



# Tehnici avansate de percepție a mediului folosind Deep Learning și estimatori probabilistici (DEEPPSENSE)

## Etapa 1 - Analiza cerințelor, pregătirea datelor și a rețelelor neuronale

**Echipa de cercetare:** Radu Dănescu, Diana Borza, Razvan Itu, Andra Petrovai, Andrei Vadan

### Cuprins

<b>Introducere - rezumatul etapei</b>	<b>2</b>
<b>Colectarea imaginilor de trafic, ingineria și selecția datelor, adnotare</b>	<b>2</b>
a. Baze de date pentru detecția obstacolelor	3
b. Baze de date pentru segmentare semantică	3
c. Definirea unui format standard pentru bazele de date	4
<b>Proiectarea, antrenarea și evaluarea rețelelor neuronale, și reglarea fină a hiperparametrilor</b>	<b>5</b>
<b>Detecția și măsurarea mișcării folosind rețele neuronale convoluționale</b>	<b>9</b>
<b>Rezultate suplimentare: Calibrarea automată a parametrilor extrinseci ai camerei</b>	<b>10</b>
<b>Rezultate suplimentare: Dezvoltarea și implementarea unui algoritm pentru monitorizarea atenției șoferului pe baza clipirilor și a direcției de privire</b>	<b>13</b>
<b>Diseminarea rezultatelor parțiale</b>	<b>14</b>
<b>Concluzii</b>	<b>15</b>
<b>Bibliografie</b>	<b>15</b>



## 1. Introducere - rezumatul etapei

Obiectivul principal al acestei etape a fost proiectarea, antrenarea și evaluarea mai multor soluții de interpretare a imaginilor din trafic folosind rețele neuronale convoluționale (*convolutional neural networks*). Rețelele neuronale convoluționale au atins performanțe foarte bune, apropiate de cele ale oamenilor într-o mulțime de probleme de clasificare automată, dar necesită un set de antrenare foarte mare. De asemenea, antrenarea lor necesită reglajul fin (*fine tuning*) a mai multor hiper-parametri: rata de învățare (*learning rate*), funcțiile de activare utilizate, *weight decay*, utilizarea regularizării, utilizarea nivelurilor (layer) de *dropout* etc.

Având în vedere aceste aspecte, activitățile principale ale acestei etape au fost:

(1) centralizarea bazelor de date existente pentru interpretarea automată a scenelor de trafic pentru a avea o bază de date suficient de mare pentru antrenarea rețele neuronale convoluționale. Toate bazele de date au fost transformate într-un format unic, bazat pe fișiere cu metadate asociate fiecărei imagini. În plus, folosind dispozitive mobile, am capturat și adnotat (folosind) formatul definit anterior o nouă bază de date cu imagini din trafic.

(2) antrenarea și reglajul fin a mai multor rețele neuronale convoluționale care au obținut performanțe ridicate în problemele de clasificare de obiecte și segmentare semantică a scenei.

(3) integrarea rezultatelor CNN într-un framework probabilistic pentru urmărirea scenei, pentru estimarea parametrilor de poziție și viteză ai obiectelor relevante.

Un alt aspect pe care ne-am focalizat a fost monitorizarea atenției șoferului. Pentru aceasta am dezvoltat un framework generic de analiză a clipirilor și a direcției de privire. Pe baza acestora se poate determina dacă șoferul este atent la drum (direcția de privire) sau dacă manifestă oboseală (frecvența clipirilor). În plus, metoda propusă are aplicații și în alte domenii, cum ar fi detecția înșelăciunii (*deceit detection*) sau în psihologie.

Rezultatele obținute în această etapă au fost publicate în două articole la o conferință internațională, respectiv un jurnal.

## 2. Colectarea imaginilor de trafic, ingineria și selecția datelor, adnotare

Pentru realizarea obiectivelor acestui proiect, vom folosi date din două surse principale: baze de date internaționale disponibile pentru comunitatea de cercetare, care au avantajul că sunt adnotate cu obiecte detectate, distanțe, și alte rezultate de referință, și date achiziționate de noi, unde putem adăuga informații suplimentare, sau putem acoperi scenarii diferite.

Bazele de date existente au fost create de diferite instituții de cercetare și educație, de aceea ele vor avea imagini de dimensiuni diferite și datele de etichetă stocate în formate diferite.



În Tabelul 2.1. se prezintă o comparație din mai multe puncte de vedere (dimensiune, scenarii) a bazelor de date existente pentru analiza automată a imaginilor de trafic.

**Tabelul 2.1.** Prezentare comparativă a bazelor de date cu imagini de trafic

Baza de date	Număr secvențe	Număr imagini	Mai multe orașe	Condiții meteo diferite	Momente diferite ale zilei	Scenarii diferite
KITTI	22	14999	Nu	Nu	Nu	Da
Cityscapes	50	5000	Da	Nu	Nu	Nu
BDD100K	100000	120000000	Da	Da	Da	Da

#### a. Baze de date pentru detecția obstacolelor

##### KITTI

Baza de date KITTI [1] a fost creată de Karlsruhe Institute of Technology din Germania și oferă diferite seturi de date cu aplicații multiple: baza de date pentru obstacole (vehicule, autobuze, pietoni, etc.), segmentare, urmărire obstacole. Baza de date cu obstacole conține un total de 7481 imagini de antrenare cu vehicule și 7518 imagini de test, conținând un total de peste 80.000 de obiecte etichetate.

##### Udacity

Baza de date oferită de Udacity [2] pentru detecția de obstacole conține imagini din traficul rutier împreună cu fișiere de tip "csv" asociate fiecărei imagini. În fișierele text sunt stocate informațiile despre locația obstacolelor în spațiul imaginii și sunt definite coordonatele dreptunghiului care le încadrează: coordonatele punctului din colțul stânga sus: (xmin, ymin) și colțul din dreapta jos: (xmax, ymax) al dreptunghiului, numele imaginii, tipul clasei etichetate. Opțional în fișiere text mai este un câmp numit "occluded" care va avea valoarea 0 dacă obstacolul nu este obstrucționat și valoarea 1 în caz contrar.

#### b. Baze de date pentru segmentare semantică

##### CityScapes

Baza de date CityScapes [3] cuprinde secvențe din diferite orașe din Germania și Elveția (ex: Aachen, Koln, Dusseldorf, Zurich). Datele sunt din traficul rutier urban în 5000 de imagini și



cu 30 de clase de obiecte diferite. Imaginile sunt stocate pe 3 canale, de dimensiune 2048 x 1024 pixeli (un raport de 2:1). În imaginile adnotate culoarea clasei drumului este: (128, 64, 128). Baza de date CityScapes conține 2975 imagini adnotate cu carosabil.

#### KITTI

Baza de date KITTI [1] adnotată cu zone de drum conține doar 289 imagini din scene din trafic în orașul german Karlsruhe. Drumul este marcat folosind culoarea: (255, 0 255) în imaginile adnotate.

#### Berkeley Deep Drive

Baza de date Berkeley Deep Drive [4] pentru segmentare conține peste 100.000 imagini etichetate. Dimensiunea imaginilor este de 1280 x 720 pixeli pe 3 canale color (RGB). Imaginile de etichetă sunt de aceeași dimensiune, unde fiecare clasă are asignată câte o culoare diferită. Pentru drum, culoarea este: (128, 64, 128). Datele sunt stocate pe 8 biți, ceea ce înseamnă că intensitatea pe fiecare canal poate avea valori cuprinse între 0 și 255. Baza de date conține un total de 6575 imagini adnotate.

### c. Definirea unui format standard pentru bazele de date

#### Format standard pentru segmentare semantica:

Prin folosirea și definirea unui format standard pentru toate cele 3 baze de date pentru segmentare, am ajuns la un total de 9839 imagini, din care 7871 de antrenare și 1968 imagini de test (validare). Aceste imagini sunt transformate în matrice de valori reale (numere în virgulă mobilă) scalate în intervalul (0...1), un format mai potrivit pentru antrenarea și validarea rețelelor neuronale.

Deoarece bazele de date au un raport diferit între înălțimea și lățimea imaginilor sursă, am ales folosirea unei dimensiuni standard. Astfel toate imaginile au fost scalate la dimensiunea 256 x 256 pixeli.

#### Format standard pentru detecția obstacolelor:

Pentru a putea procesa toate bazele de date pentru obstacole, am ales definirea unei perechi imagine - fișier text. Toate bazele de date existente pentru detecție de obstacole au fost convertite în formatul comun definit în această secțiune. În cazul în care bazele de date nu dispun de toate informațiile senzoriale, au fost scrise valori de 0. De asemenea, în această etapă am implementat și un sistem propriu de achiziție, cu același format de date, folosind dispozitive mobile rulând sistemul de operare Android. Achiziția imaginilor este realizată prin camera foto a dispozitivului, iar restul de date senzoriale sunt preluate din senzorii disponibili pe dispozitivul mobil. Imaginile sunt salvate în format RGB pe 8 biți.



Fișierele text au următoarea structură:

- timestamp: UNIX timestamp - timpul achiziționării imaginii
- accX, accY, accZ: accelerația gravitațională în jurul axelor x, y și z
- gyroX, gyroY, gyroZ: rata de rotație a dispozitivului în jurul axelor x, y și z (radiani/sec)
- magX, magY, magZ: rotația din senzorul geo-magnetic în jurul axelor x, y și z
- yaw: unghiul de rotație în jurul axei x (rotație laterală a vehiculului de-a lungul axei drumului)
- pitch: unghiul de înclinație a dispozitivului mobil
- roll: unghiul de rotație în jurul axei y (rotația dispozitivului mobil pe axa verticală a drumului)
- yawRate: rata de rotație în jurul axei x (aceeași cu gyroX)
- speedMs: viteza de deplasare (metri / secundă)
- lat, long: coordonata de latitudine și longitudine obținută din senzorul GPS

Cele 3 unghiuri de rotație (yaw, pitch și roll) au fost obținute din valorile senzorului geo-magnetic și cel de accelerație gravitațională. Acest format de date poate fi accesat și procesat “offline” pe sisteme desktop.

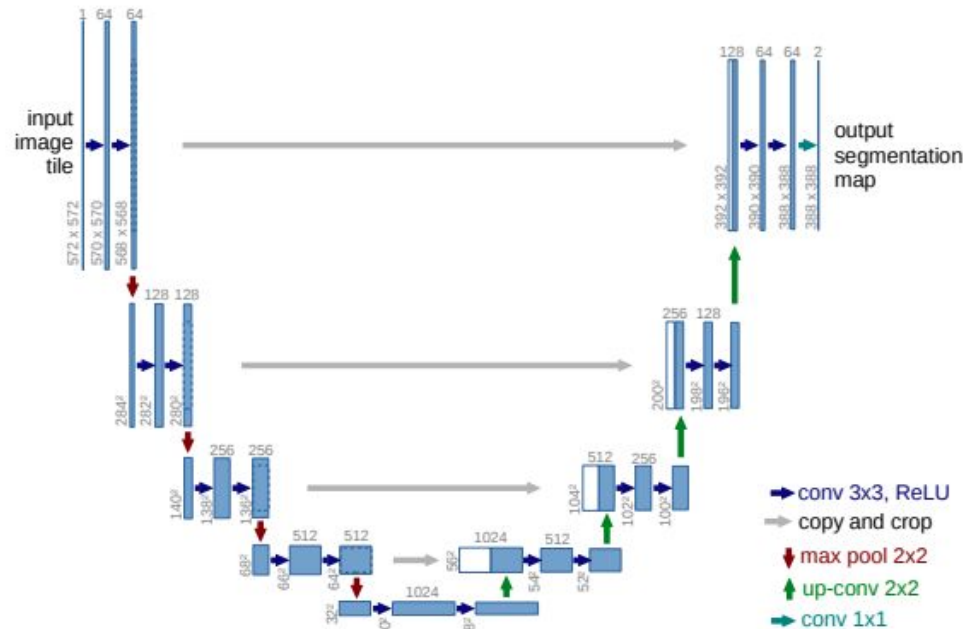
### **3. Proiectarea, antrenarea și evaluarea rețelelor neuronale, și reglarea fină a hiperparametrilor**

Inteligența artificială reprezintă un domeniu activ de interes și cercetare, mai ales în ultimii ani datorită popularității rețelelor neuronale convoluționale și apariției unor metode de a le antrena folosind hardware accesibil (plăci video grafice). În cadrul acestei etape am proiectat și antrenat o rețea neuronală convoluțională pentru segmentarea semantică a imaginilor din trafic. În această fază am ales segmentarea zonei de drum (“driveable road area”) din imagini. Deși în această etapă am procesat și adus la un format comun și bazele de date pentru detecție de obstacole, în prima fază am implementat doar partea de detecție de carosabil, urmând ca în etapele următoare să implementăm o altă rețea neuronală pentru a detecta obstacolele.

Pentru segmentare am ales o arhitectură de rețea de tipul “encoder-decoder”. Am ales folosirea unei rețele U-NET [5] modificată. Prima parte, cea de codificare, are rolul de a extrage informații relevante despre trăsături și caracteristici ale imaginii de intrare, în timp ce reduce dimensiunea spațială (datorită aplicării operațiilor de convoluție cu filtre având dimensiune de nucleu variabilă). Informațiile despre localizare sunt astfel pierdute, dar prin introducerea unor legături directe între nivelurile de codificare și cele de decodificare se va rezolva problema localizării trăsăturilor. Partea de decodificare, prin folosirea operațiilor de deconvoluție, are rolul



de a reconstrui harta de segmentare. În rețelele de tip U-NET, partea de decodificare are același număr de niveluri ca și cea de codificare (figura 1).



**Figura 1.** Arhitectura clasică a rețelei U-NET. Sursa: <https://arxiv.org/pdf/1505.04597.pdf>

Rețeaua de tip U-NET folosită în această etapă are ca intrare perechi de imagini de dimensiune 256 x 256 pixeli. Predicția (ieșirea rețelei), adică imaginea segmentată, va fi de aceeași dimensiune. Perechile de imagini de intrare sunt stocate pe disc astfel: imaginea originală (cu scena de trafic) este de tip întreg pe 8 biți pe 3 canale de culoare, în timp ce imaginea de etichetă (mască) este pe 1 canal pe 8 biți și are valoarea intensității 255 pentru carosabil, iar 0 pentru restul. Intensitățile din imaginile de etichetă sunt normalizate în intervalul [0, 1] pentru o antrenare mai eficientă. O pereche de imagini de intrare pentru rețea este ilustrată în figura 2.



**Figura 2.** Date de intrare pentru rețea: stânga imaginea originală, centru: imaginea de etichetă, dreapta: rezultatul aplicării etichetei ca mască pe imaginea originală



În procesul de antrenare am ales minimizarea funcției de cost “intersect over union”. Implementarea funcției de cost a fost făcută prin utilizarea coeficientul Sorensen-Dice (“dice loss”). Am utilizat o variantă modificată a coeficientului Sorensen, formula de calcul fiind:

$$IoU_{loss} = \frac{2(predicted_{output} \times label)}{predicted_{output} + label}$$

**Ecuatia 1.** Calculul funcției de cost “IoU” modificata.

Antrenarea unei rețele convoluționale este facilitată prin utilizarea mai multor date de antrenare. Pentru extinderea bazelor de date existente am ales să aplicăm procesul de augmentare a datelor. Am folosit următoarele tehnici de augmentare: variații aleatorii ale intensității imaginilor de intrare (în spațiul de culoare HSV), translații și rotații aleatorii ale imaginilor, scalarea lor și oglindirea pe orizontală (“horizontal flipping”). Pentru antrenarea și predicția folosind rețele neuronale am folosit limbajul de programare Python și software-ul TensorFlow [6] și Keras [7]. Pre-procesarea imaginilor și generarea bazelor de date modificate, iar apoi generarea imaginilor augmentare am realizat-o folosind software-ul OpenCV [8] pentru Python.

Rețeaua aleasă are următoarea structură: 5 straturi de codificare, 1 strat central și 5 pentru decodificare. Un strat de codificare are următoarele operații:

- convoluție
- “batch normalization”
- funcție de activare ReLU
- convoluție
- “batch normalization”
- activare ReLU
- “max pooling”

Stratul central al rețelei are aceleași operații ca și cel de codificare, cu excepția operației de “max pooling” de la final.

Stratul de decodificare are următoarele operații:

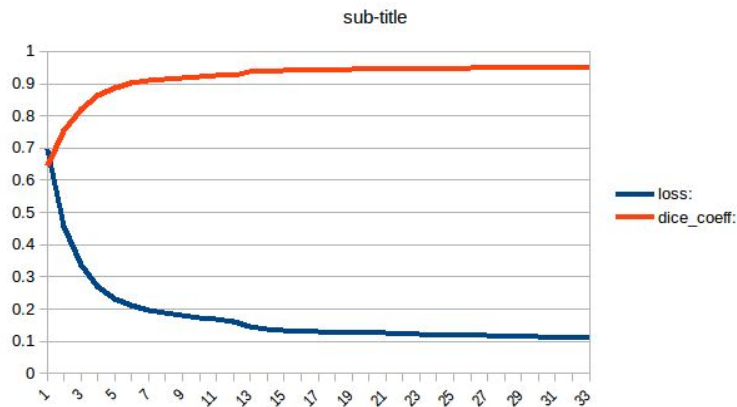
- deconvoluție (“up sampling”)
- concatenare (legatura directă cu stratul echivalent din codificare)
- deconvoluție
- “batch normalization”
- activare ReLU
- deconvoluție
- “batch normalization”



- activare ReLU
- deconvoluție
- “batch normalization”
- activare ReLU

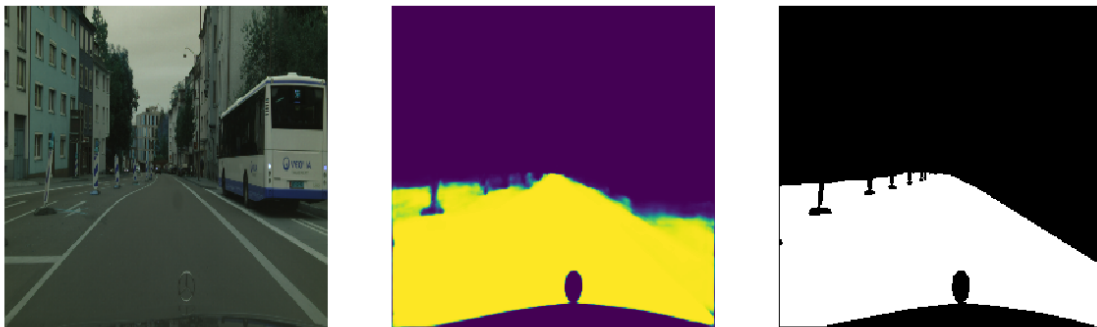
Predicția finală va fi rezultatul aplicării a unei convoluții cu nucleu 1 x 1 și a unei activări sigmoide.

Pentru antrenare am folosit o dimensiune a lotului de 16 și un număr maxim de 50 de epoci. Timpul de antrenare pentru o epocă este de 172 secunde (în medie) pe un sistem desktop echipat cu 2 plăci video Nvidia 1080 Ti. O analiză a antrenării este ilustrată în figura 3, unde se observă minimizarea funcției de cost (“loss”) și convergența coeficientului Dice spre valoarea 1 (în metrică IoU utilizată, valoarea 1 reprezintă o suprapunere perfectă a datelor).



**Figura 3.** Funcția de cost (“loss”) scade, iar coeficientul Sorensen-Dice (“dice coeff”) converge spre valoarea 1 după antrenarea timp de 33 epoci.

Un exemplu de rezultat al segmentării folosind rețeaua neuronală este ilustrat în figura 4:



**Figura 4.** Exemple de predicție a rețelei. Stânga: imaginea de intrare color, centru: rezultatul segmentării (predicția), dreapta: imaginea etichetată (“ground truth”).





## 4. Detecția și măsurarea mișcării folosind rețele neuronale convoluționale

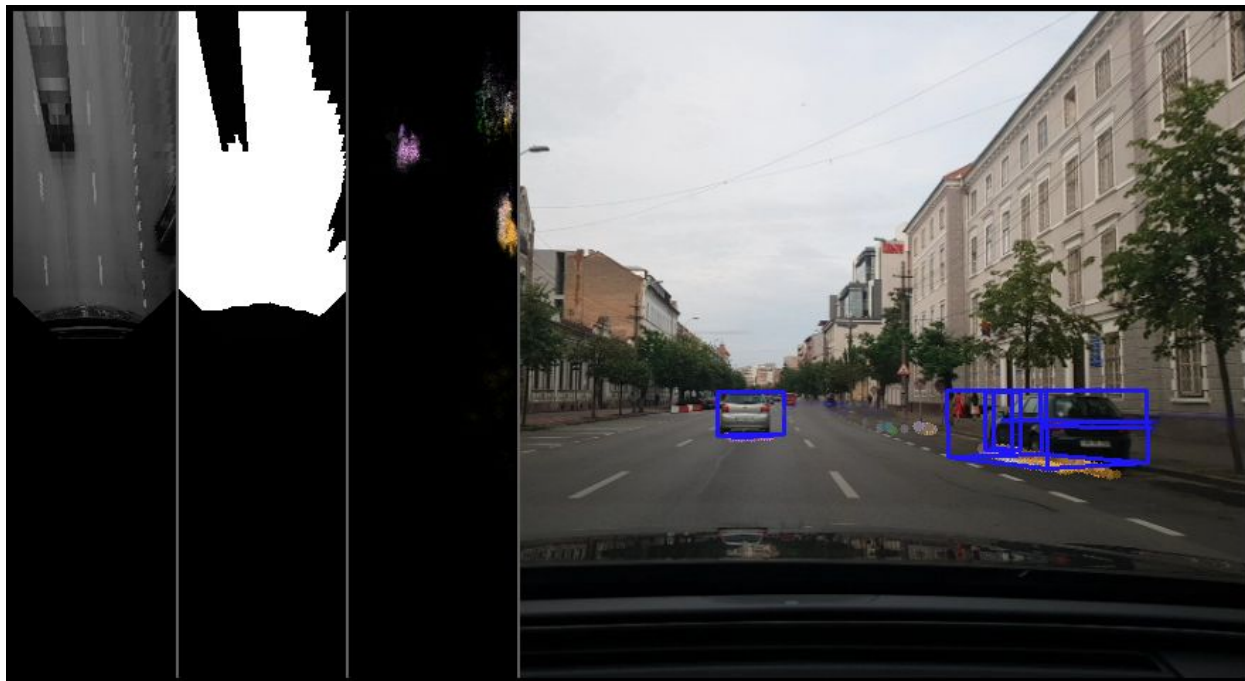
Imaginile segmentate folosind tehnicile descrise în secțiunea anterioară sunt integrate într-un framework probabilistic, bazat pe o hartă dinamică a scenei din jurul autovehiculului.

Din detecția zonei libere de drum, putem considera că tot ce nu este drum reprezintă un posibil obstacol și poate fi interpretat și urmărit în timp. Astfel, această soluție este utilă pentru a detecta orice tip de obstacol, nefiind limitată doar la vehicule sau pietoni. Pentru implementarea acestei funcționalități am ales folosirea unei hărți de ocupare bazată pe un filtru de particule. Aceste abordări pot fi folosite pentru estimarea sistemelor neliniare, multi-modale.

Harta de ocupare a scenei bazată pe filtre de particule este o reprezentare dinamică a scenei unde efectul de perspectivă este eliminat. În implementarea din cadrul acestui proiect, obstacolele sunt compuse din particule, fiecare având o poziție și un vector de viteză. Harta de ocupare este împărțită în celule de 20 x 20 cm, iar probabilitatea unei celule de a fi ocupată este dată de numărul de particule din acea celulă. Totodată, particulele pot avea un rol dublu: de a forma ipoteze pentru posibile obiecte, sau pot fi considerate elemente de bază a mediului cu proprietatea că particulele se pot deplasa dintr-o celulă în alta (astfel reușind o reprezentare dinamică a scenei). Partea de predicție a stării este realizată prin deplasarea particulelor.

Obstacolele din scenă sunt reprezentate de un set de particule, unde fiecare particulă are o poziție în harta de ocupare, o componentă de viteză dar și o variabilă pentru "vârsta" particulei din momentul creării ei. Acest parametru de "vârsta" este folosit în procesul de validare și estimare a vitezei pentru particulă, deoarece sunt luate în considerare numai particulele care "supraviețuiesc" pentru mai multe cadre. Probabilitatea unei celule de a fi ocupată este estimată ca fiind raportul dintre numărul de particule din celula (poziția lor coincide cu poziția celulei) și numărul total de particule permise într-o singură celulă " $N_c$ ". Acest parametru este constant și a fost setat la 200 de particule totale într-o celulă. Alegerea acestei valori reprezintă un compromis între acuratețe și performanța sistemului: un număr mare de particule într-o celulă înseamnă că se mențin mai multe ipoteze de viteză simultan, iar sistemul de urmărire va gestiona mai bine obstacolele care deplasează rapid în scenă, dar viteza de execuție a algoritmului va scădea. Viteza unei celule din harta de ocupare este calculată ca fiind viteza medie a particulelor din acea celulă.

Primul pas al algoritmului de urmărire este cel de predicție și este aplicat pe fiecare particulă. Poziția unei particule este modificată în funcție de viteza ei proprie și în funcție de parametri de mișcare ai ego-vehiculului (viteza și rata de rotație "yaw rate"). Al doilea pas al algoritmului este procesul de actualizare, bazat pe harta binară de ocupare a celulelor creată prin procesarea imaginii segmentate.



**Figura 5.** Framework-ul probabilistic de percepție a scenei.

Figura 5 ilustrează etapele de procesare din cadrul framework-ului probabilistic. În imaginea din stânga se pot observa etapele procesării și crearea hărții de ocupare dinamică: prima imagine reprezintă scena cu efectul de perspectivă eliminat (IPM), în a doua este ilustrată imaginea segmentată IPM (obținută din rețeaua neuronală convoluțională), a 3-a imagine ilustrează particulele și viteza lor codificată în funcție de vectorii de viteză (albastru - obiect în mișcare, galben - obiect static). În imaginea din dreapta este prezentat rezultatul procesării și grupării particulelor și a celulelor în obstacole.

## 5. Rezultate suplimentare: Calibrarea automată a parametrilor extrinseci ai camerei

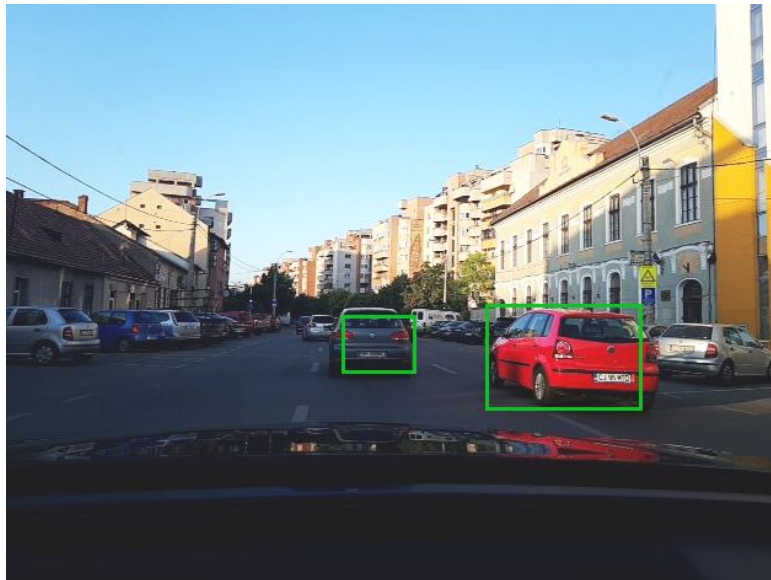
Unul dintre dezideratele acestui proiect este de a oferi un sistem de percepție a mediului care să poată funcționa cu resurse senzoriale limitate, inteligența artificială suplinind absența unor senzori costisitori dar foarte preciși, precum laserscannerele sau radarele. Șoferul va putea monta o cameră sau un dispozitiv mobil în spatele parbrizului, camera fiind orientată spre traficul din exterior, dar pentru a putea utiliza informațiile de la această cameră trebuie ca trăsăturile din imaginile achiziționate să poată fi puse în legătură cu trăsăturile 3D ale lumii reale. Acest lucru implică procesul de calibrare, ceea ce un utilizator obișnuit ar prefera să nu fie obligat să facă.

În mod tradițional, calibrarea se face în medii controlate, în condiții de laborator, de obicei prin măsurarea manuală a unor obiecte de referință. Dacă ne referim la scenariul în care



utilizatorul pur și simplu montează camera în mașină, dorind apoi să pornească la drum, condițiile de laborator nu pot fi îndeplinite, și astfel avem nevoie de calibrare automată.

Pentru calibrarea automată, avem nevoie de repere cu dimensiuni 3D cunoscute, precum lățimea unei benzi de circulație sau dimensiunea unor obstacole cunoscute. Noi am ales varianta a doua, și ne vom folosi de avantajele rețelelor neuronale pentru a selecta doar obstacolele de tip autovehicul din imagini necalibrate. Soluția pentru detecția obstacolelor este bazată pe o CNN pre-antrenată cu exemple formate din perechi imagini-obstacole reprezentate ca un dreptunghi în spațiul imagine. Am folosit API-ul TensorFlow pentru simplificarea procesului de antrenare, și am re-antrenat o soluție existentă de detecție a obstacolelor din imagini singulare numită MobileNet [10], care este o rețea proiectată în mod explicit pentru a avea un număr redus de parametri de antrenare, și care necesită putere de calcul redusă pentru predicție, putând astfel fi utilizată în timp real pe dispozitive mobile. Un exemplu de detecție automată a acestor obiecte este dat în figura 6.



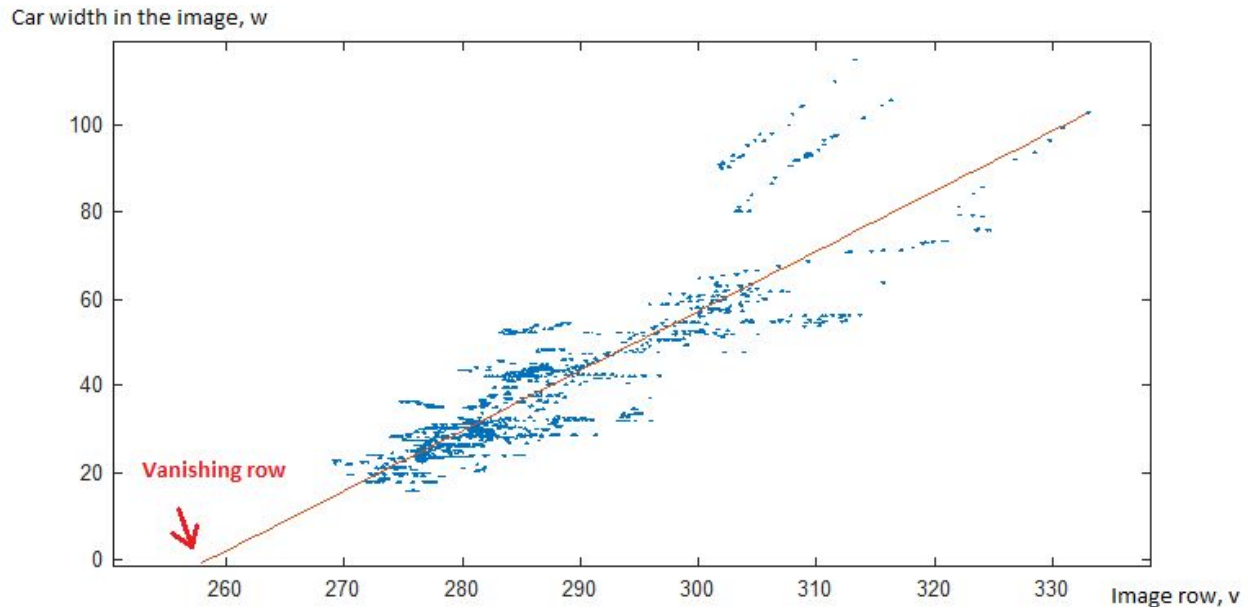
**Figura 6.** Exemplu de autovehicule detectate automat de MobileNet.

În funcție de distanța unui autovehicul față de noi, dimensiunea lui aparentă în imagine va fi variabilă. Distanța în lumea reală va avea o legătură directă cu coordonata rândului din imagine, iar obiectul va avea dimensiunea zero (va dispărea) când coordonata limitei lui inferioare va fi egală cu linia orizontului (linia de fugă, “vanishing line”). De fapt, între coordonata rândului limitei inferioare a obiectului și dimensiunea lui în imagine va exista o relație liniară, parametrii dreptei fiind direct influențați de parametrii extrinseci ai camerei video, unghiul de aplecare (pitch) și înălțimea camerei față de sol.

Datorită faptului că suprafața drumului nu e întotdeauna plană, dimensiunea vehiculelor nu este cunoscută cu exactitate, și detecția nu este perfectă (vehiculul nu este încadrat perfect), nu vom obține o dreaptă exactă pentru toate vehiculele din scenă. Colectând toate detecțiile



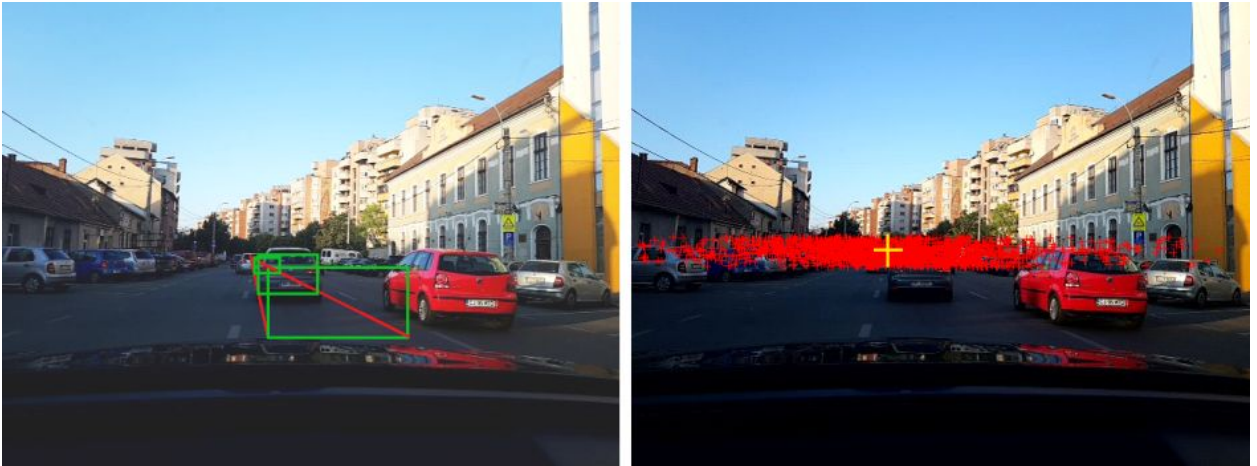
dintr-o secvență de 5 minute, se va obține un grafic de genul celui prezentat în figura 7. Pentru identificarea dreptei care se potrivește cel mai bine acestor date, se folosește metoda RANSAC.



**Figura 7.** Exemplu de autovehicule detectate automat de MobileNet.

Din parametrii dreptei obținute vom estima înălțimea și unghiul de aplecare al camerei folosind o metodă bazată pe filtrul Kalman extins, unde starea estimată este un vector format din înălțimea camerei și unghiul de pitch, iar măsurătoarea este dată de două puncte de tip (coordonată rând, lățime) de pe dreapta estimată. Filtrul converge cam după 4-5 iterații.

Pentru măsurarea unghiului de orientare laterală a camerei față de axul longitudinal al mașinii, se folosesc perechi de detecții obținute în cadre succesive, presupunând că ele aparțin aceluiași obiect. Traectoria obiectului trebuie să convergă spre punctul de fugă, care este poziționat pe linia orizontului, pe care am estimat-o deja. În realitate, fără un algoritm de urmărire complex, nu putem garanta că detecțiile aparțin aceluiași obiect, dar vom utiliza din nou o abordare statistică, și vom selecta mediana ipotezelor acestor puncte de fugă, precum în figura 8.



**Figura 8.** Principiul detecției punctului de fugă (stânga), și selectarea punctului de fugă din multiple ipoteze (dreapta).

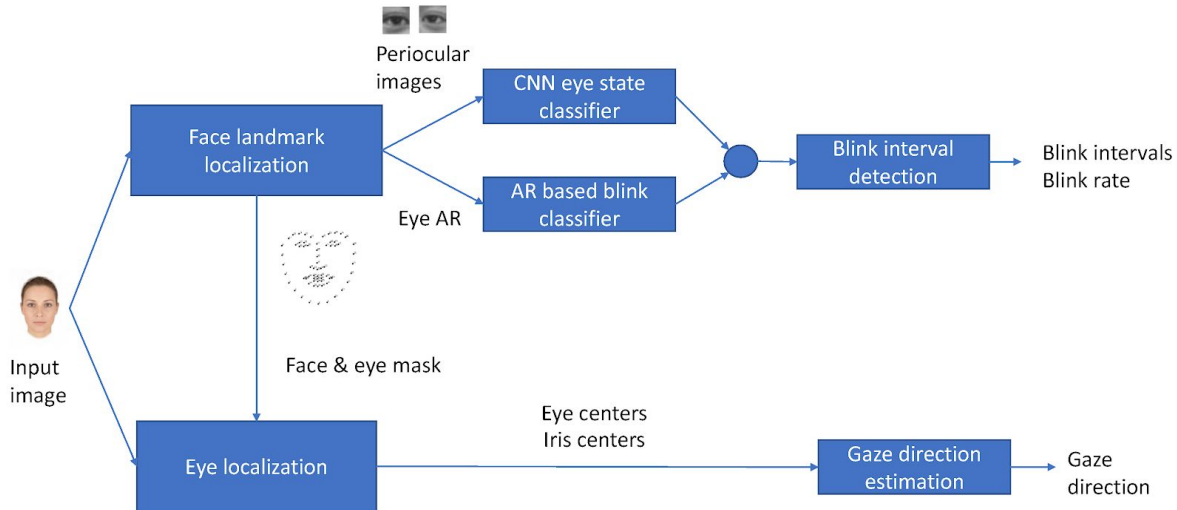
Coordonata coloanei punctului de fugă este direct legată de unghiul de orientare laterală a camerei. O poziție perfect centrală a acestui punct înseamnă o aliniere perfectă a camerei cu axul autovehiculului propriu, iar orice deviație de la acest punct este dată de unghiul de orientare  $\psi$ , astfel:

$$\tan \psi = (x_0 - w/2)/f,$$

unde  $x_0$  este coordonata coloană a punctului de fugă,  $w$  este lățimea imaginii, și  $f$  este distanța focală a camerei în pixeli.

## 6. Rezultate suplimentare: Dezvoltarea și implementarea unui algoritm pentru monitorizarea atenției șoferului pe baza clipirilor și a direcției de privire

În cadrul acestei etape am dezvoltat un sistem automat de analiză a mișcărilor oculare care calculează mai multe metrice oculomotoare, cu implicații importante într-o varietate de aplicații: monitorizarea atenției, estimarea nivelului de oboseală, interacțiune om calculator etc. Arhitectura sistemului propus este ilustrată în Figura 9: metoda propusă se bazează pe două module principale: modulul de analiză a clipirilor și modulul de analiză a direcției de privire.



**Figura 9.** Arhitectura soluției de analiză a clipirilor și a direcției de privire

La început, se utilizează o bibliotecă software (*dlib* [9]) pentru a detecta fețele umane și 68 de puncte de interes de pe față, inclusiv colțurile ochilor și mai multe puncte de pe pleoape. Pe baza acestor puncte se stabilesc regiunile de interes pentru metoda propusă.

Modulul de analiză a clipirilor combină răspunsul a doi clasificatori pentru a determina starea ochilor (închiși sau deschiși) în fiecare cadru al secvenței video. Primul clasificator se bazează pe punctele de interes detectate de librăria *dlib*, mai precis analizează raportul de aspect (*aspect ratio*) dintre lățimea și înălțimea conturului ochilor. Cel de-al doilea clasificator este o rețea neuronală convoluțională care analizează regiunea perioculară pentru a determina starea ochiului.

Apoi, utilizând o metodă rapidă de localizare a centrului irisurilor, modulul de detecție a direcției de privire corelează poziția irisurilor cu centrul și colțurile ochilor pentru a determina punctul în care privește subiectul.

Pe baza acestor trăsături extrase automat, se calculează mai multe metrici care pot fi utilizate pentru a monitoriza nivelul de atenție al șoferului: rata clipirilor, direcția de privire, recunoașterea mișcărilor oculare.

## 7. Diseminarea rezultatelor parțiale

Rezultatele obținute în această etapă au fost publicate în următoarele jurnale/conferințe internaționale:

- a. Diana Borza, Razvan Itu, Radu Danescu, "[In the Eye of the Deceiver: Analyzing Eye Movements as a Cue to Deception](#)", Journal of Imaging, Vol. 4, No. 10, 2018, Art. No. 120. [Indexat ISI]



- b. Razvan Itu, Radu Danescu, "[Machine Learning Based Automatic Extrinsic Calibration of an Onboard Monocular Camera for Driving Assistance Applications on Smart Mobile Devices](#)", 2018 International Forum on Advanced Microsystems for Automotive Applications, pp. 16-28.[**Capitol carte Springer**]

## 8. Concluzii

În cadrul acestei etape am reușit aducerea la un format comun a bazelor de date pentru detecția de obstacole și a bazelor de date pentru segmentarea semantică a scenelor din trafic. Am reușit, de asemenea și implementarea și antrenarea unei rețele neuronale convoluționale pentru a segmenta semantic zona de drum din imaginile din traficul rutier. Acesta abordare a fost implementat utilizând biblioteci software precum TensorFlow și Keras și folosind limbajul de programare Python. Rezultatele preliminare obținute prin analiza imaginilor folosind rețelele neuronale au fost integrate într-un framework probabilistic pentru urmărirea componentelor statice și dinamice ale scenei, obținând o estimare preliminară a poziției și mișcării obiectelor din scenă. În mod suplimentar, am realizat, folosind rezultatele CNN, un sistem de calibrare a parametrilor extinseci ai camerei, element valoros pentru realizarea unui sistem de percepție ce necesită o cât mai mică intervenție din partea utilizatorului. Totodată, în această etapă am obținut dezvoltat un framework generic pentru analiza mișcărilor oculare în vederea monitorizării atenției șoferului. Sistemul propus determină momentele în care subiectul a clipit, precum și rata de clipire. De asemenea, sistemul include o metodă simplă de estimare a direcției de privire pe baza deplasamentului relativ dintre centrele ochilor și centrele irișilor. Pe baza acestor metrici se poate propune un sistem care atenționează șoferul atunci când ele indică o creștere a oboselii.

Rezultatele preliminare au fost publicate în jurnale și conferințe internaționale.

## 9. Bibliografie

[1] - Geiger, A, Lenz, P, Urtasun, R, "Are we ready for autonomous driving?", Computer Vision and Pattern Recognition Conference (CVPR), pp. 3354–3361, 2012.

[2] - Udacity Vehicle Dataset, online:

<https://github.com/udacity/self-driving-car/tree/master/annotations>

[3] - M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The Cityscapes Dataset for Semantic Urban Scene Understanding," in Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016.



- [4] - Yu F., et. al, "BDD100K: A Diverse Driving Video Database with Scalable Annotation Tooling", preprint arxiv:1805.04687, 2018, online: <https://arxiv.org/abs/1805.04687>.
- [5] - O. Ronneberger, P. Fischer, T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation", Medical Image Computing and Computer-Assisted Intervention (MICCAI), Springer, LNCS, Vol. 9351: 234-241, 2015, available at arXiv:1505.04597.
- [6] - TensorFlow, online: <https://www.tensorflow.org/>
- [7] - Keras, online: <https://keras.io/>
- [8] - OpenCV, online: <https://www.opencv.org/>
- [9] - King, D.E., "Dlib-ml: A machine learning toolkit", J. Mach. Learn. Res. 2009, 10, 1755 - 1758.
- [10] - Howard, A, et al, "Mobilenets: Efficient convolutional neural networks for mobile vision applications", 2017, arXiv:1704.04861(preprint).