



POLITECNICO DI TORINO

Telecommunications Engineering

Information Theory and Signal Processing Applications

Master Degree Dissertation

Nearest Neighbour Search using binary clustered Neural Networks

Applied to object retrieval and classification

Candidate

Demetrio Ferro

Advisor

Guido Montorsi

Internship Student at

Télécom Bretagne





Co-Advisor:

Claude Berrou

Internship Tutor:

Vincent Gripon

Télécom Bretagne

Brest, FRANCE

This work has been funded in part by the ERC under the EU Seventh Framework Programme (FP7 / 2007 - 2013) / Grant agreement n° 290901.

PRESENTATION OVERVIEW

1. Problem Statement

2. Vector Quantization

3. Neural Networks

4. Training Stage

5. Query Stage

6. Empirical Results

7. Conclusions

PROBLEM STATEMENT

NEAREST NEIGHBOUR SEARCH

The problem of searching for the Nearest Neighbour is formulated as:

“Given a collection of data points and a query point in a hyper-dimensional metric space, find the data point that is closest to the query one.”

[Kevin Beyer, 1998]

The closeness metric is generally considered as **Euclidean Distance**.

Motivations can be found in its wide set of applications, mainly related to the domains of Data Processing and Machine Learning:

→ **Object Retrieval** → **Classification**

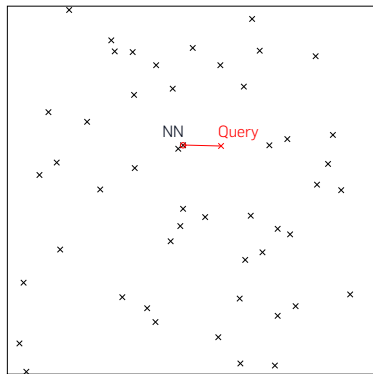
Pattern Recognition

Computer Vision

Computational Statistics

Data Mining.

PROBLEM STATEMENT



Example of Nearest Neighbour search query over a two-dimensional set of uniform distributed points.

Running an **Exhaustive Search**:

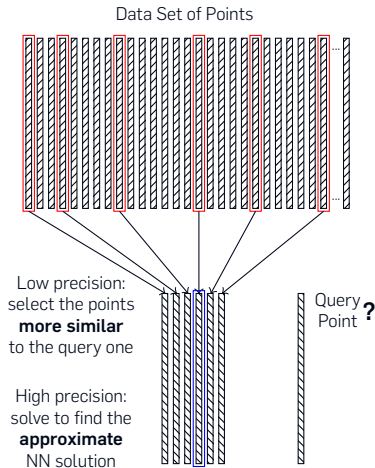
- Computing Euclidean Distances,
- Retrieving the minimum one.

The **Computational Complexity** is:

- Linear with the number of data points, considered to be N ,
- Linear with the dimensionality of the search space d .

It may hence be denoted as $\mathcal{O}(N \cdot d)$.

PROPOSED SOLUTION



The proposed solution is structured as:

- **Training Stage**

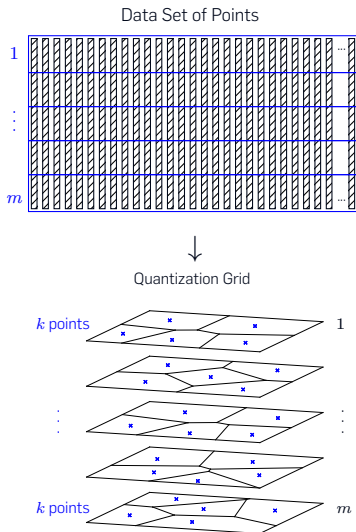
- Coarse Vector Quantization
acts on the data dimensionality,
- Neural Networks Learning
allows a quick data access.

- **Query Stage**

- Neural Networks Polling
acts on the data cardinality.
- Fine Vector Quantization
Despite a higher compression
rate, it works at low dispersion.

VECTOR QUANTIZATION

VECTOR QUANTIZATION



The quantization grid is computed empirically over the data set of points.

- **Splitting over m layers**

→ The search space is split in contiguous orthogonal sections.

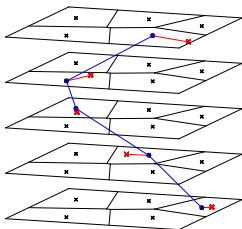
- **Clustering the layers to k points**

→ Each data points section is processed to determine the best set of **cluster centroids**.

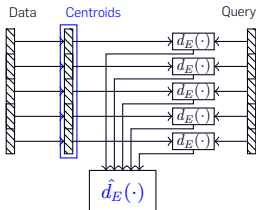
The optimality criteria considered is the minimum average dispersion one, Its achievement is sought in the use of a specific procedure known as **k -means**.

PRODUCT QUANTIZATION

Quantization Pattern



Asymmetric Distance Computation



Product Quantization (PQ) is a state-of-the-art technique to solve Approximate NN (ANN) search:

- **Data Points Approximation**

→ Associate them centroids sets called **Quantization Patterns**.

- **Distance Metric Approximation**

→ The Euclidean Distance is computed as **Sum of Distances** over lower-dimensional layers.

The scan search can be accelerated by using a Look Up Table approach.

NEURAL NETWORKS

ARTIFICIAL NEURAL NETWORKS

Artificial Neural Networks are a family of computational models able to accelerate the **data access** with a neuro-inspired approach.

Based on the biological model of **synaptic interaction**, it concerns a system of interconnected neural cells exchanging information messages.

The connections have numeric weights that can be tuned based on experience, making them capable of **learning** and **retrieving** messages.

They may be represented in two equivalent ways:

Weighted, Undirected Graphs



Nodes represent neurons, whereas edges are their interconnections.

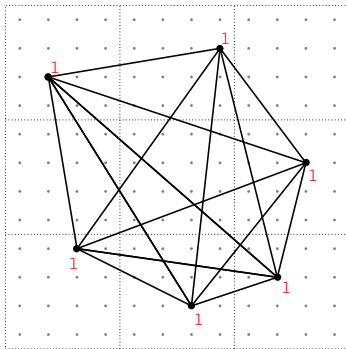
They always generate fully connected polyhedra.

Connections Matrices

0	1	0	0	0
0	0	0	0	0
1	0	1	1	1
1	0	0	0	0
1	1	0	0	0

Each of the elements expresses the activation of a cell, hence the weight of the connections between couple of nodes.

ADOPTED MODEL



Peculiar properties of the adopted model, can be expressed by:

→ **Binary Connections Weights**

Ones in the matrix establish active nodes, i.e. they identify only the existence of edges among nodes.

→ **Clustered neural activations**

The connections matrix is partitioned in such a way that only unique activations occur within each cluster.

The activation co-occurrences within clusters identify a fully connected sub-graph, or **Connections Pattern**.

FUNCTIONALITIES

Learning Rule

Considering \mathbf{z} to be one of the **binary messages** in the **learning set** \mathcal{Z} with a clustered structure, the **adjacency connections matrix** $W(\mathcal{Z})$ is generated as:

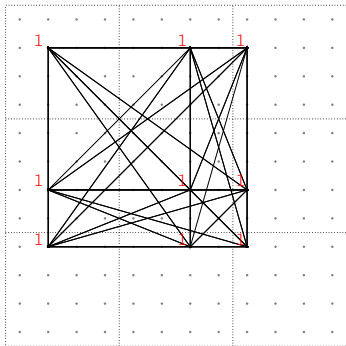
$$W(\mathcal{Z}) = g \left[\sum_{\mathbf{z} \in \mathcal{Z}} (\mathbf{z} \cdot \mathbf{z}^T) \right], \quad g(\xi) = \begin{cases} 1 & \text{if } \xi > 0 \\ 0 & \text{otherwise} \end{cases}, \quad \xi \in \mathbb{Z}^*$$

Polling Rule

Being \mathbf{z}_0 a **binary query message**, still matching the clustering structure, the polling rule for the neural network trained by the set \mathcal{Z} is expressed by the **score**:

$$s(\mathbf{z}_0, \mathcal{Z}) = \frac{\mathbf{z}_0^T \cdot W(\mathcal{Z}) \cdot \mathbf{z}_0}{(\mathbf{z}_0^T \cdot \mathbf{z}_0)^2}, \in [0, 1] \subset \mathbb{R}$$

MOTIVATIONS OF USE



The choice of a binary, clustered model is motivated by many reasons:

- Binary connection weights
Reduced **computational complexity** and preserve **memory occupancy**.
- Clustered nodes structure
The learning process acts as **patterns overlapping**.
Clusters prevent from generating **false connections patterns**.
- Adjacency connections matrix
For its **symmetry**, it carries a lot of **redundancy**, making it robust to **noise or memory faults**.

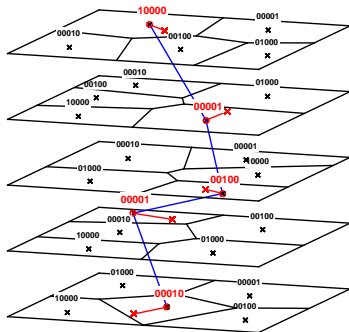
TRAINING STAGE

TRAINING STAGE OUTLINE

Set of operations executed offline, preparatory to the Query stage.

- Run a first **Coarse Vector Quantization**
 - Build the Quantization Grids over m_c **layers**, each of them composed of k_c **centroids**, obtained by running k -means;
 - Assign each data point to a **Quantization Pattern**.
- Go to the **Neural Networks Learning Stage**
 - Convert the quantized data points to **binary messages**, so that they can be learnt by binary neural networks.
 - Identify **L learning sets** of binary messages to be learnt.
 - applying the **learning rule** over all of the neural networks.

CLUSTERED BINARY LABELLING



$z = [10000 \ 00001 \ 00100 \ 00001 \ 00010]$

Quantization
Patterns



Connections
Patterns

Centroids
Layers



Networks
Clusters

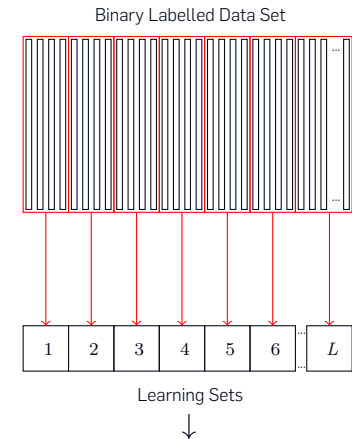
It is performed a **Binary Labelling**:

- All of the centroids, over each layer are associated to **unitary Hamming weight**, from the canonical basis of order k .
- For the network structure, only a **single activation** is allowed for each cluster.

Once all of the vectors are labelled, the whole search space is mapped to its binary version.

At this point, it is foreseen the use of a preprocessing step, aiming to partition the search space in a convenient way.

PREPROCESSING FOR ALLOCATION



$$W(\mathcal{Z}_l) = g \left[\sum_{\mathbf{z} \in \mathcal{Z}_l} (\mathbf{z} \cdot \mathbf{z}^T) \right], \quad \forall l \in \llbracket 1, L \rrbracket$$

The binary labelled data set of points is partitioned in **L disjoint learning sets**.

It is used a greedy approach, aiming to:

→ **Portion of ones in the matrices**

Minimize the mutual Hamming distance within each set, hence the diversity of messages.

→ **Constant diversity**

The diversity of messages has to be uniformly distributed among the L learning sets $\mathcal{Z}_l, \forall l \in \llbracket 1, L \rrbracket$.

→ **Constant cardinality**

The partition has to be fair, hence $|\mathcal{Z}_l| = n = N/L, \forall l \in \llbracket 1, L \rrbracket$, in such a way that cardinality reductions do not depend on the specific query.

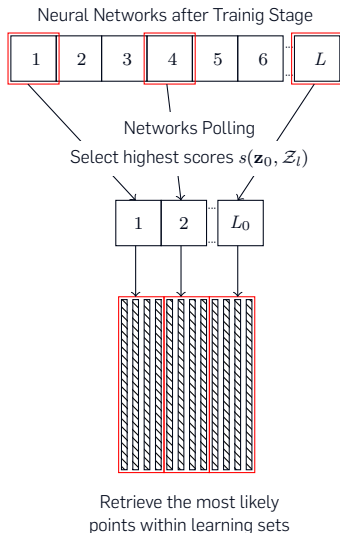
QUERY STAGE

QUERY STAGE OUTLINE

Set of operations executed on the run, when a query is performed.

- Resort to a **Neural Networks Polling**
 - Compute the likelihood scores for all of the L networks;
 - Select the **highest L_0 scores** of likelihood;
 - Retrieve the L_0 most likely **learning sets**.
- Apply a last **Fine Vector Quantization**
 - Associated data points to a **fine** grid composed of m_f **layers**, with k_f **centroids** each;
 - Approximate distance metrics with **reduced dispersion**.

NETWORKS POLLING



The selection of most likely sets is supposed to be performed as:

→ **Compute scores of likelihood**

For the query vector to be in the set \mathcal{Z}_l , learnt by the l^{th} neural network.

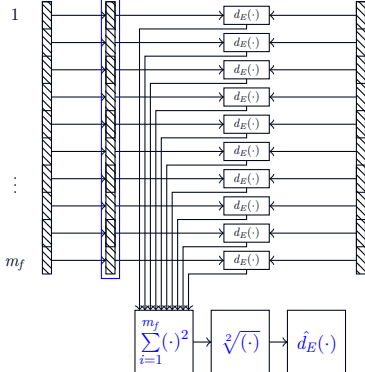
$$s(\mathbf{z}_0, \mathcal{Z}_l) = \frac{\mathbf{z}_0^T \cdot W(\mathcal{Z}_l) \cdot \mathbf{z}_0}{(\mathbf{z}_0^T \cdot \mathbf{z}_0)^2}, \forall l \in \llbracket 1, L \rrbracket$$

→ **Select the highest L_0 scores**

That is, the vectors within the most likely learning sets.

Once the data points are retrieved, it is applied a better approximation to find the Nearest Neighbour with more precision.

Data Point	Centroids	Query Point
1	2	3
4	5	6
7	8	9
10	11	12
13	14	15
16	17	18
19	20	21
22	23	24
25	26	27
28	29	30
31	32	33
34	35	36
37	38	39
40	41	42
43	44	45
46	47	48
49	50	51
52	53	54
55	56	57
58	59	60
61	62	63
64	65	66
67	68	69
70	71	72
73	74	75
76	77	78
79	80	81
82	83	84
85	86	87
88	89	90
91	92	93
94	95	96
97	98	99
100	101	102
103	104	105
106	107	108
109	110	111
112	113	114
115	116	117
118	119	120
121	122	123
124	125	126
127	128	129
130	131	132
133	134	135
136	137	138
139	140	141
142	143	144
145	146	147
148	149	150
151	152	153
154	155	156
157	158	159
160	161	162
163	164	165
166	167	168
169	170	171
172	173	174
175	176	177
178	179	180
181	182	183
184	185	186
187	188	189
190	191	192
193	194	195
196	197	198
199	200	201
202	203	204
205	206	207
208	209	210
211	212	213
214	215	216
217	218	219
220	221	222
223	224	225
226	227	228
229	230	231
232	233	234
235	236	237
238	239	240
241	242	243
244	245	246
247	248	249
250	251	252
253	254	255
256	257	258
259	260	261
262	263	264
265	266	267
268	269	270
271	272	273
274	275	276
277	278	279
280	281	282
283	284	285
286	287	288
289	290	291
292	293	294
295	296	297
298	299	300
301	302	303
304	305	306
307	308	309
310	311	312
313	314	315
316	317	318
319	320	321
322	323	324
325	326	327
328	329	330
331	332	333
334	335	336
337	338	339
340	341	342
343	344	345
346	347	348
349	350	351
352	353	354
355	356	357
358	359	360
361	362	363
364	365	366
367		



- **Quantizing the data set**
Each of the data points is associated to a fine quantization pattern.
- **Distance Approximation**
The distance among points is computed as the sum of distances over lower-dimensional sections.

It consists of iterative runs of the same technique, where solution obtained at each iteration is removed from the search space.

EMPIRICAL RESULTS

PERFORMANCES

→ Metric of Evaluation

The performances are measured in terms of success rate for **estimated solution to match the exhaustive one**, in a set of recalls of order r .

$$\eta_r = \begin{cases} 1 & \text{if } r\text{-th rank NN matches} \\ & \text{the exhaustive solution} \\ 0 & \text{otherwise} \end{cases}$$

→ Computational Complexity

Fine Product Quantization

$$\mathcal{O}(k_f d + m_f N)$$

Distance Computation
→ Fine Quantization

Binary Clustered Networks

$$\mathcal{O}(k_f d + k_c d + p_1 (m_c k_c)^2 L + m_f L_0 n)$$

Coarse, Fine Quantization ← Scores Computation
Distance Computation ←

TESTED APPLICATIONS

Object Retrieval

Pictures querying within a large data collection, assuming the shape of SIFTs, i.e. extracted features local descriptors.

TEXMEX Group - INRIA, Rennes.

Classification

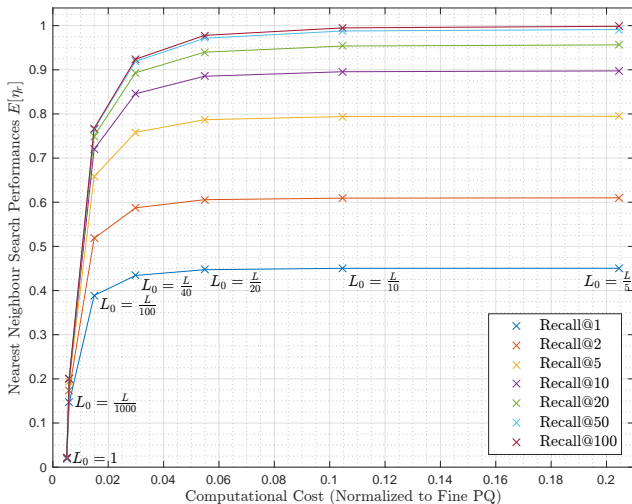
Identify labels associated to a large set of pictures containing handwritten digits, ranging from 0 to 9.

MNIST Database - Yann LeCun

In both cases, **all of the parameters of interest are trained** to work at the best settings compromise.

The achieved results are shown in the next couple of slides.

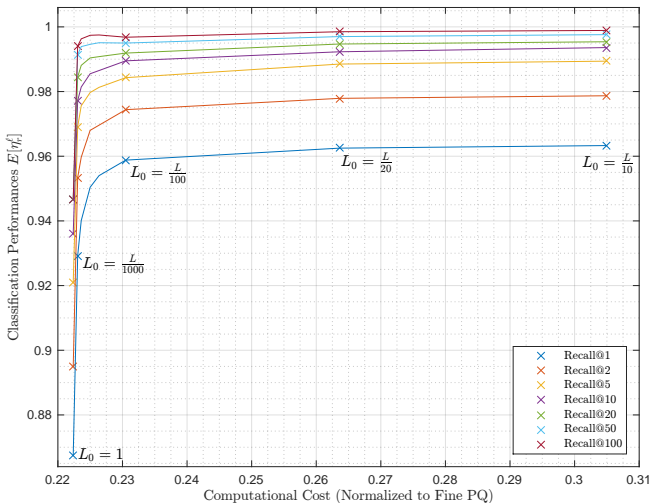
PERFORMANCES: OBJECT RETRIEVAL



Parameters list:

- Dimensionality
 $d = 128$
- Cardinality
 $N = 1000000$
- Coarse PQ
 $k_c = 2, m_c = 32$
- Fine PQ
 $k_f = 16, m_f = 256$
- Number of networks
 $L = 10000$

PERFORMANCES: CLASSIFICATION



Parameters list:

- Dimensionality
 $d = 784$
- Cardinality
 $N = 60000$
- Coarse PQ
 $k_c = 2, m_c = 32$
- Fine PQ
 $k_f = 16, m_f = 256$
- Number of
networks $L = 2000$

CONCLUSIONS

CONCLUSIONS

This work illustrates the interest of using the introduced neural networks model to accelerate ANN search over real applications.

Future works may include different new proposals, such as:

- Getting rid of the **Fine search** stage by implementing message retrieval strategies based with a Bayesian approach;
- Trying to use a **Hierarchical approach**, consisting in cascading different granularity layers;
- Proposing **refined allocation strategies** to choose which vector should be stored in which networks;
e.g. Group testing, hence the selection of non-disjoint sets

Thank you for your kind attention.

BIBLIOGRAPHY



Demetrio Ferro, Vincent Gripon, and Xiaoran Jiang.

Nearest neighbour search using binary neuroinspired associative memories.
In Neural Information Processing Systems, November 2015.
Submitted.



Demetrio Ferro, Vincent Gripon, and Xiaoran Jiang.

Nearest neighbour search using binary neuroinspired associative memories.
In Computer Vision and Pattern Recognition, November 2015.
Submitted.



Hervé Jégou, Matthijs Douze, and Cordelia Schmid.

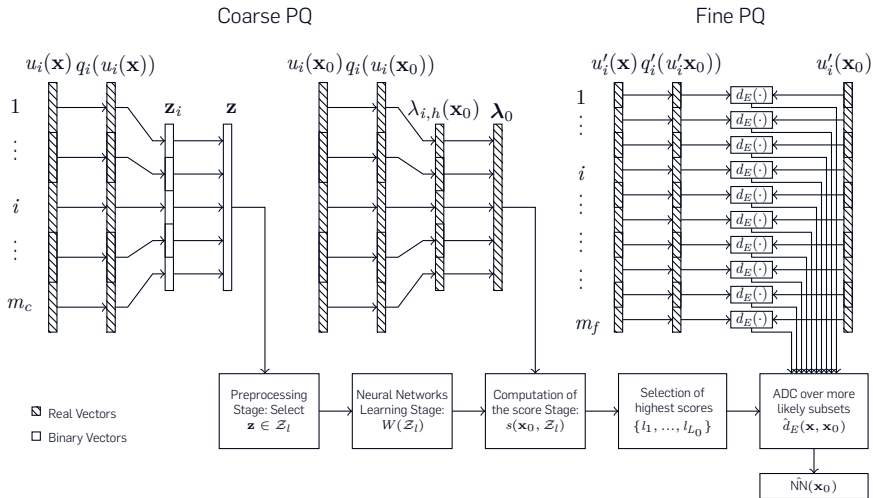
Product Quantization for Nearest Neighbour Search.
IEEE Transac. on Pattern Analysis and Machine Intelligence, January 2011.



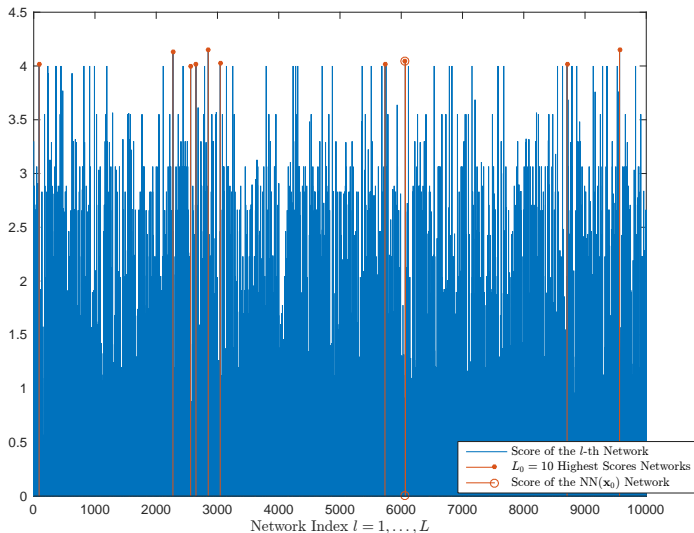
Chendi Yu, Vincent Gripon, Xiaoran Jiang, and Hervé Jégou.

Neural associative memories as accelerators for binary vector search.
In Proceedings of Cognitive, March 2015.
To appear.

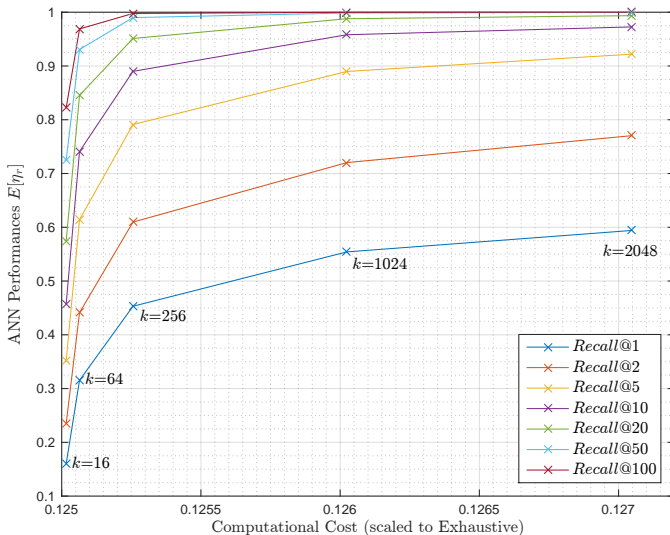
BACKUP - OVERALL TECHNIQUE SCHEME



BACKUP - RELIABILITY OF THE SCORES



BACKUP - SELECTING PARAMETERS: PQ PERFORMANCES



BACKUP - SELECTING PARAMETERS: CAOARSE GRANULARITY

