

	Models	Notes		Dataset(s) for benchmarking			Parameters (M)	Performance (Speed / FLOPS)	OpenVino Optimized Speed	Pros	Cons	Efficiency Remarks	Ideal use cases for the models	Reference resource/ notebooks	Link to research paper	Code Complexity	Code Template
Image Classification																	
					TOP 1 ACCURACY	TOP 5 ACCURACY											
	Vision Transformer			ImageNet	88.55%		632			(State of the art) High Accuracy with Less Computation Time for Training compared to Noisy Student	It does not use convolutional neural networks, instead vision transformer uses small patches to apply self attention; the efficacy of using small patches is debatable at the time of writing as it might miss edges	Accuracy is the most critical metric Semi supervised learning where getting labelled data is hard			<a href="https://arxiv.org/pdf/2010.11929v1.pdf">https://arxiv.org/pdf/2010.11929v1.pdf</a>	Hard	
	NoisyStudent (EfficientNet-L2)			ImageNet	88.40%	98.70%	480			Augments labeled dataset with unlabeled examples	Adding "noise" requires careful consideration and configuration ; computationally very expensive				<a href="https://arxiv.org/pdf/1911.04263v4.pdf">https://arxiv.org/pdf/1911.04263v4.pdf</a>	Hard	
	EfficientNet-B7			ImageNet	84.4		66			Uses Neural Architecture Search to optimize architecture					<a href="https://arxiv.org/abs/1905.11946">https://arxiv.org/abs/1905.11946</a>	Medium	
	EfficientNet-B4			ImageNet	82.6		19			Uses Neural Architecture Search to optimize architecture; excellent accuracy compared to speed			Both accuracy and performance are critical Both accuracy and performance are critical		<a href="https://arxiv.org/abs/1905.11946">https://arxiv.org/abs/1905.11946</a>	Medium	
	EfficientNet-B2			ImageNet	79.8		9.2			Uses Neural Architecture Search to optimize architecture			ML on the edge where accuracy is important!		<a href="https://arxiv.org/abs/1905.11946">https://arxiv.org/abs/1905.11946</a>	Medium	
	EfficientNet-B0			ImageNet	76.3		5.3			Uses Neural Architecture Search to optimize architecture; highly optimized for ML on edge with increased accuracy Uses filters of varying sizes to detect different image features for greater accuracy					<a href="https://arxiv.org/abs/1905.11946">https://arxiv.org/abs/1905.11946</a>		
	Inception V3			ImageNet	78.8	95.10%	55.8			Highly optimized for performance on ML on the edge applications; great backbone for CV on the edge tasks	Computationally expensive; better models exist				<a href="https://arxiv.org/pdf/1602.07281v2.pdf">https://arxiv.org/pdf/1602.07281v2.pdf</a>	Easy	
	MobileNet v2			ImageNet	71.88	90.29	3.4			Very fast and hence great for performance critical applications; works out of the box	Not as robust as larger models		ML on the edge ML on the edge, low latency		<a href="https://arxiv.org/pdf/2006.10100v1.pdf">https://arxiv.org/pdf/2006.10100v1.pdf</a>	Medium	
	SqueezeNet			ImageNet	58.1	80.42	1.25				Poor performance when used for classification and detection of small objects Increased complexity of architecture due to skip connections;	Highly efficient			<a href="https://arxiv.org/abs/1602.07360">https://arxiv.org/abs/1602.07360</a>	Medium	
	ResNet 18			ImageNet	69.76	90.08	11.174			Most tried and tested method. Works well.	Implementation of Batch normalization layers since ResNet heavily depends on it Increased complexity of architecture due to skip connections;				<a href="https://arxiv.org/pdf/1512.03385.pdf">https://arxiv.org/pdf/1512.03385.pdf</a>	Easy	
	ResNet 34			ImageNet	73.3	91.42	21.282			Most tried and tested method. Works well.	Implementation of Batch normalization layers since ResNet heavily depends on it Increased complexity of architecture due to skip connections;				<a href="https://arxiv.org/pdf/1512.03385.pdf">https://arxiv.org/pdf/1512.03385.pdf</a>	Easy	
	ResNet 50			ImageNet	76.15	92.87	25.6			Most tried and tested method. Works well.	Implementation of Batch normalization layers since ResNet heavily depends on it Increased complexity of architecture due to skip connections;				<a href="https://arxiv.org/pdf/1512.03385.pdf">https://arxiv.org/pdf/1512.03385.pdf</a>	Easy	
	ResNet 101			ImageNet	79.20%	94.70%	42.513			Most tried and tested method. Works well.	Implementation of Batch normalization layers since ResNet heavily depends on it Computationally expensive and gives much lower accuracy in comparison with much smaller models				<a href="https://arxiv.org/pdf/1512.03385.pdf">https://arxiv.org/pdf/1512.03385.pdf</a>	Easy	
	VOG 19			ImageNet	74.50%	92.00%	144			Easy to understand architecture	Computationally expensive and gives much lower accuracy in comparison with much smaller models	Low efficiency	Google deepdream		<a href="https://arxiv.org/pdf/1409.1556v5.pdf">https://arxiv.org/pdf/1409.1556v5.pdf</a>	Easy	
	VOG 16			ImageNet	74.40%	91.90%	138			Easy to understand architecture		Low efficiency	Google deepdream		<a href="https://arxiv.org/pdf/1409.1556v5.pdf">https://arxiv.org/pdf/1409.1556v5.pdf</a>	Easy	
Metrics Legend																	
Top1	The model is considered to have classified a given image correctly if the target label is the model's top prediction																
Top5	The model is considered to have classified a given image correctly if the target label is the model's top 5 predictions																
Semantic segmentation					MEAN IOU												
	FFN			Encoder (imageNet) Encoder (imageNet)			23.15M	Nil	Nil	i) Same network can perform Image Segmentation, Object detection and Pose Estimation Learns without any significant increase in number of parameters, fast.	i) Heavy network cannot be used for real time inference Not the most accurate, compromises accuracy for less params and faster fps.	High performan	<a href="https://github.com/NicolasBois/">https://github.com/NicolasBois/</a>	<a href="http://presentations.cocodataset.org/COCO17-Stuff-FAIR.pdf">http://presentations.cocodataset.org/COCO17-Stuff-FAIR.pdf</a>	Easy		
	Linknet			Cityscapes	84.00%		11.5M	21.2G	Nil	Exploits the impact of global contextual information in semantic segmentation by combining attention mechanism and spatial pyramid pooling to extract precise dense features.	i) Even though a much lighter network than using the Inception or VGG backbones, this is still heavy for real-time inference	Mld performance	<a href="https://github.com/NicolasBois/">https://github.com/NicolasBois/</a>	<a href="https://arxiv.org/abs/1707.03718">https://arxiv.org/abs/1707.03718</a>	Easy		
	PAN			Cityscapes	80.20%		21.4M	Nil	Nil	Exploits the capability of global context information by differentiation-based context aggregation through our pyramid pooling	i) Heavy model; needs a powerful system to run It is quite slow because the network must be run separately for each patch, and there is a lot of redundancy due to overlapping patches.	Where you need	<a href="https://github.com/NicolasBois/">https://github.com/NicolasBois/</a>	<a href="https://arxiv.org/abs/1612.01105">https://arxiv.org/abs/1612.01105</a>	Easy		
	PSPNet			Cityscapes	80.20%		21.4M	Nil	Nil	This network can localize the training data in terms of patches is much larger than the number of training images		When you have 1	<a href="https://github.com/NicolasBois/">https://github.com/NicolasBois/</a>	<a href="https://arxiv.org/abs/1505.04597">https://arxiv.org/abs/1505.04597</a>	Medium		
U-Net				Brain MRI segmentation			77.6M	Nil	Nil								
Metrics Legend																	
Mean IOU * IOU = true_positive / (true_positive + false_positive + false_negative) PCK@0.5: Probability of 50% pose points matching																	
Human Pose Estimation					PCK@0.5												
	MSPN			MPII Human Pose	92.60%			9.6G	Nil	i) Light model that runs fast on standard chips as well i) Maintain high res. representations throughout the the whole process ii) Fuse multi-res representations repeatedly	i) Model is heavier compared to others in the segment				<a href="https://arxiv.org/abs/1901.00148v4">https://arxiv.org/abs/1901.00148v4</a>	Medium	
	HRNet-W32			MPII Human Po	92.30%		28.5 M	16.0G	Nil						<a href="https://arxiv.org/abs/1902.09212v4">https://arxiv.org/abs/1902.09212v4</a>	Easy	
	Pyramid Residual Modules			MPII Human Pose	92.00%				Nil						<a href="https://arxiv.org/abs/1706.01101v4">https://arxiv.org/abs/1706.01101v4</a>	Hard	
	Multi-Context Attention			MPII Human Po	91.50%				Nil						<a href="https://arxiv.org/abs/1702.07432v1">https://arxiv.org/abs/1702.07432v1</a>	Hard	
	CU-Net			MPII Human Po	91.20%		(4 Stack) 3.9M (8) 7.9 M (16) 15.9 M		Nil	i) Introduces and uses dense u-net that uses ~70% lower parameters compared to traditional U-Net stack ii) Since this is a U-Net the input and output size if the same and there is no need to upscaling					<a href="https://arxiv.org/abs/1809.02194v2">https://arxiv.org/abs/1809.02194v2</a>	Medium	
	Stacked hourglass + Inception-ResNet		Pretrained weights not available. Code for running ready.	MPII Human Po	91.20%				Nil						<a href="https://arxiv.org/abs/1706.02607v2">https://arxiv.org/abs/1706.02607v2</a>	Hard	
	EfficientPose IV			MPII Human Po	91.20%		6.56 M		Nil	Very small orders of magnitude smaller than comparable models	i) Single person only ii) Relies on upscaling which can reduce the effectiveness.	Real Time Inference			<a href="https://arxiv.org/abs/2004.12186v1">https://arxiv.org/abs/2004.12186v1</a>	Easy	
	EfficientPose RT			MPII Human Po	88.40%		0.46 M		--	Order of mag. smaller than IV version	--	Real Time Inference			--	--	
Activity Recognition																	
Metrics Legend																	
3-Fold Accuracy: The data is split into three parts and 2 are used for training while the third is used for testing. This accuracy is the mean of three.																	
Action Recognition					3-FOLD ACCURACY												
	R2+1D BERT		Only following weights are available: <a href="https://arxiv.org/abs/2004.02457">https://arxiv.org/abs/2004.02457</a> <a href="https://arxiv.org/abs/2004.02457v2">https://arxiv.org/abs/2004.02457v2</a> <a href="https://arxiv.org/abs/2004.02457v3">https://arxiv.org/abs/2004.02457v3</a> <a href="https://arxiv.org/abs/2004.02457v4">https://arxiv.org/abs/2004.02457v4</a>	UCF101	98.69%		66.67 M	152.97 G	Nil	i) Uses BERT for Bi-directional embedding ii) BERT output is standard and thus can be utilised with other tasks such as image/video captioning, etc	i) Super heavy model				<a href="https://arxiv.org/abs/2009.01232v3">https://arxiv.org/abs/2009.01232v3</a>	Medium	

