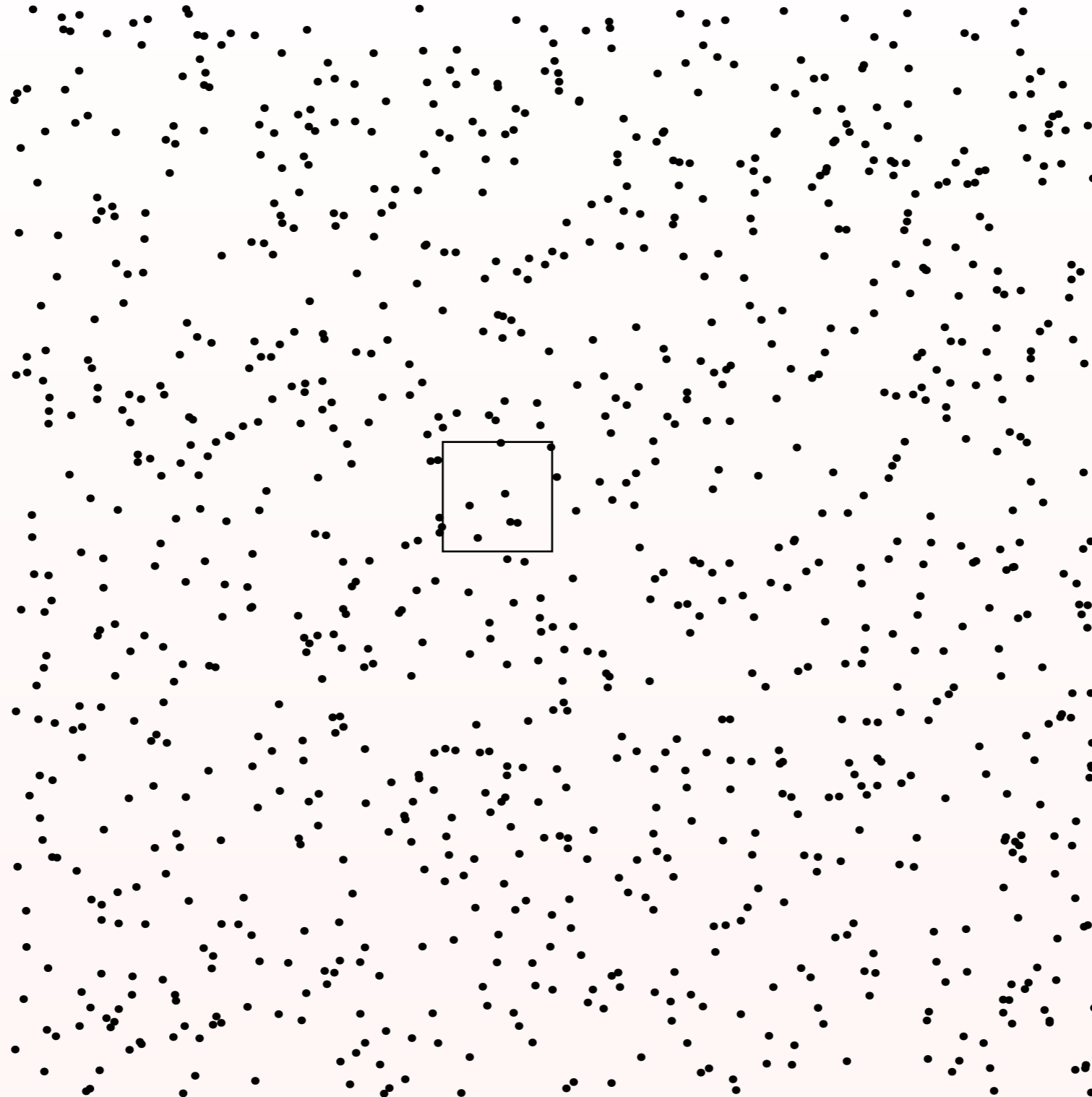
Luminosity, L (L_{Sun})



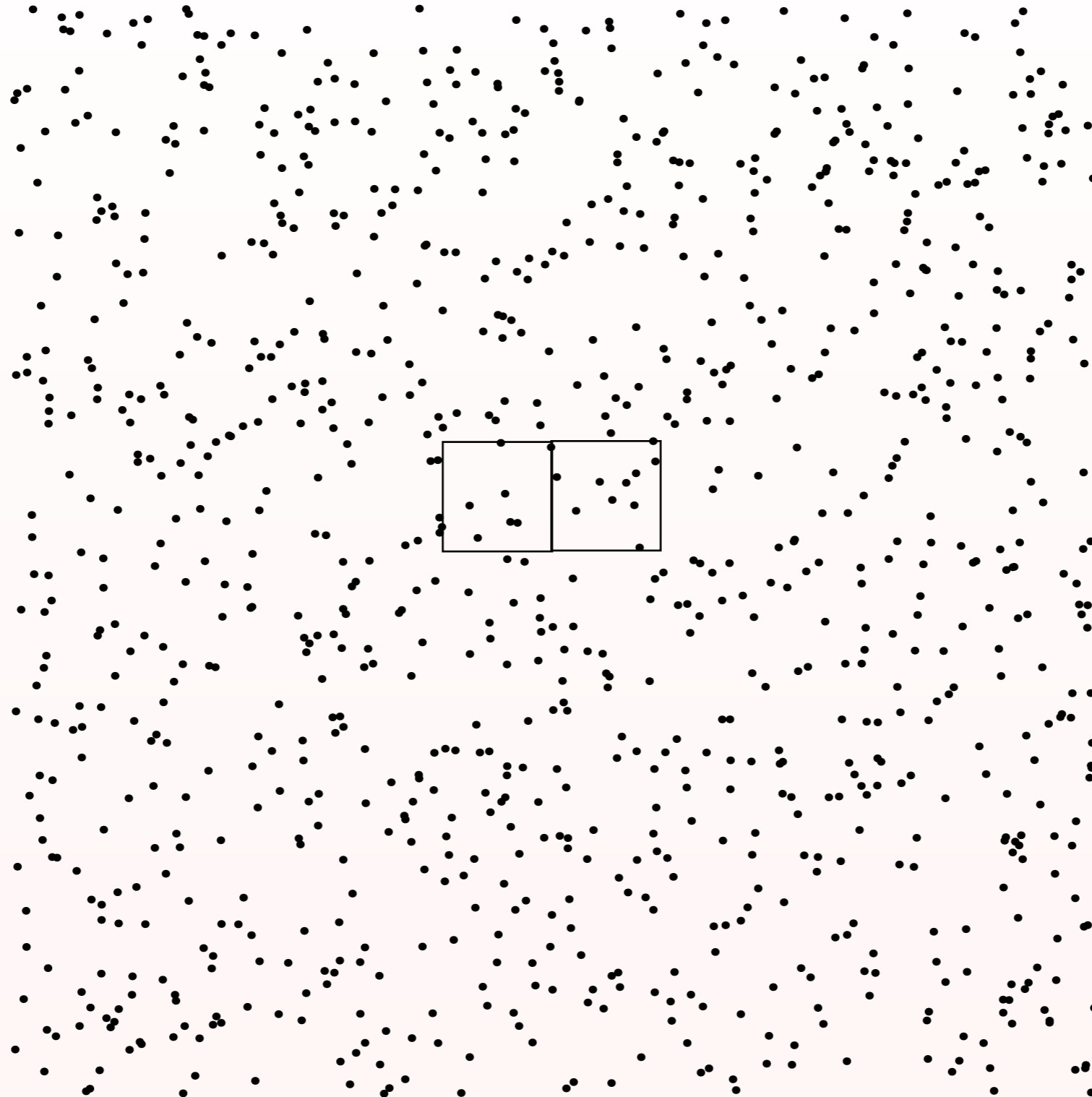
Grid search



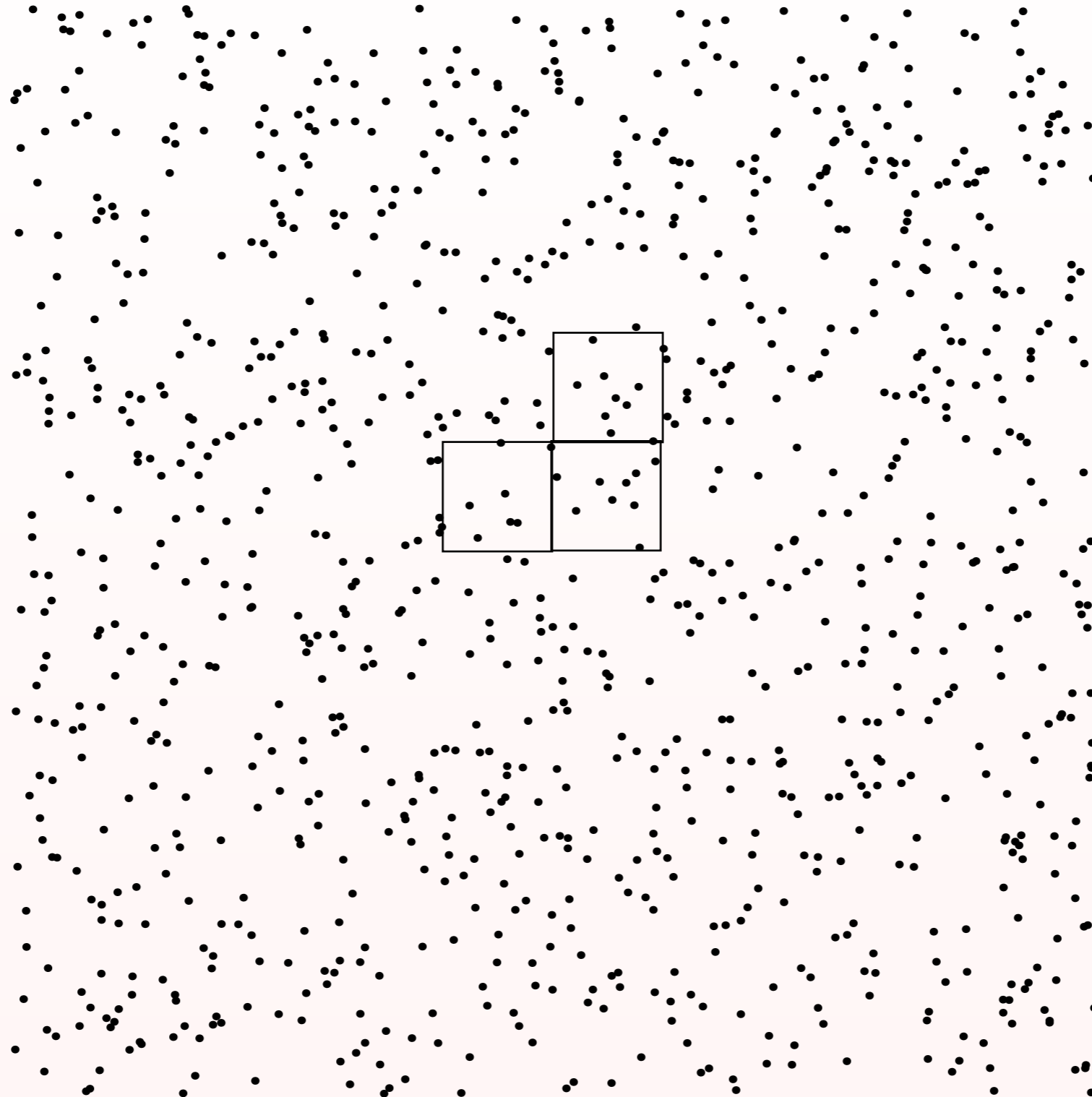
Grid search



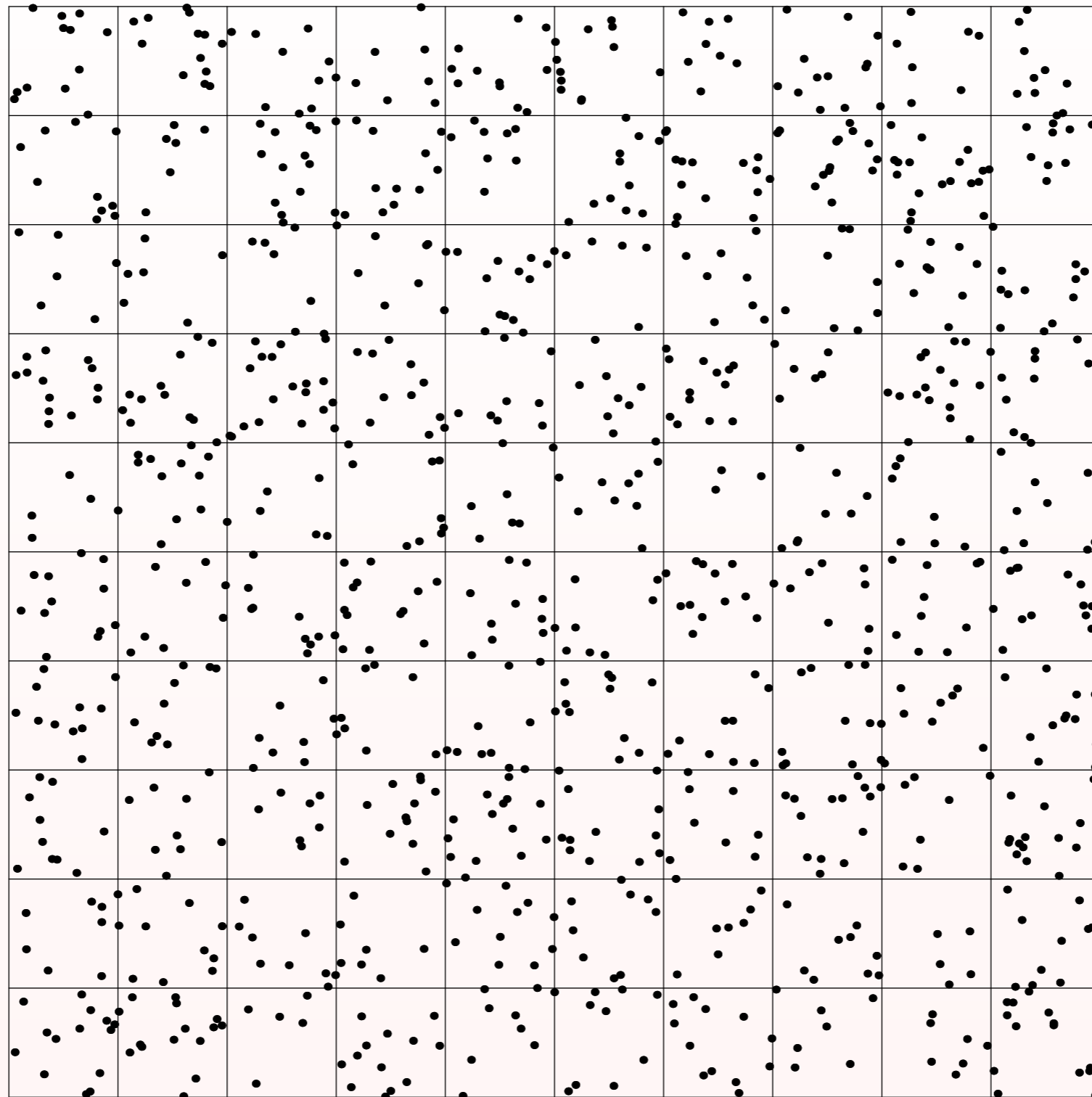
Grid search

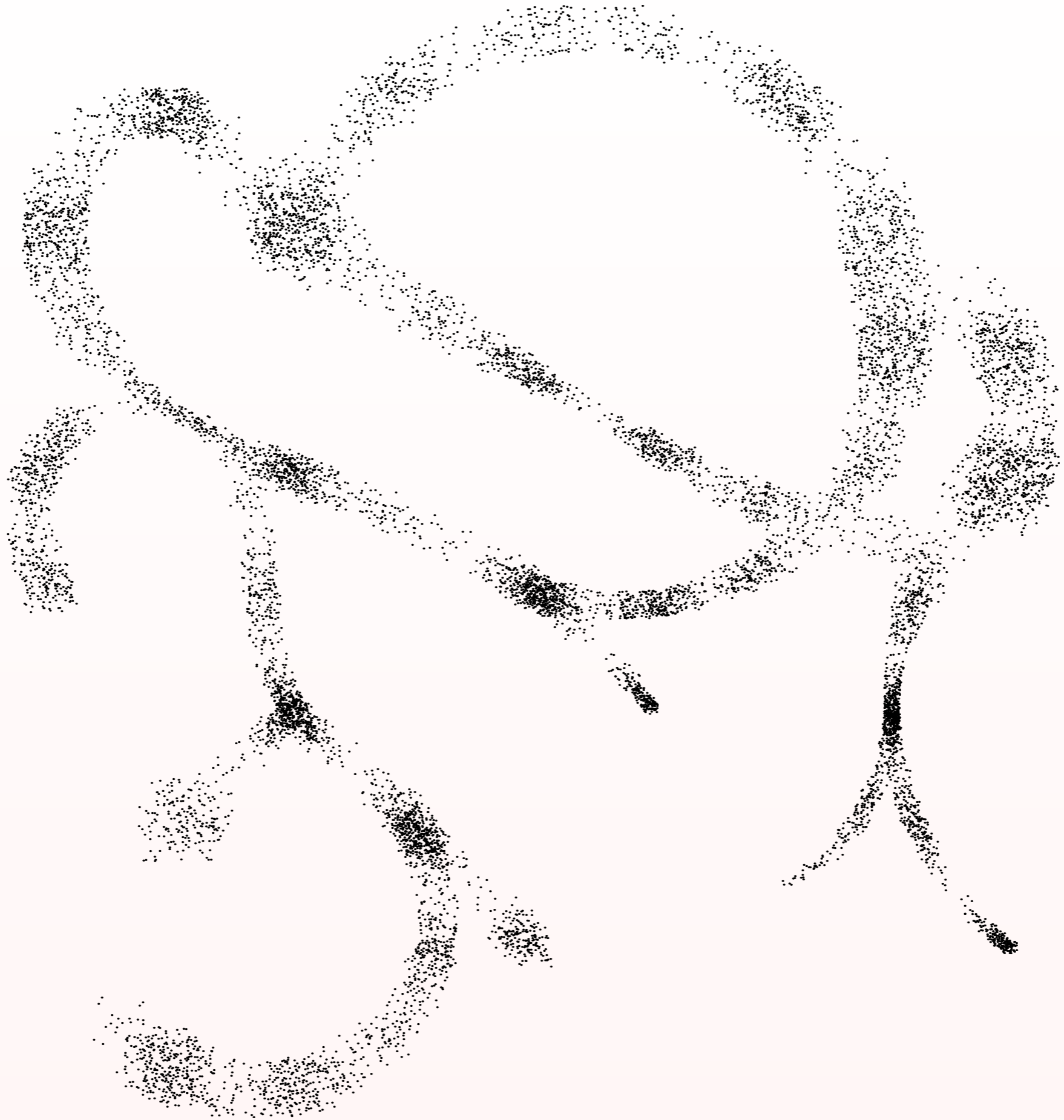


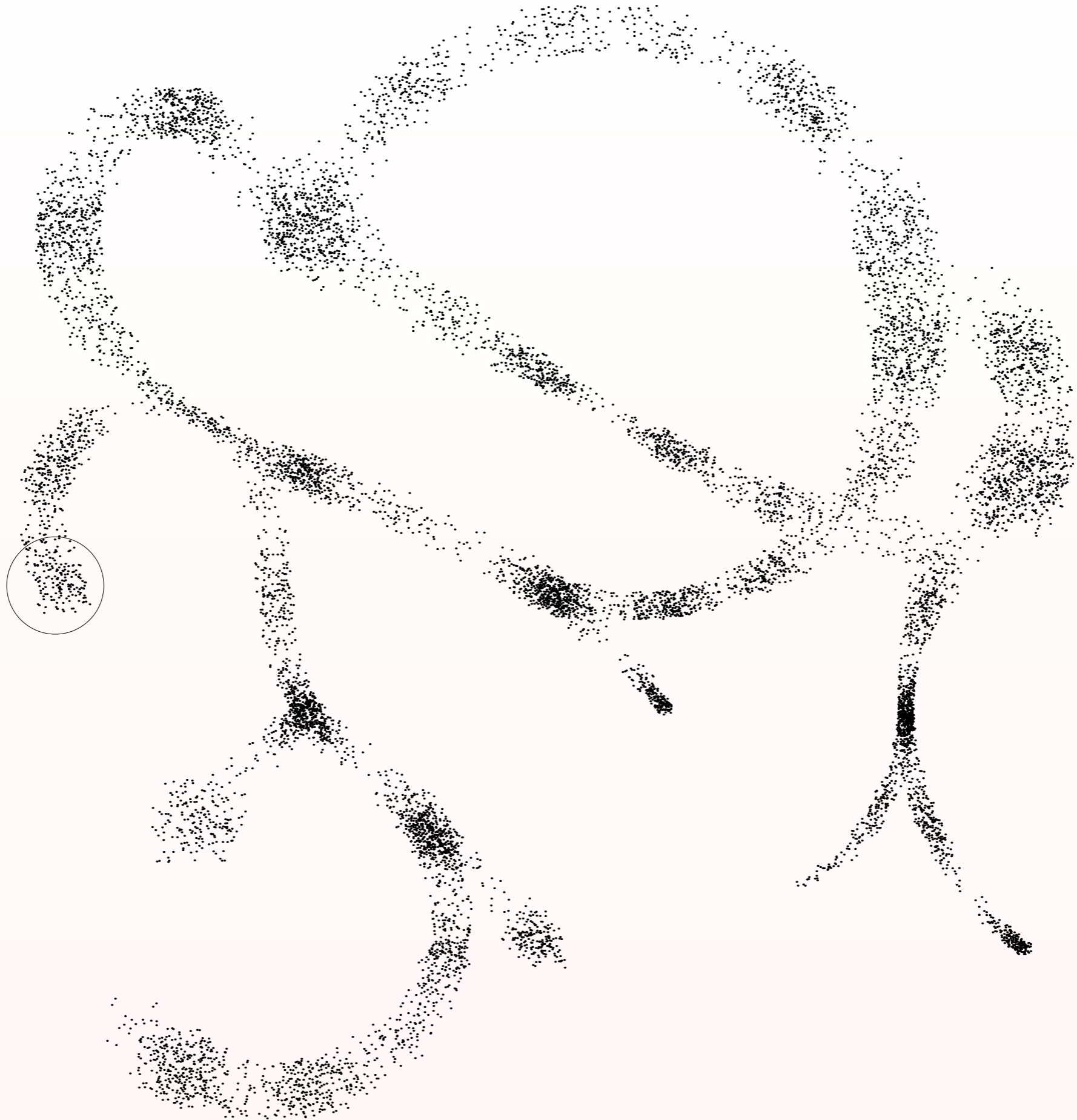
Grid search

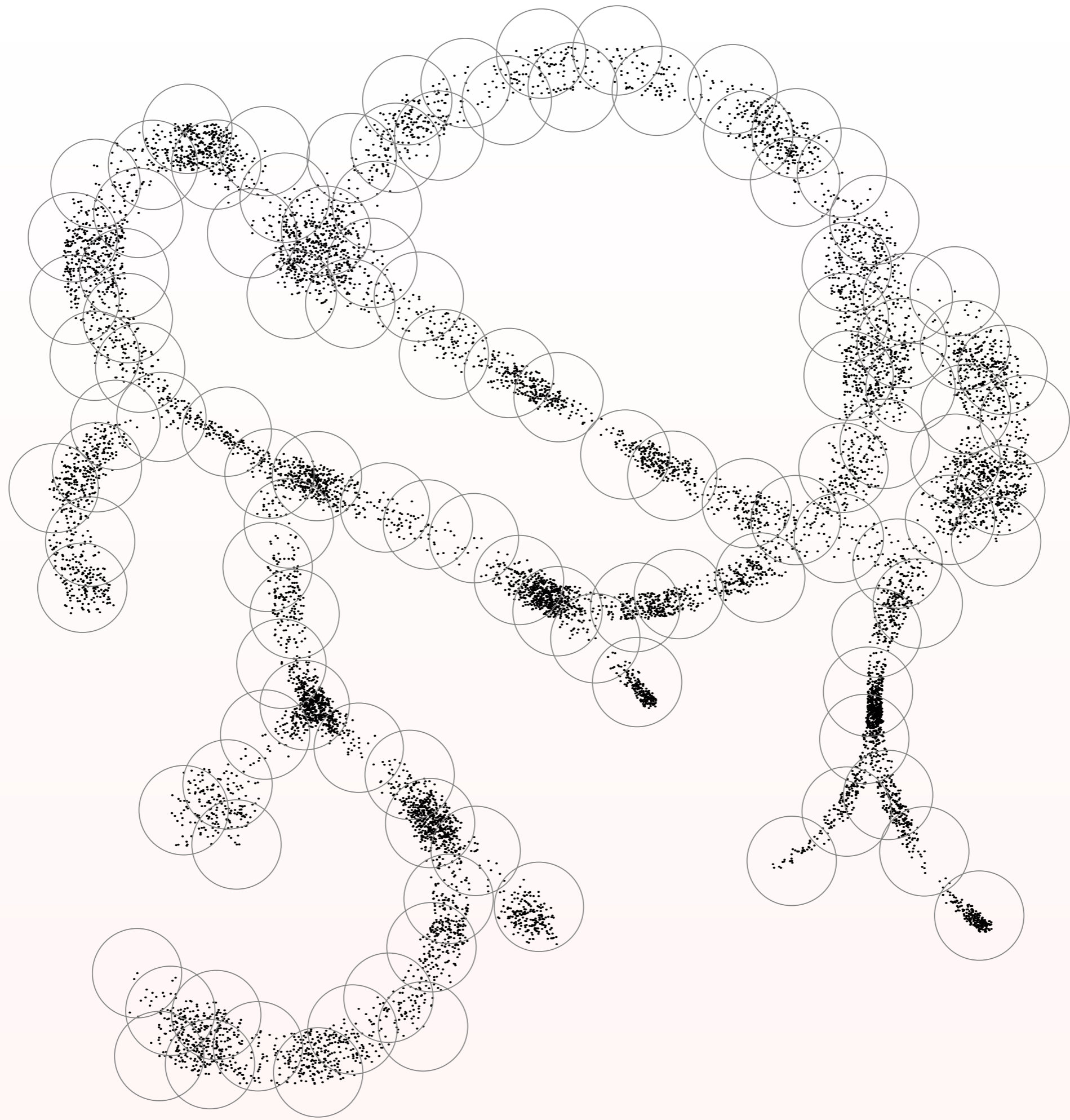


Grid search

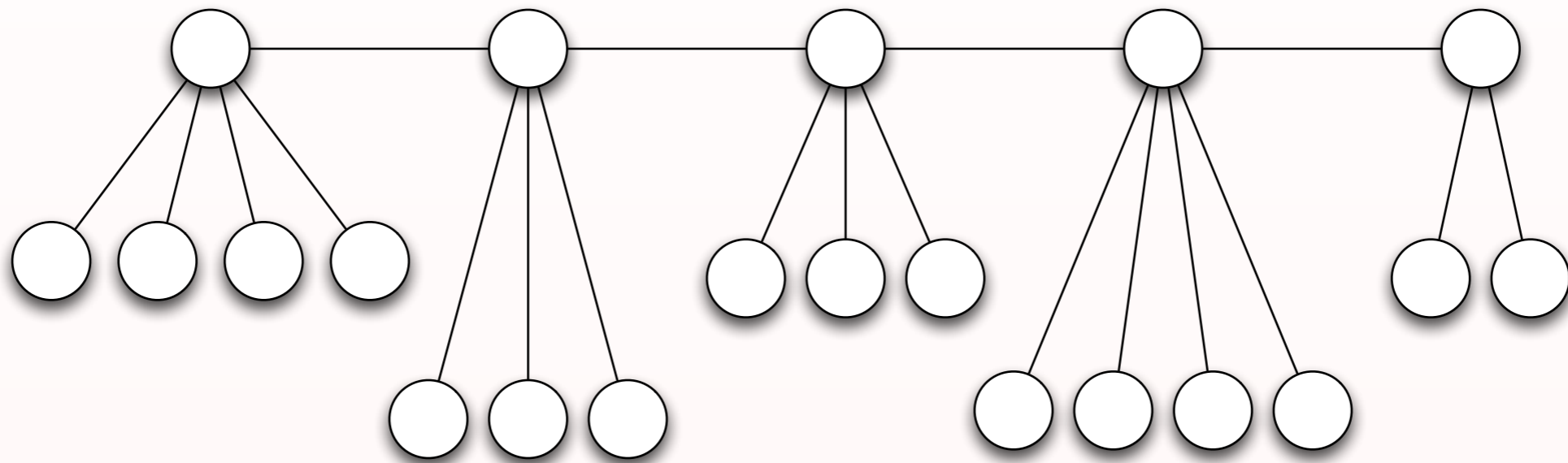




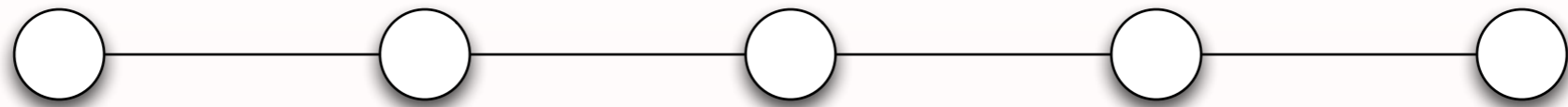




The search tree

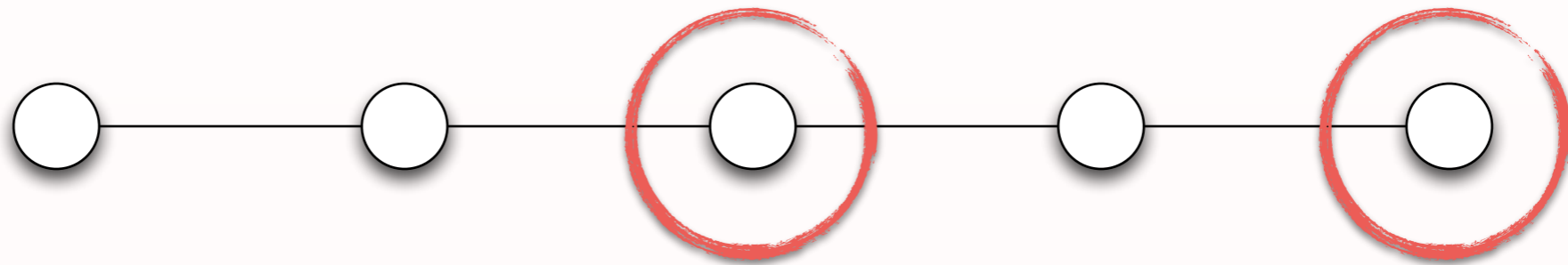


The search tree



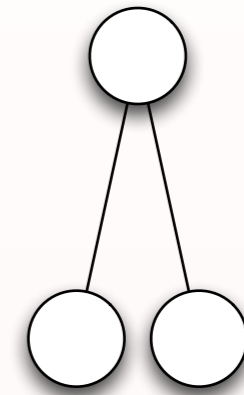
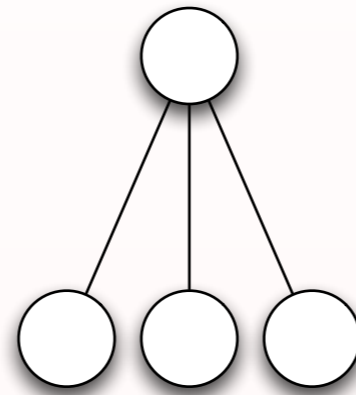
coarse search

The search tree



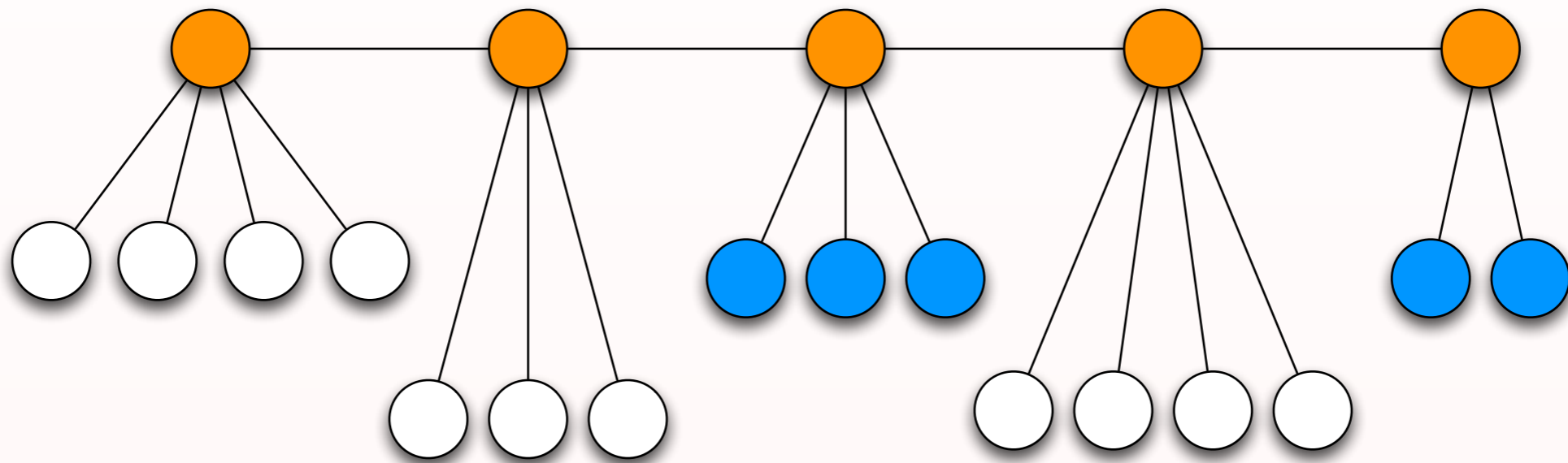
coarse search

The search tree



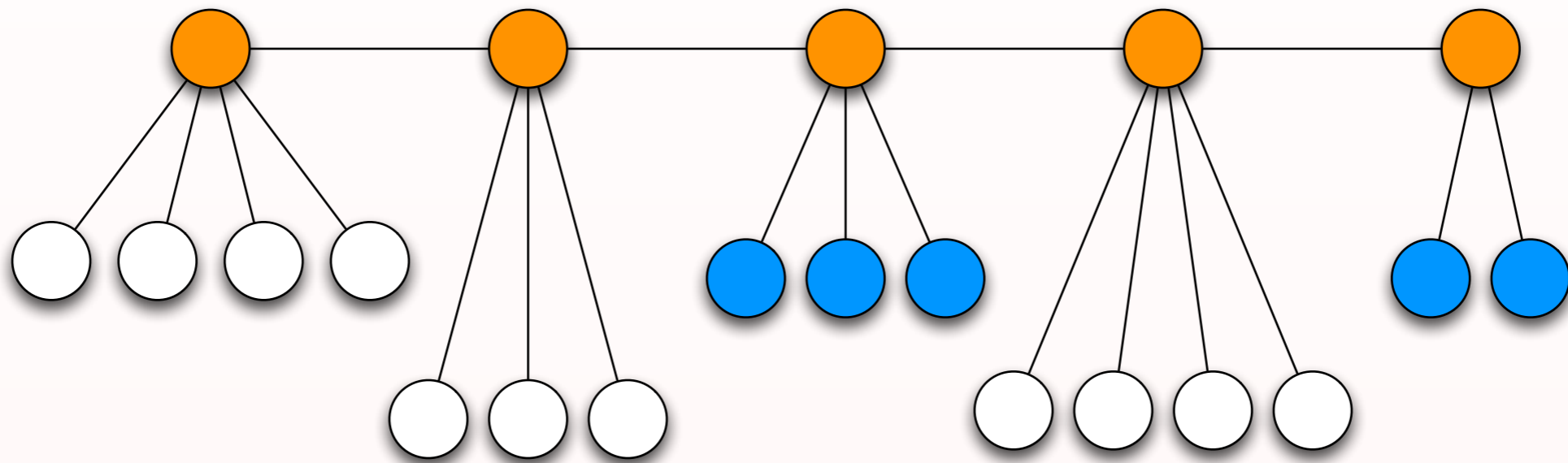
fine search

The search tree



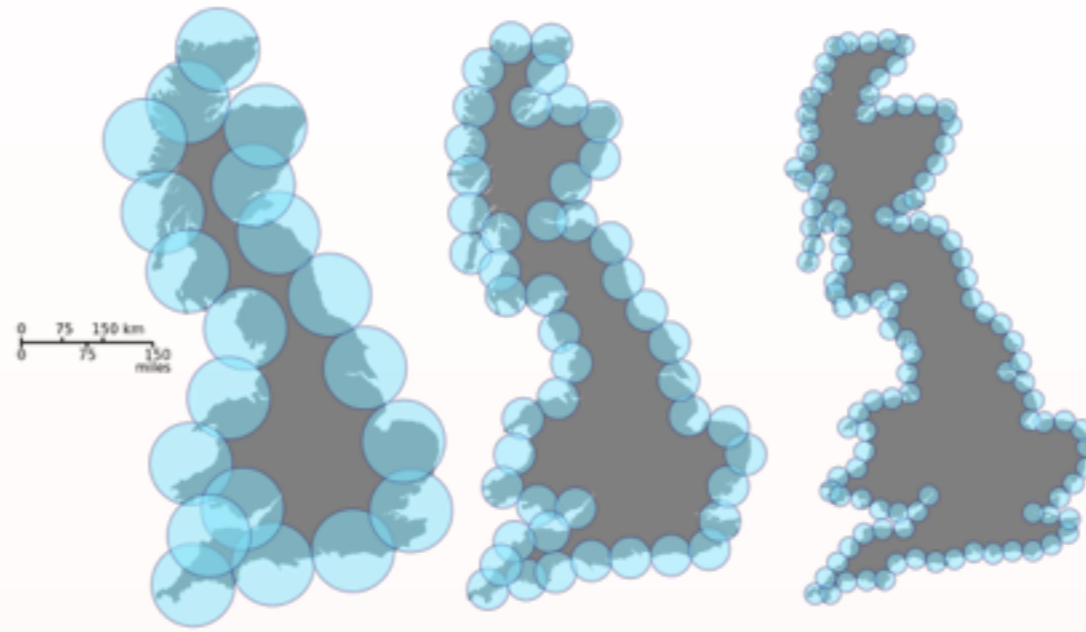
$$O(\textit{coarse} + \textit{fine})$$

The search tree

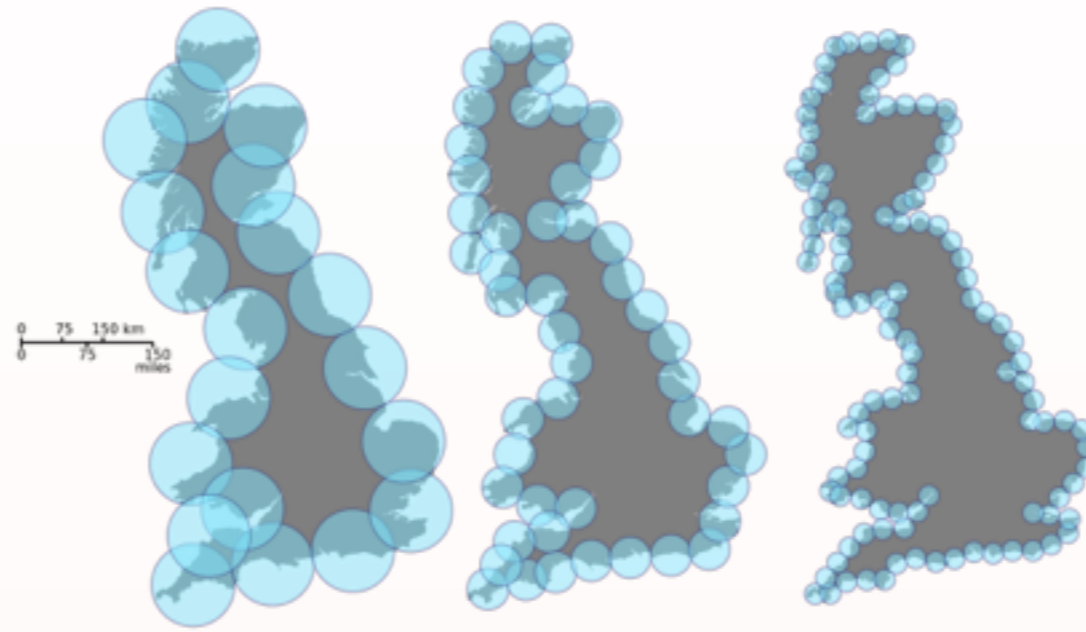


$$O(k + \textit{fine})$$

Metric entropy



Metric entropy



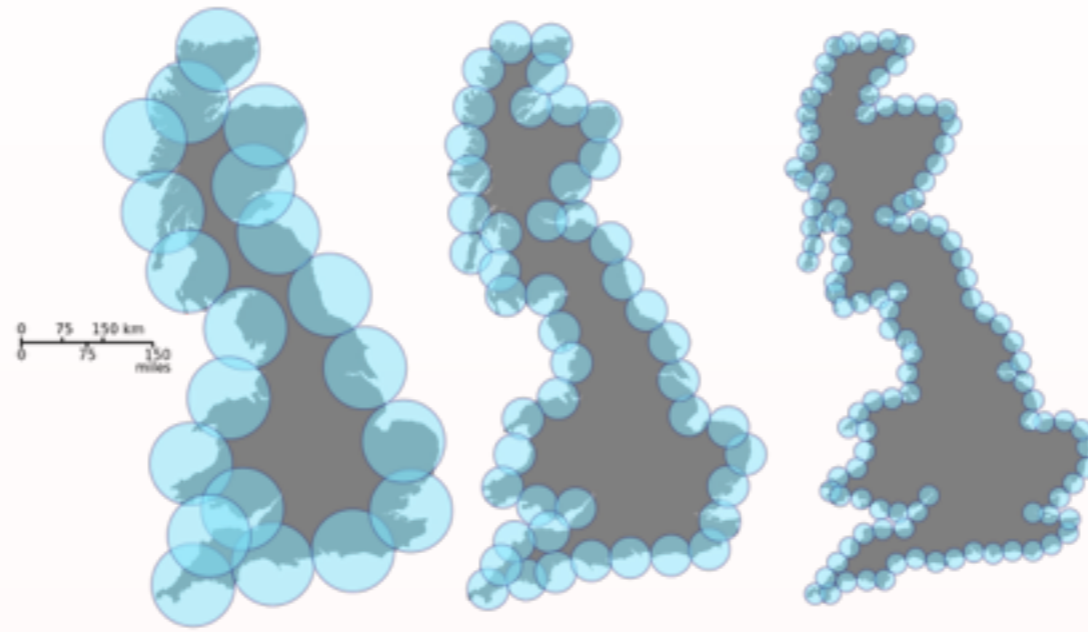
$$N_{r_c}^{ent}(D) :=$$

Metric entropy



$N_{r_c}^{ent}(D) :=$ largest number of points $x_1, \dots, x_n \in D$

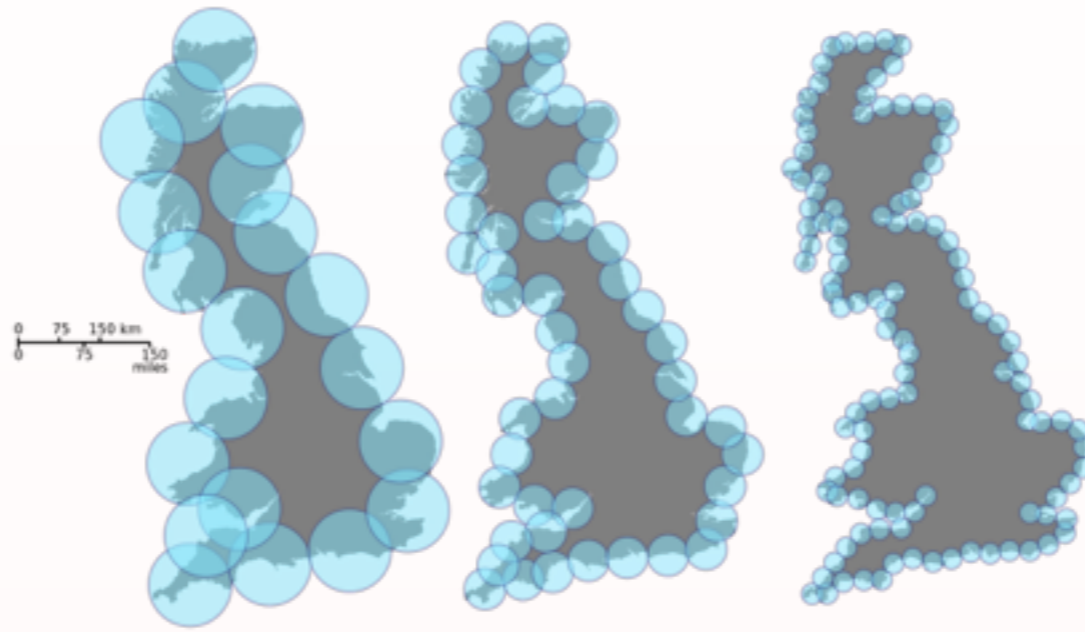
Metric entropy



$N_{r_c}^{ent}(D) :=$ largest number of points $x_1, \dots, x_n \in D$

s.t. $\|x_i - x_j\| \geq \rho, \forall i \neq j$

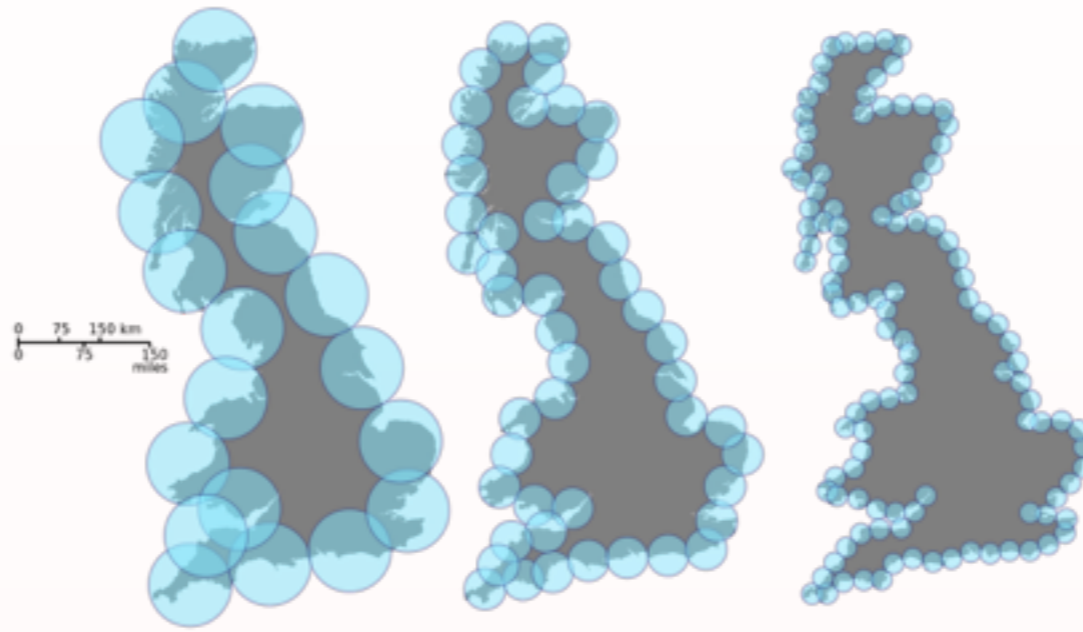
Metric entropy



$N_{r_c}^{ent}(D) :=$ largest number of points $x_1, \dots, x_n \in D$

s.t. $\|x_i - x_j\| \geq \rho, \forall i \neq j$ for some radius ρ

Metric entropy



$N_{r_c}^{ent}(D) :=$ largest number of points $x_1, \dots, x_n \in D$

s.t. $\|x_i - x_j\| \geq \rho, \forall i \neq j$ for some radius ρ

$$\dim_{\text{Minkowski}}(D) := \lim_{\rho \rightarrow 0} \frac{\log N_{\rho}(D)}{\log 1/\rho}$$

Discrete fractal dimension

Discrete fractal dimension

$$d := \arg \max_{d^*} \{N_\rho(D) \propto \rho^{d^*} \mid \rho \in [\rho_1, \rho_2]\}$$

d is the dimension that maximizes the metric entropy

Discrete fractal dimension

$$d := \arg \max_{d^*} \{N_\rho(D) \propto \rho^{d^*} \mid \rho \in [\rho_1, \rho_2]\}$$

d is the dimension that maximizes the metric entropy

$$\rho = r_c$$

Discrete fractal dimension

$$d := \arg \max_{d^*} \{N_\rho(D) \propto \rho^{d^*} \mid \rho \in [\rho_1, \rho_2]\}$$

d is the dimension that maximizes the metric entropy

$$\rho = r_c$$

k clusters

Discrete fractal dimension

$$d := \arg \max_{d^*} \{N_\rho(D) \propto \rho^{d^*} \mid \rho \in [\rho_1, \rho_2]\}$$

d is the dimension that maximizes the metric entropy

$$\rho = r_c$$

k clusters

$$k \leq N_{r_c}^{ent}(D)$$

Discrete fractal dimension

$$d := \arg \max_{d^*} \{N_\rho(D) \propto \rho^{d^*} \mid \rho \in [\rho_1, \rho_2]\}$$

d is the dimension that maximizes the metric entropy

$$\rho = r_c$$

k clusters

$$k \leq N_{r_c}^{ent}(D)$$

$$O(k + fine)$$

Discrete fractal dimension

$$d := \arg \max_{d^*} \{N_\rho(D) \propto \rho^{d^*} \mid \rho \in [\rho_1, \rho_2]\}$$

d is the dimension that maximizes the metric entropy

$$\rho = r_c$$

k clusters

$$k \leq N_{r_c}^{ent}(D)$$

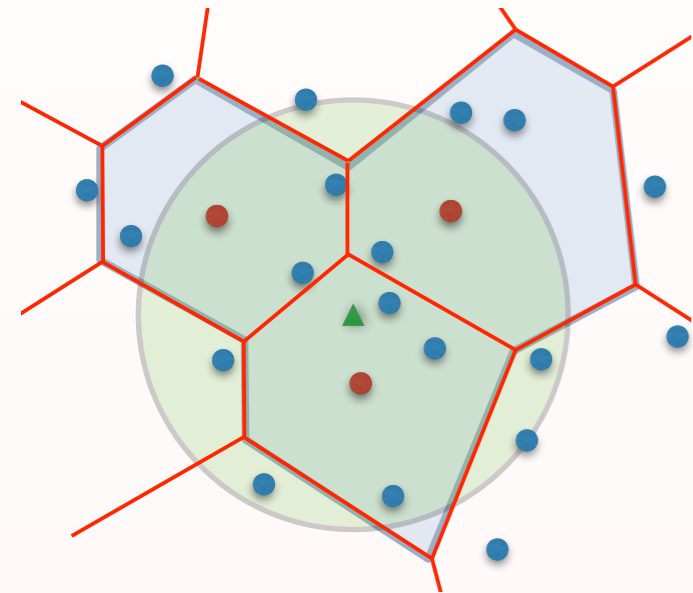
$$O(N_{r_c}^{ent}(D) + fine)$$

Fine search bounds

$$O(N_{r_c}^{ent}(D) + |F|)$$

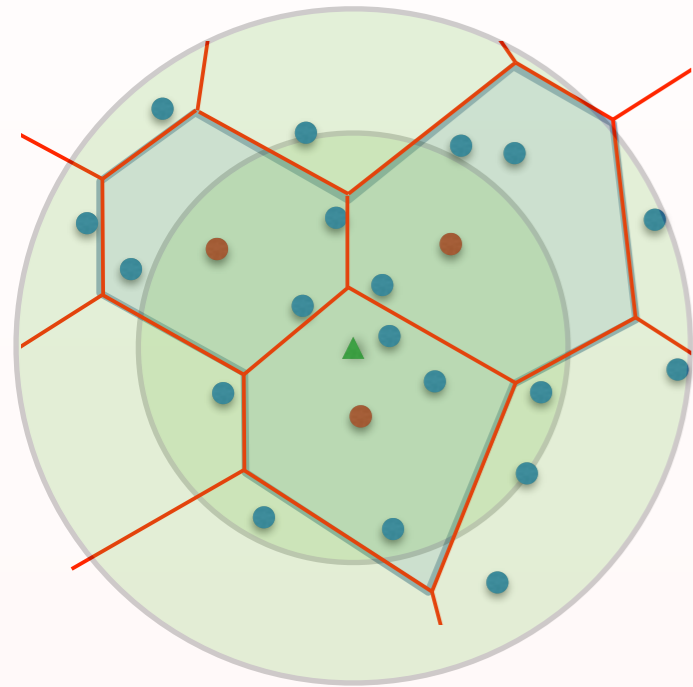
Fine search bounds

$$O(N_{r_c}^{ent}(D) + |F|) \quad F = \bigcup_{c \in B_C(q, r+r_c)} c$$



Fine search bounds

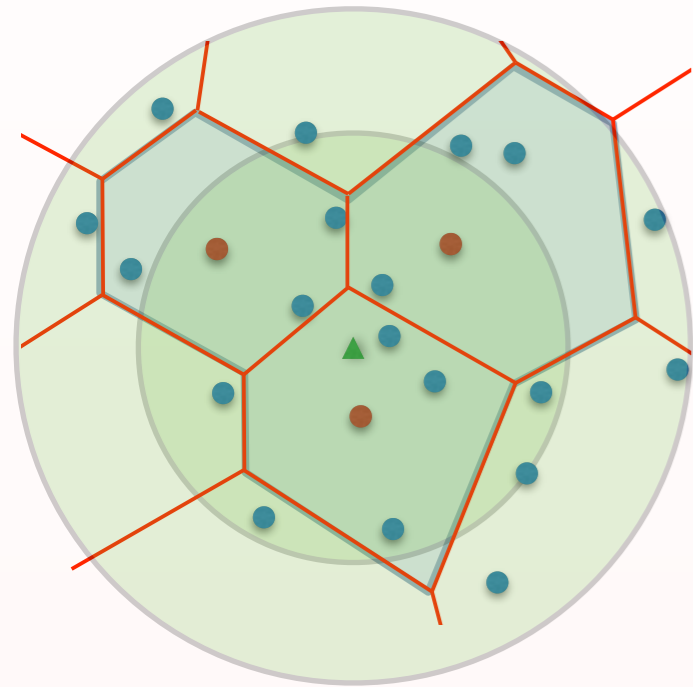
$$O(N_{r_c}^{ent}(D) + |F|) \quad F = \bigcup_{c \in B_C(q, r + r_c)} c$$



$$F \subset B_D(q, r + 2r_c)$$

Fine search bounds

$$O(N_{r_c}^{ent}(D) + |F|) \quad F = \bigcup_{c \in B_C(q, r + r_c)} c$$

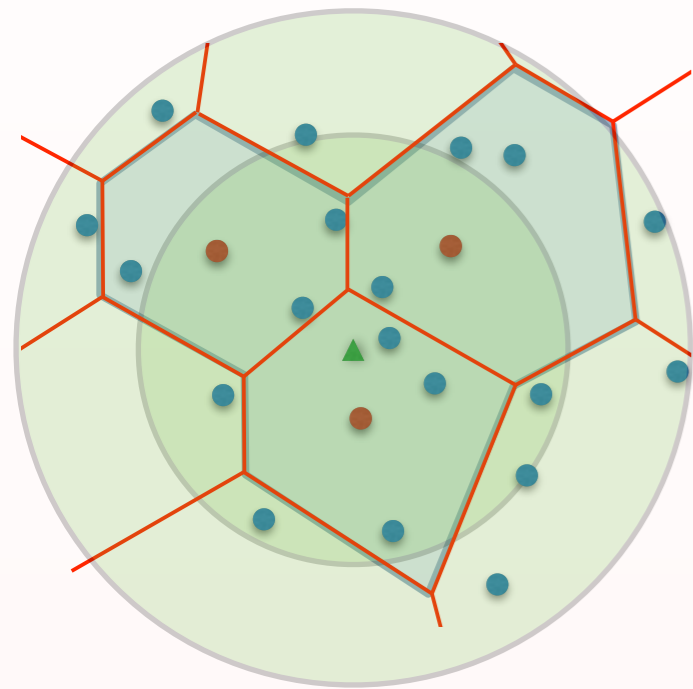


$$F \subset B_D(q, r + 2r_c)$$

$$|F| \leq |B_D(q, r + 2r_c)|$$

Fine search bounds

$$O(N_{r_c}^{ent}(D) + |F|) \quad F = \bigcup_{c \in B_C(q, r+r_c)} c$$

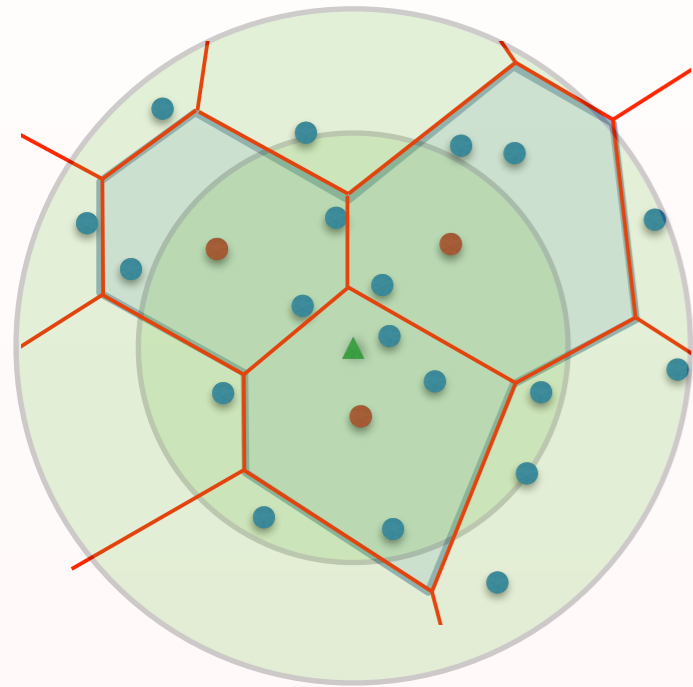


$$F \subset B_D(q, r + 2r_c)$$

$$|F| \leq |B_D(q, r + 2r_c)| \sim |B_D(q, r)| \left(\frac{r + 2r_c}{r} \right)^d$$

Fine search bounds

$$O(N_{r_c}^{ent}(D) + |F|) \quad F = \bigcup_{c \in B_C(q, r+r_c)} c$$



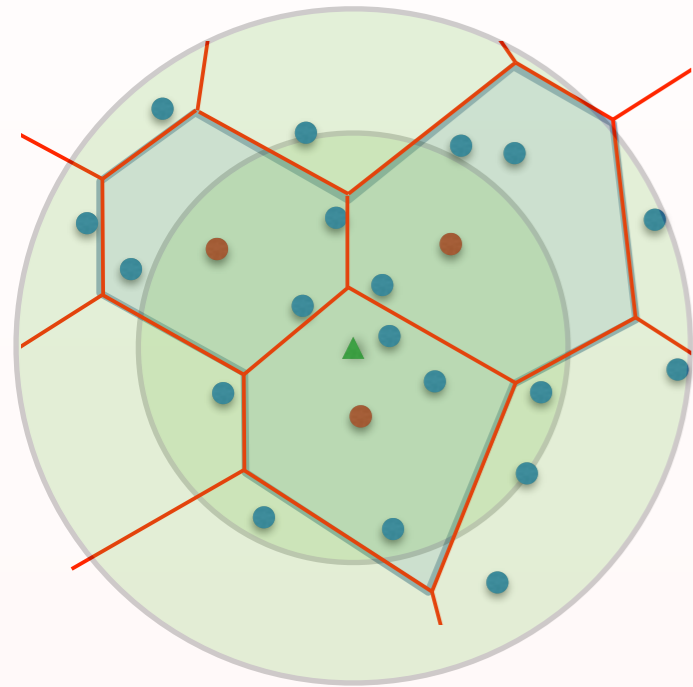
$$F \subset B_D(q, r + 2r_c)$$

$$|F| \leq |B_D(q, r + 2r_c)| \sim |B_D(q, r)| \left(\frac{r + 2r_c}{r} \right)^d$$

Doubling search radius increases hits by at most 2^d

Fine search bounds

$$O(N_{r_c}^{ent}(D) + |F|) \quad F = \bigcup_{c \in B_C(q, r+r_c)} c$$



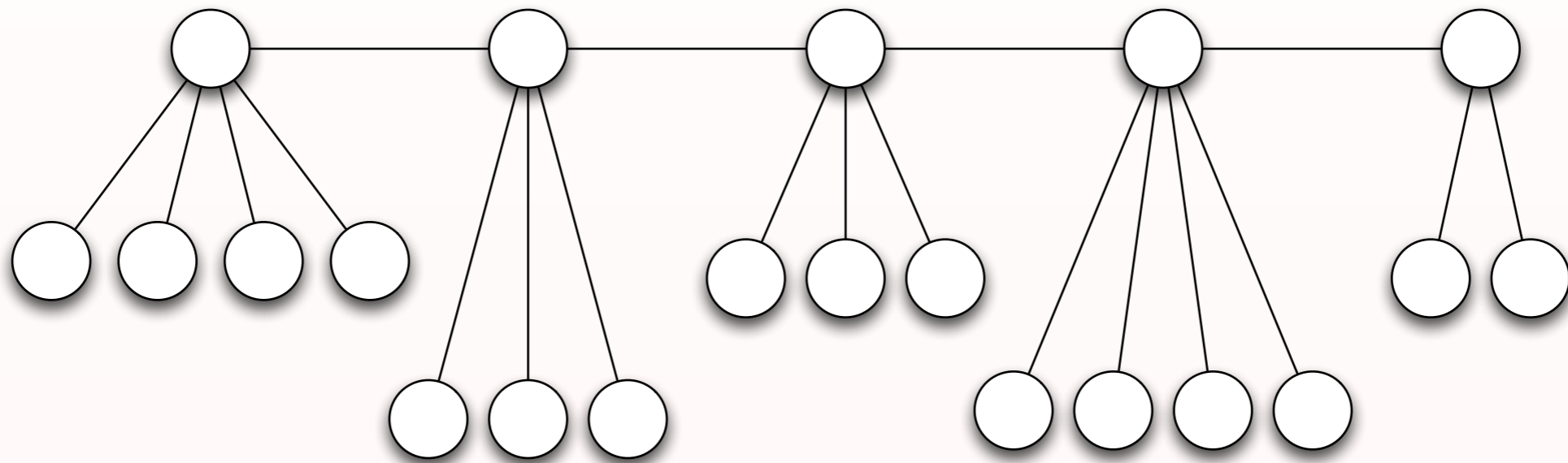
$$F \subset B_D(q, r + 2r_c)$$

$$|F| \leq |B_D(q, r + 2r_c)| \sim |B_D(q, r)| \left(\frac{r + 2r_c}{r} \right)^d$$

Doubling search radius increases hits by at most 2^d

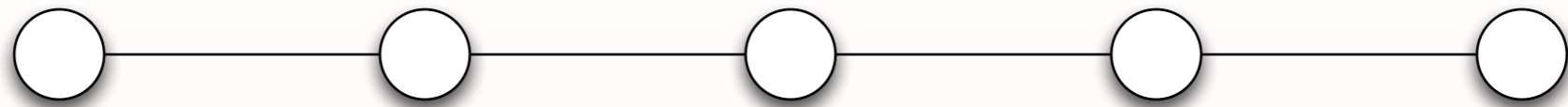
$$O \left(N_{r_c}^{ent}(D) + |B_D(q, r)| \left(\frac{r + 2r_c}{r} \right)^d \right)$$

Entropy-scaling data structure



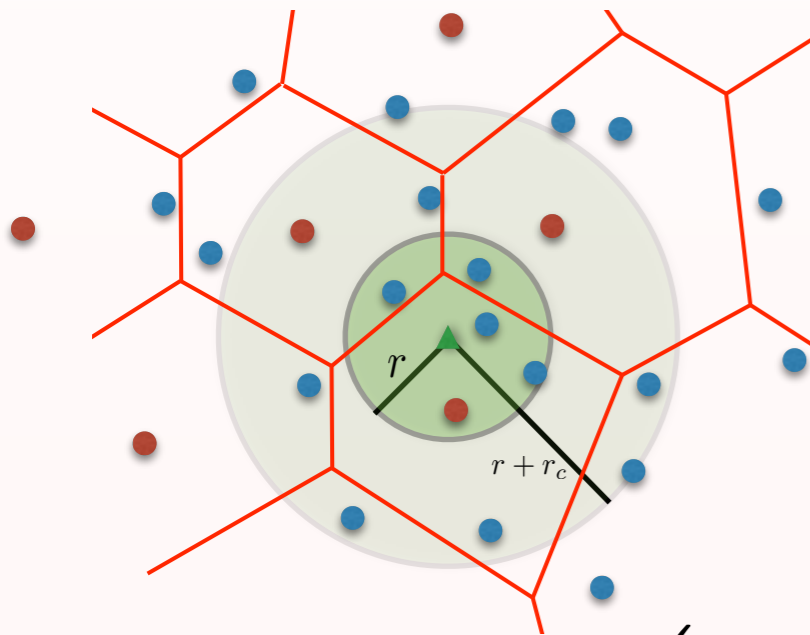
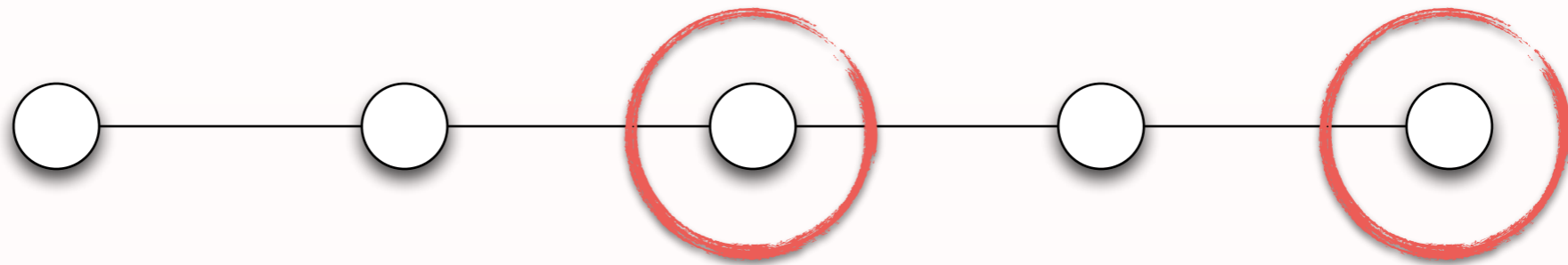
$$O \left(N_{r_c}^{ent}(D) + |B_D(q, r)| \left(\frac{r + 2r_c}{r} \right)^d \right)$$

Entropy-scaling data structure



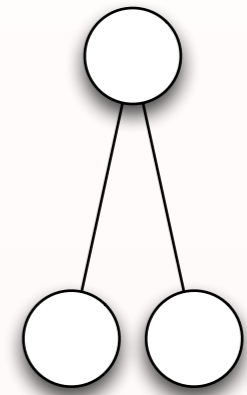
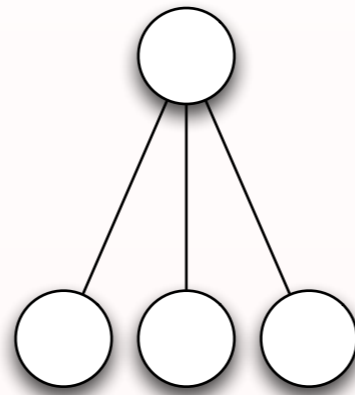
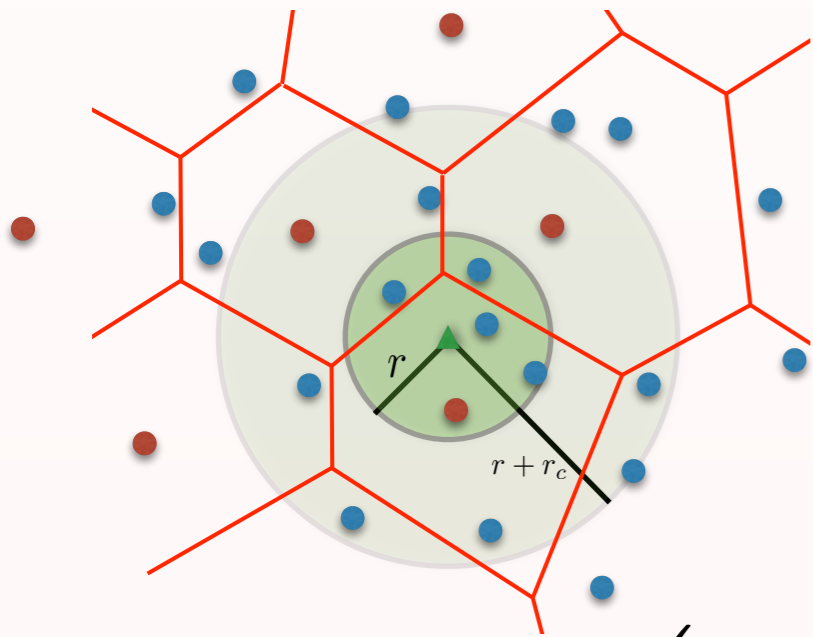
$$O \left(N_{r_c}^{ent}(D) + |B_D(q, r)| \left(\frac{r + 2r_c}{r} \right)^d \right)$$

Entropy-scaling data structure



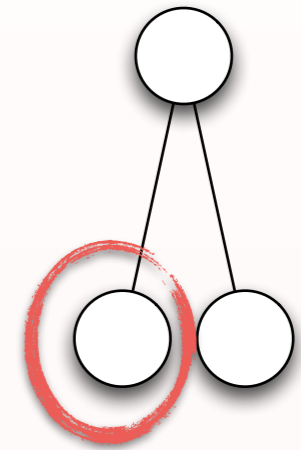
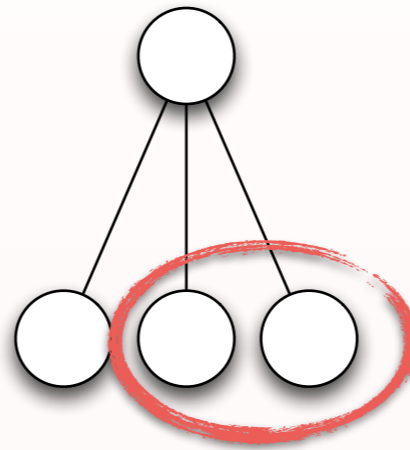
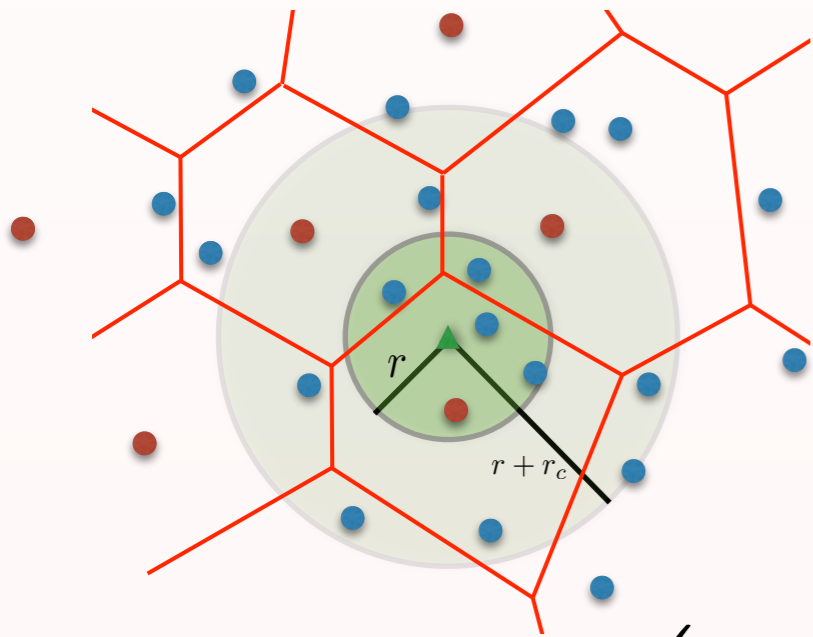
$$O \left(N_{r_c}^{ent}(D) + |B_D(q, r)| \left(\frac{r + 2r_c}{r} \right)^d \right)$$

Entropy-scaling data structure

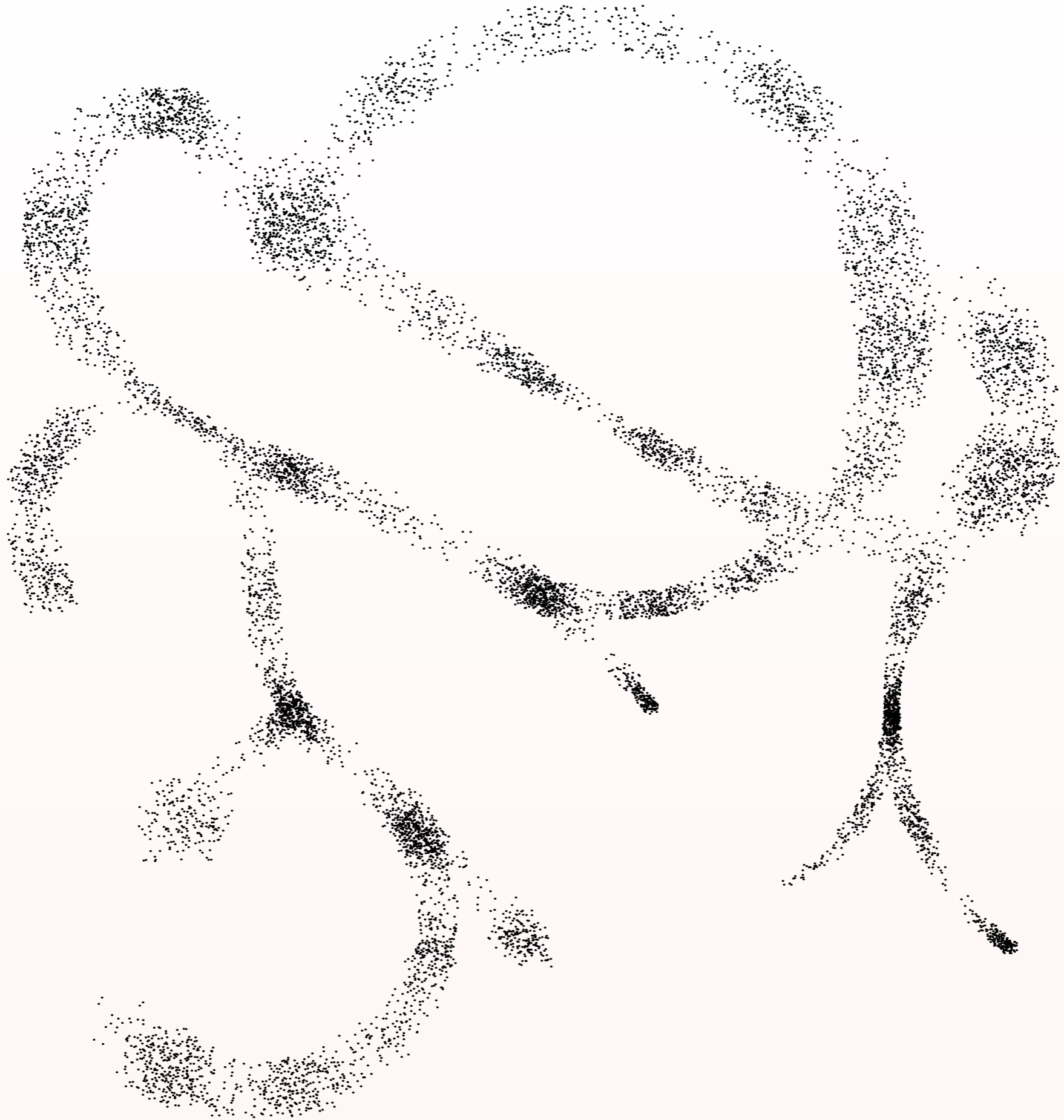


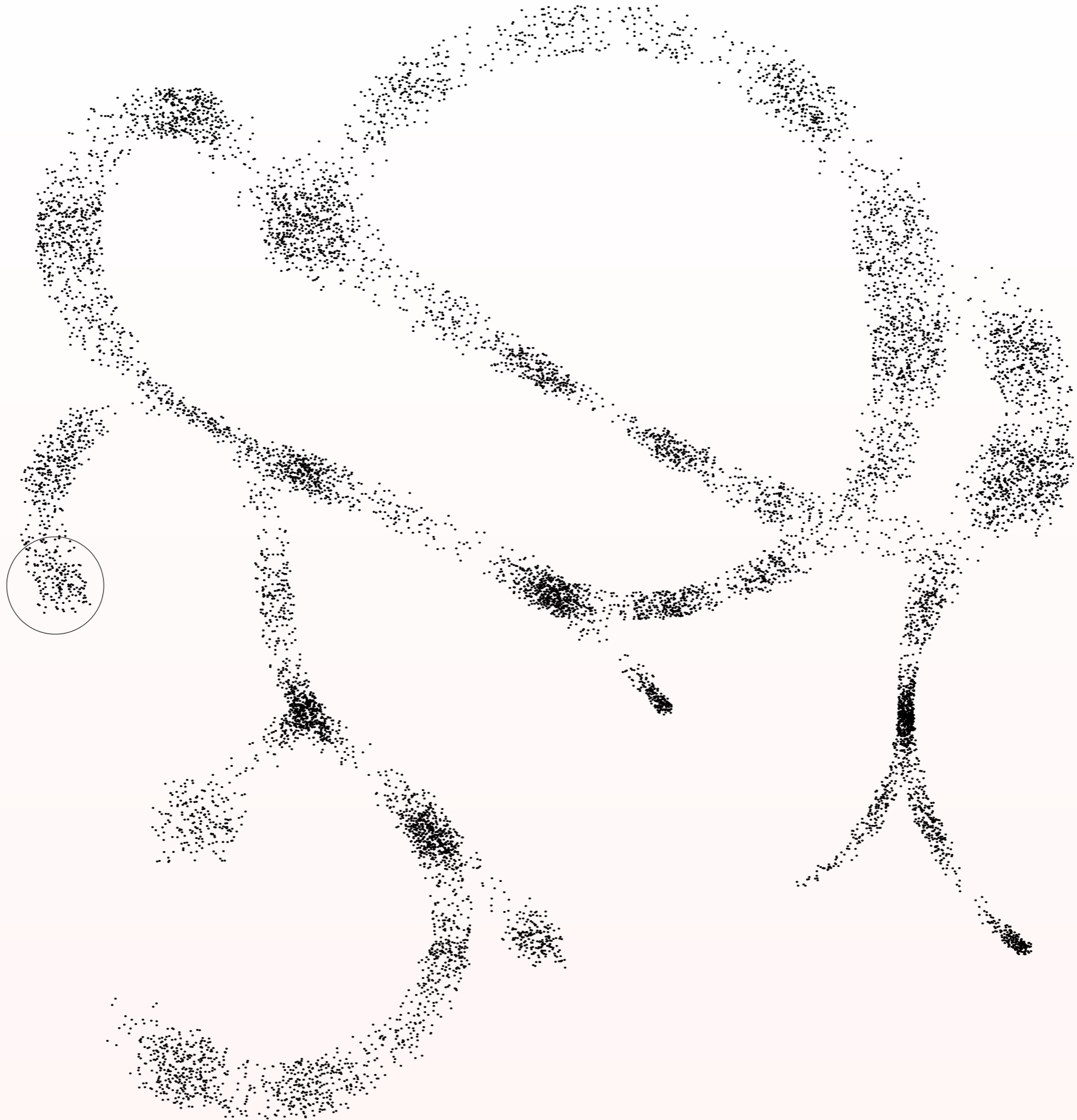
$$O \left(N_{r_c}^{ent}(D) + |B_D(q, r)| \left(\frac{r + 2r_c}{r} \right)^d \right)$$

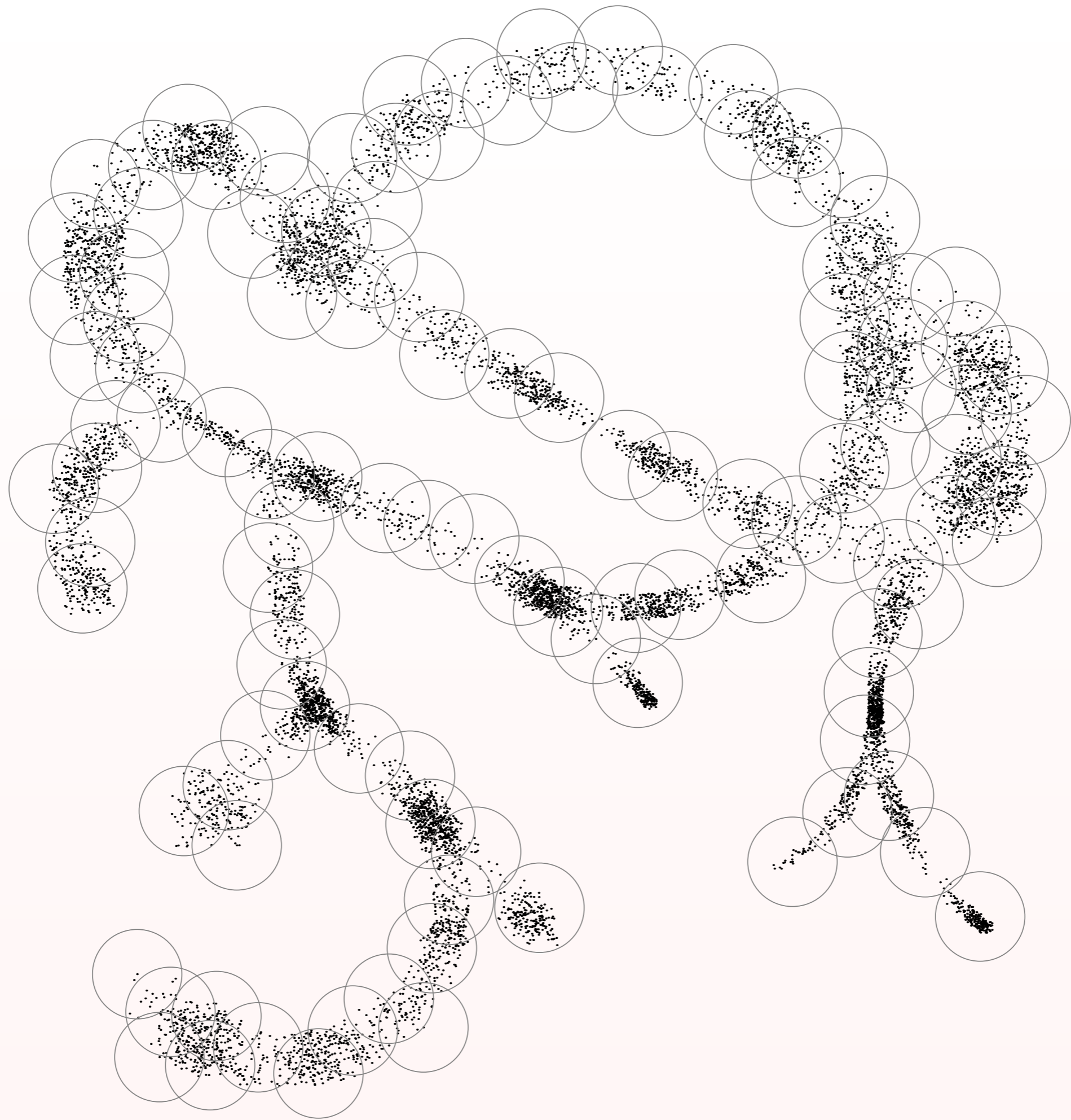
Entropy-scaling data structure

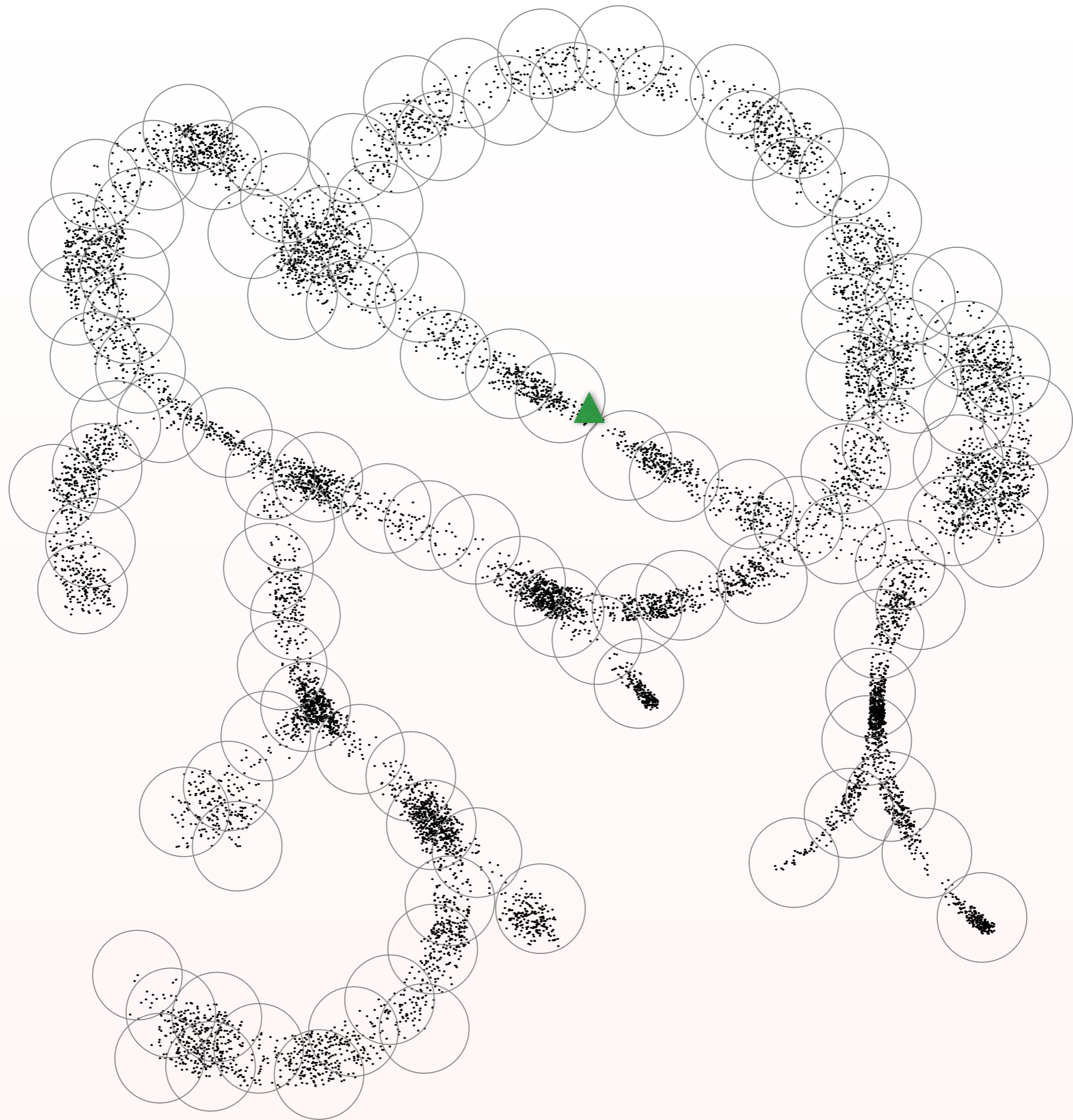


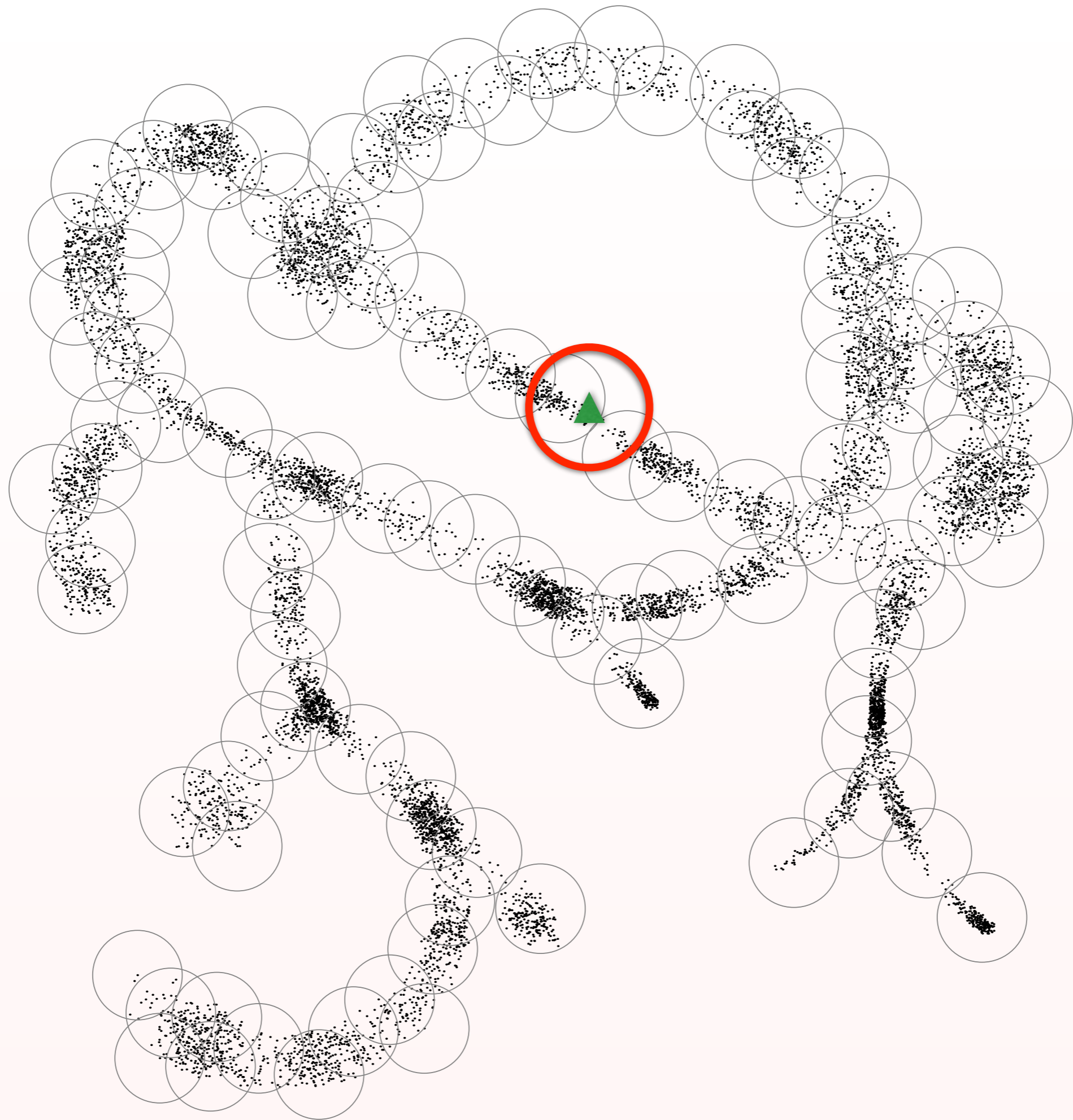
$$O \left(N_{r_c}^{ent}(D) + |B_D(q, r)| \left(\frac{r + 2r_c}{r} \right)^d \right)$$

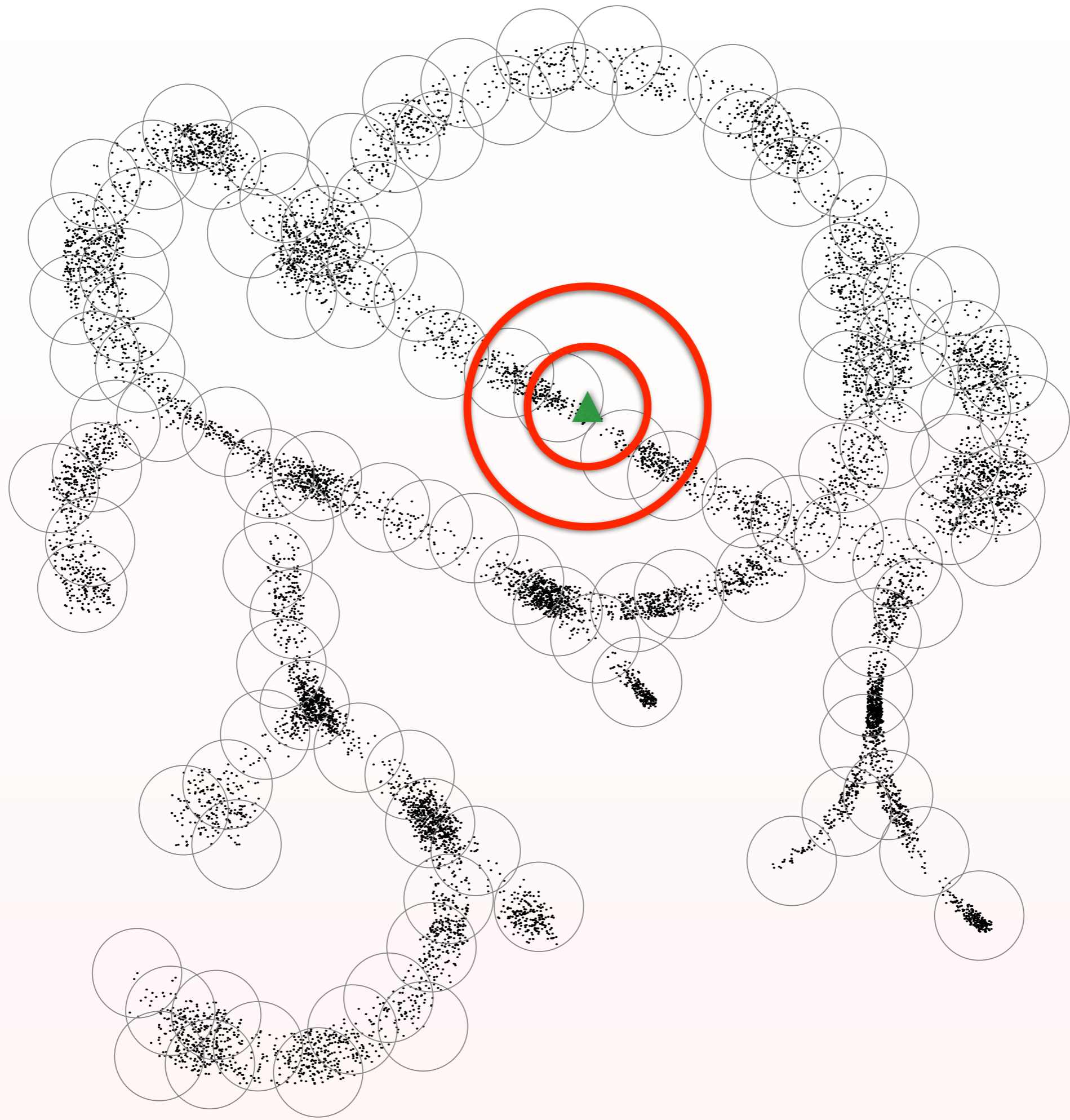


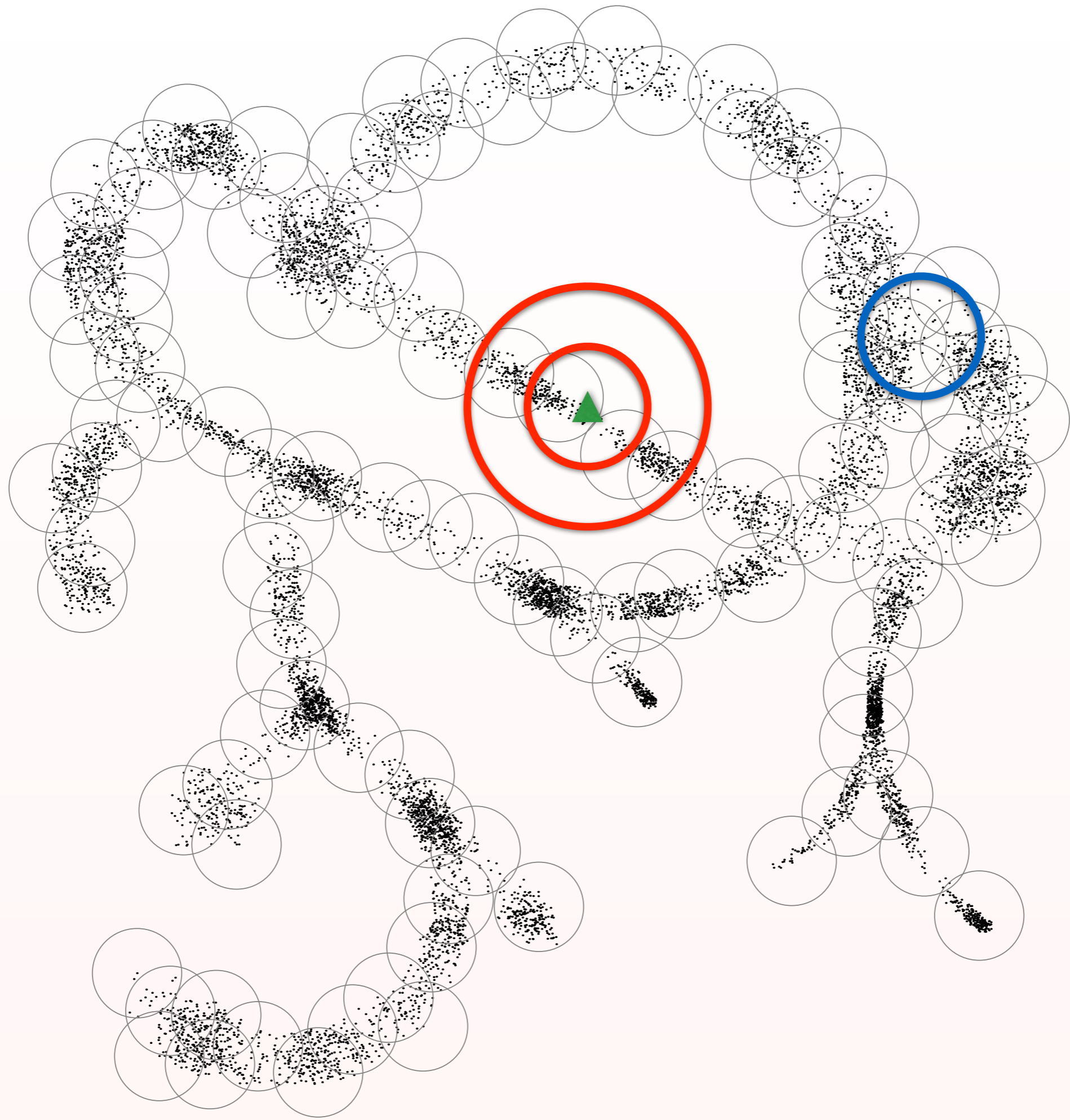


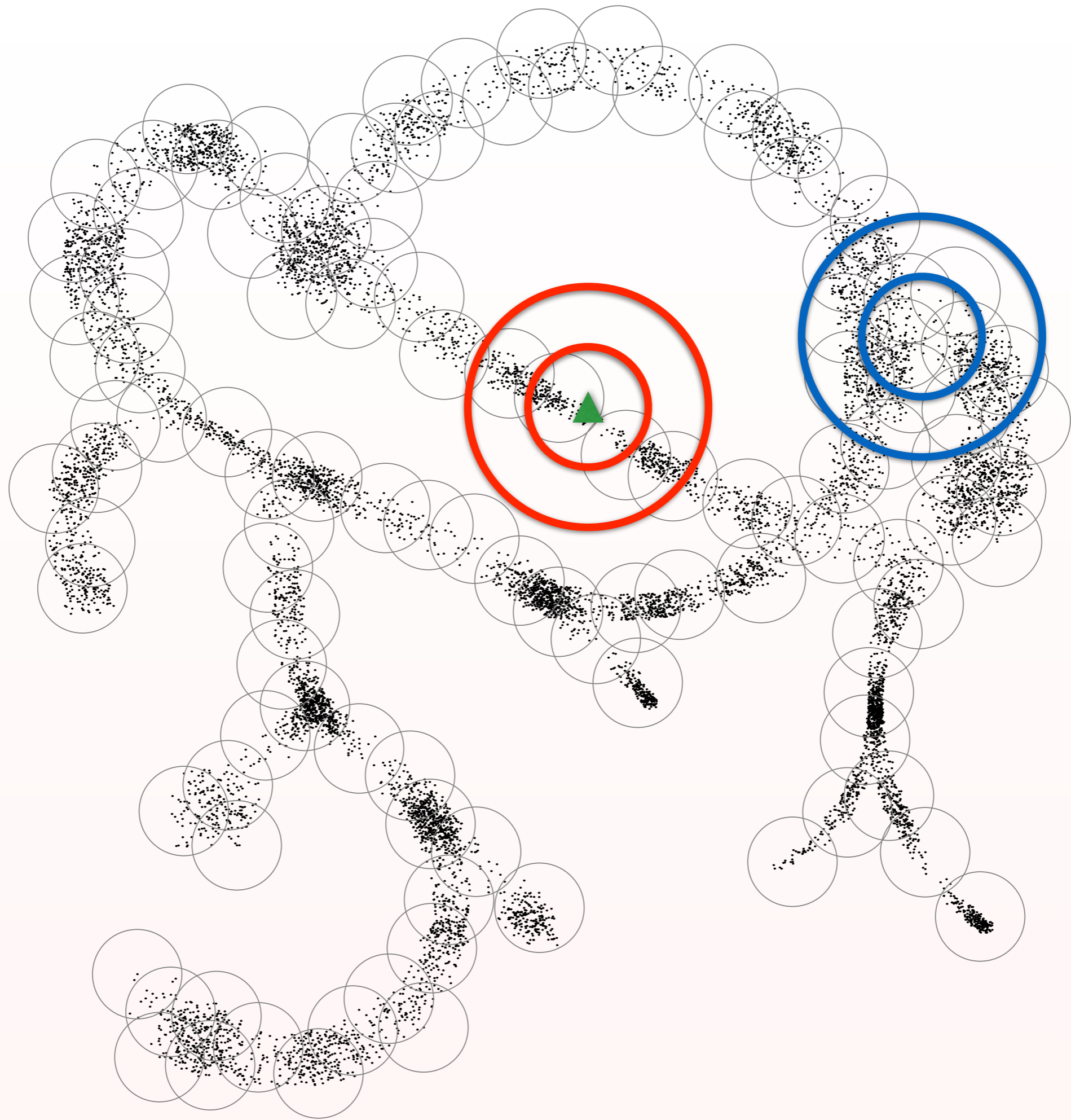












Summary

Summary

Lossless compression bounded by Shannon

Summary

Lossless compression bounded by Shannon

Lossy compression can go further

Summary

Lossless compression bounded by Shannon

Lossy compression can go further

Clever preprocessing of reads can boost standard compressors

Summary

Lossless compression bounded by Shannon

Lossy compression can go further

Clever preprocessing of reads can boost standard compressors

All comes down to *structure of the data*:

Summary

Lossless compression bounded by Shannon

Lossy compression can go further

Clever preprocessing of reads can boost standard compressors

All comes down to *structure of the data*:

minimize bits/base

Summary

Lossless compression bounded by Shannon

Lossy compression can go further

Clever preprocessing of reads can boost standard compressors

All comes down to *structure of the data*:

minimize bits/base

identify common substrings

Summary

Lossless compression bounded by Shannon

Lossy compression can go further

Clever preprocessing of reads can boost standard compressors

All comes down to *structure of the data*:

minimize bits/base

identify common substrings

find correlations among quality scores

Going further

Going further

Many more challenges:

Going further

Many more challenges:

New sequencing technologies

Going further

Many more challenges:

New sequencing technologies

longer, more error-prone reads

Going further

Many more challenges:

New sequencing technologies

longer, more error-prone reads

different error models

Going further

Many more challenges:

New sequencing technologies

longer, more error-prone reads

different error models

Can we get the benefits of compressive acceleration with a more general model