

## 0.1 OpenLearn documentazione

Questa parte della relazione è dedicata all'illustrazione e al commento della struttura e delle parti principali del componente principale del progetto, la libreria di machine learning sviluppata per la creazione, l'addestramento e l'utilizzo dei modelli neurali. La libreria è stata sviluppata per un utilizzo general-purpose, quindi non finalizzata alla risoluzione dello specifico problema proposto, il quale è stato utilizzato per testare la risposta di varie strutture neurali alla risoluzione di un problema con soluzione nota ma attraverso l'elaborazione di dati di ingresso fortemente affetti da rumore.

Nella documentazione verranno mostrate soltanto le caratteristiche principali della libreria, le classi, le strutture e i metodi più importanti effettivamente utilizzati benché la libreria ne presenti altri predisposti per future revisioni ed espansioni.

### 0.1.1 Le librerie importate

Le prime librerie importate sono le librerie per l'utilizzo del linguaggio C/CUDA, proprietario di Nvidia, per la programmazione delle schede della stessa azienda.

```
1 //////////////////////////////////////////////////CUDA INCLUDES////////////////////////////////////
2 #include "cuda.h"
3 #include "cuda_runtime.h"
4 #include <device_functions.h>
5 #include "device_launch_parameters.h"
```

Listing 1: librerie cuda

Sono state importate poi diverse librerie standard del linguaggio C++, principalmente:

la classe vector è stata usata per l'incapsulamento e la manipolazione dei dati.

le classi \*stream sono state utilizzate per il salvataggio e il caricamento dei modelli su file di testo.

la libreria windows per l'utilizzo delle funzioni di timing ad alta precisione.

```
1
2 #include "stdafx.h"
3 #include <iostream>
4 #include <stdio.h>
5 #include <fstream>
6 #include <sstream>
7 #include <string>
8 #include <windows.h>
9 #include <limits>
10 #include <math.h>
11 #include <vector>
12 #include <list>
13 #include <time.h>
14 #include <algorithm>
15 #include <assert.h>
```

Listing 2: librerie STD

### 0.1.2 classi di creazione e addestramento modelli su CPU

Di seguito sono illustrate le strutture che formano lo scheletro dei modelli, tali strutture sono state pensate per essere completamente modulari e per permettere la realizzazione di diversi modelli, dalle più semplici feed-forward, alle più complesse ricorrenti. Tale approccio permette la massima maneggevolezza della struttura sacrificando l'efficienza che deriva da una computazione matriciale.

```
1 //prototipo struttura neuron
2 struct Neuron;
3 //dichiarazione puntatore a neuron
4 typedef Neuron *ptNeuron;
5
6 //La struttura arc, presente in array nella struttura Neuron
7 //conserva i dati delle connessioni e il puntatore alla struttura Neuron puntata
8 struct arc {
9     ptNeuron target = nullptr;
10     float weight = 0;
11     float oldDelta = 0;
12     //bool enabled = true;
13 };
14
15 //La struttura interArc è stata creata per l'addestramento MLP supportato da un rete ricorrente
16 //fornisce un arco di interconnessione tra i neuroni di 2 reti
```

```

17 struct interArc {
18     ptNeuron target = nullptr;
19     ptNeuron base = nullptr;
20 };
21
22 //La struttura Neuron è l'unità base delle strutture delle reti e ne conserva tutti i parametri
23 struct Neuron {
24     //OutArcs è il vettore di connessioni ai neuroni degli strati successivi
25     // to initialize: = new vector<arc>(10);
26     vector<arc> OutArcs;
27     // numero di archi in uscita
28     u_int numOutArcs = 0;
29     //numero di archi in ingresso
30     u_int numInArcs = 0;
31     //indice riga del neurone
32     u_int layer = 0;
33     //indice della colonna del neurone
34     u_int column = 0;
35     //peso del bayes
36     float bayes = 0.01f;
37     //ultima variazione del peso
38     float oldBayesDelta = 0;
39     // vettore delle interconnessioni temporali
40     //vector<float> timeBayes;
41     // vettore contenente la percentuale di influenza relativa ad ogni input
42     vector<float> influenceInput;
43     // vettore contenente la percentuale di influenza relativa all'errore retropropagato da ogni
44     // ↳ output
45     vector<float> influenceOutput;
46     //potenziale attuale del neurone
47     float output = 0;
48     //somma in valore assoluto degli input del neurone
49     float absOutSum = 0;
50     //somma in valore assoluto delle variazioni dei pesi
51     float absDeltaSum = 0;
52     //errore di retropropagazione
53     float BPerr = 0;
54     //ogni neurone è contraddistinto da un indice unico che si riferisce alla sua posizione
55     int neurIdx = 0;
56 };
57
58 //La struct Layer è stata realizzata per contenere un layer di neuroni
59 struct Layer {
60     vector<Neuron> Neurons;
61     u_int numNeurons = 0;
62 };
63
64 //Omap conserva la scala di conversione tra l'output
65 struct Omap {
66     float maxValue = 1;
67     float minValue = 0;
68 };
69
70 // struttura necessaria per l'inserimento di un nuovo neurone
71 struct conMap {
72     u_int startLyr;
73     u_int startCol;
74     u_int arcRef;
75     u_int targetLyr;
76     u_int targetCol;
77 };
78
79 //conserva una serie temporale prima dell'accumulazione in una struct Dataset
80 struct timeSeries {
81     list<float> evento;
82 };
83
84 //struttura contenente il singolo esempio
85 struct example {

```

```

85     vector<float> input;
86     vector<float> Doutput;
87 };
88
89 //Struttura contenente l'intero Dataset di una rete
90 struct Dataset {
91     vector<example> trainingSet; // vettore esempi del training set
92     vector<example> validationSet; // vettore esempi del validation set
93     float triningErr = 0; //percentuale di errore training set
94     float validationErr = 0; //percentuale di errore validation set
95 };

```

Listing 3: struct strutturali del modello

La classe DatasetCore è stata realizzata per permettere l'estrazione dei dataset da vari formati, riportandoli nello standard di utilizzato dalle classi che li utilizzano per effettuare gli addestramenti.

```

1 class DatasetCore {
2 public:
3     // lista dei dataset
4     list<Dataset> Datasets;
5
6     // costruttore DatasetCore
7     DatasetCore() { ... }
8
9     //////////////////////////////////MANIPOLAZIONE DATASET////////////////////////////////////
10    //funzione che consente la lettura di serie temporali da csv e le ristruttura in un dataset
11    void readTimeSeriesCsv(string filename, int outStep, int inEx, float trainingPerc) { ... }
12
13    //estrae un dataset dalla lista interna della classe, la parte training se il flag e true
14    // o la parte validation se il flag è false
15    vector<example> getDataset(int n, bool training = true) { ... }
16    //////////////////////////////////////
17
18    //////////////////////////////////ALTRE FUNZIONI////////////////////////////////////
19    //metodo interno per contare le righe di un file
20    int countrRow(string filename) { ... }

```

Listing 4: class DatasetCore

La classe network fornisce le variabili e i metodi generali che la maggior parte degli oggetti rete utilizza.

```

1 class Network {
2
3 public:
4     // vettore di struct layer
5     vector<Layer> Layers;
6     // vettore di esempi per l'apprendimento
7     vector<example> examples;
8     // vettore contenente i valori di rimappatura dell'output della rete
9     vector<Omap> map;
10
11    // nome del file contenente la struttura della rete
12    string genoma = "";
13    // Layer nella rete compresi strato input e strato output
14    u_int nLayers = 0;
15    // numero totale neuroni nella rete
16    u_int nNeurons = 0;
17    //numero totale di tutti gli archi presenti nella rete
18    u_int nArc = 0;
19
20    //costruttore
21    Network(string filename) {
22        genoma = filename;
23    }
24
25    //////////////////////////////////FUNZIONE COSTRUZIONE RETE DA FILE////////////////////////////////////
26    //dato il nome del file contenete la struttura carica l'oggetto corrispondente
27    void getNetParams() { ... }
28    //////////////////////////////////////
29

```

```

30 //////////////////////////////////////////////////FUNZIONE COSTRUZIONE DATASET DA FILE////////////////////////////////////
31 //dato il file contenente un dataset prestrutturato correttamente lo carica nel modello
32 void getDataset(string filename) { ... }
33 //////////////////////////////////////////////////
34
35 //////////////////////////////////////////////////SALVA RETE SU FILE////////////////////////////////////
36 //salva la rete sul file specificato
37 void saveNet(string filename = "") { ... }
38 //////////////////////////////////////////////////
39
40 //////////////////////////////////////////////////SALVA DATASET SU FILE////////////////////////////////////
41 //salva il dataset sul file specificato
42 void saveDataset(string filename) { ... }
43 //////////////////////////////////////////////////
44
45 ////////////////////////////////////////////////// FUNZIONI MODIFICA RETE////////////////////////////////////
46 //elimina un arco dati i parametri diriferimento
47 void deleteArc(int Nlayer, int Ncolumn, int targetLayer, int targetColumn) { ... }
48
49 //elimina un neurone, compresi tutti gli archi ad esso connesso
50 void deleteNeuron(int Nlayer, int Ncolumn) { ... }
51
52 //aggiunge un arco dati i riferimenti al neurone di partenza e di arrivo
53 void addArc(int Nlayer, int Ncolumn, int targetLayer, int targetColumn) { ... }
54 //////////////////////////////////////////////////
55
56 //////////////////////////////////////////////////FUNZIONI VARIE////////////////////////////////////
57 float DsigOut(int Layer, int Neuron) { ... }
58 float sigmoid(float x) { return 1 / (1 + exp(-x)); } // Sigmoide
59 float logit(float x) { return log(x / (1 - x)); } // funzione sigmoide inversa
60 float gaussian(float x, float mu, float var) { ... }
61 void WeightsStamp(string mode = "a") { ... }
62     //m - stampa le medie dei pesi di ogni layer
63     //a - stampa tutti i pesi della rete con alcuni parametri di apprendimento
64     //w - stampa tutti i pesi con il riferimento riga colonna al target
65     //fc - stampa le medie dei gruppi di pesi tra due layer
66
67 //funzione sigmoide
68 void sigLayer(int lyr) { ... }
69
70 //applica il bayes all'output di ogni neurone del dato layer
71 void bayesLayer(int lyr, bool absSum = false) { ... }
72
73 // esegue il reset del potenziale di tutti i neuroni della rete
74 void resetPotential() { ... }
75
76 // esegue il reset della sommatoria di ogni input in valore assoluto di ogni neurone
77 void resetAbsSumPotenzial() { ... }
78
79 // esegue il reset della sommatoria di ogni input in valore assoluto di ogni neurone
80 void resetAbsSumDelta() { ... }
81
82 // esegue il reset dell'errore retropropagato in ogni neurone
83 void resetBPerr() { ... }
84
85 //resetta il valore ID nei neuroni di una rete (necessario per modifiche strutturali)
86 void resetNeuronsID() { ... }
87
88 // salva su vettore i riferimenti numerici delle connessioni verso il layer specificato
89 vector<conMap> saveConsTowardsLyr(int Layer) { ... }
90
91 // ricarica i riferimenti numerici delle connessioni verso il layer specificato
92 void loadConsTowardsLyr(vector<conMap> con) { ... }
93
94 //elimina il contenuto l'oggetto dataset dell'oggetto rete
95 void ClearDataset() { ... }
96
97 //generazione di una serie storica del seno (DEBUG)
98 void genTestDataset(int nExe, int nIn, int nOut, float step, int type, float offset) { ... }

```

```

99
100 //esegue il settaggio della mappatura di output dell'oggetto
101 void setNetMap(float max, float min) { ... }
102
103 //restituisce il valore normalizzato dell'output della rete
104 float reverseMap(int neur, float val) { ... }
105
106 // crea un vettore di n elementi successivi e li disordina
107 // creazione di una tabella di accesso casuale per un secondo vettore
108 vector<u_int> casualVector(int in, int start = 0) { ... }
109
110 //esegue il mescolamento degli elementi all'interno di un oggetto vector
111 /*vector<T, A>*/
112 template<typename T, typename A>
113 void shackleVector(vector<T, A> const& vec) { ... }
114
115 //ricalcola gli ID dei neuroni dopo il reset (da eseguire dopo delle modifiche strutturali)
116 void refreshNeurIdx() { ... }
117
118 //applica all'intero dataset un valore di offset
119 void datasetOffset(float offset) { ... }
120 ///////////////////////////////////////////////////
121
122 ///////////////////////////////////////////////////ACCESSO A VARIABILI PRIVATE//////////////////////////////////////
123 int numLayers() { ... }
124 int numNeurons(int Layer) { ... }
125 int numCon(int Layer, int Neuron) { ... }
126 int numInCon(int Layer, int Neuron) { ... }
127 int getConTargetLyr(int Layer, int Neuron, int Arc) { ... }
128 int getConTargetCol(int Layer, int Neuron, int Arc) { ... }
129 float getWeight(int Layer, int Neuron, int Arc) { ... }
130 float getDeltaWeight(int Layer, int Neuron, int Arc) { ... }
131 float getOutput(int Layer, int Neuron) { ... }
132 float getBPerr(int Layer, int Neuron) { ... }
133 ptNeuron getTarget(int Layer, int Neuron) { ... }
134 ptNeuron getConTarget(int Layer, int Neuron, int Conn) { ... }
135 int getOutConTargetID(ptNeuron base, ptNeuron target) { ... }
136 ///////////////////////////////////////////////////
137
138 ///////////////////////////////////////////////////WINDOWS HIGH PRECISION TIMING//////////////////////////////////////
139 BOOL WINAPI QueryPerformanceCounter(_Out_ LARGE_INTEGER *lpPerformanceCount);
140 BOOL WINAPI QueryPerformanceFrequency(_Out_ LARGE_INTEGER *lpFrequency);
141 inline long long PerformanceCounter() noexcept
142 {...}
143 inline long long PerformanceFrequency() noexcept
144 {...}
145 ///////////////////////////////////////////////////
146 };

```

Listing 5: class Network

La classe MLP (multi-layer perceptron) figlia della classe Network completa il corredo della madre con tutta una serie di strumenti e variabili che supporta l'addestramento e l'utilizzo di questa tipologia di reti.

```

1
2 class MLP : public Network {
3
4 public:
5     // tempo di esecuzione medio in millisecondi
6     float NetPerformance = 0;
7
8     //errore percentuale medio associato alla rete
9     float NetErrPercent = 0;
10
11     //costruttore
12     MLP(string filename) :Network(filename) {};
13
14     ///////////////////////////////////////////////////FUNZIONI CREAZIONE RETE//////////////////////////////////////
15     //permette la costruzione rapida di una rete MLP quadrata con dimensioni date

```

```

16 void qubeNet(int Nlayers, int Ncolumns, int input, int output, bool c, float initValue = 0.01f
    ↪ ) { ... }
17
18 //CREAZIONE RETE QUADRATA COMPLETAMENTE CONNESSA
19 //permette la costruzione di una rete MLP quadrata con dimensioni date in cui ogni neurone
20 //è connesso a tutti i neuroni di tutti gli strati successivi.
21 //(a parità di neuroni tale modello è più pesante ma tipicamente ha un apprendimento più
    ↪ rapido)
22 void qubeNetFC(int Nlayers, int Ncolumns, int input, int output, bool c, float initValue = 0.
    ↪ 01f) { ... }
23
24 //CREAZIONE RETE CUSTOM
25 //permette la creazione di una rete MLP con la larghezza di ogni layer variabile
26 //indicando tali larghezze con un apposito vettore dato come argomento
27 void customNet(int Nlayers, vector<int> Ncolumns, float conFill) { ... }
28 ///////////////////////////////////////////////////
29
30 ///////////////////////////////////////////////////STIMOLAZIONE RETE//////////////////////////////////////
31 //procedura di propagazione dell'informazione
32 //dati i riferimenti al vettore contenete gli input e ad un vettore di output esegue
33 //il processamento degli input e carica l'output derivante sul vettore indicato
34 void inputNet(vector<float> &input, vector<float> &output) { ... }
35
36 //esegue una propagazione dell'informazione salvando lo storico di propagazione degli input
37 //tale funzione viene utilizzata per l'apprendiment strutturale
38 void inputNetProfiler(vector<float> &input, vector<float> &output) { ... }
39 ///////////////////////////////////////////////////
40
41 ///////////////////////////////////////////////////FUNZIONI DI ADDESTRAMENTO MLP//////////////////////////////////////
42 //Algoritmo di addestramento Back-propagation su tutto il dataset caricato
43 void BP(int iter, float eps, float beta, float ErrPercent) { ... }
44
45 //esecuzione del backpropagation per un solo esempio
46 //specificando il particolare esempio
47 void oneBP(float eps, float beta, example e) { ... }
48 ///////////////////////////////////////////////////
49
50 ///////////////////////////////////////////////////ALTRE FUNZIONI MLP//////////////////////////////////////
51 //inizializza i vettori di benchmark presenti nei neuroni per l'apprendimento strutturale
52 void initVectorsProfiler() { ... }
53
54 //resetta a zero tutti i vettori di profilazione esclusi layer input e output
55 void resetVectorsProfiler(bool inInfl, bool outInfl) { ... }
56
57 //stampa a schermo l'influenza degli input per ogni uscita
58 void stampInputInfluences(bool all = false) { ... }
59
60 //stampa a schermo l'influenza degli errori degli output per ogni ingresso
61 void stampOutputErrorPropagation(bool all = false) { ... }
62
63 //dal riferimento al nodo target e dall'indice del vettore restituisce il puntatore al neurone
    ↪ base
64 ptNeuron basePosfromtarget(ptNeuron target, int k) { ... }
65
66 //dato un neurone e l'indice di un suo arco restituisce l'indice di quella connessione all'
    ↪ interno del //vettore influenceInput all'interno del neurone target
67 int idBaseConReftoTargetInfl(ptNeuron base, int arc) { ... }
68 ///////////////////////////////////////////////////
69
70 };

```

Listing 6: class MLP

### 0.1.3 altre classi non utilizzate nel progetto

Le seguenti classi sono state iniziate e verranno completate in futuro, La classe Hopfield supporta (ancora in parte) la creazione l'addestramento e l'utilizzo delle reti ricorrenti, ovvero reti con connessioni tra i neuroni che si retropropagano, ovvero che portano potenziali di uscita di neuroni nello strato output e hidden al tempo (t), ad altri neuroni nel tempo (t+1).

La classe structural learning è stata realizzata e testata, ma presenta notevoli inefficienze che non permettono di utilizzarla in questa forma su modelli di grandi dimensioni, tale classe comunque è stata realizzata per effettuare un addestramento che mira a modificare la struttura della rete durante l'addestramento eliminando le connessioni che risultano essere mediamente quelle che persistono nell'introduzione di errore rispetto alle altre.

```
1 class Hopfield : public Network { ... };
2
3 class StructuralLearning { ... };
```

Listing 7: altre classi

### 0.1.4 Classi e kernel per l'interfacciamento con la GPU

In questa sezione sono stati implementati i kernel della GPU, ovvero le funzioni scritte in C/CUDA, che vengono eseguite sulla scheda grafica. tali funzioni sono contraddistinte dalla keyword `__global__`, e richiedono come argomenti gli indirizzi di memoria corrispondenti ai dati precaricati sulla global memory della GPU. Il richiamo dei kernel dal codice C++ richiedono la sintassi speciale:

**CUDAApplyWeightCorrections** «<numOfBlocksA, ThxBlock>> (...);

si vede infatti la sequenza di caratteri «<A,B>>» dove A è il numero di blocchi su cui deve essere eseguita la funzione mentre B rappresenta il numero di thread per ogni blocco.

```
1 //////////////////////////////////////////////////CUDA Kernels////////////////////////////////////
2
3 //resetta il valore di una variabile all'interno della scheda grafica
4 __global__ void CUDAResetVar(float *val) {
5     *val = 0;
6 }
7
8 //applica ad ogni arco della rete la correzione del peso
9 __global__ void CUDAApplyWeightCorrections(float eps, float *NeuronOut, float *BPerr, float *
    ↪ weights, int *ArcIn, int *ArcOut, int nArcs) {
10     unsigned int i = (blockIdx.x * blockDim.x) + threadIdx.x;
11     if (i < nArcs) {
12         weights[i] += -eps * BPerr[ArcIn[i]] * NeuronOut[ArcOut[i]];
13     }
14 }
15
16 //applica le correzioni bayesiane dei neuroni
17 __global__ void CUDAApplyBayesCorrections(float eps, float *BPerr, float *Bayes, int startN, int
    ↪ endN) {
18     unsigned int i = startN + (blockIdx.x * blockDim.x) + threadIdx.x;
19     if (i <= endN) {
20         Bayes[i] += -eps * BPerr[i];
21     }
22 }
23
24 //retropropaga l'errore nella rete
25 __global__ void CUDAPropagationErr(float *BPerr, float *weights, float *NeuronOut, int *ArcIn,
    ↪ int *ArcOut, int startA, int endA) {
26     unsigned int i = startA + (blockIdx.x * blockDim.x) + threadIdx.x;
27
28     //retropropaga l'errore dai neuroni successivi
29     if (i <= endA) {
30         //BPerr[ArcOut[i]] += BPerr[ArcIn[i]] * weights[i];
31         atomicAdd(&BPerr[ArcOut[i]], BPerr[ArcIn[i]] * weights[i]);
32     }
33 }
34
35 //moltiplica l'errore delle uscite per la derivata puntuale della sigmoide
36 __global__ void CUDAAoutDiff(float *BPerr, float *NeuronOut, int startN, int endN) {
37     int i = startN + (blockIdx.x * blockDim.x) + threadIdx.x;
38     if (i <= endN) {
39         BPerr[i] *= NeuronOut[i] * (1 - NeuronOut[i]);
40     }
41 }
```

```

42
43 //calcola l'errore dei neuroni dello strato output
44 __global__ void CUDAoutputErr(float *NeuronOut, int OutputRef, int numNeurons, int inputN, float
    ↳ *BPerr, float *examples, int exampleRef, float *mapMaxOut, float *mapMinOut, float *
    ↳ MeanErr) {
45
46     unsigned int i = (OutputRef) + (blockIdx.x * blockDim.x) + threadIdx.x; //indice di
    ↳ scorrimento vettori: NeuronOut, BPerr,
47     unsigned int e = (exampleRef + inputN) + (blockIdx.x * blockDim.x) + threadIdx.x; //indice di
    ↳ scorrimento vettori: examples
48     unsigned int m = (blockIdx.x * blockDim.x) + threadIdx.x; // indice di scorrimento vettori:
    ↳ mapMaxOut, mapMinOut
49     //if (i == 0) *MeanErr = 0;
50     if (i < numNeurons) {
51
52         float delta = mapMaxOut[m] - mapMinOut[m];
53         BPerr[i] = (NeuronOut[i] - ((examples[e] - mapMinOut[m]) / delta)) * NeuronOut[i] * (1 -
    ↳ NeuronOut[i]); // formula valida solo per i neuroni di uscita
54         //atomicAdd(MeanErr, (abs(((NeuronOut[i] * delta) + mapMinOut[m]) - examples[e]) / examples
    ↳ [e]))*100.0f);
55         atomicAdd(MeanErr, abs(((NeuronOut[i] * delta) + mapMinOut[m]) - examples[e]) / examples[e])
    ↳ * 100.0f);
56         // calcolo l'errore percentuale sulla singola uscita e lo sommo
57         //questo rappresenta uno dei punti più inefficienti dell'algoritmo,
58         //la funzione AtomicAdd infatti blocca la parallizzazione dell'algoritmo per tutti
59         //i thread che cercheranno di accedere a tale variabile riportando solo per questi
60         //l'esecuzione ad un modello sequenziale mettendoli in attesa
61     }
62 }
63
64 //resetta un dato vettore
65 __global__ void CUDAresetVector(float *vect, int size) {
66     unsigned int i = (blockIdx.x * blockDim.x) + threadIdx.x;
67     if (i < size) vect[i] = 0.0f;
68 }
69
70 //imposta i valori di output dei neuroni di input al valore dell'esempio del dataset
71 __global__ void CUDAs setInput(float *NeuronOut, int inputN, int exampleRef, float *example) {
72     unsigned int i = (blockIdx.x * blockDim.x) + threadIdx.x;
73     if (i < inputN) NeuronOut[i] = example[exampleRef + i];
74 }
75
76 //imposta i valori di output dei neuroni di input al valore specificato
77 __global__ void CUDAs setSingleInput(float *NeuronOut, int inputN, float *example) {
78     unsigned int i = (blockIdx.x * blockDim.x) + threadIdx.x;
79     if (i < inputN) NeuronOut[i] = example[i];
80 }
81
82 //applica la sigmoide ai potenziali dei neuroni in un dato intervallo (uno strato)
83 __global__ void CUDAsigLayer(float *NeuronOut, int start, int end) {
84     unsigned int i = start + (blockIdx.x * blockDim.x) + threadIdx.x;
85     if (i <= end) {
86         NeuronOut[i] = 1 / (1 + expf(-NeuronOut[i]));
87     }
88 }
89 //aggiunge all'output del neurone il contributo del bayes
90 __global__ void CUDAbayesInput(float *NeuronOut, float *Bayes, int start, int end) {
91     unsigned int i = start + (blockIdx.x * blockDim.x) + threadIdx.x;
92     if (i <= end) {
93         NeuronOut[i] += Bayes[i];
94     }
95 }
96
97 //propaga l'informazione dai neuroni dello strato input a quello di output
98 __global__ void CUDAlayerInput(float *weights, int *ArcIn, int *ArcOut, float *NeuronOut, int
    ↳ start, int end) {
99     unsigned int i = start + (blockIdx.x * blockDim.x) + threadIdx.x;
100     if (i <= end) {
101         atomicAdd(&NeuronOut[ArcOut[i]], NeuronOut[ArcIn[i]] * weights[i]);

```



```

102 //addizione bloccante non permette ad altri thread di sovrascrivere il valore
103 // finche l'operazione non è completata
104 }
105 }
106
107 ///////////////////////////////////////////////////

```

Listing 8: cuda kernels

Questa classe è l'interfaccia di collegamento con la GPU, la quale conserva i metodi di caricamento dei dati strutturali della rete, e gli esempi di addestramento sulla GPU.

```

1 //api di interfacciamento alla GPU
2 class CUDACore {
3 public:
4 //Device specs struct (contiene le specifiche della GPU)
5 cudaDeviceProp prop;
6 int GpuID = 0;
7
8 //struttura contenente i puntatori alle aree di memoria contenenti i parametri della rete nella
9   ↪ GPU
10 struct devNetParams {
11     float *weights = 0;
12     int *ArcIn = 0;
13     int *ArcOut = 0;
14     float *examples = 0;
15     float *NeuronOut = 0;
16     float *Bayes = 0;
17     float *BPerr = 0;
18     float *mapMaxOut = 0;
19     float *mapMinOut = 0;
20     int *NeurInLyr = 0;
21     int *priority = 0;
22     float *MeanErr = 0;
23     float *InputRT = 0;
24 }gpuNetParams;
25
26 vector<float> weights; //pesi della rete
27 vector<int> ArcIn; //target dell'n-esimo arco
28 vector<int> ArcOut; //base dell'n-esimo arco
29 vector<float> NeuronOut; //vettore contenente l'output dei neuroni
30 vector<float> Bayes; //vettore contenente i bayes dei neuroni
31 vector<float> BPerr; // vettore contenete gli errori retropropagati
32 vector<float> mapMaxOut; //vettore contenente il massimo valore degli output
33 vector<float> mapMinOut; //vettore contenente il minimo valore degli output
34 vector<int> priority; // vettore contenente i punti di sincronizzazione dei thread
35 vector<int> NeurInLyr; //vettore contenente gli indici dell'ultimo neurone di ogni layer
36 vector<float> examples; //vettore degli esempi
37 float MeanErr = 0; //variabile contenente l'errore medio percentuale della rete
38 int inputN, outputN; //passo di esecuzione elementi del vettore esempi
39
40 CUDACore(int nGpu) {
41     GpuID = nGpu;
42     checkCuda(cudaGetDeviceProperties(&prop, nGpu)); // carica lo struct cudaDeviceProp prop con
43       ↪ le caratteristiche della GPU con indice 0
44 }

```

Listing 9: classe di interfaccia alla GPU

Le seguenti funzioni interne della classe di interfaccia alla GPU, rappresentano il metodo di conversione della struttura della rete dalla classe MLP alla classe CudaCore, e viceversa, tale conversione porta il modello da una struttura reticolare a puntatori in una struttura a vettori lineare.

Tale struttura divide i vari parametri della rete in vettori lineari, ogni tipo di elemento viene raggruppato in un unico vettore, mentre le connessioni precedentemente rappresentate da puntatori divengono vettori a loro volta che fungono da tabelle di accesso ordinate, scorrendo dal primo elemento all'ultimo si trovano tutti gli elementi su cui eseguire sequenzialmente le operazioni da svolgere per eseguire la propagazione. La complessità di tale operazione sta proprio nello svolgimento della struttura della rete per poter eseguire con efficienza la computazione sulla scheda grafica, la quale deve essere più "machine-friendly" possibile, data l'enorme quantità di operazioni che deve essere svolta per eseguire l'addestramento. Sono indicate anche altre funzioni di conversione per Hopfield e per il dataset il quale ha bisogno anch'esso di serializzazione per l'utilizzo nelle successive funzioni.

```

1 //copia del modello da MLP class a CudaCore class
2 void cudaNetCopyMLP (MLP *pt) {
3     cout << "copying the net into CUDACore.." << endl;
4     weights.resize(pt->nArc);
5     ArcIn.resize(pt->nArc);
6     ArcOut.resize(pt->nArc);
7     NeuronOut.resize(pt->nNeurons);
8     Bayes.resize(pt->nNeurons);
9     BPerr.resize(pt->nNeurons);
10    priority.resize(pt->nLayers + 1);
11    NeurInLyr.resize(pt->nLayers + 1);
12    mapMaxOut.resize(pt->map.size());
13    mapMinOut.resize(pt->map.size());
14    inputN = pt->numNeurons(0);
15    outputN = pt->numNeurons(pt->nLayers - 1);
16
17    int NeuronIdx = 0;
18    int ArcIdx = 0;
19    vector<int> neurons(pt->nLayers);
20    //carico il vettore di mappatura dell'output della rete
21    for (int i = 0; i < pt->map.size(); i++) {
22        mapMaxOut[i] = pt->map[i].maxValue;
23        mapMinOut[i] = pt->map[i].minValue;
24    }
25
26    NeurInLyr[0] = -1; // setto il primo valore
27    priority[0] = -1; // setto il primo valore
28
29    //carico i parametri della rete
30    for (int i = 0; i < pt->nLayers; i++) {
31
32        for (int j = 0; j < pt->numNeurons(i); j++) {
33
34            Bayes[NeuronIdx] = pt->getTarget(i, j)->bayes;
35
36            for (int k = 0; k < pt->numCon(i, j); k++) {
37
38                weights[ArcIdx] = pt->getTarget(i, j)->OutArcs[k].weight;
39                ArcIn[ArcIdx] = pt->getTarget(i, j)->OutArcs[k].target->neurIdx;
40                ArcOut[ArcIdx] = pt->getTarget(i, j)->neurIdx;
41                ArcIdx++;
42            }
43            NeuronIdx++;
44        }
45        NeurInLyr[i + 1] = NeuronIdx - 1; // salvo l'indice dell'ultimo neurone del layer corrente
46        priority[i + 1] = ArcIdx - 1; // salvo l'indice dell'ultimo arco del layer corrente
47    }
48 }
49
50 //copia del modello da CudaCore class a MLP class
51 void cudaNetPasteMLP (MLP *pt) {
52     int idx = 0;
53     int Nidx = 0;
54     for (int i = 0; i < pt->nLayers; i++) {
55         for (int j = 0; j < pt->numNeurons(i); j++) {
56             pt->getTarget(i, j)->bayes = Bayes[Nidx++];
57             for (int k = 0; k < pt->numCon(i, j); k++) {
58                 pt->getTarget(i, j)->OutArcs[k].weight = weights[idx++];
59             }
60         }
61     }
62 }
63
64 void cudaNetCopyHopfield(Hopfield* pt) { ... }
65
66 void cudaNetCopyExamples (MLP *pt) { ... }

```

Listing 10: classe di interfaccia alla GPU

## 0.1.5 metodi che lanciano i kernel sulla GPU

i seguenti metodi interagiscono direttamente con la GPU allocando memoria su di essa, caricando i dati, e lanciando i kernel che eseguono le effettive operazioni sulla GPU.

```
1
2 //////////////////////////////////////////////////CUDA Kernel functions////////////////////////////////////
3 //esegue le operazioni di allocamento memoria e preparazione al lancio del kernel di
4   ↳ propagazione della rete
5 cudaError_t hostCUDAtraningNet(float eps, int Niter, int ThxBlock) {
6     cout << "learning is started!" << endl;
7     //host variables
8     float *Cweights = &weights[0];
9     int *CArcIn = &ArcIn[0];
10    int *CArcOut = &ArcOut[0];
11    float *CNeuronOut = &NeuronOut[0];
12    float *CBayes = &Bayes[0];
13    float *CBPerr = &BPerr[0];
14    float *CmapMaxOut = &mapMaxOut[0];
15    float *CmapMinOut = &mapMinOut[0];
16    float *Cexamples = &examples[0];
17    int *CNeurInLyr = &NeurInLyr[0];
18    int *Cpriority = &priority[0];
19    float *CMeanErr = &MeanErr;
20
21    //device variables
22    float *dev_weights = 0;
23    int *dev_ArcIn = 0;
24    int *dev_ArcOut = 0;
25    float *dev_examples = 0;
26    float *dev_NeuronOut = 0;
27    float *dev_Bayes = 0;
28    float *dev_BPerr = 0;
29    float *dev_mapMaxOut = 0;
30    float *dev_mapMinOut = 0;
31    int *dev_NeurInLyr = 0;
32    int *dev_priority = 0;
33    float *dev_MeanErr = 0;
34
35    //int ThxBlock = 1024;
36
37    cudaError_t cudaStatus;
38
39    // Choose which GPU to run on, change this on a multi-GPU system.
40    cudaStatus = cudaSetDevice(GpuID);
41    if (cudaCheckStatus(cudaStatus) == true) goto Error;
42
43    // Allocate GPU buffers for vectors
44    cudaStatus = cudaMalloc((void**)&dev_weights, weights.size() * sizeof(float));
45    if (cudaCheckStatus(cudaStatus) == true) goto Error;
46    cudaStatus = cudaMalloc((void**)&dev_ArcIn, ArcIn.size() * sizeof(float));
47    if (cudaCheckStatus(cudaStatus) == true) goto Error;
48    cudaStatus = cudaMalloc((void**)&dev_ArcOut, ArcOut.size() * sizeof(float));
49    if (cudaCheckStatus(cudaStatus) == true) goto Error;
50    cudaStatus = cudaMalloc((void**)&dev_NeuronOut, NeuronOut.size() * sizeof(float));
51    if (cudaCheckStatus(cudaStatus) == true) goto Error;
52    cudaStatus = cudaMalloc((void**)&dev_Bayes, Bayes.size() * sizeof(float));
53    if (cudaCheckStatus(cudaStatus) == true) goto Error;
54    cudaStatus = cudaMalloc((void**)&dev_BPerr, BPerr.size() * sizeof(float));
55    if (cudaCheckStatus(cudaStatus) == true) goto Error;
56    cudaStatus = cudaMalloc((void**)&dev_mapMaxOut, mapMaxOut.size() * sizeof(float));
57    if (cudaCheckStatus(cudaStatus) == true) goto Error;
58    cudaStatus = cudaMalloc((void**)&dev_mapMinOut, mapMinOut.size() * sizeof(float));
59    if (cudaCheckStatus(cudaStatus) == true) goto Error;
60    cudaStatus = cudaMalloc((void**)&dev_examples, examples.size() * sizeof(float));
61    if (cudaCheckStatus(cudaStatus) == true) goto Error;
62    cudaStatus = cudaMalloc((void**)&dev_NeurInLyr, NeurInLyr.size() * sizeof(int));
63    if (cudaCheckStatus(cudaStatus) == true) goto Error;
64    cudaStatus = cudaMalloc((void**)&dev_priority, priority.size() * sizeof(int));
```

```

64     if (cudaCheckStatus(cudaStatus) == true) goto Error;
65     cudaStatus = cudaMalloc((void**)&dev_MeanErr, sizeof(float));
66     if (cudaCheckStatus(cudaStatus) == true) goto Error;
67
68
69     // Copy input vectors from host memory to GPU buffers.
70     cudaStatus = cudaMemcpy(dev_weights, Cweights, weights.size() * sizeof(float),
71                             ↪ cudaMemcpyHostToDevice);
72     if (cudaCheckStatus(cudaStatus) == true) goto Error;
73     cudaStatus = cudaMemcpy(dev_ArcIn, CArcIn, ArcIn.size() * sizeof(int),
74                             ↪ cudaMemcpyHostToDevice);
75     if (cudaCheckStatus(cudaStatus) == true) goto Error;
76     cudaStatus = cudaMemcpy(dev_ArcOut, CArcOut, ArcOut.size() * sizeof(int),
77                             ↪ cudaMemcpyHostToDevice);
78     if (cudaCheckStatus(cudaStatus) == true) goto Error;
79     cudaStatus = cudaMemcpy(dev_NeuronOut, CNeuronOut, NeuronOut.size() * sizeof(float),
80                             ↪ cudaMemcpyHostToDevice);
81     if (cudaCheckStatus(cudaStatus) == true) goto Error;
82     cudaStatus = cudaMemcpy(dev_Bayes, CBayes, Bayes.size() * sizeof(float),
83                             ↪ cudaMemcpyHostToDevice);
84     if (cudaCheckStatus(cudaStatus) == true) goto Error;
85     cudaStatus = cudaMemcpy(dev_BPerr, CBPerr, BPerr.size() * sizeof(float),
86                             ↪ cudaMemcpyHostToDevice);
87     if (cudaCheckStatus(cudaStatus) == true) goto Error;
88     cudaStatus = cudaMemcpy(dev_mapMaxOut, CmapMaxOut, mapMaxOut.size() * sizeof(float),
89                             ↪ cudaMemcpyHostToDevice);
90     if (cudaCheckStatus(cudaStatus) == true) goto Error;
91     cudaStatus = cudaMemcpy(dev_mapMinOut, CmapMinOut, mapMinOut.size() * sizeof(float),
92                             ↪ cudaMemcpyHostToDevice);
93     if (cudaCheckStatus(cudaStatus) == true) goto Error;
94     cudaStatus = cudaMemcpy(dev_examples, Cexamples, examples.size() * sizeof(float),
95                             ↪ cudaMemcpyHostToDevice);
96     if (cudaCheckStatus(cudaStatus) == true) goto Error;
97     cudaStatus = cudaMemcpy(dev_NeurInLyr, CNeurInLyr, NeurInLyr.size() * sizeof(int),
98                             ↪ cudaMemcpyHostToDevice);
99     if (cudaCheckStatus(cudaStatus) == true) goto Error;
100    cudaStatus = cudaMemcpy(dev_priority, Cpriority, priority.size() * sizeof(int),
101                            ↪ cudaMemcpyHostToDevice);
102    if (cudaCheckStatus(cudaStatus) == true) goto Error;
103    cudaStatus = cudaMemcpy(dev_MeanErr, CMeanErr, sizeof(float), cudaMemcpyHostToDevice);
104    if (cudaCheckStatus(cudaStatus) == true) goto Error;
105
106
107    //////////////////////////////////lancio dei kernel all'interno della gpu/////////////////////////////////
108    int startA = 0;
109    int endA = 0;
110    int startN = 0;
111    int endN = 0;
112    int numLayerArcs = 0;
113    int numLayerNeur = 0;
114    int numOfBlocksMax = 0;
115    int numOfBlocksA = 0;
116    int numOfBlocksN = 0;
117    int numOfBlocksOut = floorf(outputN / ThxBlock) + 1;
118    int exampleRef = 0;
119    int outputRef = NeuronOut.size() - outputN;
120    long long t0 = 0, t1 = 0;
121    long long t0in = 0, t1in = 0;
122    double elapsedMilliseconds = 0;
123    double elapsedInMilliseconds = 0;
124
125    for (int it = 0; it < Niter; it++) { //scorro le iterazioni
126
127        t0 = PerformanceCounter();
128
129        for (int t = 0; t < (examples.size() / (inputN + outputN)); t++) { //scorro gli esempi
130            //imposto il riferimento per l'esempio di input
131            exampleRef = t * (inputN + outputN);

```

```

122 t0in = PerformanceCounter();
123 //resetto il vettore contenente lo stato di attivazione dei neuroni
124 numOfBlocksA = (floorf(NeuronOut.size() / ThxBlock) + 1);
125 CUDAResetVector <<<numOfBlocksA, ThxBlock >>> (dev_NeuronOut, NeuronOut.size());
126
127 //imposto i valori di input ai neuroni dello strato input
128 numOfBlocksA = (floorf(inputN / ThxBlock) + 1);
129 CUDASetInput <<<numOfBlocksA, ThxBlock >>> (dev_NeuronOut, inputN, exampleRef,
    ↪ dev_examples);
130
131 //propagazione dell'input nella rete
132
133 startA = 0; // indice di partenza dei vettori archi
134 endA = 0; // ultimo indice dei vettori archi
135 startN = 0; // indice di partenza dei vettori neuroni
136 endN = 0; // ultimo indice dei vettori neuroni
137
138 for (int i = 0; i < priority.size() - 1; i++) { //NB non viene applicata la sigmoide
    ↪ allo strato di input eventualmente correggi
139
140     startA = priority[i] + 1;
141     endA = priority[i + 1];
142
143     if (i < priority.size() - 2) {
144         startN = NeurInLyr[i + 1] + 1;
145         endN = NeurInLyr[i + 2];
146     }
147
148     numLayerArcs = endA - startA + 1;
149     numLayerNeur = endN - startN + 1;
150
151     numOfBlocksA = floorf(numLayerArcs / ThxBlock) + 1;
152     numOfBlocksN = floorf(numLayerNeur / ThxBlock) + 1;
153
154     if (i < priority.size() - 2) {
155         //propago l'output dei neuroni al prossimo/i layer
156         CUDALayerInput <<<numOfBlocksA, ThxBlock >>> (dev_weights, dev_ArcIn, dev_ArcOut,
            ↪ dev_NeuronOut, startA, endA);
157         //applico il contributo dei bayes all output dei neuroni del layer corrente
158         CUDABayesInput <<< numOfBlocksN, ThxBlock >>> (dev_NeuronOut, dev_Bayes, startN,
            ↪ endN);
159         //applico la sigmoide allo stato di attivazione dei neuroni
160         CUDASigLayer <<<numOfBlocksN, ThxBlock >>> (dev_NeuronOut, startN, endN);
161     }
162 }
163
164 tlin = PerformanceCounter();
165 elapsedInMilliseconds += ((tlin - t0in) * 1000.0) / PerformanceFrequency();
166
167
168 //resetto il vettore contenente l'errore dei neuroni
169 numOfBlocksN = (floorf(BPerr.size() / ThxBlock) + 1);
170 CUDAResetVector <<<numOfBlocksN, ThxBlock >>> (dev_BPerr, BPerr.size());
171
172 CUDAResetVar <<<1, 1 >>> (dev_MeanErr);
173 CUDAOuputErr <<<numOfBlocksOut, ThxBlock >>> (dev_NeuronOut, outputRef, NeuronOut.size
    ↪ (), inputN, dev_BPerr, dev_examples, exampleRef, dev_mapMaxOut, dev_mapMinOut,
    ↪ dev_MeanErr);
174 cudaMemcpy(CMeanErr, dev_MeanErr, sizeof(float), cudaMemcpyDeviceToHost);
175
176 MeanErr += *CMeanErr / outputN;
177
178
179 //retropropagazione dell'errore
180
181 for (int i = priority.size() - 2; i > 1; i--) {
182
183     startA = priority[i - 1] + 1;
184     endA = priority[i];

```

```

185     startN = NeurInLyr[i - 1] + 1;
186     endN = NeurInLyr[i];
187
188     numLayerArcs = endA - startA + 1;
189     numLayerNeur = endN - startN + 1;
190
191     numOfBlocksA = floorf(numLayerArcs / ThxBlock) + 1;
192     numOfBlocksN = floorf(numLayerNeur / ThxBlock) + 1;
193     //numOfBlocksMax = maxOf(numOfBlocksA, numOfBlocksN);
194
195     CUDAPropagationErr <<<numOfBlocksA, ThxBlock >>> (dev_BPerr, dev_weights,
196         ↪ dev_NeuronOut, dev_ArcIn, dev_ArcOut, startA, endA);
197     CUDAOutDiff <<<numOfBlocksN, ThxBlock >>> (dev_BPerr, dev_NeuronOut, startN, endN);
198     cudaStatus = cudaMemcpy(CBPerr, dev_BPerr, BPerr.size() * sizeof(float),
199         ↪ cudaMemcpyDeviceToHost);
200     if (cudaCheckStatus(cudaStatus) == true) goto Error;
201     copy(CBPerr, CBPerr + BPerr.size(), BPerr.begin());
202 }
203
204 //applico a ogni peso la sua correzione
205
206 startN = NeurInLyr[1] + 1; // la correzione dei bais va applicata dal primo layer
207 ↪ nascosto in poi
208 endN = NeurInLyr[NeurInLyr.size() - 1];
209
210 numLayerNeur = endN - startN + 1;
211
212 numOfBlocksA = floorf(weights.size() / ThxBlock) + 1;
213 numOfBlocksN = floorf(numLayerNeur / ThxBlock) + 1;
214
215 CUDAApplyWeightCorrections <<<numOfBlocksA, ThxBlock >>> (eps, dev_NeuronOut, dev_BPerr,
216     ↪ dev_weights, dev_ArcIn, dev_ArcOut, weights.size());
217 CUDAApplyBayesCorrections <<<numOfBlocksN, ThxBlock >>> (eps, dev_BPerr, dev_Bayes,
218     ↪ startN, endN);
219 }
220
221 t1 = PerformanceCounter();
222 elapsedMilliseconds = ((t1 - t0) * 1000.0) / PerformanceFrequency(); // calcolo il tempo
223 ↪ di esecuzione di una iterazione di addestramento (tutto il set)
224 MeanErr = MeanErr / (examples.size() / (inputN + outputN)); //calcolo l'errore percentuale
225 ↪ medio sul dataset
226 elapsedInMilliseconds = elapsedInMilliseconds / (examples.size() / (inputN + outputN));
227 cout << "Iterazione: " << it << " " << MeanErr << " %Err " << "execution time:" <<
228     ↪ elapsedMilliseconds << "ms" << endl;
229 cout << "mean InputTime: " << elapsedInMilliseconds << "ms" << endl;
230 printNetSpecs();
231 MeanErr = 0;
232 }
233
234 cudaStatus = cudaMemcpy(Cweights, dev_weights, weights.size() * sizeof(float),
235     ↪ cudaMemcpyDeviceToHost);
236 if (cudaCheckStatus(cudaStatus) == true) goto Error;
237 copy(Cweights, Cweights + weights.size(), weights.begin());
238
239 cudaStatus = cudaMemcpy(CBayes, dev_Bayes, Bayes.size() * sizeof(float),
240     ↪ cudaMemcpyDeviceToHost);
241 if (cudaCheckStatus(cudaStatus) == true) goto Error;
242 copy(CBayes, CBayes + Bayes.size(), Bayes.begin());
243 //checkpoint di errore (se la GPU richiama un qualunque errore ripare da qui)
244 Error:
245
246 //libero la memoria nella scheda grafica
247 cudaFree(dev_weights);
248 cudaFree(dev_ArcIn);
249 cudaFree(dev_ArcOut);
250 cudaFree(dev_NeuronOut);
251 cudaFree(dev_examples);
252 cudaFree(dev_BPerr);
253 cudaFree(dev_mapMaxOut);

```

```

244     cudaFree(dev_mapMinOut);
245     cudaFree(dev_priority);
246     cudaFree(dev_NeurInLyr);
247
248     //ritorno lo stato della GPU
249     return cudaStatus;
250 }
251
252 //esegue il caricamento nella gpu dei parametri della rete
253 cudaError_t hostCUAuploadNetParams() {
254
255     cudaError_t cudaStatus;
256
257     cudaStatus = cudaSetDevice(GpuID);
258     if (cudaCheckStatus(cudaStatus) == true) goto Error;
259
260     //host variables
261     float *Cweights = &weights[0];
262     int *CArcIn = &ArcIn[0];
263     int *CArcOut = &ArcOut[0];
264     float *CNeuronOut = &NeuronOut[0];
265     float *CBayes = &Bayes[0];
266     float *CBPerr = &BPerr[0];
267     float *CmapMaxOut = &mapMaxOut[0];
268     float *CmapMinOut = &mapMinOut[0];
269     float *Cexamples = &examples[0];
270     int *CNeurInLyr = &NeurInLyr[0];
271     int *Cpriority = &priority[0];
272     float *CMeanErr = &MeanErr;
273
274     // Choose which GPU to run on, change this on a multi-GPU system.
275     cudaStatus = cudaSetDevice(GpuID);
276     if (cudaCheckStatus(cudaStatus) == true) goto Error;
277
278     // Allocate GPU buffers for vectors
279     cudaStatus = cudaMalloc((void**)&gpuNetParams.weights, weights.size() * sizeof(float));
280     if (cudaCheckStatus(cudaStatus) == true) goto Error;
281     cudaStatus = cudaMalloc((void**)&gpuNetParams.ArcIn, ArcIn.size() * sizeof(float));
282     if (cudaCheckStatus(cudaStatus) == true) goto Error;
283     cudaStatus = cudaMalloc((void**)&gpuNetParams.ArcOut, ArcOut.size() * sizeof(float));
284     if (cudaCheckStatus(cudaStatus) == true) goto Error;
285     cudaStatus = cudaMalloc((void**)&gpuNetParams.NeuronOut, NeuronOut.size() * sizeof(float));
286     if (cudaCheckStatus(cudaStatus) == true) goto Error;
287     cudaStatus = cudaMalloc((void**)&gpuNetParams.Bayes, Bayes.size() * sizeof(float));
288     if (cudaCheckStatus(cudaStatus) == true) goto Error;
289     cudaStatus = cudaMalloc((void**)&gpuNetParams.BPerr, BPerr.size() * sizeof(float));
290     if (cudaCheckStatus(cudaStatus) == true) goto Error;
291     cudaStatus = cudaMalloc((void**)&gpuNetParams.mapMaxOut, mapMaxOut.size() * sizeof(float));
292     if (cudaCheckStatus(cudaStatus) == true) goto Error;
293     cudaStatus = cudaMalloc((void**)&gpuNetParams.mapMinOut, mapMinOut.size() * sizeof(float));
294     if (cudaCheckStatus(cudaStatus) == true) goto Error;
295     cudaStatus = cudaMalloc((void**)&gpuNetParams.examples, examples.size() * sizeof(float));
296     if (cudaCheckStatus(cudaStatus) == true) goto Error;
297     cudaStatus = cudaMalloc((void**)&gpuNetParams.NeurInLyr, NeurInLyr.size() * sizeof(int));
298     if (cudaCheckStatus(cudaStatus) == true) goto Error;
299     cudaStatus = cudaMalloc((void**)&gpuNetParams.priority, priority.size() * sizeof(int));
300     if (cudaCheckStatus(cudaStatus) == true) goto Error;
301     cudaStatus = cudaMalloc((void**)&gpuNetParams.InputRT, inputN * sizeof(float));
302     if (cudaCheckStatus(cudaStatus) == true) goto Error;
303     cudaStatus = cudaMalloc((void**)&gpuNetParams.MeanErr, sizeof(float));
304     if (cudaCheckStatus(cudaStatus) == true) goto Error;
305
306
307     // Copy input vectors from host memory to GPU buffers.
308     cudaStatus = cudaMemcpy(gpuNetParams.weights, Cweights, weights.size() * sizeof(float),
309                             ↪ cudaMemcpyHostToDevice);
309     if (cudaCheckStatus(cudaStatus) == true) goto Error;
310     cudaStatus = cudaMemcpy(gpuNetParams.ArcIn, CArcIn, ArcIn.size() * sizeof(int),
311                             ↪ cudaMemcpyHostToDevice);

```



```

311     if (cudaCheckStatus(cudaStatus) == true) goto Error;
312     cudaStatus = cudaMemcpy(gpuNetParams.ArcOut, CArcOut, ArcOut.size() * sizeof(int),
        ↪ cudaMemcpyHostToDevice);
313     if (cudaCheckStatus(cudaStatus) == true) goto Error;
314     cudaStatus = cudaMemcpy(gpuNetParams.NeuronOut, CNeuronOut, NeuronOut.size() * sizeof(float)
        ↪ , cudaMemcpyHostToDevice);
315     if (cudaCheckStatus(cudaStatus) == true) goto Error;
316     cudaStatus = cudaMemcpy(gpuNetParams.Bayes, CBayes, Bayes.size() * sizeof(float),
        ↪ cudaMemcpyHostToDevice);
317     if (cudaCheckStatus(cudaStatus) == true) goto Error;
318     cudaStatus = cudaMemcpy(gpuNetParams.BPerr, CBPerr, BPerr.size() * sizeof(float),
        ↪ cudaMemcpyHostToDevice);
319     if (cudaCheckStatus(cudaStatus) == true) goto Error;
320     cudaStatus = cudaMemcpy(gpuNetParams.mapMaxOut, CmapMaxOut, mapMaxOut.size() * sizeof(float)
        ↪ , cudaMemcpyHostToDevice);
321     if (cudaCheckStatus(cudaStatus) == true) goto Error;
322     cudaStatus = cudaMemcpy(gpuNetParams.mapMinOut, CmapMinOut, mapMinOut.size() * sizeof(float)
        ↪ , cudaMemcpyHostToDevice);
323     if (cudaCheckStatus(cudaStatus) == true) goto Error;
324     cudaStatus = cudaMemcpy(gpuNetParams.examples, Cexamples, examples.size() * sizeof(float),
        ↪ cudaMemcpyHostToDevice);
325     if (cudaCheckStatus(cudaStatus) == true) goto Error;
326     cudaStatus = cudaMemcpy(gpuNetParams.NeurInLyr, CNeurInLyr, NeurInLyr.size() * sizeof(int),
        ↪ cudaMemcpyHostToDevice);
327     if (cudaCheckStatus(cudaStatus) == true) goto Error;
328     cudaStatus = cudaMemcpy(gpuNetParams.priority, Cpriority, priority.size() * sizeof(int),
        ↪ cudaMemcpyHostToDevice);
329     if (cudaCheckStatus(cudaStatus) == true) goto Error;
330     cudaStatus = cudaMemcpy(gpuNetParams.MeanErr, CMeanErr, sizeof(float),
        ↪ cudaMemcpyHostToDevice);
331     if (cudaCheckStatus(cudaStatus) == true) goto Error;
332
333     if (false) {
334         Error:
335         //libero la memoria nella scheda grafica
336         cudaFree(gpuNetParams.weights);
337         cudaFree(gpuNetParams.ArcIn);
338         cudaFree(gpuNetParams.ArcOut);
339         cudaFree(gpuNetParams.NeuronOut);
340         cudaFree(gpuNetParams.examples);
341         cudaFree(gpuNetParams.BPerr);
342         cudaFree(gpuNetParams.mapMaxOut);
343         cudaFree(gpuNetParams.mapMinOut);
344         cudaFree(gpuNetParams.priority);
345         cudaFree(gpuNetParams.NeurInLyr);
346         cout << "ERRORE: libero la memoria della gpu. " << endl;
347     }
348
349     return cudaStatus;
350 }
351
352 //esegue il download dalla gpu dei parametri della rete
353 cudaError_t hostCUDAdownloadNetParams() {
354
355     cout << "downloading net params from gpu.." << endl;
356
357     float *Cweights = &weights[0];
358     float *CBayes = &Bayes[0];
359
360     cudaError_t cudaStatus;
361
362     cudaStatus = cudaSetDevice(GpuID);
363     if (cudaCheckStatus(cudaStatus) == true) goto Error;
364
365     cudaStatus = cudaMemcpy(Cweights, gpuNetParams.weights, weights.size() * sizeof(float),
        ↪ cudaMemcpyDeviceToHost);
366     if (cudaCheckStatus(cudaStatus) == true) goto Error;
367
368     cudaStatus = cudaMemcpy(CBayes, gpuNetParams.Bayes, Bayes.size() * sizeof(float),

```



```

369     ↪ cudaMemcpyDeviceToHost);
370 if (cudaCheckStatus(cudaStatus) == true) goto Error;
371
372 if (false) {
373 Error:
374     //libero la memoria nella scheda grafica
375     cudaFree(gpuNetParams.weights);
376     cudaFree(gpuNetParams.ArcIn);
377     cudaFree(gpuNetParams.ArcOut);
378     cudaFree(gpuNetParams.NeuronOut);
379     cudaFree(gpuNetParams.examples);
380     cudaFree(gpuNetParams.BPerr);
381     cudaFree(gpuNetParams.mapMaxOut);
382     cudaFree(gpuNetParams.mapMinOut);
383     cudaFree(gpuNetParams.priority);
384     cudaFree(gpuNetParams.NeurInLyr);
385     cout << "ERRORE: libero la memoria della gpu. " << endl;
386 }
387
388 return cudaStatus;
389 }
390
391 //esegue l'input della rete gia addestrata prendendo in input l'esempio dato
392 cudaError_t hostCUDAInputNet(float *input, int ThxBLOCK) {
393     //importante verificare che l'input abbia la stessa dimansione dell'input della rete
394
395     cudaError cudaStatus;
396
397     cudaStatus = cudaSetDevice(GpuID);
398     if (cudaCheckStatus(cudaStatus) == true) goto Error;
399
400     ////////////lancio dei kernel all'interno della gpu//////////
401     //int ThxBLOCK = 1024;
402     int startA = 0;
403     int endA = 0;
404     int startN = 0;
405     int endN = 0;
406     int numLayerArcs = 0;
407     int numLayerNeur = 0;
408     int numOfBlocksMax = 0;
409     int numOfBlocksA = 0;
410     int numOfBlocksN = 0;
411     int numOfBlocksOut = floorf(outputN / ThxBLOCK) + 1;
412     int outputRef = NeuronOut.size() - outputN;
413     long long t0in = 0, t1in = 0;
414     double elapsedInMilliseconds = 0;
415
416     t0in = PerformanceCounter();
417     //resetto il vettore contenente lo stato di attivazione dei neuroni
418     numOfBlocksA = (floorf(NeuronOut.size() / ThxBLOCK) + 1);
419     CUDAResetVector <<<numOfBlocksA, ThxBLOCK >>> (gpuNetParams.NeuronOut, NeuronOut.size());
420
421     cudaStatus = cudaMemcpy(gpuNetParams.InputRT, input, inputN * sizeof(float),
422         ↪ cudaMemcpyHostToDevice);
423     if (cudaCheckStatus(cudaStatus) == true) goto Error;
424
425     //imposto i valori di input ai neuroni dello strato input
426     numOfBlocksA = (floorf(inputN / ThxBLOCK) + 1);
427     CUDASetSingleInput <<<numOfBlocksA, ThxBLOCK>>> (gpuNetParams.NeuronOut, inputN,
428         ↪ gpuNetParams.InputRT);
429
430     //propagazione dell'input nella rete
431
432     startA = 0; // indice di partenza dei vettori archi
433     endA = 0; // ultimo indice dei vettori archi
434     startN = 0; // indice di partenza dei vettori neuroni
435     endN = 0; // ultimo indice dei vettori neuroni

```

```

435     for (int i = 0; i < priority.size() - 1; i++) { //NB non viene applicata la sigmoide allo
436         ↪ strato di input eventualmente correggi
437
438         startA = priority[i] + 1;
439         endA = priority[i + 1];
440
441         if (i < priority.size() - 2) {
442             startN = NeurInLyr[i + 1] + 1;
443             endN = NeurInLyr[i + 2];
444         }
445
446         numLayerArcs = endA - startA + 1;
447         numLayerNeur = endN - startN + 1;
448
449         numOfBlocksA = floorf(numLayerArcs / ThxBlock) + 1;
450         numOfBlocksN = floorf(numLayerNeur / ThxBlock) + 1;
451
452         if (i < priority.size() - 2) {
453             CUDALayerInput <<<numOfBlocksA, ThxBlock >>> (gpuNetParams.weights, gpuNetParams.ArcIn,
454                 ↪ gpuNetParams.ArcOut, gpuNetParams.NeuronOut, startA, endA); //propago l'output dei
455                 ↪ neuroni al prossimo/i layer
456             CUDABayesInput <<<numOfBlocksN, ThxBlock >>> (gpuNetParams.NeuronOut, gpuNetParams.Bayes
457                 ↪ , startN, endN); //applico il contributo dei bayes all output dei neuroni del
458                 ↪ layer corrente
459             CUDASigLayer <<<numOfBlocksN, ThxBlock >>> (gpuNetParams.NeuronOut, startN, endN); //
460                 ↪ applico la sigmoide allo stato di attivazione dei neuroni
461         }
462     }
463
464     //copio l'output dei neuroni dello strato output nella memoria della cpu
465     cudaStatus = cudaMemcpy(&NeuronOut[0] + outputRef, gpuNetParams.NeuronOut + outputRef,
466         ↪ outputN * sizeof(float), cudaMemcpyDeviceToHost); //TODO da errore e non carica il
467         ↪ vettore trovare il BUG
468     if (cudaCheckStatus(cudaStatus) == true) goto Error;
469
470     tlin = PerformanceCounter();
471     elapsedInMilliseconds = ((tlin - t0in) * 1000.0) / PerformanceFrequency();
472
473     //////////////////////////////////visualizzazione dell'esempio////////////////////////////////////
474
475     float delta;
476     cout << "input time: " << elapsedInMilliseconds << " ms" << endl;
477
478     for (int on = 0; on < outputN; on++) {
479         delta = mapMaxOut[on] - mapMinOut[on];
480         cout << "Y" << on << ": " << (NeuronOut[NeuronOut.size() - outputN + on] * delta) +
481             ↪ mapMinOut[on] << endl;
482     }
483     cout << endl;
484
485     //////////////////////////////////////////////////////////////////////
486
487     if (false) {
488     Error:
489         //libero la memoria nella scheda grafica
490         cudaFree(gpuNetParams.weights);
491         cudaFree(gpuNetParams.ArcIn);
492         cudaFree(gpuNetParams.ArcOut);
493         cudaFree(gpuNetParams.NeuronOut);
494         cudaFree(gpuNetParams.examples);
495         cudaFree(gpuNetParams.BPerr);
496         cudaFree(gpuNetParams.mapMaxOut);
497         cudaFree(gpuNetParams.mapMinOut);
498         cudaFree(gpuNetParams.priority);
499         cudaFree(gpuNetParams.NeurInLyr);
500     }
501
502     return cudaStatus;

```

```

495 }
496 //////////////////////////////////////////////////
497
498 //////////////////////////////////////////////////CUDA UTILITY////////////////////////////////////
499 //verifica la corretta esecuzione di un operazione
500 inline cudaError_t checkCuda(cudaError_t result){ ... }
501
502 //verifica la corretta esecuzione di un operazione restituendo un bool
503 bool cudaCheckStatus(cudaError_t cudaStatus) { ... }
504
505 //stampa a schermo le principali proprietà della scheda
506 void printDeviceSpecs() { ... }
507
508 //stampa i parametri della rete che vengono passati alla scheda
509 void printNetSpecs() { ... }
510
511 //calcola il peso del modello
512 float sizeOfModel(string mesureUnit = "B") { ... }
513
514 template<typename T, typename A>
515 float sizeOfVector(vector<T, A> const& vect,string mesureUnit = "B") { ... }
516 //////////////////////////////////////////////////
517 };

```

Listing 11: classe di interfaccia alla GPU