

Edge Devices	Tensorflow Version	Model (CNN, NLP)	Dataset	Quantization	OS (Ubuntu/L4T/Docker)	CUDA version
Nvidia Jetson TX1, Nvidia Jetson TX2, Nvidia Jetson Nano, Nvidia Jetson Xavier, Raspberry pi 4 + Coral TPU	2.5	Mobilenet V1	ImageNet	Tensorflow FP32, EdgeTPU-tflite, TensorLite, TensorRT	18.04 / 32.6.1 / 20.10	10.2.300
		Mobilenet V2				
		Inception V3				
		YOLO V5	COCO	Tensorflow FP32		
		RNN	IMDb			
		LSTM				
		BERT				
DistilBERT	GLUE SST-2					

Image Classification

MobileNet V1

Edge device	Setup						Metric									
	OS (Ubuntu/L4T/ Docker)	Model	Dataset	Optimization Platform	Quantization	batch size	Accuracy	Model load time	Data load time	Inference task time	Inference time	IPS (model, dataset load + inference)	IPS (inference only)	Flop (giga)	parameter (m)	Flop / parameter (m)
Google Coral TPU	18.04 / 20.10	YOLO V5	mp4	edge tpu tfLite - TFv2	INT8	32										
Google Coral TPU	18.04 / 20.10	YOLO V5	mp4	edge tpu tfLite - TFv2	INT8	64										
Google Coral TPU	18.04 / 20.10	YOLO V5	mp4	edge tpu tfLite - TFv2	INT8	128										
Nvidia Jetson TX1	18.04 / 32.6.1 / 20.10	YOLO V5	mp4	tfLite - TFv2	FP16	1		5.404241323	0.0789109726	84.69850264	2.210868567	0.4213810761	0.4523053248	16.4	7.2	2277.777778
Nvidia Jetson TX2	18.04 / 32.6.1 / 20.10	YOLO V5	mp4	tfLite - TFv2	FP16	1		5.126217842	0.03652477264	76.50713611	1.998791111	0.4652878221	0.500302405	16.4	7.2	2277.777778
Nvidia Jetson Xavier	18.04 / 32.6.1 / 20.10	YOLO V5	mp4	tfLite - TFv2	FP16	1		3.832612276	0.03427410126	56.88416886	1.485840628	0.6255035382	0.6730196908	16.4	7.2	2277.777778
Nvidia Jetson Nano	18.04 / 32.6.1 / 20.10	YOLO V5	mp4	tfLite - TFv2	FP16	1		7.65151906	0.06461715698	96.5831039	2.523046274	0.3643363073	0.396346276	16.4	7.2	2277.777778
Nvidia Jetson TX1	18.04 / 32.6.1 / 20.10	YOLO V5	mp4	tfLite - TFv2	INT8	1		5.189818859	0.0424387455	47.69496274	1.239058526	0.7179670452	0.8070643792	16.4	7.2	2277.777778
Nvidia Jetson TX2	18.04 / 32.6.1 / 20.10	YOLO V5	mp4	tfLite - TFv2	INT8	1		5.158523798	0.03731636865	42.02915525	1.091719676	0.8046586446	0.9159860544	16.4	7.2	2277.777778
Nvidia Jetson Xavier	18.04 / 32.6.1 / 20.10	YOLO V5	mp4	tfLite - TFv2	INT8	1		3.864009142	0.03376603127	57.95437908	1.514623171	0.6143662537	0.6602302268	16.4	7.2	2277.777778
Nvidia Jetson Nano	18.04 / 32.6.1 / 20.10	YOLO V5	mp4	tfLite - TFv2	INT8	1		6.017657042	0.04754137993	56.25687218	1.462034784	0.6097358389	0.6639782549	16.4	7.2	2277.777778

Image Classification

MobileNet V1

Edge device	Setup		Metric										
	Model	Dataset	batch size	Accuracy	Model load time	Data load time	Inference task time	Inference time	IPS (model, dataset load + inference)	IPS (inference only)	Flop (giga)	parameter (m)	Flop / parameter (m)
Nvidia Jetson Xavier	MobileNet V1	imagenet 1000	1	0.93	38.88625741	1.037235498	65.82744575	0.06412833238	9.456155237	15.59373155	1.15	4.25	270.5882353
Nvidia Jetson Xavier	MobileNet V1	imagenet 1000	2	0.93	38.01113439	1.054350138	32.73214889	0.03176783514	13.92799128	31.47838043	2.3	4.25	541.1764706
Nvidia Jetson Xavier	MobileNet V1	imagenet 1000	4	0.93	39.47181606	1.041965723	18.81584477	0.01820332026	16.85487822	54.93503303	4.59	4.25	1080
Nvidia Jetson Xavier	MobileNet V1	imagenet 1000	8	0.93	38.30404139	1.030306578	12.16145039	0.01117571183	19.41897256	85.35546233	9.18	4.25	2160
Nvidia Jetson Xavier	MobileNet V1	imagenet 1000	16	0.93	38.36752224	1.058549881	8.744111538	0.008018319332	20.92574659	124.7144144	18.4	4.25	4329.411765
Nvidia Jetson Xavier	MobileNet V1	imagenet 1000	32	0.93	38.1481607	1.017022848	8.580517292	0.00614240556	21.44679822	162.8026659	36.7	4.25	8635.294118
Nvidia Jetson Xavier	MobileNet V1	imagenet 1000	64	0.93	38.18341804	1.039550543	8.765650034	0.005155852996	21.33827334	193.9543274	73.4	4.25	17270.58824
Nvidia Jetson Xavier	MobileNet V1	imagenet 1000	128	0.93	37.66842175	1.027097225	10.46026754	0.004670296796	20.83164779	214.1191543	1.47E+02	4.25	3.46E+04
Nvidia Jetson Nano	MobileNet V1	imagenet 1000	1	0.93	71.08252144	2.345131636	156.859957	0.1526636727	4.342385849	6.550346801	1.15	4.25	270.5882353
Nvidia Jetson Nano	MobileNet V1	imagenet 1000	2	0.93	71.12680101	2.0991745	112.5438108	0.1095119245	5.38299902	9.131425683	2.3	4.25	541.1764706
Nvidia Jetson Nano	MobileNet V1	imagenet 1000	4	0.93	70.64579511	2.058488846	85.60306263	0.07203072262	6.31681362	13.88296499	4.59	4.25	1080
Nvidia Jetson Nano	MobileNet V1	imagenet 1000	8	0.93	69.93868399	2.046374798	62.63308215	0.05643707871	7.42839512	17.71884766	9.18	4.25	2160
Nvidia Jetson Nano	MobileNet V1	imagenet 1000	16	0.93	70.81421661	2.031164169	68.55348921	0.05673198426	7.128752148	17.62674113	18.4	4.25	4329.411765
Nvidia Jetson Nano	MobileNet V1	imagenet 1000	32	0.93	70.34334397	2.068148136	107.1531734	0.08601080929	5.702673581	11.62644566	36.7	4.25	8635.294118
Nvidia Jetson Nano	MobileNet V1	imagenet 1000	64								73.4	4.25	17270.58824
Nvidia Jetson Nano	MobileNet V1	imagenet 1000	128								1.47E+02	4.25	3.46E+04

MobileNet V2

Edge device	Setup		Metric										
	Model	Dataset	batch size	Accuracy	Model load time	Data load time	Inference task time	Inference time	IPS (model, dataset load + inference)	IPS (inference only)	Flop (giga)	parameter (m)	Flop / parameter (m)
Nvidia Jetson Xavier	MobileNet V2	imagenet 1000	1	0.906	66.94901466	1.029478643	92.66682172	0.09103277969	6.224886781	10.98505399	0.615	3.53	174.2209632
Nvidia Jetson Xavier	MobileNet V2	imagenet 1000	2	0.906	66.94992638	1.034856796	48.19621539	0.04725036407	8.607235995	21.16385809	1.23	3.53	348.4419263
Nvidia Jetson Xavier	MobileNet V2	imagenet 1000	4	0.906	67.81008959	1.019435406	26.64441657	0.02604661059	10.47391131	38.39271127	2.46	3.53	696.8838527
Nvidia Jetson Xavier	MobileNet V2	imagenet 1000	8	0.906	66.61496758	1.004122734	15.34189391	0.01490839934	12.05383741	67.0762821	4.92	3.53	1393.767705
Nvidia Jetson Xavier	MobileNet V2	imagenet 1000	16	0.906	66.90273213	1.034251213	10.25050116	0.009799251008	12.89260094	102.0486157	9.84	3.53	2787.535411
Nvidia Jetson Xavier	MobileNet V2	imagenet 1000	32	0.906	66.81817079	1.010535479	8.577705383	0.007023798302	13.40196461	142.3731088	19.7	3.53	5580.736544
Nvidia Jetson Xavier	MobileNet V2	imagenet 1000	64	0.906	67.01150298	1.009224176	8.856770515	0.00599429966	13.31964638	166.82516	39.4	3.53	11161.47309
Nvidia Jetson Xavier	MobileNet V2	imagenet 1000	128	0.906	67.0898447	1.022275925	10.45342684	0.005248077912	13.03367568	190.5459516	78.7	3.53	22294.61756
Nvidia Jetson Nano	MobileNet V2	imagenet 1000	1	0.906	127.2817204	2.542519569	253.3297305	0.2492911539	2.609913198	4.011373787	0.615	3.53	174.2209632
Nvidia Jetson Nano	MobileNet V2	imagenet 1000	2	0.906	125.8207233	2.181444604	173.7648814	0.1711566045	3.31381096	5.842602468	1.23	3.53	348.4419263
Nvidia Jetson Nano	MobileNet V2	imagenet 1000	4	0.906	125.2173266	2.570510387	158.1518993	0.1498946307	3.497234534	6.67135304	2.46	3.53	696.8838527
Nvidia Jetson Nano	MobileNet V2	imagenet 1000	8	0.906	125.2875311	2.136475801	147.7138095	0.08967195725	3.634533248	11.15175837	4.92	3.53	1393.767705
Nvidia Jetson Nano	MobileNet V2	imagenet 1000	16	0.906	125.7341638	2.186259031	88.78350043	0.07893936454	4.651493628	12.66795097	9.84	3.53	2787.535411
Nvidia Jetson Nano	MobileNet V2	imagenet 1000	32								19.7	3.53	5580.736544
Nvidia Jetson Nano	MobileNet V2	imagenet 1000	64								39.4	3.53	11161.47309
Nvidia Jetson Nano	MobileNet V2	imagenet 1000	128								78.7	3.53	22294.61756

Inception V3

Edge device	Setup		Metric										
	Model	Dataset	batch size	Accuracy	Model load time	Data load time	Inference task time	Inference time	IPS (model, dataset load + inference)	IPS (inference only)	Flop (giga)	parameter (m)	Flop / parameter (m)
Nvidia Jetson Xavier	Inception V3	imagenet 1000	1	0.919	132.2537916	1.042407036	147.0621588	0.1454411669	3.566860513	6.879632405	11.5	23.8	483.1932773
Nvidia Jetson Xavier	Inception V3	imagenet 1000	2	0.919	132.1550634	1.044339657	76.47736812	0.07555388546	4.769238864	13.23558668	22.9	23.8	962.1848739
Nvidia Jetson Xavier	Inception V3	imagenet 1000	4	0.919	131.8163593	1.031774044	82.13434196	0.04199156785	4.651537333	23.8143049	45.9	23.8	1928.571429
Nvidia Jetson Xavier	Inception V3	imagenet 1000	8	0.919	132.4758158	1.018234968	35.21715236	0.0256972863	5.927262872	38.91394505	91.8	23.8	3857.142857
Nvidia Jetson Xavier	Inception V3	imagenet 1000	16	0.919	132.8492031	1.030791521	35.47679901	0.02195885707	5.95192651	45.53971079	1.84E+02	23.8	7.73E+03
Nvidia Jetson Xavier	Inception V3	imagenet 1000	32	0.919	132.9358473	1.021033049	37.76926231	0.02061955165	5.96297753	48.49765974	3.67E+02	23.8	1.54E+04
Nvidia Jetson Xavier	Inception V3	imagenet 1000	64	0.919	133.0718038	1.035824299	49.38590672	0.03192617209	5.581181505	31.32226429	7.34E+02	23.8	3.08E+04
Nvidia Jetson Xavier	Inception V3	imagenet 1000	128	0.919	126.1647036	1.033347368	76.42906809	0.05759925744	5.028794448	17.36133493	1.47E+03	23.8	6.18E+04
Nvidia Jetson Nano	Inception V3	imagenet 1000	1								11.5	23.8	483.1932773
Nvidia Jetson Nano	Inception V3	imagenet 1000	2								22.9	23.8	962.1848739
Nvidia Jetson Nano	Inception V3	imagenet 1000	4								45.9	23.8	1928.571429
Nvidia Jetson Nano	Inception V3	imagenet 1000	8								91.8	23.8	3857.142857
Nvidia Jetson Nano	Inception V3	imagenet 1000	16								1.84E+02	23.8	7.73E+03
Nvidia Jetson Nano	Inception V3	imagenet 1000	32								3.67E+02	23.8	1.54E+04
Nvidia Jetson Nano	Inception V3	imagenet 1000	64								7.34E+02	23.8	3.08E+04
Nvidia Jetson Nano	Inception V3	imagenet 1000	128								1.47E+03	23.8	6.18E+04

Edge device	Model	Optimization Platform	Quantization	Dataset	Data file size	batch size	Accuracy	Model load time	Data load time	Inference task time	Inference time	IPS (model, dataset load + Inference)	IPS (inference only)
Nvidia Jetson Xavier	MobileNet V1	TensorFlow 2.5V	FP32	raw image1000	107M	1	0.93	34.7032392	0.977602005	54.16603136	0.0526967833	11.13002465	18.97649035
							0.8154304199	3.17361927	6.146800041	126.9653633	0.1261841333	7.337522529	7.924926644
	MobileNet V2						0.906	61.97093296	0.9851565361	84.89456916	0.08355212736	6.763571484	11.96857616
							0.6952713623	6.596752882	6.228187799	127.2615538	0.126822253	7.138446886	7.88505325
	Inception V3						0.919	122.310173	0.9865400791	142.1121132	0.134392396	3.767768645	7.440897178
							0.8360081787	23.30724597	5.745856285	140.1731172	0.1397072463	5.909249781	7.157824855

Edge device	Model	Optimization Platform	Quantization	Dataset	Data file size	batch size	Accuracy	Model load time	Data load time	Inference task time	Inference time	IPS (model, dataset load + Inference)	IPS (inference only)
Nvidia Jetson Xavier	MobileNet V1	TensorFlow 2.5V	FP32	raw image 1000	107M	1	0.8154304199	3.17361927	6.146800041	126.9653633	0.1261841333	7.337522529	7.924926644
Nvidia Jetson TX2		0.8154304199	4.838519812				10.63910866	174.9827492	0.1745162234	5.250435876	5.730126291		
Rpi + Coral TPU	Edgetpu tflite - TFv2	INT8	0.7466875				2.74034667	0.4634003639	35.30190253	0.005293801611	25.97021505	188.9001654	
Nvidia Jetson Xavier	TensorFlow 2.5V	FP32	0.6952713623				6.596752882	6.228187799	127.2615538	0.126822253	7.138446886	7.88505325	
Nvidia Jetson TX2	MobileNet V2	TensorFlow 2.5V	FP32				0.6952713623	6.596752882	6.228187799	127.2615538	0.126822253	7.138446886	7.88505325
Rpi + Coral TPU	Edge tpu tflite - TFv2	INT8	0.7055976563				2.697169304	0.5129849911	36.1754241	0.005687035799	25.39000418	175.8385274	
Nvidia Jetson Xavier	Inception V3	TensorFlow 2.5V	FP32	0.8360081787	23.30724597	5.745856285	140.1731172	0.1397072463	5.909249781	7.157824855			
Nvidia Jetson TX2		0.8360081787	10.70541787	9.359740734	242.2576783	0.2415484159	3.812096621	4.139956771					
Rpi + Coral TPU	Edge tpu tflite - TFv2	INT8	0.8341835938	2.763403893	0.4765644073	102.1026423	0.0697758156	9.492834802	14.3316132				

Edge device	Model	Optimization Platform	Quantization	Dataset	Data file size	batch size	Accuracy	Model load time	Data load time	Inference task time	Inference time	IPS (model, dataset load + Inference)	IPS (inference only)
Nvidia Jetson Xavier	MobileNet V1	TensorFlow 2.5V	FP32	raw image 1000	107M	1	0.8154304199	3.099321604	6.734873295	124.3262188	0.1238430877	7.453763537	8.074734075
						2	0.8212866821	3.112899542	6.372859955	64.14415121	0.06390857816	13.58143709	15.64735171
						4	0.8119281006	3.146010876	5.918588161	34.61884713	0.03449584007	22.89196681	28.98900267
						8	0.7822409668	3.150560379	5.697333813	20.99315047	0.02093475342	33.51089117	47.76746017
						16	0.7599161542	3.130496979	5.50075531	12.26452923	0.01223503852	47.85654938	81.73247665
						32	0.7634943724	3.129932642	5.502838135	9.182964563	0.009167948961	56.13015578	109.0756509
						64	0.7335840464	3.16045928	5.483443022	7.301475286	0.007291737795	62.71409971	137.1415194
						128	0.7278847694	3.080677748	5.424599171	6.390872002	0.006385187864	67.13144487	156.6124633
						1	0.6952713623	4.216674566	6.683954	127.1585252	0.1266751056	7.243271978	7.894210907
						2	0.6922662354	4.171766043	6.275687456	67.42864656	0.0671857481	12.8409101	14.88410903
						4	0.7005634766	4.249879122	5.887251377	36.27891159	0.03615801263	21.54427553	27.65638727
						8	0.6724165649	4.20552969	5.581227541	21.01236987	0.02095461965	32.46845265	47.72217377
	16	0.6743112594	4.295340776	5.516176462	13.38004827	0.01335142159	43.11912447	74.89839177					
	32	0.656946063	4.197406054	5.540759563	10.32744312	0.01030998278	49.83651445	96.99337251					
	64	0.6237511635	4.209680796	5.468185425	8.486898899	0.008474980354	55.05163394	117.9943738					
	128	0.5953434706	4.233690262	5.471747637	7.474613428	0.007468005419	58.2070438	133.9045627					
	Inception V3	TensorFlow 2.5V	FP32	raw image 1000	107M	1	0.8360081787	6.934961081	6.631296873	138.0621166	0.1376137908	6.595071686	7.266713565
						2	0.8371211548	7.027062654	6.197374582	73.80433631	0.07357701612	11.49045263	13.5912008
						4	0.8426464844	6.988181114	5.9403162	41.27468991	0.04115663791	18.4490996	24.29741716
						8	0.8290883179	7.03455162	5.73078537	25.18580437	0.02512357044	26.34966866	39.80325974
						16	0.832520379	7.061951637	5.630872726	17.21933675	0.01718498349	33.4312187	58.19033812
						32	0.8306828737	7.025388718	5.583262444	13.13729382	0.01311532331	38.84106802	76.24669074
						64	0.8323142529	7.02260494	5.546610355	10.34937906	0.01033916235	44.25121431	96.71963416
						128	0.8211693764	7.003484726	5.51829505	9.343114614	0.009336360216	45.73541414	107.108121

	Edge Device	Model	Dataset	batch size	Accuracy	Model load time	Data load time	Inference task time	Inference time	IPS (model, dataset load + inference)	IPS (inference only)	model size
old	Nvidia Jetson TX2	Mobilenet V1	imagenet 1000	1	0.93	59.92601061	1.6062181	99.50904369	0.09781955647	6.209582935	10.22290466	1.3MB
new	Nvidia Jetson TX2	Mobilenet V1	imagenet 1000	1	0.93	2.744750023	1.727652311	98.29226208	0.09653062463	9.730919374	10.3594067	h5 : 17MB / json : 44KB
old	Nvidia Jetson TX2	Mobilenet V2	imagenet 1000	1	0.906	103.5648978	1.562505245	157.4630725	0.1557951264	3.808201653	6.41868602	2.5MB
new	Nvidia Jetson TX2	Mobilenet V2	imagenet 1000	1	0.906	4.72168088	2.054520845	159.603534	0.1578788443	6.010341644	6.333970867	h5 : 14MB / json : 80KB
old	Nvidia Jetson TX2	Inception V3	imagenet 1000	1	0.919	190.2802024	1.556982994	254.3616292	0.2526663547	2.241152661	3.957788529	4.4M
new	Nvidia Jetson TX2	Inception V3	imagenet 1000	1	0.919	9.547210455	1.754220247	250.801995	0.2490842218	3.815285038	4.014706321	h5 : 92MB / json : 148KB

edge device	model	saved file	model size	model load time
Nvidia Jetson TX2	Mobilenet V1	pb (weight, network)	1.3MB	59.92601061
	Mobilenet V1	h5(weight) / json (network)	17MB / 44KB	2.744750023
	Mobilenet V2	pb (weight, network)	2.5MB	103.5648978
	Mobilenet V2	h5(weight) / json (network)	14MB / 80KB	4.72168088
	Inception V3	pb (weight, network)	4.4M	190.2802024
	Inception V3	h5(weight) / json (network)	92MB / 148KB	9.547210455

model	Optimization Platform	Input size	batch size	Accuracy	Model load time	Data load time	Inference task time	Inference time	IPS (model, dataset load + inference)	IPS (inference only)	model size
Coral-Mobilenet V1	edge tpu tflite - TFv2	224x224x3	1	0.7466875	3.13670969	0.6347105503	49.95246148	0.01541767024	18.61369596	64.86064267	4.6M
Coral-Mobilenet V2	edge tpu tflite - TFv2	224x224x3	1	0.7055976563	3.117435455	1.01078105	50.14280176	0.01652078829	18.42604086	60.52979934	4.2M
Coral-Mobilenet V3	edge tpu tflite - TFv1	224x224x3	1	0.783796875	3.19194746	0.4738194942	190.9825015	0.1694372175	5.137471851	5.901891064	12M
Coral-Inception V3	edge tpu tflite - TFv1	299x299x3	1	0.9133398438	3.178961277	0.442332983	577.8533888	0.5539650221	1.71976533	1.805168124	25M
custom-Mobilenet v1	edge tpu tflite - TFv2	224x224x3	1								
custom-Mobilenet V2	edge tpu tflite - TFv2	224x224x3	1	0.6149140625	3.109596014	0.4486391544	40.48142219	0.01313121428	22.706807	76.15441943	4.1M
custom-Inception V3	edge tpu tflite - TFv2	299x299x3	1	0.8341835938	3.23638773	1.589375019	616.6254816	0.5810301368	1.609136693	1.721081122	24M

TRT vs TF-Lite													
Edge Device	model	Optimization Platform	Quantization	Dataset	batch size	Accuracy	Model load time	Data load time	Inference task time	Inference time	IPS (model, dataset load + inference)	IPS (inference only)	model size
Nvidia Jetson TX2	MobileNet V1	TensorRT	FP32	imagenet 1000	1	0.93	175.3601644	0.7383546829	51.67811322	0.04575625014	4.390261077	21.85493778	19MB
Nvidia Jetson TX2	MobileNet V2	TensorRT	FP32	imagenet 1000	1	0.906	190.3930771	0.7177202702	51.81776929	0.04641133618	4.116433942	21.54646003	17MB
Nvidia Jetson TX2	Inception V3	TensorRT	FP32	imagenet 1000	1	0.919	875.0414793	0.8669939041	142.8614163	0.05401475263	0.9815757357	18.5134607	97MB
Nvidia Jetson TX2	MobileNet V1	TensorRT	FP16	imagenet 1000	1	0.93	177.7084451	0.7184529305	62.83537316	0.05922599053	4.14486264	16.88447911	19MB
Nvidia Jetson TX2	MobileNet V2	TensorRT	FP16	imagenet 1000	1	0.905	191.834491	0.7160468102	63.0419898	0.05860773921	3.91247506	17.06259299	17MB
Nvidia Jetson TX2	Inception V3	TensorRT	FP16	imagenet 1000	1	0.917	913.6164947	0.8568546772	91.71054363	0.07264208722	0.9938536383	13.76612427	97MB
Nvidia Jetson TX2	MobileNet V1	TensorRT	INT8	imagenet 1000	1	0.93	174.7160046	0.7275948524	9.512676716	0.005189748764	5.406674982	192.6875549	19MB
Nvidia Jetson TX2	MobileNet V2	TensorRT	INT8	imagenet 1000	1	0.908	190.4329388	0.6990947723	9.935600758	0.005090838671	4.973444673	196.4312886	17MB
Nvidia Jetson TX2	Inception V3	TensorRT	INT8	imagenet 1000	1								
memory limit													
Nvidia Jetson TX2	MobileNet V1	TensorLite	FP32	imagenet 1000	1	0.4974281616	0.01261377335	0.3975148201	103.4773719	0.08632320921	9.625797088	11.58437006	16.1MB
Nvidia Jetson TX2	Inception V3	TensorLite	FP32	imagenet 1000	1	0.5728826904	0.01034545898	0.3892166615	89.72137737	0.0724643943	11.09620035	13.79988075	13.3MB
Nvidia Jetson TX2	Inception V3	TensorLite	FP32	imagenet 1000	1								90.9MB
Segmentation fault (core dumped)													
Nvidia Jetson TX2	MobileNet V1	TensorLite	FP16	imagenet 1000	1	0.7657109375	0.004939079285	0.4116418362	99.55065703	0.0817362262	10.00327728	12.23447725	8.1MB
Nvidia Jetson TX2	MobileNet V2	TensorLite	FP16	imagenet 1000	1	0.573520752	0.006615638733	0.4141659737	88.37753868	0.0712310019	11.26147428	14.03883103	6.7MB
Nvidia Jetson TX2	Inception V3	TensorLite	FP16	imagenet 1000	1	0.9696956177	0.03152155876	0.9304687977	596.7319081	0.5765084247	1.67309722	1.734580029	45.5MB
Nvidia Jetson TX2	MobileNet V1	TensorLite	INT8	imagenet 1000	1	0.4970228271	0.007335424423	0.4081280231	103.8265989	0.08669070854	9.593056561	11.53526159	4.4MB
Nvidia Jetson TX2	MobileNet V2	TensorLite	INT8	imagenet 1000	1	0.6178867188	0.005393505096	0.4070675373	90.02196455	0.07223293441	11.05773596	13.84410045	3.8MB
Nvidia Jetson TX2	Inception V3	TensorLite	INT8	imagenet 1000	1	0.8316484375	0.01927137375	0.4072015285	641.392735	0.6221491756	1.558071163	1.607331552	23.2MB

TF vs TRT													
Edge Device	model	Optimization Platform	Quantization	Dataset	batch size	Accuracy	Model load time	Data load time	Inference task time	Inference time	IPS (model, dataset load + inference)	IPS (inference only)	model size
Nvidia Jetson TX2	MobileNet V1	TensorFlow 2.5V	FP32	imagenet 1000	1	0.93	59.92601061	1.6062181	99.50904369	0.09781955647	6.209582935	10.22290466	
Nvidia Jetson TX2	MobileNet V2	TensorFlow 2.5V	FP32	imagenet 1000	1	0.906	103.5648978	1.562505245	157.4630725	0.1557951264	3.808201653	6.41868602	
Nvidia Jetson TX2	Inception V3	TensorFlow 2.5V	FP32	imagenet 1000	1	0.919	190.2802024	1.556982994	254.3616292	0.2526663547	2.241152661	3.957788529	
Nvidia Jetson TX2	MobileNet V1	TensorRT	FP32	imagenet 1000	1	0.93	175.3601644	0.7383546829	51.67811322	0.04575625014	4.390261077	21.85493778	
Nvidia Jetson TX2	MobileNet V2	TensorRT	FP32	imagenet 1000	1	0.906	190.3930771	0.7177202702	51.81776929	0.04641133618	4.116433942	21.54646003	
Nvidia Jetson TX2	Inception V3	TensorRT	FP32	imagenet 1000	1	0.919	875.0414793	0.8669939041	142.8614163	0.05401475263	0.9815757357	18.5134607	

Edge Device	Model	Quantization	Accuracy	IPS (inference only)	Model size
Nvidia Jetson TX2	MobileNet V1	FP32	0.93	21.85493778	19MB
	MobileNet V1	FP16		16.88447911	
	MobileNet V1	INT8		192.6875549	
	MobileNet V2	FP32	0.906	21.54646003	17MB
	MobileNet V2	FP16	0.905	17.06259299	
	MobileNet V2	INT8	0.908	196.4312886	
	Inception V3	FP32	0.919	18.5134607	97MB
	Inception V3	FP16	0.917	13.76612427	
	Inception V3	INT8	memory limit		

Model load time	Data load time	Inference task time	Inference time	IPS (model, dataset load + inference)
59.92601061	1.6062181	99.50904369	0.09781955647	6.209582935
103.5648978	1.562505245	157.4630725	0.1557951264	3.808201653
190.2802024	1.556982994	254.3616292	0.2526663547	2.241152661
175.3601644	0.7383546829	51.67811322	0.04575625014	4.390261077
190.3930771	0.7177202702	51.81776929	0.04641133618	4.116433942
875.0414793	0.8669939041	142.8614163	0.05401475263	0.9815757357

Edge Device	Model	Quantization	Optimization Platform	Accuracy	IPS (inference only)
Nvidia Jetson TX2	MobileNet V1	FP32	TensorFlow 2.5V	0.93	10.22290466
	MobileNet V1		TensorRT		21.85493778
	MobileNet V2		TensorFlow 2.5V	0.906	6.41868602
	MobileNet V2		TensorRT	21.54646003	
	Inception V3		TensorFlow 2.5V	0.919	3.957788529
	Inception V3		TensorRT		18.5134607

single GPU 할당 없이 추론 성능 비교

		IPS (model, dataset load + inference)	IPS (inference only)
	container1	11.47055448	24.48302847
	container2	21.59180339	24.15756832
	container3	21.68663102	24.42456412
	container4	21.62289646	24.44595701
	container5	22.05667186	25.30064568

single GPU 할당 여부에 따른 추론 성능 비교

		IPS (model, dataset load + inference)	IPS (inference only)
	with-gpu container	27.71799156	31.57648431
	without-gpu container	27.52136015	31.26643088

		ggaman.com/vram: 25	IPS (model, dataset load + inference)	IPS (inference only)		
		container1	4.753275123	7.603665533	avg=	7.595620248
		container2	4.725597928	7.447515023		
		container3	4.612298272	7.662968067		
		container4	4.561359302	7.668332368		
		ggaman.com/vram: 50	IPS (model, dataset load + inference)	IPS (inference only)		
		container1	4.778885324	7.655733805	avg=	7.659350936
		container2	4.62523149	7.662968067		

Text Classification

RNN

Edge device	Setup					Metric										
	OS (Ubuntu/L4T/ Docker)	Model	Dataset	Optimization Platform	Quantization	batch size	Accuracy	Model load time	Data load time	Inference time	Inference time (avg)	SPS	SPS(inf)	Flop (giga)	parameter (m)	Flop / parameter (m)
Nvidia Jetson Nano	18.04 / 32.6.1 / 20.10	DistiBERT	GLUE SST-2	TensorFlow 2.5V	FP32	64	0.9106	80.0597	2.5090	200.0222	0.2294	3.0857	4.3595	6.7	66	101.5151515
Nvidia Jetson Nano	18.04 / 32.6.1 / 20.10	DistiBERT	GLUE SST-2	TensorFlow 2.5V	FP32	128	0.9106	80.0597	2.3074	200.3049	0.2297	3.0848	4.3534	6.7	66	101.5151515
Google Coral TPU	18.04 / 32.6.1 / 20.10	DistiBERT	GLUE SST-2	TFLite	INT8	1								6.7	66	101.5151515
Google Coral TPU	18.04 / 32.6.1 / 20.10	DistiBERT	GLUE SST-2	TFLite	INT8	2								6.7	66	101.5151515
Google Coral TPU	18.04 / 32.6.1 / 20.10	DistiBERT	GLUE SST-2	TFLite	INT8	4								6.7	66	101.5151515
Google Coral TPU	18.04 / 32.6.1 / 20.10	DistiBERT	GLUE SST-2	TFLite	INT8	8								6.7	66	101.5151515
Google Coral TPU	18.04 / 32.6.1 / 20.10	DistiBERT	GLUE SST-2	TFLite	INT8	16								6.7	66	101.5151515
Google Coral TPU	18.04 / 32.6.1 / 20.10	DistiBERT	GLUE SST-2	TFLite	INT8	32								6.7	66	101.5151515
Google Coral TPU	18.04 / 32.6.1 / 20.10	DistiBERT	GLUE SST-2	TFLite	INT8	64								6.7	66	101.5151515
Google Coral TPU	18.04 / 32.6.1 / 20.10	DistiBERT	GLUE SST-2	TFLite	INT8	128								6.7	66	101.5151515

BERT

Nvidia Jetson TX1	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	1	0.9317	111.6667	0.8395	121261.4603	4.8505	0.2060	0.2062	13.39	110	121.7272727
Nvidia Jetson TX1	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	2	0.9317	113.7475	0.9170	115239.0107	4.6096	0.2167	0.2169	13.39	110	121.7272727
Nvidia Jetson TX1	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	4	0.9317	113.7475	0.8231	114094.7305	4.5638	0.2189	0.2191	13.39	110	121.7272727
Nvidia Jetson TX1	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	8	0.9221	113.7475	0.9000	112788.7456	4.5115	0.2214	0.2217	13.39	110	121.7272727
Nvidia Jetson TX1	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	16								13.39	110	121.7272727
Nvidia Jetson TX1	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	32								13.39	110	121.7272727
Nvidia Jetson TX1	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	64								13.39	110	121.7272727
Nvidia Jetson TX1	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	128								13.39	110	121.7272727
Nvidia Jetson TX2	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	1	0.9317	99.3734	0.6965	106923.4484	4.2769	0.2336	0.2338	13.39	110	121.7272727
Nvidia Jetson TX2	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	2	0.9317	99.3734	0.2575	103062.0319	4.1225	0.2423	0.2426	13.39	110	121.7272727
Nvidia Jetson TX2	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	4	0.9317	99.3734	0.2343	102395.8735	4.0958	0.2439	0.2442	13.39	110	121.7272727
Nvidia Jetson TX2	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	8	0.9317	99.3734	0.2458	102382.2849	4.0953	0.2439	0.2442	13.39	110	121.7272727
Nvidia Jetson TX2	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	16	0.9317	99.3734	0.1973	102362.6500	4.0945	0.2440	0.2442	13.39	110	121.7272727
Nvidia Jetson TX2	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	32	0.9317	99.3734	0.1970	102202.6855	4.0881	0.2444	0.2446	13.39	110	121.7272727
Nvidia Jetson TX2	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	64	0.9317	99.3734	0.1952	101898.2051	4.0759	0.2451	0.2453	13.39	110	121.7272727
Nvidia Jetson TX2	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	128	0.9317	99.3734	0.4065	101542.3949	4.0617	0.2460	0.2462	13.39	110	121.7272727
Nvidia Jetson Xavier	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	1	0.9317	66.1882	0.1192	6262.0173	0.2505	3.9505	3.9923	13.39	110	121.7272727
Nvidia Jetson Xavier	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	2	0.9317	66.1882	0.1531	6082.0433	0.2433	4.0661	4.1105	13.39	110	121.7272727
Nvidia Jetson Xavier	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	4	0.9317	66.1882	0.1452	6022.0454	0.2409	4.1062	4.1514	13.39	110	121.7272727
Nvidia Jetson Xavier	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	8	0.9317	66.1882	0.1506	5962.0398	0.2385	4.1471	4.1932	13.39	110	121.7272727
Nvidia Jetson Xavier	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	16	0.9317	66.1882	0.1283	5898.9912	0.2360	4.1909	4.2380	13.39	110	121.7272727
Nvidia Jetson Xavier	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	32	0.9317	66.1882	0.1160	5789.2058	0.2316	4.2695	4.3194	13.39	110	121.7272727
Nvidia Jetson Xavier	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	64	0.9317	66.1882	0.1142	5789.6695	0.2316	4.2691	4.3180	13.39	110	121.7272727
Nvidia Jetson Xavier	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	128	0.9317	66.1882	0.1143	5841.9759	0.2337	4.2314	4.2794	13.39	110	121.7272727
Nvidia Jetson Nano	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	1								13.39	110	121.7272727
Nvidia Jetson Nano	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	2								13.39	110	121.7272727
Nvidia Jetson Nano	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	4								13.39	110	121.7272727
Nvidia Jetson Nano	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	8								13.39	110	121.7272727
Nvidia Jetson Nano	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	16								13.39	110	121.7272727
Nvidia Jetson Nano	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	32								13.39	110	121.7272727
Nvidia Jetson Nano	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	64								13.39	110	121.7272727
Nvidia Jetson Nano	18.04 / 32.6.1 / 20.10	BERT	IMDb	TensorFlow 2.5V	FP32	128								13.39	110	121.7272727
Google Coral TPU	18.04 / 32.6.1 / 20.10	BERT	IMDb	TFLite	INT8	1								13.39	110	121.7272727
Google Coral TPU	18.04 / 32.6.1 / 20.10	BERT	IMDb	TFLite	INT8	2								13.39	110	121.7272727
Google Coral TPU	18.04 / 32.6.1 / 20.10	BERT	IMDb	TFLite	INT8	4								13.39	110	121.7272727
Google Coral TPU	18.04 / 32.6.1 / 20.10	BERT	IMDb	TFLite	INT8	8								13.39	110	121.7272727
Google Coral TPU	18.04 / 32.6.1 / 20.10	BERT	IMDb	TFLite	INT8	16								13.39	110	121.7272727
Google Coral TPU	18.04 / 32.6.1 / 20.10	BERT	IMDb	TFLite	INT8	32								13.39	110	121.7272727
Google Coral TPU	18.04 / 32.6.1 / 20.10	BERT	IMDb	TFLite	INT8	64								13.39	110	121.7272727
Google Coral TPU	18.04 / 32.6.1 / 20.10	BERT	IMDb	TFLite	INT8	128								13.39	110	121.7272727