

Mobile Deep Visual Detection and Recognition

Ondertitel (facultatief)

Thijs VERCAMMEN

Promotor(en): Prof. dr. ir. Toon Goedéme

Co-promotor(en): Ing. Floris De Feyter

Masterproef ingediend tot het behalen van
de graad van master of Science in de
industriële wetenschappen: Elektronica-ICT
ICT

Academiejaar 2021 - 2022

©Copyright KU Leuven

Zonder voorafgaande schriftelijke toestemming van zowel de promotor(en) als de auteur(s) is overnemen, kopiëren, gebruiken of realiseren van deze uitgave of gedeelten ervan verboden. Voor aanvragen i.v.m. het overnemen en/of gebruik en/of realisatie van gedeelten uit deze publicatie, kan u zich richten tot KU Leuven Technologicampus De Nayer, Jan De Nayerlaan 5, B-2860 Sint-Katelijne-Waver, +32 15 31 69 44 of via e-mail iiw.denayer@kuleuven.be.

Voorafgaande schriftelijke toestemming van de promotor(en) is eveneens vereist voor het aanwenden van de in deze masterproef beschreven (originele) methoden, producten, schakelingen en programma's voor industrieel of commercieel nut en voor de inzending van deze publicatie ter deelname aan wetenschappelijke prijzen of wedstrijden.

Voorwoord

Het voorwoord vul je persoonlijk in met een appreciatie of dankbetuiging aan de mensen die je hebben bijgestaan tijdens het verwezenlijken van je masterproef en je hebben gesteund tijdens je studie.

Samenvatting

De (korte) samenvatting, toegankelijk voor een breed publiek, wordt in het Nederlands geschreven en bevat **maximum 3500 tekens**. Deze samenvatting moet ook verplicht opgeladen worden in KU Lokaal.

Abstract

Het extended abstract of de wetenschappelijke samenvatting wordt in het Engels geschreven en bevat **500 tot 1.500 woorden**. Dit abstract moet **niet** in KU Loket opgeladen worden (vanwege de beperkte beschikbare ruimte daar).

Keywords: Voeg een vijftal keywords in (bv: Latex-template, thesis, ...)

Inhoudsopgave

Voorwoord	iii
Samenvatting	iv
Abstract	v
Inhoud	vii
Figurenlijst	viii
Tabellenlijst	ix
Symbolenlijst	x
Lijst met afkortingen	xi
1 Situering en doelstelling	1
1.1 Situering	1
1.2 Probleemstelling	1
1.3 Doelstellingen	1
2 Herkenning en Detectie Algemeen	3
2.1 Deep learning-gebaseerde herkenningssystemen	3
2.1.1 Herkenning	3
2.1.2 convolutioneel neuraal netwerk (CNN)	4
2.1.3 Trainen van een CNN	5
2.2 Deep learning-gebaseerde detector	6
2.2.1 Two-stage detector	6
2.2.2 One-stage detector	7
3 Herkenning en detectie implementatie op mobiel platform	9

3.1	CNN architecturen voor mobiele platformen	9
3.1.1	MobileNet	10
3.1.2	EfficientNet	10
3.1.3	TinyYOLO	10
3.2	Frameworks	10
3.2.1	Tensorflow	10
3.2.2	PyTorch	11
3.2.3	MMDetection	12
3.2.4	Core ML	13
3.2.5	Mobile AI Compute Engine (MACE)	13
3.2.6	KotlinDL	13
3.2.7	paddle lite	13
3.2.8	Detectron	13
3.2.9	Open Neural Network Exchange (ONNX)	13
3.3	Implementatie op mobiele platformen	13
3.4	Optimalisaties van neurale netwerken voor snelheid en bestadsgrootte	14
3.4.1	Pruning	15
3.4.2	Parameter quantisatie	15
3.4.3	Convolutionele filter compressie	16
3.4.4	Matrix factorisatie	16
A	Uitleg over de appendices	19

Lijst van figuren

2.1	CNN met 2 convolutie lagen en 2 pooling lagen en een fully-connected layer	4
2.2	Convolutie laag waarbij een filter wordt herleid tot een output feature.	4
2.3	ReLu, waarbij het maximum wordt genomen van 0 en de input waarde.	5
2.4	R-CNN	6
2.5	Faster R-CNN	7
2.6	YOLO waarbij de input is opgedeeld in een $S \times S$ rooster. En waarbij bounding box voorspellingen zijn gedaan.	8
2.7	One-stage detector met VGG net backbone	8
3.1	Implementatie flow van een getraind TensorFlow model naar de applicatie	11
3.2	Implementatie flow van een getraind PyTorch model naar de applicatie met code . .	12
3.3	CNN voor en na pruning	15

Lijst van tabellen

Lijst van symbolen

Maak een lijst van de gebruikte symbolen. Geef het symbool, naam en eenheid. Gebruik steeds SI-eenheden en gebruik de symbolen en namen zoals deze voorkomen in de hedendaagse literatuur en normen. De symbolen worden alfabetisch gerangschikt in opeenvolgende lijsten: kleine letters, hoofdletters, Griekse kleine letters, Griekse hoofdletters. Onderstaande tabel geeft het format dat kan ingevuld en uitgebreid worden. Wanneer het symbool een eerste maal in de tekst of in een formule wordt gebruikt, moet het symbool verklaard worden. Verwijder deze tekst wanneer je je thesis maakt.

b	Breedte	$[mm]$
A	Oppervlakte van de dwarsdoorsnede	$[mm^2]$
c	Lichtsnelheid	$[m/s]$

Lijst van afkortingen

CNN Rols ReLu

Hoofdstuk 1

Situering en doelstelling

1.1 Situering

Tegenwoordig wordt deep learning steeds meer en meer gebruikt om beeldverwerking problemen op te lossen. Via neurale netwerken kunnen we met meer en betere features werken om de afbeeldingen te analyseren. Maar veel van deze modellen hebben behoorlijk wat rekenkracht en geheugen nodig om te werken. Ook is er steeds meer interesse naar real-time toepassingen waarvan het resultaat zo snel mogelijk beschikbaar moet zijn. Dit wordt moeilijk bij veel hedendaagse systemen waarbij de foto eerst genomen moet worden en vervolgens door een computer geanalyseerd moet worden, omdat hedendaagse systemen veel rekenwerk en geheugen vragen. In deze masterproef wordt er onderzocht of de computer kan weggelaten worden en de afbeelding meteen door het mobiel apparaat geanalyseerd kan worden. Dus er moet onderzocht worden hoe een bestaand model aangepast kan worden om efficiënt te werken op een mobiel apparaat. Hierbij moet vooral rekening gehouden worden met de rekenkracht en geheugen van het mobiele apparaat.

1.2 Probleemstelling

Mobiele apparaten zijn kleine toestellen met beperkt geheugen en beperkte rekenkracht. In deze masterproef wordt er onderzocht hoe het rekenwerk beperkt kan worden zodat het resultaat real-time geleverd kan worden. Er gaat ook onderzocht worden hoe alle data efficiënt kan worden opgeslagen op het toestel.

1.3 Doelstellingen

Het uiteindelijke doel van deze masterproef is er voor zorgen dat een bestaand deep learning model aangepast kan worden zodat dit real-time resultaten kan geven op een mobiel apparaat. Dit gebeurt aan de hand van de volgende stappen:

- grondig begrijpen van een deep learning herkenningssysteem

- grondig begrijpen van een deep learning detectiesysteem
- Onderzoeken welke technieken er gebruikt kunnen worden om bestaande systemen op een mobiel apparaat te implementeren.
- onderzoeken voor optimalisaties voor een herkenningssysteem
- onderzoeken voor optimalisaties voor een detectiesysteem
- gevonden technieken testen en analyseren
- werkend prototype applicatie ontwerpen voor een mobiel apparaat

Hoofdstuk 2

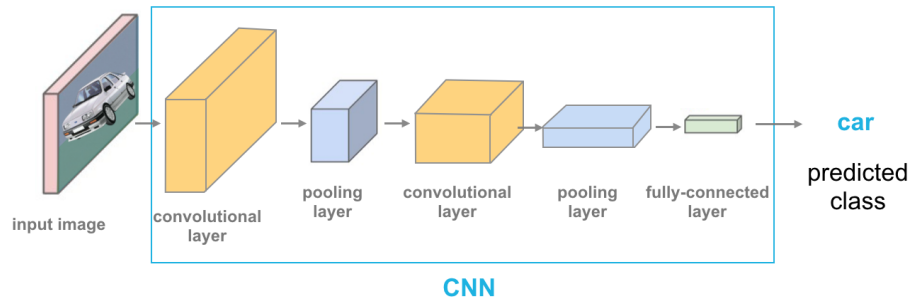
Herkenning en Detectie Algemeen

2.1 Deep learning-gebaseerde herkenningssystemen

Herkenningssystemen voorspellen wat de klasse van een object is in een afbeelding. Dus het herkennen van objecten in digitale afbeeldingen zonder deze te lokaliseren of aan te duiden. Bij herkenningssystemen is er geen of weinig overlap tussen de trainingsafbeeldingen en de inputafbeeldingen. Bijvoorbeeld bij een gezichts herkenningssysteem wordt er een algemeen herkenningssystemen ontworpen dat gezichten herkent, en niet elk individueel gezicht herkent. Voor een herkenningssysteem is er een goed getraind netwerk nodig dat input afbeeldingen omzet in features. Er moet een database zijn met daarin de gegevens van de objecten die men wilt herkennen. Vervolgens hebben is er ook een methode nodig om features van het neurale netwerk te vergelijken met de gegevens in de database om het juiste object te herkennen.

2.1.1 Herkenning

Wanneer er een getraind CNN is kan er een herkenningssysteem ontwikkeld worden. Als men bepaalde objecten in een afbeelding wil ontdekken gaat men met behulp van een CNN de afbeelding omzetten in een embedding. Embeddings Koehrsen (2018) zijn vector representaties die kunnen worden vergeleken in een embedding space, waar gelijkaardige objecten dicht bij elkaar liggen. De embedding van de input afbeelding wordt vergeleken met de embeddings die zich in een galerij bevinden. Jiang et al. (2019) geeft aan dat met behulp van een query kunnen er gelijkaardige objecten uit de galerij gehaald worden om deze te gaan vergelijken in een embedding space. De galerij is een database/verzameling met gekende embeddings/ID's van de objecten die men wilt herkennen. Een query is een embedding van de input waarvan het label niet gekend is. Gelijkaardige embeddings kunnen gezocht worden via de nearest neighbour techniek, waar we naar de klasse van de dichtsbijzijnde buur gaan kijken.

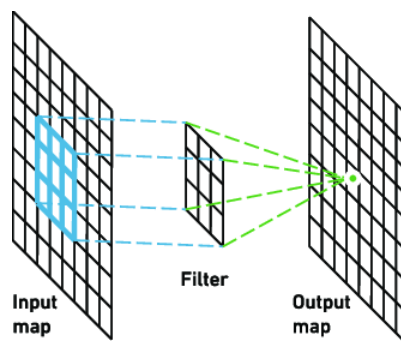


Figuur 2.1: CNN met 2 convolutie lagen en 2 pooling lagen en een fully-connected layer

2.1.2 convolusioneel neuraal netwerk (CNN)

De belangrijkste bouwsteen van een herkenningssysteem is een goed getrainde CNN, in dit hoofdstuk wordt het CNN besproken dat beschreven wordt door Jiang et al. (2019). Het algemeen model van een CNN is weergegeven in figuur 2.1. In tegenstelling tot fully connected netwerken wordt bij een CNN de gewichten gedeeld over verschillende locaties om zo het aantal parameters te verminderen.

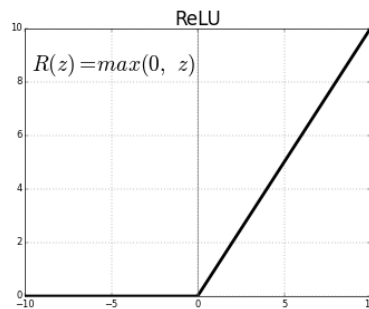
Het belangrijkste deel van een CNN zijn de convolutielagen (figuur 2.2) waarbij men een kernel/filter over de input laat gaan wat als output een feature map genereert. Een kernel bestaat uit set van gewichten die met de input worden vermenigvuldigd, deze kernel wordt over de input afbeelding geschoven. Al de pixels binnen het veld van de kernel worden gereduceerd tot een enkele waarde. CNN leren verschillende features met verschillende kernels in parallel. Waardoor de matrices met feature mappen steeds kleiner worden maar ook dieper worden. Een andere factor van een convolutie laag is de stride, deze waarde geeft aan met hoeveel pixels de kernel telkens moet doorschuiven. Een CNN bestaat uit een opeenvolging van een aantal convolutie lagen die steeds meer high-level features extraheren. Hoe meer convolutielagen een netwerk telt hoe meer features er uit de input worden gehaald, maar hoe trager het netwerk is.



Figuur 2.2: Convolutie laag waarbij een filter wordt herleid tot een output feature.

Elke convolutie laag wordt gevolgd door een niet-lineaire activatie functie, de meest gebruikt functie hiervoor is de rectified linear unit (ReLU) (figuur 2.3). De ReLU wordt vaak gebruikt omdat deze eenvoudig is, kan exact 0 weergeven en ziet er lineair uit. $\text{Max}(0, x)$ is de ReLU bewerking, dus

er wordt verdergegaan met 0 of de input waarde. Zonder een niet-lineaire activatie functie kan het CNN herleid worden tot 1 convolutie laag die geen high-level features kan extraheren. Andere mogelijkheden voor Lineaire activatie functies zijn: Sigmoid en Tangens hyperbolicus maar deze functies vragen meer rekenwerk.



Figuur 2.3: ReLu, waarbij het maximum wordt genomen van 0 en de input waarde.

Een volgende bouwsteen is de pooling laag waarbij het aantal samples in de feature map wordt verlaagt. De meest voorkomende methode is max-pooling waarbij er verder wordt gegaan met de maximum waarde in een bepaalde regio. Het doel van een pooling laag is om het aantal parameters te verminderen en zo ook het rekenwerk te verminderen. Er kan ook gebruik gemaakt worden van avarage pooling waarbij er verder wordt gegaan met de gemiddelde waarde van een regio. Er is ook minimal pooling waarbij er verder wordt gegaan met de minimum waarde.

Op het einde van elk CNN volgen er meestal 1 of meerdere fully connected lagen. Deze lagen connecteren elke input van één laag met elke activatie eenheid van de volgende laag. Dit zorgt voor meer parameters en meer rekenwerk waardoor deze lagen een vertragende factor vormen. De fully connected lagen zorgen voor een classificatie op basis van de features van de convolutie lagen.

2.1.3 Trainen van een CNN

Het trainen van een CNN bestaat uit het leveren van veel voorbeelden aan het netwerk. Op basis van het resultaat van deze voorbeelden worden telkens de gewichten van de kernels aangepast, zodat er steeds een beter resultaat wordt geleverd.

De loss functie geeft de error van de voorspelling weer tijdens het trainen van een neuraal netwerk. Op basis van de loss functie gaat men via de stochastic gradiënt decent de gewichten van het netwerken bijstellen zodat bij een volgende trainingsinput de loss functie een beter resultaat geeft. Bij de stochastic gradiënt decent wordt er per batch/group trainingsvoorbeelden de gewichten bijgesteld. De gradienten worden berekend door de loss af te leiden naar de gewichten via de ketting regel, en de gewichten worden bijgesteld volgens de tegengestelde gradiënt.

De learning rate bij het trainen van een CNN beïnvloed de grootte van de stap waarmee de gewichten worden bijgesteld. Hoe kleiner de learning rate hoe langer het trainen van een CNN duurt. Maar als de learning rate te hoog is kan het resultaat een slecht getraind netwerk zijn, omdat de

veranderingen op de gewichten te groot is om een beter resultaat te verkrijgen.

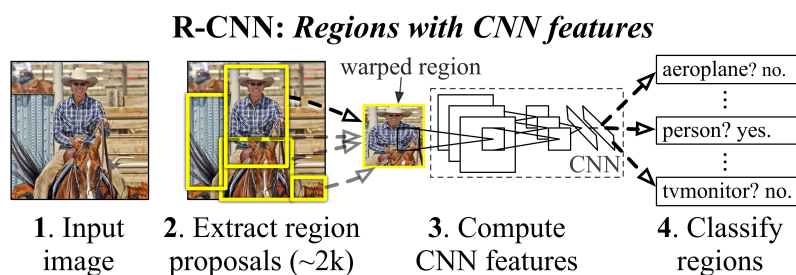
2.2 Deep learning-gebaseerde detector

Object detectie is het lokaliseren en classificeren van objecten in een afbeelding, waarbij de objecten aangeduid worden met een Bounding box. Door gebruik te maken van CNN kunnen er vrij nauwkeurige object detectoren ontworpen worden. Object detectie maakt voornamelijk gebruik van twee methodes: de single-stage detector en de two-stage detector.

2.2.1 Two-stage detector

Zoals de naam zegt bestaat deze methode uit 2 niveaus. Het eerste deel worden er Regions of Interest (Rols) gecreëerd, dit is het filteren van regio's waarbij de kans groot is dat deze een object bevatten. Het tweede deel classificeert en verfijnt de lokalisatie van de Rols die in het eerste deel gecreëerd werden. Dit gebeurt door elk van de Rols door een CNN te voeren. Region-based Convolutional Neural Network (R-CNN) Girshick (2015) is het basis principe van de two-stage detectoren weergegeven in figuur 2.4. Hierbij wordt met een region proposal algoritme regio's uit de afbeelding gefilterd waar de kans groot is dat er objecten op staan. R-CNN bestaat uit 3 stappen:

1. Via een selective search algoritme Uijlings et al. (2013) worden er ongeveer 2000 mogelijke regio's met objecten geselecteerd.
2. Elke mogelijke regio wordt omgezet naar een feature vector via een CNN.
3. Elke regio wordt vervolgens geclassificeerd met een klas-specifieke SVMs.

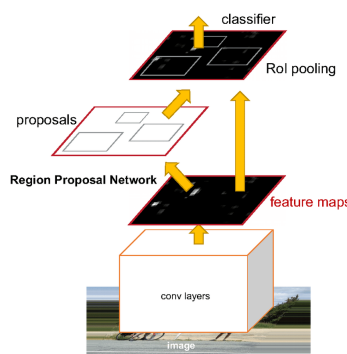


Figuur 2.4: R-CNN

R-CNN is een trage detector vermits elke Rols door een CNN moet gaan. Deze methode is geëvolueerd tot de veel snellere methode Faster R-CNN Ren et al. (2016). Hierbij wordt de afbeelding door een CNN behandeld en vervolgens wordt er gebruik gemaakt van een Region Proposal Network (RPN) dit is weergegeven in figuur 2.5. Het RPN gaat zoals bij R-CNN regio's uit de afbeelding filteren waar de kans groot is dat er objecten opstaan. Maar het RPN werkt sneller en levert betere resultaten dan het region proposal algoritme.

Het RPN is een deep fully convolutional network dat per input een set van regio's als output geeft. Elk van deze regio's heeft een objectness score wat een maat is voor het object t.o.v. de achtergrond in de afbeelding. Om een region proposal te genereren wordt het RPN over de feature map geschoven die gegenereerd is door het voorgaande CNN. Op elke sliding window locatie worden er meerdere regio voorspellingen gedaan. Deze voorspelling wordt gedaan door verschillende anchor boxes in een sliding window te evalueren.

Vervolgens worden Rols omgezet naar een feature vector met vaste lengte door RoI pooling. Elk van deze features gaat door een set van fully-connected lagen die 2 lagen als output heeft. een softmax laag die de klasse voorspelt, en een bounding box regressie laag die de bounding box voorspelt.



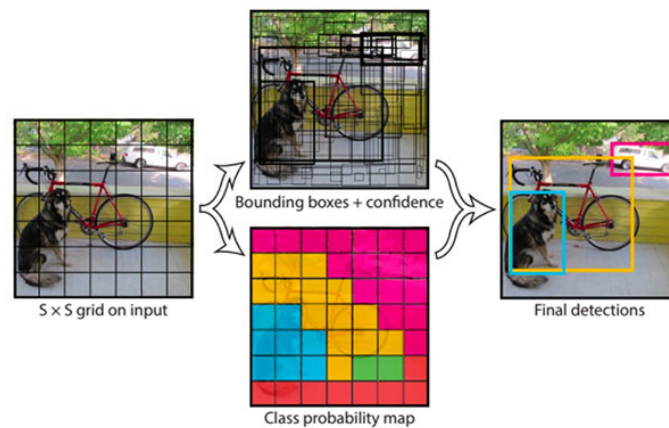
Figuur 2.5: Faster R-CNN

2.2.2 One-stage detector

Bij one-stage detectoren gebeurt object detectie in één keer met één neurale netwerk. Dus er is geen region proposal niveau meer zoals bij de two-stage detector. Deze detector gebruiken minder geheugen en rekenkracht t.o.v. two-stage detectoren. One-stage detectoren zijn sneller dan two-stage detectoren omdat ze alles in één keer doen, maar kunnen wat in nauwkeurigheid verliezen t.o.v. two-stage detectoren. De twee veel gebruikte technieken van one-stage detectie zijn: You Only Look Once (YOLO) en Single Shot Detection (SSD).

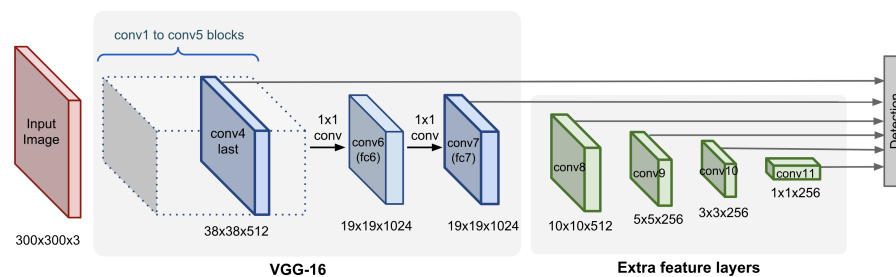
YOLO Redmon et al. (2016) verdeelt de afbeelding in een $S \times S$ rooster zoals in figuur 2.6 te zien is. De cel waarin het middelpunt van het object valt is verantwoordelijk voor de object detectie. Elke cel voorspelt K aantal bounding boxes en een score die aangeeft hoe zeker het model is dat een bepaalde bounding box een object bevat. Deze score wordt bepaald met de Intersection Of Union (IOU) tussen de voorspelde box en de ground truth box. Vervolgens is er nog een methode nodig om de overbodige bounding boxes te verwijderen. Een eerste mogelijkheid is door enkel bounding boxes te tekenen waarvan de voorspelling boven een threshold ligt. Een andere methode is non-maxima suppression deze methode zorgt ervoor dat elk object maar één bounding box heeft. Deze techniek houdt enkel de bounding box over met de beste voorspelling en onderdrukt de rest van de bounding boxes.

SSD Liu et al. (2016) is een one-stage detector (figuur 2.7) waarbij een afbeelding door verschil-



Figuur 2.6: YOLO waarbij de input is opgedeeld in een $S \times S$ rooster. En waarbij bounding box voorspellingen zijn gedaan.

lende convolutielagen gaat, wat als resultaat feature mappen op verschillende schalen oplevert. Op elke locatie van deze feature mappen van verschillende schalen wordt een vaste set van bounding boxes geëvalueerd. Voor elk van deze boxen wordt de zekerheid dat het een object bevat voorspeld. Op het einde wordt non maximum suppression gebruikt om de finale voorspelling te maken. Het netwerk van een SSD bestaat uit een basis netwerk dat gevormd wordt door een standaard classificatie netwerk zonder de fully-connected lagen. Vervolgens worden er extra convolutie lagen toegevoegd, wat het model toelaat om voorspellingen te doen op verschillende schalen.



Figuur 2.7: One-stage detector met VGG net backbone

Hoofdstuk 3

Herkenning en detectie implementatie op mobiel platform

Dit hoofdstuk zal gaan over het implementeren van deep learning herkenningssystemen en detectiesystemen op een mobiel platform. In dit hoofdstuk wordt er besproken welke technologieën er gebruikt kunnen worden om een CNN te implementeren op een mobiel platform. Dit gaan een aantal frameworks zijn die de programmeur in staat stelt om een bestaand model te implementeren op een mobiel apparaat. Deze frameworks gaan vergeleken worden om te kijken van welk framework het best gebruik wordt gemaakt voor een bepaalde toepassing. Er wordt ook onderzocht hoe een bestaand herkenningssystemen en detectiesystemen geoptimaliseerd kan worden zodat dit gebruikt kan worden op een mobiel platform. Bij het uitvoeren van een neurale netwerk op een mobiel apparaat zal men rekening moeten houden met de volgende zaken: gelimiteerde rekenkracht en beschikbaar geheugen. Ook moet er rekening gehouden worden met een beperkte batterij, want CNN's gebruiken veel bandbreedte en voeren veel berekeningen uit wat meer energie verbruikt. Dus er zal onderzocht moeten worden welke methodes men kan gebruiken om het geheugen van het model, het aantal bewerkingen en het energieverbruik te kunnen verbeteren.

3.1 CNN architecturen voor mobiele platformen

Dit deel gaat over CNN architecturen die specifiek ontworpen zijn voor mobiele en embedded toestellen. Deze netwerk architecturen streven naar zo weinig mogelijk parameters en een zo snel mogelijke uitvoering zonder een te groot effect op de accuraatheid te hebben. Het overlopen van deze netwerken voor de masterproef is iets minder relevant vermits we een reeds getraind CNN willen optimaliseren. Dus dat netwerk is reeds getraind, maar er kan wel gekeken worden welke technieken deze architecturen gebruiken zodat het model gebruikt kan worden op een mobiel platform.

3.1.1 MobileNet

MobileNet is voornamelijk opgebouwd uit diepe afzonderlijke convoluties (depthwise separable convolution), om het aantal berekeningen te verminderen. Deze convoluties bestaan uit een depthwise convolutie, waarbij één filter over elk input kanaal gaat waardoor er afzonderlijke feature maps ontstaan. Vervolgens gebeurt er een pointwise convolutie die met een 1x1 convolutie de outputs samenvoegt. Deze manier van werken zorgt voor een grootte vermindering in het aantal parameters met een klein effect op de accuraatheid. VGG-16 heeft 138 miljoen parameters en een accuraatheid van 71.5%, MobileNet heeft 4.2 miljoen parameters met een accuraatheid van 70.6%. Ondertussen is er al een MobileNetV2 ... en ook een MobileNetV3

3.1.2 EfficientNet

3.1.3 TinyYOLO

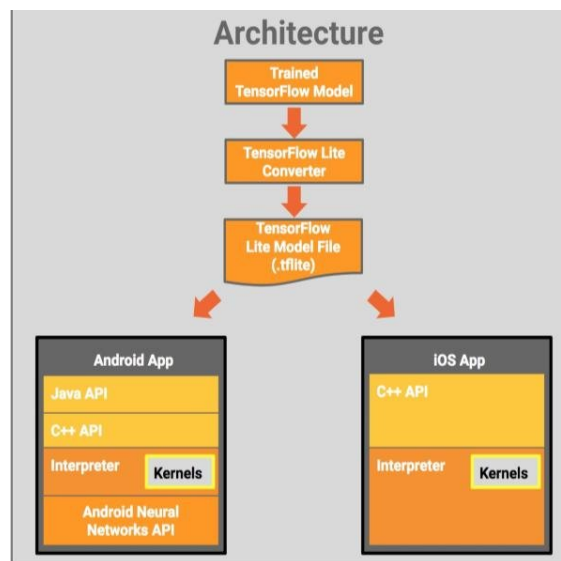
3.2 Frameworks

Om machine learning modellen te ontwerpen en te trainen kan er gebruik gemaakt worden van frameworks. Deze frameworks geven de programmeur een set van tools die hun in staat stelt om op een overzichtelijke en flexibele manier machine learning modellen te ontwerpen en trainen. In deze paragraaf worden enkele van de meest gebruikte frameworks besproken. Dit zijn enkele van de meest gebruikte frameworks, maar zeker niet allemaal.

3.2.1 Tensorflow

TensorFlow (Abadi et al. (2016)) is ontworpen door Google en is een open source library voor machine learning implementaties. TensorFlow focust op het trainen en het deployen van neurale netwerken. Ook ondersteunt TensorFlow meerdere programmeertalen zoals: Python, Java en C. Door de introductie van de Keras API is TensorFlow meer gebruiksvriendelijk geworden. Keras is een framework dat bovenop TensorFlow gebruikt kan worden, waarmee machine learning modellen kunnen ontworpen worden op een overzichtelijke manier. Zo heeft TensorFlow een gebruiksvriendelijke API voor eenvoudige projecten en meer uitgebreide tools voor complexe projecten. Door gebruik te maken van TensorBoard kan data op een flexibele manier gevisualiseerd worden tijdens runtime. TensorFlow biedt ook goede ondersteuning op gebied van deployment van een productie model. Nog een voordeel van TensorFlow is dat het een grootte community achter zich heeft, omdat dit een veel gebruikt framework is.

Maar één van de belangrijkste aspecten van TensorFlow voor deze masterproef is de mogelijkheid om een model te optimaliseren naar TensorFlow Lite (figuur: 3.1), dit model is bruikbaar op een mobiel toestel. TensorFlow Lite zorgt ervoor dat het model een lage latency heeft en een kleine binaire grootte. Het converteren kan makkelijk worden uitgevoerd via de volgende lijnen code in Python.



Figuur 3.1: Implementatie flow van een getraind TensorFlow model naar de applicatie

Listing 3.1: Converteren van TensorFlow naar een TensorFlow Lite model

```
import tensorflow as tf

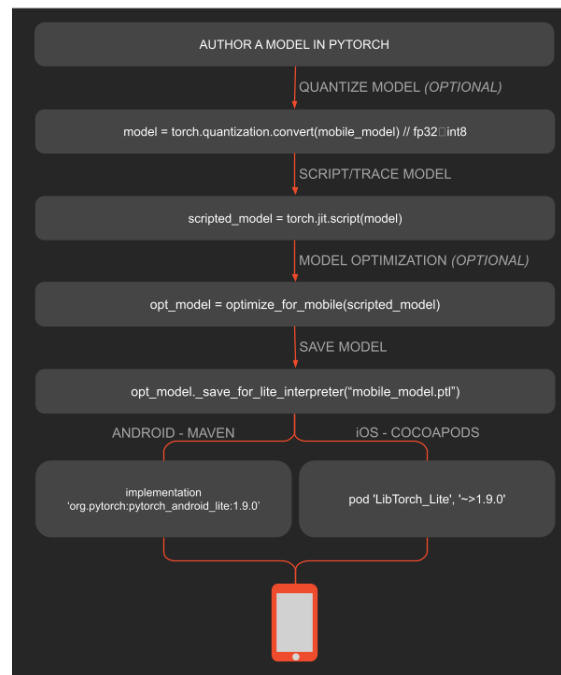
converter = tf.lite.TFLiteConverter.from_saved_model(saved_model_dir)
tflite_model = converter.convert()
```

Het TensorFlow Lite framework geeft ook de mogelijkheid om verdere optimalisaties zoals pruning en quantisatie uit te voeren na het converteren. Het TensorFlow Lite model is ook compatibel met android en IOS. Bij Android kan TensorFlow gebruik maken van de Neural Network API (NNAPI) ... die beschikbaar is vanaf Android 8.1. Deze API kan gebruikt worden om netwerk modellen te versnellen met de GPU, DSP en NPU.

3.2.2 PyTorch

PyTorch (Li et al. (2020)) is ontworpen door facebook en wordt zoals TensorFlow ook gebruikt voor machine learning implementaties. Dit is een python gebaseerd framework dat focust op flexibiliteit, maar deze extra flexibiliteit zorgt voor meer lijnen code. Door zijn flexibiliteit is het gemakkelijk om nieuwe functionaliteiten toe te voegen door bestaande code aan te passen of nieuwe code toe te voegen. PyTorch maakt gebruik van externe tools zoals TensorBoard om data te visualiseren. Dit framework wordt gebruikt om een machine learning model te ontwerpen en trainen. Het model dat is ontworpen kan vervolgens gebruikt worden gebruikt om een applicatie te ontwerpen.

Zoals TensorFlow biedt PyTorch ook de mogelijk om het model te optimaliseren naar PyTorch Mobile. In figuur: 3.2 is te zien welke lijnen code er moeten worden uitgevoerd om een PyTorch model te converteren naar een PyTorch mobile model. Maar PyTorch mobile zit nog in zijn beta fase, dus hierbij kunnen er onverwachte complicaties optreden. Een ander framework is Caffe2 dit



Figuur 3.2: Implementatie flow van een getraind PyTorch model naar de applicatie met code

framework focust ook op mobiele implementaties, maar het framework is ondertussen geïntrigeerd met PyTorch. Ook is het PyTorch framework compatibel met android en IOS.

PyTorch wordt voornamelijk gebruikt om te experimenteren met neurale netwerken. Bedrijven en onderzoekers gebruiken dit framework vooral om een CNN experimenteel op te bouwen en te trainen. TensorFlow wordt voornamelijk gebruikt om een model effectief in gebruik te nemen. De TensorFlow Lite conversie geeft meer zekerheid t.o.v. de PyTorch Mobile conversie vermits deze nog in zijn Beta fase zit.

3.2.3 MMDetection

MMDetection maakt deel uit van OpenMMLab en is een open source object detectie toolbox gebaseerd op PyTorch. De toolbox bevat gewichten van meer dan 200 voorgetrainde modellen. Via modulair ontwerp kan het detectie framework opgesplitst worden in verschillende componenten. Met de verschillende componenten kan een eigen detectie model worden gemaakt door de verschillende componenten te combineren. MMDetection ondersteund ook veel verschillende detectie methodes zoals: YOLO, Faster R-CNN, etc. MMDetection biedt geen ondersteuning om een bestaand model te optimaliseren naar een model voor mobiele toepassingen. Dus dit model zal geconverteerd moeten worden naar een framework waarbij het optimaliseren naar een mobiel model wel mogelijk is.

3.2.4 Core ML

Core ML is het Apple framework om machine learning tools te integreren in een applicatie. Dit kan een model zijn van Create ML het machine learning framework van Apple zelf. Maar Core ML biedt ondersteuning om modellen te converteren van TensorFlow, PyTorch en ONNX naar Core ML. Uiteraard is dit framework enkel van toepassing voor Apple, en in deze masterproef wordt er vooral gefocust op een Android implementatie.

3.2.5 Mobile AI Compute Engine (MACE)

Het MACE framework is ontworpen door Xiaomi en dient specifiek voor mobiele toepassingen van neurale netwerken op Android, IOS, Linux en Windows. MACE biedt ondersteuning voor verschillende frameworks zoals: TensorFlow, Caffe en ONNX. Het model kan geïmplementeerd worden op Android, IOS en Linux. Het MACE framework bespaart geheugen door de core library zo klein mogelijk te maken door het aantal externe dependencies te minimaliseren. Het Winograd algoritme ... wordt gebruikt om convolutie bewerkingen te versnellen en verbeterd op deze manier de latency van het CNN model.

3.2.6 KotlinDL

3.2.7 paddle lite

3.2.8 Detectron

3.2.9 Open Neural Network Exchange (ONNX)

Om één van deze frameworks te gebruiken voor een mobiele implementatie kan er gebruik gemaakt worden van de bovenstaande frameworks. Maar de kans bestaat dat het bestaande model niet gemodelleerd is in één van de bovenstaande frameworks. Daarom is er een methode nodig die een willekeurig model omzet in één van de bovenstaande modellen. ONNX (2017) biedt de mogelijkheid om verschillende tools/frameworks samen te laten werken. Hierbij wordt een bestaand model geconverteerd naar een ONNX model, en dit model kan op zijn beurt geconverteerd worden naar het gewilde framework. ONNX ondersteunt niet elk machine learning framework, maar toch wel de meest bekende. Het ONNX framework biedt zelf ook de mogelijkheid om een model te deployen en te optimaliseren voor mobiel gebruik.

3.3 Implementatie op mobiele platformen

Uit paragraaf 3.1 kan er afgeleidt worden dat sommige frameworks ondersteuning bieden voor het optimaliseren van het machine learning model naar een lichtere versie. Dus het bestaande CNN model zou naar één van deze frameworks geconverteerd moeten worden zodat er gebruik gemaakt kan worden van de mobiele optimalisatie die dat framework ondersteunt. Niet elk framework zal

even compatibel zijn met het bestaande model, er zal dus gekeken moeten worden tussen welke frameworks een conversie mogelijk is. Ook zal niet elk framework de optimalisatie voor mobiele platformen op dezelfde manier toepassen. Dus het converteren van het standaard model naar het mobiele model zal voor elk framework een ander resultaat geven. In deze paragraaf zal er besproken worden welke frameworks er compatibel zijn met elkaar en welke frameworks het beste optimaliseren voor een mobile implementatie.

We hebben al gezien dat voor het converteren van een model niet veel lijnen code nodig zijn. Maar dit is enkel het geval als het model in het zelfde framework is ontworpen. Als het model in een ander framework is ontworpen en getraind is de eerste stap van het process om het CNN model naar het gewilde framework te converteren. ONNX dient goed als overkoepelend framework dat er voor zorgt dat verschillende frameworks compatibel zijn met elkaar. Volgens de website zijn er 23 frameworks die naar het ONNX framework geconverteerd kunnen worden en een CNN kunnen ontwerpen en trainen. Om dan geconverteerd te worden naar een framework dat het model kan omzetten naar een mobile versie. Een ander manieren om een model te gebruiken over verschillende frameworks in plaats van ONNX hangt af van de compatibiliteit tussen frameworks. Bepaalde frameworks zullen zelf een methode hebben om modellen van andere frameworks in te laden.

In de paper geschreven door Luo et al. (2020) wordt PyTorch Mobile vergeleken met TensorFlow Lite voor verschillende netwerk architecturen en een model getraind met dezelfde data. Voor alle netwerk architecturen in deze paper geeft de optimalisatie naar TensorFlow Lite de kleinste bestand grootte van het model. Uit deze paper is ook af te leiden dat de optimalisatie naar een mobiel model niet alleen afhankelijk is van het framework maar ook van de netwerk architectuur. Zo geeft TensorFlow Lite volgens Luo et al. (2020) betere latency resultaten voor de zwaardere netwerken (ResNet50, InceptionV3, DenseNet121) dan PyTorch Mobile. Maar PyTorch Mobile heeft op zijn beurt wel een betere latency voor SqueezeNet en MobileNetV2. Dus uit deze paper kunnen we afleiden dat TensorFlow Lite het beste de bestands grootte verkleint, maar dat de netwerk architectuur ook een rol speelt.

Febvay (2020) vergelijkt TensorFlow Lite met MACE voor verschillende neurale netwerken (SqueezeNet, MobileNetV1/V2). Hierbij geeft TensorFlow Lite het beste resultaat, TensorFlow Lite gaf een Top-1 resultaat van 69,19% en MACE gaf een top-1 resultaat van 66.84% bij MobileNetV1. Ook voor de latency gaf TensorFlow Lite in de meeste gevallen de beste resultaten buiten bij het gebruik van 4 of 6 CPU cores, dan gaf MACE betere resultaten.

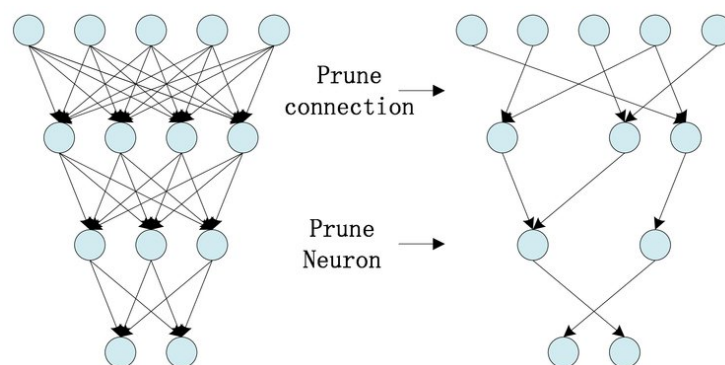
3.4 Optimalisaties van neurale netwerken voor snelheid en bestands-grootte

In deze paragraaf wordt er onderzocht welke optimalisaties er kunnen worden toegepast om de accuraatheid, snelheid en gebruikt geheugen te verbeteren. Maar het optimaliseren van een bepaalde factor zal vaak negatieve gevolgen hebben voor een andere factor, dit zal meestal de accuraatheid zijn. Dus er zal een goede balans gevonden moeten worden tussen de optimalisatie en

de negatieve gevolgen op de andere factoren.

3.4.1 Pruning

Pruning is de eerste stap van de Deep compression methode voorgesteld door Han et al. (2016). Bij het trainen van een CNN hebben bepaalde gewichten een grotere invloed op het resultaat. Andere gewichten hebben weinig tot geen invloed op het resultaat. Maar alle gewichten worden steeds berekend ongeacht hun invloed op het resultaat. Bij pruning worden de gewichten met een kleine invloed op het resultaat verwijderd dit is weergegeven op figuur 3.3. Waardoor er geen berekeningen meer moeten uitgevoerd worden voor de verwijderde gewichten. Deze parameters moeten dan ook niet meer worden bijgehouden waardoor het CNN model minder geheugen in beslag neemt. Eerst wordt het CNN op een normale manier getraind zonder pruning. Vervolgens worden al de kleine gewichten onder een bepaalde threshold verwijderd. volgens Han et al. (2016) wordt voor VGG-16 het aantal parameters met factor 13 verminderd, voor AlexNet met een factor 9. Deze methode heeft zeer weinig tot geen effect op de accuraatheid.



Figuur 3.3: CNN voor en na pruning

3.4.2 Parameter quantisatie

Het quantiseren en delen van gewichten is een tweede methode voorgesteld door Han et al. (2016). Een CNN bestaat uit miljoenen gewichten, en de waarde van elke van deze gewichten moeten op het systeem worden opgeslagen. De default representatie van een waarde wordt opgeslagen als een floating point nummer wat 4 bytes in beslag neemt. Dus voor miljoenen parameters hebben de gewichten veel schijfruimte nodig. Een mogelijke oplossing hiervoor is quantiseren van gewichten, waarbij de getal representatie van de gewichten wordt veranderd naar fixed point. Hierbij worden de waarden van gewichten beperkt tot een set van beschikbare waardes. Waarbij de waardes éénmalig worden opgeslagen en al de gewichten refereren naar een waarde van de vaste set met waardes. Hoe kleiner de set met waardes is hoe minder geheugen er in beslag wordt genomen, maar een kleinere set van waardes zorgt ook voor een mindere accuraatheid. Dus de grootte van de set moet goed worden gekozen zodat er niet te veel geheugen wordt gebruikt met een accepteerbare

daling in accuratie. Han et al. (2016) past vervolgens Huffman encoding toe die een compressie uitvoert op de gekwantiseerde parameters.

3.4.3 Convolutionele filter compressie

een andere methode voorgesteld is Compressed Convolutional Filters. Hierbij wordt de grootte van de kernel verkleind om het aantal parameters en rekenwerk te verminderen. Maar door de kernels te verkleinen daalt de accuraatheid van het CNN.

3.4.4 Matrix factorisatie

Hierbij worden grootte en complexe matrices opgesplitst in verschillende kleinere en simpelere matrices.

Bibliografie

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., Kudlur, M., Levenberg, J., Monga, R., Moore, S., Murray, D. G., Steiner, B., Tucker, P., Vasudevan, V., Warden, P., Wicke, M., Yu, Y., and Zheng, X. (2016). TensorFlow: A system for large-scale machine learning. *arXiv:1605.08695 [cs]*. arXiv: 1605.08695.
- Febvay, M. (2020). Low-level Optimizations for Faster Mobile Deep Learning Inference Frameworks. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 4738–4742, Seattle WA USA. ACM.
- Girshick, R. (2015). Fast R-CNN. *arXiv:1504.08083 [cs]*. arXiv: 1504.08083.
- Han, S., Mao, H., and Dally, W. J. (2016). Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding. *arXiv:1510.00149 [cs]*. arXiv: 1510.00149.
- Jiang, X., Hadid, A., Pang, Y., Granger, E., and Feng, X. (2019). *Deep Learning in Object Detection and Recognition*, edited by Xiaoyue Jiang, Abdenour Hadid, Yanwei Pang, Eric Granger, Xiaoyi Feng. Springer Singapore : Imprint: Springer, Singapore, 1st ed. 2019. edition.
- Koehrsen, W. (2018). Neural Network Embeddings Explained.
- Li, S., Zhao, Y., Varma, R., Salpekar, O., Noordhuis, P., Li, T., Paszke, A., Smith, J., Vaughan, B., Damania, P., and Chintala, S. (2020). PyTorch distributed: experiences on accelerating data parallel training. *Proceedings of the VLDB Endowment*, 13(12):3005–3018.
- Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.-Y., and Berg, A. (2016). SSD: Single Shot MultiBox Detector. volume 9905, pages 21–37.
- Luo, C., He, X., Zhan, J., Wang, L., Gao, W., and Dai, J. (2020). Comparison and Benchmarking of AI Models and Frameworks on Mobile Devices. *arXiv:2005.05085 [cs, eess]*. arXiv: 2005.05085.
- ONNX (2017). ONNX Tutorials. original-date: 2017-11-15T18:59:18Z.
- Redmon, J., Divvala, S., Girshick, R., and Farhadi, A. (2016). You Only Look Once: Unified, Real-Time Object Detection. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 779–788. ISSN: 1063-6919.

- Ren, S., He, K., Girshick, R., and Sun, J. (2016). Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *arXiv:1506.01497 [cs]*. arXiv: 1506.01497.
- Uijlings, J. R., R, van de Sande, K. E., A, Gevers, T., Smeulders, A. W., and M (2013). Selective Search for Object Recognition. *International Journal of Computer Vision*, 104(2):154–171. Num Pages: 154-171 Place: New York, Netherlands Publisher: Springer Nature B.V.

Bijlage A

Uitleg over de appendices

Bijlagen worden bij voorkeur enkel elektronisch ter beschikking gesteld. Indien essentieel kunnen in overleg met de promotor bijlagen in de scriptie opgenomen worden of als apart boekdeel voorzien worden.

Er wordt wel steeds een lijst met vermelding van alle bijlagen opgenomen in de scriptie. Bijlagen worden genummerd met een drukletter A, B, C,...

Voorbeelden van bijlagen:

Bijlage A: Detailtekeningen van de proefopstelling

Bijlage B: Meetgegevens (op USB)

FACULTEIT INDUSTRIËLE INGENIEURSWETENSCHAPPEN
CAMPUS DE NAYER SINT-KATELIJNE-WAVER
J. De Nayerlaan 5
2860 SINT-KATELIJNE-WAVER, België
tel. + 32 15 31 69 44
iiw.denayer@kuleuven.be
www.iw.kuleuven.be

