

1 Prefazione

1.1 scopo del progetto

L'obiettivo di questo progetto è la realizzazione di un modello di machine learning in grado di apprendere le relazioni tra i dati fruiti da una scheda embeded dotata di accelerometro, giroscopio e magnetometro, e la posizione spaziale di un braccio umano. Il problema che ci si è posti di risolvere poteva essere risolto in modo deterministico ma per lo scopo del progetto si voleva risolvere il problema tramite un algoritmo di ML, per poter verificare l'efficacia di questo metodo per la risoluzione di un problema con soluzione nota.

1.2 introduzione generale alle tecniche utilizzate

Per la realizzazione del progetto è stata realizzata una scheda con microcontrollore arduino, provvista di bluetooth, accelerometro, magnetometro, alimentata a batteria. È stata realizzata una libreria in C++/CUDA per la creazione, l'addestramento e l'utilizzo di una rete neurale con supporto per CPU e GPU. È stato realizzato un secondo programma C++, per l'acquisizione dei dati della scheda e per l'acquisizione dei punti dello scheletro forniti dal kinect, lo stesso programma si occupa della sincronizzazione di tali dati e della creazione di un dataset necessario per il processo di addestramento della rete neurale. In fine è stato realizzato un programma C# che sfrutta il game engine Unity per la realizzazione dell'ambiente grafico, che rende possibile la visualizzazione dell'output della rete, muovendo un manichino secondo gli input della scheda. Quest'ultima operazione è stata effettuata mettendo in comunicazione il programma C++ che si occupa dell'acquisizione dei dati della scheda, e che ne esegue l'input nel modello addestrato, per poi passare i dati elaborati tramite comunicazione socket interna al motore grafico.

2 Richiami di teoria

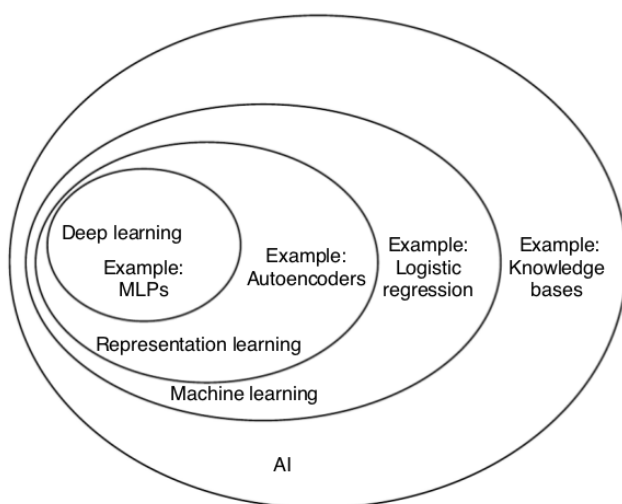
2.1 Dal machine learning al deep learning

Oggi l'intelligenza artificiale è un fiorente campo di ricerca, con l'obiettivo di risolvere una grande varietà di problemi, che per essere risolti attraverso la programmazione classica necessitano di una grande quantità di conoscenze pregresse, spesso non disponibili.

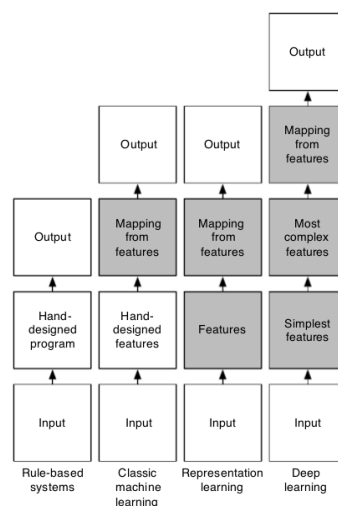
I Programmi basati sul paradigma dell'intelligenza artificiale si propongono di superare questi ostacoli acquisendo direttamente queste conoscenze dai dati grezzi, tale capacità è nota come machine learning. Sotto questa famiglia di algoritmi si trovano altri sottogruppi quali il representation learning e all'interno di quest'ultimo il deep learning.

Il deep learning rispetto ai metodi più classici, tipicamente in grado di riconoscere relazioni lineari (come ad esempio il noto SVM o support vector machine), si propone come alternativa per l'apprendimento di funzioni a molte variabili non lineari anche molto complesse. Le applicazioni più comuni comprendono il riconoscimento di oggetti nelle immagini, delle parole in tracce audio o anche per le traduzioni multilingue.

Il termine "deep learning" deriva proprio dalla capacità di questi modelli di riuscire a cogliere delle relazioni molto "profonde" tra i dati di ingresso e di uscita, approssimando con un certo errore, dipendente dal caso specifico, la funzione che dati gli input restituisce l'output desiderato, purtroppo però presentano il grande problema di non rendere disponibile in maniera chiara le relazioni apprese.



famiglie di algoritmi



differenze generali tra le diverse tipologie di algoritmi

2.1.1 Le reti feed-forward (MLP)

Uno dei modelli della famiglia del deep learning più comuni e semplici da utilizzare sono le reti neurali feed-forward o multi layer perceptron, le quali sono assimilabili a modelli di regressione statistica.

Le reti neurali sono un modello matematico molto semplificato del cervello, e per parlarne è necessario qualche riferimento al modello biologico.

Il neurone biologico

L'unità base del cervello è rappresentata dai neuroni (Figura 2), i quali a loro volta sono composti dai dendriti, dal soma, dall'assone e dalle sinapsi.

I dendriti sono l'apparato di input del neurone, attraverso il quale riceve i segnali dall'ambiente o da altri neuroni, tali segnali caricano il neurone, che accumula un potenziale elettrico all'interno del soma. In determinate condizioni o al raggiungimento di una certa soglia di potenziale il neurone emette un impulso lungo l'assone giungendo alle sinapsi, poste alla fine dell'assone, le quali sono connesse con altri neuroni o con altri apparati del corpo, come i muscoli.

Questa particolare cellula oltre a scambiare impulsi è in grado di accumulare informazioni attraverso l'ispessimento o l'assottigliamento della guaina mielinica, la quale è presente lungo l'assone. Le variazioni dello spessore della guaina comportano un aumento o una diminuzione della resistenza elettrica dell'assone determinando una variazione dell'impulso percepito dai neuroni connessi alle sinapsi a parità di potenziale rilasciato.

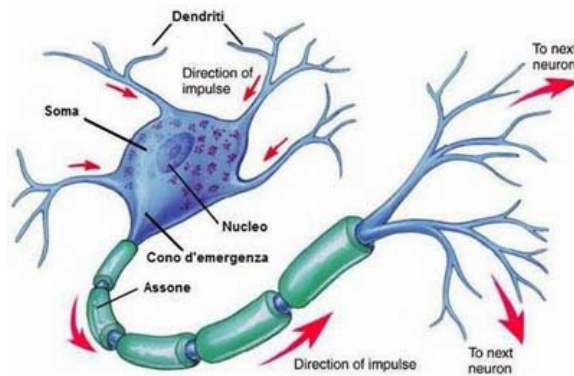


Figura 2: Il neurone biologico

Il perceptrone

Il neurone artificiale è composto da una serie di connessioni in ingresso, come i dendriti nel caso del neurone biologico, questa volta però il compito di immagazzinare informazione viene svolto da queste connessioni in input e non più dalle connessioni in output. Ad ogni input del neurone è associato un coefficiente che moltiplica il potenziale in ingresso, a questo punto gli ingressi pesati giungono al nodo sommatore, il corrispettivo del soma, che ne accumula la sommatoria. Lo step successivo è il blocco contenente la funzione di trasferimento, una funzione che prende in input il potenziale del neurone e rende disponibile il risultato come output finale o come input per i neuroni successivi. La funzione di trasferimento del neurone può essere di vari tipi, può variare infatti anche tra i neuroni della medesima rete, alcune delle più comuni sono:

- (a) la funzione gradino.
- (b) la funzione lineare con saturazione
- (c) la funzione sigmoide tra 0 e 1
- (d) la funzione sigmoide tra -1 e 1

Volendo esprimere il perceptrone in termini matematici si ottiene la seguente equazione:

$$y = f(P) = f(\sum_i W_i \cdot X_i)$$

y - output del perceptrone

f - funzione di trasferimento del perceptrone

P - potenziale del neurone

W_i - peso o coefficiente dell' i -esimo ingresso

X_i - potenziale d'ingresso dell' i -esimo neurone

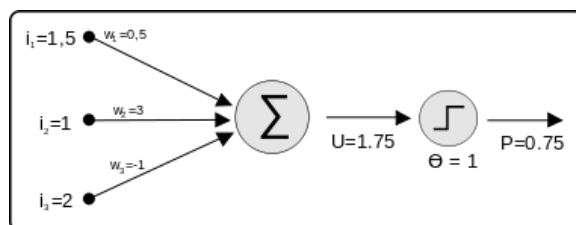


Figura 3: Il perceptrone

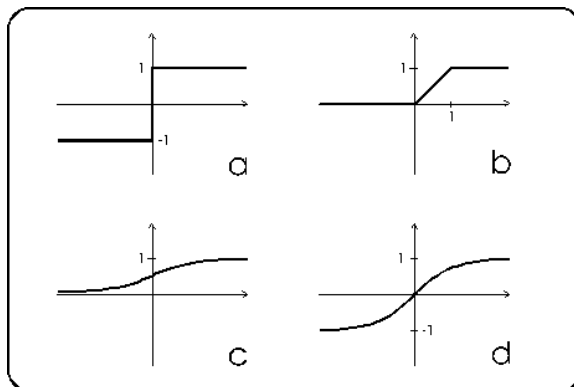
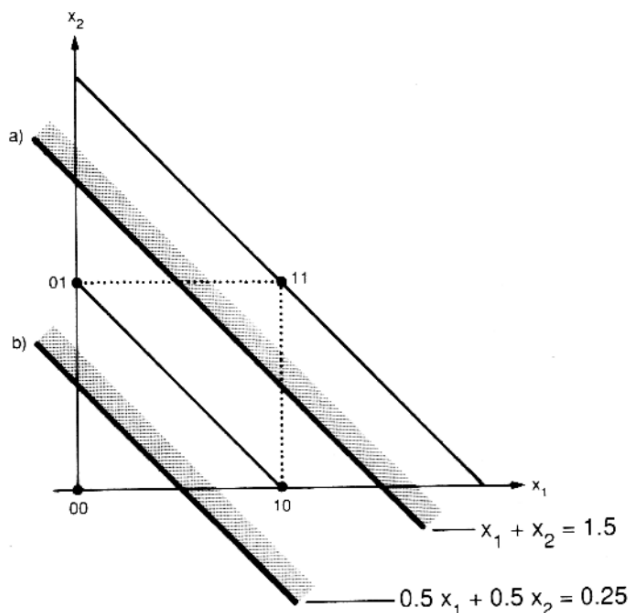


Figura 4: Le funzioni di trasferimento

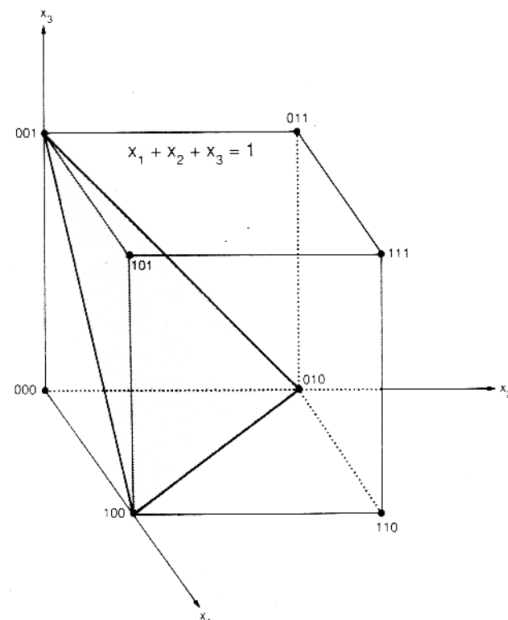
Una delle caratteristiche più importanti del perceptrone è la sua capacità di compiere separazioni lineari. Tali separazioni avvengono in uno spazio a n dimensioni, dove n corrisponde al numero di variabili in ingresso al neurone, perciò si avrà ad esempio per un perceptrone a due input la capacità di separare i valori su \mathbb{R}^2 attraverso una retta, nel caso \mathbb{R}^3 si potranno separare i valori dello spazio tridimensionale attraverso un piano, e così via.

Le reti multi-strato e la separazione non lineare

Le reti multi-strato sono composte da strati di perceptron, interconnessi tra loro in modo che gli output degli strati precedenti siano connessi con gli input dei perceptron successivi. Tale struttura "profonda" permette di rappresentare funzioni non lineari molto complesse, tanto più sono gli strati della rete.



esempio di linee di separazione per neurone a 2 ingressi binari
(a - funzione OR)
(b - funzione AND)

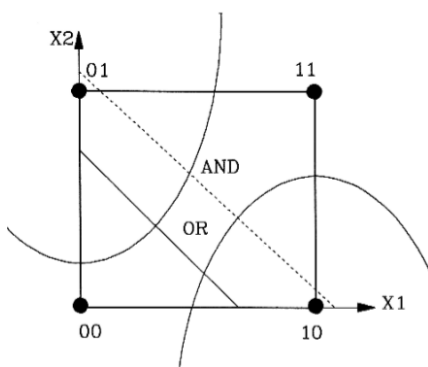


esempio di piano di separazione lineare
a 3 ingressi per un neurone binario

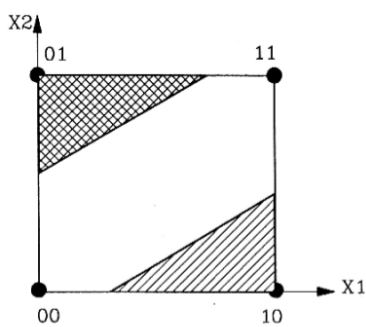
Uno degli esempi che meglio chiarisce questa problematica è la funzione xor nel seguente esempio, la quale per semplicità è rappresentata con neuroni binari a soglia.

Il problema dello XOR

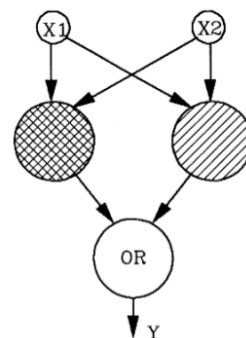
Lo xor è una delle funzioni non lineari più semplice, nella quale diventa evidente l'impossibilità di approssimare tale funzione con un solo strato di neuroni. Il grafico illustra come un neurone singolo possa effettuare la separazione lineare che riproduce la funzione and e or, ma come si nota lo xor necessita di due parabole per la separazione e non più di rette, tale separazione si effettua utilizzando una rete ad almeno due strati.



Esempio di linee di separazione per and or e xor



MLP in grado di rappresentare la funzione xor



2.1.2 La funzione di addestramento: Il back-propagation

Come già detto la struttura della rete neurale deve apprendere una data funzione da dei dati.

Tale operazione si effettua presentando degli input alla rete ai quali questa risponderà con degli output, i quali dovranno essere confrontati con degli output corretti per il set di input presentato, per calcolare l'errore nell'output della rete.

L'errore che la rete commette ad ogni tentativo viene utilizzato per calcolare attraverso una prescelta funzione, nel nostro caso il BP, dei coefficienti di correzione che vanno applicati ai pesi della rete per avvicinare il suo output a quello desiderato.

Come dice il nome "retropropagazione" in italiano, l'errore calcolato dall'output viene applicato ai neuroni precedenti e poi a quelli ancora precedenti fino ad arrivare al primo strato, tale procedura può essere effettuata in diversi modi, calcolando gli errori su un determinato numero di esempi del dataset, accumulando le correzioni ed applicandole successivamente, tale tecnica è attualmente la più utilizzata ed è detta "mini-batch", ma è possibile anche effettuare la correzione esempio per esempio, e questo è detto addestramento in-line, in fine tra i metodi più comuni c'è anche la possibilità di eseguire tutti gli esempi del dataset e solo alla fine di un'intero ciclo di addestramento applicare le correzioni. questi tre metodi ovviamente risultano validi ma a seconda del problema da risolvere possono

avere prestazioni migliori o peggiori in modo difficilmente prevedibile.

Dopo questa introduzione generale all'addestramento delle reti passiamo a dare una vista più accurata dell'algoritmo utilizzato nel progetto che è principalmente pensato per reti mlp con neuroni ad ingressi continui e con funzione di trasferimento sigmoidea o ad arcotangente.

Il BP cerca di minimizzare l'errore quadratico medio, relativo a un training set di esempi dati, scendendo lungo la superficie d'errore e seguendo la direzione di massima pendenza, in cerca di una "valle" sufficientemente profonda.

Consideriamo una rete con n input X_i ($i = 1, 2, \dots, n$), uno strato nascosto di q neuroni Z_k ($k = 1, 2, \dots, q$) e uno strato output di m neuroni Y_j ; ($j = 1, 2, \dots, m$). Il training set sia costituito da p esempi o casi C_r , ($r = 1, 2, \dots, p$). Tutti i neuroni abbiano la stessa funzione di trasferimento.

L'errore quadratico medio dello strato output è:

$$E = \frac{1}{2} \sum_j \sum_r (Y_{rj} - D_{rj})^2 \quad (1)$$

dove Y_j è l'output del neurone Y_j alla presentazione dell'esempio C_r e D_{rj} è il suo valore desiderato. Per semplificare la notazione, eliminiamo l'indice r e minimizziamo l'errore quadratico medio ad ogni presentazione di esempio (modalità on-line), anziché dopo ogni ciclo completo di presentazione del training set o epoca (modalità batch). Avremo allora:

$$E = \frac{1}{2} \sum_j (Y_j - D_j)^2 \quad (2)$$

Si applica a questo punto la regola del gradiente (fig. Propagazione a):

$$\Delta W_{jk} = -\eta \frac{\partial E}{\partial W_{jk}} = -\eta \frac{\partial E}{\partial Y_j} \frac{\partial Y_j}{\partial P_j} \frac{\partial P_j}{\partial W_{jk}} = -\eta (Y_j - D_j) f'(P_j) Z_k \quad (3)$$

$$\frac{\partial E}{\partial Y_j} = (Y_j - D_j) \quad (4)$$

$$\frac{\partial Y_j}{\partial P_j} = f'(P_j) \quad (5)$$

$$\frac{\partial P_j}{\partial W_{jk}} = \frac{\partial (\sum_k W_{jk} Z_k)}{\partial W_{jk}} = Z_k \quad (6)$$

Ponendo a questo punto:

$$\delta_j = (Y_j - D_j) f'(P_j) \quad (7)$$

$$\Delta W_{jk} = -\eta \delta_j Z_k \quad (8)$$

La formula (3) viene impiegata per aggiornare i pesi sinattici delle connessioni tra strato nascosto e strato output, analogamente, per quanto riguarda le connessioni tra strato input e strato nascosto si utilizza la formula (fig. Propagazione b):

$$\Delta W_{ki} = -\eta \frac{\partial E}{\partial W_{ki}} = -\eta \frac{\partial E}{\partial Z_k} \frac{\partial Z_k}{\partial P_k} \frac{\partial P_k}{\partial W_{ki}} = -\eta \frac{\partial E}{\partial Z_k} f'(P_k) X_i \quad (9)$$

Confrontando questa formula con la precedente formula (3), al posto degli errori noti $(Y_j - D_j)$ troviamo le derivate $\frac{\partial E}{\partial Z_k}$, che dobbiamo ora calcolare retropropagando l'errore dallo strato output a ogni neurone nascosto:

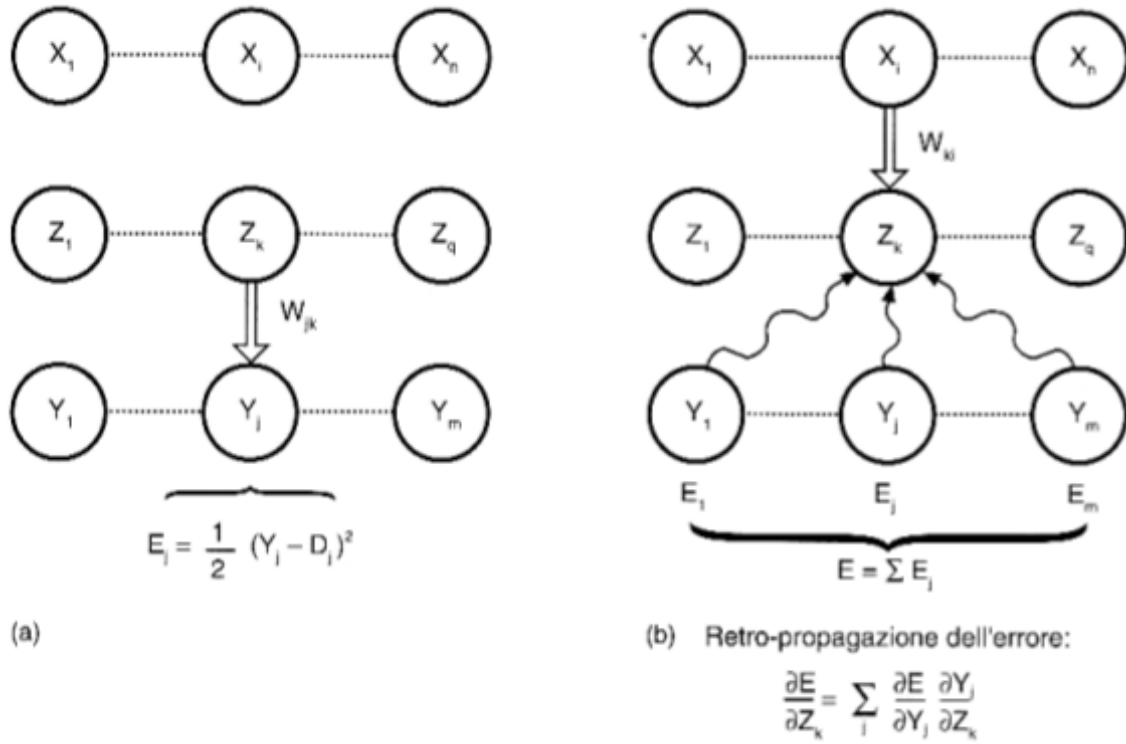
$$\frac{\partial E}{\partial Z_k} = \sum_j \left(\frac{\partial E}{\partial Y_j} \frac{\partial Y_j}{\partial P_j} \frac{\partial P_j}{\partial Z_k} \right) = \sum_j (Y_j - D_j) f'(P_j) W_{jk} \quad (10)$$

o anche, applicando la formula (7) :

$$\frac{\partial E}{\partial Z_k} = \sum_j \delta_j W_{jk} \quad (11)$$

In fine otteniamo le equazioni:

$$\Delta W_{ki} = -\eta f'(P_k) \left(\sum_j \delta_j W_{jk} \right) X_i \quad (12)$$



Propagazione dell'errore nella rete

che ponendo:

$$\delta_k = f'(P_k) \left(\sum_j \delta_j W_{jk} \right) \quad (13)$$

la quale assume la stessa forma della (8) :

$$\Delta W_{ki} = -\eta \delta_k X_i \quad (14)$$

Le formule (12) (13) (14) sono valide non solo per aggiornare i pesi sinattici tra strato nascosto (l'unico nella rete che abbiamo ipotizzato) e strato input ma anche, in reti più generali, per le connessioni tra due strati nascosti consecutivi. Se si adotta la funzione di trasferimento sigmoide si ha:

$$Y = f(P) = \frac{1}{1 + e^{kP}} \quad (15)$$

$$f'(P) = kf(P)(1 - f(P)) = kY(1 - Y) \quad (16)$$

Dove k è proporzionale alla pendenza della sigmoide nel suo punto di flessione, normalmente si pone $k = 1$. A questo punto otteniamo la rielaborazione delle formule (8) (12):

$$\Delta W_{jk} = -\eta k (Y_j - D_j) Y_j (1 - Y_j) Z_k \quad (17)$$

$$\Delta W_{ki} = -\eta k Z_k (1 - Z_k) \left(\sum_j \delta_j W_{jk} \right) X_i \quad (18)$$

infine i pesi sinattici vengono aggiornati, ad ogni tempo t dell'intervallo 1, 2, ..., (t-1), t, (t+1), ...ecc con:

$$W_{jk}(t+1) = W_{jk}(t) + \Delta W_{jk} \quad (19)$$

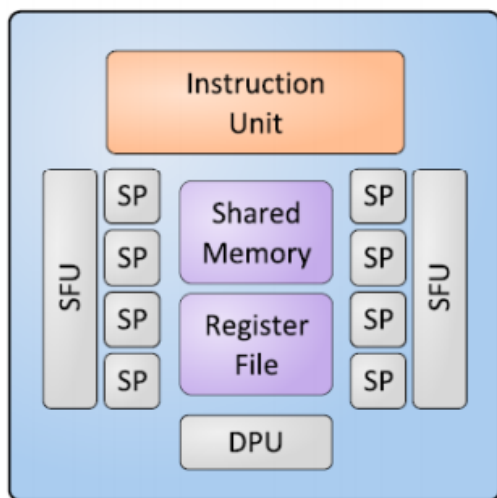
$$W_{ki}(t+1) = W_{ki}(t) + \Delta W_{ki} \quad (20)$$

2.2 Il computo parallelo

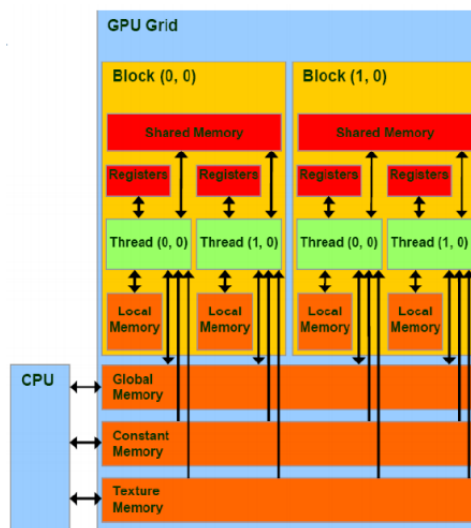
L'incremento di prestazioni dei calcolatori fa sempre più affidamento all'ausilio delle GPU, le quali vengono utilizzate sempre di più per operazioni non prettamente grafiche, le quali presentano un'architettura che punta alla massimizzazione del numero di unità di elaborazione e non più all'incremento della frequenza di lavoro.

Tale approccio non è sempre il migliore, infatti la loro efficacia è strettamente legata al tipo di applicazioni, in particolare alla parallelizzabilità del codice, ovvero alla dipendenza di determinate istruzioni di avere a disposizione i dati elaborati da istruzioni precedenti.

Questa architettura trova la massima efficacia quando un grande numero di istruzioni è completamente indipendente dalle altre.



Struttura di uno stream-multiprocessor



Struttura logica della GPU

2.2.1 L'architettura della GPU

La struttura della GPU è definita "SIMT" (single instruction, multiple thread), tale aggettivo è dato in quanto la GPU è realizzata per eseguire un medesimo programma su un gran numero di processori contemporaneamente, i quali però ricevono dati diversi da computare. Tipicamente i processori presenti su una GPU sono più piccoli e semplici di quelli di una CPU e anche molto meno performanti ma il loro grande numero e la strutturazione delle memorie presenti al suo interno permette per il codice parallelizzabile di ottenere tempi di esecuzione estremamente minori rispetto a quelli in una CPU.

Uno dei componenti più importanti in una GPU è lo stream-multiprocessor, il quale può essere considerato un processore vettoriale, una GPU solitamente ne contiene diversi, e compongono il modello fisico della GPU, i quali a loro volta sono composti principalmente da tre unità di calcolo:

- **SP (scalar processor)** - E' un'unità logica che esegue operazioni in virgola fissa e mobile di un solo thread per volta, è importante inoltre sapere che gli SP di uno stesso SM hanno una memoria detta shared condivisa che gli permette di scambiarsi dati con estrema velocità a tempo di esecuzione.
- **SFU (special functions unit)** - E' un'unità logica di supporto agli SP che esegue a livello HW funzioni complesse come le funzioni goniometriche e trascendenti che appesantirebbero eccessivamente gli SP.
- **DPU (double precision unit)** - Quest'ultima si occupa delle operazioni su numeri a 64 bit, i double.

Una delle peculiarità più importanti della GPU è la differenza tra il modello fisico e quello logico, infatti il modello al quale il programmatore deve far riferimento per ottimizzare l'esecuzione del codice è una via di mezzo tra quello logico e quello che la macchina implementa fisicamente, cercando il bilanciamento migliore possibile tra comprensibilità dei dati e ottimizzazione del codice. Il modello logico. Il programmatore infatti ha la possibilità di interagire e modellare questa struttura composta da oggetti logici detti griglie (GRID), tali griglie possono avere tre dimensioni e a loro volta sono composte da altri oggetti logici detti blocchi, i quali possono avere fino a tre dimensioni, che in fine sono composti da thread. Ogni SM generalmente processa un blocco alla volta, mettendo a disposizione la sua shared memory tra i thread dello stesso blocco, tra i vari blocchi e griglie è possibile lo scambio di informazioni esclusivamente attraverso la global memory, molto più grande delle cache degli SM ma anche molto più lenta. Un'ultimo argomento da trattare riguardo il legame tra la struttura logica e quella hardware della GPU è "il flusso di esecuzione" ovvero il metodo con cui lo scheduler, ovvero l'organo interno che si occupa della ripartizione delle operazioni, strutturate dal programmatore nella griglia, ai vari SM. Tale operazione ha come unità base il warp, ovvero un numero n dipendente dalla specifica architettura, nel nostro caso 32, di thread appartenenti allo stesso blocco che vengono eseguiti contemporaneamente in maniera "lookstepped" (ovvero in sincronizzazione forzata). L'esecuzione del warp si conclude e si passa al successivo solo nel momento in cui tutti i suoi thread arrivano a fine esecuzione indipendentemente dalla "branch" di programma che hanno seguito.

2.2.2 L'utilizzo della GPU per il deep-learning

Nel deep learning come illustrato nella sezione precedente, si effettuano moltissimi calcoli, i quali sono per gran parte scorrelati tra loro, caratteristica questa fondamentale per la parallelizzazione delle operazioni. L'utilizzo di calcolatori paralleli per l'addestramento delle reti neurali consente di completare addestramenti per modelli molto grandi in tempi molto ridotti. Confrontando l'addestramento di una rete su una CPU desktop rispetto ad una GPU da gaming si possono incrementare le prestazioni da 200 a 500 volte, ciò significa che un addestramento di un'ora su GPU, può richiedere su una CPU un tempo che può andare da una a tre settimane.

(Nel nostro caso da un i5-4670k abbiamo avuto un incremento di prestazioni di 250 volte su una GPU Nvidia gtx760 asus, con un programma che pecca fortemente di ottimizzazione).

2.3 Il kinect

Il kinect è un dispositivo per la rilevazione del movimento, prodotto dall'azienda Microsoft e commercializzato insieme alla console videoludica Xbox. Esso può tuttavia essere interfacciato ad un comune pc tramite l'apposito connettore. Nel kinect è presente una fotocamera 1080p RGB, una fotocamera infrarossi per calcolare la profondità e 4 microfoni. Con questo set di sensori ed il suo SDK (software development kit), consente l'utilizzo dei dati prodotti per scopi più generici, una tra le quali l'utilizzo delle gestture per controllare il pc o la ricostruzione delle espressioni facciali. Questi dati sono resi disponibili attraverso l'apposito SDK e possono essere usati in generici programmi C++, C# ed altri linguaggi di programmazione, nello specifico per il linguaggio C++ si usa la libreria "Kinect.h".



Figura 8: Il kinect.

2.3.1 Modello dello scheletro

In particolare ci siamo interessati alla capacità del kinect di poter ricostruire la posizione delle parti del corpo nello spazio, infatti utilizzando la fotocamera infrarossi è in grado di rilevare 24 punti diversi nel corpo come indicato in figura 9.

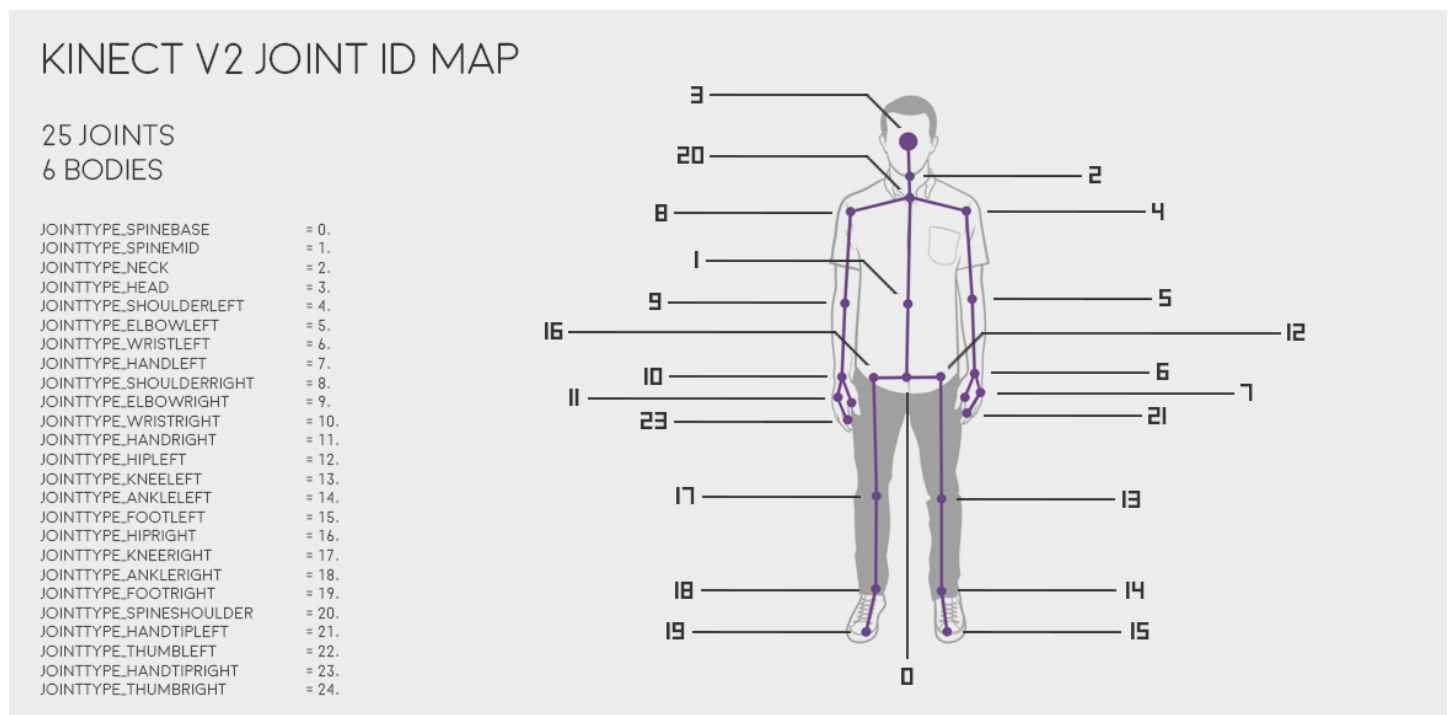


Figura 9: La mappa dei giunti che il kinect è in grado di rilevare.

I dati di posizione e profondità restituiti dal kinect fanno riferimento alla terna di orientazione indicata in figura 10. Per in nostri scopi abbiamo utilizzato delle procedure già definite dalla libreria che permettono di acquisire la posizione nello spazio di ogni singolo giunto dello scheletro e il suo stato di tracciamento, se la sua posizione è stata interpolata o è stato effettivamente tracciato.



Figura 10: Orientazione degli assi rispetto al kinect.

2.4 Scheda embedded

Per questo progetto abbiamo realizzato una scheda embedded che comprendesse l'imu, un modulo bluetooth, la batteria con il suo regolatore e un arduino.

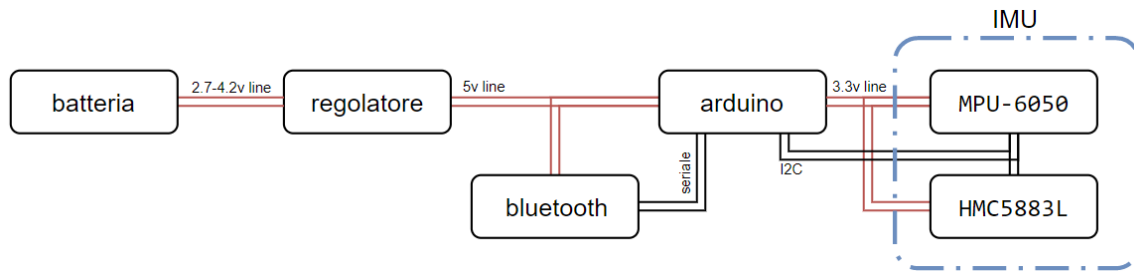


Figura 11: Schema scheda embedded

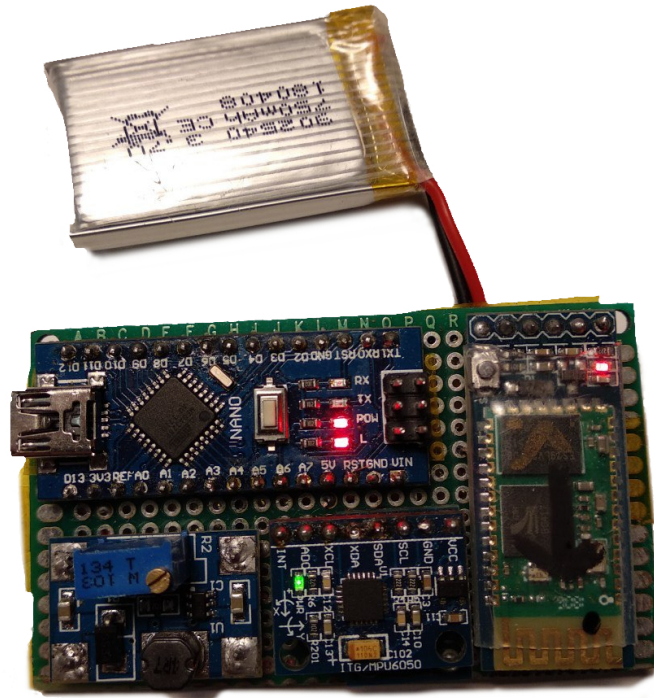


Figura 12: Scheda embedded

La scheda è stata realizzata su millefori saldando e wrappando i diversi componenti.

2.4.1 Arduino

Arduino è una piattaforma hardware open-source dotata di microcontrollore e tutto il suo ecosistema, questo lo rendono estremamente utile per realizzare progetti che non richiedono specifiche particolari senza doveresi realizzare una scheda apposita. Le schede arduino possono essere programmate attraverso l'ide proprietario "Arduino IDE", che oltre ad includere la compatibilità con tutte le schede arduino offre alcuni strumenti utili come il monitor seriale.

Nel nostro progetto abbiamo usato la scheda "Arduino nano" (fig.13), si alimenta a 5 V e include un regolatore di tensione da 5 V a 3.3 V, una seriale, l'I2C ed altri pin general purpose.

Il regolatore interno è stato usato per alimentare i moduli HMC5883L e MPU-6050 che costituiscono l'imu della scheda, la seriale è stata utilizzata per comunicare con il modulo bluetooth HC-05, il bus di comunicazione I2C è stato usato per interfacciarsi con l'imu ed un pin GPIO è stato usato per controllare lo stato del modulo bluetooth.

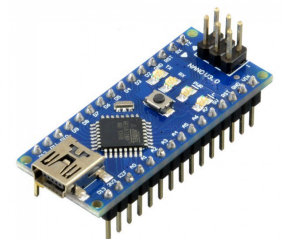


Figura 13: Arduino nano

2.4.2 Bluetooth

Il modulo bluetooth utilizzato è "HC-05", ha 6 pin per interfacciarsi con le altre periferiche: 2 di alimentazione, 2 per la seriale, state e key.

Questo modulo ha un regolatore da 5 V a 3.3 V integrato, quindi va alimentato a 5 V.

Per comunicare con arduino usa la seriale (RS-232) con livello logico 3.3 V, ciò comporta la necessità di inserire un partitore sul pin rx del modulo (quindi il pin tx di arduino).

Il pin state è stato usato per far conoscere ad arduino quando il modulo ha stabilito la connessione con il



Figura 14: HC-05

computer.

Il pin key svolge una funzione particolare poichè permette di avviare il modulo in modalità "command mode", in questa modalità si possono usare gli AT commands per programmare il modulo (es. cambiare nome al dispositivo, oppure cambiare il baud rate della seriale).

I comandi AT e la loro descrizione si trovano sul datasheet del modulo.

Per avviare l'HC-05 in modalità AT si deve collegare il modulo all'alimentazione tenendo premuto il pulsante che si trova sulla scheda, dopodichè si può procedere ad inviare gli at command tramite seriale con un baud rate di 38400Bd. Per fare ciò abbiamo collegato la seriale del modulo con 2 pin di arduino con supporto PWM ed abbiamo usato la libreria software serial per poter avere un'altra seriale (la principale connessa al pc e la secondaria connessa al modulo).

Dopodichè basta programmare arduino in modo che reindirizzi ciò che gli viene scritto dal pc al modulo, e viceversa. il programma usato è il seguente:

```
1 #include <SoftwareSerial.h>
2
3 SoftwareSerial BTSerial(10, 11); // RX | TX
4
5 void setup()
6 {
7   Serial.begin(9600);
8   Serial.println("Enter AT commands:");
9   BTSerial.begin(38400); // HC-05 default speed in AT command mode
10 }
11
12 void loop()
13 {
14   // Keep reading from HC-05 and send to Arduino Serial Monitor
15   if (BTSerial.available())
16     Serial.write(BTSerial.read());
17
18   // Keep reading from Arduino Serial Monitor and send to HC-05
19   if (Serial.available())
20     BTSerial.write(Serial.read());
21 }
```

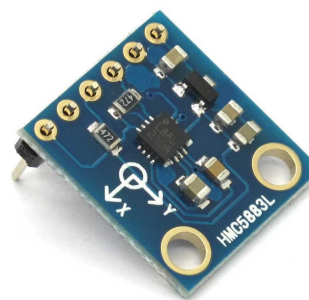
Così facendo abbiamo impostato il baud rate a 115 200 Hz.

2.4.3 L'imu

L'imu (inertial measurement unit) permette di misurare le forze ad esso applicate e l'orientazione dello stesso. Questo viene solitamente fatto combinando i dati di accelerometro, magnetometro e giroscopio. In particolare l'accelerometro misura le accelerazioni, da cui in condizioni di moto inerziale si può estrarre il vettore gravità sui 3 assi determinando quindi l'angolazione rispetto al suolo; Il magnetometro rileva invece il campo magnetico terrestre su 3 assi, dando così indicazione della direzione "nord"; Infine il giroscopio restituisce le accelerazioni angolari. Per questo specifico progetto si sono utilizzati i moduli commerciali "MPU-6050" (fig.15a) e "HMC5883L" (fig.15b), rispettivamente come accelerometro più giroscopio e magnetometro.



(a) MPU-6050



(b) HMC5883L

Figura 15: l'IMU utilizzata in questo progetto

Questi dispositivi comunicano con arduino attraverso il protocollo I2C, sono stati quindi connessi ai pin analogici 5 e 4 di arduino. In particolare l'MPU-6050 integra un accelerometro su 3 assi ed un giroscopio su 3 assi, che vengono convertiti in digitale da 6 ADC a 16 bits. Inoltre può essere programmato su diverse precisioni, il giroscopio tra $\pm 250^\circ/\text{s}$ e $\pm 2000^\circ/\text{s}$, l'accelerometro tra $\pm 2\text{ g}$ e $\pm 16\text{ g}$. Questo modulo ha anche un sensore di temperatura che per questo progetto non è stato usato. Supporta l'I2C fino a 400 kHz.

2.4.4 Regolatore di tensione

Questo regolatore di tensione è del tipo switching e permette di elevare la tensione da 2 V-24 V a 2 V-28 V con un picco massimo di 2 A. La tensione di uscita viene impostata girando il trimmer e qualunque sia la tensione in ingresso l'uscita rimarrà al valore impostato purchè si rispettino i limiti massimi e la tensione di ingresso sia sempre minore di quella di uscita. Nel nostro caso è stato molto utile perchè la batteria a litio ha una tensione che oscilla tra 2.7 V quando è scarica e 4.2 V quando è carica e questo integrato provvede a stabilizzare l'alimentazione a 5 V.



Figura 16: SX1308

3 Programmi

3.1 Overview generale e main server

Introduciamo i programmi con uno schema generale della relazione logico temporale tra i programmi: per prima cosa abbiamo dovuto acquisire un dataset ed usarlo per addestrare la rete neurale

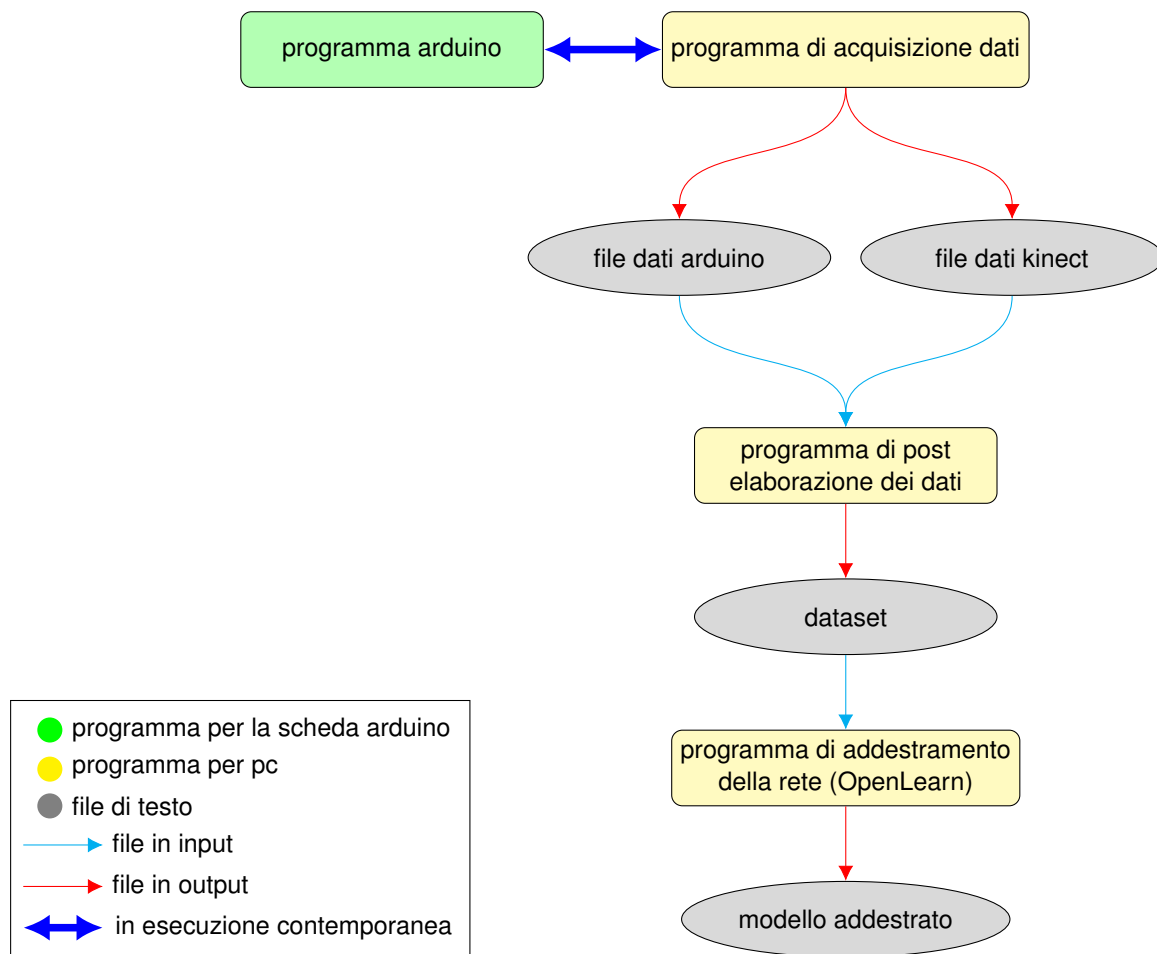


Figura 17: i programmi usati per ottenere i pesi della rete addestrata (?)

una volta addestrata la rete si possono usare i pesi così ottenuti per visualizzare in tempo reale il pupo che si muove(???)

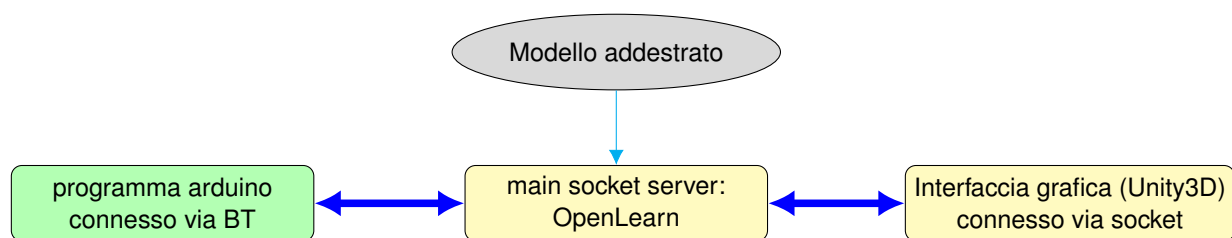


Figura 18: una volta addestrata la rete il manichino viene controllato dal suo output in real-time, controllato dai dati ricevuti via BT da arduino. Nella seguente immagine il data-flow va da sinistra verso destra.

Il programma principale è stato realizzato in C++, con l'ausilio della libreria boost per l'implementazione dei thread e del server socket, in oltre integra la libreria OpenLearn per il caricamento del modello preaddestrato e il processamento dei dati real-time, in tale programma è presente un thread che si occupa parallelamente della ricezione dei dati da arduino, tali dati vengono caricati nel circular buffer il quale è eseguito in un altro thread. I dati contenuti nel circular buffer vengono successivamente presi dal thread che esegue il controllo della rete, la quale elabora i dati che vengono poi immagazzinati e inviati al client unity su richiesta.

3.2 Programma arduino

Questo programma si occupa di acquisire i dati dai sensori attraverso l'I2C, elaborarli e scriverli sulla seriale, in modo che il modulo bluetooth li invii al pc. Il programma si scompone nella classe "MPU0_6050" che gestisce l'accelerometro, nella classe "QMC5883" che gestisce il magnetometro, nella classe "serial" per gestire la seriale e nelle 2 funzioni standard di arduino "setup" e "loop". Inoltre abbiamo usato la libreria "Wire.h" per l'I2C.

Nel "setup" si trova soltanto l'inizializzazione alla seriale per comunicare con il modulo "HC-05", il baud rate impostato è lo stesso che è stato settato nel modulo tramite la modalità AT:

```
1 void setup()
2 {
3   Serial.begin(115200);
4 }
```

Listing 1: funzione "setup"

Nella funzione "loop" si trovano le varie dichiarazioni alle altre classi, un while per aspettare che il modulo "HC-05" sia connesso al pc ed il ciclo infinito che provvede ad acquisire i dati dai sensori e scriverli sulla seriale:

```
1 void loop()
2 {
3   // chiamo i costruttori delle 2 classi che gestiscono i sensori e provvedo al loro settaggio.
4   QMC5883 mm;
5   MPU_6050 aa;
6   aa.setting();
7
8   // accendo un led per il debug
9   digitalWrite(13, HIGH);
10
11  // aspetto finchè il pin "state" del modulo "HC-05" non diventa alto
12  while (!digitalRead(3))
13  {
14    delay(10);
15  }
16
17  // quando il bluetooth è connesso posso procedere ad inizializzare e sincronizzare la seriale
18  serial ser;
19  ser.sinc();
20
21  // spengo il led
22  digitalWrite(13, LOW);
23
24  while (true)
25  {
26    // acquisisco i dati del magnetometro e li normalizzo
27    mm.get_data();
28    mm.normalize(100);
29
30    // acquisisco i dati dell'accelerometro
31    aa.get_data();
32
33    // li sposto in degli array
34    float acc_xyz[] = {aa.get_AcX(), aa.get_AcY(), aa.get_AcZ()};
35    float g_xyz[] = {aa.get_GyX(), aa.get_GyY(), aa.get_GyZ()};
36    float magn[] = {mm.getX(), mm.getY(), mm.getZ()};
37
38    // li scrivo sulla seriale in modo che il modulo bluetooth li trasmetta al pc.
39    ser.send_data(acc_xyz, g_xyz, magn, aa.get_temp());
40  }
41 }
```

Listing 2: funzione "loop"

Qui sono state usate le classi "MPU_6050", "QMC5883" e "serial".

Le classi "MPU_6050" e "QMC5883" si occupano di impostare i registri dei 2 moduli ed inoltre hanno metodi che permettono l'estrazione delle misurazioni da loro effettuate.

L'impostazione dei moduli è stata eseguita facendo riferimento ai rispettivi datasheet.

```
1 class MPU_6050
2 {
```

```

3 public:
4     int16_t AcX, AcY, AcZ, Tmp, GyX, GyY, GyZ;
5     float normalized_AcX, normalized_AcY, normalized_AcZ, normalized_GyX, normalized_GyY,
      ↪ normalized_GyZ;
6     float rescaled_AcX, rescaled_AcY, rescaled_AcZ, rescaled_GyX, rescaled_GyY, rescaled_GyZ;
7
8     const uint8_t MPU = 0x68; // I2C address of the MPU-6050
9
10    // costruttore della classe, accende l'MPU_6050
11    MPU_6050() {...}
12
13    // prede i dati dall'MPU_6050 e le sposta sulle variabili di classe.
14    void get_data() {...}
15
16    // funzioni per ritornare i valori dell'accelerometro
17    float get_AcX() { return (AcX*1.0); }
18    float get_AcY() {...}
19    float get_AcZ() {...}
20
21    // funzioni per ritornare i valori del giroscopio
22    float get_GyX() { return (GyX*1.0); }
23    float get_GyY() {...}
24    float get_GyZ() {...}
25
26    // funzioni per ritornare i valori normalizzati dell'accelerometro
27    float get_normalized_AcX() { return normalized_AcX; }
28    float get_normalized_AcY() {...}
29    float get_normalized_AcZ() {...}
30
31    // funzioni per ritornar i valori normalizzati del giroscopio
32    float get_normalized_GyX() { return normalized_GyX; }
33    float get_normalized_GyY() {...}
34    float get_normalized_GyZ() {...}
35
36    // funzioni per calcolare e ritornare la norma dell'accelerometro o del magnetometro
37    float norma(float x, float y, float z) {...}
38    float get_norma_Gy() { return norma(GyX, GyY, GyZ); }
39    float get_norma_Ac() {...}
40
41    // funzione per ritornare la temperatura
42    float get_temp() { return Tmp; }
43
44    // funzioni per normalizzare l'accelerometro o il magnetometro
45    void normalize_Ac(int Max) {...}
46    void normalize_Gy(int Max) {...}
47
48    // funzione che scrive le impostazioni all'MPU_6050
49    void setting() {...}
50 };

```

Listing 3: classe "MPU_6050"

La seguente è la classe di gestione della bussola 3D

```

1 class QMC5883
2 {
3     public:
4         uint8_t add = 0x0D;
5         int nowX = 0;
6         int nowY = 0;
7         int nowZ = 0;
8         float rescaled_x = 0;
9         float rescaled_y = 0;
10        float rescaled_z = 0;
11
12        // il costruttore provvede ad inizializzare il modulo e settarne i vari registri
13        QMC5883() {...}
14
15        // funzioni per ritornare i valori x-y-z del magnetometro
16        float getX() { return (nowX*1.0); }

```

```

17 float getY(){...}
18 float getZ(){...}
19
20 // funzioni per ritornare i valori del magnetometro dopo averli riscalati
21 float getX_rescaled(){ return rescaled_x; }
22 float getY_rescaled(){...}
23 float getZ_rescaled(){...}
24
25 // funzione che provvede a leggere i regisrti del modulo per acquisire i dati
26 // spostandoli nelle variabili di classe
27 void get_data(){...}
28
29 // funzioni per calcolare la norma e ritornare i valori normalizzati
30 float norma(float x, float y, float z){...}
31 float get_norma(){ return norma(nowX, nowY, nowZ); }
32 void normalize(int Max){...}
33 };

```

Listing 4: classe "QMC5883"

Infine la classe serial si occupa di scomporre i dati di tipo float restituiti dai moduli in array di 4 byte che poi vengono trasmessi al modulo bluetooth tramite la seriale.

```

1 class serial
2 {
3     public:
4
5     // funzione che permette di sincronizzarsi con il pc
6     void sinc(){...}
7
8     // union per poter scompattare un float in 4 byte
9     union Scomp_float{...};
10
11     // funzione che permette di inviare un intero float scompattandolo in
12     // 4 byte che vengono trasmessi serialmente
13     void send_float(float n){...}
14
15     // funzione che permette di riceve un float come sopra
16     float receive_float(){...}
17
18     // funzione per l'invio di un solo carattere
19     void send_char(char ch){...}
20
21     // funzione per la ricezione di un singolo carattere
22     char receive_char(){...}
23
24     // funzione che provvede ad inviare i dati necessari.
25     void send_data(float* acc_xyz, float* g_xyz, float* magn, float temp){...}
26 };

```

Listing 5: classe "serial"

Per effettuare l'operazione di serializaziopne dei float in array di byte abbiamo fatto ricorso all'uso di una struttura dati chiamata "union":

```

1 union Scomp_float
2 {
3     float f;
4     int i;
5     unsigned char byte_s[4];
6 };

```

Listing 6: classe "serial"

Questa struttura dati ha la particolarità che a differenza di una struct, in cui i dati sono memorizzati in modo contiguo nella memoria, l'union memorizza tutti i dati della struttura nelle stesse locazioni di memoria.

```

1 struct es_struct
2 {
3     int a;
4     unsigned char b[4];
5 }

```

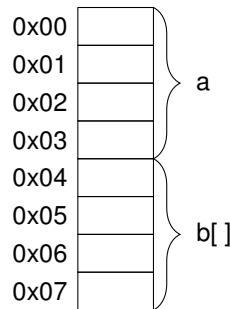
(a) Struct generico

```

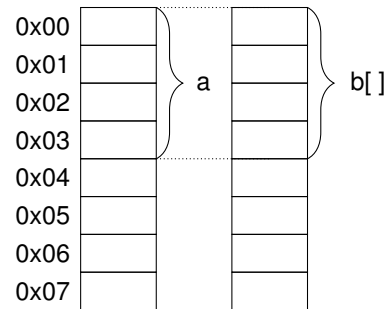
1 union es_union
2 {
3     int a;
4     unsigned char b[4];
5 }

```

(c) Union generico



(b) Rappresentazione dei dati dello struct memorizzati in ram



(d) Rappresentazione dei dati dell'union memorizzati in ram

Figura 19: Confronto tra struct ed union

Come si può vedere dalla figura 19 lo struct (figura 19a) memorizza i dati in ram facendo partire l'intero dall'indirizzo 0x00 e l'array di byte dall'indirizzo 0x04 (figura 19c), mentre l'union (figura 19d) sovrappone i 2 dati nelle stesse locazioni di memoria (figura 19). Ciò permette ad esempio di leggere un intero (o un qualunque altro tipo di dato) come un insieme di byte: con riferimento alla figura 19d se pongo:

```

1 es_union U;
2 U.a = 0x0514FFAA;

```

Listing 7: classe "serial"

e successivamente vado a vedere il valore di b[] ottengo i valori:

```

1 U.b[0] -> 0x05
2 U.b[1] -> 0x14
3 U.b[2] -> 0xFF
4 U.b[3] -> 0xAA

```

Listing 8: classe "serial"

che sono esattamente i byte costituenti U.a. Ciò ci ha permesso di inviare un float come sequenza di byte al pc, nella classe serial sul pc c'è esattamente lo stesso tipo di struttura per ricostruire il float.

3.3 programma per l'acquisizione dei dati dal kinect e da arduino

Questo programma è stato usato per acquisire i dati da arduino e dal kinect e scriverli in dei file di testo, ce vengono successivamente elaborati.

3.3.1 librerie usate

La lista completa delle librerie usate è:

```

1 // libreria usata da visual studio, da togliere in caso si usi un altro ide
2 #include "pch.h"
3
4 // librerie standard per i file stringhe e altro
5 #include <cstdlib>
6 #include <stdlib.h>
7 #include <fstream>
8 #include <iomanip>
9 #include <iostream>
10 #include <ostream>
11 #include <string>
12 #include <math.h>
13
14 // librerie create da noi
15 #include "real_time.h"

```



```

16 #include "kin_file_manager.h"
17 #include "ard_file_manager.h"
18 #include "data_structure.h"
19
20 // libreria per usare i thread
21 #include <boost/thread.hpp>
22
23 // altre librerie per log, alcuni tipi di dato, un buffer FIFO
24 #include <boost/log/core.hpp>
25 #include <boost/call_traits.hpp>
26 #include <boost/circular_buffer.hpp>
27 #include <boost/container/vector.hpp>
28
29 // altre librerie per il logger
30 #include <boost/log/attributes.hpp>
31 #include <boost/log/attributes/scoped_attribute.hpp>
32 #include <boost/log/expressions.hpp>
33 #include <boost/log/sinks/sync_frontend.hpp>
34 #include <boost/log/sinks/text_ostream_backend.hpp>
35 #include <boost/log/sources/basic_logger.hpp>
36 #include <boost/log/sources/record_ostream.hpp>
37 #include <boost/log/sources/severity_logger.hpp>
38 #include <boost/log/utility/setup/common_attributes.hpp>
39 #include <boost/log/utility/setup/console.hpp>
40 #include <boost/log/utility/setup/file.hpp>
41
42 #include <boost/smart_ptr/make_shared_object.hpp>
43 #include <boost/smart_ptr/shared_ptr.hpp>
44
45 // altre librerie per i thread
46 #include <boost/thread/condition_variable.hpp>
47 #include <boost/thread/mutex.hpp>
48 #include <boost/thread/thread.hpp>
49
50 // libreria per il tempo
51 #include <boost/date_time/posix_time/posix_time.hpp>
52
53 // libreria fatta da noi per gestire la seriale
54 #include "Serial_handler.h"
55
56 // librerie di Windows e kinect
57 #include <Windows.h>
58 #include <Kinect.h>
59 #include <Shlobj.h>
60
61 // namespace standard
62 using namespace std;

```

Listing 9: librerie usate

3.3.2 Librerie custom

Per alleggerire il file "main" abbiamo spostato delle classi in file esterni, che sono a tutti gli effetti delle librerie a se stanti. Abbiamo costruito librerie per: la gestione del tempo, i file manager, le strutture dati e la gestione della seriale.

3.3.2.1 Libreria per misurare il tempo

In "real_time.h" ci sono alcune funzioni per prendere il tempo e misurarlo:

```

1 #pragma once
2
3 #ifndef _REAL_TIME_H_
4 #define _REAL_TIME_H_
5
6 #include <boost/chrono.hpp>
7 #include <boost/timer/timer.hpp>
8
9 class real_time
10 {

```

```

11 private:
12     boost::timer::cpu_timer timer;
13     boost::timer::cpu_times t;
14
15 public:
16     /// <summary>
17     /// start the time counting
18     /// </summary>
19     /// <returns></returns>
20     void start();
21
22     /// <summary>
23     /// return elapsed time in millisecond
24     /// </summary>
25     /// <returns></returns>
26     float stop();
27
28     /// <summary>
29     /// get the time in millisecond
30     /// </summary>
31     /// <returns></returns>
32     uint64_t get_curr_time();
33
34 };
35
36 #endif // #ifndef _REAL_TIME_H_

```

Listing 10: librerie usate

3.3.2.2 File manager per il kinect

In "kin_file_manager.h" ci sono le funzioni per gestire il file di testo relativo al kinect:

```

1 #pragma once
2
3 #ifndef _KIN_FILE_MANAGER_H_
4 #define _KIN_FILE_MANAGER_H_
5
6 #include <iostream>
7 #include <ostream>
8 #include <fstream>
9 #include <string>
10
11 // in questa libreria sono dichiarate le varie strutture dati
12 #include "data_structure.h"
13
14 class kin_file_manager
15 {
16 private:
17     std::fstream f;
18     std::ios::_Openmode mode;
19     unsigned long int line_id = 0;
20
21 public:
22     // il costruttore provvede a settare il nome del file e la modalit\`a di apertura
23     kin_file_manager(std::string file_name, std::ios::_Openmode mode);
24
25     // questa funzione provvede a scrivere un singolo blocco di dati nel file
26     void write_data_line(kinect_data dat);
27
28     // funzione che legge un blocco di dati dal file e lo sposta nella struttura dati
29     // se si \`e raggiunta la fine del file la funzione ritorna 0, altrimenti ritorna 1.
30     bool read_data_line(kinect_data* dat);
31
32     // distruttore della classe, chiude il file
33     ~kin_file_manager()
34     {
35         f.close();
36     }

```

```

37
38 // funzione che chiude il file
39 void close()
40 {
41     f.close();
42 }
43
44 };
45
46 #endif // #ifndef _KIN_FILE_MANAGER_H_

```

Listing 11: librerie usate

3.3.2.3 File manager per arduino

La libreria "ard_file_manager.h" che contiene le funzioni per la gestione del file di testo relativo ai dati di arduino:

```

1 #pragma once
2
3 #ifndef _ARD_FILE_MANAGER_
4 #define _ARD_FILE_MANAGER_
5
6 #include <iostream>
7 #include <ostream>
8 #include <fstream>
9 #include <string>
10 #include "data_structure.h"
11
12
13 class ard_file_manager
14 {
15 private:
16     std::fstream f;
17     std::ios::_Openmode mode;
18     unsigned long int line_id = 0;
19
20 public:
21     ard_file_manager(std::string file_name, std::ios::_Openmode mode);
22
23     void write_data_line(arduino_data dat);
24
25     bool read_data_line(arduino_data* dat);
26
27     ~ard_file_manager()
28     {
29         f.close();
30     }
31
32     void close()
33     {
34         f.close();
35     }
36 };
37
38 #endif // #ifndef _ARD_FILE_MANAGER_

```

Listing 12: librerie usate

Questa libreria è strutturata in modo totalmente simile a quella per la gestione dei file del kinect trattata nel paragrafo 3.3.2.2 con l'unica differenza che la struttura dati utilizzata è relativa ad arduino e non al kinect.

3.3.2.4 strutture dati

La libreria "data_structure.h" è fondamentale e definisce le strutture dati utilizzate nel resto del programma:

```

1 #pragma once
2
3

```

[illegible]

```

68 float acc_xyz[3];
69 float gy_xyz[3];
70 float magn_xyz[3];
71 float temp;
72
73 // tempo esatto di acquisizione dei dati
74 uint64_t frame_time = 0;
75
76 // ?
77 uint64_t contatore = 0;
78
79 // funzione che permette di stampare il set di dati attualmente immagazzinato a schermo
80 void print_data();
81 };
82
83 // struttura template che gestisce un set di dati del dataset
84 template<typename type_in, std::size_t N, typename type_out, std::size_t M>
85 struct dataset_data
86 {
87     boost::array<type_in, N> in;
88     boost::array<type_out, M> out;
89
90     // stampa i dati attualmente immagazzinati
91     void print_data();
92 };
93
94 #endif // #ifndef _DATA_STRUCTURE_H_

```

Listing 13: librerie usate

3.3.2.5 Gestione seriale

Infine, l'ultima libreria fatta da noi ospita la classe che gestisce la seriale per la comunicazione tramite bluetooth con arduino:

```

1 #pragma once
2
3 #ifndef _SERIAL_HADLER_H_
4 #define _SERIAL_HADLER_H_
5
6 #include <iostream>
7 #include <string>
8 #include <boost/asio.hpp>
9 #include <boost/bind.hpp>
10 #include <boost/asio/serial_port.hpp>
11
12
13 class Serial
14 {
15 private:
16     // dichiaro la porta seriale che verra inizializzata nel costruttore
17     boost::asio::serial_port* port;
18
19     // per ricomporre i float ricevuti uso un union
20     union Scomp_float
21     {
22         float n_float;
23         uint32_t n_int;
24         uint8_t n_bytes[4];
25     };
26
27 public:
28     // nel costruttore si inizializza la seriale e si effettua la connessione al dispositivo
29     // ↳ bluetooth
30     Serial(std::string com);
31
32     // effettua una sincronizzazione iniziale tra arduino ed il pc
33     void sinc();
34
35     // riceve un float da arduino

```

```

35 float receive_float();
36
37 // metodo per inviare un carattere
38 void send_char(char ch);
39
40 // metodo per ricevere un carattere
41 char receive_char();
42
43
44 // riceve i dati di acc, gy, magn da arduino
45 // il metodo ritorna
46 //     -> 0 se la trasmissione \e andata a buon fine
47 //     -> 1 se \e fallita
48 int receive_data(float* acc_xyz, float* g_xyz, float* magn, float* temp);
49
50 ~Serial()
51 {
52     port->close();
53 }
54 };
55
56 #endif // #ifndef _SERIAL_HADLER_H_

```

Listing 14: librerie usate

Per trasmettere il blocco di dati da arduino al pc abbiamo sviluppato un "protocollo" di comunicazione, cioè i vari dati sono separati da sei tag di controllo che permettono di verificare se il dato è stato trasmesso correttamente o no.

```

1 void send_data(float* acc_xyz, float
    ↳ * g_xyz, float* magn, float
    ↳ temp)
2 {
3     send_char('a');
4     for (int i = 0; i < 3; i++)
5         send_float(acc_xyz[i]);
6     send_char('g');
7     for (int i = 0; i < 3; i++)
8         send_float(g_xyz[i]);
9     send_char('m');
10    for (int i = 0; i < 3; i++)
11        send_float(magn[i]);
12    send_char('t');
13    send_float(temp);
14    send_char(4);
15 }

```

(a) codice arduino

```

1 int Serial::receive_data(float*
    ↳ acc_xyz, float* g_xyz, float*
    ↳ magn, float* temp)
2 {
3     if (receive_char() != 'a')
4         return 1;
5     for (int i = 0; i < 3; i++)
6         acc_xyz[i] = receive_float();
7     if (receive_char() != 'g')
8         return 1;
9     for (int i = 0; i < 3; i++)
10        g_xyz[i] = receive_float();
11    if (receive_char() != 'm')
12        return 1;
13    for (int i = 0; i < 3; i++)
14        magn[i] = receive_float();
15    if (receive_char() != 't')
16        return 1;
17    temp[0] = receive_float();
18    if (receive_char() != 4)
19        return 1;
20    return 0;
21 }

```

(b) codice che gira sul pc

Figura 20: codice per l'invio e la ricezione dei dati necessari

Nel codice in figura 20a troviamo la funzione per inviare i dati da arduino, mentre in figura 20b troviamo il codice per ricevere gli stessi dati sul pc. Come si può vedere viene inviato un carattere di controllo prima di ogni dato, questo carattere serve per controllare l'integrità della comunicazione e deve essere ricevuto dal pc. Se ciò non accade significa che nella trasmissione è andato perso 1 byte e quindi arduino ed il pc vanno risincronizzati.

3.3.3 programma principale dell'acquisizione dati

Oltre a queste classi/strutture integrate in delle librerie separate ci sono altre classi/altro nel main, infatti abbiamo: "namespace logger", "class bounded_buffer", "class kinect_class", "class arduino_class", "class ard_handler", "class kin_handler", "void start_acquire_data" e "void acquire_data".

3.3.3.1 Logger

Per prima cosa è bene spiegare il logger, che viene usato nel resto del programma. Questa libreria serve esclusivamente per il debug e ha attributi specifici per lavorare con più thread contemporaneamente in modo asincrono. Serve fondamentalmente a stampare un testo a schermo, ma una volta settate le impostazioni si può aggiungere a questo testo il nome del thread il tempo esatto in cui il messaggio viene stampato ed evitare che 2 thread stampino contemporaneamente nella console (se lo facessero il testo si mischierebbe e non si capirebbe nulla). La documentazione approfondita su come funziona questa libreria si può trovare qui. Nel nostro caso i settaggi usati sono: un id che identifica il thread, un tag che identifica il nome del thread, un tag che identifica il livello di severità, il tempo assoluto e la stringa da stampare. Ad esempio se questa riga è posta all'inizio del main e il main thread viene chiamato "main":

```
1 BOOST_LOG_SEV(slg, logger::normal) << "hello log";
```

Listing 15: librerie usate

Allora viene stampato a schermo un messaggio formattato in questo modo:

```
1 id->(0x24a8): main < normal >   time->[00:00:00.000404]  -> hello log
```

Listing 16: librerie usate

Ciò consente di debuggare facilmente operazioni asincrone su più thread. Il codice relativo hai settaggi necessari per ottenere quel tipo di formattazione è il "namespace logger":

```
1 namespace logger
2 {
3     // definizione dei diversi livelli di severità
4     enum severity_level
5     {
6         normal,
7         warning,
8         error,
9         critical
10    };
11
12    // definizione degli attributi usati
13    BOOST_LOG_ATTRIBUTE_KEYWORD(severity, "Severity", severity_level)
14    BOOST_LOG_ATTRIBUTE_KEYWORD(tag_attr, "Tag", std::string)
15    BOOST_LOG_ATTRIBUTE_KEYWORD(timeline, "Timeline", boost::log::attributes::timer::value_type)
16    BOOST_LOG_ATTRIBUTE_KEYWORD(tread_id, "Tread_id", boost::thread::id)
17    BOOST_LOG_ATTRIBUTE_KEYWORD(tread_name, "tread_name", std::string)
18
19    // provvede a sostituire il numero del severity_level con la relativa stringa
20    std::ostream& operator<< (std::ostream& strm, severity_level level){...}
21
22    // inizializzazione della formattazione del logger
23    void init(){...}
24
25    // definizione dei vari parametri
26    boost::log::sources::severity_logger< severity_level > f_init(){...}
27    void attr_thread(...){...}
28    void attr_time(...){...}
29    void attr_tag(...){...}
30    void attr_thread_name(...){...}
31};
```

Listing 17: librerie usate

Inoltre per funzionare correttamente richiede che all'inizio del main siano posti alcuni settaggi:

```
1 // inizializzazione logger
2 logger::init();
3 auto slg = logger::f_init();
4 attr_thread(&slg);
5 attr_time(&slg);
6
```



```

7 // imposto il nome del main thread in "main"
8 logger::attr_thread_name(&slg, "main");

```

Listing 18: librerie usate

3.3.3.2 Buffer FIFO

Un'altra classe usata nel resto del programma è "class bounded_buffer", questa classe ospita funzioni specifiche per immagazzinare temporaneamente dei dati gestendone l'ingresso e l'uscita dal buffer. Fondamentalmente costituisce un registro FIFO indipendente dal tipo di dato che si utilizza. Questo tipo di buffer viene comunemente usato in caso si ha la struttura thread che produce dati - thread che consuma dati: il produttore acquisisce i dati e li inserisce nel buffer, il consumatore estrae i dati e li processa. La dimensione del buffer viene specificata quando lo si inizializza e deve essere sufficientemente grande in modo da garantire la non perdita di dati e il non rallentamento del thread produttore. Nel nostro caso abbiamo adottato questo tipo di buffer per poter garantire l'asincronicità delle acquisizioni tra arduino e il kinect, così facendo si garantisce la massima velocità tra le acquisizioni dei dati. Il codice relativo a questa classe è stato quasi interamente preso dagli esempi della libreria boost, qui ed è:

```

1 template <class T>
2 class bounded_buffer
3 {
4 public:
5
6     typedef boost::circular_buffer<T> container_type;
7     typedef typename container_type::size_type size_type;
8     typedef typename container_type::value_type value_type;
9     typedef typename boost::call_traits<value_type>::param_type param_type;
10
11     // nel costruttore si imposta la capacità del buffer
12     explicit bounded_buffer(size_type capacity) : m_unread(0), m_container(capacity) {}
13
14     // questa funzione permette di aggiungere un elemento al buffer
15     void push_front(param_type item){...}
16
17     // questa funzione permette di estrarre un elemento dal buffer
18     void pop_back(value_type* pItem){...}
19
20     // restituisce la percentuale di riempimento del buffer
21     float percentage_of_filling(){...}
22
23     // svuota il buffer
24     void flush(){...}
25
26 private:
27     bounded_buffer(const bounded_buffer&); // Disabled copy constructor
28     bounded_buffer& operator = (const bounded_buffer&); // Disabled assign operator
29
30     bool is_not_empty() const { return m_unread > 0; }
31
32     bool is_not_full() const { return m_unread < m_container.capacity(); }
33
34     size_type m_unread;
35     container_type m_container;
36     boost::mutex m_mutex;
37     boost::condition_variable m_not_empty;
38     boost::condition_variable m_not_full;
39 };

```

Listing 19: librerie usate

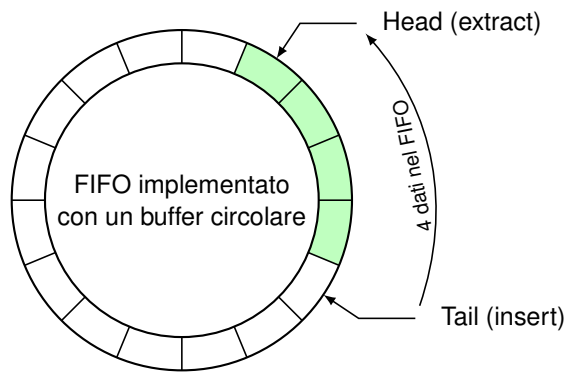


Figura 21: Schema di funzionamento circular buffer

La classe nel listato 19 utilizza un buffer circolare per implementare il registro FIFO, il funzionamento di questo buffer è descritto dalla figura 21: ogni volta che si aggiunge un dato viene occupata una locazione in testa al buffer ed ogni volta che si rimuove un dato viene liberata una locazione in coda, cioè permette di riusare continuamente le stesse locazioni di memoria.

3.3.3.3 Classe gestione kinect

La classe "kinect_class" provvede a leggere i dati dal kinect ed aggiungerli al suo corrispettivo registro FIFO. L'estrazione dei dati dal kinect è stata realizzata sulla base dell'esempio body basic D2D integrato nell'sdk browser fornito da microsoft. Il codice relativo a questa classe è:

```

1 class kinect_class
2 {
3 public:
4
5     // Current Kinect
6     IKinectSensor* m_pKinectSensor;
7     ICoordinateMapper* m_pCoordinateMapper;
8
9     // Body reader
10    IBodyFrameReader* m_pBodyFrameReader;
11
12    real_time t;
13
14    bounded_buffer<kinect_data>* FIFO_kin;
15
16    kinect_class(bounded_buffer<kinect_data>* data_FIFO) {...}
17
18    void start()
19    {
20
21        bool flag = false;
22        while (true)
23        {
24            {...}
25
26            // prendo il tempo
27            uint64_t t1 = t.get_curr_time();
28
29            // acquisisco i dati e li sposto in data
30            flag = Update(&data);
31
32            // se l'acquisizione \e andata a buon fine procedo
33            if (flag)
34            {
35                // prendo un'altra volta il tempo
36                uint64_t t2 = t.get_curr_time();
37                // cos'i da poter impostare che il tempo in cui sono stati acquisiti i dati sia
38                // la media tra il tempo prima di cominciare l'acquisizione ed il tempo dopo che essa \
39                // e finita
40                data.frame_time = (t1 + t2) / 2;
41
42                // inserisco i dati cos'i acquisiti nel registro fifo

```

```

42     FIFO_kin->push_front(data);
43
44     // faccio aspettare un po di tempo a questo trhead perch'\e altrimenti il kinect da
45         ↳ errore,
46     // dato che non supporta un acquisizione a pi'\u di 30fps
47     boost::this_thread::sleep(boost::posix_time::milliseconds(35)); al kinect scoppia
48 }
49
50 {...}
51
52 }
53 }
54
55 // inizializza il kinect, ritorna 1 se ha avuto successo, 0 altrimenti
56 int InitializeDefaultSensor(){...}
57
58 // acquisisce i dati e li passa alla funzione per processarli
59 bool Update(kinect_data* data)
60 {
61     {...}
62
63     flag = ProcessBody(BODY_COUNT, ppBodies, data);
64
65     {...}
66
67     return flag;
68 }
69
70
71 // questa funzione provvede a spostare i dati da joints a data sistemandoli nella loro
72     ↳ struttura dati,
73 // una spiegazione di come ci\o viene fatto verra data nel capitolo successivo
74 bool handle_data(Joint* joints, kinect_data* data){...}
75
76 // questa funzione separa le diverse persone rilevate e trasferisce le informazioni
77 // relative ai giunti della prima persona alla funzione handle_data
78 bool ProcessBody(int nBodyCount, IBody** ppBodies, kinect_data* data)
79 {
80     {...}
81
82     for (int i = 0; i < nBodyCount; ++i)
83     {
84         {...}
85
86         Joint joints[JointType_Count];
87
88         hr = pBody->GetJoints(_countof(joints), joints);
89
90         return handle_data(joints, data);
91
92         {...}
93     }
94
95     {...}
96 }
97
98
99 // provvede a spegnere il kinect e deallocare le sue risorse nel caso la classe venga
100     ↳ distrutta
101 ~kinect_class(){...}
102 template<class Interface>
103 inline void SafeRelease(Interface *& pInterfaceToRelease){...}
104 };

```

Listing 20: librerie usate

3.3.3.4 Classe gestione arduino

La classe "arduino_class" provvede a caricare continuamente i dati provenienti dal kineck nel registro fifo, infatti la funzione "start" viene lanciata come thread e continua a caricare i dati fino al termine del programma:

```
1 class arduino_class
2 {
3 public:
4     Serial* ser;
5     real_time t;
6     bounded_buffer<arduino_data>* FIFO_ard;
7
8     // il costruttore provvede ad inizializzare la seriale e sincronizzare i 2 dispositivi
9     arduino_class(bounded_buffer<arduino_data>* data_FIFO, string com_port)
10    {
11        ser = new Serial(com_port);
12        cout << "serial ok" << endl;
13        FIFO_ard = data_FIFO;
14        ser->sinc();
15    }
16
17    // questa funzione viene lanciata all'avvio del thread e provvede ad acquisire i dati ed
18    // ↳ inserirli nel registro fifo
19    void start()
20    {
21        while (true)
22        {
23            arduino_data data;
24
25            uint64_t t1 = t.get_curr_time();
26            ser->receive_data(data.acc_xyz, data.gy_xyz, data.magn_xyz, &data.temp);
27            uint64_t t2 = t.get_curr_time();
28
29            // per ottenere un accurata misurazione del tempo viene preso prima e dopo l'acquisizione e
30            // ↳ poi fatta la media dei valori ottenuti
31            data.frame_time = (t1 + t2) / 2;
32
33            FIFO_ard->push_front(data);
34        }
35    }
36};
```

Listing 21: librerie usate

3.3.3.5 Gestore di arduino e del kinect

Dopo i 2 thread producer si hanno i 2 thread consumer: "class ard_handler" e "class kin_handler", che si occupano rispettivamente di salvare su file di testo le informazioni di arduino e del kinect.

```
1 /// <summary>
2 /// handle the data from arduino and write it on a file
3 /// </summary>
4 class ard_handler
5 {
6 private:
7     bounded_buffer<arduino_data>* ard;
8     ard_file_manager* f;
9
10 public:
11     // il costruttore provvede ad inizializzare la classe per gestire i file di testo di arduino
12     ard_handler(bounded_buffer<arduino_data>* ard, string file_name)
13     {
14         this->ard = ard;
15         f = new ard_file_manager(file_name, std::ios::out);
16     }
17
18     // questa funzione \e quella che viene effettivamente lanciata all'avvio del thread e
19     // ↳ provvede a scaricare il buffer scrivendo i dati sul file fino al raggiungimento della
20     // ↳ quota predichiarata "tot_esempi"
```

```

19 void start()
20 {
21     // inizializzazione logger
22     auto slg = logger::f_init();
23     attr_thread(&slg);
24     logger::attr_thread_name(&slg, "ard_handler");
25
26     BOOST_LOG_SEV(slg, logger::normal) << "hello log";
27
28     for (int i = 0; i < tot_esempi; i++)
29     {
30         arduino_data data_ard;
31         ard->pop_back(&data_ard);
32
33         f->write_data_line(data_ard);
34     }
35 }
36
37 // segnale sulla console quando il thread ha finito il suo lavoro
38 BOOST_LOG_SEV(slg, logger::normal) << "end";
39 }
40 };

```

Listing 22: librerie usate

La classe "kin_handler" esegue esattamente le stesse operazioni della precedente con la sola differenza che usa "kin_file_manager" al posto di "ard_file_manager".

```

1  /// <summary>
2  /// handle the data from the kinect and write it in the file
3  /// </summary>
4  class kin_handler
5  {
6  private:
7      bounded_buffer<kinect_data>* kin;
8      kin_file_manager* f;
9
10 public:
11     kin_handler(bounded_buffer<kinect_data>* kin, string file_name){...}
12
13     void start(){...}
14 };

```

Listing 23: librerie usate

3.3.3.6 Funzioni di lancio dei thread per l'acquisizione ed il salvataggio dei dati

Dopo tutte queste classi si hanno 2 funzioni che si occupano di inizializzare le altre classi, lanciare tutti i thread necessari ed acquisire dall'utente i nomi dei file etc.

```

1 void start_acquire_data(string out_ard, string out_kin, string com_port)
2 {
3     // in inizializzo il logger per questa funzione
4     auto slg = logger::f_init();
5     attr_thread(&slg);
6     logger::attr_thread_name(&slg, "main");
7
8     // dichiaro i 2 buffer FIFO
9     bounded_buffer<arduino_data> data_FIFO_ard(20);
10    bounded_buffer<kinect_data> data_FIFO_kin(20);
11
12    // chiamo i costruttori delle classi che gestiscono il kinect e arduino
13    // al costruttore di arduino gli devo passare la com a cui il bluetooth \e collegato
14    arduino_class a(&data_FIFO_ard, com_port);
15    kinect_class k(&data_FIFO_kin);
16
17    // chiamo i costruttori delle classi per la scrittura nei file di testo
18    ard_handler a_h(&data_FIFO_ard, out_ard);
19    kin_handler k_h(&data_FIFO_kin, out_kin);

```

```

20
21 BOOST_LOG_SEV(slg, logger::normal) << "costructor end" << std::endl;
22
23 // lancio i thread che "producono i dati"
24 boost::thread ard([&a]() { a.start(); });
25 boost::thread kin([&k]() { k.start(); });
26
27 BOOST_LOG_SEV(slg, logger::normal) << "producer thread launced" << std::endl;
28
29 // lancio i thread che "consumano dati"
30 boost::thread ard_h([&a_h]() { a_h.start(); });
31 boost::thread kin_h([&k_h]() { k_h.start(); });
32
33 BOOST_LOG_SEV(slg, logger::normal) << "consumer thread launched" << std::endl;
34
35 // aspetto all'infinito (finch\`e il programma non viene chiuso
36 kin.join();
37 ard.join();
38 kin_h.join();
39 ard_h.join();
40 }

```

Listing 24: librerie usate

infine c'è la funzione che si occupa di acquisire i dati dall'utente e passarli alla funzione appena vista.

```

1 void acquire_data()
2 {
3     string ard_file_mane = "";
4     string kin_file_mane = "";
5     string arduino_com_port = "";
6
7     // acquisisco il nome dei 2 file che conterranno i dati di arduino e i dati del kinect
8     cout << "insert the file names: " << endl << "\t arduino -> ";
9     std::getline(std::cin, ard_file_mane);
10    cout << "\t kinect -> ";
11    std::getline(std::cin, kin_file_mane);
12
13    // acquisisco il numero della com a cui \`e collegato il bluetooth del computer
14    cout << "insert the arduino com port -> : ";
15    std::getline(std::cin, arduino_com_port);
16
17    // lancio la funzione passandogli i dati appena acquisiti
18    start_aquire_data(ard_file_mane, kin_file_mane, arduino_com_port);
19 }

```

Listing 25: librerie usate

3.3.3.7 Main

Infine il main chiama semplicemente la funzione appena trattata.

```

1 int main()
2 {
3     logger::init();
4
5     auto slg = logger::f_init();
6     attr_thread(&slg);
7     attr_time(&slg);
8     logger::attr_thread_name(&slg, "main");
9
10    BOOST_LOG_SEV(slg, logger::normal) << "hello log";
11
12    acquire_data();
13
14    return 0;
15 }

```

Listing 26: librerie usate

3.3.4 Riassunto

Riassumendo la struttura generale di questo programma si ha:

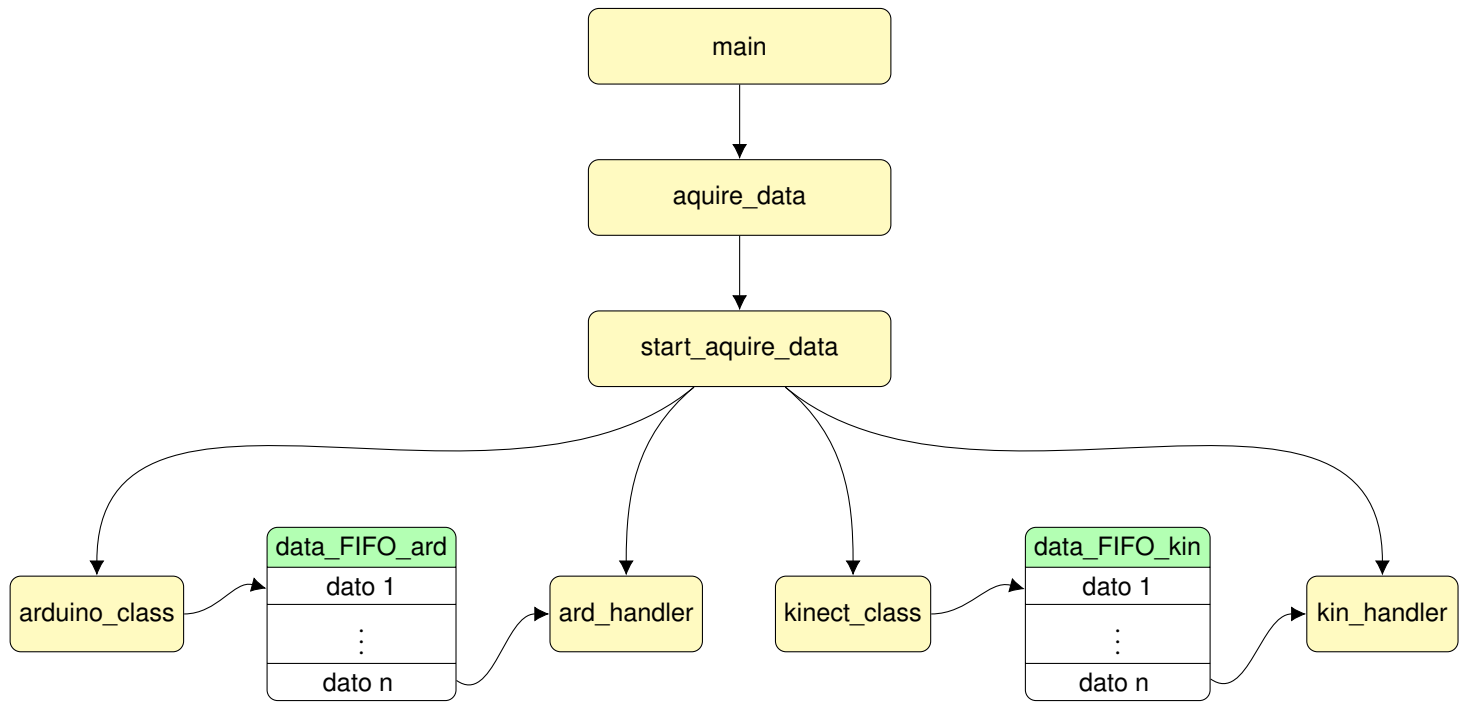


Figura 22: programmi che si occupano dell'immagazzinamento dei dati per la creazione del dataset

Cioè il "main" chiama "aquire_data" che chiama a sua volta "start_aquire_data". Dopodichè vengono lanciati i 4 thread ("arduino_class", "ard_handler", "kinect_class", "kin_handler"). I thread che si occupano di prelevare i dati dai dispositivi sono "arduino_class" e "kinect_class", che inseriscono i dati nel corrispettivo buffer FIFO. Infine i thread che si occupano di scrivere i dati su i file di testo sono "ard_handler" e "kin_handler", che scrivono i dati prelevati dal buffer nel corrispettivo file di testo.

3.4 Programma che si occupa di post-elaborare i dati

Dopo aver acquisito i dati di arduino e del kinect bisogna effettuare una post elaborazione per ottenere coppie di valori riferiti a circa lo stesso istante temporale. Il programma usato per fare ciò riprende molte cose dal codice del programma "prg_get_data", le librerie usate sono le stesse ed anche le classi "kin_file_manager", "kin_file_manager" e "data_structure.h". L'unica classe nuova è "make_dataset":

```
1 class make_dataset
2 {
3 private:
4     // dichiaro i puntatori alle classi per la gestione dei file cos'i da renderle disponibili
5     // per tutta la classe
6     ard_file_manager *f_a;
7     kin_file_manager *f_k;
8
9     // delta di accettazione in millisecondi
10    const int range = 35;
11
12    // dichiaro 2 liste per conservare i dati
13    list<arduino_data> a_dat;
14    list<kinect_data> k_dat;
15
16    // devo dichiarare il numero di ingressi e di uscite della rete
17    static const std::size_t n_in = 9;
18    static const std::size_t n_out = 4;
19
20    // dichiaro una lista di tipo dataset_data che verr'a riempita con i dati accoppiati
21    list<dataset_data<float, n_in, float, n_out>> dataset;
22
23    // questa funzione si occupa di scrivere il dataset sul file
24    void write_dataset(string dataset_file_name){...}
25
26    // restituisce 1 solo se i 2 tempi che gli vengono passati differiscono meno della soglia
27    // impostata sopra
```



```

26 bool check_sync(uint64_t t1, uint64_t t2){...}
27
28 // questo metodo sposta i dati dalle strutture "arduino_data" e "kinect_data" alla struttura "
   ↳ dataset" dichiarata globalmente nella classe
29 void transfer_data_to_dataset(arduino_data* ard_iter, kinect_data* kin_iter){...}
30
31 // restituisce la differenza temporale tra i tempi passati alla funzione
32 uint32_t time_distance(uint32_t t1, uint32_t t2){...}
33
34 // sincronizza i dati caricati in "a_dat", "k_dat" e li sposta in "dataset" attraverso la
   ↳ funzione "transfer_data_to_dataset". Questa funzione non tiene conto dei "cloni", nella
   ↳ sezione successiva vi \\'e una spiegazione pi\'u dettagliata.
35 void sync_with_clones(){...}
36
37 // stessa cosa della funzione precedente ma eliminando i cloni.
38 void sync_without_clones(){...}
39
40 // carica i dati dai file usando "f_a" e "f_k", sposta i dati acquisiti in "a_dat" e "k_dat"
41 void load_data(){...}
42
43 public:
44 // l'unica funzione accessibile \\'e il costruttore che provvede ad inizializzare i file
   ↳ manager, caricare i dati, sincronizzarli (con o senza cloni) e scrivere il file del
   ↳ dataset con i dati appena generati.
45 make_dataset(string ard_file_name, string kin_file_name, string dataset_file_name, bool
   ↳ want_clones)
46 {
47     f_a = new arduino_file_manager(ard_file_name, std::ios::in);
48     f_k = new kinect_file_manager(kin_file_name, std::ios::in);
49
50     load_data();
51
52     if (want_clones)
53     {
54         sync_with_clones();
55     }
56     else
57     {
58         sync_without_clones();
59     }
60
61     write_dataset(dataset_file_name);
62 }
63
64 // distruttore della classe
65 ~make_dataset(){...}
66
67
68 };

```

Listing 27: classe make_dataset

In questa classe le 2 funzioni principali sono "sync_with_clones" e "sync_without_clones" che provvedono a costruire le coppie di dati ed a salvarle nella struttura dati "dataset". Analizziamo la funzione "sync_with_clones":

```

1 void sync_with_clones()
2 {
3     uint64_t time_old_good_pair = 0;
4
5     uint64_t time_old_ard = 0;
6     uint64_t time_old_kin = 0;
7
8     auto ard_iter = a_dat.begin();
9     auto kin_iter = k_dat.begin();
10
11     uint32_t ard_cont = 0;
12     uint32_t kin_cont = 0;
13
14     uint32_t old_ard_cont = 0;
15     uint32_t old_kin_cont = 0;

```

```

16
17 while (true)
18 {
19     // controllo la sincronia
20     if (check_sync(ard_iter->frame_time, kin_iter->frame_time))
21     {
22         // aggiungo il dato al dataset
23         transfer_data_to_dataset(&(*ard_iter), &(*kin_iter));
24
25         // stampo delle info utili a capire la qualit a dei dati acquisiti
26         cout << abs((long long) (ard_iter->frame_time - kin_iter->frame_time)) << "   kin: " <<
            ↳ kin_cont << " ard: " << ard_cont << endl;
27         cout << "time since another good pair: " << ((ard_iter->frame_time + kin_iter->frame_time)
            ↳ / 2 - time_old_good_pair) << endl;
28         cout << "ard frame rate: " << (1.0f / ((ard_iter->frame_time - time_old_ard) / 1000.0f)) <
            ↳ < "Hz" << endl;
29         cout << "kin frame rate: " << (1.0f / ((kin_iter->frame_time - time_old_kin) / 1000.0f)) <
            ↳ < "Hz" << endl;
30
31         // salvo i tempi
32         time_old_good_pair = (ard_iter->frame_time + kin_iter->frame_time) / 2;
33         time_old_ard = ard_iter->frame_time;
34         time_old_kin = kin_iter->frame_time;
35
36         // controllo se ci sono pi u dati che si accoppiano con uno solo (solo debug)
37         if ((kin_cont == old_kin_cont) || (ard_cont == old_ard_cont))
38         {
39             cout << "dati multipli nello stesso range" << endl;
40         }
41         else
42         {
43             old_kin_cont = kin_cont;
44             old_ard_cont = ard_cont;
45         }
46
47         cout << endl;
48
49     }
50     else
51     {
52         cout << "dato scartato" << endl << endl;
53     }
54
55     // incremento chi tra i 2 dati  e il pi u vecchio
56     if (ard_iter->frame_time >= kin_iter->frame_time)
57     {
58         kin_iter++;
59         kin_cont++;
60     }
61     else
62     {
63         ard_iter++;
64         ard_cont++;
65     }
66
67     // se sono arrivato alla fine esco dal ciclo
68     if ((kin_iter == k_dat.end()) || (ard_iter == a_dat.end())) { break; }
69 }
70 }

```

Listing 28: funzione sync_with_clones

Questa prima funzione effettua il pair dei dati senza per  tener conto che, dati i differenti frame rate di arduino e del kinect, un singolo dato del kinect potrebbe accoppiarsi con pi  di un dato di arduino e viceversa. Questi "cloni" interferiscono con il processo di apprendimento della rete e quindi si   passati alla funzione "sync_without_clones" che effettua la stessa operazione ma tenendo conto che potrebbero esserci pi  dati di un tipo accoppiati con uno dell'altro e salva solo la coppia che ha la differenza di tempo minore.

```

1 void sync_without_clones()
2 {
3     uint64_t time_old_good_pair = 0;

```

```

4
5 uint64_t time_old_ard = 0;
6 uint64_t time_old_kin = 0;
7
8 auto ard_iter = a_dat.begin();
9 auto kin_iter = k_dat.begin();
10
11 uint32_t ard_cont = 0;
12 uint32_t kin_cont = 0;
13
14 uint32_t old_ard_cont = 0;
15 uint32_t old_kin_cont = 0;
16
17 kinect_data k_tmp = *kin_iter;
18 arduino_data a_tmp = *ard_iter;
19
20 for(;;)
21 {
22     // controllo la sincronia
23     if (check_sync(ard_iter->frame_time, kin_iter->frame_time))
24     {
25         // stampo i tempi per valutare la qualit   del dataset
26         cout << "common " << abs((long long) (ard_iter->frame_time - kin_iter->frame_time)) << "
                ↳ kin: " << kin_iter->contatore << " ard: " << ard_iter->contatore << endl;
27         cout << "ard frame rate: " << (1.0f / ((ard_iter->frame_time - time_old_ard) / 1000.0f)) <
                ↳ < "Hz" << endl;
28         cout << "kin frame rate: " << (1.0f / ((kin_iter->frame_time - time_old_kin) / 1000.0f)) <
                ↳ < "Hz" << endl;
29
30         // salvo i tempi attuali
31         time_old_ard = ard_iter->frame_time;
32         time_old_kin = kin_iter->frame_time;
33
34         // se il contatore del kinect non    variato dall'ultima volta significa che ci sono pi  
                ↳ u dati riferiti ad arduino che si accoppiano ad un solo dato del kinect
35         if (kin_cont == old_kin_cont)
36         {
37
38             cout << "multi ard in the same kin" << endl;
39             if (time_distance(kin_iter->frame_time, ard_iter->frame_time) < time_distance(k_tmp.
                ↳ frame_time, a_tmp.frame_time))
40             {
41                 a_tmp = *ard_iter;
42             }
43         }
44         // se il contatore dei arduino non    variato dall'ultima volta significa che ci sono pi  
                ↳   u dati riferiti al kinect che si accoppiano ad un solo dato di arduino
45         else if (ard_cont == old_ard_cont)
46         {
47             cout << "multi kin in the same ard" << endl;
48             if (time_distance(kin_iter->frame_time, ard_iter->frame_time) < time_distance(k_tmp.
                ↳ frame_time, a_tmp.frame_time))
49             {
50                 k_tmp = *kin_iter;
51             }
52         }
53            se entrambi sono variati significa che la sequenza    finita e posso salvare i dati
54         else
55         {
56             cout << "sequence reset" << endl;
57
58             // salvo la coppia salvata
59             transfer_data_to_dataset(&a_tmp, &k_tmp);
60
61             // stampo i tempi per le valutazioni
62             cout << "best " << abs((long long) (a_tmp.frame_time - k_tmp.frame_time)) << " kin: " <<
                ↳ k_tmp.contatore << " ard: " << a_tmp.contatore << endl;
63             cout << "time since another good pair: " << ((a_tmp.frame_time + k_tmp.frame_time) / 2 -
                ↳ time_old_good_pair) << endl;

```

```

64     time_old_good_pair = (a_tmp.frame_time + k_tmp.frame_time) / 2;
65
66     // salvo i dati correnti e i contatori correnti
67     old_kin_cont = kin_cont;
68     old_ard_cont = ard_cont;
69     a_tmp = *ard_iter;
70     k_tmp = *kin_iter;
71 }
72
73 cout << endl;
74
75 }
76 else
77 {
78     cout << "dato scartato" << endl << endl;
79 }
80
81 // incremento chi tra i 2 dati \e il pi\u vecchio
82 if (ard_iter->frame_time >= kin_iter->frame_time)
83 {
84     kin_iter++;
85     kin_cont++;
86 }
87 else
88 {
89     ard_iter++;
90     ard_cont++;
91 }
92
93 // se sono arrivato alla fine dei dati esco
94 if ((kin_iter == k_dat.end()) || (ard_iter == a_dat.end())) { break; }
95 }
96 }

```

Listing 29: funzione sync_without_clones

La funzione che si occupa di acquisire i nomi dei file e altro dall'utente è "post_elaborate":

```

1 void post_elaborate()
2 {
3     string ard_file_mane = "";
4     string kin_file_mane = "";
5     string dataset_file_name = "";
6     string tmp = "";
7     bool clones = false;
8
9     // acquisisco i nomi dei file dei dati di arduino e del kinect
10    cout << "insert the input file names: " << endl
11        << "\t arduino -> ";
12    std::getline(std::cin, ard_file_mane);
13    cout << "\t kinect -> ";
14    std::getline(std::cin, kin_file_mane);
15
16    // acquisisco il nome del file in cui verr\ a scritto il dataset
17    cout << "insert the ouptut file name -> : ";
18    std::getline(std::cin, dataset_file_name);
19
20    // chiedo se si vuole il dataset con o senza cloni
21    cout << "do you want clones in the data? [y/n] -> : ";
22    std::getline(std::cin, tmp);
23
24    if (tmp == "y")
25    {
26        clones = true;
27    }
28
29    // chiamo la classe make_dataset passandogli i dati appena acquisiti
30    make_dataset(ard_file_mane, kin_file_mane, dataset_file_name, clones);
31 }

```

3.5 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.

3.5.1 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 31: 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 32: librerie STD

3.5.2 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;

```

```

81 };
82
83 //struttura contenente il singolo esempio
84 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 trainingErr = 0; //percentuale di errore training set
94     float validationErr = 0; //percentuale di errore validation set
95 };

```

Listing 33: 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 counRow(string filename) { ... }

```

Listing 34: 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 del'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 di riferimento
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)); } // Sigmoid
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 35: 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 36: class MLP

3.5.3 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 37: altre classi

3.5.4 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 bayes 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 CUDAOutDiff(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 }
```

```

41 }
42
43 //calcola l'errore dei neuroni dello strato output
44 __global__ void CUDAAoutputErr(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 CUDASetInput(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 CUDASetSingleInput(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[ArcIn[i]], NeuronOut[ArcOut[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 38: 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 39: 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 40: classe di interfaccia alla GPU

3.5.5 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    //////////////////////////////////lancio dei kernel all'interno della gpu////////////////////////////////////
107    int startA = 0;
108    int endA = 0;
109    int startN = 0;
110    int endN = 0;
111    int numLayerArcs = 0;
112    int numLayerNeur = 0;
113    int numOfBlocksMax = 0;
114    int numOfBlocksA = 0;
115    int numOfBlocksN = 0;
116    int numOfBlocksOut = floorf(outputN / ThxBlock) + 1;
117    int exampleRef = 0;
118    int outputRef = NeuronOut.size() - outputN;
119    long long t0 = 0, t1 = 0;
120    long long t0in = 0, t1in = 0;
121    double elapsedMilliseconds = 0;
122    double elapsedInMilliseconds = 0;
123
124    for (int it = 0; it < Niter; it++) { //scorro le iterazioni
125
126        t0 = PerformanceCounter();
127
128        for (int t = 0; t < (examples.size() / (inputN + outputN)); t++) { //scorro gli esempi
129            //imposto il riferimento per l'esempio di input
130            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 //resetto il vettore contenente l'errore dei neuroni
168 numOfBlocksN = (floorf(BPerr.size() / ThxBlock) + 1);
169 CUDAResetVector <<<numOfBlocksN, ThxBlock >>> (dev_BPerr, BPerr.size());
170
171 CUDAResetVar <<<1, 1 >>> (dev_MeanErr);
172 CUDASetOutputErr <<<numOfBlocksOut, ThxBlock >>> (dev_NeuronOut, outputRef, NeuronOut.size()
    ↪ (), inputN, dev_BPerr, dev_examples, exampleRef, dev_mapMaxOut, dev_mapMinOut,
    ↪ dev_MeanErr);
173 cudaMemcpy(CMeanErr, dev_MeanErr, sizeof(float), cudaMemcpyDeviceToHost);
174
175 MeanErr += *CMeanErr / outputN;
176
177 //retropropagazione dell'errore
178
179 for (int i = priority.size() - 2; i > 1; i--) {
180
181     startA = priority[i - 1] + 1;
182     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);
310     if (cudaCheckStatus(cudaStatus) == true) goto Error;
311     cudaStatus = cudaMemcpy(gpuNetParams.ArcIn, CArcIn, ArcIn.size() * sizeof(int),
312                             ↪ 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, tlin = 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
465     float delta;
466     cout << "input time: " << elapsedInMilliseconds << " ms" << endl;
467
468     for (int on = 0; on < outputN; on++) {
469         delta = mapMaxOut[on] - mapMinOut[on];
470         cout << "Y" << on << ": " << (NeuronOut[NeuronOut.size() - outputN + on] * delta) +
471             ↪ mapMinOut[on] << endl;
472     }
473     cout << endl;
474
475     //////////////////////////////////////
476
477     if (false) {
478         Error:
479         //libero la memoria nella scheda grafica
480         cudaFree(gpuNetParams.weights);
481         cudaFree(gpuNetParams.ArcIn);
482         cudaFree(gpuNetParams.ArcOut);
483         cudaFree(gpuNetParams.NeuronOut);
484         cudaFree(gpuNetParams.examples);
485         cudaFree(gpuNetParams.BPerr);
486         cudaFree(gpuNetParams.mapMaxOut);
487         cudaFree(gpuNetParams.mapMinOut);
488         cudaFree(gpuNetParams.priority);
489         cudaFree(gpuNetParams.NeurInLyr);
490     }
491
492     return cudaStatus;
493 }
494
495 //////////////////////////////////////
496
497 //////////////////////////////////////CUDA UTILITY////////////////////////////////////
498 //verifica la corretta esecuzione di un operazione
499 inline cudaError_t checkCuda(cudaError_t result){ ... }
500
501 //verifica la corretta esecuzione di un operazione restituendo un bool
502 bool cudaCheckStatus(cudaError_t cudaStatus) { ... }
503
504 //stampa a schermo le principali proprietà della scheda

```



```

497 void printDeviceSpecs() { ... }
498
499 //stampa i parametri della rete che vengono passati alla scheda
500 void printNetSpecs() { ... }
501
502 //calcola il peso del modello
503 float sizeofModel(string mesureUnit = "B") { ... }
504
505 template<typename T, typename A>
506 float sizeofVector(vector<T, A> const& vect, string mesureUnit = "B") { ... }
507 //////////////////////////////////////
508 };

```

Listing 41: classe di interfaccia alla GPU

3.6 Interfaccia grafica

L'interfaccia grafica dell'applicativo è stata realizzata in Unity, un popolare motore grafico per videogiochi e applicazioni android e ios. E' stato realizzato un rig (un manichino) attraverso blender con uno scheletro interfacciabile con Unity, in modo da poterne controllare i movimenti attraverso l'inserimento degli angoli dei vari "ossi" dello scheletro. Lo script dell'interfaccia grafica è stato scritto in C# e si basa su un client che richiede continuamente al server i dati di uscita della rete neurale.

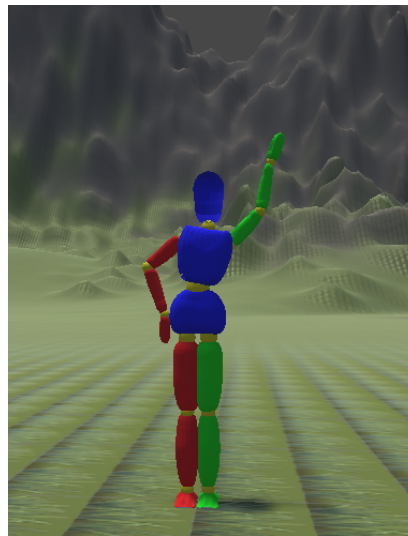


Figura 23: rig utilizzato per la renderizzazione

3.6.1 classe syncSocketServer

Questa parte illustra la classe utilizzata per far effettuare allo script in unity la richiesta al server centrale dei dati elaborati.

3.6.1.1 librerie utilizzate

tale classe è stata implementata nella parte unity e ha richiesto l'importazione di UnityEngine, oltre alle librerie per la comunicazione socket e per la manipolazione dei byte ricevuti dal main server.

```

1 using System;
2 using System.Net;
3 using System.Net.Sockets;
4 using System.Runtime.InteropServices;
5 using System.Text;
6 using System.Collections;
7 using System.Collections.Generic;
8 using UnityEngine;

```

Listing 42: librerie class client C#

3.6.1.2 classe client C# per la comunicazione

In questa sezione è illustrato la classe che permette la comunicazione con il main server. Si noti come l'utilizzo della classe in unity richiede che essa sia figlia della classe MonoBehaviour, interna di unity.

```
1
2 public class SyncSocketClient : MonoBehaviour {
3     //struttura union per la conversione Bytes -> float
4     [StructLayout(LayoutKind.Explicit)]
5     public struct Ubones
6     {
7         [FieldOffset(0)] public float fBones;
8         [FieldOffset(0)] public byte strBones1;
9         [FieldOffset(1)] public byte strBones2;
10        [FieldOffset(2)] public byte strBones3;
11        [FieldOffset(3)] public byte strBones4;
12    }
13
14    //Oggetti socket
15    public IPEndPoint ipHostInfo;
16    public IPAddress ipAddress;
17    public IPEndPoint remoteEP;
18    public Socket sender;
19    private const int number_of_joints = 2;
20    private const int number_of_float = number_of_joints * 3;
21
22
23    // alloco la struttura di conversione
24    public Ubones[] bones = new Ubones[number_of_float]; //4 joints
25
26    // dichiaro il vettore di buffer
27    public byte[] bytes;
28
29    //settaggio parametri e inizzializzazione comunicazione
30    public SyncSocketClient(int Port){ ... }
31
32    //chiusura comunicazione
33    public void CloseConnection(){ ... }
34
35    // funzione per la ricezione e la conversione float
36    public void ReciveFloat(){ ... }
37
38    //funzione per la ricezione della stringa dal server
39    public string ReciveString(){ ... }
40
41    //trasmissione stringa
42    public void sendMessage(string msg){ ... }
43
44    //funzione di comunicazione, invia il codice di richiesta, riceve i byte li carica nell'
45    ↪ union
46    public void TestComunication(){ ... }
47
48    //funzione per la reinizializzazione delle variabili bones
49    public void initBones(){ ... }
50 }
```

Listing 43: class client C#

3.6.2 classe di controllo del manichino

In questa sezione viene illustrata la classe interna al motore grafico realizzata per il caricamento dei dati ricevuti dal server nella struttura di dell'oggetto

3.6.2.1 librerie della classe controllo

```
1 using System.Collections.Generic;
2 using System.IO;
```

```

3 using System.Text;
4 using UnityEngine;

```

Listing 44: Librerie classe controllo C#

3.6.2.2 corpo della classe controllo

Anche questa classe è stata resa figlia di MonoBehaviour, il settaggio delle variabili viene effettuato dal metodo Start() richiamato da Unity all'avvio del programma, in questo caso però l'utilizzo della classe è richiamato attivamente da Unity attraverso il metodo Update() che viene richiamato una volta per ogni aggiornamento dello schermo.

```

1 public class ProvaControllo : MonoBehaviour {
2     //dichiarazione della classe socketClient sincrono
3     //e dell'oggetto di controllo del manichino
4     private SyncSocketClient socket;
5     public GameObject[] bones;
6
7     private List<string> Tags = new List<string>();
8
9     /*{ 4, 6, 3, 5, 8; */ // indici degli oggetti da aggiornare nella lista GameObject
10    public int[] neededObjects = new int[] {3, 5};
11
12    //ref kinect{ 4, 5, 8, 9; } // indici degli oggetti da aggiornare nella struct Ubones della
13    ↪ com socket
14    public int[] neededBonesF = new int[] {0, 1};
15
16    //vettori da tre per il passaggio degli angoli
17    public Vector3 BraccioDx;
18    public Vector3 avanBraccioDx;
19    public Vector3 BraccioSx;
20    public Vector3 avanBraccioSx;
21
22    // Metodo di inizializzazione della classe, richiamato all'avvio
23    void Start () {
24
25        //attivazione della comunicazione
26        socket = new SyncSocketClient(9999);
27
28        //inizializzazione della struttura di conversione byte/float
29        socket.initBones();
30
31        //creazione della lista dei tag
32        Tags.Add("Root"); //0
33        Tags.Add("Bacino"); //1
34        Tags.Add("schiena"); //2
35        Tags.Add("BraccioDx"); //3
36        Tags.Add("BraccioSx"); //4
37        Tags.Add("AvanbraccioDx"); //5
38        Tags.Add("AvanbraccioSx"); //6
39        Tags.Add("testa"); //7
40        Tags.Add("schiena"); //8
41
42        //inizializzazione dell'array di oggetti del gioco
43        bones = new GameObject[10];
44
45        //link degli oggetti creati alle ossa del manichino
46        int i = 0;
47        foreach (string Tag in Tags) {
48            bones[i] = GameObject.FindGameObjectWithTag(Tag);
49            i++;
50        }
51
52        // Update is called once per frame
53        void Update () {
54
55            //richiesta al server centrslle dei nuovi dati e aggiornamento
56            //della struttura interna alla classe socket
57            socket.TestCommunication();

```

```

58
59 //caricamento dei dati ricevuti sulla struttura unity del manichino
60 for (int i = 0; i < 2; i++)
61 {
62     bones[neededObjects[i]].transform.rotation = Quaternion.Euler(socket.bones[
        ↪ neededBonesF[i]*3].fBones, socket.bones[(neededBonesF[i] * 3) + 1].fBones,
        ↪ socket.bones[(neededBonesF[i] * 3) + 2].fBones);
63 }
64 }

```

Listing 45: class client C#

4 Conclusioni

Il progetto in fine si è concluso con esito negativo per quanto riguarda la buona generalizzazione del modello, principalmente a causa della cattiva qualità della bussola 3D, la quale anche durante i test sull'output non filtrato restituiva valori non coincidenti con la realtà. In compenso la libreria sviluppata per l'addestramento dei modelli ha dato buoni risultati in termini di ottimizzazione dell'addestramento dei modelli, incrementando su GPU la velocità di addestramento di ben 250 volte rispetto alla CPU, risultati soddisfacenti considerando che il programma è affetto fortemente da "race-condition", una problematica che affligge gli approcci al calcolo che si basano fortemente su multithreading. Tale problematica si presenta quando processi paralleli cercano di eseguire calcoli che mirano a modificare una stessa variabile contemporaneamente, e tale variabile è anche un input di queste funzioni, in particolare un processo prende tale variabile come input copiandola, un altro thread la modifica, e il primo la sovrascrive con il suo risultato non considerando la modifica fatta dal secondo. Per far fronte a tale problema è stato necessario utilizzare le funzioni atomiche che causano un blocco della variabile e mettono in attesa i thread che cercano di utilizzarla, causando una serializzazione delle operazioni. Un possibile altro approccio sarebbe stato effettuare queste operazioni con uno specifico ordine che mirava a fare più operazioni su più variabili evitando però la race-condition, con tale ottimizzazione si possono raggiungere velocità di computazione molto maggiori, in quanto un solo thread in attesa in un warp in cui tutti gli altri hanno concluso il proprio lavoro, li blocca tutti, e non permette allo stream-multiprocessor di lanciarne un altro. Abbiamo comunque deciso di non intraprendere questa strada in quanto richiedeva molto tempo e ingenti modifiche alla struttura del programma. In oltre con modelli a molti neuroni per layer tendeva ad attenuarsi in quanto si verificavano meno collisioni. Concludendo il problema non è stato risolto con i risultati sperati, ma la libreria realizzata è perfettamente funzionante e pronta per essere applicata per l'apprendimento di nuove task.