

HybridNets: End-to-End Perception Network

Thanh Dat Vu Viet Hoai Bao Ngo Thesis Defense

Advisor Assoc. Prof. Duy Hung Phan

Introduction

Main Problem

Literature Review

Proposed Methods

Experiments

Discussion

Introduction

Proposed Methods

Network Architecture

Loss Function

Training Strategy

Experiments

Discussion

Introduction

Proposed Methods

Experiments

Cost Computation

Quantitative Results

Video Comparison

Discussion

Introduction

Proposed Methods

Experiments

Discussion

Community

The advantages of HybridNets

The limitations of HybridNets

Introduction

Proposed Methods

Experiments

Discussion

Conclusion

Implications For Future Work

Introduction



Two separate models for object detection and semantic segmentation

Cost Computation



Cost computation between two separate models are different

Synchronous Problem

Assume that Segmentation Model is faster than Object Detection Model



Decision Making

Assume that Segmentation Model is faster than Object Detection Model



Our purpose

- Build End-to-End Perception Network for Multi-tasks.
- Our work focus on Autonomous Driving Perception.





Road







• The base feature extractor is made of a series of convolutional layers which are shared between all tasks, and the extracted features are used as input to task-specific output heads.



¹Zhanpeng Zhang et al., Facial landmark detection by deep multi-task learning, ECCV 2014

• Each task has a separate network, but cross-stitch units combine information from parallel layers of the different task networks with a linear combination is called Share Tunk problem.



 1 lshan Misra et al., FCross-stitch networks for multi-task learning, IEEE Conference on Computer Vision and Pattern Recognition 2016

 Instead of combining information from different task networks with a linear combination of parallel features (as in Cross-Stitch networks (Misra et al., 2016)), NDDR-CNN uses concatenation and a 1x1 convolution to fuse features from separate task networks.



¹Gao et al., Nddr-cnn: Layerwise feature fusing in multi-task cnns by neural discriminative dimensionality reduction, IEEE Conference on Computer Vision and Pattern Recognition 2019

• Preliminary predictions are made for four tasks, then these predictions are re-combined and used to compute final, refined predictions for two output tasks.



 $^1{\rm Xu}$ et al., Pad-net: Multi-tasks guided prediction-and-distillation network for simultaneous depth estimation and scene parsing, IEEE Conference on Computer Vision and Pattern Recognition 2018

MultiNet

• Limitation: Fixed input size due to cell-based method.



¹Marvin Teichmann et al., MultiNet: Real-time Joint Semantic Reasoning for Autonomous Driving, IEEE Intelligent Vehicles Symposium 2018

DLT-Net

• Limitation: Shared context tensor from drivable area to other decoders, making finetuning of specific tasks harder.



¹Yeqiang Qian, John M. Dolan, Ming Yang, DLT-Net: Joint Detection of Drivable Areas, Lane Lines, and Traffic Object, IEEE Transactions on Intelligent Transportation Systems 2020

YOLOP

• Limitation: Two different segmentation heads for drivable area and lane line.



 $^{^{1}\}mathrm{Dong}$ Wu et al., YOLOP: You Only Look Once for Panoptic Driving Perception, arXiv 2021

Proposed Methods

HybridNets Architecture



Backbone Network



EfficientNet



¹Mingxing Tan, Quoc V. Le, EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks, ICML 2019





 $^{^1 {\}rm Mingxing}$ Tan, Ruoming Pang, Quoc V. Le, EfficientDet: Scalable and Efficient Object Detection, CVPR 2020

Feature Network Design



 $^{^1 {\}rm Mingxing}$ Tan, Ruoming Pang, Quoc V. Le, EfficientDet: Scalable and Efficient Object Detection, CVPR 2020





Why this research used Bi-FPN instead of FPN, PA-Net, NAS-FPN?



When fusing features with different resolutions, a common way is to first resize them to first resize them to the same resolution and then simply sum them up treating all input features equally without distinction, as named

Weighted Feature Fusion

HybridNets Architecture Overview



Detection Head



Anchor Boxes is calculated by Anchor Scales and Anchor Ratios based on K-Means algorithm.

Anchor Boxes Calculation



Base scale = 1.25 for each feature map level Anchor Scales = [2**0, 2**0.70, 2**1.32] Anchor Ratios = [(0.62, 1.58), (1.0, 1.0), (1.58, 0.62)]

Anchor Boxes Visualization



- This method work wells with small and huge objects.
- Become a good base to help network convergence fast.

Segmentation Head



• Why our work up scale on P2 Level and add P2 feature map from backbone into output feature map?

Object Detection Loss Function

$$\begin{split} \mathcal{L}_{all} &= \alpha \mathcal{L}_{det} + \beta \mathcal{L}_{seg} \\ \\ \mathcal{L}_{det} &= \alpha_1 \mathcal{L}_{class} + \alpha_2 \mathcal{L}_{obj} + \alpha_3 \mathcal{L}_{box} \\ \\ \mathbf{FL}(p_{\mathbf{t}}) &= -\alpha_{\mathbf{t}} (1 - p_{\mathbf{t}})^{\gamma} \log(p_{\mathbf{t}}). \\ \\ \\ \mathbf{Smooth}_{L1}(\mathbf{x}) &= \begin{cases} \delta_1 x^2 & if \quad x < \delta_2 \\ x - \delta_1 \\ \mathbf{x} = \delta b_p \cdot |b_x - \delta_x| + |b_y - \delta_y| + |b_w - \delta_w| + |b_h - \delta_h| \end{cases} \end{split}$$

Segmentation Loss Function









Training Strategy Overview



1 (11.0)

Experiments

Dataset

• Dataset: BDD100K.



- Merge four classes: Car, truck, bus, train into a single class (vehicle).
- Merge two drivable area {direct, alternative} into drivable, re-labelled two lane line annotations into a central one line.
- Basic augmentation techniques: rotating, flipping, HSV shifiting, Mosaic, Mixup.

HybridNets Training Details

• Data augmentation:







Phase 2



Phase 3

Phase 1: Mosaic + Mixup + Flip + Affine Phase 2: Flip Phase 3: Flip + Random(Blur, Gray, Brightness)

• Training stats:

- Training time: 750 hours (1 GPU 3090 RTX)
- GPU memory usage: 24 GB
- Batch size: 16

• Precision and Recall



• IoU metric for Drivable Area Segmentation and Lane Line Segmentation.



- AP = $\sum_{k=0}^{k=n-1}$ [Recalls(k) Recalls(k+1)] * Precisions(k).
- Recalls(n) = 0, Precisions(n) = 1 where n = Number of thresholds.



• mAP50 is used for Traffic Object Detection task.



Cost computation



Parameters FLOPs



Why HybridNets has less cost computation than YOLOP?

Depthwise Separable Convolutions



- In depth-wise convolution, we use each filter channel only at one input channel.
- Point-wise Convolution is a type of convolution that uses a 1x1 kernel: a kernel that iterates through every single point.

Traffic Object Detection Task



Traffic Object Detection

Recall MAP50

Traffic Object Detection Visualization

False Positive False Negative



a) YOLOP



b) HybridNets

Lane Segmentation Task



Lane Line Segmentation

Lane Segmentation Visualization









a) YOLOP



b) HybridNets

Drivable Segmentation Task



Drivable Area Segmentation

Drivable Segmentation Visualization

False Positive False Negative



a) YOLOP



b) HybridNets

Comparison Between HybridNets and Prior SOTA YOLOP



Discussion

Research

- Show once again that multi-tasking is an underexplored field with untapped potential in autonomous driving.
- Currently, the preprint on arXiv have yet to receive a citation.

Production

- Create a scalable network on a wide range of embedded systems, from industrial-grade edge computing to off-the-shelf mobile phones.
- Porting to TensorRT, allowing streamlined user experience.

Community Interaction



Github statistics as of April 10th

Community Contribution

- Github.com/PINTO0309
- Provided a script to convert model to TFLite, ONNX, OpenVINO, CoreML, TFJS.
- Stress-tested on a foggy video clip with real-time performance (ONNX 384x512, Input 720x1280, FP16).



Community Contribution

- Github.com/iwatake2222
- Created a demo environment entirely in C++, using ported model from PINTO0309.
- Stress-tested on Pixel 4a, only achieved 6 FPS maximum (TFLite 384x640, Input 720x1280, FP16).



- Github.com/ibaiGorordo
- Created a Python script to infer Youtube videos with customizable bird eye view, using ported model from PINTO0309.



Advantages

- Vehicle's perception in a single pass.
- Scalable thanks to EfficientNet backbone.
- Switchable backbone.



Limitations

- Fixed to 2 heads for 2 specific tasks of object detection and segmentation.
- Lack of ablation studies.





Future Work

• Perform various tasks related to perception.



- Perform various tasks related to perception.
- Improve parameters and FLOPs of network for edge devices.
- Build a decoder network detecting 3D Object Detection with only one input.

Questions?