

# UL-VIO: Ultra-lightweight Visual-Inertial Odometry with Noise Robust Test-time Adaptation

Jinho Park<sup>1</sup>, Se Young Chun<sup>2</sup>, Mingoo Seok<sup>1</sup>  
<sup>1</sup>Columbia University, <sup>2</sup>Seoul National University



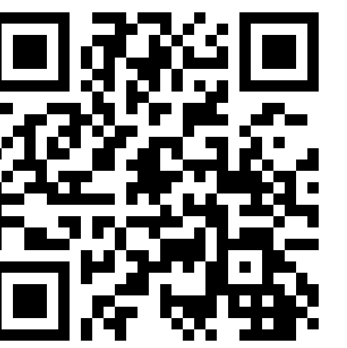
VLSILab  
@COLUMBIA UNIV



Sponsored by:

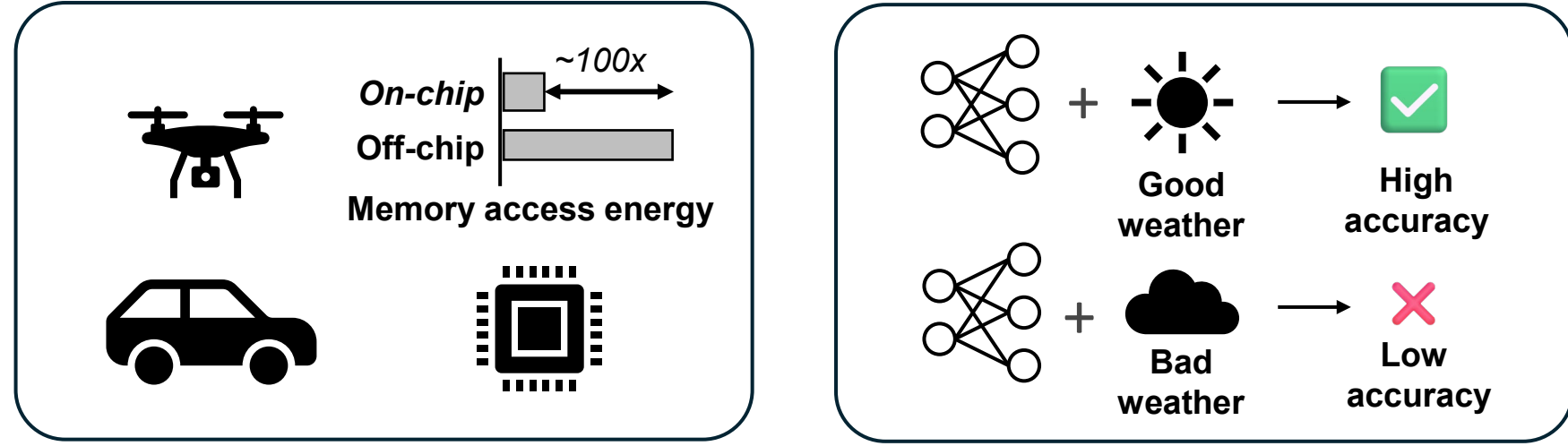


Project page



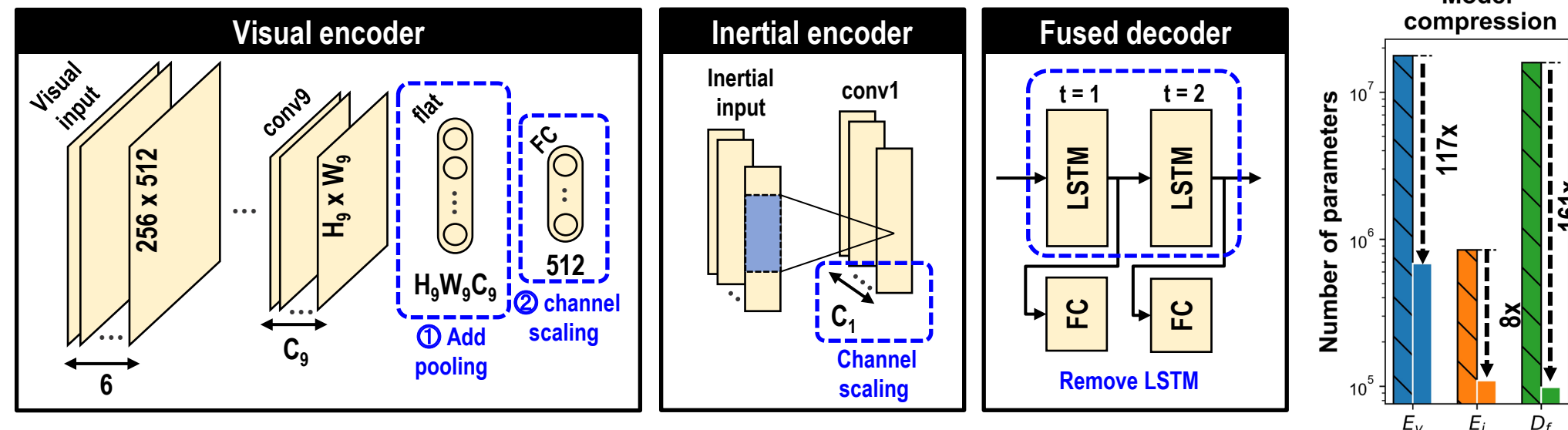
Linkedin

## Motivation for Lightweight & Robust Network



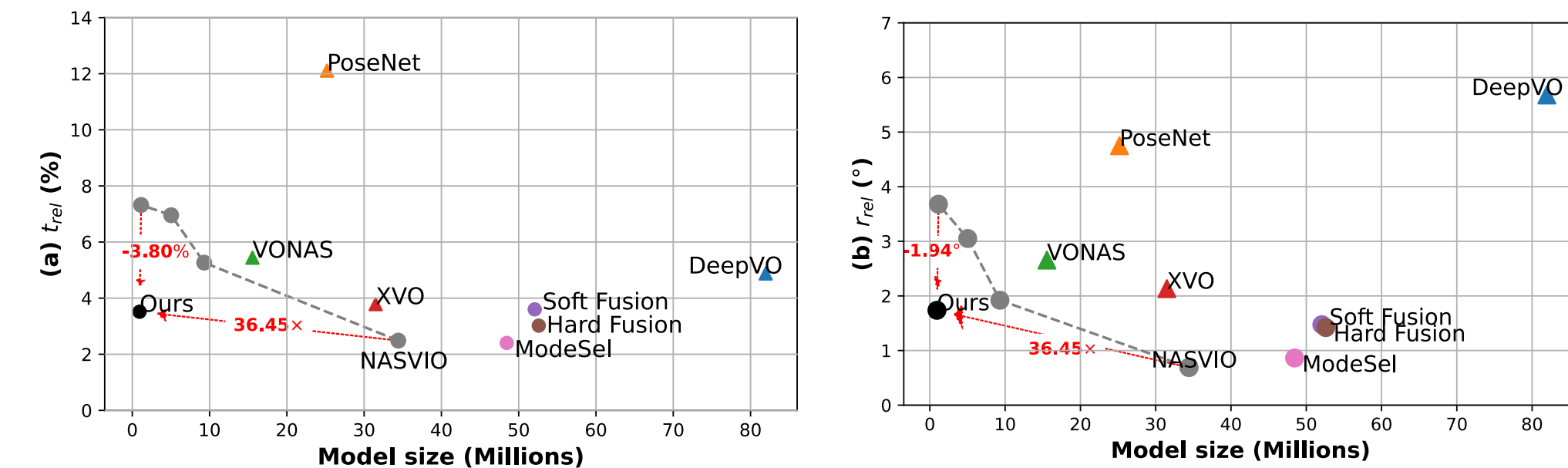
- Need for **lightweight** neural network to be hosted entirely by *on-chip* memory, which is of few MB in modern processors
- Robust to **environmental factors** and **noises** inducing distribution shifts

## Model Compression



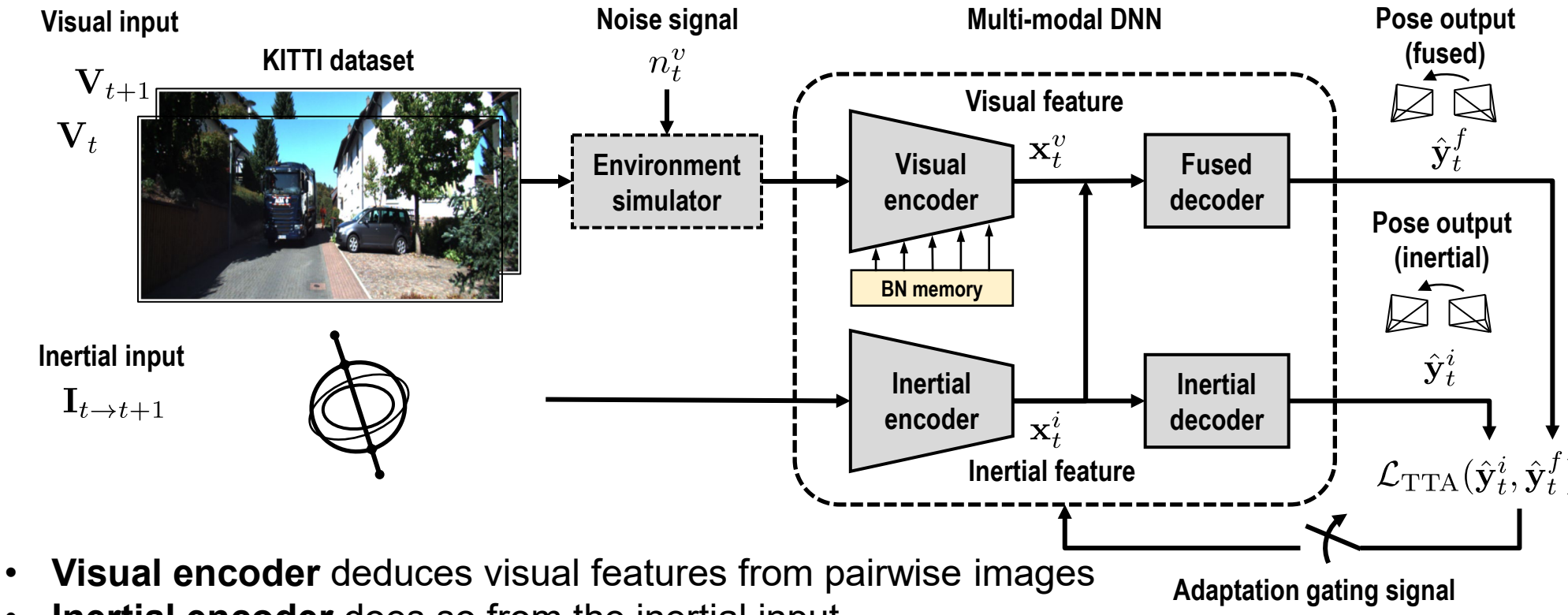
- Channel scaling** – model parameter quadratically proportional to channel size in  $E_v$ ,  $E_i$ ,  $D_f$
- Addition of **AveragePool** prior to FullyConnected (FC) to reduce feature size in  $E_v$
- Replace LSTM w/ FC** layer to reduce the model size in  $D_f$

## Model Compression (results)



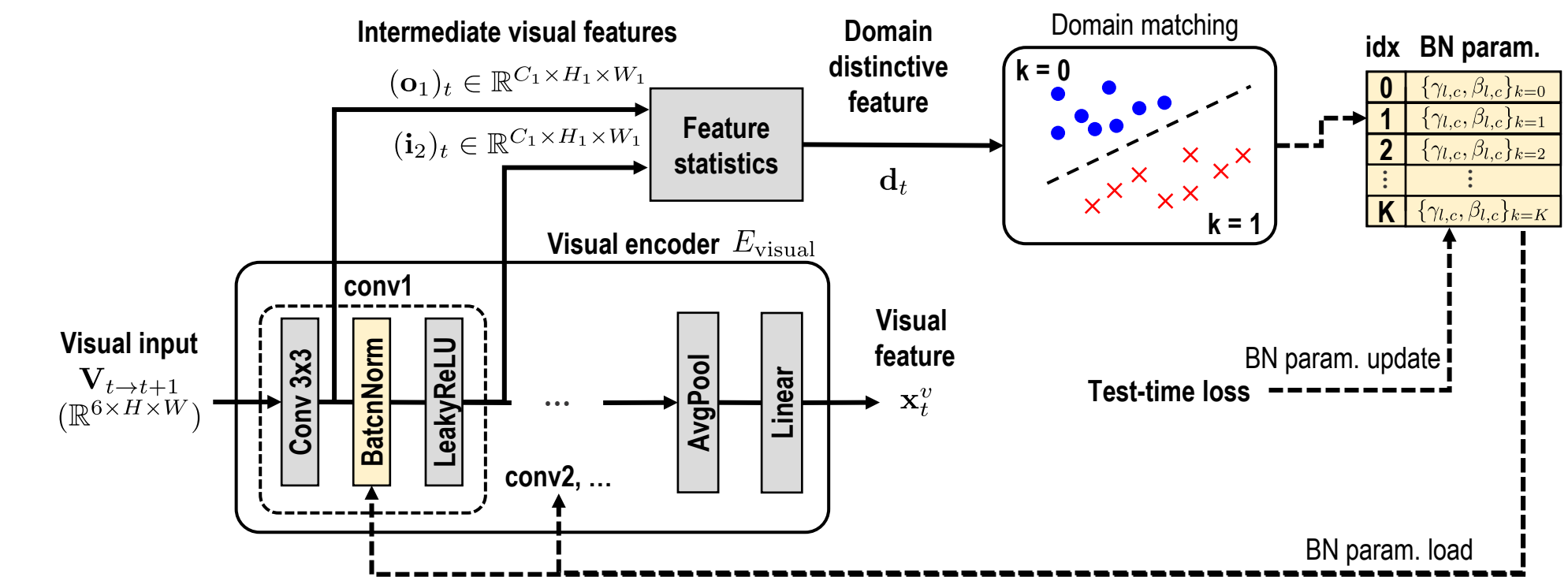
- We achieve **36x model size reduction** w/ a minute (1%) increase in absolute pose estimation error
- Latest **Apple A16** and **Qualcomm Snapdragon** CPUs possess only a **few MB** of *on-chip* memory

## Multi-modal Consistency-based TTA



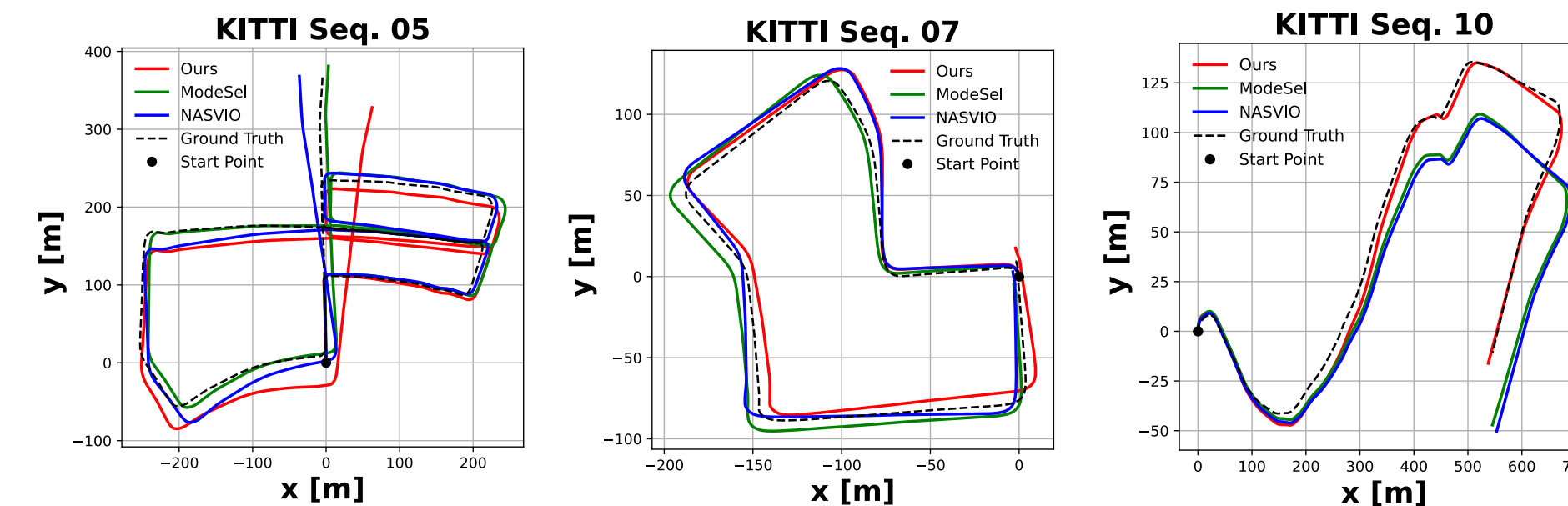
- Visual encoder** deduces visual features from pairwise images
- Inertial encoder** does so from the inertial input
- Decoder** predicts pose transformation from fused features

## Dictionary-based Visual Encoder Adaptation

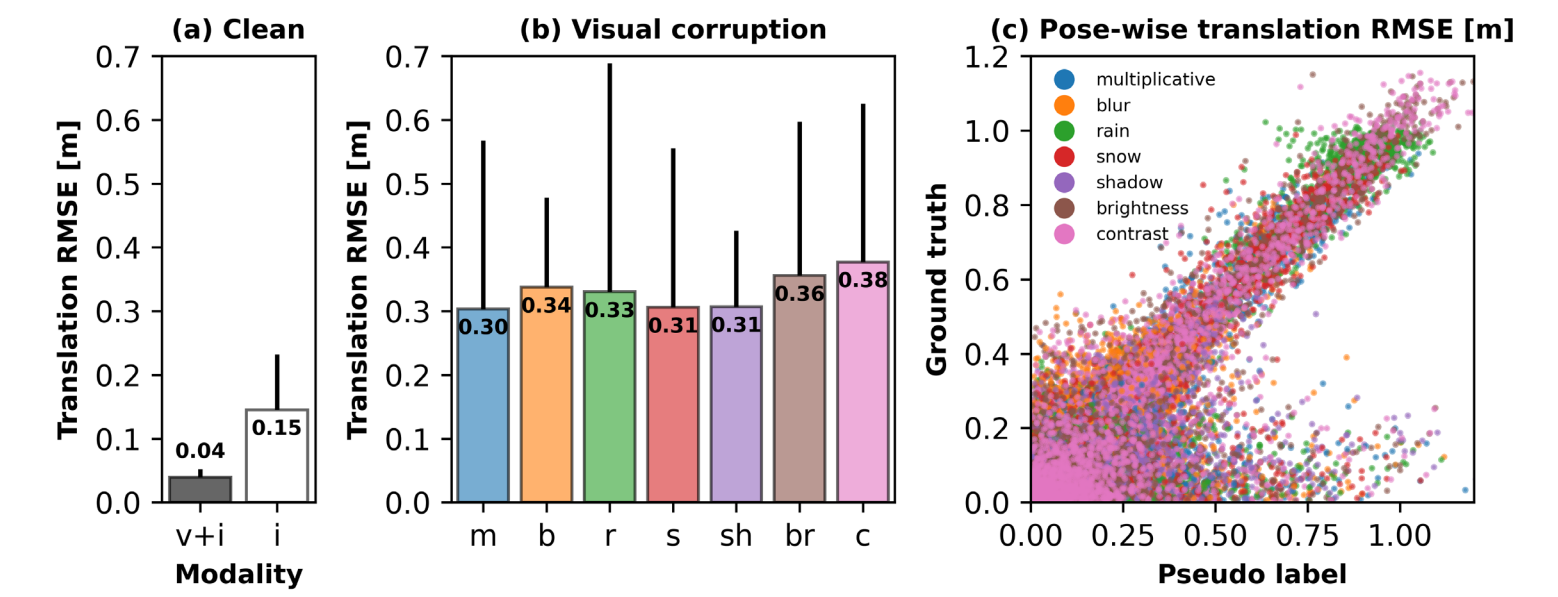


- Forward-path generate **visual feature**  $x^v$ , which will be used for decoder in the deeper layers
- First layer conv output  $o_1$  is used to generate **domain distinctive feature**  $d_t$
- Domain  $k$  is found using  $l_2$  distance with the proxies  $\{d^k\}$
- BatchNorm** parameter corresponding to the domain  $k$  is updated

## Trajectory Results for UL-VIO



## Motivation for Multi-modal Consistency



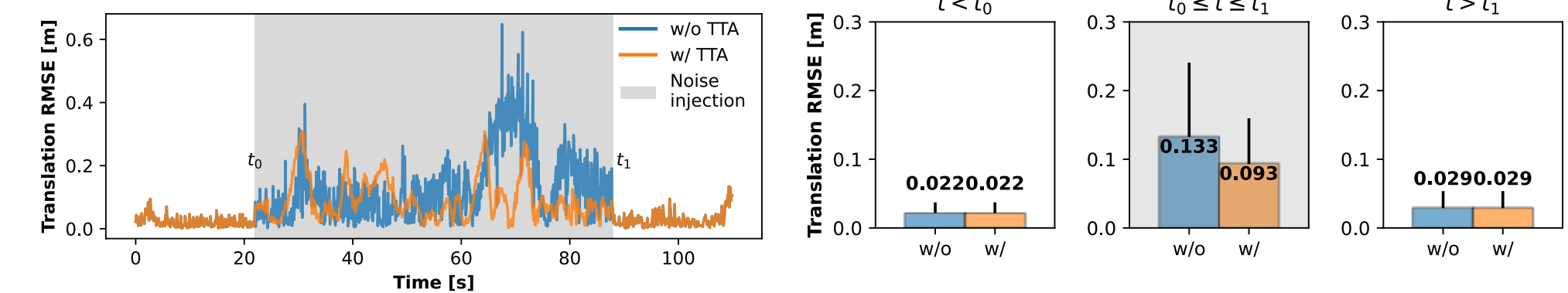
- While using **visual feature** along with inertial is preferred, inertial-only result become reliable under noise
- Inertial-only pose output (**pseudo label**) is highly correlated with the **ground-truth** target pose

## Against Fine-tuned Baselines

Model	Average pose-wise $t_{rmse}$ [m]							
	Clean	Multi.	Blur	Rain	Snow	Shadow	Bright.	Cont.
Source	0.059	0.154	<u>0.261</u>	0.176	0.191	0.203	0.226	<u>0.250</u>
Fine-tuned with adver. noise (FT)								
Multi.	0.099	<u>0.129</u>	0.394	0.227	0.372	0.192	0.299	0.331
Blur	0.115	0.176	0.263	0.193	0.247	0.184	0.242	0.261
Rain	0.289	0.325	0.372	<b>0.095</b>	0.394	0.311	0.525	0.531
Snow	0.091	0.148	0.319	0.263	<u>0.183</u>	0.208	0.369	0.450
Shadow	0.085	<b>0.112</b>	0.322	0.179	0.243	<b>0.121</b>	<u>0.221</u>	0.252
Bright.	0.091	0.151	0.312	0.177	0.226	0.185	0.233	0.278
Cont.	0.093	0.150	0.330	0.197	0.219	0.184	0.237	0.273
TTA (ours)	-	0.156	<b>0.230</b>	<u>0.143</u>	<b>0.172</b>	<u>0.155</u>	<b>0.193</b>	<b>0.212</b>

- Except for one case, e.g., multiplicative noise, our TTA method has the **best or 2nd best accuracy**

## Continual TTA



- Translation RMSE quickly increases when the **visual corruption** is applied
- TTA** mitigates the increase in the error to some extent

## Continual TTA w/ Dynamic Noise Shifts

Time	Seq. 05				Seq. 07				Seq. 10				Avg.
Noise	Blur	Rain	Snow	Con.	Blur	Rain	Snow	Con.	Blur	Rain	Snow	Con.	
Baseline	0.118	0.121	<b>0.103</b>	0.166	0.127	0.153	0.110	0.191	0.137	0.134	<b>0.120</b>	0.167	0.137
TTA	<b>0.112</b>	<b>0.107</b>	0.110	<b>0.107</b>	<b>0.101</b>	<b>0.108</b>	<b>0.106</b>	<b>0.104</b>	<b>0.123</b>	<b>0.124</b>	0.121	<b>0.127</b>	<b>0.113</b>
ddf acc.	97.9	100	100	100	98.2	100	100	100	98.8	100	100	100	99.6

- We inject different types of noise {blur, rain, snow, contrast} in a continuous manner for KITTI