

Topological approach for finding nearest neighbor sequence in time series

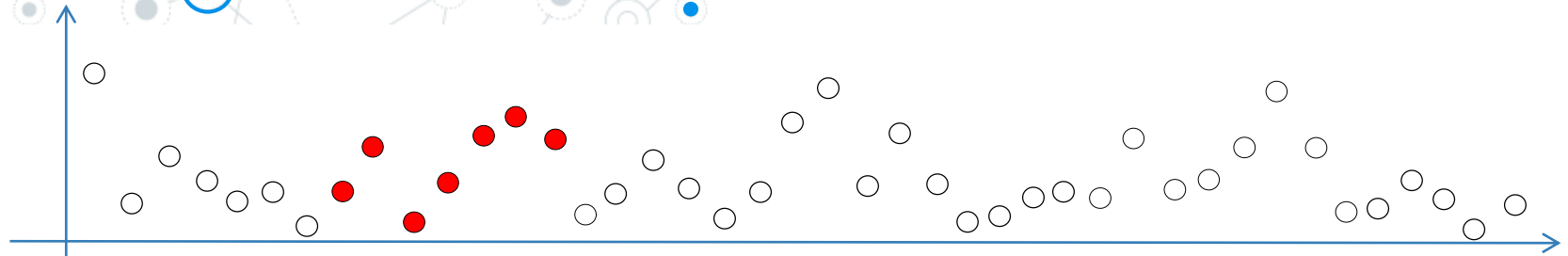
Paolo Avogadro, Matteo Dominoni

Dipartimento di Informatica, Sistemistica e Comunicazione

Università degli Studi di Milano-Bicocca

paolo.avogadro@unimib.it

Motivation



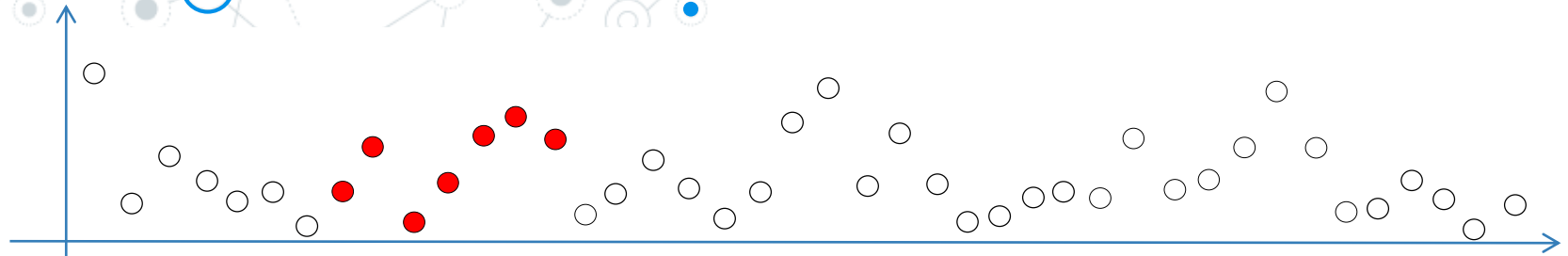
Given a **sequence**, knowing which is its (Euclidean) **nearest neighbor** provides important information, regarding:

anomalies

Recurrent patterns

clusters

Motivation



Given a **sequence**, knowing which is its (Euclidean) **nearest neighbor** provides important information, regarding:

anomalies

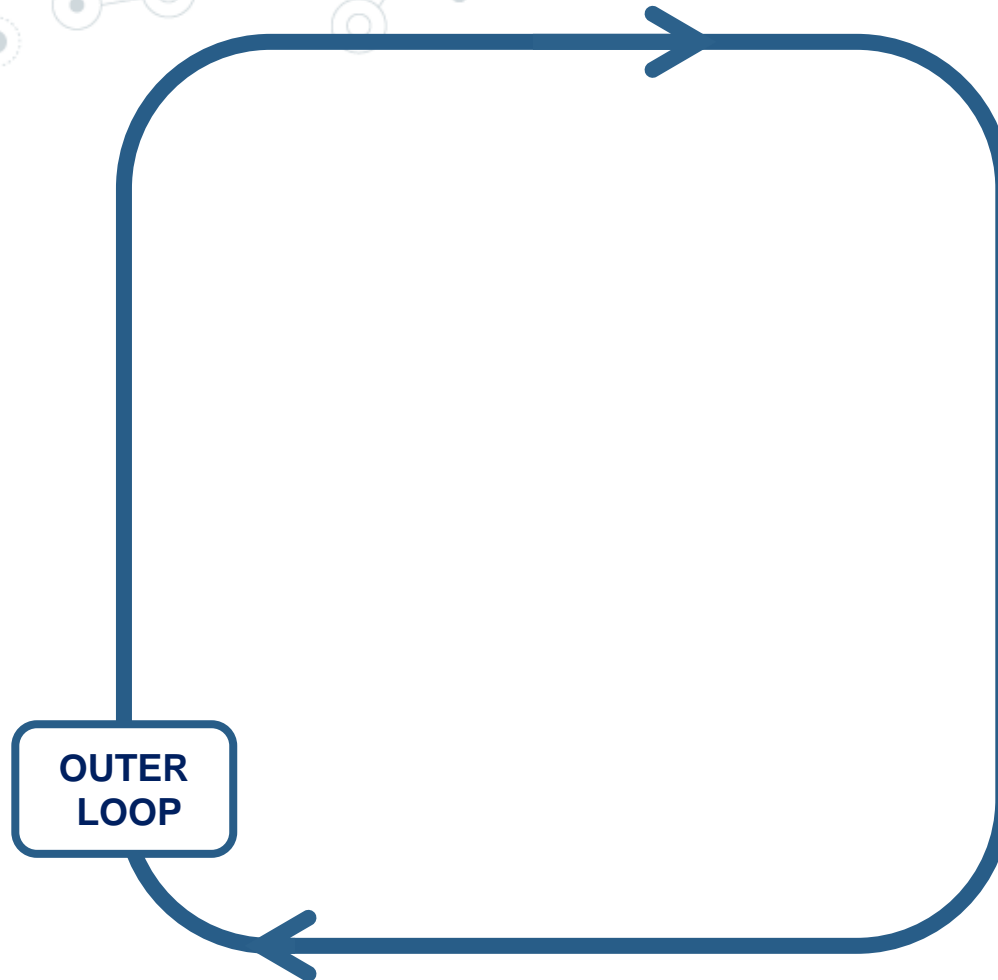
Recurrent patterns

clusters

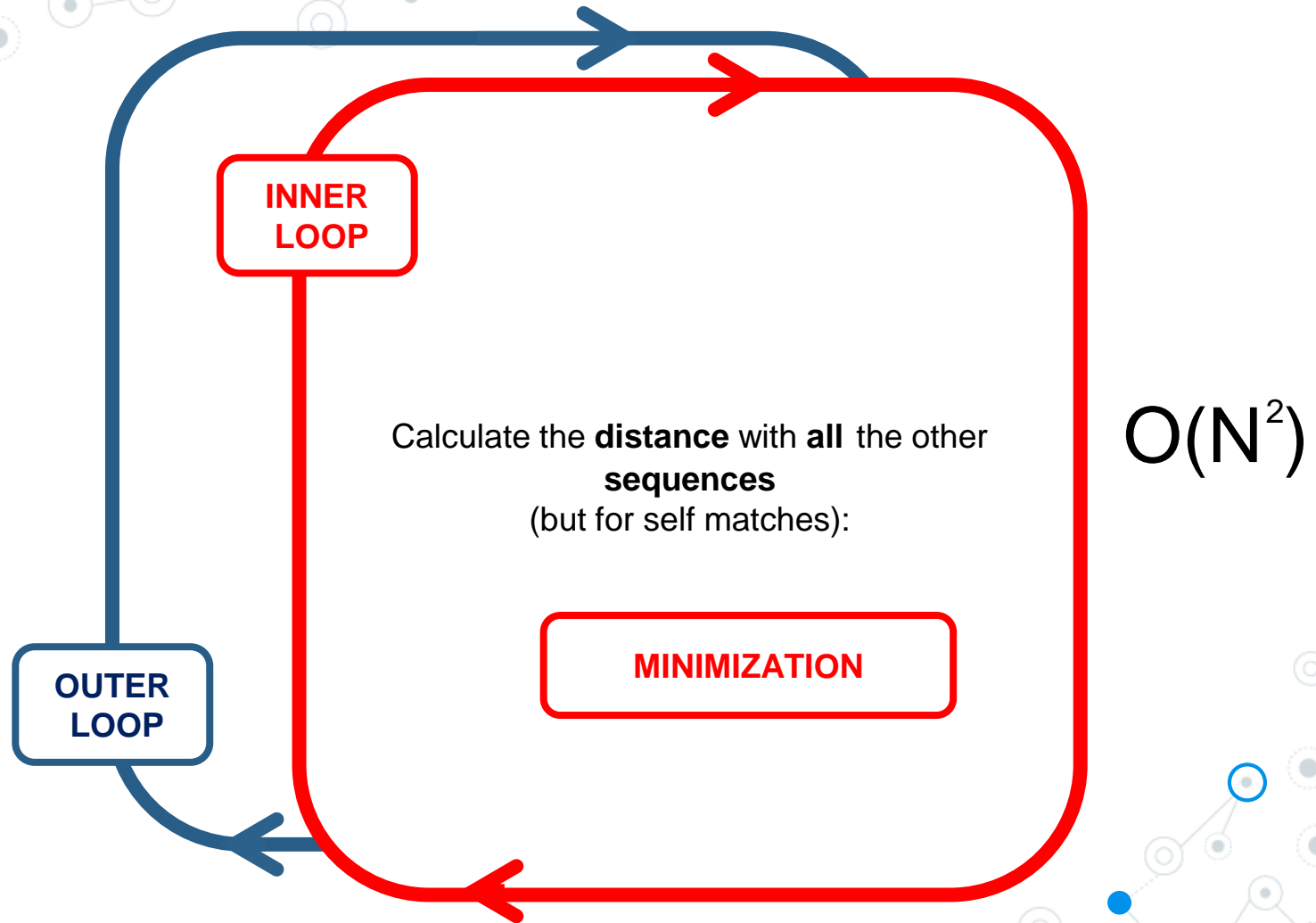
Aim

A fast and accurate **algorithm** to find an approximate **nearest neighbor distance (nnd)** profile for **all** the sequences of a time series

In order to obtain the **nnd** of **all** the sequences one needs **2** nested **loops**.

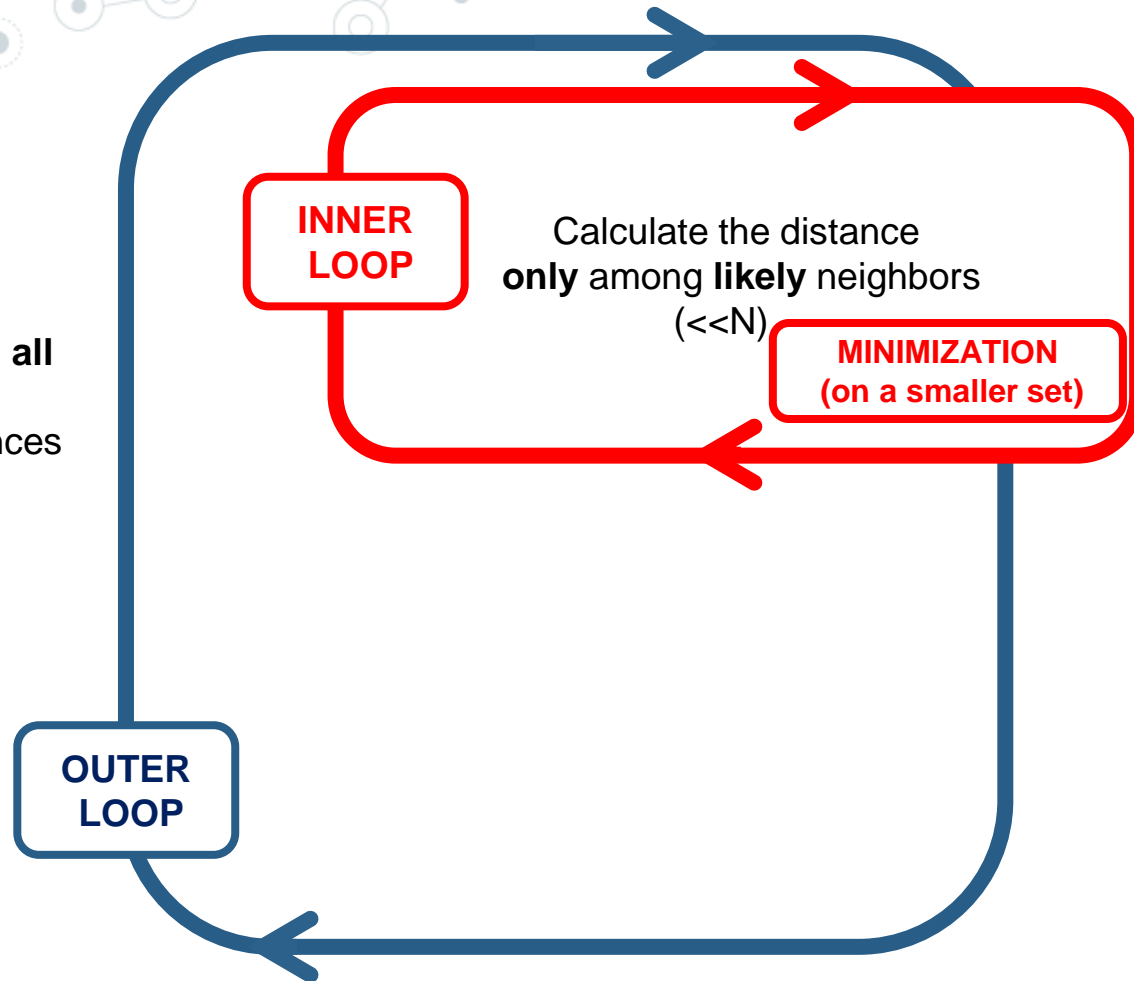


In order to obtain the **nnd** of **all** the sequences one needs **2** nested **loops**.

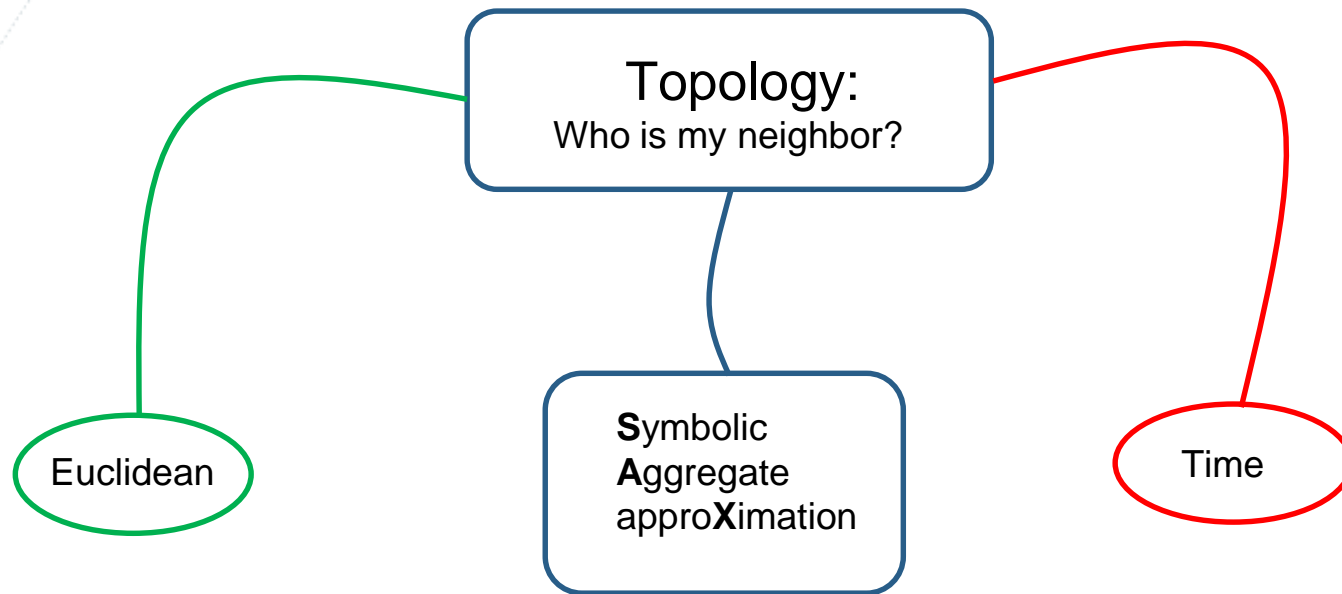


IDEA

In order to obtain the **nnd** of **all** the sequences one needs **2** nested **loops**.

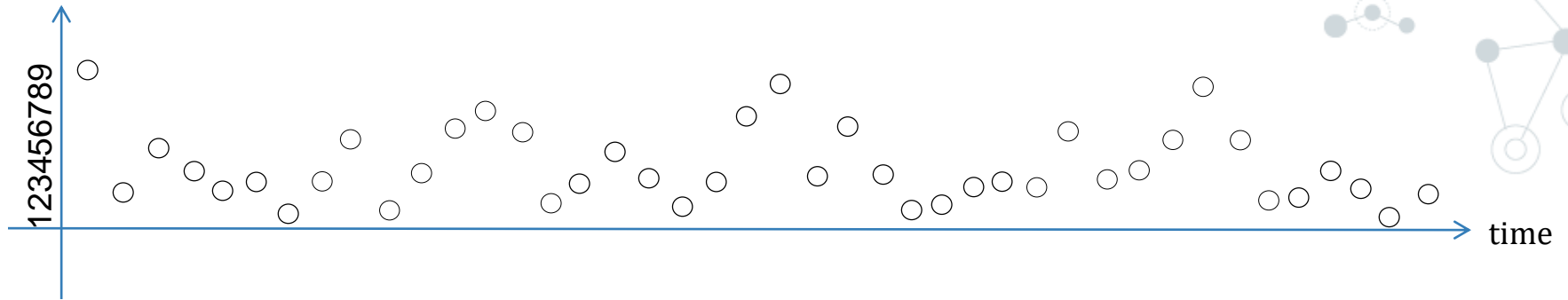


Let's **prune** the **inner loop**, thanks to **topologies**

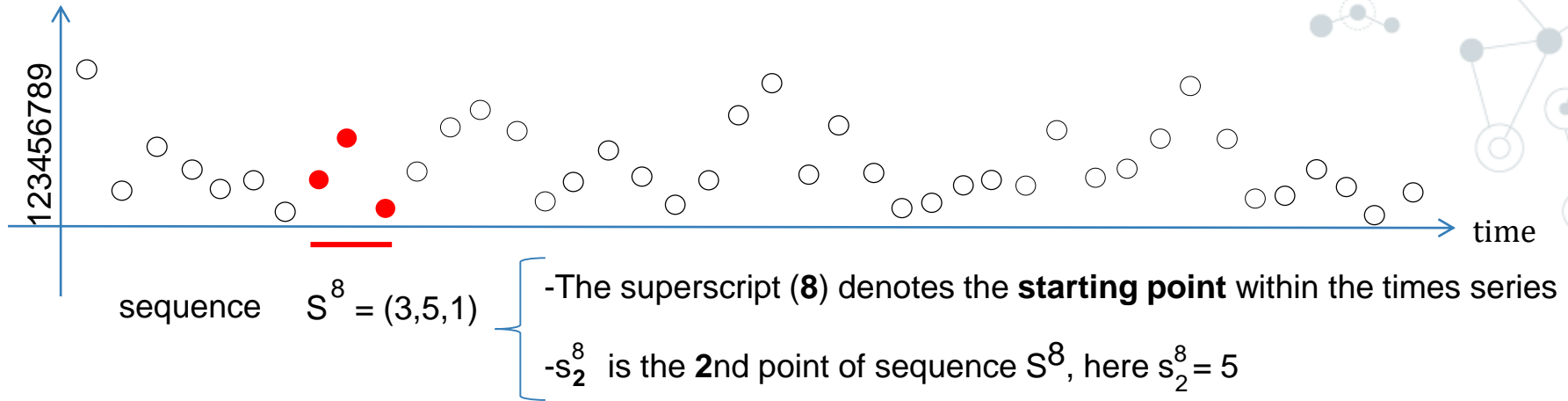


Topology

Euclidean

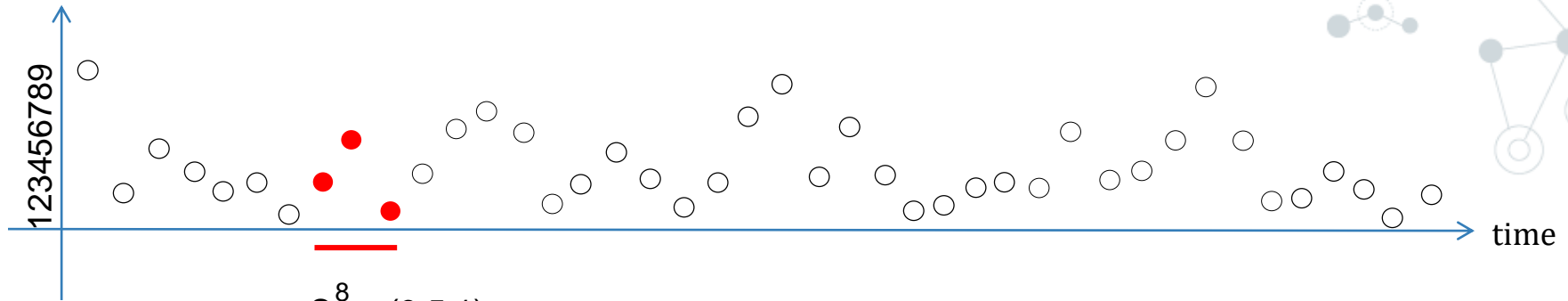


Euclidean

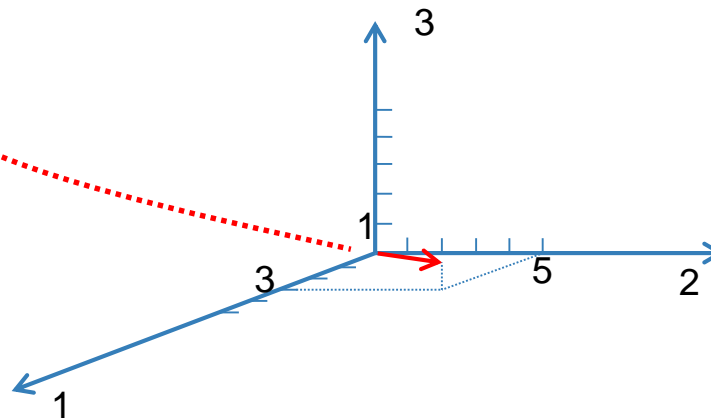


Topology

Euclidean

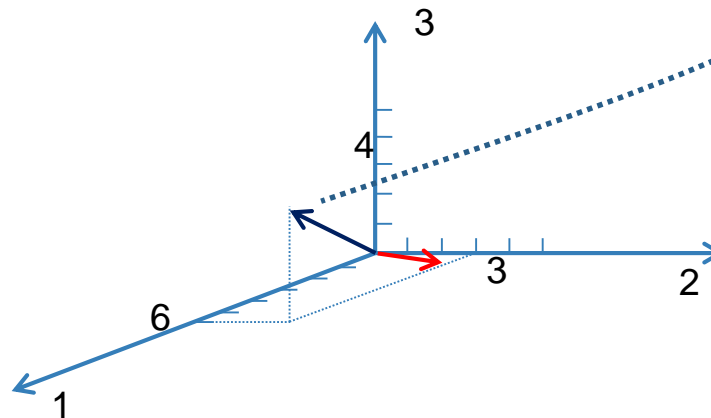
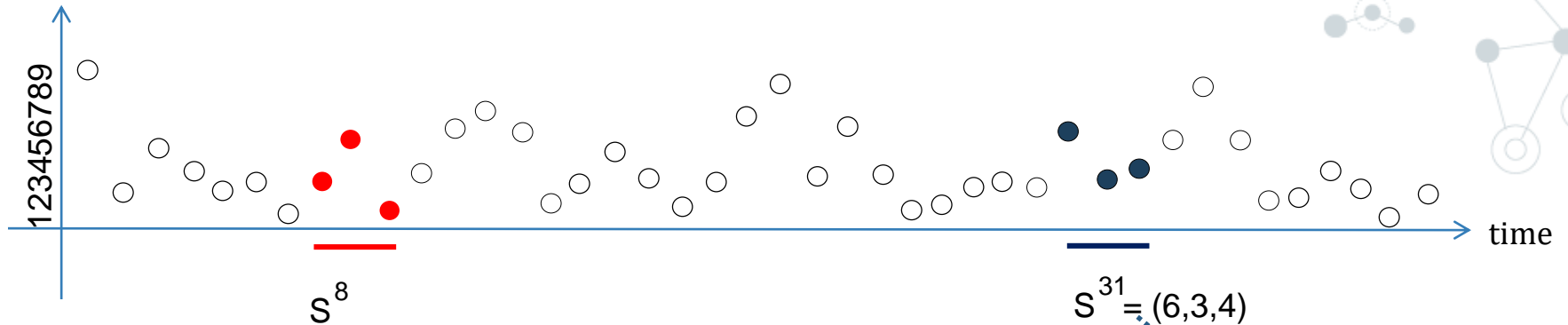


$$S^8 = (3, 5, 1)$$



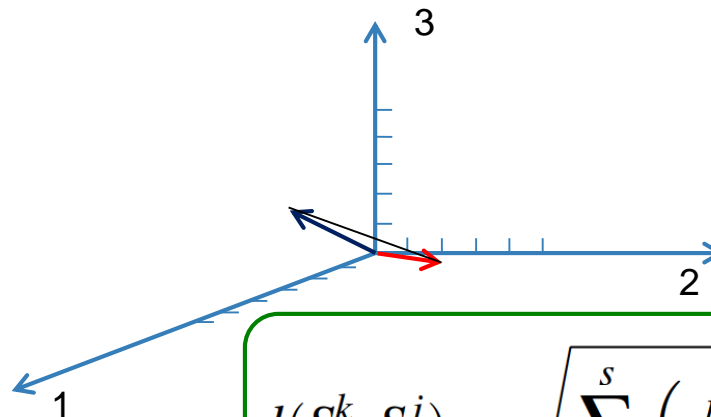
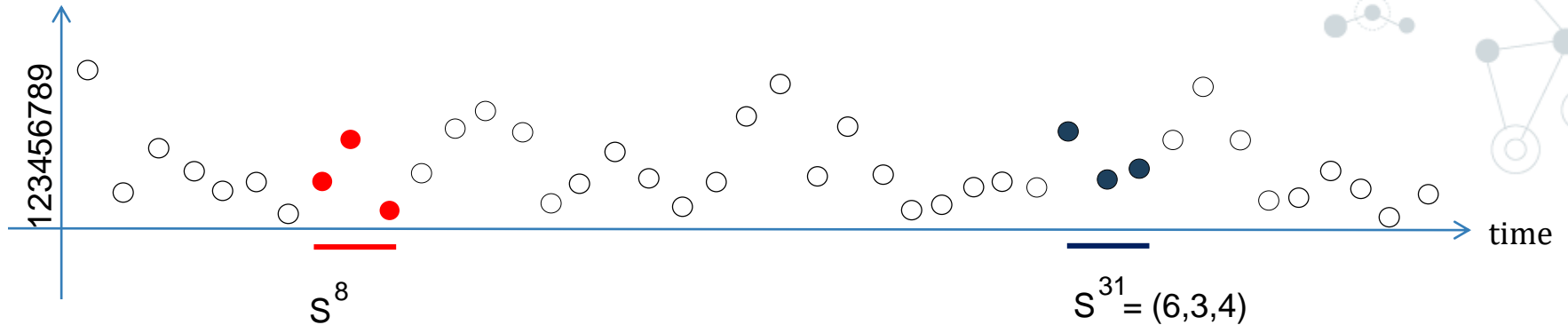
Topology

Euclidean



Topology

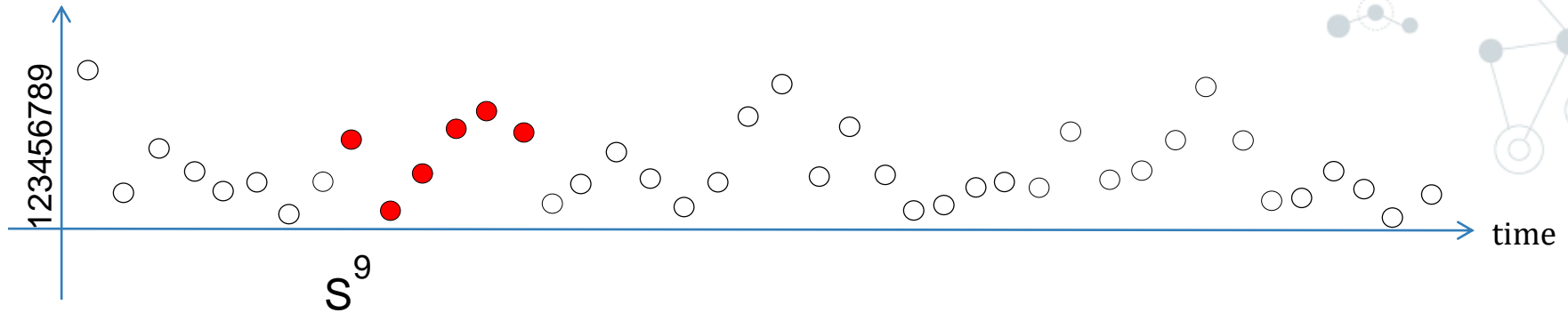
Euclidean



$$d(S^k, S^j) = \sqrt{\sum_{n=1}^s (s_n^k - s_n^j)^2}$$

Another neighborhood:

Time



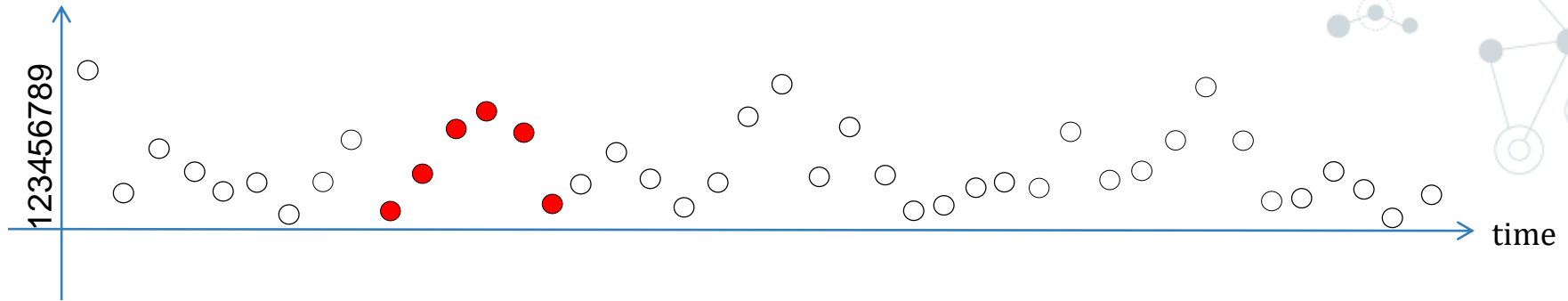
Another neighborhood:

Time



Another neighborhood:

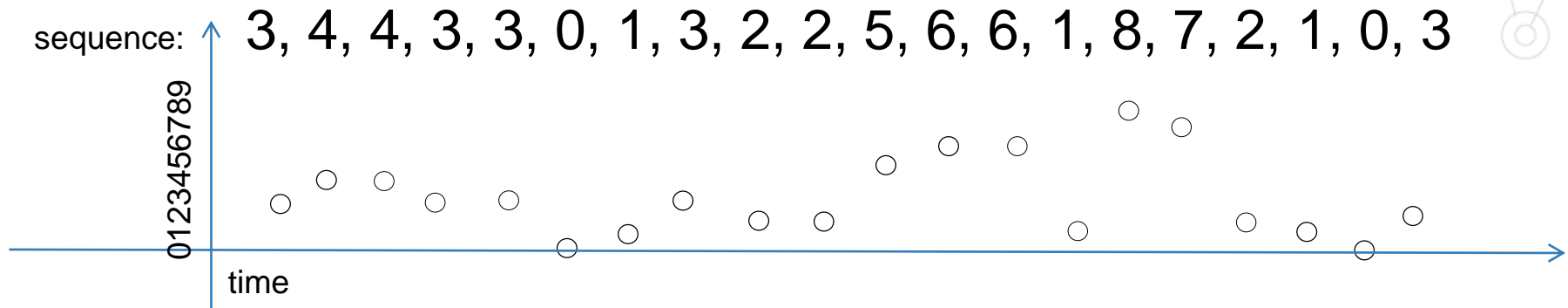
Time



$$d_t(S^k, S^j) = |k - j|$$

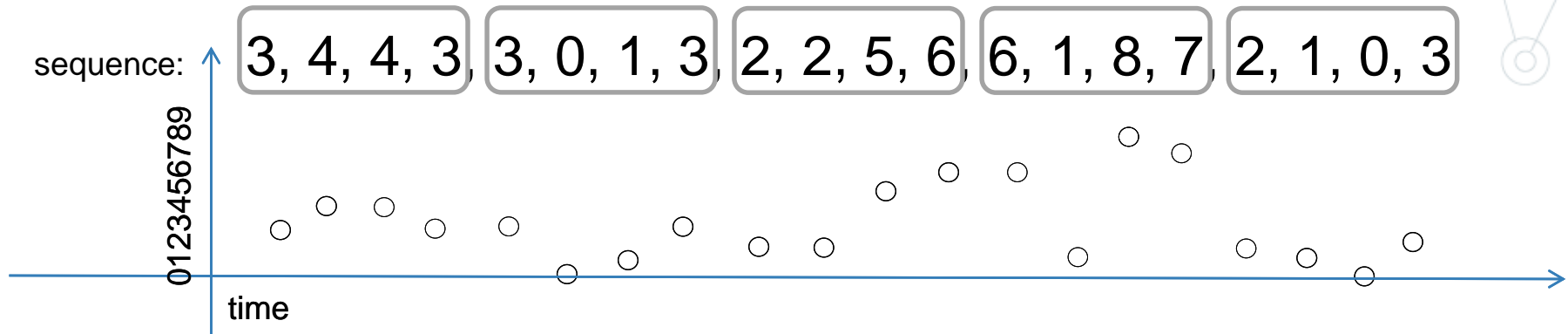
Symbolic
Aggregate
approximation

Topology



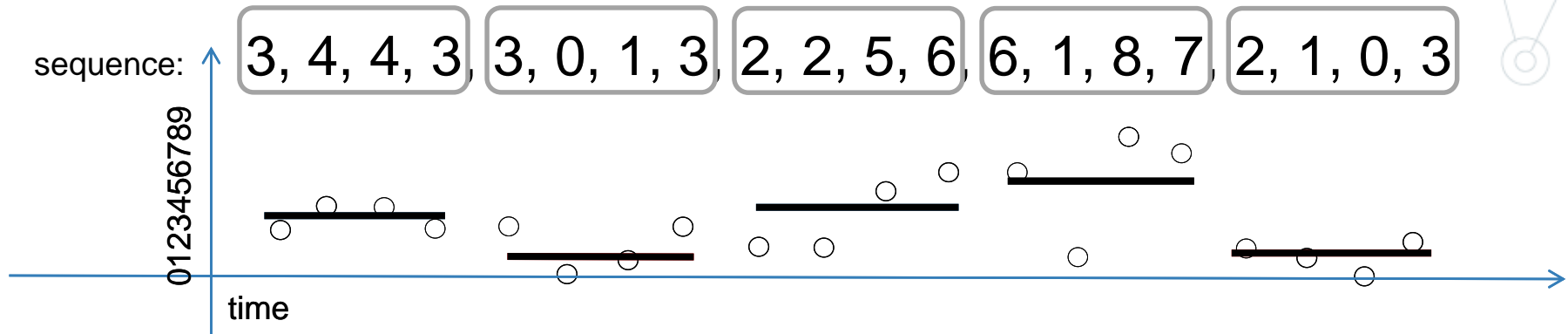
Symbolic
Aggregate
approximation

Topology



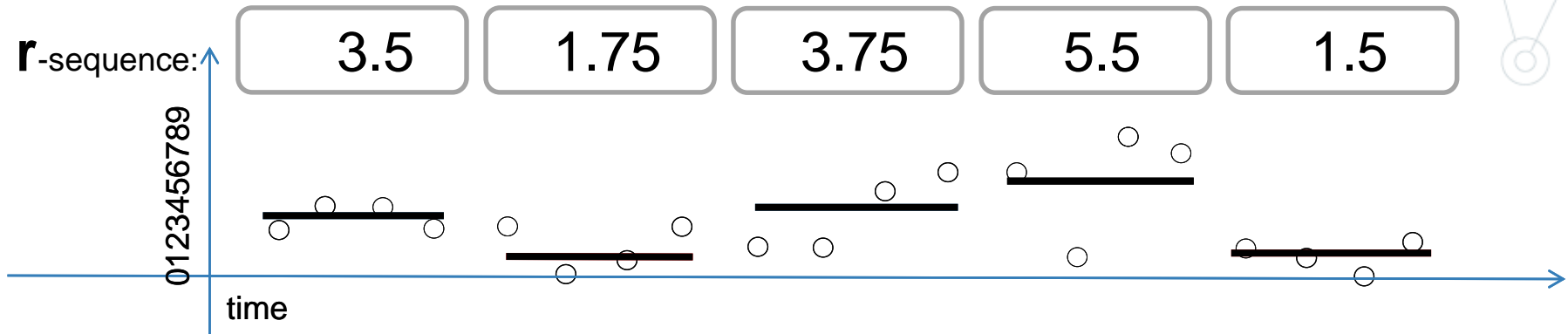
divide the sequence in parts

Symbolic
Aggregate
approximation



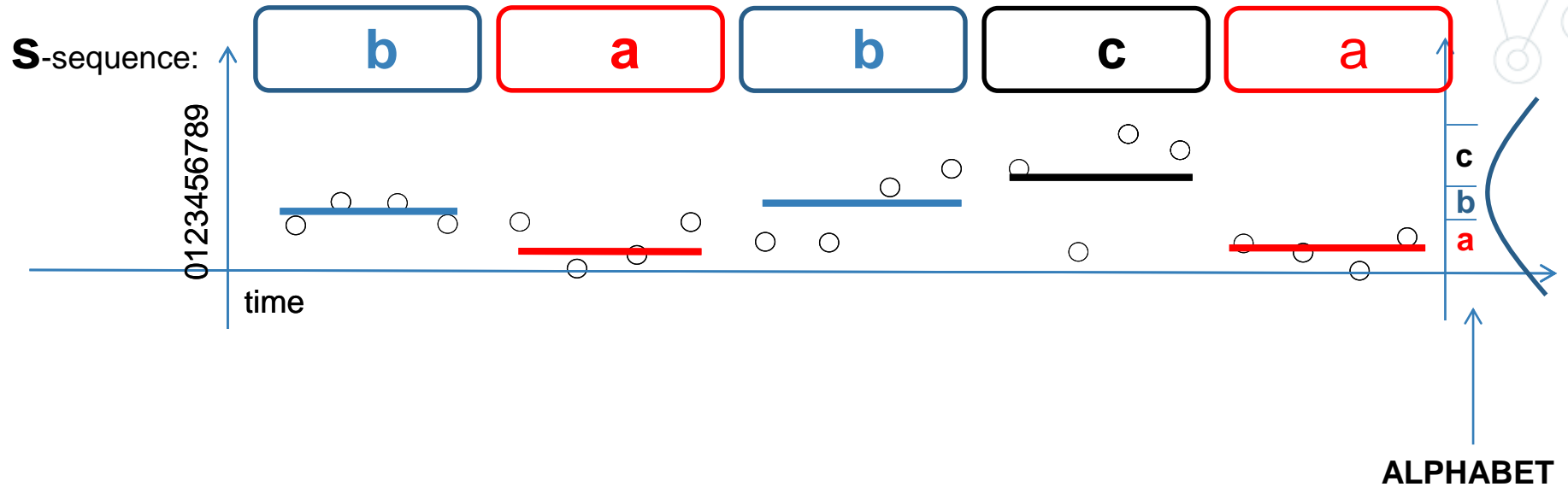
Symbolic
Aggregate
approximation

Topology



Piecewise Aggregate Approximation (PAA)

Symbolic
Aggregate
approximation



Symbolic
Aggregate
approximation

sequence

34433013225661872103

with 2 parameters:
- the number of sub-sequences
- the number of symbols

Symbolic
Aggregate
approximation

sequence

34433013225661872103

Symbolic sequence

babca

A different form of
Neighborhood:

Symbolic
Aggregate
approximation

Topology

Many sequences

S^3

S^{16}

S^{97}

S^{684}

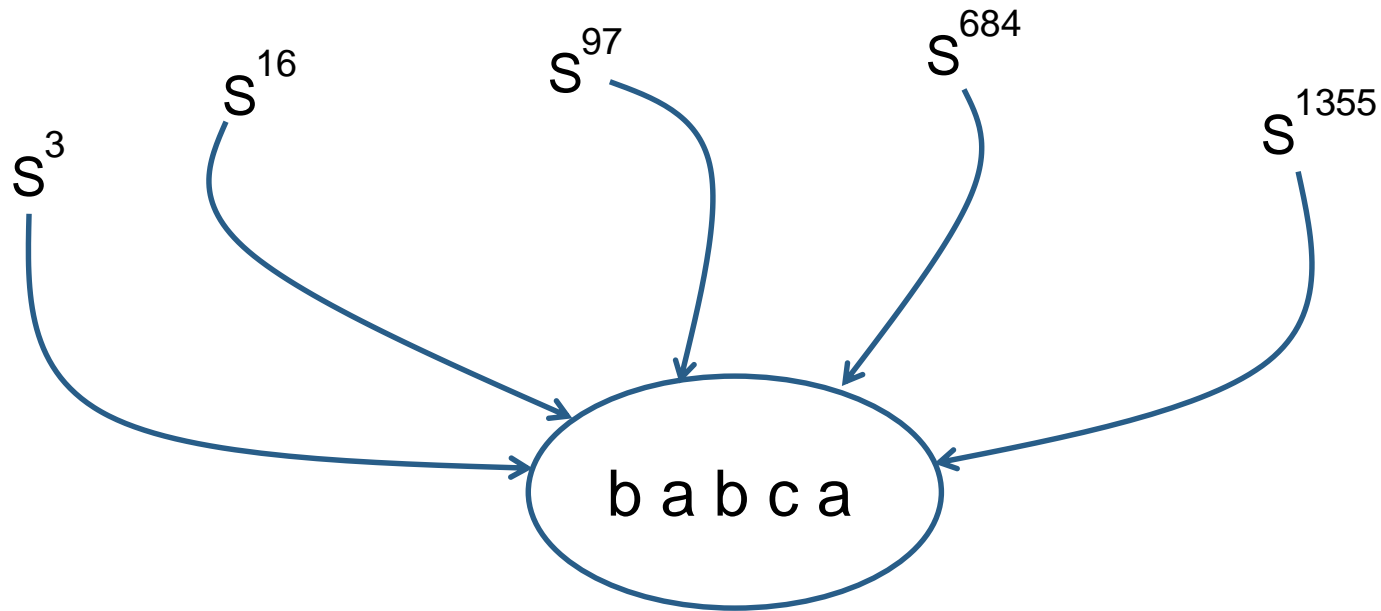
S^{1355}

Symbolic Aggregate approximation

Topology

A different form of
Neighborhood:

Many sequences



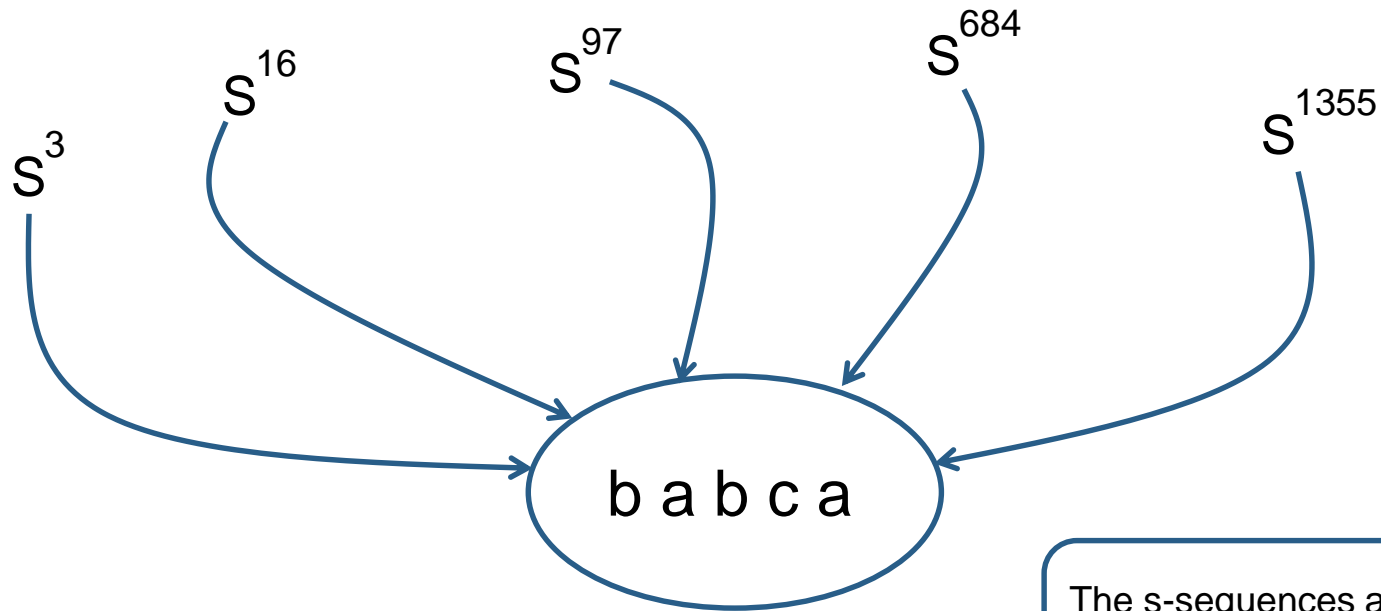
Symbolic-sequence

Symbolic Aggregate approximation

Topology

A different form of
Neighborhood:

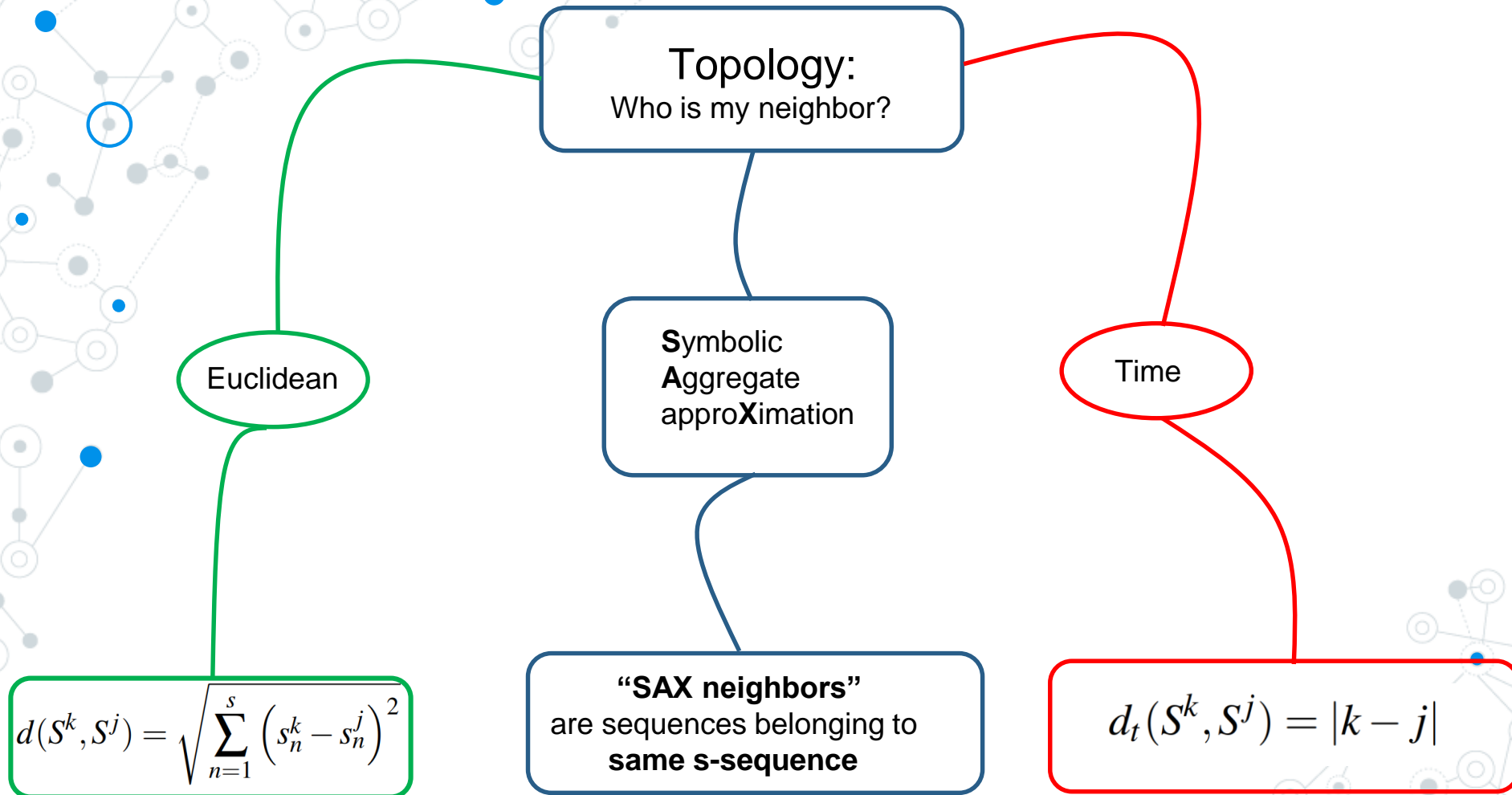
Many sequences

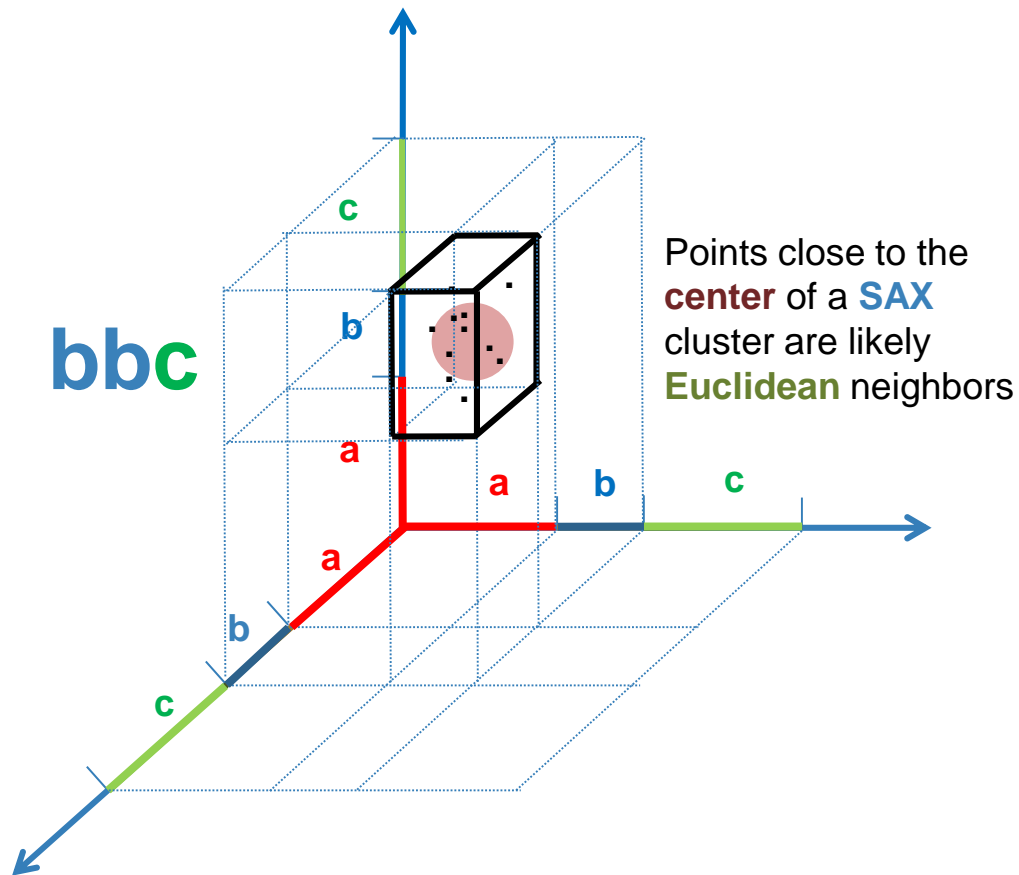


Symbolic-sequence

The s-sequences are
natural **clusters** of
sequences!

Topology





Topology

Diagram illustrating the body-centered cubic (bcc) unit cell. The unit cell is a cube with a smaller cube inside, representing the bcc structure. The axes are labeled a , b , and c . The unit cell is labeled **bcc**.

For points close to the **border**, the **SAX-Euclidean** connection is more loose!

HOT SAX, an
algorithm for
quickly
finding **the**
sequence with the
highest *nnd*

SAX **EUCLIDEAN**

Algorithm

Algorithm

SAX EUCLIDEAN

HOT SAX, an algorithm for quickly finding the sequence with the highest *nnd*

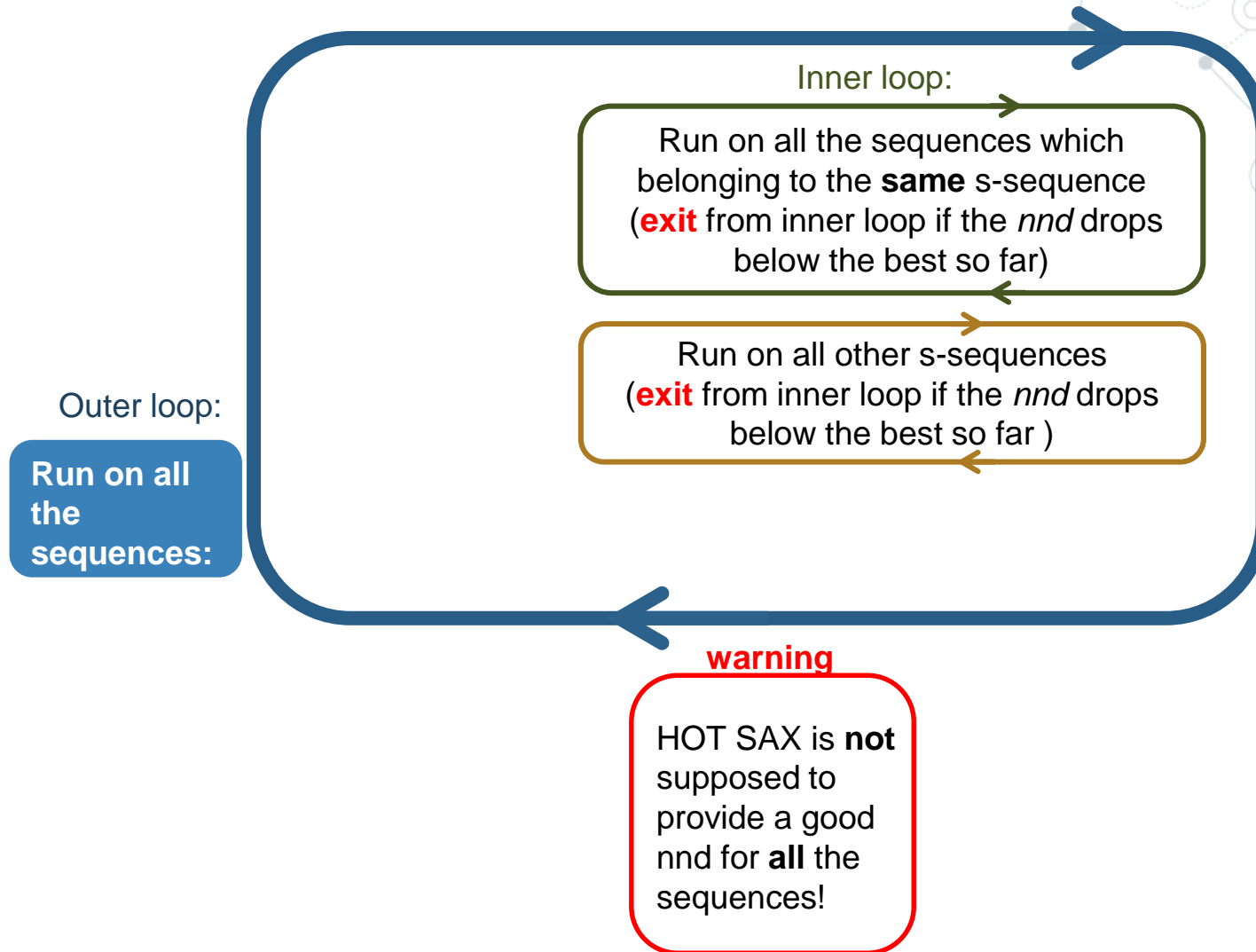
Outer loop:

Run on all the sequences:

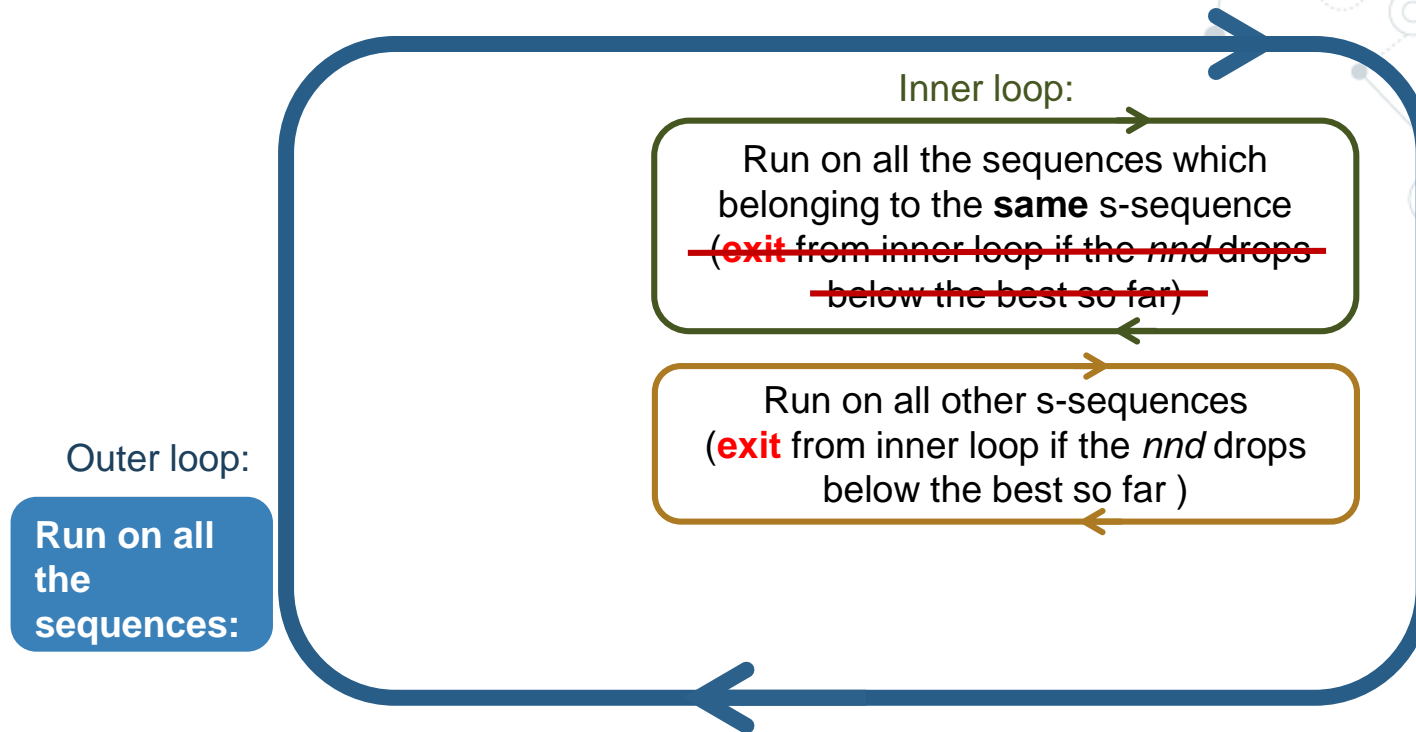
Inner loop:

Run on all the sequences which belonging to the **same** s-sequence
(**exit** from inner loop if the *nnd* drops below the best so far)

Run on all other s-sequences
(**exit** from inner loop if the *nnd* drops below the best so far)

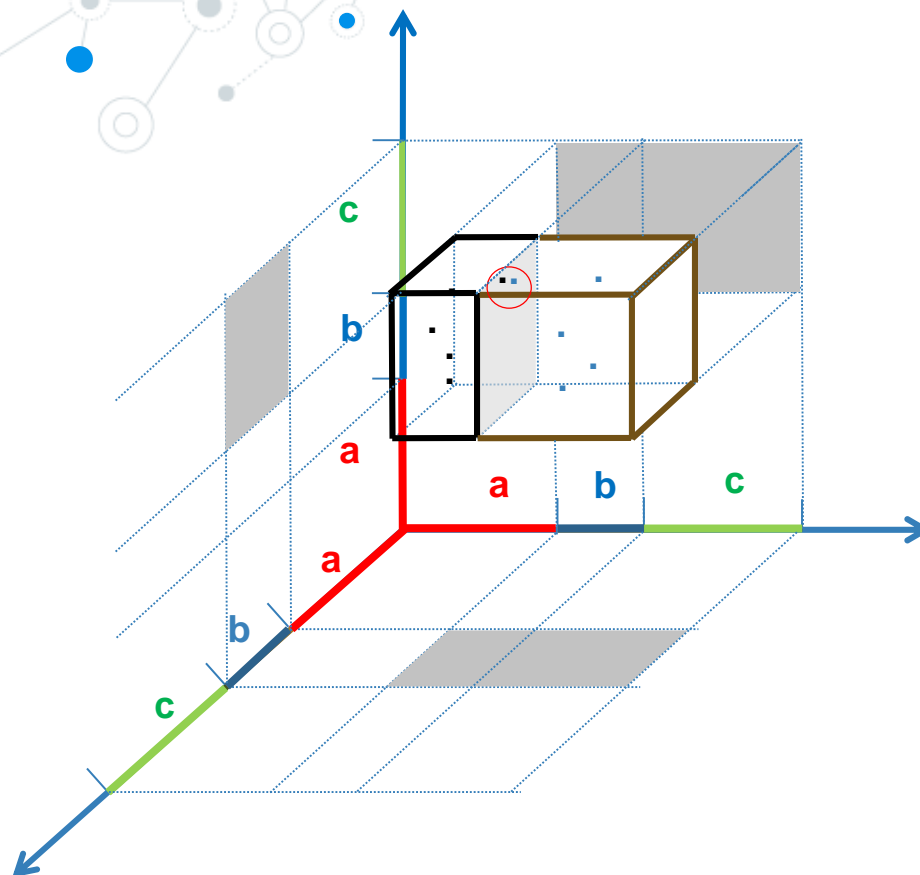


Algorithm

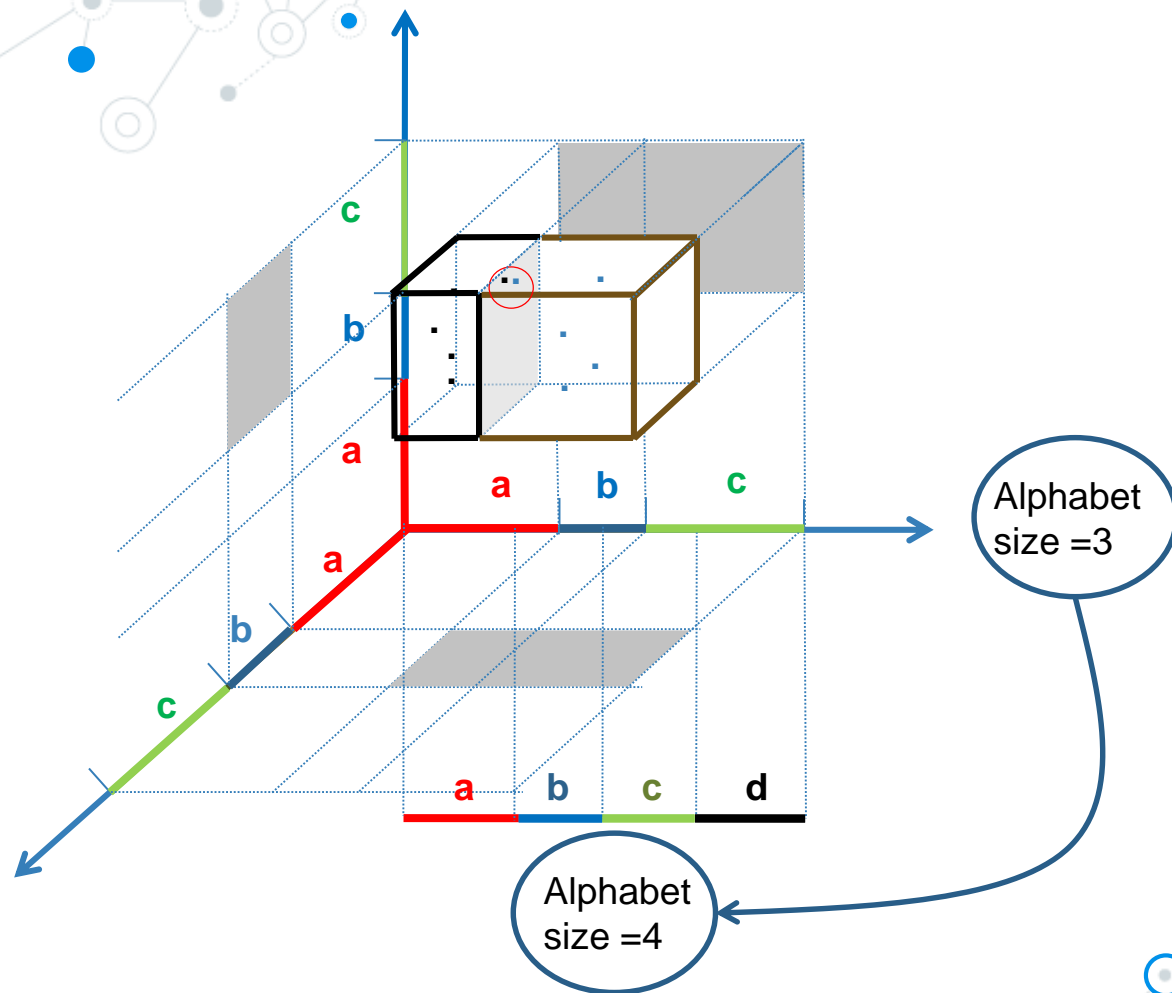


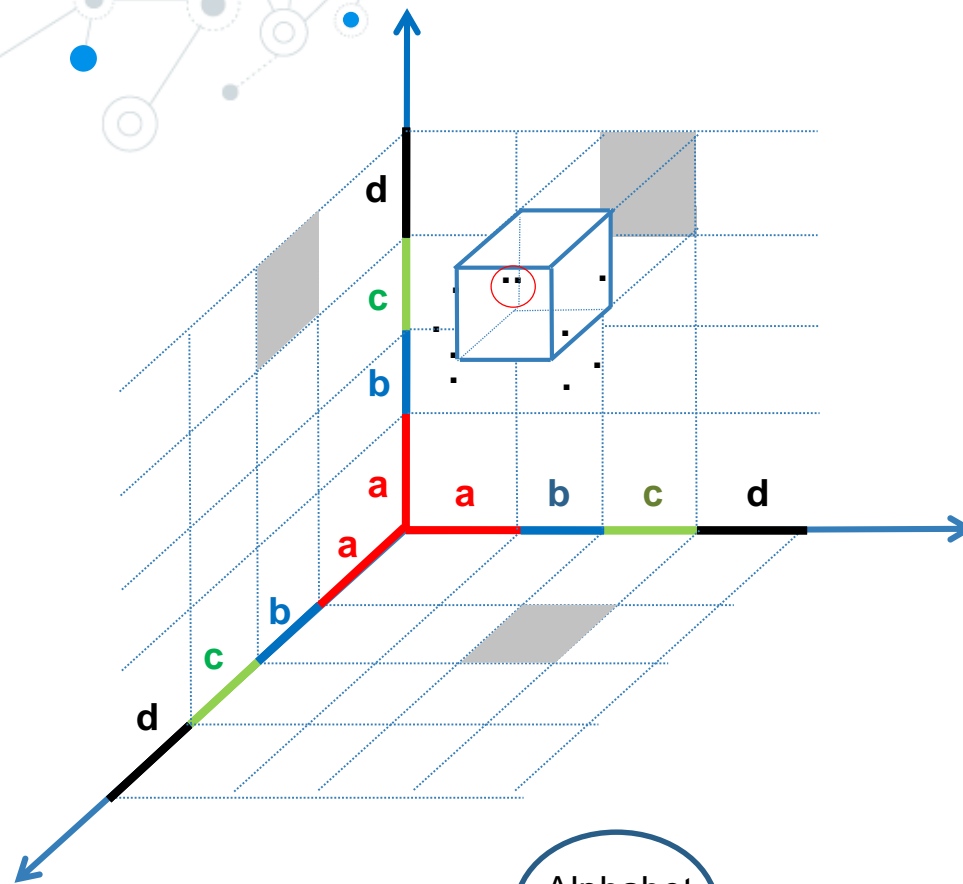
1

Let's force a full run on the same cluster!



Alphabet
size =3





Alphabet
size =4

Algorithm

2

Let's force a full run on the same cluster, but with alphabet +1

Outer loop:

Run on all the sequences

Inner loop:

Run on all the sequences which belonging to the **same** s-sequence

Run on all the sequences which belonging to the **same** s-sequence (alphabet + 1)

Run on all other s-sequences (**exit** from inner loop if the *nnd* drops below the best so far)

Unfortunately as the **dimension** of the clusters increases,
more and more points are doomed to be **close to the borders**!

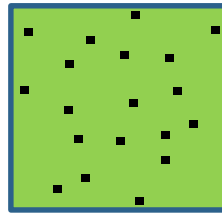
CURSE

s-sequences  hypercubes

Unfortunately as the **dimension** of the clusters increases more and more points are doomed to be **close to the borders**!

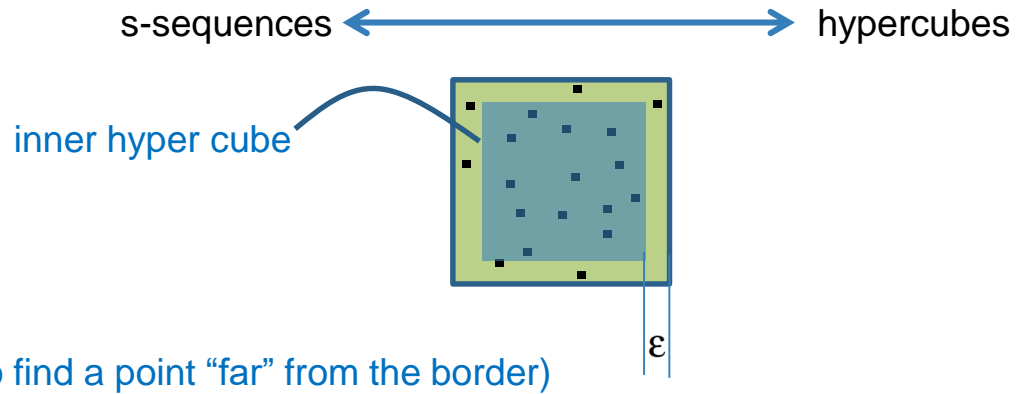
CURSE

s-sequences  hypercubes



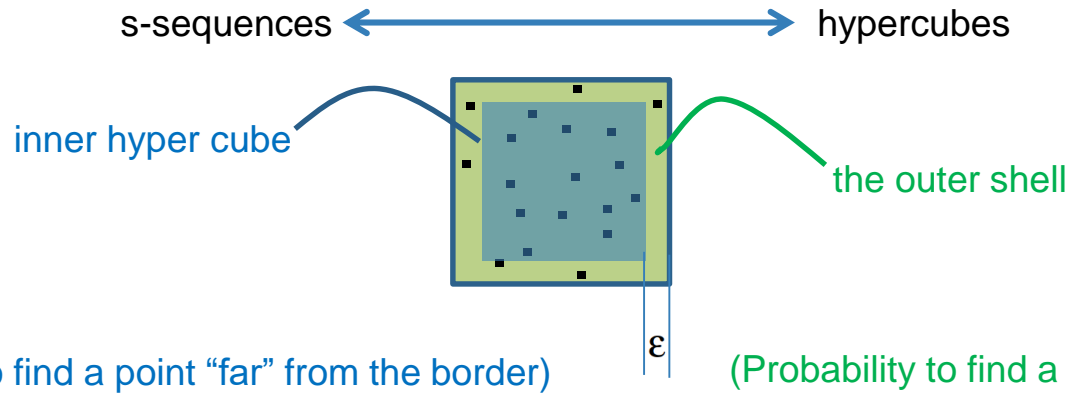
Unfortunately as the **dimension** of the clusters increases more and more points are doomed to be **close to the borders**!

CURSE



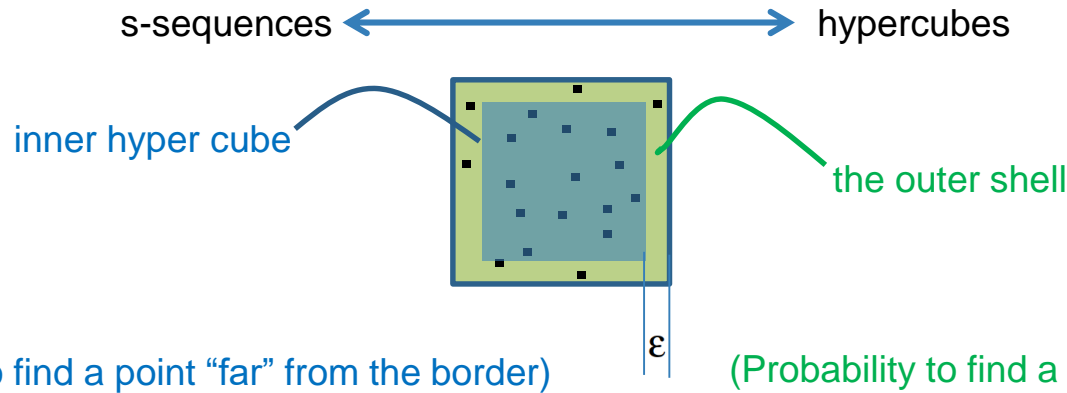
Unfortunately as the **dimension** of the clusters increase more and more points are doomed to be **close to the borders**!

CURSE



Unfortunately as the **dimension** of the clusters increase more and more points are doomed to be **close to the borders**!

CURSE

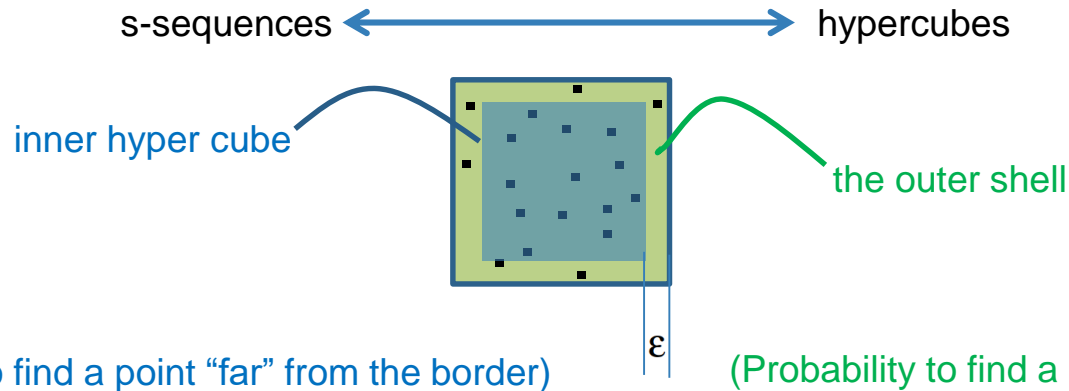


$$\frac{\text{inner volume}}{\text{volume}} = \left(\frac{l - 2\varepsilon}{l} \right)^n \xrightarrow{n \rightarrow \infty} 0$$

For higher dimensions most of the volume of a hypercube is close to its borders!

Unfortunately as the **dimension** of the clusters increase more and more points are doomed to be **close to the borders**!

CURSE

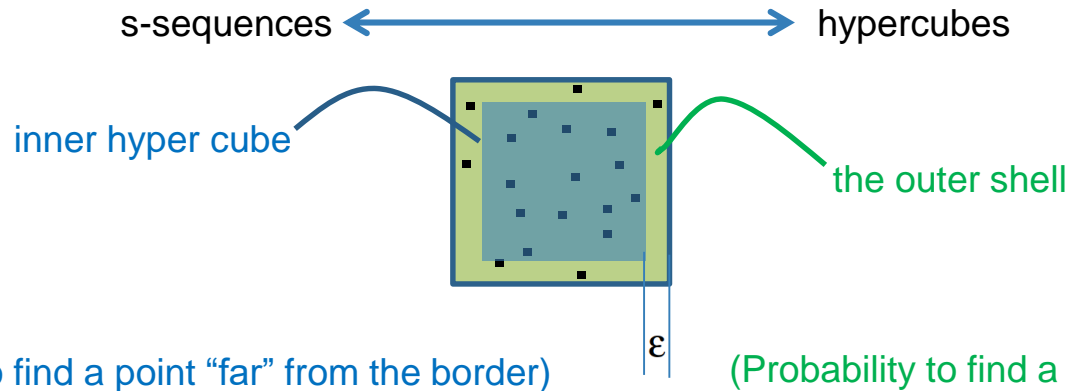


Curse of dimensionality

As n grows the **SAX-Euclidean** connection neighborhood becomes looser

Unfortunately as the **dimension** of the clusters increase more and more points are doomed to be **close to the borders**!

CURSE



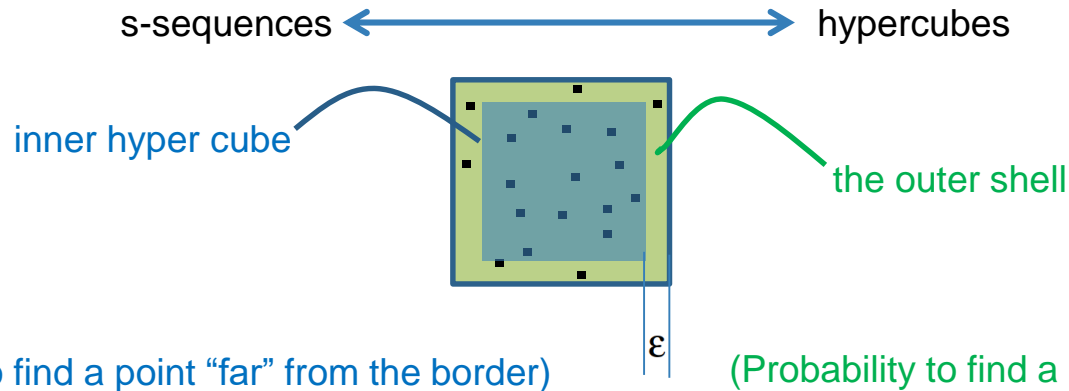
Curse of dimensionality

As n grows the **SAX-Euclidean** connection neighborhood becomes looser

There are too many close cluster!

Unfortunately as the **dimension** of the clusters increase more and more points are doomed to be **close to the borders**!

CURSE



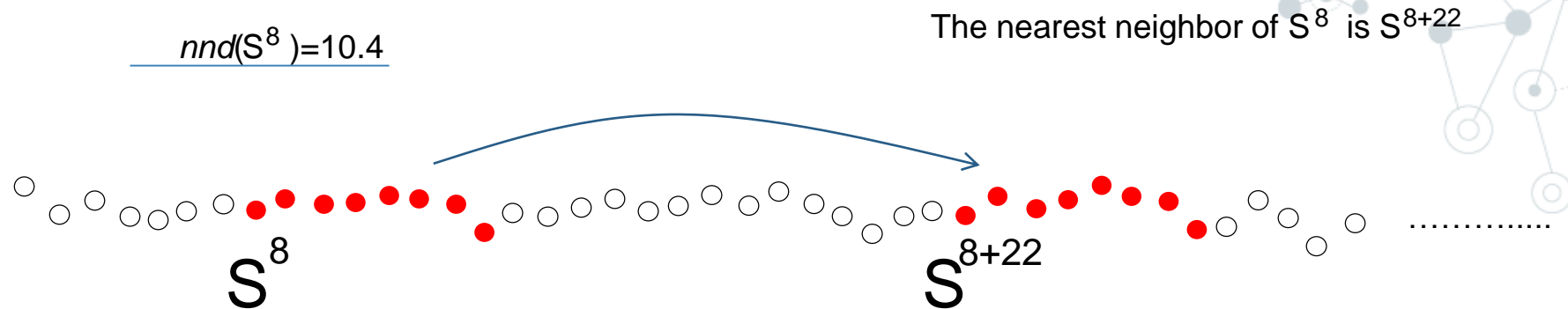
Curse of dimensionality

As n grows the **SAX-Euclidean** connection neighborhood becomes looser

There are too many close cluster!

Next?

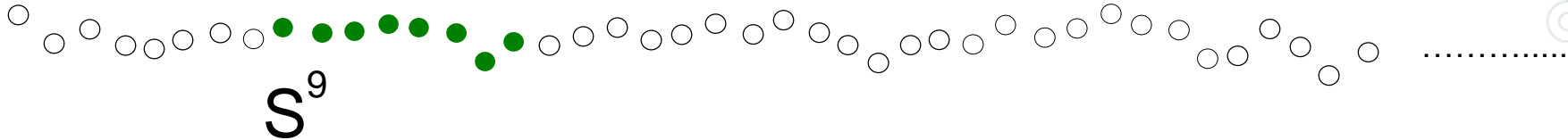
Solution: time topology



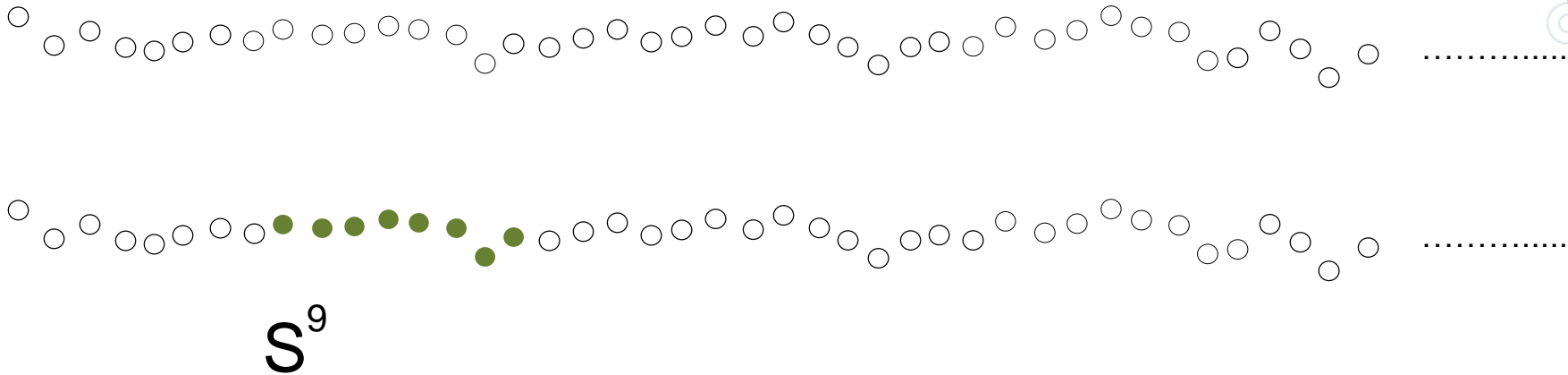
Solution: time topology

The nearest neighbor of S^9 is S^{453}

$$nnd(S^9)=17.8$$

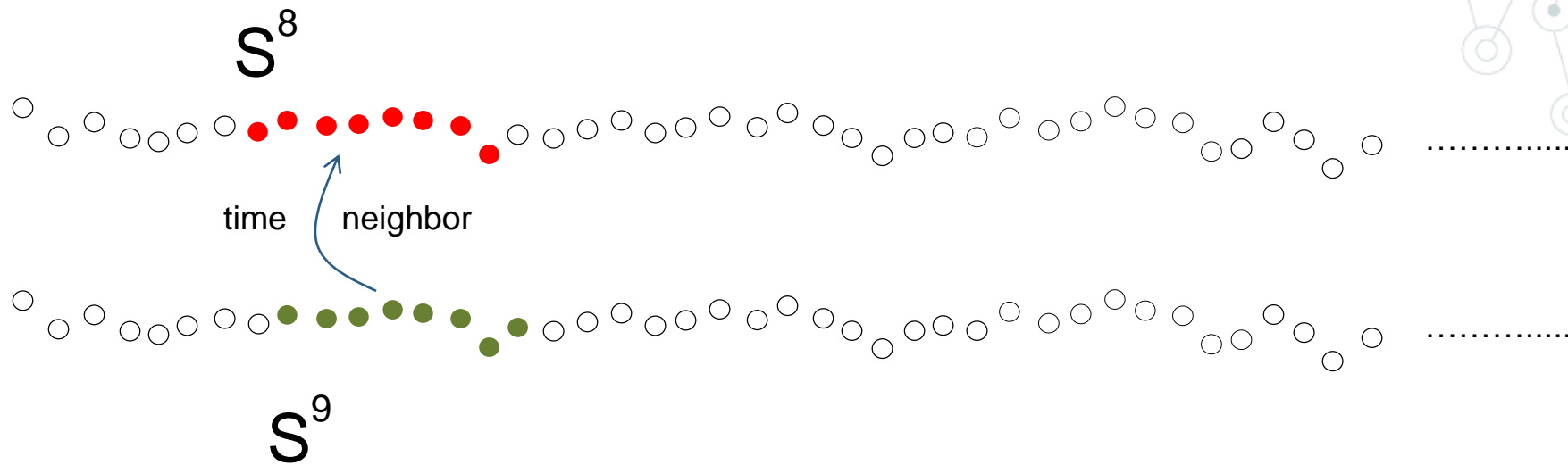


Solution: time topology

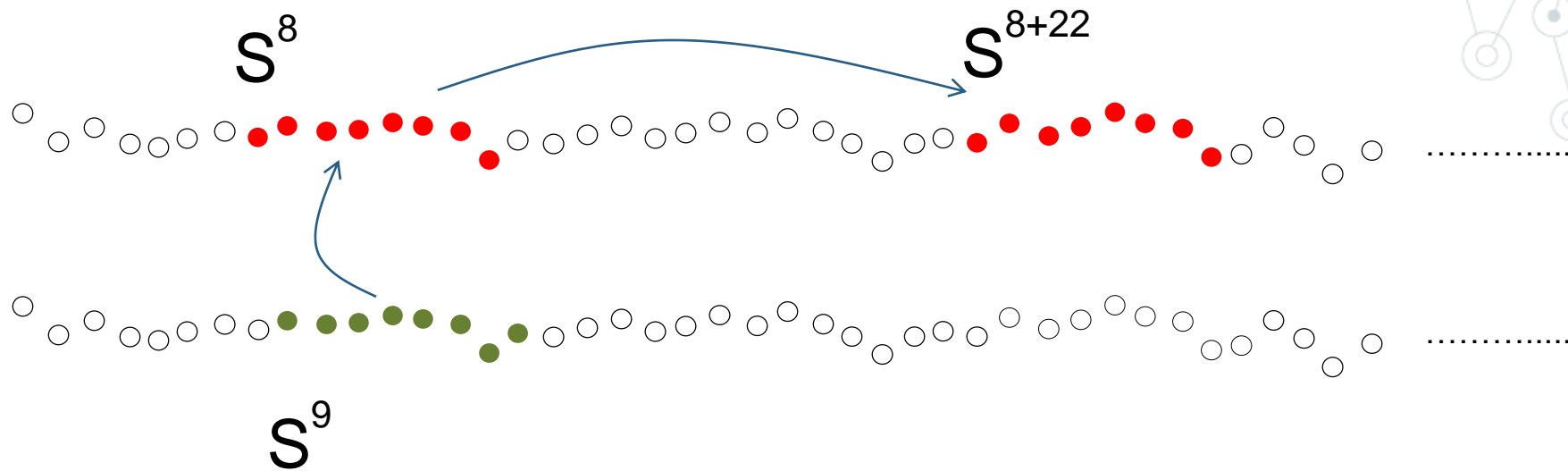


Solution: time topology

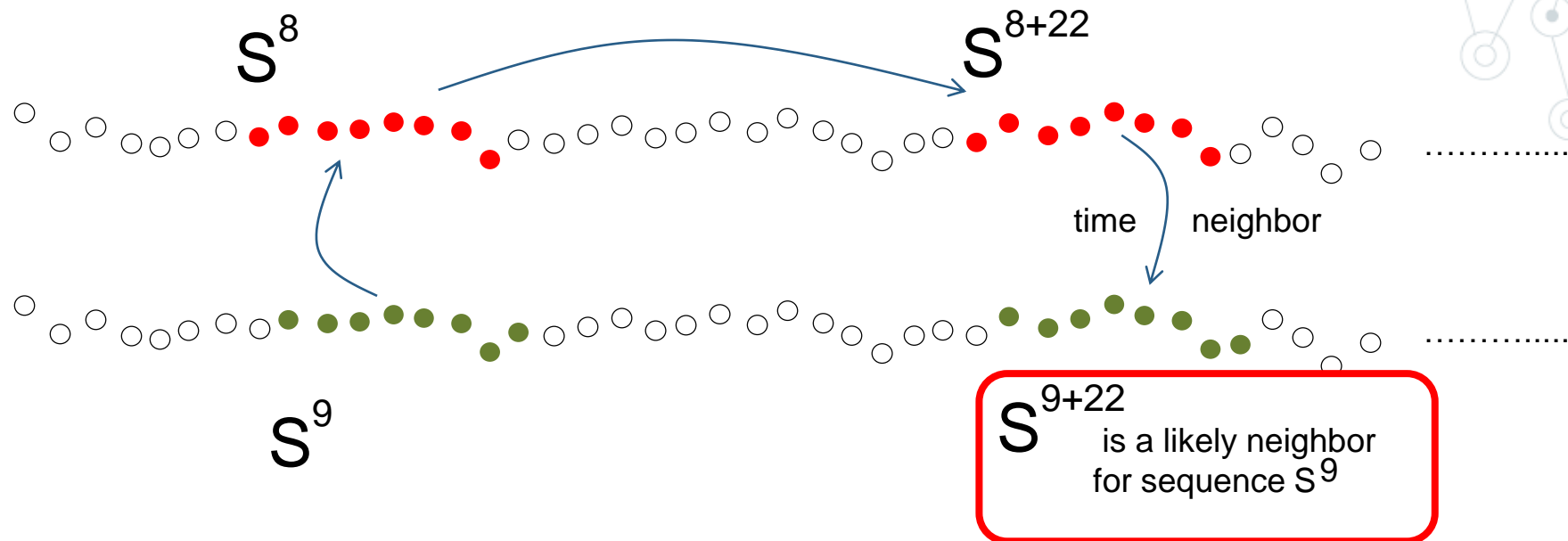
Since $nnd(S^8) \ll nnd(S^9)$



Solution: time topology

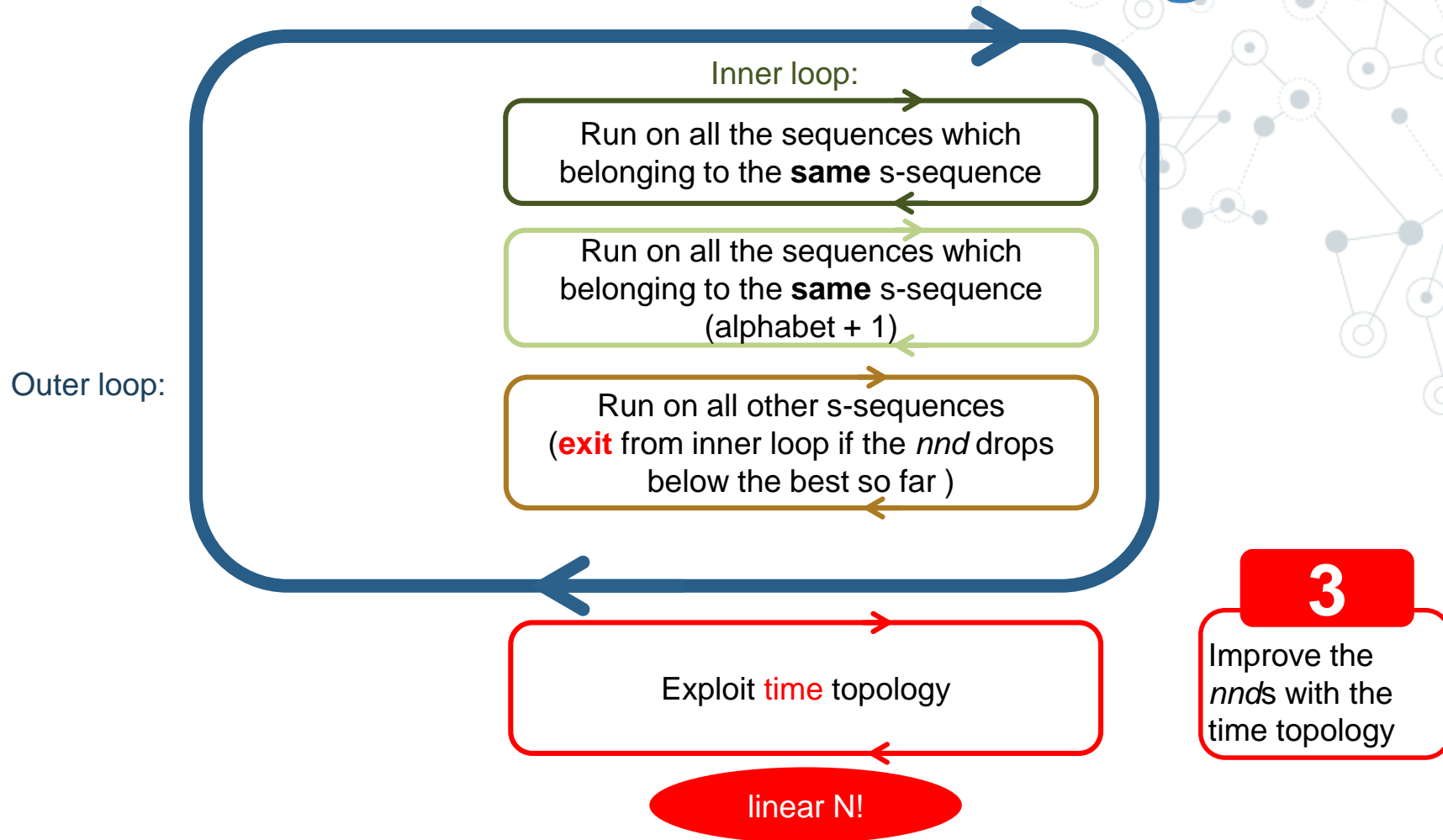


Solution: time topology

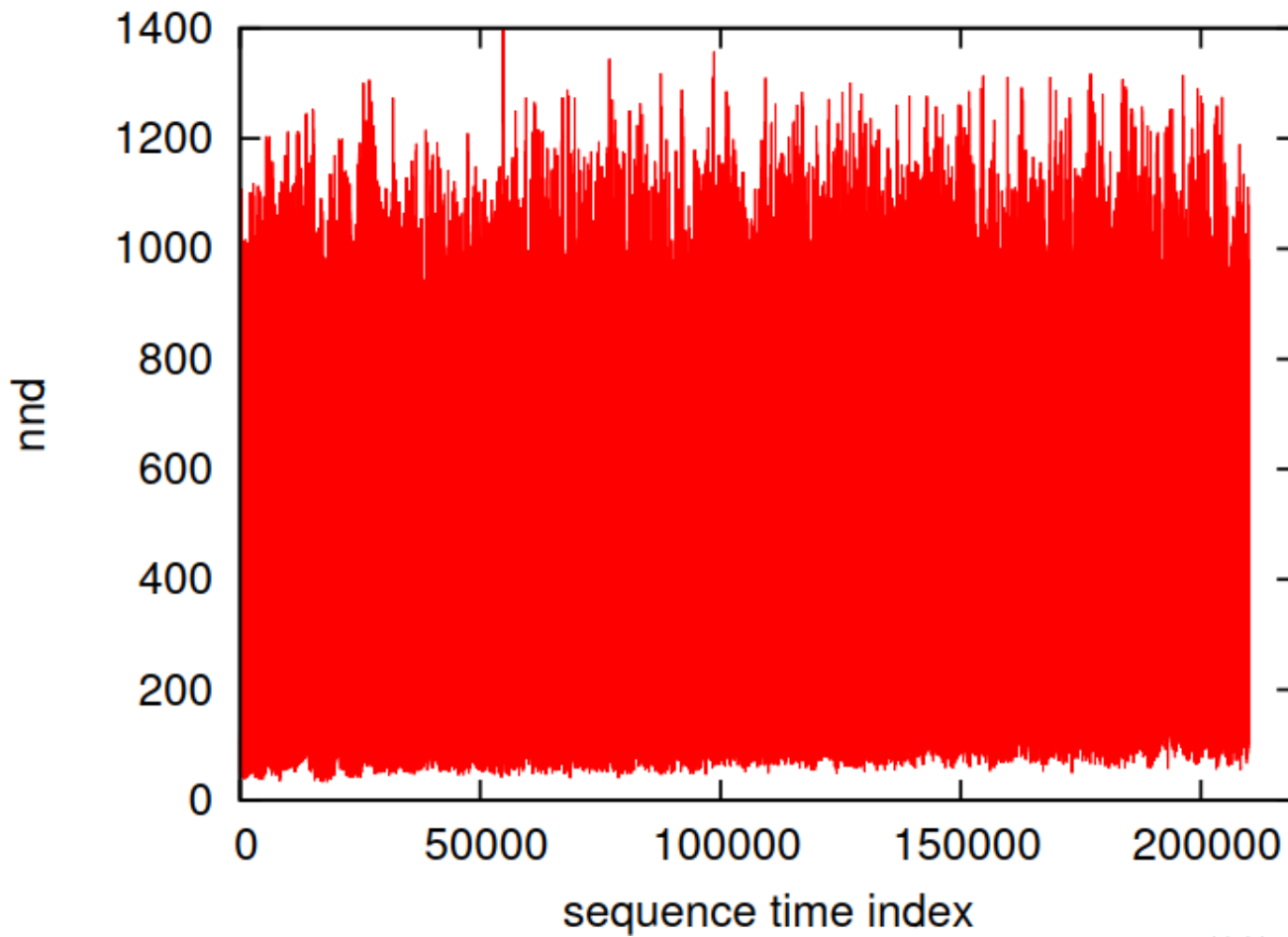


$$nnd(S^{i+1}) \leq d_2(S^{i+1}, S^{i+k+1}) = nnd(S^i) + (p_{i+s} - p_{i+s+k})^2 - (p_i - p_{i+k})^2$$

Algorithm



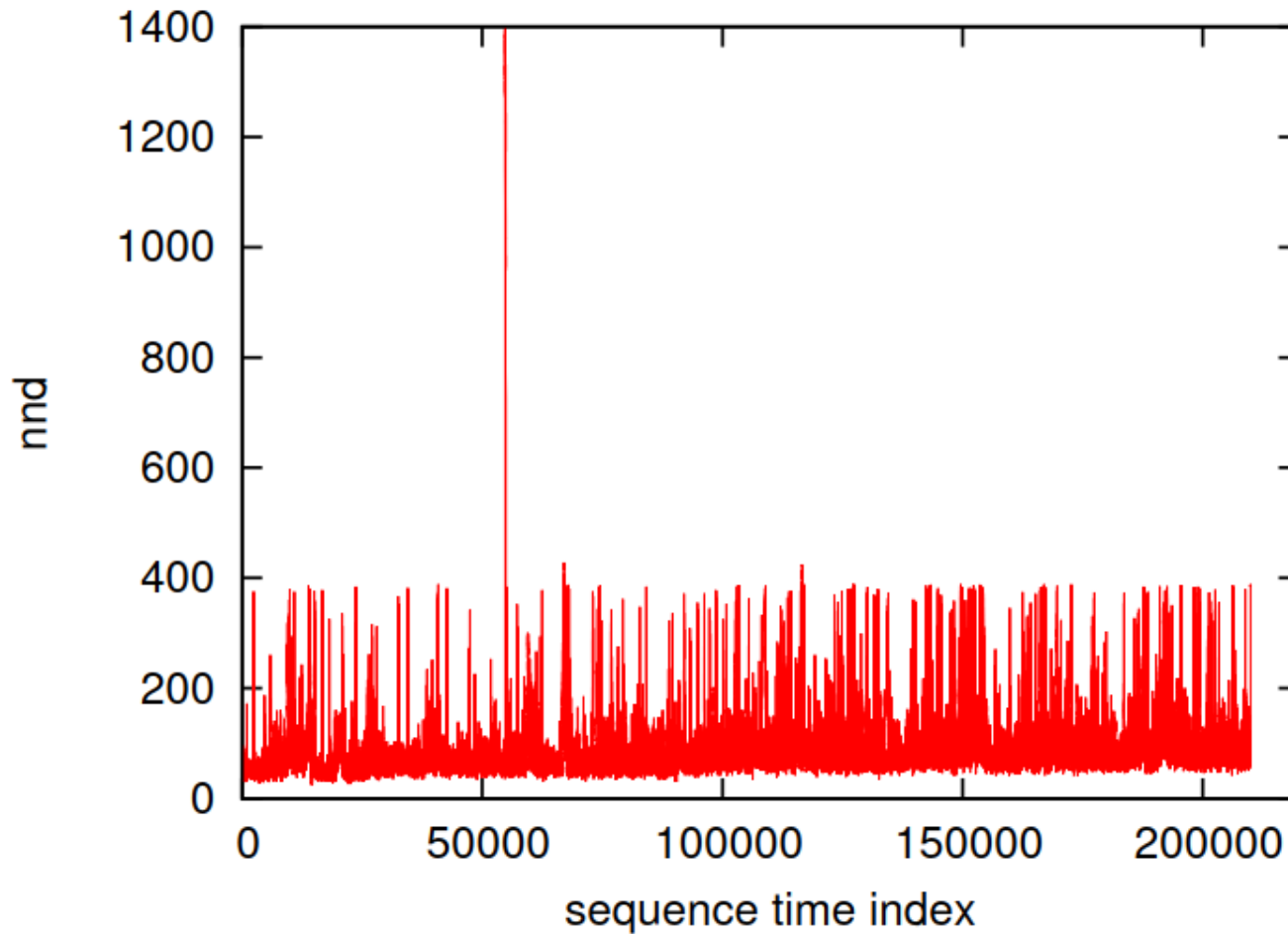
VALIDATION



0

The approx.
nnds obtained
with HOT SAX

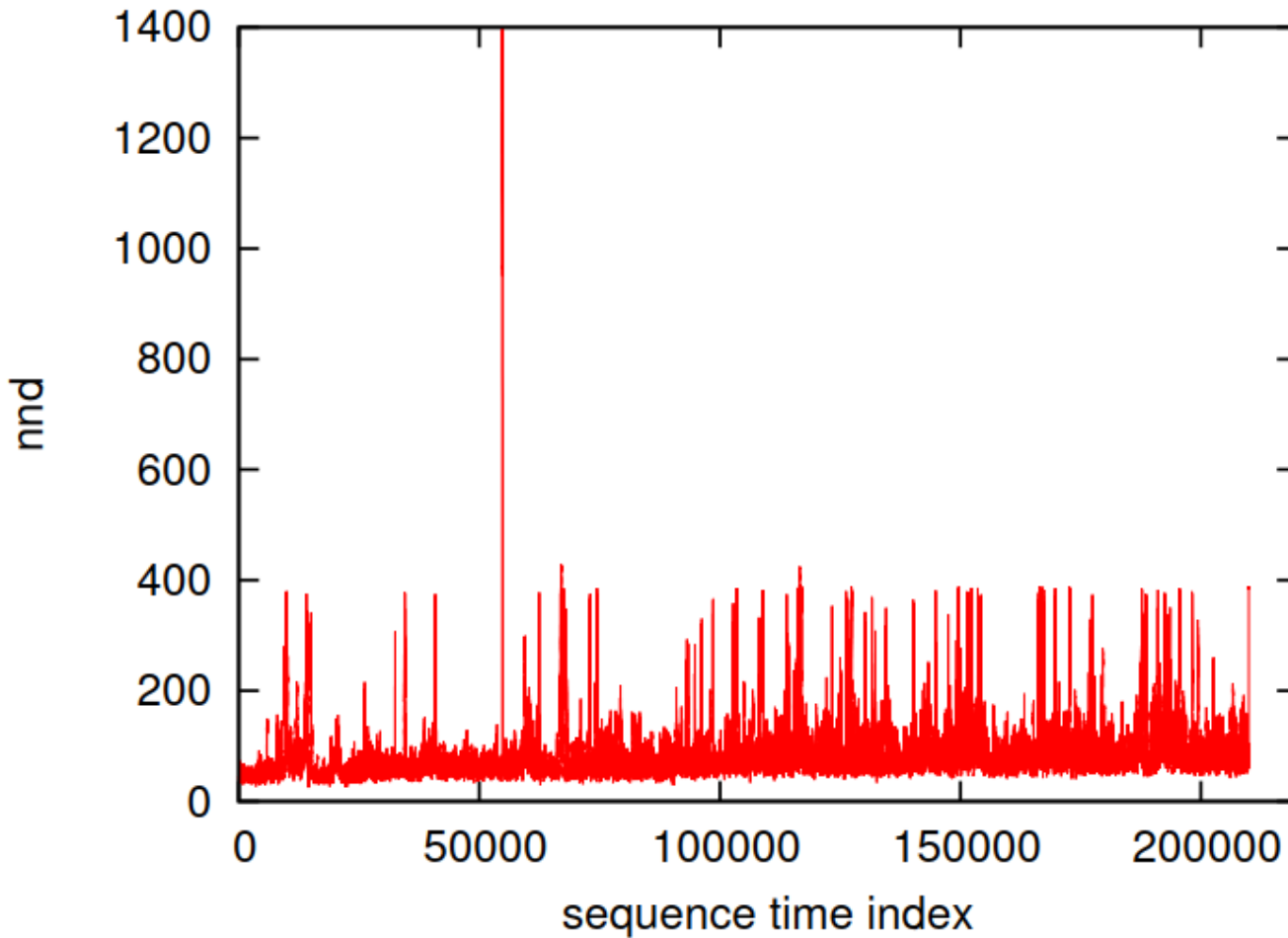
VALIDATION



1

Let's force a full run on the same cluster!

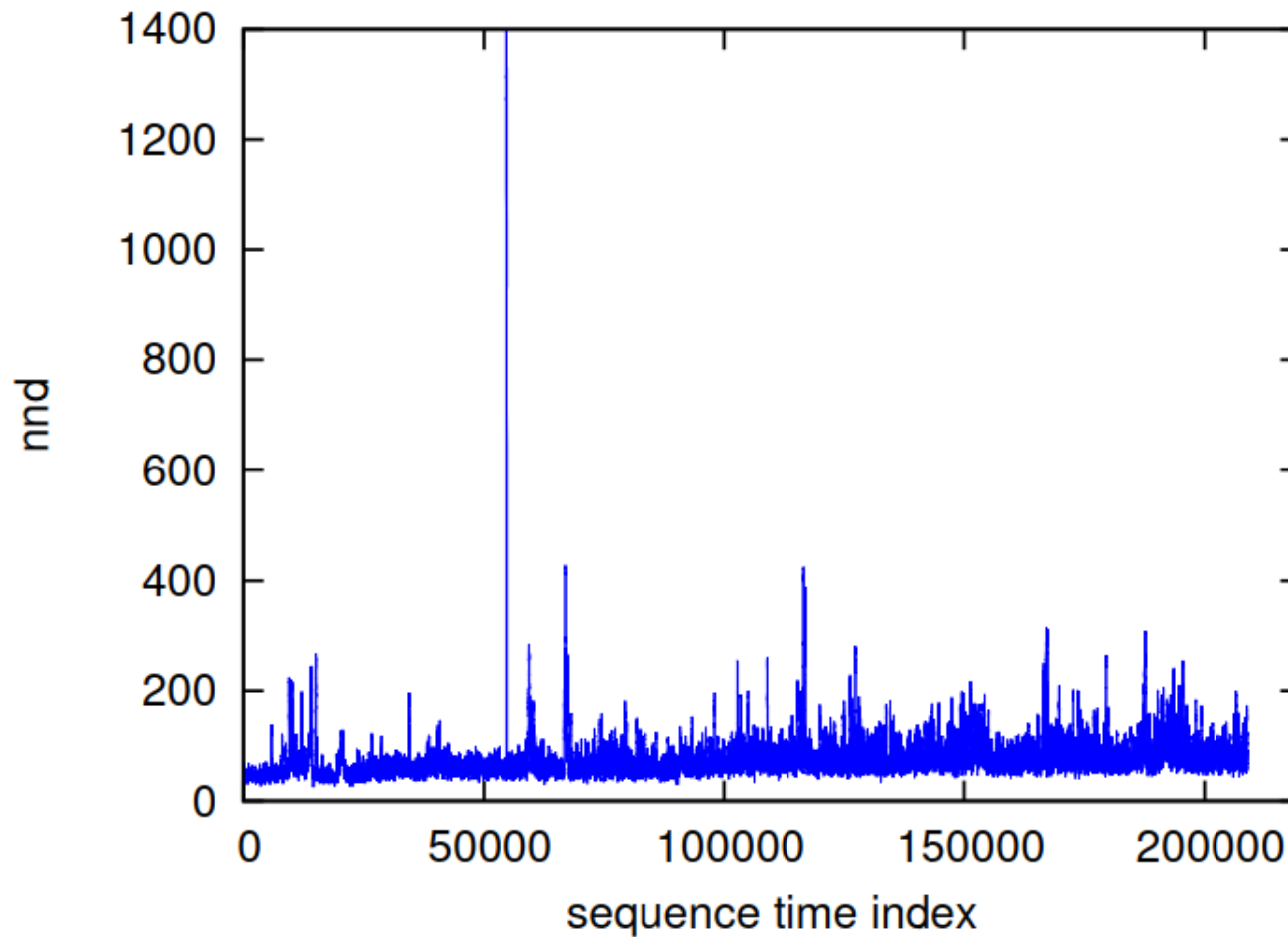
VALIDATION



2

Let's force a full run on the same cluster, but with alphabet +1

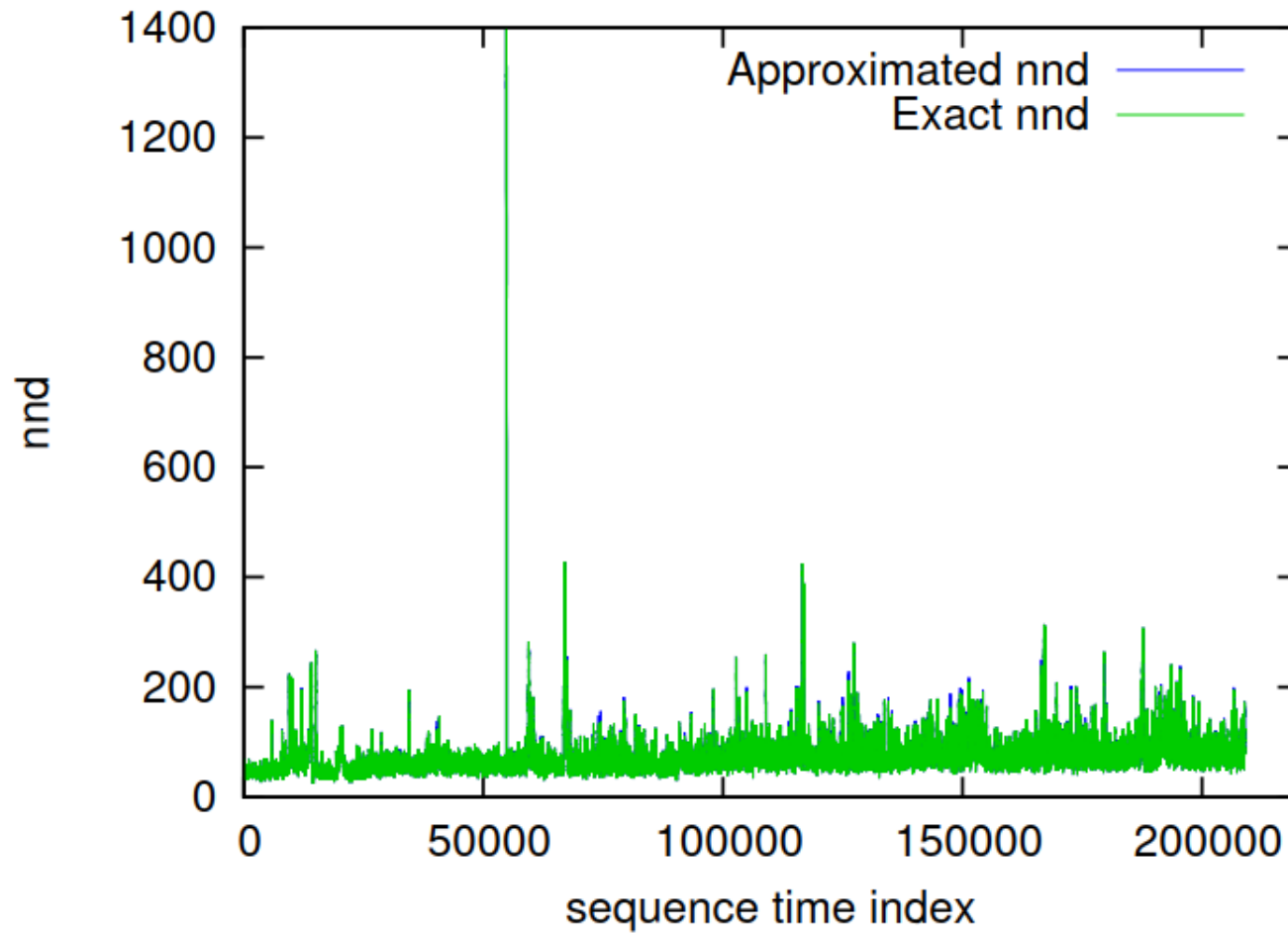
VALIDATION



3

Improve the *nnds* with the time topology

VALIDATION



VALIDATION

file name	% of exact <i>nnds</i>	Err	speedup
<i>sel0606</i> ₂	99.6	0.05	83
<i>sel0606</i> ₃	99.5	0.05	116
<i>sel102</i> ₂	98.8	0.05	40
<i>sel102</i> ₃	99.7	0.06	26
<i>sel123</i> ₂	98.8	0.05	32
<i>sel123</i> ₃	99.2	0.04	14
<i>bidmc15</i> ₂	99.6	0.15	5
<i>bidmc15</i> ₃	99.8	0.09	30
<i>bidmc15</i> ₄	98.4	0.05	40
<i>bidmc15</i> ₅	98.9	0.06	58

VALIDATION

Err = Average Relative Error

file name	% of exact <i>nnds</i>	Err	speedup
<i>sel0606</i> ₂	99.6	0.05	83
<i>sel0606</i> ₃	99.5	0.05	116
<i>sel102</i> ₂	98.8	0.05	40
<i>sel102</i> ₃	99.7	0.06	26
<i>sel123</i> ₂	98.8	0.05	32
<i>sel123</i> ₃	99.2	0.04	14
<i>bidmc15</i> ₂	99.6	0.15	5
<i>bidmc15</i> ₃	99.8	0.09	30
<i>bidmc15</i> ₄	98.4	0.05	40
<i>bidmc15</i> ₅	98.9	0.06	58

N_a = number of sequences for which only an approximate *nnd* has been found

$$Err = \frac{1}{N_a} \sum_{i=1}^N \frac{nnd_a(S^i) - nnd(S^i)}{nnd(S^i)}$$

approximate *nnd*

exact *nnd*

VALIDATION

file name	% of exact <i>nnds</i>	Err	speedup
<i>sel0606</i> ₂	99.6	0.05	83
<i>sel0606</i> ₃	99.5	0.05	116
<i>sel102</i> ₂	98.8	0.05	40
<i>sel102</i> ₃	99.7	0.06	26
<i>sel123</i> ₂	98.8	0.05	32
<i>sel123</i> ₃	99.2	0.04	14
<i>bidmc15</i> ₂	99.6	0.15	5
<i>bidmc15</i> ₃	99.8	0.09	30
<i>bidmc15</i> ₄	98.4	0.05	40
<i>bidmc15</i> ₅	98.9	0.06	58

Conclusions

1. By interweaving **Euclidean**, **SAX**, and **Time** topologies, (starting from a discord search algorithm) we obtain a new algorithm for finding an **approximate** *nnd* profile for all the sequences of a time series

Conclusions

1. By interweaving **Euclidean**, **SAX**, and **Time** topologies, (starting from a discord search algorithm) we obtain a new algorithm for finding an **approximate *nnd*** profile for all the sequences of a time series
2. The evaluation we performed returns more than 98% of the **exact *nnds***, and the values of the approximate *nnds* are within 15% in respect to the exact ones.

Conclusions

1. By interweaving **Euclidean**, **SAX**, and **Time** topologies, (starting from a discord search algorithm) we obtain a new algorithm for finding an **approximate *nnd*** profile for all the sequences of a time series
2. The evaluation we performed returns more than 98% of the **exact *nnds***, and the values of the approximate *nnds* are within 15% in respect to the exact ones.
3. The speedups are between **1 and 2 orders** of magnitude in respect to the brute force

Conclusions

1. By interweaving **Euclidean**, **SAX**, and **Time** topologies, (starting from a discord search algorithm) we obtain a new algorithm for finding an **approximate** *nnd* profile for all the sequences of a time series
2. The evaluation we performed returns more than 98% of the **exact nnds**, and the values of the approximate *nnds* are within 15% in respect to the exact ones.
3. The speedups are between **1 and 2 orders** of magnitude in respect to the brute force
4. The algorithm **assures** to find the sequences with the highest ***nnd*** (*discords*)

Conclusions

1. By interweaving **Euclidean**, **SAX**, and **Time** topologies, (starting from a discord search algorithm) we obtain a new algorithm for finding an **approximate** *nnd* profile for all the sequences of a time series
2. The evaluation we performed returns more than 98% of the **exact nnds**, and the values of the approximate *nnds* are within 15% in respect to the exact ones.
3. The speedups are between **1 and 2 orders** of magnitude in respect to the brute force
4. The algorithm **assures** to find the sequences with the highest ***nnd*** (*discords*)

Future

1. Use faster algorithms to calculate the distance between sequences
2. Add the second *nnd*, third *nnd*, etc.