

#### Learning *Multi-Scene* **Absolute Pose Regression** with Transformers Position 0 1 2 3 MLP Transformer Selected Scene Index CNN Orientation 0 1 2 3 ► MLP ►a Transformer Bar-Ilan **Ron Ferens**

Yoli Shavit

Yosi Keller

# **Camera Pose Regression**

#### **Camera Pose Estimation**



**Camera Pose** 



Input

Output

#### **Absolute Camera Pose Regression**

A learning-based method for solving the camera pose estimation problem



### **Camera Pose Estimation at Inference Time**





**Localization Pipelines** 

**Absolute Pose Regressors (APRs)** 

### **Camera Pose Estimation at Inference Time**



✓ Fast (order of magnitude)
✓ Light-weight
✓ Standalone



**Localization Pipelines** 

**Absolute Pose Regressors (APRs)** 

### The Cons of Single-Scene APRs

✓ Fast

x Less accurate

- ✓ Light-weight
- ✓ Standalone

x Trained per scene



For localizing images from N scenes we need to train, deploy and choose from N models

# Learning Multi-Scene Pose Regression with Transformers



We extract visual features using a convolutional backbone and then encode, project and flatten activation maps











 $\hat{S}_x$  and  $\hat{S}_q$  are learned parameters representing task uncertainty

# **Comparison with MSPN**

Median position and orientation errors of our method and a recent multi-scene approach (MSPN) for the **CambridgeLandmarks** (top) and **7Scenes** (bottom) datasets

Method	K. College	Old Hospital	Shop Facade	St. Mary
MSPN 3	1.73/3.65	2.55/4.05	2.92/7.49	2.67/6.18
MS-Transformer (ours)	0.83/1.47	1.81/2.39	0.86/3.07	1.62/3.99

Method	Chess	Fire	Heads	Office	Pumpkin	Kitchen	Stairs
MSPN[3]	<b>0.09</b> /4.76	0.29/10.5	0.16/13.1	<b>0.16</b> /6.8	0.19/5.5	<b>0.21</b> /6.61	0.31/11.63
MS-Transformer (ours)	0.11/ <b>4.66</b>	0.24/9.6	0.14/12.19	0.17/ <b>5.66</b>	0.18/4.44	0.17/5.94	0.26/8.45

# **Comparison with Single-Scene APRs**

Mean of median position and orientation errors and the respective ranks of our method, MSPN and state-of-the-art single-scene APRs on the CambridgeLandmarks (left) and 7Scenes (right) datasets.

Method	Average [m/deg]	Ranks	Method	Average [m/deg]
Single-scene APRs			Single-scene APRs	
PoseNet [17]	2.09/6.84	10/11	PoseNet [17]	0.44/10.4
BayesianPN [15]	1.92/6.28	8/10	BayesianPN [15]	0.47/9.81
LSTM-PN 35	1.30/5.52	2/9	LSTM-PN 35	0.31/9.86
SVS-Pose [21]	1.33/5.17	3/7	GPoseNet [8]	0.31/9.95
GPoseNet [8]	2.08/4.59	6/3	PoseNet-Learnable [16]	0.24/7.87
PoseNet-Learnable [16]	1.43/2.85	5/2	GeoPoseNet [16]	0.23/8.12
GeoPoseNet [16]	1.63/2.86	6/3	MapNet [7]	0.21/7.78
MapNet [7]	1.63/3.64	6/5	IRPNet [29]	0.23/8.49
IRPNet [29]	1.42/3.45	4/4	AttLoc 36	0.20/7.56
Multi-scene APRs			Multi-scene APRs	
MSPN 3	2.47/5.34	11/8	MSPN [3]	0.20/8.41
MS-Transformer (Ours)	1.28/2.73	1/1	MS-Transformer (Ours)	0.18/7.28

Ranks

10/11

 $\frac{11/8}{8/9}$ 

8/8

 $\frac{7}{4}$  5/5

 $\frac{4}{3}$ 5/7

2/2

2/6

1/1

### From Multi-Scene to Multi-Dataset

APR Method	CambridgeLand.	7Scenes
	[m/deg]	[m/deg]
Single-scene [16]	1.43/2.85	0.24/7.87
Multi-scene (Ours)	1.28/2.73	0.18/7.28
Multi-dataset (Ours)	1.50/ 2.57	0.22/6.78

Our method is able to learn multiple scenes from datasets with different scales and properties

## **Robustness and Scalability**

Encoder/Decoder	Position	Orientation	Transformer Dimension	Position	Orientation
# Layers	[meters]	[degrees]		[meters]	[degrees]
2	0.19	7.48	64	0.18	8.06
4	0.18	6.94	128	0.19	7.56
6	0.18	7.28	256	0.18	7.28
8	0.18	6.92	512	0.18	7.19

Backbone	Position	Orientation
	[meters]	[degrees]
Resnet50	0.19	8.6
EfficientNetB0	0.18	7.28
EfficientNetB1	0.17	7.26

Our method maintains state-of-the-art performance across architectural choices

Num. Scenes	Runtime [ms]		Memory [Mb]		-
Num. Layers	2	6	2	6	
1	18.8	34.6	40.8	74.6	-
4	18.8	35	40.8	74.6	
7	19.2	35.2	40.8	74.6	
10	19.2	35.2	40.8	74.6	The memory footprint for a
100	19.6	35.4	41.0	74.8	1000 scenes with a single scene
500	21.0	41.0	41.8	75.6	1000 seenes with a single seene
1000	27.0	58.6	42.8	76.7	APR approach is <b>~5000Mb</b>

# Conclusion

# We propose a novel transformer-based approach for multi-scene absolute pose regression

- Two Transformer Encoders separately attend to position- and orientationinformative image cues
- Two Transformer Decoders attend to scene-specific information

#### Our approach is shown to provide a new state-of-the-art APR accuracy

- Outperforming single and multi-scene APRs across indoor and outdoor benchmarks
- Demonstrating robustness to specific architecture choices



X

MLP

► MLP

Bar-Ilan

#### Learning *Multi-Scene* **Absolute Pose Regression** Position ▶ 0 1 2 with Transformers Transformer https://github.com/yolish/multiscene-pose-transformer Selected Scene Index CNN Query Image

Orientation

Transformer





Yoli Shavit

Ron Ferens Yosi Keller