

UNIVERSITATEA NATIONALA DE STIINTA SI TEHNOLOGIE  
POLITEHNICA BUCURESTI  
FACULTATEA DE AUTOMATICĂ ȘI CALCULATOARE  
DEPARTAMENTUL DE CALCULATOARE

## PROIECT DE DIPLOMA

Orientare in Spatiu folosind ORB-SLAM  
BUCUREȘTI

Alfred Andrei Pietraru

**Coordonator științific:**

Prof. dr. ing. Anca Morar

NATIONAL UNIVERSITY OF SCIENCE AND TECHNOLOGY  
POLITEHNICA BUCHAREST  
FACULTY OF AUTOMATIC CONTROL AND COMPUTERS  
COMPUTER SCIENCE AND ENGINEERING DEPARTMENT

## DIPLOMA PROJECT

Spatial Orientation using ORB-SLAM  
BUCHAREST

Alfred Andrei Pietraru

**Thesis advisor:**

Prof. dr. ing. Anca Morar

# CUPRINS

<b>1</b>	<b>Introducere</b>	<b>1</b>
1.1	Context . . . . .	1
1.2	Problema . . . . .	1
1.3	Obiective . . . . .	1
1.4	Solutia propusa . . . . .	2
1.5	Rezultatele obtinute . . . . .	2
1.6	Structura lucrarii . . . . .	2
<b>2</b>	<b>Cerinte si Motivatie</b>	<b>3</b>
2.1	Motivatie . . . . .	3
2.2	Cerinte Functionale si Nonfunctionale . . . . .	4
<b>3</b>	<b>Studiu de piata</b>	<b>5</b>
3.1	State of the art Visual SLAM . . . . .	5
<b>4</b>	<b>Solutie propusa</b>	<b>8</b>
4.1	Achizitia datelor . . . . .	9
4.2	Extragerea trasaturilor . . . . .	9
4.3	Harta punctelor din spatiu . . . . .	12
4.4	Asociere puncte din spatiu cu feature-uri ORB . . . . .	13
4.5	Optimizare Estimare Pozitie Initiala . . . . .	14
4.5.1	Proiectarea MapPoint in planul imaginii . . . . .	14
4.5.2	Motion Only Bundle Adjustment . . . . .	15

4.6	Crearea unui cadru cheie . . . . .	16
4.6.1	Optimizare harta locala . . . . .	17
4.6.2	Reteaua Neurala FastDepth . . . . .	19
<b>5</b>	<b>Detalii de implementare</b>	<b>23</b>
5.1	Limbaje de programare si librarii folosite . . . . .	23
5.2	Mediu de lucru si principalele clase . . . . .	25
5.3	Pipeline antrenare FastDepth . . . . .	37
<b>6</b>	<b>Evaluare</b>	<b>39</b>
6.1	Setul de date TUM RGBD Dataset . . . . .	39
6.2	Metrici utilizate . . . . .	40
<b>7</b>	<b>Concluzii</b>	<b>43</b>
<b>8</b>	<b>Referinte</b>	<b>45</b>

## SINOPSIS

Această lucrare reprezintă o reimplementare a algoritmului ORB-SLAM2, utilizat pentru localizarea și cartografierea unui mediu interior necunoscut. Sunt abordate concepte precum extragerea de trăsături folosind descriptorul ORB, realizarea sarcinilor de feature matching între cadre, utilizarea algoritmului Bundle Adjustment pentru optimizare, implementat prin biblioteca Ceres, precum și relocalizarea poziției camerei prin vectori obținuți cu metoda bag-of-words, implementată în DBOW2. Algoritmul este testat pe setul de date TUM RGB-D, iar rezultatele obținute sunt comparabile cu valorile de tip ground truth. Se explorează, de asemenea, integrarea unei rețele neuronale de tip FastDepth pentru estimarea matricii de adâncime în cadrul fluxului de procesare al algoritmului. Implementarea este realizată în C++, având un număr de linii de cod de aproximativ trei ori mai mic decât cel din implementarea oficială. De asemenea, algoritmul este optimizat pentru a funcționa pe perioade lungi de timp datorită unui management eficient al memoriei.

## ABSTRACT

This work represents a reimplementation of the ORB-SLAM2 algorithm, used for localization and mapping of an unknown indoor environment. Concepts such as feature extraction using the ORB descriptor, performing feature matching tasks between frames, utilizing the Bundle Adjustment algorithm for optimization implemented through the Ceres library, as well as camera position relocalization through vectors obtained with the bag-of-words method implemented in DBOW2 are addressed. The algorithm is tested on the TUM RGB-D dataset, and the obtained results are comparable with ground truth values. The integration of a FastDepth neural network for depth matrix estimation within the algorithm's processing pipeline is also explored. The implementation is carried out in C++, having a number of lines of code approximately three times smaller than that of the official implementation. Additionally, the algorithm is optimized to function over long periods of time due to efficient memory management.

# 1 INTRODUCERE

## 1.1 Context

SLAM, Simultaneous Localization and Mapping reprezinta o clasa de algoritmi de planificare si control al miscarii unui agent prin mediu pentru a construi un model al spatiului cat mai apropiat de realitate. Acestea au castigat atentia publicului in ultimii ani, lucru care a condus la dezvoltarea numeroaselor variante care exista la momentul curent pe piata, fiecare adaptat pentru mediul si tipul de senzori folositi. O atentie deosebita a fost acordata sistemelor de tip Visual SLAM din mai multe motive: camerele video sunt unul dintre cele mai comune tipuri de senzori, exista o multitudine de tehnici de Computer Vision pentru procesarea imaginilor iar filozofia pe care acesti algoritmi o urmeaza este asemanatoare cu modul in care creierul uman interpreteaza mediul inconjurator: sunt alese un set de puncte din spatiu care vor fi considerate referinte iar unghiul din care acestea sunt observate poate oferi informatii despre pozitia agentului in mediu. Astazi, cei mai populari algoritmi de tip Visual SLAM imbina domenii precum Machine Learning, Computer Vision, Robotica si Matematica pentru a crea sisteme robuste, capabile sa indeplineasca o varietate de sarcini.

## 1.2 Problema

Creati un sistem capabil sa exploreze un mediu de interior necunoscut. Acesta trebuie sa creeze o harta a zonei parcurse, sa reconstruiasca traseul estimand pozitia camerei pentru fiecare cadru citit, sa fie tolerant la erori si sa poata opera pentru perioade de timp indelungate.

## 1.3 Obiective

Obiectivele principale ale lucrarii sunt:

- studierea, configurarea si implementarea unui sistem de tip ORB-SLAM2

- testarea performantelor folosind seturi de date consacrate, utilizarea unui video realizat de mine pentru etapa de evaluare
- testarea unei implementari in care o camera tip RGBD sa fie inlocuita cu o camera monoculara traditionala si o retea neurala, aceasta sarcina presupunand alegerea unei arhitecturi potrivite si compararea celor 2 metode
- prezentarea problemelor intalnite in etapa de dezvoltare si a unor directii de imbunatatire

## **1.4 Solutia propusa**

Lucrarea propune implementarea algoritmului ORB-SLAM2 adaptat pentru camerele de tip RGBD si testarea acestuia pe setul de date TUM RGBD. Se vor analiza aspecte precum acuratetea traiectoriei si indeplinirea conditiilor de functionare in timp real. De asemenea, se va incerca inlocuirea matricei de adancime utilizata de camera RGBD, cu o harta de distante calculata de catre reseaua neurala FastDepth.

## **1.5 Rezultatele obtinute**

Rezultatele au aratat ca implementarea algoritmului ORB-SLAM2 folosind o camera tip RGBD da rezultate bune in medii indoor si ca poate functiona in timp real cu viteze de aproximativ 10-15 cadre pe secunda. Utilizarea unei retele neurale pentru estimarea distantei nu a functionat, algoritmul fiind capabil sa functioneze pentru cel mult 50 de cadre.

## **1.6 Structura lucrarii**

Lucrarea este structurată în mai multe capitole: introducerea oferă contextul lucrării, definește problema abordată, obiectivele și soluția propusă. Capitolul 2, descrie cerintele functionale, nonfunctionale si motivatia. Capitolul 3 realizeaza un studiu de piata asupra metodelor curente, separandu-le in 3 categorii. Capitolul 4 descrie la nivel conceptual solutia propusa: algoritmii folositi si componentele logice. Capitolul 5 detaliaza modul in care este realizata evaluarea si rezultatele obtinute. Ultimul capitol prezinta cateva impresii personale despre acest proiect: lucruri pe care le-as fi imbunatatit, problemele intalnite si directiile viitoare.

## 2 CERINTE SI MOTIVATIE

### 2.1 Motivatie

Filmele, cartile si jocurile pe calculator ne prezinta un viitor al omenirii in care roboti inteligenti indeplinesc sarcinile din viata de zi cu zi sau cele care ar putea pune in pericol siguranta omului. Exista zeci de filme care descriu acest modul in care ar arata un robot inteligent, capabil sa se adapteze la mediu si sa interactioneze cu omul. Desi in momentul de fata suntem departe de a crea un framework suficient de complex pentru un asemenea agent, consider ca algoritmi din categoria Visual SLAM sunt un pas in directia corecta. Imi este greu sa imi imaginez un robot care sa poata simula compartamentul uman si sa nu fie capabil sa se deplaseze si sa inteleaga mediul in care se afla. Pentru noi, aceste lucruri sunt adanc inradacinate in modul in care functioneaza creierul, dar pentru un calculator, a fost nevoie de aproape 20 de ani de cercetare pentru a crea algoritmi suficienti de complexi pentru a indeplini niste sarcini minimale de orientare cum ar fi capacitatea de invatare a mediului si de pozitionare a agentului in spatiu. Chiar si cu algoritmi de tip SLAM dezvoltati pana in acest moment, exista numeroase aplicatii practice:

- realizarea sarcinilor din viata de zi cu zi: cazul robotilor de curatenie sau a celor care transporta obiecte in interiorul unei cladiri
- in aplicatii medicale, ca de exemplu asistenta pentru persoanele nevazatoare
- in aplicatii militare: cartografierea zonelor necunoscute: in interiorul cladirilor sau in medii ostile unde nu este acces la un sistem de coordonate globale cum ar fi pozitia data de un sistem GPS
- in aplicatii industriale, inspectii asupra instalatiei sau depozitelor, detectarea unor erori si raportarea zonei in care au fost observate

Exista numeroase aplicatii pentru sistemele SLAM, doar toate pleaca de la principiul ca agentul trebuie sa creeze o harta a mediului si sa inteleaga care este pozitia acestuia. Pe masura ce aceste sisteme vor evolua va creste si complexitatea sarcinilor pe care le pot indeplini.



## 2.2 Cerinte Functionale si Nonfunctionale

Din punct de vedere al cerintelor functionale, algoritmul ORB-SLAM va primi un video realizat cu o camera de tip RGBD si va return 2 fisiere text, unul va contine estimarea pozitiei pentru fiecare cadru in parte iar celalalt va avea salvata harta mediului inconjurator, alcatuita dintr-un nor de puncte in spatiu si cadrele cheie asociate acestora. Algoritmul va avea o interfata grafica minimala alcatuita din 2 ferestre. Cea din stanga va contine o reprezentare pentru cadrul curent procesat, acesta va avea culoarea albastru, celelalte cadre cheie salvate in harta vor avea verde si cu rosu punctele din spatiu. Toate acestea vor alcatui impreuna harta mediului inconjurator. In fereastra din dreapta va fi afisat fiecare cadru in format alb negru, iar cu rosu vor fi marcate feature-urile detectate de algoritmul ORB. Cadrele cheie consecutive sunt conectate intre ele prin intermediul unei drepte de culoare neagra. Cadrele cheie si modul in care sunt unite intre ele vor recompune traseul realizat de camera in video.

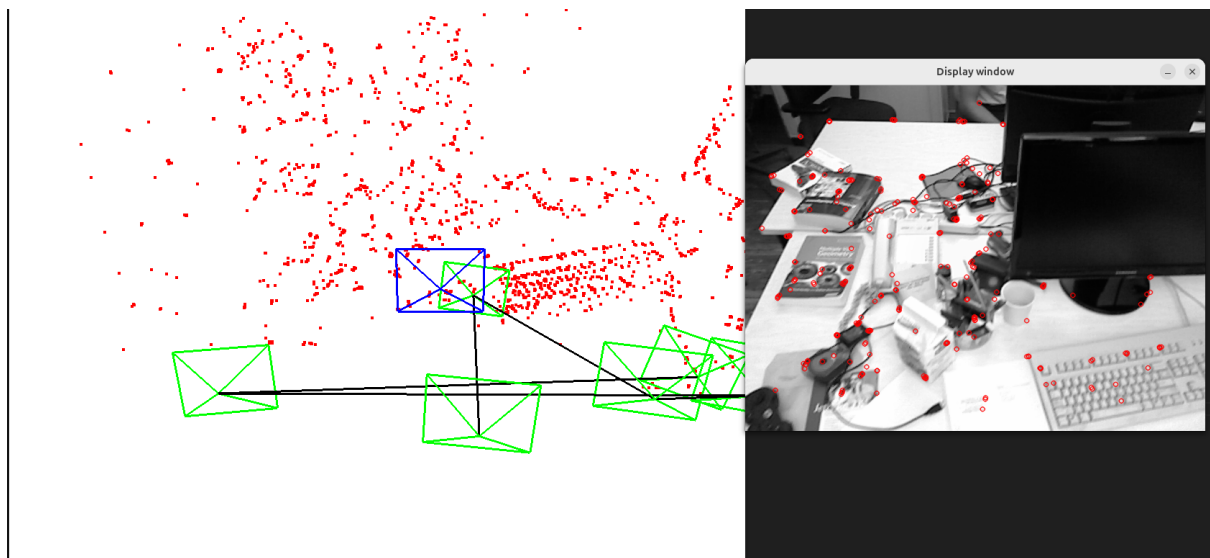


Figura 1: Interfata grafica, stanga reprezentarea hartii, dreapta extragere feature-uri cu ORB

Ca cerinte nonfunctionale, algoritmul trebuie sa mearga in timp real, sa proceseze intre 10-15 cadre pe secunda si sa poata fi folosit in sisteme embedded in care unitatea centrala de procesare are cel mult 2 core-uri si nu poate folosi GPU-ul. Sistemul trebuie sa fie rezistent la erorile de estimare pentru fiecare imagine primita, sa optimizeze atat harta cat si traseul realizat si sa aiba capacitatea de relocalizare in situatia in care urmarirea cadru cu cadru esueaza. Mediul in care poate sa opereze este unul static, de mici dimensiuni si nu poate accesa coordonatele globale ale pozitiei sale prin intermediul tehnologiilor precum GPS.

## 3 STUDIU DE PIATA

### 3.1 State of the art Visual SLAM

Cei mai noi algoritmi de visual SLAM functioneaza acum folosind tehnici de deep learning. In continuare vor fi prezentate lucrarile care au reprezentat SOA, pana la inceputul anului 2025, grupate in categorii in functie de modul in care sunt folosite tehnicile de deep learning. Am considerat potrivita impartirea pe 3 nivele a algoritmilor, in functie gradul de utilizare al tehnicilor de deep learning pentru realizarea operatiile specifice sistemelor SLAM:

1. Algoritmi care se bazeaza fundamental pe tehnici de deep learning pentru a functiona, DPV-SLAM, ESLAM.
2. Algoritmi care sunt la granita dintre metodele clasice si cele deep learning, in care doar anumite componente sunt imbunatatite cu ajutorul retelelor neurale: Light-SLAM, HFNet-SLAM.
3. Algoritmii clasici, care nu folosesc deloc retele neurale: ORB-SLAM3, SVO.

Deep Patch Visual SLAM (DPV-SLAM), este un sistem SLAM care foloseste deep neural networks. Acesta imparte operatiile care trebuie realizate in 2 categorii: frontend-ul care realizeaza sarcina de visual odometry cu ajutorul unui sistem derivat din Deep Patch Visual Odometry (DPVO) si partea de backend alcatuita din 2 metode de loop closure: proximity loop closure si classical loop closure. Algoritmul are nevoie intre 5-6 GB de memorie pe GPU pentru a putea rula. O alta problema o reprezinta proximity loop closure. Aceasta functioneaza cu ajutorul unei harti foarte dense de feature-uri obtinute cu ajutorul metodei de optical flow fiind imposibil de folosit in timp real fara a folosi GPU.

ESLAM sau Efficient Dense Visual SLAM using Neural Implicit Maps este un sistem de SLAM monocameră RGB-D care folosește o combinație între o harta densa 3D, reprezentata de o retea neurala implicita si un backend optimizat geometric pentru estimarea matricei de pozitie

a camerei. Avantajele acestei implementari sunt ca produce o harta densa si detaliata si poate reconstrui detalii chiar si in zone partial observate. Problema acestei implementari este ca necesita un GPU si resurse mari de calcul si nu este potrivit pentru dispozitivele embedded.

Light-SLAM este construit pornind de la aceeasi filozofie ca si ORB-SLAM2, partea de backend reprezentata de local mapping, adica optimizarea hartii create si loop closure, recunoasterea zonelor prin care a mai trecut algoritmul si inchiderea buclelor traiectoriei, acestea sunt realizate folosind metode clasice. Extragerea de keypoint-uri, descriptori si sarcina de matching intre descriptorii a doua imagini consecutive este realizata de 2 retele neurale. Acest sistem poate functiona in timp real daca se poate folosi un GPU, dar cea mai mare problema o reprezinta faptul ca retele neurale nu sunt capabile sa gaseasca feature-uri cu acuratete suficient de buna in zone care nu seamana cu ceea ce a intalnit in setul de date de antrenare, astfel algoritmul nu are garantia ca va functiona in situatii critice.

HFNet-SLAM este o metoda construita pe baza ORB-SLAM3 si folosindu-se de arhitectura HF-Net, avand straturile de convolutie separate in depthwise convolution si pointwise convolution, asemanator cu modul in care este gandit Mobile\_Net. In loc sa foloseasca 2 retele neurale precum Light-SLAM, aceasta foloseste una singura, atat pentru extragere keypoint-urilor si a descriptorilor cat si pentru feature-urile globale, folosite in sarcinile de loop closure. Pe langa problema retelei neurale care trebuie sa ruleze pe GPU si a feature-urilor instabile extrase din imagini pentru zone care nu au fost intalnite in setul de antrenare, algoritmul calculeaza pentru fiecare cadru in parte feature-urile sale globale lucru care adauga un overhead computational inutil. De asemenea reseaua neurala nu extrage keypoint-urile pe mai multe nivele, obtinandu-se prea putine puncte pentru a mentine sistemul stabil in mediile slab texturate.

ORB-SLAM3 este continuare implementarii algoritmului ORB-SLAM2 pe care l-am ales eu. Acesta a aparut in 2021 si pana in acest moment este cea mai complexa si completa metoda de a estima traiectoria camerei si a reconstrui o harta de puncte a mediului inconjurator folosind doar metode clasice. In comparatie cu precedesorul acestuia, implementarea de ORB-SLAM3 foloseste datele obtinute de la Inertial Measurement Unit (IMU) si optimizeaza rezultatele

primite folosind tehnica de Maximum a Posteriori Estimation (MAP). In comparatie cu predecesorul sau, versiunea curenta algoritmului lucreaza cu multiple harti locale, respectiv noruri de puncte in spatiu. In momentul in care ORB-SLAM3 pierde orientarea si trebuie sa execute o relocalizare, sistemul genereaza o noua harta pentru a mentine fluida procesarea cadru cu cadru. In situatia in care tracker-ul recunoaste zona in care a ajuns, incearca sa uneasca hartile intre ele pentru a reconstrui intreg mediul. Am considerat ca zona pe care o parcurge agentul nostru este de mici dimensiuni si nu ar avea nevoie de un sistem atat de complex de interconectare al hartilor generate iar utilizarea acestuia ar adauga un overhead nejustificat. In plus, nu avem acces la datele ce apartin componentei IMU.

SVO sau Semi-Direct Visual Odometry for Monocular and Multi-Camera Systems este un exemplu de algoritm tip SLAM care foloseste doar 2 thread-uri: unul responsabil de urmarirea cadru cu cadru si celalalt pentru optimizarea hartii. Acesta foloseste gradientii pixelilor in imagini pentru a crea feature-uri, nu doar contururile obiectelor. In cazul algoritmilor din familia ORB-SLAM care incearca sa optimizeze eroarea de proiectie a punctelor din spatiu, aici se folosesc metode directe si trebuie minimizata eroarea fotometrica a pixelilor aflati in apropierea conturilor obiectelor. Este printre cei mai rapizi algoritmi de SLAM, procesand peste 100 de cadre pe secunda pe un CPU, dar genereaza o harta cu mult prea putine puncte care rareori poate fi refolosita, nu exista capacitate de relocalizare si este dificil de extins. Fiind poate a doua cea mai buna optiune dupa ORB-SLAM2.

In ciuda faptului ca algoritmul ORB-SLAM2 a aparut in 2017, in continuare ramane un exemplu de sistem bine gandit, cu multe posibilitati de extindere si capacitate de a fi adaptat la cerintele din zilele noastre. Indeplineste cu succes toate cerintele functionale si nonfunctionale pe care sistemul ar trebui sa le aiba: poate fi folosit in real time, implementarea procesand aproximativ 15 cadre pe secunda, creaza o harta mediului inconjurator pe care o poate optimiza, are capacitate de relocalizare si corecteaza erorile de estimare care apar in timp prin mecanismul de loop closure. Acesta poate rula exclusiv pe CPU, fiind potrivit atat pentru vehicule la sol, dar si pentru drone. Nu are nevoie de o estimare a pozitiei globale, putand fi folosit in medii ostile unde terenul este complet necunoscut.

## 4 SOLUTIE PROPUSA

Solutia mea propune implementarea algoritmului ORB-SLAM2. Acesta are 2 scopuri fundamentale:

- sa estimeze pentru fiecare cadru in parte matricea de pozitie si orientare a camerei, reconstruind astfel traseul parcurs in timpul functionarii algoritmului
- sa creeze o harta locala a mediului inconjurator pentru a memora zonele prin care a mai trecut si pentru a imbunatatii estimarea traiectoriei

Matricea de pozitie si orientare a camerei (pose matrix) are dimensiuni  $4 \times 4$  si are formatul prezentat mai jos, unde  $R$  reprezinta matricea de rotatie  $3 \times 3$ , iar  $t$  este vectorul coloana de dimensiune 3, reprezentand translatia fata de punctul de origine  $(0,0,0)$ . Aceasta mai este denumita si matricea de conversie din sistemul de coordonate global (world space) in sistemul de coordonate al camerei (camera space) si este notata in implementarea mea ca  $T_{cw}$ . Inversa acestei matrice notata  $T_{wc}$  realizeaza operatia de conversie dintre cele 2 sisteme de coordonate in sens opus.

$$T_{cw} = \begin{bmatrix} R & t \\ 0 & 1 \end{bmatrix}, \quad T_{wc} = \begin{bmatrix} R & -R^t t \\ 0 & 1 \end{bmatrix} \quad (1)$$

Algoritmul primeste ca date de intrare: sursa de la care va obtine imaginile de tip RGB pe care va trebui sa le prelucreze, acestea pot sa provina atat de la un video, set de date consacrat sau chiar in timp real direct de la camera, parametrii de distorsiune a imaginii si matricea parametrilor interni ai camerei, avand dimensiunea  $3 \times 3$  si notata in mod traditional cu  $K$ . Aceasta contine 4 constante importante: distanta focala a camerei pe axa  $x$  si pe  $y$   $f_x, f_y$  si  $c_x, c_y$  reprezentand coordonatele centrului imaginii. Aceasta matrice trebuie modificata de fiecare data cand este schimbata camera cu care se realizeaza filmarea, sau cand se fac operatii de modificare a dimensiunii imaginilor fata de modul in care ar fi acestea extrase natural. Matricea are urmatoarea forma:

$$K = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \quad (2)$$

Algoritmul va returna un fisier text in care se vor afla estimarile matricilor de pozitie impreuna cu timestamp-ul asociat pentru fiecare cadru in parte in ordine cronologica. Rezultatul poate fi comparat cu fisiere care contin valorile reale si care respecta acelasi format pentru a verifica corectitudinea algoritmului. Diagrama UML prezinta interactiunea dintre principale componente descrise din punct de vedere functional, dar si flow-ul natural al algoritmului. In continuare voi detalia logica fiecarei componente din punct de vedere al algoritmilor folositi, ale valorilor de intrare si de iesire asociate acestora.

## 4.1 Achizitia datelor

Scopul acestei componente este sa citeasca imaginea de tip RGB de la camera, sa extraga matricea de adancime asociata cadrului curent, de exemplu: prin intermediul unei camere stereo, a unei camere de tip RGBD sau cu ajutorul unei retele neurale si sa creeze o estimare initiala pentru pozitia curenta a camerei pe baza masuratorilor anterioare. In viitor, o alta functie a acestei componente ar putea fi extragerea datelor de la instrumente de masura precum giroscop sau accelerometru, pentru a obtine informatii suplimentare cu privire la orientarea si distanta efectuata de catre camera care ar putea imbunatatii considerabil estimarea initiala a pozitiei.

## 4.2 Extragerea trasaturilor

Ca date de intrare aceasta componenta primeste doar imaginea de tip RGB, si extrage aproximativ 1000 de trasaturi si descriptori asociati acestora. Trasaturile sunt zone de interes in imagine care pot fi folosite pentru a detecta obiecte sau gasi asocieri intre cadrele consecutive. Acestea mai sunt numite si keypoint-uri in literatura de specialitate iar librarii precum OpenCV au structuri de date dedicate pentru acestea. O trasatura poate fi interpretata matematic ca o zona in care apare o schimbare brusca a gradientului culorii. De cele mai multe ori, astfel de variatii se regasesc in zonele de frontiera dintre obiecte, deoarece apare o diferenta de culoare

si implicit una de intensitate luminoasa. Zonele slab texturate, cum ar fi cerul sau peretii in interiorul unei cladiri nu au zone care sa poata fi usor de comparat: aproape toti pixelii de pe suprafata respectiva arata similar si este greu de estimat de unde a fost extras keypoint-ul respectiv. In schimb, o camera complet mobilata ar fi o zona puternic texturata iar un algoritm de detectie de keypoint-uri ar putea sa gaseasca usor 1000 de trasaturi pe care sa le foloseasca. Daca algoritmul nu reuseste sa gaseasca suficiente keypoint-uri pentru a face urmarirea intre cadre, de obicei minim 500, urmarirea cadru cu cadru nu poate continua. Din aceasta cauza algoritmul de ORB-SLAM2 da rezultate eronate in zonele slab texturate. Daca algoritmul de extragere functioneaza corect iar traiectoria camerei este una stabila, fara schimbari bruste ale directiei de deplasare, trasaturi similare ar trebui sa fie observate in ambele imagini. Asocierile dintre ele, ne pot da informatii despre modul in care s-a deplasat camera intre cele 2 cadre. Problema este ca aceste keypoint-uri nu pot fi comparate direct intre ele, din aceasta cauza ne folosim de descriptori. Acestia sunt vectori de diferite dimensiuni care trebuie sa surprinda informatia esentiala observata in zona respectiva din imagine, in mod ideal descriptorii ar trebui sa ramana invariabili la operatiile de redimensionare si rotatie aplicate pe keypoint-uri. Algoritmului Oriented Fast and Rotated Brief (ORB) este folosit pentru extragerea de keypoint-uri si descriptori. A fost creat in anul 2011 ca alternativa pentru alti algoritmi de extragere de feature-uri precum SIFT si SURF. Motivul pentru care acesta a ajuns atat de popular se datoreaza mai multor factori:

- Este mult mai rapid decat SIFT si SURF fiind mult mai potrivit pentru sisteme in timp real si pentru dispozitive embedded.
- la momentul realizarii lucrarii stiintifice ORB-SLAM2 atat SIFT cat si SURF se aflau sub protectia drepturilor de autor, ORB nu avea o astfel de restrictie
- ORB este invariant din punct de vedere al rotatiei
- Foloseste descriptori binari, care pot fi usor de comparat folosind distanta Hamming, aceasta converteste in biti vectorul de elemente obtinute,

Implementarea algoritmului ORB poate fi separata in 2 componente, calcularea keypoint-urilor si cea a descriptorilor. Pasii pe care ii urmeaza algoritmul sunt realizati intr-un for loop. ORB extrage feature-uri la diferite dimensiuni ale imaginii, pentru a crea trasaturi mai robuste la modificarea distantei. Numarul de executii al buclei for, este acelasi cu numarul de resize-uri pe care trebuie sa le aplice algoritmul.

1. Calcularea keypoint-urilor folosind algoritmul FAST-9.
2. Selectarea celor mai potrivite keypoint-uri folosind Harris Corner Measure, Trasaturile sunt sortate in ordine descrescatoare si sunt selectate primele N cele mai potrivite
3. Pentru fiecare keypoint se calculeaza orientarea acestuia folosind intensitatea centroidului, dupa aceasta operatie avem toate informatiile necesare despre keypoint-uri.
4. Inainte de a calcula descriptorii se aplica o operatie de smoothing Gaussian pentru fiecare zona selectata de un keypoint avand  $31 \times 31$  de pixeli, folosind un kernel cu o dimensiune de  $5 \times 5$ .
5. se calculeaza descriptorii de tip steer BRIEF, avand valorile modificate dupa unghiul dat de orientare. Operatia de modificare a orientarii duce la o variatie redusa a valorilor bitilor din descriptori si la o corelatie puternica intre acestia, facand descriptorii ineficienti
6. Se obtine rBRIEF o varianta optimizata a algoritmul steer BRIEF, prin alegerea bitilor despre care se stie ca au varianta mare si grad scazut de corelatie intre ei.

In cazul algoritmul FAST-9, cifra 9 vine de la diametrul ferestrei circulare in care se face compararea intre valoarea intensitatii pixelului si centru. Acest algoritm primeste ca parametru imaginea si pragul pe care trebuie sa il depaseasca diferenta de intensitate intre pixeli pentru a fi considerat un keypoint. De cele mai multe ori, informatia data de keypoint-uri este redundanta, pentru a selecta un numar restrans de trasaturi, de preferat cele mai expresive, se foloseste Harris Corner Measure. Pentru a calcula orientarea vom defini notiunea de centroid  $C$  care este diferit de centrul zonei determinata de keypoint  $O$ . Vectorul  $\vec{OC}$  va fi cel care va da unghiul  $\theta$  al keypoint-ului pe care il vom obtine direct din urmatoarea formula, unde  $I(x, y)$  reprezinta intensitatea luminoasa a pixelului cu coordonate  $(x, y)$ .

$$m_{pq} = \sum_{x,y} x^p y^q I(x, y), \quad \theta = \text{atan2}(m_{01}, m_{10}) \quad (3)$$

In etapele 5 si 6 se realizeaza calcularea descriptorilor: acestia vor avea forma binara si o lungime finala de 256 de biti. Compararea lor se va realiza folosind distanta Hamming. Cu cat 2 descriptori au o valoarea mai mica a acestei distante, cu atat mai similari sunt. Valorile descriptorilor sunt asociate pe baza unui test binar in care este comparata intensitatea a 2 puncte din planul imaginii. Problema este ca descriptorii BRIEF sunt sensibili la schimbarile de rotatie, din aceasta cauza, prin rotirea coordonatelor pixelilor cu unghiul  $\theta$  al orientarii se obtine steered BRIEF. Pentru a obtine rBRIEF, au fost invatate in offline prin aplicarea



unui algoritm de tip Greedy, care teste de verificare a intensitatii au cea mai mare variatie, si primele 256 dintre acestea au fost alese pentru a alcatui descriptorul.

### 4.3 Harta punctelor din spatiu

Unul dintre scopurile fundamentale ale algoritmului de ORB-SLAM2, pe langa cel de estimare al traseului camerei este cel de creare a hartii locale a mediului inconjurator. Problema este ca, in comparatie cu versiuni mai avansate ale acestui algoritm, special modificate pentru o reconstructie cat mai fidela a mediului, algoritmul nostru trebuie sa functioneze pentru un sistem embedded care nu are capacitate de procesare suficient de mare, fiind nevoit astfel sa simuleze mediul printr-un nor de puncte cu o densitate redusa (sparse). Cele 2 sarcini sunt dependente una de cealalta, fiecare element din norul de puncte actioneaza ca o referinta, o caracteristica a mediului care ar trebui sa fie observata de fiecare data cand punctul se afla in frustum-ul camerei. De exemplu: presupunem ca avem o imagine in care este observata in totalitate o masa in interiorul unei incaperi. ORB va identifica aproape instantaneu feature-urile (colturile mesei) si teoretic, indiferent de modul in care ne-am rotit in jurul mesei, aceleasi feature-uri ar trebui sa fie observate de fiecare data, mai mult de atat, considerand ca mediul este static, acestea sunt mereu asociate cu acelasi punct din spatiu, devenind astfel o referinta pe baza careia putem estima modul in care s-ar deplasa camera. In literatura de specialitate, aceste puncte din spatiu sunt denumite MapPoint-uri iar functionalitatea corecta a algoritmului depinde strict de modul in care aceste MapPoint-uri sunt observate cadru cu cadru. Un astfel de punct in spatiu este creat dintr-un keypoint, dar nu vom avea nevoie de toate punctele din regiunea respectiva si vom considera ca centrul este punctul cel mai semnificativ, avand coordonate  $x$  si  $y$ , si distanta fata de camera fiind estimata ca fiind  $d$ . Ne vom folosi de matricea transformarii din coordonatele camerei in coordonatele globale si de parametrii interni ai camerei  $f_x$ ,  $f_y$  distanta focala, si  $c_x$ ,  $c_y$  coordonatele centrului imaginii. Vectorul coloana cu 3 dimensiuni reprezinta pozitia in spatiu a feature-ului gasit in cadrul curent, astfel am creat primul MapPoint. Ca alternativa, pentru a nu lucra cu matrici de dimensiuni  $4 \times 4$  putem folosi  $R_{wc}$  reprezentand matricea de rotatie si  $t_{wc}$  vectorul de translatie.

$$\begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} = T_{wc} * \begin{bmatrix} \frac{x-c_x}{f_x} * d \\ \frac{y-c_y}{f_y} * d \\ d \\ 1 \end{bmatrix}, \quad \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = R_{wc} * \begin{bmatrix} \frac{x-c_x}{f_x} * d \\ \frac{y-c_y}{f_y} * d \\ d \end{bmatrix} + t_{wc} \quad (4)$$

In literatura de specialitate MapPoint-urile sunt considerate ca fiind niste ancore (landmark) pozitionate dinamic de catre algoritm, acestea sunt asociate cu un anumit cadru cheie si ne vor ajuta in optimizarea matricei de pozitie dar si pentru sarcina de relocalizare si de memorare a zonelor cunoscute.

#### 4.4 Asociere puncte din spatiu cu feature-uri ORB

Ca date de intrare avem feature-urile si descriptorii extrasi din imagine, matricea de adancime si harta de MapPoint-uri. Scopul acestei componente este sa gaseasca cat mai multe asocieri de 1:1 intre feature-uri si MapPoint-uri. Intr-un caz ideal fiecare feature gasit ar trebui sa aiba asociat un MapPoint, dar in realitate nu se poate intampla acest lucru din 2 motive: imperfectiuni ale algoritmului ORB de detectie ale feature-urilor: acesta nu garanteaza ca acelasi feature va fi gasit de fiecare data pentru cadre consecutive si faptul ca modelul isi schimba orientarea, facand ca MapPoint-urile aflate la limita campului vizual al camerei sa nu mai poata fi observate. Un MapPoint este un feature al unui cadru anterior, proiectat in spatiu. In final, aceasta componenta realizeaza tot o comparare de feature-uri intre cadrul curent, si multiple cadre anterioare. Aceasta operatie de comparare se realizeaza prin intermediul distantei Hamming dintre descriptori, cu cat valoarea obtinuta este mai mica, cu atat cele 2 feature-uri sunt mai asemanatoare. Exista mai multe tipuri de algoritmi folositi pentru feature matching, dar cel folosit in implementarea curenta este Brute Force Feature Matching optimizat. Acest algoritm primeste ca date de intrare 2 seturi de feature-uri si incearca sa gaseasca asocieri intre ele. Asocierele sunt facute cu ajutorul descriptorilor, se calculeaza distanta Hamming iar daca valoarea obtinuta este minima, perechea respectiva de feature-uri se considera ca a fost corect asociata. Pentru ORB-SLAM2 o potrivire intre 2 keypoint-uri arata ca ele se refera la exact acelasi punct din spatiu, observat din 2 imagini diferite. Daca  $N$  este numarul de feature-uri din primul set,  $M$  numarul de feature-uri din al doilea set si  $D$ , dimensiunea descriptorului, in cazul nostru ORB este 32, complexitatea algoritmului devine

$O(N * M * D)$ . Destul de costisitor de folosit pentru un sistem in timp real, mai mult de atat este predispus la erori, compararea feature-urilor nu tine cont de locatia acestora in imagine, obtinandu-se astfel asocieri care matematic par corecte, dar ele nu au sens din punct de vedere logic. Pentru a rezolva aceasta problema si a reduce complexitatea temporala se stabileste o fereastră circulară de dimensiune prestabilită in jurul punctului de proiectie unde se pot cauta feature-uri. O data ce 2 keypoint-uri au fost considerate ca facand referinta la acelasi punct din spatiu, cadrului curent ii este asociat un nou MapPoint.

## 4.5 Optimizare Estimare Pozitie Initiala

Aceasta componenta primeste ca data de intrare estimarea pozitiei curente a camerei  $T_{cw}$ , si o asociere bijectiva intre feature-urile gasite in imagine si punctele care exista la momentul respectiv in spatiu. Ca date de iesire vom avea doar matricea pozitiei curente a camerei optimizata. Daca asocierile intre feature-uri si MapPoint-uri sunt perfecte, ar trebui ca proiectia punctului din spatiu pe imagine sa se suprapuna pe centrul keypoint-ului. Rareori se petrece acest lucru in practica, iar distanta dintre proiectia unui MapPoint si coordonatele centrului feature-ului reprezinta eroarea de asociere. Pentru a minimiza aceasta eroare, exista 2 optimizari care se pot face: prima este modificarea valorilor matricei de pozitiei, iar cea de-a doua este modificarea coordonatelor din spatiu ale MapPoint-ului. Inainte de a prezenta algoritmul de optimizare folosit, voi arata modul in care se proiecteaza un MapPoint in plan.

### 4.5.1 Proiectarea MapPoint in planul imaginii

Aceasta operatie de proiectie poate fi vazuta ca aplicarea unui functii  $\pi(\cdot)$  ce primeste ca date de intrare coordonatele globale ale punctului, iar ca rezultat va returna coordonatele omogene in planul imaginii. Aceasta transformare se petrece in 2 etape:

1. conversia din sistemul de coordonate globale in sistemul de coordonate al camerei
2. conversia din sistemul de coordonate al camerei in sistemul de coordonate al imaginii

In prima etapa putem folosi coordonatele omogene, pentru a face conversia in mod direct. Alternativ, putem extrage din matricea de pozitie  $T_{cw}$  atat matricea de rotatie  $R_{cw}$  cat si vectorul coloana de translatie  $t_{cw}$ .

$$\mathbf{X}_{camera} = \mathbf{T}_{cw} \cdot \begin{bmatrix} \mathbf{X}_w \\ 1 \end{bmatrix}, \quad \mathbf{T}_{cw} = \begin{bmatrix} \mathbf{R}_{cw} & \mathbf{t}_{cw} \\ \mathbf{0}^T & 1 \end{bmatrix}, \quad \mathbf{X}_{camera} = \mathbf{R}_{cw} \cdot \mathbf{X}_w + \mathbf{t}_{cw} \quad (5)$$

Matricea  $\mathbf{T}_{cw}$  este utilizata atat pentru a descrie pozitia si orientarea in spatiu cat si pentru a schimba din sistemul de coordonate global in cel al camerei. In sistemul de coordonate global, un punct se afla la exact aceeasi valoare indiferent de pozitia camerei care il priveste, in sistemul de coordonate al camerei, pozitia unui MapPoint o sa difere de fiecare data. In etapa a doua MapPoint-ul este in sistemul de referinta al camerei, coordonatele fiind reprezentate prin vectorul coloana  $\mathbf{X}_{camera}$ . Vom considera a 3-a valoare a acestui vector  $Z_c$ . Aceasta reprezinta distanta dintre planul camerei si punctul pe care il analizam.  $Z_c$  ne spune daca punctul respectiv poate fi observat in imagine. Daca valoarea  $Z_c$  este mai mica sau egala cu 0, inseamna ca punctul se proiecteaza in spatele camerei, facandu-l invalid. In situatia in care  $Z_c$  este mai mare decat 0, vom realiza conversia in coordonatele omogene ale imaginii cu ajutorul urmatoarei formule,  $u$  fiind asociat axei x si  $v$  fiind asociat axei y. Daca valorile  $u$  si  $v$  au valori mai mari ca 0 si mai mici decat dimensiunea imaginii, vectorul coloana  $[u, v, 1]$  este rezultatul cautat.

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \mathbf{K} \cdot \begin{bmatrix} \frac{X_c}{Z_c} \\ \frac{Y_c}{Z_c} \\ 1 \end{bmatrix}, \quad \mathbf{K} = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \quad (6)$$

## 4.5.2 Motion Only Bundle Adjustment

Algoritmul folosit in aceasta etapa se numeste Motion Only Bundle Adjustment. Acesta modifica doar matricea pozitiei curente a camerei. Coordonatele punctelor din spatiu sunt considerate ca fiind constante. Algoritmul este unul iterativ, minimizand o functie de cost. Forma generala a functiei de cost este suma erorilor de proiectie pentru toate perechile (feature, MapPoint). Iar formula generala este aceasta.

$$\mathbf{R}_{cw}, \mathbf{t}_{cw} = \min_{\mathbf{R}_{cw}, \mathbf{t}_{cw}} \sum_{i=1}^N \rho(\|\mathbf{x}_i - \mathbf{K} \cdot (\mathbf{R}_{cw} \cdot \mathbf{X}_i + \mathbf{t}_{cw})\|^2) \quad (7)$$

In aceasta formula,  $\mathbf{x}_i$  reprezinta coordonatele omogene feature-ului in sistemul de coordonate al imaginii iar  $\mathbf{X}_i$  reprezinta coordonatele globale ale MapPoint-ului pentru care calculam

eroarea de proiectie. Simbolul  $\rho(\cdot)$  reprezinta functia Huber pentru scalarea valorilor de eroare. Daca o asociere intre un feature si un MapPoint nu este potrivita, diferenta dintre centrul feature-ului si proiectia MapPoint-ului este mai mare decat un prag prestabilit. Aceasta diferenta, lasata nemodificata, ar destabiliza algoritmul. Iar o astfel de problema este usor de observat, daca modificarea matricei de pozitie duce la variatii enorme a orientarii sau a translaticii intre 2 cadre consecutive, atunci cel mai probabil asocierile intre feature-uri si MapPoint-uri aveau valori eronate. Termenul de *outlier* este folosit pentru a descrie o pereche incorecta. Functia de pierdere Huber reduce valoarea acestor outlier-ere permitandu-le in acelasi timp sa faca parte din algoritmul de optimizare. In acest fel algoritmul devine mai robust si capabil ajunga la o valoare optima in mai putine iteratii. Mai jos este prezentata formula matematica a functiei Huber Loss, unde  $\delta$  reprezinta un numar real pozitiv, toleranta a erorii de proiectie.

$$\rho(s) = \begin{cases} \frac{1}{2}s^2 & \text{if } |s| \leq \delta \\ \delta(|s| - \frac{1}{2}\delta) & \text{if } |s| > \delta \end{cases} \quad (8)$$

In urma executiei algoritmului obtinem matricea de pozitie optimizata, mai mult de atat, stim care dintre perechile (feature, MapPoint) au avut statutul de outlier si le putem elimina pentru a nu influenta in mod negativ functionalitatea algoritmului.

## 4.6 Crearea unui cadru cheie

Aceasta componenta primeste ca date de intrare absolut toate informatiile procesate de pana acum pentru cadrul curent: imaginea de tip rgb, matricea de adancime, punctele cheie, descriptorii, asocierile (feature, MapPoint) si matricea estimarii pozitiei. Toate acestea impreuna vor alcatui un cadru cheie care va fi salvat in memorie. Salvarea pozitiilor cadrelor anterioare ne poate ajuta in 2 feluri. Putem folosi doar 2 cadre anterioare pentru a estima pozitia celui care urmeaza bazandu-ne pe legea inertiei. Consider ca o data inceputa deplasarea camerei intr-o anumita directie, este foarte probabil ca aceea miscare sa fie mentinuta si la urmatorul cadru. Fie  $T_{cw}$  matricea de pozitie pentru cadrul la care vrem sa estimam deplasarea, iar  $T_{cw1}, T_{cw2}$  matricile de pozitie a celor 2 cadre imediat predecesoare. Formula de estimare a pozitiei curente este:

$$\mathbf{T}_{cw} = \mathbf{T}_{cw1} \cdot (\mathbf{T}_{cw2}^{-1} \cdot \mathbf{T}_{cw1}) \quad (9)$$

Al doilea motiv pentru care avem nevoie de cadre cheie este recreerea mediului si a traseului parcurs. Incercam sa salvam numarul minim de cadre necesare pentru a reproduce harta de MapPoint-uri a mediului inconjurator. Un cadru cheie nou (KeyFrame) aduce cu sine MapPoint-uri noi, extrase din feature-urile gasite in imaginea respectiva. Functionarea corecta a urmarii cadru cu cadru, este determinata de numarul de MapPoint-uri gasite in imaginea curenta in comparatie cu un cadru de referinta. In momentul in care numarul de puncte cheie gasite in imaginea curenta scade sub un anumit prag, stim ca este necesar un nou cadru cheie care: sa stabilizeze urmarirea, sa introduca noi MapPoint-uri, si sa ajute la optimizarea intregii harti a mediului.

#### 4.6.1 Optimizare harta locala

Harta locala este alcatuita din KeyFrame-uri si MapPoint-uri. Pentru a optimiza harta trebuie sa adaugam noi puncte de tip MapPoint si noi KeyFrame-uri in ea. Pentru a valida conexiunile care deja exista. Modul in care sunt create si sterse punctele urmeaza o abordare numita survival of the fittest. La fiecare nou KeyFrame adaugat, sunt create in aproximativ 100 de noi MapPoint-uri si inserate in harta. Acestea vor fi supuse unui test care sa evalueze cat de usor sunt observate feature-urile pe care le reprezinta. Din punct de vedere matematic, un MapPoint este observat de un KeyFrame daca proiectia acestuia in imagine este un vector valid in sistemul de coordonate al KeyFrame-ului respectiv. Cu cat mai multe Keyframe-uri observa acelasi MapPoint, cu atat mai stabil este punctul respectiv din spatiu. Intr-un caz ideal, ar trebui ca orice MapPoint creat sa fie stabil. De cele mai multe ori nu se intampla acest lucru din cauza erorilor de feature matching. Astfel, doar cele mai evidente feature-uri raman salvate in harta pana la finalul algoritmului. Scopul acestei componente este adaugarea KeyFrame-urilor noi, eliminarea celor redundante si testarea stabilitatii tuturor MapPoint-urilor create. Pentru a salva Keyframe-ul curent in harta urmatoarele operatii trebuie urmate:

1. Folosind harta de adancime, sunt selectate cele mai apropiate  $N$  feature-uri care nu au un MapPoint asociat si au valoarea adancimii mai mare ca 0. Coordonatele acestor feature-uri sunt proiectate in spatiu pentru a obtine noi MapPoint-uri.

2. KeyFrame-ul curent este comparat cu alte cadre cheie, pentru a vedea cu cine imparte cele mai multe puncte comune. Keyframe-urile sunt stocate in harta intr-o structura de tip graf neorientat unde nodurile sunt cadrele cheie iar arcele sunt numarul de MapPoint-uri comune dintre ele.
3. sunt eliminate punctele cheie redundante sau care au fost observate in prea putine cadre pentru a fi luate in considerare.
4. se executa un algoritm numit Local Bundle Adjustment in care KeyFrame-urile care au cele mai multe puncte comune cu cadrul curent analizat vor avea matricea de pozitie optimizata si coordonatele globale ale MapPoint-urilor asociate acestora.

Local Bundle Adjustment este similar cu Motion Only Bundle Adjustment. In continuare vorbim de un algoritm iterativ care incearca sa minimizeze o functie de cost, folosind metoda scaderii gradientului. Diferenta este ca optimizarea se aplica pe mai mult de un cadru cheie: atat matricea pozitiei si punctele din spatiu asociate (MapPoint-urile) vor fi optimizate. Avem urmatoarele etape:

1. Se creeaza lista de cadre mobile. Plecand de la cadrul curent, se vor selecta toti vecinii de gradul 1 si 2 din graful neorientat stocat in harta. Aceste Keyframe-uri sunt considerate *mobile* deoarece matricea lor de pozitie se va modifica.
2. Se creeaza lista de MapPoint-uri ale caror coordonate vor fi optimizate. Fiecare Keyframe din multimea cadrelor mobile observa un numar de puncte in spatiu, toate aceste puncte vor fi folosite de catre algoritmul de optimizare.
3. Se creeaza lista de cadre fixe pentru care matricea de pozitie nu se va modifica. Pentru fiecare punct din lista de MapPoint-uri ce vor fi optimizate se va itera prin lista de KeyFrame-uri care observa acel MapPoint. Daca un KeyFrame apartine multimii de cadre mobile va fi ignorat, iar daca nu, va fi adaugat in lista de cadre fixe. Acestea sunt incluse in algoritm pentru a garanta ca modificarea coordonatelor unui MapPoint nu va strica asocierea (feature, MapPoint) in cadrele care nu vor avea matricea de pozitie modificata.

Lucrarea stiintifica care sta la baza ORB-SLAM2, implementeaza deja functia de cost pe care algoritmul de Local Bundle Adjustment o foloseste. Pentru a intelege mai usor formula matematica, aceasta trebuie privita de la dreapta la stanga.  $E_{k,j}$  reprezinta eroarea de proiectie

a unui MapPoint pe feature-ul asociat. Indicele  $k$  apartine KeyFrame-ului,  $j$  reprezinta ordinul perechii (feature, MapPoint) pentru care calculam eroarea. Simbolul  $\rho(\cdot)$  este asociat functiei Huber, folosita pentru a ameliora efectele perechilor de tip outlier.  $X_k$  este o notatie pentru multimea tuturor asocierilor (feature, MapPoint) pentru un Keyframe  $k$ . Suma erorilor tuturor perechilor este calculata pentru toate cadrele fixe si mobile. Parametrii care vor fi optimizati sunt: coordonatele MapPoint-urilor selectate de catre algoritmul  $X_i$  cat si matricile de pozitie pentru cadrele mobile  $K_l$ . Algoritmul optimizeaza valorile pana cand ajunge la o valoare de minim sau pentru un numar de iteratii.

$$\{\mathbf{X}_i, \mathbf{R}_l, \mathbf{t}_l \mid i \in \mathcal{P}_L, l \in \mathcal{K}_L\} = \arg \min_{\mathbf{X}_i, \mathbf{R}_l, \mathbf{t}_l} \sum_{k \in \mathcal{K}_L \cup \mathcal{K}_F} \sum_{j \in \mathcal{X}_k} \rho(E_{kj}) \quad (10)$$

$$E_{kj} = \|\mathbf{x}_j - K \cdot (\mathbf{R}_k \mathbf{X}_j + \mathbf{t}_k)\| \quad (11)$$

## 4.6.2 Reteaua Neurala FastDepth

Retele Neurale Artificiale sunt o tehnica des intalnita in Machine Learning pentru a rezolva sarcini complexe pentru care nu exista solutii algoritmice clar definite sau implementarea acestora este mult prea costisitoare. Retele neurale definesc o functie nonliniara care gaseste o corespondenta intre un set multivariat de date de intrare  $x$  si un set multivariat de date de iesire  $y$ , modificand un set de parametri  $\phi$ ,  $f(x, \phi) = y$ . Aceasta este alcatuita dintr-un numar enorm de elemente de procesare care contin parametrii functiei  $\phi$ , conectate intre ele intr-o structura de tip graf si dispuse pe straturi. Cele mai importante fiind: stratul de intrare si de iesire, unde se stabileste forma generala pe care trebuie o sa respecte datele care vor parcurge reteaua si modul in care va arata rezultatul obtinut. Celelalte nivele sunt denumite straturi ascunse. Acestea fac prelucrarea informatiei primite de la straturile anterioare si o transmit mai departe. Spunem ca o retea neurala invata din datele primite, daca isi modifica parametrii  $\phi$  astfel incat sa reprezinte cu mai multa acuratete corespondenta intre datele de intrare  $x$  si cele de iesire  $y$ . FastDepth este o arhitectura de retea neurala folosita pentru estima adancimii in imagini. Aceasta primeste o imagine de tip RGB al interiorului unei incaperi si returneaza o matrice cu valori in intervalul  $0m - 10m$  estimand pentru fiecare pixel in parte distanta de la planul de proiectie al imaginii pana la punctul din spatiu surprins de fotografie. Scopul nostru este antrenarea unei retele neurale care sa produca o matrice de adancime cu valori cat mai



aproprate de distanta reala la care se afla obiectele fata de camera. ORB-SLAM2 foloseste o camera tip RGBD / Stereo care descrie cu foarte mare acuratete distanta pana intr-un anumit punct din spatiu, dar creeaza o matrice rara de valori, suprafetele lucioase sau cele care nu au putut fi clar observate vor avea adancimea 0 pentru a arata ca distanta nu a putut fi corect estimata in pixelii respectivi. Un motiv pentru care retelele neurale sunt o alternativa buna este ca ele vor avea o estimare pentru fiecare pixel din imagine. Reteaua FastDepth pare sa realizeze o pseudosegmentare a zonelor din imagine, identifica conturul obiectelor si atribuie valori ale distantei asemanatoare pentru pixelii ce apartin aceleasi entitati. Arhitecturile de mari dimensiuni nu pot fi folosite in timp real fara a utiliza un GPU, dar nu este cazul si pentru arhitectura FastDepth care poate procesa aproximativ 100 de cadre pe secunda. De asemenea consuma o cantitate redusa de memorie, parametrii retelei pot fi stocati intr-un fisier ONNX ce ocupa mai putin de 8 MB, fiind usor de integrat intr-un dispozitiv embedded.

Un posibil dezavantaj al acestei arhitecturi este limitarea de 10m, fiind nepotrivit de folosit afara, dar ideal pentru un spatiu inchis de mici dimensiuni. Un alt dezavantaj este ca valorile approximate vor avea o acuratete mai slaba decat cele obtinute de camerele Stereo/RGBD.

Exista mai multe filozii cand vine vorba de modul in care ar trebui sa arate arhitectura retelelor neurale si operatiile pe care ar trebui sa le realizeze fiecare strat. Feed forward neural network a fost printre primele arhitecturi definite. Elementele de procesare sunt dispuse pe straturi, si fiecare strat primeste input-ul de la stratul precedent si transmite output-ul la stratul imediat urmator. Informatia circula liniar, de la intrarea in retea pana la finalul acesteia. Abordarea s-a dovedit eficienta in situatiile in care era nevoie de retele neurale de mici dimensiuni, cu un numar redus de straturi si parametrii. In momentul in care crestea complexitatea, abordarea de feed forward neural network devenea greu de antrenat si dadea rezultate mai slabe. Pentru a rezolva aceasta problema au aparut arhitecturile de tip residual network. Acestea au aplicatii in procesarea imaginilor, unde datele de intrare au dimensiuni mari si este nevoie de multe nivele pentru a extrage suficiente informatii. Principiul de functionare este utilizarea unor straturi reziduale denumite si skip connections, in care rezultatul unui strat este salvat si transmis ca data de intrare la un alt nivel decat la cel imediat urmator. Abordarea aceasta pastreaza din informatiile initiale ale datelor de intrare in straturile viitoare stabilizand antrenarea. O alta arhitectura des intalnita este cea de encoder-decoder folosita in numeroase aplicatii practice in ceea ce priveste imaginile: sarcini de colorare a imaginilor gri, reconstructie a imaginilor care contin parti lipsa si de generare de imagini: un exemplu fiind Variational Auto Encoder.

Pe langa tipurile de arhitecturi propuse, exista mai multe categorii de straturi in retele neurale. Primele folosite erau cele fully connected unde fiecare element de procesare era conectat cu toate celelalte elemente de procesare din stratul urmator. Matematic, operatia poate fi vazuta ca o inmultire de matrici, o operatie costisitoare, iar utilizarea exclusiva a straturilor complet conectate creea o retea neurala incapabila sa reprezinte functii nonliniare, scazand capacitatea de generalizare. De cele mai multe ori straturile liniare sunt folosite impreuna cu functii de activare nonliniare precum ReLU sau Sigmoid dar in continuare ramane problema numarului mare de parametrii care trebuie antrenati. Din aceasta cauza au fost create straturile convolutionale care folosesc mai putini parametrii si au aplicabilitate in procesarea imaginilor. Principiul teoretic pe care se bazeaza este ca pixelii alaturati in imagine au aceeasi semnificatie, reprezentand acelasi feature. Operatia de convolutie trebuie realizata pe o zona a imaginii iar modificarea parametrilor afecteaza output-ul generat de mai multi pixeli. Se stabileste un kernel, o matrice de mici dimensiuni, in FastDepth folosindu-se kernel-uri de (3, 3), acestea vor stoca parametrii  $w_{mn}$  pe care reteaua neurala ii va antrena pentru stratul convolutional. In formula  $h_{ij}$  reprezinta intensitatea pixelului dupa calculul operatiei de convolutie, iar  $x_{ij}$  este valoarea intensitatii pixelului de pe coloana  $i$  si linia  $j$ . Litera  $a$  reprezinta functia de activare folosita iar  $\beta$  este o valoare numerica denumita bias. Acesta poate fi modificat in timpul antrenarii si creste capacitatea de generalizare a functiei de convolutie.

$$h_{ij} = a \left[ \beta + \sum_{m=1}^3 \sum_{n=1}^3 \omega_{mn} x_{i+m-2, j+n-2} \right] \quad (12)$$

Stratul de convolutie este in continuare prea costisitor pentru a crea o arhitectura in timp real de mari dimensiuni. Presupunem ca avem un vector de intrare pentru un strat de convolutie cu dimensiunile  $[d_{in}, h, w]$  unde  $d_{in}$  este numarul de canale,  $h$  inaltimea si  $w$  latimea vectorului. Kernelul folosit are dimensiunile  $[k, k, d_{in}, d_{out}]$ , unde  $d_{out}$  este numarul de canale rezultate in urma convolutiei. In total se vor executa  $h \cdot w \cdot d_{in} \cdot d_{out} \cdot k \cdot k$  operatii. Pentru a rezolva aceasta problema a fost creat un strat numit Depthwise Separable Convolutions, obtinut prin compunerea a 2 straturi de convolutie, unul numit depthwise convolution, iar celalalt pointwise convolution. Aceasta abordare creste viteza de procesare si imbunatateste acuratetea in sarcini de clasificare pentru seturi de date precum ImageNet ILSVRC2012. In cazul depthwise convolution, fiecare canal al datelor de intrare este procesat de un singur kernel al stratului de convolutie. Pointwise convolution uneste printr-o combinatie liniara rezultatul procesarii fiecarui canal. Complexitatea temporală obtinuta astfel este de:  $h \cdot w \cdot d_{in} \cdot (k^2 + d_{out})$ .

FastDepth foloseste tehnica de skip connections, straturile finale primind ca date de intrare valorile calculate de straturile aflate la inceput, si urmeaza o arhitectura encoder-decoder. Encoder-ul transforma datele de intrare intr-o forma mai compacta asemeni unei operatii de arhivare. Aceasta este realizata folosind o alta retea neurala numita Mobile\_Net si ulterior Mobile\_Netv2 care reduce numarul de parametri si creste viteza de procesare fara a impacta acuratetea. Partea de decoder este alcatuit din 5 straturi de tip depthwise convolution fiecare urmate de o interpolare liniara care dubleaza dimensiunea rezultatului, ultimul strat fiind un pointwise convolution care uneste canalele obtinute si returneaza matricea de adancime. Pentru Mobile\_Netv2 exista parametrii preantrenati in biblioteca Pytorch pe setul de date ImageNet fiind un motiv in plus de a folosi aceasta arhitectura in dezvoltarea FastDepth. Mobile\_Netv2 foloseste atat depthwise convolution cat si pointwise convolution intr-un strat numit Inverted Residual, acesta fiind alcatuit din urmatoarele componente unde  $t$  este factorul de multiplicare al numarului de canale,  $s$  este parametrul de stride, determina daca se micsoreaza numarul de canale,  $h, w, d$  sunt dimensiunile matricei de intrare.

**DE INSERAT AICI IMAGINEA**

## 5 DETALII DE IMPLEMENTARE

### 5.1 Limbaje de programare si librarii folosite

Implementarea este realizata in C++17. Pentru management-ul librariilor si al codului folosesc CMake 3.28.3. Acesta imi permite sa grupez in foldere codul scris de mine si face operatia de linking automat cu binarele pachetelor folosite. Librariile principale sunt OpenCV 4.9.0, Ceres 2.2.0, Eigen 3.4.0, DBoW2 si ultima versiune de Sophus pana la data de ianuarie 2025. In comparatie cu alte librarii care inca mai trec prin diverse update-uri, Sophus a intrat intr-o etapa de mentenanta, dezvoltarea efectiva a acestuia fiind finalizata din iunie 2024. O prima problema pe care am intalnit-o a fost gasirea unei versiuni compatibile de C++ cu toate aceste pachete. Am incercat mai multe variante printre care C++11, C++14, C++17 si C++20. Preferinta mea ar fi fost sa folosesc o versiune cat mai noua cu putinta, dar care sa poata fi compatibila cu toate librariile mentionate. C++11 si C++14 nu erau compatibile cu Ceres, versiunea minima pentru aceasta librarie era C++17. C++20 si C++23 nu era compatibil cu Sophus si cu Eigen, iar ambele librarii sunt fundamentale deoarece implementeaza metode puternic optimizate de a lucra cu matrici iar API-ul lor era mai simplu decat cel din OpenCV. Singura optiune ramasa a fost C++17 care era incompatibila cu DBoW2. Libraria folosea o versiune mai veche a functiei throw pentru erori. In momentul in care am eliminat aceasta directiva, am putut recompila codul ca librarie. Bibliotecile utilizate sunt urmatoarele:

OpenCV este o librarie de computer vision. Contine implementari ale algoritmilor de extragere de trasaturi precum FAST, ORB, SIFT, SURF, API-uri pentru procesare video: citirea unui video cadru cu cadru, procesarea de imagini: aplicarea de filtre, transformarea in grayscale, eliminarea distorsiunii cauzata de camera. Foarte importante sunt structurile ce abstractizeaza matricile si parametrii prin care se indentifica trasaturile: `cv::Mat` si `cv::KeyPoint`. Structura `KeyPoint` este fundamentala pentru implementarea algoritmului deoarece stocheaza nume-roase informatii despre zona pe care o reprezinta: orientarea acesteia, coordonatele centrului si nivelul la care a fost observat, parametrii de care am avut nevoie in fiecare etapa de procesare a cadrelor. Pe langa aceste lucruri, OpenCV are un modul dedicat pentru citirea parametrilor

retelelor neurale din fisierele care urmeaza un format de tip ONNX, fiind o alternativa potrivita daca vreau sa utilizez un model doar pentru sarcini de inferenta.

Ca librarie de optimizare am avut de ales intre Ceres si g2o. In implementarea oficiala g2o era cel folosit. Motivul principal fiind ca permite abstractizarea parametrilor care trebuie optimizati si a relatiilor dintre acestia sub forma unui graf neorientat. API-ul de g2o permite activarea si dezactivarea anumitor noduri, pentru a face implementarea mai robusta impotriva perechilor de tip outlier, si pentru a putea reintroduce noduri eliminate temporar in graful de optimizare. Ceres din pacate nu permite acest lucru. O data create conditiile initiale acestea pot fi dezactivate si nu mai este permisa reutilizarea lor in aceeaasi problema de optimizare. La finalizarea algoritmului memoria folosita de catre noduri este eliberata. Cu toate acestea, Ceres are un API usor de utilizat si are o viteza comparativa cu cel din g2o.

Folosesc libraria Eigen deoarece este mai simplu API-ul de calcul cu matrici decat cel din OpenCV. Pentru a accesa elementele unei matrici in OpenCV se foloseste o referinta la vectorul de date facand accesarea elementelor mult mai nesigura iar verificarea indicelui este facuta la runtime. In cazul matricilor din Eigen, accesarea elementelor si operatiile cu matrici sunt verificate la compile time, prevenind astfel erorile inainte de a rula programul.

Sophus este o librarie care imi permite sa lucrez cu algebra de tip Lie. In loc de a vedea estimarile pozitiei ca pe niste matrici de  $4 \times 4$ , le pot vedea ca pe un vector alcatuit din 7 elemente. Primii 4 parametrii alcatuiesc un quaternion, aceasta fiind o exprimare vectoriala a unei matrici  $3 \times 3$  de rotatie, iar ultimii 3 parametrii reprezinta un vector de translatie. Biblioteca implementeaza operatii care imi permit sa lucrez cu acesti vectori, care fac parte dintr-un grup numit  $se(3)$  si garanteaza ca rezultatul obtinut este scalat corespunzator pentru a face parte in continuare din aceeaasi categorie.

DBoW2 este o metoda de tip bag of words pentru compararea imaginilor intre ele. Este utilizat pentru operatii precum feature matching intre imagini consecutive, relocalizari si recunoasterea zonelor prin care a trecut pentru a inchide buclele create de mai multe cadre cheie salvate in harta. Acesta este alcatuit dintr-o structura de tip arbore. Fiecare nivel este obtinut din realizarea unui algoritm de clusterizare a descriptorilor de tip ORB ca de exemplu kmeans++, separarea tuturor descriptorilor in functie de centroizi si reluarea aceleasi operatii in fiecare dintre clusterelor nou create. Nodurile de tip frunza sunt alcatuite dintr-un singur descriptor. Construirea arborelui se realizeaza intr-o etapa offline. In cazul libreriei DBoW2, setul de date folosit a fost Bovis 2008-09-01. Au fost alese 10K imagini iar pentru fiecare cadru in parte

extrasi 1000 descriptori ORB. Acestia au fost folositi pentru a crea un arbore de adancimea 6 iar numarul de clustere pe care le creeaza fiecare iteratie a algoritmului kmeans++ este de 10. Pe ultimul strat exista un milion de frunze, si tot aceeasi lungime o va avea si vectorul de feature-uri care va reprezenta o imagine. Fiecare descriptor va primi o valoare numerica numita greutate, invers proportionala cu frecventa pe care o are acesta. Cu cat este mai rar un anumit descriptor, cu atat este mai util pentru a diferenta o imagine de multe altele. Scopul principal al librăriei este sa primeasca ca data de intrare descriptorii ORB ai unei imagini si sa calculeze vectorul sau bag-of-words. Vectorul bag-of-words este alcatuit in principal din valori de 0. Fiecare descriptor al imaginii parcurge arborele de la radacina spre frunze, parcurgerea realizandu-se prin calcularea distantei Hamming dintre descriptor si toate nodurile de pe un anumit nivel, si alegerea nodului cu distanta Hamming minima. Nodul frunza la care va ajunge va avea asociat un index, in cazul de fata cu valori de la 0 la un milion. La acelasi index va fi modificata valoarea din vectorul bag-of-words in valoarea greutatii descriptorului stocat in arbore. Apelul de biblioteca returneaza de asemenea un vector de feature-uri in care fiecare element este o pereche de forma *(int, vector\_descriptori)*, primul element este indexul clusterului de la nivelul 4 al arborelui DBOW2, iar cel de-al doilea element reprezinta un vector de descriptori din imaginea curenta care se potrivesc in acelasi cluster. Doua imagini pot fi comparate intre ele prin intermediul acestui vector de feature-uri, lucru care va fi detaliat in descrierea clasei OrbMatcher. In implementarea ORB-SLAM2, nu este practica creerea unui vector de tip bag of words cu un milion de elemente, mai ales ca majoritatea valorilor sunt 0, asa ca o reducere a dimensionalitatii vectorului ar creste viteza de calcul a sistemului. Din aceasta cauza, compararea descriptorilor se realizeaza doar pana la nivelul 4 in arbore, vectorul bow avand doar 1000 de elemente, iar cel de feature-uri avand acelasi numar de elemente cu numarul de descriptori.

## 5.2 Mediu de lucru si principalele clase

Structura de fisiere este una simpla, in folderul radacina se regaseste fisierul de CmakeLists.txt care va fi interpretat de utilitarul cmake pentru a genera automat Makefile-ul. Acest Makefile va contine regulile de build si de clean pentru proiect. In fisierul main.cpp vor fi initializate componentele si se va putea selecta pe care dintre cele 2 seturi de date se va aplica algoritmul. Aceste seturi de date contin de fapt cadrele dintr-un video facut cu o camera RGBD Micro-

soft Kinetic impreuna cu matricile de adancime si pozitiile acestora in spatiu pentru fiecare cadru in parte. Aceste seturi de date sunt suficient de complexe pentru a permite evaluarea functionarii algoritmului de ORB-SLAM2. Tot in main.cpp, se va realiza citirea fisierului ORBvoc.txt, acesta contine datele pe care le va folosi clasa ORBVocabulary pentru a initializa arborele folosit de biblioteca DBOW2.

Tot in folderul radacina se regaseste si fisierul fast\_depth.onnx, in care este stocata arhitectura si parametrii retelei neurale FastDepth pentru estimarea adancimii. In folderul de include se afla antetele claselor pe care le voi implementa si in folderul de src se regaseste codul de C++ ce implementeaza logica programului. Am observat ca separarea codului in acest fel este o practica des intalnita in proiectele de mari dimensiuni si garanteaza flexibilitate in includerea dependintelor intre fisiere. Algoritmul ORB-SLAM2 este unul complex, depinzand de o multitudine de parametri care pot influenta acuratetea. Cei mai importanti sunt cei corelati cu camera. In fisierul config.yaml se regaseste matricea  $K$ , parametrii de distorsiune ai imaginii si alte constante pe care le-am considerat ca fiind niste hiperparametrii ai algoritmului. Acestia vor trebui modificati in functie de mediul in care va rula ORB-SLAM2 pentru a garanta functionarea corecta.

Clasa TumDatasetReader este responsabila de achizitia de date si de scrierea in fisier a traiectoriei pe care o estimeaza algoritmul cadru cu cadru. Achizitia de date presupune citirea din memorie a matricei RGB, convertirea acesteia in grayscale pentru o procesare mai rapida de catre algoritmul ORB si de obtinerea hartii de adancime pentru cadrul respectiv. Acest lucru poate fi realizat in 2 feluri: matricea de distante este citita din setul de date TUM RGBD si a fost inregistrata cu o camera RGBD tip Microsoft Kinetic, sau se foloseste reseaua neurala FastDepth care estimeaza in timp real distanta pentru fiecare pixel din cadrul curent. Imaginea RGB si harta de adancime vor fi transmise ca parametrii clasei Tracker. TumDatasetReader stocheaza estimarile pozitiilor camerei pentru fiecare cadru in parte. Cadrele cheie, cele salvate in clasa Map, vor avea matricea de pozitie stocata nealterat in memorie, ele sunt deja relative fata de primul cadru citit. Pentru celelalte cadre, matricea de pozitie salvata in clasa TumDatasetReader este relativa la un cadru cheie, de preferat ultimul cadru cheie creat pana la citirea imaginii curente. Motivul pentru care se realizeaza salvarea pozitiilor in acest fel, este ca doar cadrele din clasa Map sunt salvate in memorie si pot fi optimizate de catre algoritmul Bundle Adjustment asa ca doar acestea ar trebui sa aiba valoarea lor salvata explicit.

Clasa MapPoint este fundamentala pentru buna functionare a algoritmului ORB-SLAM2. Aceasta este formata cu ajutorul unui KeyPoint si al unui KeyFrame asociat acestuia. In etapa anterioara, am prezentat modul in care se face proiectia coordonatelor unui punct cheie in spatiu, acestea devenind coordonatele globale ale MapPoint-ului pe care il creem. Punctului din spatiu i se asociaza de asemenea descriptorul acelui keypoint care l-a creat, pentru compararea ulterioara cu alte KeyPoint-uri din alte imagini. Un MapPoint are nevoie de un vector de orientare, acesta ajuta in verificarea proprietatii unui MapPoint de a fi sau nu vizibil dintr-un KeyFrame. Pentru a calcula acest vector de orientare prima data se determina coordonatele globale ale centrului camerei pentru cadru cheie care a creat acel KeyPoint, acest lucru se realizeaza in felul urmator:

$$\begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix} = -R_{wc}^t * t_{wc}, \quad T_{wc} = \begin{bmatrix} R_{wc} & t_{wc} \\ 0 & 1 \end{bmatrix} \quad (13)$$

Normalizarea diferentei intre coordonatele globale ale centrului camerei si ale MapPoint-ului creeaza vectorul de orientare. Acesta poate fi modificat, daca se constata ca mai multe cadre observa acelasi punct. In situatia respectiva, vectorul de orientare final va fi media aritmetica a celorlalti vectori de orientare individuali.

Clasa Feature, aceasta componenta nu exista in implementarea oficiala a ORB-SLAM, dar am considerat ca utilizarea acesteia ar simplifica codul. Extinde clasa KeyPoint fara a o mosteni explicit, are asociata distanta, extrasa din matricea de adancime, descriptorul si o valoare de tip boolean care arata daca punctul este monocular sau stereo. Aceasta clasificare se obtine prin compararea adancimii cu o valoare foarte apropiata de 0. In cazul implementarii mele, daca distanta este mai mica decat  $1e-2$ , consider ca valoarea estimata de camera RGBD ori de retea neurala nu este corecta si ca punctul respectiv este monocular, altfel consider ca este stereo. Pentru fiecare KeyPoint extras, se va crea o instanta a clasei Feature care va fi stocata direct in KeyFrame. Fiecare Feature are setat pe null la initializare o referinta la un obiect de tip MapPoint. Pentru a garanta functionarea in real time a algoritmilor, asocierea (Feature, MapPoint) trebuie sa poata fi accesata in  $O(1)$ . Vectorul de elemente Feature, impreuna cu o structura de tip dictionar unde cheia va fi MapPoint si valoarea de tip Feature, vor face acest lucru posibil. Singura problema este ca cele 2 structuri incearca sa reprezinte



aceleasi corelatii, in cazul vectorului am indicele unui Feature drept cheie si incerc sa accesez MapPoint-ul asociat, iar in cazul dictionarului, am referinta unui MapPoint si incerc sa obtin adresa unui Feature. Ambele structuri trebuie sa contina aceleasi perechi, altfel comportamentul algoritmului devine nedefinit.

Clasa KeyFrame, contine mai multe elemente legate de cadrul curent. Pentru a mentine functionarea sistemului in timp real, trebuie sa stocam in memorie rezultatele calculelor noastre. In aceasta clasa se vor regasi matricea de adancime, vectorul de instante ale clasei Feature, vectorul de trasaturi calculat de metoda bag-of-words implementata in DBOW2 si cadrul initial, convertit in format grayscale ce va fi folosit ulterior pentru afisarea in timp real a performantelor algoritmului. In interiorul constructorului acestei clase sunt mai multe operatii realizate, majoritatea necesare pentru a creste eficienta accesarii datelor. De exemplu: vectorul de tip Feature in medie contine 1000 de elemente care nu sunt sortate. In situatia in care proiectam un MapPoint in plan ar trebui sa comparam coordonatele proiectiei cu pozitia fiecarui Feature in parte pentru a stabili care este cel mai apropiat. O modalitate de a rezolva acest lucru este segmentarea suprafetei in  $K$  zone, in cazul meu am ales  $K = 100$ , fiecare reprezentand o portiune din imaginea initiala, avand asociate referintele Feature-urilor care se gasesc pe suprafata respectiva. In acest fel, in functie de zona in care este proiectat un MapPoint, vom stii ce Feature are o posibilitate mare de a corespunde, reducand astfel numarul de comparatii. Considerand ca toate valorile de tip Feature sunt dispuse in mod egal pe suprafata imaginii atunci complexitatea devine,  $O(N/K)$  unde  $N$  reprezinta numarul de Feature-uri iar  $K$  reprezinta numarul de zone in care a fost impartita imaginea. Constructorul este responsabil de initializarea structurilor de tip Feature, partitionarea lor in functie de coordonatele in imagine, si de memorarea estimarii curente a pozitie camerei si a centrului camerei in coordonate globale. Tot in aceasta clasa se regaseste structura de tip dictionar (MapPoint, Feature), care va fi adaptata pe tot parcursul algoritmului. Alta metoda importanta este: *get\_vector\_keypoints\_after\_reprojection*. Aceasta primeste ca date de intrare coordonatele proiectiei unui MapPoint, valoarea ferestrei de proiectie, si octava minima si maxima. Octavele reprezinta nivelul la care a fost observat un Keypoint in imagine si o estimare grosiera a distantei dintre camera si punctul din spatiu observat. Acesta poate sa aiba valori intre 0 si 7 inclusiv si ne spune de cate ori s-a facut resize la imagine pentru a surprinde o anumita trasatura. De exemplu: daca un Keypoint are valoarea octavei 0, inseamna ca algoritmul a

detectat-o in imaginea nemodificata. Daca ar fi 1, atunci dimensiunea imaginii a fost reduasa o singura data cu 0.8 din valoarea initiala si asa mai departe. Feature-ul care are asociat un MapPoint trebuie sa aiba valori ale octavei apropiate intre ele. Daca aceasta situatie nu s-ar respecta, ar introduce erori de estimare a distantei, ne asteptam ca Feature-uri corespondente, sa fie approximate la aceeasi distanta. Altfel ar putea inseamna ca cele 2 puncte din spatiu sunt diferite. Daca mediul are o structura simetrica, de exemplu: o sala de clasa cu bancile aliniate una in fata celeilalte, algoritmul ar putea observa 2 colturi ce apartin de 2 mese diferite, daca nu ar avea aceasta separare pe baza octavei, urmatorul cadru care observa aceleasi mese ar putea sa asocieze eronat punctele intre ele, afectand estimarea pozitiei. Fereastra de proiectie reprezinta cat de departe poate sa fie Feature-ul de coordonatele punctului de proiectie ale unui MapPoint pentru a fi considerata corecta o asociere intre cele 2 elemente. In functie de dimensiunea ferestrei aceasta poate intersecta 1, 2 sau 4 subsectiuni din cele 100 in care este impartita imaginea. Problema cea mai mare pe care am avut-o cu clasele KeyFrame si MapPoint era dependenta circulara. MapPoint-urile aveau nevoie de un KeyFrame si un Feature pentru a fi create si trebuia sa mentina o lista a KeyFrame-urilor care observa MapPoint-ul respectiv. In cazul KeyFrame-ului, acesta trebuie sa pastreze referinte asupra tuturor MapPoint-urilor pe care le observa. Pentru a rezolva aceasta problema am folosit o clasa aditionala care face operatii cu cele 2 structuri si am folosit forward declaration.

Clasa Map implementeaza harta pe care o foloseste algoritmul ORB-SLAM2. Aceasta este responsabila de stocarea corecta a KeyFrame-urilor, a MapPoint-urilor si rezolva problema dependentei circulare a celor 2 clase. Aici am implementate metodele de adaugare/stergere a unui MapPoint dintr-un KeyFrame. De asemenea, clasa MapPoint contine referinte la toate KeyFrame-urile care o observa. Aceste referinte sunt adaugate / sterse de catre 2 metode care se regasesc aici. Clasa Map creaza o structura de tip graf ponderat neorientat, in care nodurile sunt reprezentate de KeyFrame-uri. Arcele arata daca exista mai mult de 15 puncte comune intre 2 KeyFrame-uri iar ponderea lor este determinata de numarul de MapPoint-uri comune. Clasa Map realizeaza operatii pe grafurile de KeyFrame-uri, adauga/sterge noduri si face interogari pentru a afla vecinii directi sau cei pe nivel 2. Am ales sa implementez aceasta structura folosind *std::unordered\_map*. Drept cheie va avea KeyFrame-ul curent iar valoarea returnata de structura de tip dictionar va fi un alt *std::unordered\_map*, ce va contine toate celelalte KeyFrame-uri cu care este direct conectata dar si ponderea conexiunii.

În acest fel accesarea vecinilor de ordinul 1 va fi o operație ce se poate realiza în timp constant. Funcția *track\_local\_map* este folosită de către clasa *Tracking*. Aceasta primește ca date de intrare cadrul curent și ultimul cadru cheie salvat. Nu returnează nimic, doar încearcă să găsească câte un *MapPoint* pentru Feature-urile care încă nu au fost corelate cu un punct din spațiu. Aceasta operație este costisitoare și funcționează în felul următor:

1. sunt cautate toate KeyFrame-urile vecine de gradul 1 și 2 cu ultimul KeyFrame adăugat
2. din aceste KeyFrame-uri sunt extrase toate MapPoint-urile observate de către ele
3. MapPoint-urile sunt proiectate și sunt cautate potriviri pentru Feature-urile care încă nu au MapPoint-uri asociate.

Pentru a nu fi necesar să calculăm de fiecare dată KeyFrame-urile vecine și harta locală de MapPoint-uri, le stochez ca variabile în interiorul clasei *Map*. Acestea vor fi modificate în momentul în care un KeyFrame este adăugat în harta. Într-un caz ideal, ar trebui ca pentru fiecare Feature adăugat să se găsească un MapPoint, dar acest lucru rareori se întâmplă. În situația în care s-au găsit mai puțin de 30 de puncte din spațiu care s-au proiectat corect în imagine, se consideră că a apărut o eroare de urmărire și algoritmul începe o etapă de relocalizare.

Clasa *OrbMatcher* este responsabilă de realizarea urmăririi feature-urilor asemănătoare între cadre consecutive. Înainte de a începe prezentarea metodelor implementate, voi descrie pipeline-ul de procesare al unui punct din spațiu pentru a fi considerat observabil de către camera. Avem o instanță a obiectului *MapPoint* *mp*, dacă una dintre operațiile prezentate eșuează, punctul respectiv este ignorat de către KeyFrame-ul curent.

1. Se proiectează coordonatele globale ale *mp* în planul imaginii folosind matricea de estimare a poziției  $T_{cw}$  și matricea parametrilor camerei  $K$ . Se verifică dacă coordonatele proiectiei sunt valide pentru imagine.
2. Se calculează distanța  $d$  de la centrul camerei la *mp*. În funcție de valoarea octavei stocată în acest MapPoint, se pot estima o limită minimă și maximă pentru  $d$ . Dacă valoarea obținută nu se încadrează în acest interval se consideră că punctul este invalid.
3. Cu ajutorul geometriei analitice se obține ecuația dreptei care unește *mp* și centrul camerei. Aceasta dreaptă și vectorul de direcție al MapPoint-ului, trebuie să creeze un unghi cu o valoare mai mare de 60 de grade pentru a fi considerat *mp* observabil.

Daca aceste 3 verificari au fost realizate cu succes se considera ca punctul poate fi observat de catre camera. Exista 2 functii responsabile de asocierile intre cadrele curente, scopul acestora este ca gaseasca corespondente intre Feature-urile din cadrul curent si MapPoint-urile din spatiu. Pentru a se gasi perechea Feature  $f$  si MapPoint  $mp$ , trebuie ca  $mp$  sa se proiecteze in vecinatatea  $f$  iar descriptorii asociati atat Feature-ului cat si al MapPoint-ului sa aiba distanta Hamming sub un prag, setat in aceasta implementare la 50. O metoda ar fi compararea tuturor Feature-urilor din spatiu, cu totalitatea MapPoint-urilor observate de cadrul curent. Dar aceasta metoda ar fi ineficienta. O alta abordare ar fi separarea Feature-urilor in clustere in functie de distanta Hamming a descriptorilor, abordare stabila dar lenta si preferabil de utilizat cand nu ne putem baza pe estimarea matricei de pozitie a cadrului anterior. Iar cealalta abordare o reprezinta clusterizarea in functie de coordonatele in imagine ale Feature-ului.

Functie *match\_\_frame\_\_reference\_\_frame* implementeaza prima metoda. Aceasta primeste ca parametru 2 vectori de feature-uri calculati de biblioteca DBoW2, unul asociat cadrului curent, pentru care estimam matricea de pozitie si unul asociat cadrului anterior, pentru care cunoastem deja matricea de pozitie si asocierile de tip (Feature, MapPoint). Elementele acestor vectori sunt de tip (*int, vector\_\_descriptori*). Daca 2 astfel de perechi au prima valoare egala intre ele, inseamna ca cei 2 vectori de descriptori fac parte din acelasi cluster, conform arborelui din biblioteca DBOW2. Fiecare descriptor are asociata o instanta a clasei Feature. In cadrul anterior, instanta poate avea sau nu, un MapPoint corespondent. Daca exista acel MapPoint se poate proiecta in imagine. Descriptorul intern al MapPoint-ului este comparat cu ceilalti descriptori din acelasi cluster din cadrul curent si se aplica testul de proportionalitate Lowe pentru a garanta ca descriptorul cu distanta Hamming minima este cel mai bun. Functia *match\_\_consecutive\_\_frames* este mai simpla si implementeaza a doua metoda. MapPoint-ul din spatiu este proiectat in imagine si toate Feature-urile dintr-o zona circulara de raza de variabila sunt considerati posibili candidati pentru a crea o asociere (Feature, MapPoint). Se calculeaza distanta Hamming intre descriptorul MapPoint-ului si cel al Feature-ului. Iar descriptorul cu distanta minima si mai mica decat un prag setat la 100 este considerat ca fiind cel mai potrivit. Feature-ul asociat acelui descriptor, va pastra o referinta a MapPoint-ului.

Clasa MotionOnlyBA implementeaza in Ceres algoritmul Motion Only Bundle Adjustment, primeste ca date de intrare KeyFrame-ul curent si returneaza matricea de pozitie optimizata. Biblioteca lucreaza cu o notiune din C++ numita functori. Acestea sunt clase/structuri pentru

care s-a facut overload la operatorul  $()$ . Clasa `BundleError` se afla din aceeaasi categorie si implementeaza functia de eroare obtinuta din proiectarea unui `MapPoint` si asocierea acestuia cu un `Feature`. Pentru a crea problema de optimizare, clasa `ceres::Problem` trebuie sa stie care parametrii trebuie optimizati si functia de eroare pe care trebuie sa o minimizeze. In cazul acestui algoritm, singurul lucru care va fi modificat este matricea de pozitie a `KeyFrame`-ului pe care o voi converti in forma  $se(3)$ , transformand-o intr-un vector de 7 elemente. Iar pentru functia de eroare, nu voi scrie explicit ca este suma erorilor de proiectie. In schimb, voi initializa pentru fiecare asociere de tip  $(Feature, MapPoint)$  cate un element al clasei `BundleError`. Algoritmul de optimizare implementat de biblioteca Ceres, va incerca in mod independent sa reduca valoarea erorii pentru fiecare pereche in parte, modificand pe rand vectorul pozitiei. Exista un motiv pentru care schimb modul in care este exprimata pozitia camerei, matricea de pozitie contine 2 componente: matricea de rotatie  $R$  si un vector de translatie  $t$ . Pentru  $t$  nu exista restrictii de modificare atata timp cat aceasta nu aduce modificari mari intre pozitile a doua cadre consecutive, orice mod in care ar varia parametrii acestui vector, in continuare semnificatia lui de vector de translatie ramane nealterata, in aceasta situatie putem spune ca parametrii sunt alterati de catre biblioteca Ceres folosind *EuclidianManifold*, mici modificari bazate pe calcularea derivatelor partiale ale acestora din functia de eroare definita in `BundleError`, asemanator modului in care sunt modificati parametrii in retele neurale. Pentru matricea de rotatie  $R$  nu se mai poate aplica aceeaasi logica. Aceasta trebuie sa faca parte din structura de tip grup numit  $SO(3)$ , adica sa respecte egalitatea  $R * R^t = R^t * R = I$  si trebuie sa reprezinte o rotatie reala pe cele 3 axe. Alterarea aleatorie a parametrilor ar duce la o matrice invalida. Din aceasta cauza, modificarea rotatiei trebuie facuta cu un anumit unghi iar acest lucru se poate realiza printr-o inmultire de 2 matrici de rotatie valide. Din pacate nu exista implementare in forma matriceala pentru schimbarea unghiului de rotatie, dar este pentru Quaternioni. Din aceasta cauza fac conversia din matrice de pozitie in vector din categoria  $se(3)$ , iar pentru primii 4 parametrii asociati rotatiei, optimizarea lor se realizeaza folosind *QuaternionManifold*. Aceasta abordare rezolva problema instabilitatii numerice si garanteaza ca rezultatul operatiei de optimizare este un element valid in  $se(3)$ , ce poate fi ulterior convertit in forma matriceala. In functie de categoria din care face parte `Feature`-ul, acesta este considerat monocular sau stereo. Functia de eroare implementata de clasa `BundleError` este identica pentru ambele, cu exceptia ca pentru punctele stereo, este verificata si distanta la care se afla punctul fata de valoarea la care a fost estimata de camera RGBD.

Pentru a preveni instabilitatea cauzata de punctele de tip outlier, functia Huber descrisa in capitolul anterior este folosita in calcularea finala a erorii de proiectie. In implementarea oficiala realizata de g2o, agloritmul de optimizare este rulat de 4 ori, si dupa fiecare executie sunt eliminate punctele de tip outlier. Experimental, am observat ca etapa de optimizare cadru cu cadru este cea mai costisitoare operatie pe care o realizeaza algoritmul de ORB-SLAM2, executia acesteia de 4 ori, nu creste semnificativ acuratetea si reduce viteza de prelucrarea la aproximativ 5 cadre pe secunda, facandu-l nepotrivit pentru un sistem in timp real. Am observat ca obtin rezultate foarte bune, ruland o singura data Motion Only Bundle Adjustment, urmat apoi de o etapa de eliminare a corelatiilor (Feature, MapPoint) de tip outlier. Daca mai putin de 3 asocieri raman, se considera ca algoritmul a acumulat prea multe erori in urmarirea cadru cu cadru si trece intr-o stare de relocalizare.

Clasa Tracker realizeaza urmarirea traiectoriei cadru cu cadru. Aceasta integreaza fiecare dintre componentele definite anterior, si este responsabila de captarea cadrului curent, transformarea acestuia in KeyFrame si luarea deciziei daca va fi salvat in Map pentru a completa harta mediului inconjurator. Pasii urmasori se executa pentru fiecare cadru in parte:

1. Se creeaza KeyFrame-ul curent.
2. Se estimeaza matricea de pozitie pe baza legii de miscare.
3. Se realizeaza asocierea intre Feature-urile (puncte 2D) din cadru curent si MapPoint-urile observate de cadru anterior (puncte 3D)
4. Pe baza asocierilor respective realizate anterior, se optimizeaza matricea de pozitie a KeyFrame-ului curent, sunt eliminate asocierile de tip outlier
5. Este proiectata harta locala pe cadrul curent, si se gasesc noi asocieri (Feature, MapPoint), se executa din nou aceeasi operatie de optimizare Motion Only Bundle Adjustment
6. Este evaluat KeyFrame-ul curent, se verifica daca trebuie salvat in clasa Map.

Cadrul curent si matricea de adancime sunt citite de TumDatasetReader. In imaginea RGB se foloseste ORB pentru a extrage un vector de KeyPoint-uri si un vector de descriptori. Acestea sunt folosite pentru a initializa un obiect de tip KeyFrame. Pentru algoritmul ORB se foloseste o versiune modificata implementata in clasa ORBextractor si este conceputa sa extraga aproximativ 1000 de puncte cheie, acestea fiind distribuite cat mai egal pe suprafata imaginii. Daca un numar foarte mare de keypoint-uri s-ar obtine din aceeasi zona, acuratetea

estimarii ar avea de suferit, pixelii din zonele aflate mai aproape de camera se misca cu o viteză mai mare decât cei aflați în depărtare, dacă am considera doar punctele dintr-o anumită zonă în realizarea estimării, am obținut variații în mișcare prea bruste / lente depinzând de locul unde s-au găsit majoritatea punctelor. Parametrii setați pentru algoritmul ORB sunt următorii: 1000 de feature-uri, factorul de scalare al imaginii este 1.2, există maxim 8 nivele, și algoritmul FAST care face extragerea inițială de KeyPoint-uri să ia în considerare zona respectivă dacă diferența de intensitate între pixeli este de la 20 în sus. Dacă în schimb, zona este slab texturată atunci poate să seteze această diferență la 7, pentru a garanta că vor fi găsite feature-uri chiar și în cele mai dezavantajoase zone din imagine. De-a lungul duratei de viață a algoritmului, clasa Tracker păstrează 4 referințe de tip KeyFrame: cadrul curent care este analizat, 2 cadre imediat anterioare care vor fi folosite la estimarea poziției și ultimul cadrul referință care a fost creat. Cadrul referință este ultimul KeyFrame adăugat în Map și indică aproximativ în ce zonă se află camera și care MapPoint-uri ar trebui să fie vizibile. Cadrele mai vechi care nu au fost salvate în Map au fost șterse pentru a reduce cantitatea de memorie folosită. Pentru ultimul KeyFrame creat urmează etapa de estimare a poziției curente, aceasta se face pe baza legii de mișcare, iar valorile matricii vor fi calculate folosindu-ne de cele 2 cadre salvate în Tracker. Important aici de observat că pentru primul cadrul citit, poziția acestuia este matricea identitate  $4 \times 4$ , acest lucru sugerând că dispozitivul care înregistrează mediul consideră că primul KeyFrame este chiar originea sistemului de coordonate, iar toate matricile de poziție viitoare sunt de fapt transformări relative față de origine. Primul KeyFrame va fi salvat întotdeauna în clasa Map și este utilizat pentru a inițializa primele puncte de tip MapPoint: pentru toate Feature-urile de tip stereo din imagine, se vor crea puncte în spațiu. Din cauza acestui mod de inițializare, ORB-SLAM2 este sensibil până la apariția următorului cadrul cheie, estimările făcute de acesta în prima etapă fiind predispuse la erori. Uneori algoritmul își pierde orientarea cu totul, fiind necesară o etapă de relocalizare, sau de reluare a execuției acestuia. ORB-SLAM3, implementează o metodă mult mai robustă de inițializare, generând mai multe hărți locale în situația în care urmărirea cadru cu cadru esuează și le unește între ele în momentul în care recunoaște o zonă pe care a vizitat-o deja. După ce a fost creat KeyFrame-ul și a fost făcută estimarea inițială a poziției, clasa OrbMatcher este folosită pentru a găsi corelații între Feature-uri și MapPoint-uri. Alegerea metodei care îndeplinește acest lucru fiind determinată de numărul de KeyFrame-uri create de la ultima relocalizare sau de la adăugarea unui nou cadrul cheie în Map. Dacă nu se vor găsi minim

15 asocieri, se va considera ca algoritmul si-a pierdut orientarea, altfel, asocierile respective vor fi utilizate de catre MotionOnlyBA pentru a realiza optimizarea pozitiei. Perechile de tip outlier vor fi eliminate si noua pozitie a KeyFrame-ului va fi returnata. Daca vor ramane mai putin de 3 asocieri se va considera, din nou, ca algoritmul si-a pierdut orientarea. In final, se foloseste clasa Map pentru a proiecta toate punctele din harta locala pe cadrul curent iar asocierile gasite vor trece din nou printr-un proces de optimizare. Daca nu se gasesc minim 50 de perechi (Feature, MapPoint) inseamna ca urmarirea cadrului curent a esuat. Altfel se trece la etapa urmatoare si se va decide daca vom stoca in Map KeyFrame-ul curent. Acest lucru se va intampla daca urmatoarele conditii vor avea loc simultan.

1. au trecut mai mult de 30 de cadre de la ultimul KeyFrame adaugat in Map
2. numarul de MapPoint-uri in cadrul curent este 25% din numarul urmarit de cadrul de referinta
3. cadrul curent are cel putin 70 de Feature-uri de tip stereo, cu distanta dintre centrul camerei si punct este mai mica de 3.2 metri si urmareste cel putin 100 de MapPoint-uri

Clasa LocalMapping este responsabila de optimizarea hartii algoritmului. Aceasta sterge/adauga KeyFrame-uri si MapPoint-uri, iar la fiecare cadru cheie nou, realizeaza operatia de Local Bundle Adjustment. Aceasta metoda optimizeaza matricile de pozitie si toate MapPoint-urile vecinilor directi si cei de categoria a doua pentru KeyFrame-ul abia adaugat. In momentul in care thread-ul de Tracking considera ca un nou cadru cheie trebuie de adaugat in harta, se executa metoda principala *local\_map*, aceasta indeplineste urmatoarele operatii:

1. Creeaza noi MapPoint-uri din primele 100 de Feature-uri de tip stereo, sortate in ordine crescatoare dupa distanta la care se afla acestea de centrul camerei
2. Adauga cadrul curent in graful de KeyFrame-uri stabilind vecinii directi ai acestuia
3. Noile MapPoint-uri create sunt adaugate intr-o lista numita *recently\_added*, pentru a iesi din aceasta lista, punctele trebuie sa treaca un test care dovedeste ca nu sunt rezultatul unui Feature eronat detectat de catre algoritmul ORB, si ca pot fi folosite cu incredere
4. Se executa operatia de *culling*, punctele sunt verificate daca sunt valide iar daca nu, memoria lor este eliberata.
5. Se foloseste operatia de triangulare pentru a crea noi MapPoint-uri din Feature-urile care se potrivesc intre ele si fac parte din cadre cheie diferite.



6. Se detecteaza entitatile de tip MapPoint care reprezinta acelasi punct din spatiu, iar una dintre referinte este stearsa pentru creste coorenta hartii si a creste ponderea conexiunii dintre KeyFrame-urile adiacente
7. Se executa operatia de KeyFrame culling, se verifica daca informatiile pe care le detine un KeyFrame, adica totalitatea valorilor de tip MapPoint pe care le detine, sunt observate si din alte cadre. Daca peste 90% din punctele observate de un anumit cadru sunt vizibile si din alte cadre, KeyFrame-ul analizat este considerat redundant si memoria lui este eliberata. Acest lucru garanteaza ca graful clasei Map, contine doar cadre esentiale pentru reprezentarea norului de puncte.

In etapa a 4-a se executa operatia de *culling*, aceasta elimina punctele care nu sunt de incredere. Singurele puncte care nu vor trece prin aceasta etapa de verificare sunt cele generate de primul KeyFrame, tot primul KeyFrame nu poate fi sters deoarece ar da peste cap sistemul de coordonate local sub care lucreaza ORB-SLAM2. Un punct este considerat de incredere daca din momentul in care a fost creat, el a fost observat in 3 cadre cheie consecutive si daca a fost observat in cel putin 25% din numarul total de cadre care au trecut de la creerea acestuia. Ambele conditii trebuie sa fie respectate simultan in momentul in care se face verificarea punctului respectiv. Politica pe care o urmeaza familia de algoritmi ORB-SLAM este sa genereze multe puncte, fara a impune restrictii, pe care apoi le va supune acestui test de relevanta.

Ultima clasa este cea de MapDrawer pe care o folosesc pentru a afisa norul de MapPoint-uri, cadrul curent analizat si pozitiile cadrelor cheie observate. Folosesc biblioteca Pangolin si OpenGL pentru desenarea fiecărei structuri, camera urmareste cadrul curent. Interfata grafica scade viteza de procesare a cadrelor dar este o modalitate eficienta de a intelege vizual ce se petrece in algoritm. Implementarea pentru interfata grafica am realizat-o spre final, cand aveam celelalte componente finalizate, lucru care a ingreunat procesul de dezvoltare deoarece lucram cu valori numerice in terminal. Acum daca as reincepe implementarea, interfata grafica ar fi printre primele lucruri pe care le-as realiza. Datorita acestei clase am reusit sa gasesc erori in modul de constructie al grafului ponderat din clasa Map si al modului in care proiectam punctele in spatiu.

### 5.3 Pipeline antrenare FastDepth

Pentru rețeaua Neurala FastDepth pipeline-ul de antrenare a fost scris folosind biblioteca Pytorch iar pentru operațiile de preprocesare folosesc biblioteca Albumentations. Setul de date pe care am făcut antrenarea se numește Nyu Depthv2 Dataset și l-am obținut de pe Kaggle. Rezultatul acestui pipeline trebuie să fie un fișier de tip ONNX cu valorile parametrilor rețelei FastDepth în urma antrenării pe setul de date. O problemă pe care am observat-o la setul de date este că pentru imaginile de antrenament, adâncimile sunt exprimate ca fiind în intervalul  $[0, 255]$ , pe când în setul de date de validare, acestea se află între  $[0, 10000]$  reprezentând valorile în milimetri ale distanțelor. O limitare a acestui set de date este că nu poate detecta distanțe mai mari de 10 metri. Dar considerând că algoritmul trebuie să funcționeze pentru încăperi de mici dimensiuni, consider că această distanță maximă nu ar trebui să reprezinte o problemă. Pentru antrenare am ales să urmez lucrarea științifică și am setat hiperparametrii:

- Optimizatorul folosit a fost implementarea din Pytorch pentru Stochastic Gradient Descent, `torch.SGD`, având un learning rate de  $1e-3$ , o valoare a momentumului de  $\beta = 0.9$  și `weight_decay = 1e-4`.
- antrenarea s-a realizat pentru 50 de epoci iar durata antrenării a fost de aproximativ 6 ore jumătate. Laptopul pe care am antrenat este un Asus TUF Gaming A15, având un procesor AMD Ryzen 7 cu o frecvență de 4.2 GHz și placa video NVIDIA GeForce RTX 2060, cu o memorie de 6GB.
- Imaginile în setul de date au o dimensiune de  $(3, 460, 640)$ . Pentru a crește viteza de procesare am modificat dimensiunile la  $(3, 256, 320)$  și am aplicat o funcție de normalizare de tip `min_max`. Ambele transformări sunt aplicate atât pe setul de date de antrenare cât și pe cel de test.
- un batch de date are dimensiune de 8
- În lucrarea FastDepth funcția de pierdere folosită este `L1Loss`, aceasta fiind suma diferențelor dintre valoarea reală și cea determinată de rețeaua neurală în modul. În implementarea mea am ales să folosesc o funcție de pierdere mai robustă conform acestei lucrări științifice.

Acuratetea a fost verificată prin compararea diferenței relative între valorile obținute prin inferență și cele reale cu un factor `RELATIVE_ERROR = 0.15`. Această operație a fost

realizata pentru fiecare pixel in parte, iar acuratetea reprezinta procentul de pixeli cu o valoare care se incadreaza in limita impusa de `RELATIVE_ERROR`. Pentru a preveni antrenarea pentru intervale lungi fara a obtine rezultate, am avut 2 metode pe care le-am implementat: o strategie de early stopping: in situatia in care valoarea acuratetii nu ar fi crescut pentru 5 epoci antrenarea ar fi fost oprita si o strategie pentru modificarea learning rate-ului in timpul antrenarii. Daca acuratetea nu crestea pentru 3 epoci valoarea parametrului sa fie redusa la 0.3 din valoarea initiala. In practica am observat ca reseaua converge aproape monotonic catre o valoare optima. Functia de pierdere primeste ca date de intrare matricea de adancime obtinuta de catre reseaua neurala si matricea cu valori reale din setul de date, denumita groundtruth, si returneaza o valoare numerica de tip double care exprima cat de departe se afla estimarea noastra de realitate. Ideea antrenarii unei retele neurale este minimizarea acestei valori. Functia de eroare este alcatuita dintr-o combinatie liniara a 3 componente diferite: `L1Loss`, `GradientEdgeLoss` si `Structural Similarity Loss`, formula matematica este:

$$loss = 0.6 \cdot L1Loss + 0.2 \cdot GradientEdgeLoss + StructuralSimilarityLoss \quad (14)$$

`Structural Similarity Loss` se asigura ca media si distributia standard pe care o urmeaza valorile estimate, se apropie de media si distributia standard a matricei groundtruth. In comparatie cu celelalte 2 componente ale functiei de pierdere care sunt aplicate la nivel de pixel, aceasta abstractizeaza rezultatele ca fiind 2 distributii Normale cu parametrii  $\mathcal{N}(\mu, \sigma^2)$  care trebuie sa se suprapuna. Principiul de functionare pentru `GradientEdgeLoss` este ca pixeli din regiuni apropiate trebuie sa aiba cam aceleasi valori de estimare ale distantei si ca diferenta intre pixeli adiacenti pe axele x si y, ar trebui sa fie identica cu cea din imaginea groundtruth. Aceasta poate fi scrisa in felul urmatoar, unde  $N$  reprezinta numarul de pixeli din imagine, iar derivata valorilor pixelilor in raport cu axa de coordonate reprezinta diferenta intre matricea imaginii initiale si aceeasi matrice avand un rand shiftat la dreapta pentru axa X notata  $\frac{\partial I}{\partial x}$  si un rand shiftat vertical pentru axa Y notata  $\frac{\partial I}{\partial y}$ .

$$L_{edges} = \frac{1}{N} \sum_{i=1}^N \left( \left| \frac{\partial I_{pred}}{\partial x} - \frac{\partial I_{true}}{\partial x} \right| + \left| \frac{\partial I_{pred}}{\partial y} - \frac{\partial I_{true}}{\partial y} \right| \right) \quad (15)$$

## 6 EVALUARE

### 6.1 Setul de date TUM RGBD Dataset

Setul de date utilizat pentru a realiza evaluarea se numeste TUM RGBD Dataset. Acesta contine numeroase subseturi, fiecare verificand un aspect diferit al implementarii, ajutand la creerea unui imagini de ansamblu cu privire la robustetea algoritmului in functie de mediu in care se lucreaza si de traiectoria pe care o urmeaza. Cele 2 subseturi pe care le-am considerat potrivite pentru implementarea sunt:

- Subsetul `rgbd_dataset_freiburg1_xyz` contine cadrele unui video de 35 de secunde, in care traiectoria este in principal alcatuita din translatii, exista foarte putine rotatii fiind ideal pentru a verifica daca estimarea pozitiei in spatiu este corect realizata.
- Subsetul `rgbd_dataset_freiburg1_rpy` are 27 de secunde si contine foarte putine translatii. Exista in schimb numeroase schimbari bruste de rotatie care reduc acuratetea imaginii captate, testand la maxim capacitatea algoritmului ORB de a extrage feature-uri. Sistemul isi schimba orientarea pe toate cele 3 axe, fiind unul dintre cele mai dificile subseturi de date pe care se poate face antrenarea. Algoritmul ORB-SLAM2 este sensibil la operatiile de rotatie, mai ales atunci cand camera isi schimba orientarea catre o zona necunoscuta. Pentru a crea harta zonei respective sunt generate numeroase KeyFrame-uri si MapPoint-uri, pe care algoritmul trebuie sa le filtreze in clasa de LocalMapping, lucru care creste complexitatea temporala si spatiala si scade acuratetea sistemului.

Videourile sunt realizate cu ajutorul unei camere RGBD Microsoft Kinect, avand frecventa de 30 de cadre pe secunda, setul de date contine imaginile de tip RGB, hartile de adancime pentru fiecare cadru in parte, vectorii de pozitie in forma  $se(3)$ , primii 3 parametrii fiind pozitia in spatiu  $(tx, ty, tz)$  iar urmatorii 4 parametrii sunt asociati matricei de rotatie, scrisa sub forma de Quaternion,  $(qw, qx, qy, qz)$  si timestamp-urile asociate momentului in care au fost inregistrate fiecare din valorile din setul de date. Cu ajutorul acestor timestamp-uri putem

crea asocieri de tip (imagine RGB, matrice de adancime, pozitie) pe care le putem transmite algoritmul ORB-SLAM2. Pozitia este considerata ca fiind valoarea ideala, groundtruth, si va fi comparata cu rezultatele obtinute. Clasa TumDatasetReader este responsabila de citirea datelor si stocarea matricilor de pozitie obtinute pentru fiecare cadru. Dupa parcurgerea intregului set de date, valorile estimate sunt salvate intr-un fisier de tip text unde vor fi comparate cu cele reale.

## 6.2 Metrici utilizate

Algoritmul ORB-SLAM2 scrie intr-un fisier estimarile matricilor de pozitie pentru fiecare cadru in parte. Pentru a realiza comparatia cu datele de tip groundtruth din setul de date, folosesc un pachet din python numit *evo*, acesta este capabil sa creeze un grafic al traiectoriei, permitand astfel o reprezentare vizuala a rezultatelor si o separare a acestora in functie de ceea ce vreau sa evaluez: viteza, translatia sau orientarea. De exemplu, figura de mai jos reprezinta variatia translatie pe fiecare dintre cele 3 axe. Cu albastru este valoarea de tip groundtruth iar cu galben este estimarea realizata de implementarea mea pentru algoritmul ORB-SLAM. Rezultatele sunt obtinute pentru subsetul de date `rgbd_dataset_freiburg1_xyz`, acesta fiind special conceput pentru a testa corectitudinea estimarii translatiei intre cadre.

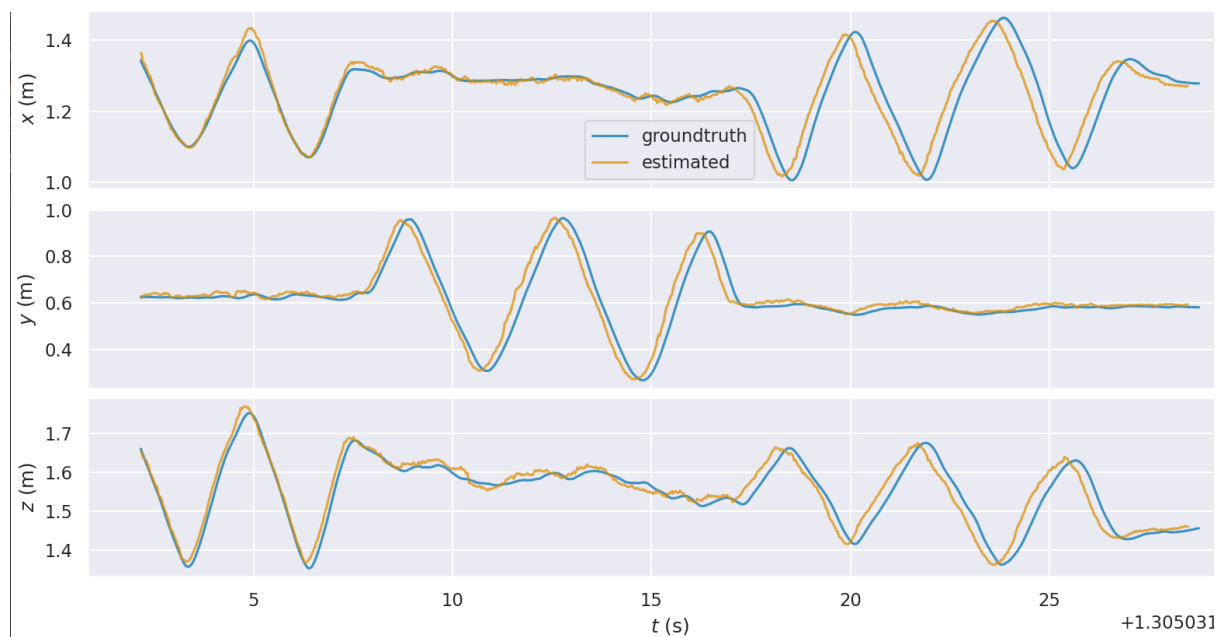


Figura 2: Graficul translatiei pe fiecare din cele 3 axe, groundtruth si estimare ORB-SLAM2

Consider ca o reprezentare grafica in care traiectoria groundtruth se suprapune exact cu ceea

ce a obținut estimarea ORB-SLAM2 poate fi considerată în mod neoficial, o metrică pe baza căreia să putem spune dacă algoritmul funcționează corect. Pentru exactitate se poate folosi: APE (Absolute Pose Error), această metrică măsoară distanța euclidiană dintre pozițiile estimate și cele reale, la fiecare moment de timp. În general valorile scorurilor APE obținute pentru ambele seturi de date sunt până în 0.05, sugerând că acuratețea este bună. Alte metrici pe care le folosesc pentru implementarea mea a algoritmului ORB-SLAM2 este numărul de secunde necesar pentru parcurgerea setului de date sau numărul de cadre pe secundă. Numărul de KeyFrame-uri create, în momentul în care sistemul nu mai poate realiza urmărirea cadru cu cadru, acesta inserează un nou KeyFrame, cu cât sunt mai puține KeyFrame-uri noi adăugate, se poate considera că traiectoria este ușor de interpretat și că sistemul este stabil. O altă metrică este legată de numărul de relocalizări pe care a trebuit să le facă algoritmul pentru a parcurge setul de date. Relocalizarea apare în situația în care urmărirea cadru cu cadru esuează și este căutat KeyFrame-ul care seamănă cel mai bine cu cadrul curent folosind vectorul de feature-uri calculat de metoda bag-of-words. Ideal, numărul necesar de relocalizări ar trebui să fie 0.

În cazul rularii implementării mele pe setul de date `rgbd_dataset_freiburg1_xyz` acesta durează în medie 67 de secunde, funcționând la aproximativ 15 cadre pe secundă, în medie este nevoie între 5-7 KeyFrame-uri noi pentru parcurgerea setului de date. Logica de relocalizare nu este deloc folosită, algoritmul fiind capabil să realizeze urmărirea cadru cu cadru. Pe setul de date `rgbd_dataset_freiburg1_rpy` a durat 73 de secunde, videoclipul având 27 de secunde, reprezintă aproximativ 10-11 cadre pe secundă. Acest lucru se datorează numeroaselor operații de optimizare a hărții pe care trebuie să le facă algoritmul deoarece sunt adăugate între 17-19 KeyFrame-uri pentru a parcurge întreg setul de date, din cauza mișcărilor bruste ale camerei care reduc considerabil claritatea imaginilor extrase. În continuare, numărul de relocalizări este 0. Mai jos, am atașat graficul care compară estimarea orientării pentru fiecare cadru, estimările fiind descompuse după cele 3 dimensiuni ale rotației. Graficul realizat de algoritmul ORB-SLAM2 pare să fie shiftat în timp față de cel real, dar să aibă aproximativ aceeași formă cu cel al valorilor de tip groundtruth. Consider că problema poate să pornească de la modul în care sunt atașate timestamp-urile pentru fiecare cadru în parte, lucru care nu are legătură directă cu modul în care este realizată implementarea, ci cu modul în care setul de date creează perechile (imagini RGB, vector poziție, matrice de adâncime). Am încercat utilizarea rețelei neurale FastDepth pentru a estima adâncimea în loc de a folosi matricea de

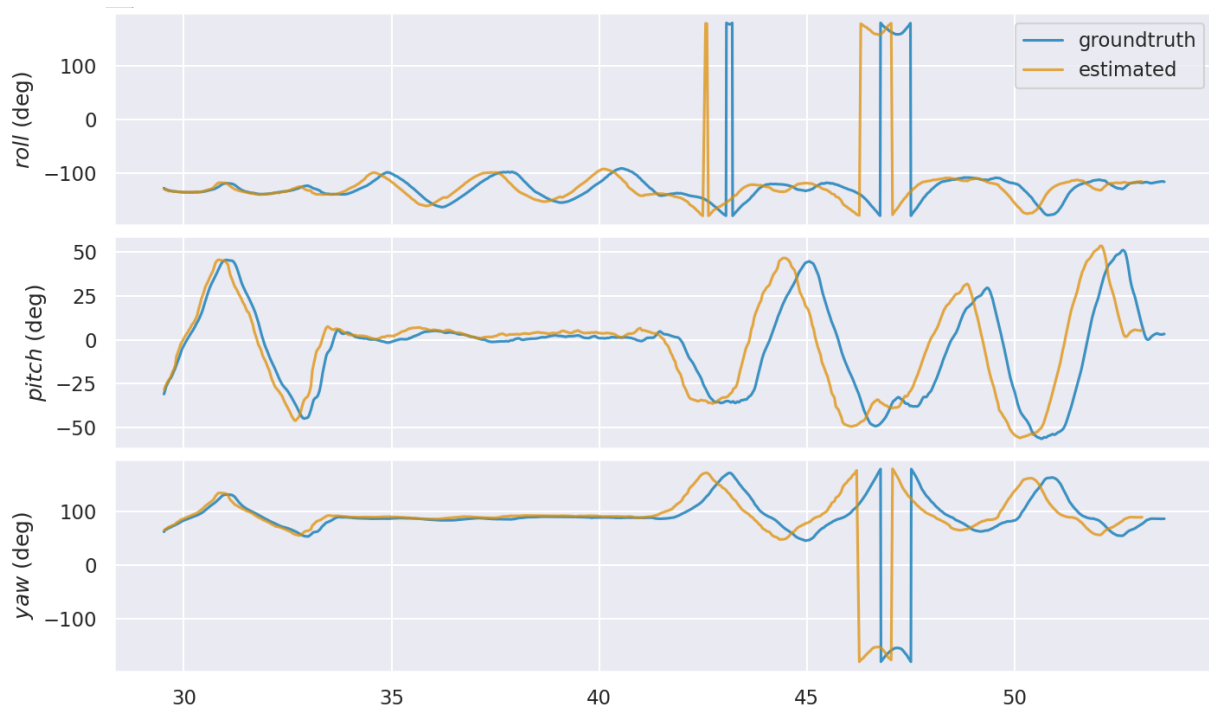


Figura 3: Graficul orientarii pentru setul de date rgbd\_dataset\_freiburg1\_rpy

distanțe a setului de date TUM RGBD. Pentru a testa arhitectura voi lua doar imaginea și vectorul de poziție simulând astfel o situație în care sistemul ar avea doar o cameră RGB. Problema cu această implementare este că rețeaua neurală nu este suficient de exactă iar estimările distanțelor între cadre consecutive în continuare variază mult. Pentru subsetul de date rgbd\_dataset\_freiburg1\_xyz implementarea intră în etapa de relocalizare în medie după primele 50-60 de cadre procesate. Sistemul este mult prea instabil pentru a realiza în mod corect urmărirea cadru cu cadru.

În etapele inițiale ale dezvoltării algoritmului am încercat utilizarea unui video realizat folosind camera telefonului pentru testarea implementării. Au fost 3 probleme pe care le-am întâlnit: nu puteam determina parametrii corecți ai camerei telefonului. Aveam nevoie de distanța focală, de coordonatele centrului imaginii și de parametrii de distorsiune. A doua problemă o reprezenta lipsa unei matrici de adâncime pentru fiecare cadru iar cea de-a treia, era lipsa unui vector de poziție pentru fiecare imagine. Chiar și în situația în care obțineam parametrii camerei telefonului folosind algoritmi implementați în OpenCV, în continuare nu puteam fi sigur dacă estimările realizate de mine sunt cele corecte. Din cauza acestor multe probleme, am ajuns la concluzia că un set de date ar fi o variantă mai potrivită.

## 7 CONCLUZII

O problema pe care am avut-o în testarea algoritmului a fost că nu am putut face funcțională implementarea inițială a ORB-SLAM2, exista conflicte între versiunile de biblioteci Eigen și g2o. Versiunea de Eigen folosită la momentul respectiv nu mai exista acum în repo-ul oficial și eu nu am reușit să o accesez pentru a testa.

Libraria Ceres este ușor de folosit pentru problemele de optimizare non-liniară, consider că cel mai mare plus pe care îl aduce lucrarea mea este că am cea mai completă implementare a algoritmului Bundle Adjustment în această bibliotecă la care se adaugă o logică de filtrare a punctelor de tip outlier și are caz separat de folosire atât pentru punctele monoculare cât și pentru cele stereo. În plus am adus optimizări la codul oficial pentru ORB-SLAM2: implementarea lor nu eliberează absolut deloc memoria pentru MapPoint-urile și KeyFrame-urile considerate invalide, eu am rezolvat această problemă, extinzând durata pentru care poate rula algoritmul și făcându-l potrivit pentru sistemele embedded. De asemenea, în implementarea oficială nu este folosită încă structura de tip dicționar, clasele principale folosite au parametri de stare care își modifică valoarea la fiecare cadru, funcțiile având efecte laterale care generează erori greu de urmărit și corectate. În total am scris 4030 de linii de cod în C++ și am modificat codul pentru numeroase funcții, făcându-l mai ușor de înțeles și menținut. ORB-SLAM2 îndeplinește 3 funcții: urmărirea cadru cu cadru, corectarea erorilor, și închiderea buclelor. Etapa de închidere a buclelor nu am reușit să o implementez din cauza apropierii termenului de predare, cu toate acestea, algoritmul obține în continuare rezultate bune pentru subseturile de date alese.

În ciuda faptului că utilizarea unei rețele neurale pentru a înlocui o camera RGBD nu a funcționat așa cum am crezut inițial, algoritmul devine instabil și își pierde complet orientarea după primele 50 de cadre, consider că o rețea neurală precum FastDepth poate fi folosită în estimarea adâncimii pentru punctele care au avut atribuită distanța 0 de către camera tip RGBD. O direcție viitoare ar fi utilizarea unui sistem care combină cele 2 abordări. Separarea straturilor convoluționale în depthwise și pointwise s-a dovedit a fi o tehnică bună pentru a crește viteza arhitecturii astfel încât să poată fi folosită în timp real.



Un lucru pe care îl regret este că nu am utilizat o interfață grafică încă de la primele etape, pentru a detecta erorile în timpul dezvoltării algoritmului. Deși scorul calculat de APE este o metrică bună pentru a vedea cât de bine se potrivesc 2 estimări ale poziției unui cadru, o simplă valoare numerică nu este suficientă pentru a avea o privire generală asupra modului în care funcționează implementarea.

Ca direcții viitoare, mă gândeam să utilizez arhitecturi de rețele neurale pentru sarcini bine delimitate, cum ar fi extragerea de feature-uri sau găsirea de corelații de tip (Feature, MapPoint), adăugarea unui modul de object detection și a unui algoritm de planificare de trasee, pentru a îndeplini sarcini simple de găsire a unor obiecte de mici dimensiuni. Până acum algoritmul a folosit doar metode clasice, ORB pentru extragere de feature-uri, Brute Force pentru matching, bag-of-words pentru relocalizare. Acum, consider că direcția pe care ar trebui să o urmeze această clasă de algoritmi de tip SLAM este una în care tehnici de Machine Learning sunt folosite pentru a crește viteza și poate acuratețea operațiilor. Această fiind și direcția încurajată de lucrările ce fac parte din state of the art până la momentul curent.

## **8 REFERINTE**