



گزارش پروژه نهایی

محمد حسین شکوهی - ۹۹۲۰۲۵۶۷

یادگیری عمیق - دکتر عمادالدین فاطمی زاده

دانشگاه صنعتی شریف

فهرست مطالب

۳	شبکه تشخیص عمق
۴	مقدمه
۴	معماری شبکه
۵	تابع Loss
۶	Data Augmentation
۶	دیتاست مورد استفاده
۶	جزئیات پیاده سازی
۷	ارزیابی
۹	شبکه تشخیص اشیاء
۱۰	استفاده از کتابخانه ی Detectron2
۱۰	مدل مورد استفاده
۱۱	مقدمه
۱۱	معماری شبکه
۱۲	ارزیابی
۱۳	اتصال دو شبکه به یکدیگر
۱۴	ارزیابی
۱۵	محیط گرافیکی
۱۵	گیت هاب

شبکه تشخیص عمق:

به منظور یافتن شبکه و مدل مناسب برای تشخیص عمق روی دیتاست NYUv2، تحقیقات گسترده ای توسط اینجانب انجام شده است. در یکی از مقاله^۱ ها، تحلیل ها و بررسی های مختلفی بر روی نحوه کارکرد مدل MonoDepth انجام شده است و از آن به عنوان یکی از برترین مدل های موجود در حوزه تشخیص عمق یاد شده است. در یکی دیگر از مقاله^۲ ها، از مدل های MonoDepth و PSMNet به عنوان مطرح ترین مدل های تشخیص عمق یاد شده است. در یک سایت^۳ معتبر، مدل های تشخیص عمق مختلف بر اساس عملکردشان روی دیتاست NYUv2 رتبه بندی شده اند. در نهایت، پس از بررسی های متعدد و سنجش میزان پیچیدگی و عملکرد مدل های مختلف، برای پیاده سازی شبکه تشخیص عمق از مدل DenseDepth^۴ استفاده شده است. لازم به ذکر است که در سایت ذکر شده، همانطور که در شکل ۱ مشاهده می شود، عملکرد مدل DenseDepth روی دیتاست NYUv2، از بین ۳۰ شبکه موجود رتبه ۱۱ را به خود اختصاص داده است که نوید از یک مدل قوی می دهد. معماری و جزئیات پیاده سازی این مدل در ادامه شرح داده خواهد شد.

Rank	Model	absolute relative error	RMSE	log 10	Delta < 1.25	Delta < 1.25^2	Delta < 1.25^3	Extra Training Data	Paper	Code	Result	Year	Tags
1	LeReS	0.09			0.916			✓	Learning to Recover 3D Scene Shape from a Single Image	GitHub	arXiv	2020	
2	GLPDepth	0.098	0.344	0.042	0.915	0.988	0.997	×	Global-Local Path Networks for Monocular Depth Estimation with Vertical CutDepth	GitHub	arXiv	2022	
3	AdaBins	0.103	0.364	0.044	0.903	0.984	0.997	×	AdaBins: Depth Estimation using Adaptive Bins	GitHub	arXiv	2020	
4	DPT-Hybrid	0.110	0.357	0.045	0.904	0.988	0.998	✓	Vision Transformers for Dense Prediction	GitHub	arXiv	2021	
5	BTS	0.110	0.392	0.047	0.885	0.978	0.994	×	From Big to Small: Multi-Scale Local Planar Guidance for Monocular Depth Estimation	GitHub	arXiv	2019	
6	LapDepth	0.110	0.393	0.047	0.885	0.979	0.995	×	Monocular Depth Estimation Using Laplacian Pyramid-Based Depth Residuals	GitHub	arXiv	2021	
7	VNL	0.111	0.416	0.048	0.875	0.976	0.994	×	Enforcing geometric constraints of virtual normal for depth prediction	GitHub	arXiv	2019	
8	SC-DepthV2	0.138	0.532	0.059	0.820	0.956	0.989	×	Auto-Rectify Network for Unsupervised Indoor Depth Estimation	GitHub	arXiv	2021	
9	U-Net		0.382					×	PhaseCam3D — Learning Phase Masks for Passive Single View Depth Estimation	GitHub	arXiv	2019	
10	DSN		0.429					✓	On Deep Learning Techniques to Boost Monocular Depth Estimation for Autonomous Navigation	GitHub	arXiv	2020	
11	DenseDepth		0.465					×	High Quality Monocular Depth Estimation via Transfer Learning	GitHub	arXiv	2018	

شکل ۱ - رتبه بندی عملکرد شبکه های مختلف روی دیتاست NYUv2

¹ [How do neural networks see depth in single images?](#)

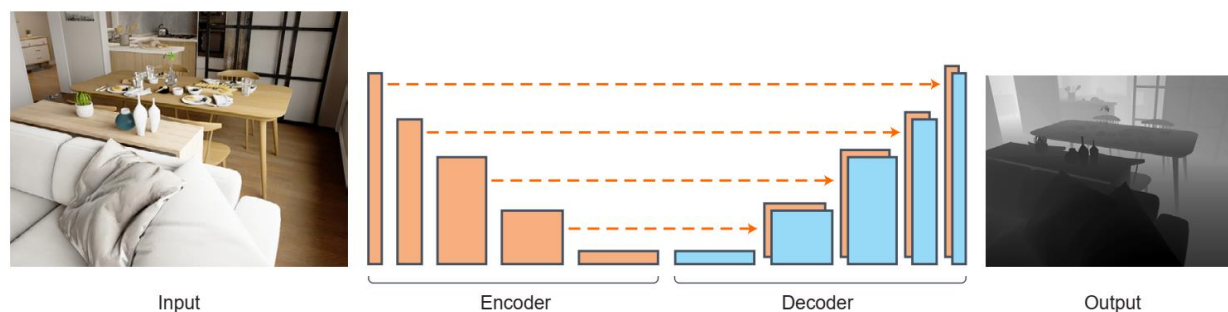
² [Joint Object Detection and Depth Estimation in Multiplexed Image](#)

³ [Monocular Depth Estimation on NYU-Depth V2](#)

⁴ [High Quality Monocular Depth Estimation via Transfer Learning](#)

مقدمه: استخراج دقیق عمق از تصاویر دارای کاربرد های متعددی در حوزه های مختلف می باشد که مدل های متعددی برای آن پیشنهاد شده است. مدل DenseDepth از یک شبکه CNN و ایده Transfer Learning برای تولید نقشه عمق تصاویر RGB با رزولوشن بالا استفاده می کند. در واقع، این مدل از یک معماری استاندارد Encoder-Decoder به همراه Augmentation های مختلف استفاده می کند تا به دقت بالایی در تشخیص عمق دست یابد. یکی از اهداف مهم این مدل، استفاده از شبکه های نسبتاً ساده برای دستیابی به دقت بالا است که سرعت یادگیری را به شدت افزایش می دهد. به منظور دستیابی به این هدف، مدل DenseDepth از ایده Transfer Learning استفاده می کند و برای بخش Encoder از شبکه های از پیش آموزش دیده با دقت بالا که در اصل برای دسته بندی تصاویر آموزش دیده اند استفاده می کند. نتایج اولیه نشان می دهد که ایده استفاده از شبکه های Encoder ساده که در اصل برای اهداف دیگری از جمله دسته بندی تصاویر آموزش دیده اند، منجر به دست یابی به دقت بالایی می شود. در ادامه با ارائه Decoder مناسب و آموزش آن بر روی دیتاست NYUv2، شبکه تکمیل شده و نقشه عمق تصاویر در خروجی تولید خواهد شد. در این مدل، همانطور که پیش تر اشاره شد، از یک شبکه Encoder-Decoder ساده به همراه Skip Connection استفاده می شود.

معماری شبکه: شکل ۲ یک نمای کلی از معماری شبکه مورد استفاده را نشان می دهد. برای بخش Encoder، از شبکه DenseNet-169 که بر روی دیتاست ImageNet از پیش آموزش دیده است استفاده می شود تا تصویر RGB ورودی را به یک بردار ویژگی ها انکود کند. در ادامه، این بردار از لایه های Up-sampling متعدد عبور می کند تا نقشه عمق در خروجی تولید شود. این لایه های Up-sampling به همراه Skip Connection های متناظر، Decoder ما را تشکیل می دهند که توسط خودمان بر روی دیتاست NYUv2 آموزش خواهد دید. لازم به ذکر است که در بخش Decoder از Batch Normalization استفاده نشده است.



شکل ۲ - ما از یک شبکه ساده Encoder-Decoder به همراه Skip Connection استفاده می کنیم. برای بخش Encoder از شبکه ای از پیش آموزش دیده DenseNet-169 استفاده می شود. بخش Decoder از لایه های کانولوشنی متعدد استفاده می کند که ورودی آن ها ترکیب Up-sample شده ی خروجی لایه ی قبل و لایه متناظر در Encoder می باشد. شبکه Decoder در ادامه بر روی دیتاست NYUv2 آموزش خواهد دید.

حال می خواهیم کمی بیشتر وارد جزئیات این معماری شده و لایه های مختلف آن را بررسی کنیم. جدول ۱ ساختار لایه به لایه ی شبکه را نشان می دهد. بخش Encoder مبتنی بر شبکه DenseNet-169 است که بر روی ImageNet آموزش دیده است. برای بخش Decoder، با یک لایه ی کانولوشنی 1×1 با همان تعداد کانال های خروجی Encoder شروع می کنیم. سپس به طور متوالی از بلاک های Up-sampling استفاده می کنیم. هر بلاک از یک لایه ی Bilinear up-sampling و در ادامه دو لایه ی کانولوشنی تشکیل شده است. در هر بلاک، ورودی اولین لایه ی کانولوشنی، حاصل ادغام خروجی لایه ی قبلی و خروجی لایه ی Pooling متناظر Encoder (Skip-connection) است. خروجی هر بلاک Up-sampling، به جز بلاک آخر، از تابع فعالساز

Leaky ReLU با پارامتر $\alpha = 0.2$ عبور می‌کند. ماتریس تصاویر ورودی تقسیم بر ۲۵۵ می‌شود تا محدوده ورودی بین ۰ و ۱ باشد و هیچ روش Normalization دیگری بر روی لایه‌ی ورودی استفاده نمی‌شود.

LAYER	OUTPUT	FUNCTION
INPUT	$480 \times 640 \times 3$	
CONV1	$240 \times 320 \times 64$	DenseNet CONV1
POOL1	$120 \times 160 \times 64$	DenseNet POOL1
POOL2	$60 \times 80 \times 128$	DenseNet POOL2
POOL3	$30 \times 40 \times 256$	DenseNet POOL3
...
CONV2	$15 \times 20 \times 1664$	Convolution 1×1 of DenseNet BLOCK4
UP1	$30 \times 40 \times 1664$	Upsample 2×2
CONCAT1	$30 \times 40 \times 1920$	Concatenate POOL3
UP1-CONVA	$30 \times 40 \times 832$	Convolution 3×3
UP1-CONVB	$30 \times 40 \times 832$	Convolution 3×3
UP2	$60 \times 80 \times 832$	Upsample 2×2
CONCAT2	$60 \times 80 \times 960$	Concatenate POOL2
UP2-CONVA	$60 \times 80 \times 416$	Convolution 3×3
UP2-CONVB	$60 \times 80 \times 416$	Convolution 3×3
UP3	$120 \times 160 \times 416$	Upsample 2×2
CONCAT3	$120 \times 160 \times 480$	Concatenate POOL1
UP3-CONVA	$120 \times 160 \times 208$	Convolution 3×3
UP3-CONVB	$120 \times 160 \times 208$	Convolution 3×3
UP4	$240 \times 320 \times 208$	Upsample 2×2
CONCAT3	$240 \times 320 \times 272$	Concatenate CONV1
UP2-CONVA	$240 \times 320 \times 104$	Convolution 3×3
UP2-CONVB	$240 \times 320 \times 104$	Convolution 3×3
CONV3	$240 \times 320 \times 1$	Convolution 3×3

جدول ۱ - سطر های نارنجی رنگ مربوط به Encoder (DenseNet-169) و سطر های آبی رنگ مربوط به Decoder می‌باشند.

تابع Loss: یک تابع هدف استاندارد که در مسائل تشخیص عمق استفاده می‌شود، اختلاف بین نقشه‌ی عمق اصلی و نقشه‌ی عمق خروجی شبکه می‌باشد. استفاده از تابع هدف مناسب نقش بسزایی در سرعت و عملکرد شبکه خواهد داشت. در این مدل، ما یک تابع هدف جدید تعریف می‌کنیم که علاوه بر محاسبه‌ی اختلاف نقشه‌ی اصلی و نقشه‌ی خروجی شبکه، وجود ناهنجاری‌های فرکانس بالا در نقشه‌ی خروجی شبکه را نیز مجازات می‌کند. این تابع هدف شامل جمع وزن دار سه جمله است:

$$L(y, \hat{y}) = \lambda L_{depth}(y, \hat{y}) + L_{grad}(y, \hat{y}) + L_{SSIM}(y, \hat{y}).$$

جمله‌ی اول همان L1 loss است که روی مقادیر عمق تعریف می‌شود:

$$L_{depth}(y, \hat{y}) = \frac{1}{n} \sum_p^n |y_p - \hat{y}_p|.$$

جمله‌ی دوم L1 loss است که روی گرادیان نقشه‌ی عمق تعریف می‌شود:

$$L_{grad}(y, \hat{y}) = \frac{1}{n} \sum_p^n |g_x(y_p, \hat{y}_p)| + |g_y(y_p, \hat{y}_p)|$$

جملات g_y و g_x به ترتیب اختلاف مولفه‌های x و y گرادیان نقشه‌ی عمق اصلی و خروجی شبکه را نشان می‌دهند.

در نهایت، جمله‌ی سوم از SSIM (Structural Similarity) استفاده می‌کند. اخیراً نشان داده شده است که وجود این جمله در تابع هدف می‌تواند به بهبود کیفیت آموزش شبکه‌های تشخیص عمق کمک کند.

$$L_{SSIM}(y, \hat{y}) = \frac{1 - SSIM(y, \hat{y})}{2}.$$

Data Augmentation: همانطور که می‌دانیم Data Augmentation یکی از روش‌های مناسب برای جلوگیری از Over-fitting در فرآیند آموزش شبکه می‌باشد. در این مدل، دو نوع Data Augmentation متفاوت مورد استفاده قرار گرفته است:

۱- Horizontal Flipping با احتمال ۰.۵

۲- جابجا کردن کانال‌های رنگ تصویر ورودی (برای مثال جابجا کردن کانال قرمز و سبز) با احتمال ۰.۲۵. ثابت شده است که این روش منجر به بهبود عملکرد شبکه می‌شود.

دیتاست مورد استفاده: همانطور که پیش‌تر اشاره شد و مطابق دستورالعمل‌های دستیاران آموزشی درس، برای آموزش شبکه از دیتاست NYUv2 با ۱۴۴۹ تصویر استفاده کردم. از این عدد، ۱۳۰۴ تصویر به آموزش و ۱۴۵ تصویر به تست اختصاص داده شده است. با توجه به معماری شبکه، نقشه‌ی عمق خروجی دارای ابعاد 320×240 می‌باشد. به هنگام آموزش، نقشه‌ی عمق واقعی را با ضریب ۲ کوچک می‌کنیم تا با ابعاد خروجی شبکه همخوانی داشته باشد. اما به هنگام تست، خروجی شبکه را با ضریب ۲ بزرگ می‌کنیم تا با نقشه‌ی عمق واقعی همخوانی داشته باشد. لازم به ذکر است که به هنگام تست، خروجی نهایی برابر است با میانگین خروجی شبکه برای تصویر اصلی و خروجی شبکه برای آینه شده‌ی تصویر اصلی.

جزئیات پیاده‌سازی: برای پیاده‌سازی شبکه‌ی DenseDepth از کد موجود در گیت‌هاب^۵ که با کتابخانه‌ی TensorFlow نوشته شده است استفاده کردیم. البته از کد ذکر شده تنها به عنوان شروع و پایه‌ی کار استفاده شده و تغییرات فراوان و متعددی در آن اعمال شده تا برای استفاده در این پروژه مناسب گردد. برای آموزش شبکه از سرویس آنلاین Google Colab با کارت گرافیک Tesla T4 و ۱۶ گیگابایت VRAM استفاده شده است. همانطور که پیش‌تر اشاره شد، بخش Encoder یک شبکه

⁵ [GitHub DenseDepth](#)

DenseNet-169 است که از قبل روی دیتاست ImageNet آموزش دیده است. وزن های بخش Decoder به صورت رندوم مقدار دهی اولیه شده اند. برای آموزش شبکه از بهینه ساز Adam و $\text{learning rate} = 0.0001$ و مقادیر پارامتر $\beta_1 = 0.9$ و $\beta_2 = 0.999$ و سایز بچ برابر با ۴ استفاده شده است. فرآیند آموزش برای تعداد epoch برابر با ۲۰ انجام شد و حدود ۳۰ دقیقه به طول انجامید.

ارزیابی: برای ارزیابی دقت شبکه روی داده های تست، از ۶ معیار معتبر برای ارزیابی دقت شبکه ها استفاده شده است. این معیار ها را در شکل ۳ مشاهده می فرمایید:

$$\begin{aligned} \text{average relative error (rel): } & \frac{1}{n} \sum_p^n \frac{|y_p - \hat{y}_p|}{y_p}; \\ \text{root mean squared error (rms): } & \sqrt{\frac{1}{n} \sum_p^n (y_p - \hat{y}_p)^2}; \\ \text{average (log}_{10}\text{) error: } & \frac{1}{n} \sum_p^n |\log_{10}(y_p) - \log_{10}(\hat{y}_p)|; \\ \text{threshold accuracy } (\delta_i): & \% \text{ of } y_p \text{ s.t. } \max\left(\frac{y_p}{\hat{y}_p}, \frac{\hat{y}_p}{y_p}\right) = \\ & \delta < thr \text{ for } thr = 1.25, 1.25^2, 1.25^3; \end{aligned}$$

شکل ۳ - معیار های مورد استفاده برای سنجش دقت شبکه

همانطور که پیش تر اشاره شد، از یک زیر مجموعه ی ۱۴۵ تایی از دیتاست برای انجام تست استفاده شده است. نتایج در جدول ۲ آورده شده است.

δ_1	δ_2	δ_3	rel	rms	\log_{10}
0.7960	0.9606	0.9931	0.1421	0.6948	0.0616

جدول ۲ - معیار های سنجش عملکرد شبکه

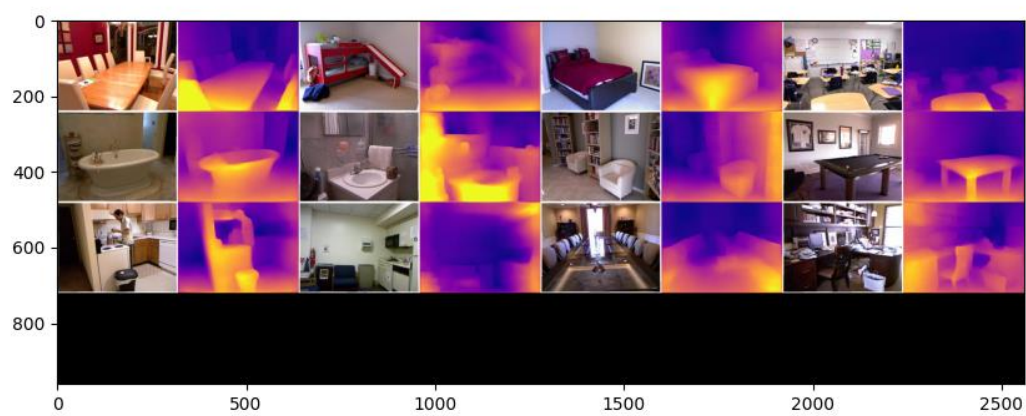
علیرغم اینکه نتایج سنجش عملکرد این شبکه که توسط اینجانب آموزش دیده است، نشان دهنده ی عملکرد بسیار خوب شبکه هستند، اما کمی نسبت به نتایج آورده شده در خود مقاله اصلی این مدل ضعیف تر هستند. این اتفاق می تواند ۲ دلیل داشته باشد:

۱- در مقاله اصلی از دیتاست NYUv2 با 50k تصویر استفاده شده است، حال آنکه ما برای آموزش شبکه تنها از ۱۳۰۴ عدد تصویر از این دیتاست استفاده کردیم. بنابراین منطقی است که عملکرد شبکه ی آموزش دیده توسط اینجانب کمی ضعیف تر باشد.

۲- برخی از روش های استفاده شده در مقاله اصلی برای Data Augmentation را به دلیل پیچیدگی زیاد حذف کردیم.

گراف کامل شبکه در فایل model_plot.png قابل مشاهده می باشد.

شکل ۴ خروجی شبکه را برای برخی از تصاویر تست نشان می دهد.



شکل ۴ - خروجی شبکه‌ی تشخیص عمق برای برخی از داده های تست

شبکه تشخیص اشیاء:

برای بخش تشخیص اشیاء مطابق پیشنهاد دستیاران آموزشی تصمیم گرفتیم از شبکه های از پیش آموزش دیده استفاده کنیم. برای مساله ی تشخیص اشیاء، کتابخانه های متعددی برای زبان پایتون نوشته شده اند که مدل های آماده و از پیش آموزش دیده ای با استفاده از انواع معماری های مشهور را در اختیار کاربر قرار می دهند تا بتواند بنا بر نیاز خود، یکی از آنها را انتخاب کرده و عملیات تشخیص اشیاء را انجام دهد. بنابراین جستجو های بسیار فراوانی انجام دادیم تا مدل مناسب برای تشخیص اشیاء انتخاب شود. بر این اساس، در یک سایت معتبر^۶، انواع الگوریتم ها و کتابخانه های مشهور برای تشخیص اشیاء ذکر شده و مزایا و معایب هر کدام شرح داده شده است. بر اساس کتابخانه های ذکر شده در این سایت، بر روی چندین کتابخانه ی که در زمینه ی تشخیص اشیاء وجود دارد تست های مختلفی انجام دادیم تا کتابخانه مناسب برای استفاده در پروژه را انتخاب کنیم. چند نمونه از این کتابخانه ها در ذیل آورده شده است:

۱- **ImageAI⁷**: این کتابخانه که با استفاده از TensorFlow نوشته شده است، مجموعه ی بزرگی از مدل ها و معماری های از پیش آموزش دیده را در اختیار کاربر می گذارد تا بتواند بر اساس نیاز های خود یکی را انتخاب کرده و عملیات تشخیص اشیاء را انجام دهد. ما این کتابخانه را نصب کردیم و تست های متعددی روی آن انجام دادیم. بر اساس تست هایی که انجام دادیم، مدل های موجود در این کتابخانه، از جمله مدل های YOLO، SSD و RetinaNet از دقت نسبتاً خوبی برخوردار بودند اما به دلیل استفاده ی این کتابخانه از پکیج های بسیار قدیمی که در ورژن های جدید پایتون دیگر در دسترس نیستند، این کتابخانه را انتخاب نکردیم.

۲- **GluonCV⁸**: این کتابخانه که با استفاده از MXNet نوشته شده است، یکی دیگر از کتابخانه های معتبر در حوزه ی تشخیص اشیاء می باشد. این کتابخانه، مدل های از پیش آموزش دیده ی متعددی را برای Image Classification، Object Detection و Semantic Segmentation در اختیار کاربر قرار می دهد. به طور خاص، این کتابخانه برای انجام Object Detection مدل های معروف از پیش آموزش دیده همچون SSD، Faster R-CNN و Yolo-v3 را در اختیار کاربر قرار می دهد. ما با استفاده از تک تک این مدل ها بر روی دیتاست NYUv2 تست های متعددی انجام دادیم. به دلیل آموزش دیدن مدل های این کتابخانه روی دیتاست های نه چندان مناسب، تشخیص اشیاء با استفاده از مدل های این کتابخانه بر روی دیتاست NYUv2 عملکرد بسیار ضعیف و دقت بسیار کمی برخوردار بود. به دلیل عملکرد غیر قابل قبول این کتابخانه بر روی دیتاست NYUv2، این مدل را نیز انتخاب نکردیم.

۳- **YOLOv3_TensorFlow⁹**: این کتابخانه با استفاده از مدل YOLOv3 عملیات تشخیص اشیاء را انجام می دهد و از دقت نسبتاً خوبی برخوردار است، اما به دلیل استفاده ی از ورژن های بسیار قدیمی و منسوخ شده ی TensorFlow، این کتابخانه را نیز انتخاب نکردیم.

۴- **Detectron2¹⁰**: این کتابخانه که توسط تیم Facebook AI Research (FAIR) توسعه داده شده است، یکی از بهترین و قوی ترین کتابخانه های موجود برای انجام عملیات تشخیص اشیاء می باشد که مجموعه ی بسیار عظیمی از مدل

⁶ [Object Detection Algorithms and Libraries](#)

⁷ [GitHub ImageAI](#)

⁸ [GitHub GluonCV](#)

⁹ [GitHub YOLOv3 TensorFlow](#)

¹⁰ [GitHub Detectron2](#)

های از پیش آموزش دیده را برای انواع عملیات Object Detection و Semantic Segmentation در اختیار کاربر قرار می‌دهد. کار کردن با این کتابخانه که با استفاده از PyTorch نوشته شده است، به دلیل منظم بودن و Documentation های بسیار خوب، برای کاربران ساده است. مدل های این کتابخانه به دلیل استفاده از دیتاست های معتبر و قوی برای آموزش مدل ها، از جمله دیتاست های COCO و ImageNet از دقت بسیار بالایی برخوردار است. بر اساس تست هایی که با استفاده از این کتابخانه روی دیتاست NYUv2 انجام دادیم، مدل های آموزش دیده این کتابخانه قادر بودند تا با دقت بسیار بالایی عملیات Object Detection را انجام دهند. به دلیل قوی بودن تیم توسعه ای این کتابخانه، سادگی کار با آن، و دقت بسیار بالای آن روی دیتاست NYUv2، در نهایت این کتابخانه برای استفاده در پروژه انتخاب شد.

استفاده از کتابخانه ی Detectron2: کار با کتابخانه ی Detectron2 بسیار ساده می‌باشد. کفایت مطابق دستورالعمل های ذکر شده در صفحه ی گیت هاب این کتابخانه، ابتدا پکیج های مورد نیاز نصب گردند و سپس خود کتابخانه نصب شود. در ادامه، با توجه به ¹¹Model Zoo ارائه شده توسط توسعه دهندگان، کاربر مدل مورد نظر خود را از بین چندین مدل از پیش آموزش دیده انتخاب می‌کند، آن را دانلود کرده و در کتابخانه ی Detectron2 وارد می‌کند و عملیات تشخیص اشیاء را به سادگی انجام می‌دهد.

مدل مورد استفاده: تا کنون در مورد کتابخانه ی مورد استفاده صحبت کردیم. هر کتابخانه، شامل مدل های از پیش آموزش دیده ی متعددی می‌شود. بنابراین، حال نوبت به آن رسیده که مدل مناسب برای تشخیص اشیاء در این پروژه را مشخص کنیم. کتابخانه ی Detectron2 مدل های از پیش آموزش دیده ی متعددی را برای Object Detection و Instance Segmentation در اختیار کاربران قرار می‌دهد که از جمله معروف ترین آن ها می‌توان به Faster R-CNN و Mask R-CNN اشاره کرد. به دلیل برخی ملاحظات که در بخش های بعدی به آن اشاره خواهد شد، تصمیم گرفتیم تا از بین مدل های موجود در Model Zoo کتابخانه ی Detectron2، از مدل Mask R-CNN استفاده کنیم که عملیات Instance Segmentation انجام می‌دهد، یعنی هم عملیات Object Detection (Box Prediction) و هم عملیات Semantic Segmentation (Mask Prediction) را به صورت توانمند انجام می‌دهد. این مدل بر روی دیتاست COCO آموزش دیده است و با توجه به تست های انجام شده، برای استفاده در این پروژه بسیار مناسب می‌باشد. حال می‌رسیم به قدم بعدی! کتابخانه ی Detectron2 برای مدل Backbone Mask R-CNN های مختلفی را پیشنهاد داده است که در ادامه به آنها اشاره می‌شود:

- ۱- FPN: از یک زیرساخت (FPN) ResNet+Feature Pyramid Network (FPN) به همراه لایه های conv و FC استاندارد برای تشخیص ماسک و تشخیص اشیاء استفاده می‌شود. این زیر ساخت بهترین Tradeoff را بین سرعت و دقت برقرار می‌کند بنابراین برای این پروژه این زیرساخت را انتخاب کردیم.
- ۲- C4: از زیرساخت ResNet conv4 به همراه لایه conv5 استفاده می‌کند.
- ۳- DC5: از زیرساخت ResNet conv5 به همراه لایه های conv و FC برای تشخیص ماسک و اشیاء استفاده می‌کند.

¹¹ [Detectron2 Model Zoo and Baselines](#)

در شکل ۵ که از Model Zoo کتابخانه برداشته شده است، جدول انواع مدل های از پیش آموزش دیده مبتنی بر Mask R-CNN را مشاهده می فرمایید:

COCO Instance Segmentation Baselines with Mask R-CNN

Name	lr sched	train time (s/iter)	inference time (s/im)	train mem (GB)	box AP	mask AP	model id	download
R50-C4	1x	0.584	0.110	5.2	36.8	32.2	137259246	model metrics
R50-DC5	1x	0.471	0.076	6.5	38.3	34.2	137260150	model metrics
R50-FPN	1x	0.261	0.043	3.4	38.6	35.2	137260431	model metrics
R50-C4	3x	0.575	0.111	5.2	39.8	34.4	137849525	model metrics
R50-DC5	3x	0.470	0.076	6.5	40.0	35.9	137849551	model metrics
R50-FPN	3x	0.261	0.043	3.4	41.0	37.2	137849600	model metrics
R101-C4	3x	0.652	0.145	6.3	42.6	36.7	138363239	model metrics
R101-DC5	3x	0.545	0.092	7.6	41.9	37.3	138363294	model metrics
R101-FPN	3x	0.340	0.056	4.6	42.9	38.6	138205316	model metrics
X101-FPN	3x	0.690	0.103	7.2	44.3	39.5	139653917	model metrics

شکل ۵ - مدل های از پیش آموزش دیده مبتنی بر Mask R-CNN برای عملیات Object Detection and Instance Segmentation

در ادامه، بر اساس مقاله^{۱۲} اصلی مدل Mask R-CNN، معماری و جزئیات این مدل را بیشتر بررسی خواهیم کرد.

مقدمه: Mask R-CNN یک شبکه ی سریع، انعطاف پذیر و تعمیم پذیر می باشد که برای انجام Instance Segmentation طراحی شده است. ما در این گزارش قصد نداریم وارد جزئیات Mask R-CNN شویم و تنها مفاهیم کلی آن را بررسی می کنیم. این شبکه نه تنها به طور بهینه اشیاء موجود در تصاویر را تشخیص می دهد، بلکه به صورت همزمان برای هر شیء تشخیص داده شده یک Segmentation Mask با کیفیت بالا نیز تولید می کند. این شبکه به نوعی شبکه ی Faster R-CNN را با اضافه کردن یک شاخه ی جدید گسترش می دهد که از آن برای تولید Segmentation Mask به طور موازی با شاخه ی موجود برای تولید Bounding Box استفاده می شود. این شاخه ی جدید در واقع یک FCN است که به هر Region of Interest (RoI) اعمال می شود و به صورت پیکسل به پیکسل Segmentation Mask تولید می کند. در شکل ۶ نمونه خروجی شبکه ی Mask R-CNN را مشاهده می فرمایید.

Mask R-CNN از نظر مفهومی بسیار ساده است: شبکه ی Faster R-CNN برای هر شیء تشخیص داده شده دو خروجی دارد، یکی برچسب کلاس و یکی آفست های Bounding Box. کافیت به این ساختار یک شاخه ی سوم اضافه کنیم که Object Mask را در خروجی تولید می کند.

معماری شبکه: برای نشان دادن انعطاف پذیری Mask R-CNN، نویسندگان مقاله معماری های مختلفی را برای این مدل پیاده سازی کرده اند. معماری این شبکه به طور کلی دو بخش دارد: (۱) معماری بخش Backbone که یک شبکه کانولوشنی بوده و بردار

¹² [Mask R-CNN](#)

ویژگی ها را از تصاویر استخراج می‌کند. ۲) بخش Head که برای تولید مشخصات Bounding Box و Mask Prediction از آن استفاده می‌شود. در ادامه، در مورد هر کدام از این بخش ها به طور جداگانه صحبت خواهد شد و معماری های مورد استفاده برای هر بخش ذکر خواهد شد.



شکل ۶- نمونه خروجی شبکه‌ی Mask R-CNN که بر روی دیتاست COCO آموزش دیده است.

برای بخش Backbone، از معماری های ResNet و ResNeXt با ۵۰ یا ۱۰۱ لایه استفاده شده است. پیاده سازی اصلی Faster R-CNN با استفاده از ResNet، ویژگی ها را از آخرین لایه‌ی کانولوشنی مرحله‌ی چهارم استخراج می‌کرد که به آن C4 می‌گوییم. بنابراین، استفاده از زیرساخت ResNet-50. استخراج ویژگی ها از C4 را ResNet-50-C4 نام‌گذاری می‌کنیم. نویسندگان مقاله همچنین از یک معماری مشهور دیگر نیز برای Backbone استفاده کرده اند که به آن Feature Pyramid Network (FPN) گفته می‌شود. شبکه‌ی Faster R-CNN با استفاده از FPN، ویژگی های RoI ها را از لول های مختلفی از Feature Pyramid با توجه به ابعاد آنها استخراج می‌کند، اما بقیه‌ی ساختار دقیقاً همانند ResNet است. استفاده از زیرساخت ResNet-FPN در شبکه‌ی Mask R-CNN منجر به بهترین نتایج چه از نظر سرعت و چه از نظر دقت شد. به همین دلیل بود که ما تصمیم گرفتیم تا از این زیرساخت در کتابخانه‌ی Detectron2 استفاده کنیم.

ارزیابی: بر اساس جدول های موجود در Model Zoo کتابخانه‌ی Detectron2، عملکرد شبکه‌ی Mask R-CNN با استفاده از زیرساخت R50-FPN و روی دیتاست COCO در جدول ۳ آورده شده است. لازم به ذکر است که عملکرد شبکه با استفاده از معیارهای AP جداگانه برای Box Prediction و Mask Prediction سنجیده شده است.

box AP	mask AP
41.0	37.2

جدول ۳- عملکرد شبکه‌ی Mask R-CNN با استفاده از زیرساخت R50-FPN

اتصال دو شبکه به یکدیگر:

برای اتصال دو شبکه به یکدیگر، حالت های مختلفی را بررسی کردیم تا بهینه ترین حالت را پیدا کنیم. بر این اساس، می توان برای اتصال دو شبکه، دو استراتژی مختلف را دنبال نمود:

۱- ابتدا شبکه ی تشخیص اشیاء، تمامی اشیاء موجود در تصویر را تشخیص می دهد، و سپس با استفاده از نقشه ی عمق تولید شده توسط شبکه ی تشخیص عمق، عدد عمق هر شیء به دست می آید.

۲- ابتدا شبکه ی تشخیص عمق، نقشه ی عمق تصاویر را تولید می کند، سپس بر اساس محدودیت هایی که کاربر روی میزان عمق قابل تشخیص اعمال نموده، بخش هایی از تصویر اصلی را Mask می کنیم و آن را به شبکه ی تشخیص اشیاء می دهیم. سپس نتایج و Bounding Box های خروجی این شبکه بر روی تصویر اصلی اعمال خواهد شد.

پس از بررسی های انجام شده، تصمیم گرفتیم تا در این پروژه از استراتژی ۱ برای ترکیب دو شبکه استفاده کنیم زیرا به نظر می رسد این استراتژی منجر به نتایج بهتری می شود. همانطور که پیش تر اشاره شد، در این استراتژی ابتدا شبکه ی تشخیص اشیاء، تمامی اشیاء موجود در تصویر را تشخیص می دهد و سپس، عمق هر تصویر با توجه به نقشه ی عمق تولید شده توسط شبکه ی تشخیص عمق قابل محاسبه است. در ادامه، بر اساس محدودیت هایی که کاربر بر روی عمق قابل تشخیص اعمال کرده، آن Bounding Box هایی که شرایط مورد نظر کاربر را ندارند حذف می شوند و Bounding Box های باقیمانده به همراه عمق بر روی تصویر اصلی رسم می شوند و خروجی نهایی تولید می شود. اما چگونه عمق هر شیء را از روی نقشه ی عمق محاسبه کنیم؟ اگر به خاطر داشته باشید، پیش تر اشاره شد که برای شبکه ی تشخیص اشیاء از Instance Segmentation استفاده کرده ایم که هم Bounding Box Prediction و هم Mask Prediction انجام می دهد. در اینجا دلیل اینکار مشخص می شود! Bounding Box ها به تنهایی اطلاعات کافی برای محاسبه ی عمق اشیاء را در اختیار ما نمی گذارند زیرا در هر Bounding Box، تنها درصدی از پیکسل ها مربوط به شیء تشخیص داده شده است و درصدی دیگر مربوط به سایر اشیاء موجود در پس زمینه هستند، بنابراین اطلاعات دقیقی از مختصات پیکسلی هر شیء در اختیار نداریم و تخمین عمق شیء از روی نقشه ی عمق بسیار دشوار و بی کیفیت می شود. اما با استفاده از Instance Segmentation، علاوه بر Bounding Box ها، Segmentation Mask ها را نیز در اختیار داریم که مختصات پیکسلی هر شیء تشخیص داده شده را به طور دقیق مشخص می کنند. به عبارتی، Segmentation Mask به ما اجازه می دهد که به طور دقیق مرز های هر شیء تشخیص داده شده را مشخص کنیم. در واقع، با استفاده از Mask R-CNN، همانطور که در شکل ۶ مشاهده می شود، علاوه بر مختصات Bounding Box ها، برای هر شیء تشخیص داده شده یک آرایه ی True و False نیز در خروجی داریم که دارای ابعاد تصویر اصلی است و مکان پیکسل های متعلق به آن شیء را در اختیار ما قرار می دهد. در شکل ۷ نمونه یک تصویر و Segmentation Mask تولید شده برای یکی از اشیاء را مشاهده می فرمایید.

با داشتن پیکسل های هر شیء، می توانیم به سادگی پیکسل های متناظر را در نقشه ی تشخیص عمق بررسی کرده و با توجه به آنها عمق شیء را تشخیص دهیم. برای تخمین عمق اشیاء از روش میانگین گیری استفاده شده است، بدین معنا که عمق هر شیء برابر است با میانگین عمق تک تک پیکسل های آن شیء روی نقشه ی تشخیص عمق:

$$D_I = \frac{1}{n} \sum_{p \in I} D_p$$

در این رابطه، D_I عمق شیء مورد نظر است، n تعداد کل پیکسل های آن شیء است که توسط Segmentation Mask مشخص می شود، و $\sum_{p \in I} D_p$ عمق تمامی پیکسل های متعلق به آن شیء را جمع می کند.



شکل ۷ - نمونه یک تصویر و Segmentation Mask تولید شده برای اسب موجود در تصویر

ارزیابی: حال می خواهیم بر اساس معیار های عملکرد هر کدام از شبکه های جداگانه، عملکرد شبکه ی نهایی را تخمین بزنیم.

قبل از ورود به جزئیات بیشتر باید یک نکته ی مهم را یادآوری کنیم: شبکه ی تشخیص اشیاء، دو معیار مختلف را برای عملکرد شبکه در اختیار ما قرار می دهد، یکی AP_{box} که دقت تولید Bounding Box را نشان می دهد، و دیگری AP_{mask} که دقت تولید Segmentation Mask را نشان می دهد. با توجه به اینکه ما در معماری خود از Segmentation Mask های تولید شده توسط شبکه ی تشخیص اشیاء (و نه از Bounding Box های تولید شده) برای تخمین عمق استفاده می کنیم، بنابراین دقت AP_{box} تأثیری در عملکرد تشخیص عمق اشیاء ندارد. به عبارتی، ما در ابتدا با استفاده از معیار هایی که در ادامه معرفی می شوند و با ترکیب کردن AP_{mask} و معیار های عملکرد شبکه ی تشخیص عمق، دقت تخمین عمق اشیاء کلی را بدست می آوریم، و سپس بین عدد بدست آمده و عدد AP_{box} مینیمم می گیریم تا دقت شبکه ی نهایی بدست آید. به عبارتی، ما دقت شبکه ی نهایی را از مینیمم گرفتن بین دقت عملیات تشخیص عمق اشیاء و عملیات تشخیص اشیاء بدست می آوریم، اما دقت عملیات تشخیص عمق اشیاء خود حاصل ترکیب دقت شبکه ی تشخیص عمق و AP_{mask} شبکه ی تشخیص شیء است. بنابراین ۲ مرحله ی زیر انجام می شود:

۱- با ترکیب دقت شبکه ی تشخیص عمق و AP_{mask} شبکه ی تشخیص اشیاء، دقت عملیات تشخیص عمق اشیاء را بدست می آوریم.

۲- بین عدد بدست آمده در مرحله ی ۱ و عدد AP_{box} مینیمم می گیریم تا دقت نهایی شبکه ی کلی بدست آید.

برای انجام مرحله ی اول، از دو روش زیر برای ترکیب استفاده می کنیم:

$$1- \gamma_1 = \delta_1 \times AP_{mask} = 0.7960 \times 0.372 = 0.296$$

$$2- \gamma_2 = e^{-rms} \times AP_{mask} = 0.186$$

بر این اساس، دقت شبکه‌ی نهایی به صورت زیر محاسبه خواهد شد:

$P_1 = \min(\gamma_1, AP_{box})$	$P_2 = \min(\gamma_2, AP_{box})$
29.6	18.6

جدول ۴ - دقت شبکه‌ی نهایی

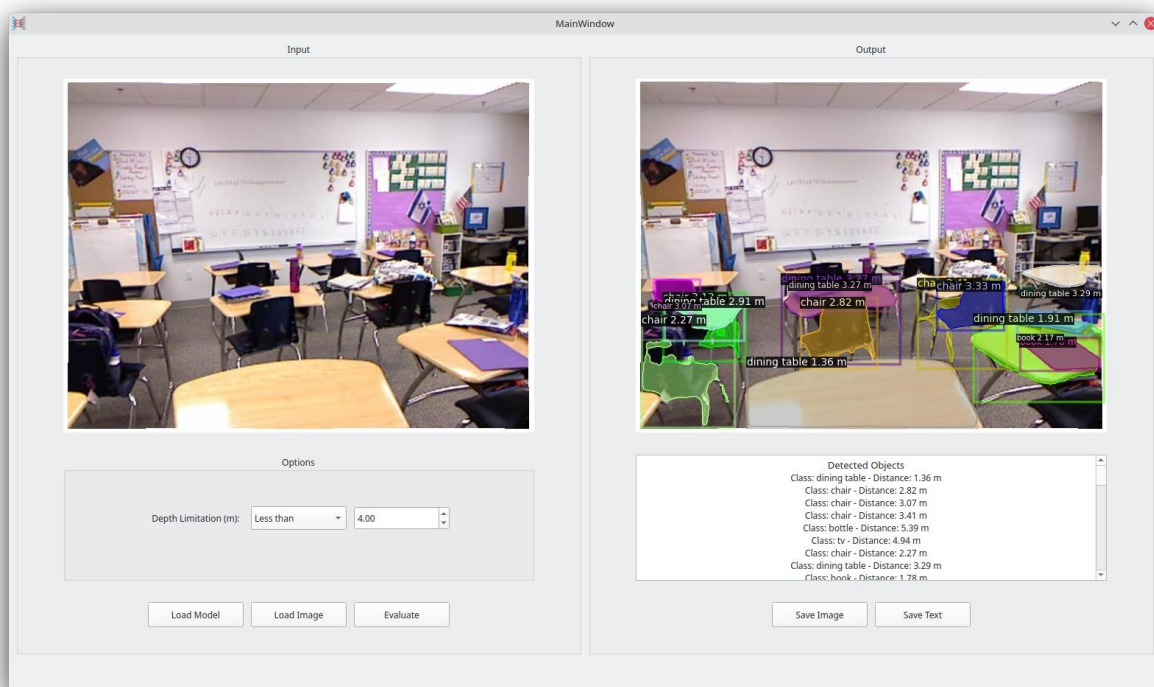
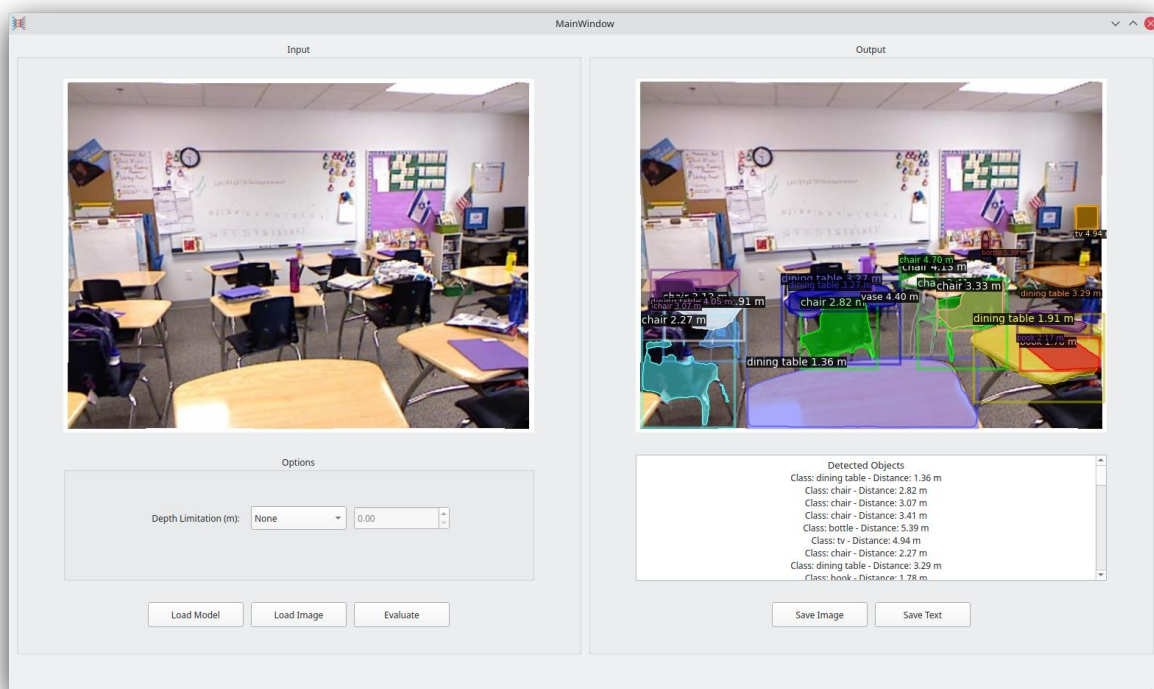
لازم به ذکر است که هر دو معیار معرفی شده برای سنجش دقت شبکه‌ی نهایی، کاملاً جدید بوده و توسط اینجانب معرفی شده اند و از هیچ مرجعی الهام گرفته نشده اند. نمونه ورودی و خروجی پروژه را در شکل ۸ مشاهده می‌کنید.

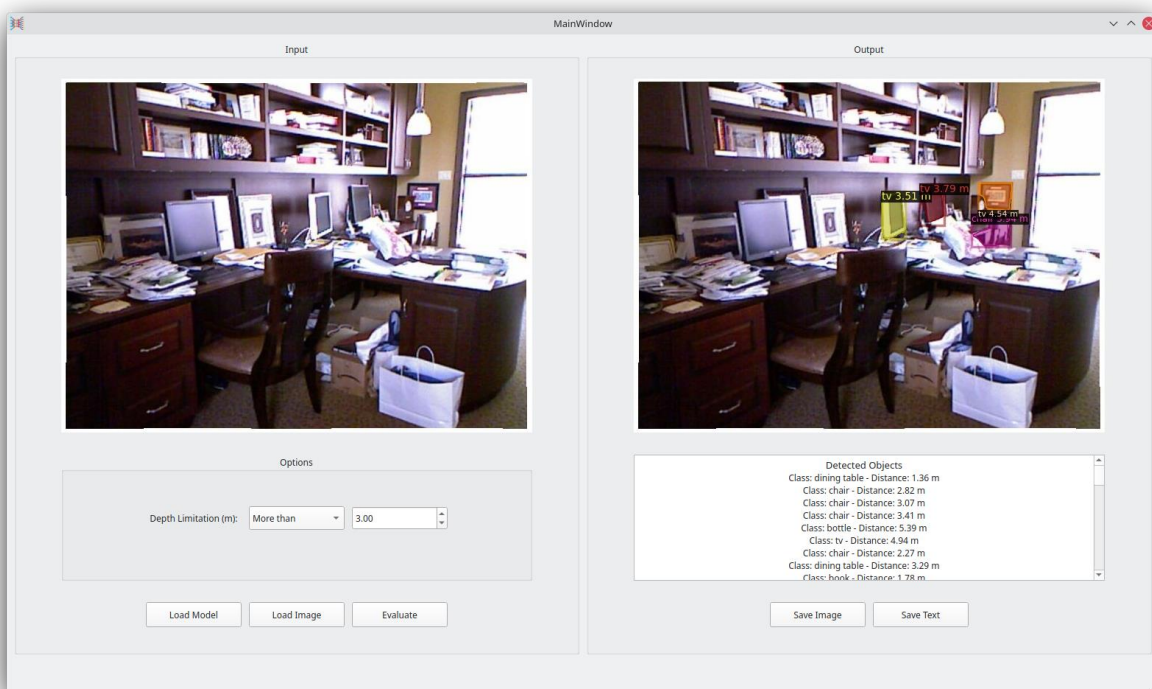
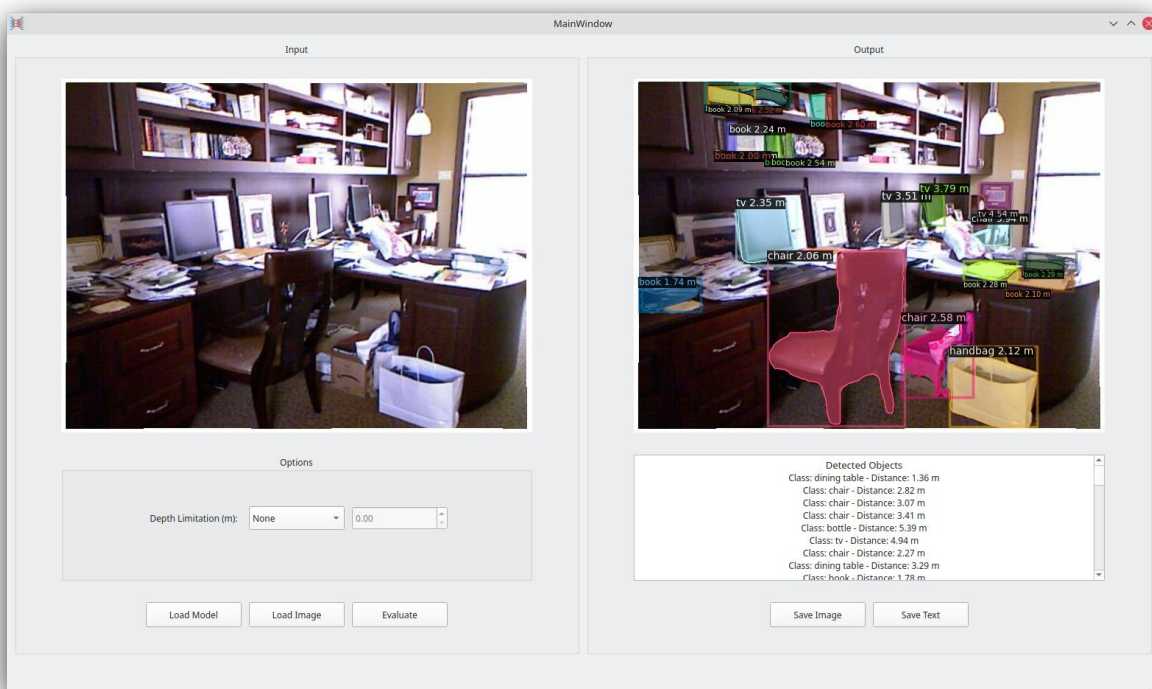
محیط گرافیکی:

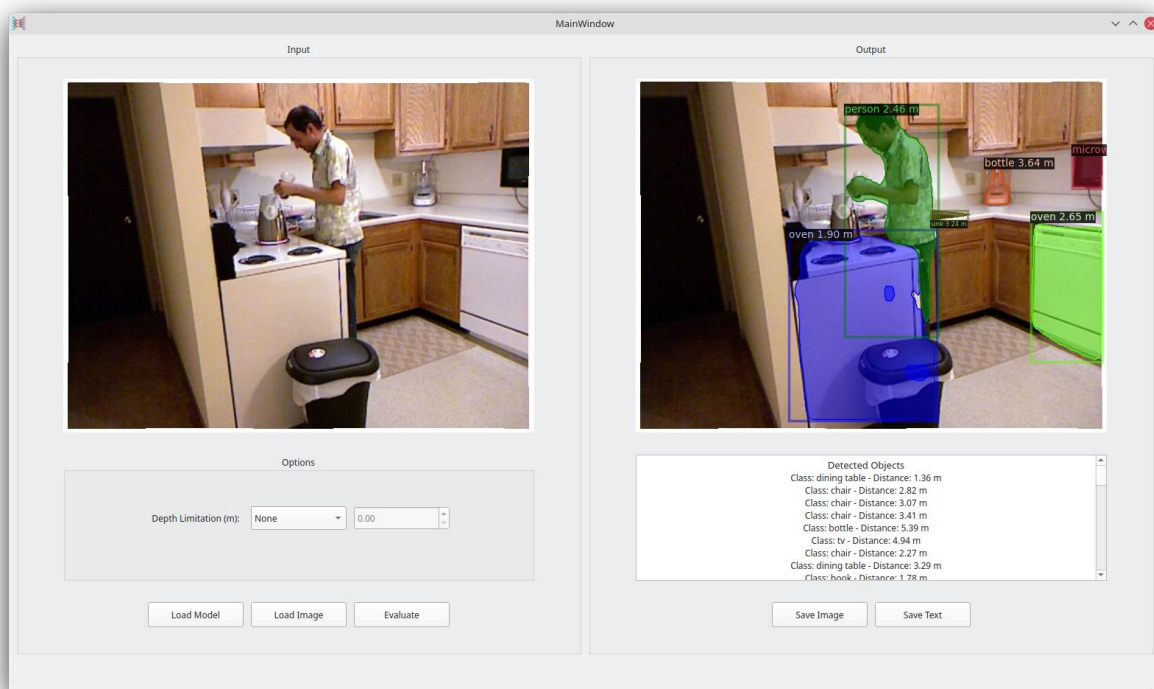
به منظور سهولت استفاده از این پروژه و همچنین جذابیت بصری آن، با استفاده از کتابخانه‌ی PyQt5 و ابزار Qt Designer، یک محیط گرافیکی تعاملی نیز برای آن ایجاد شده است. این محیط گرافیکی، مدل از پیش آموزش دیده و یک تصویر را از کاربر به عنوان ورودی دریافت می‌کند و با انجام عملیات Object Detection و Depth Estimation بر روی تصویر ورودی، تمامی اشیاء موجود و عمق هر شیء را مشخص کرده و در تصویر خروجی اعمال می‌کند. همچنین، تمامی اشیاء تشخیص داده شده و عمق آنها به صورت متنی جدای از تصویر نیز نمایش داده می‌شوند. آموزش قدم به قدم استفاده از محیط گرافیکی و سایر جزئیات در فایل README موجود در گیت هاب به طور کامل شرح داده شده است. نمونه ورودی و خروجی پروژه را در شکل ۸ مشاهده می‌کنید.

گیت هاب:

این گزارش به همراه تمامی فایل ها و کد های پروژه، نحوه‌ی نصب کتابخانه های مورد نیاز و سایر آموزش های مورد نیاز برای کار با آن در Repository گیت هاب موجود می‌باشد که جناب آقای بهراد احمدپور به آن اضافه شده اند. در صورتی که نیاز بود شخص دیگری نیز به Repository گیت هاب اضافه شود با ایمیل shokuhi.mh@gmail.com یا آیدی تلگرام @Mhosein_shokuhi در ارتباط باشید.







شکل ۸ - چند نمونه ورودی و خروجی پروژه