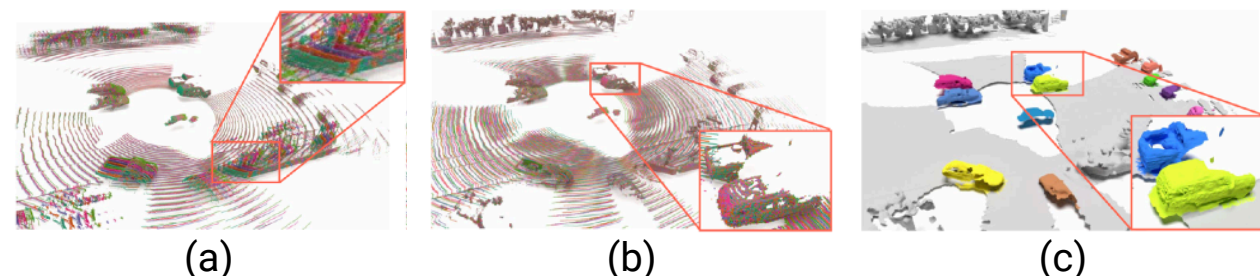


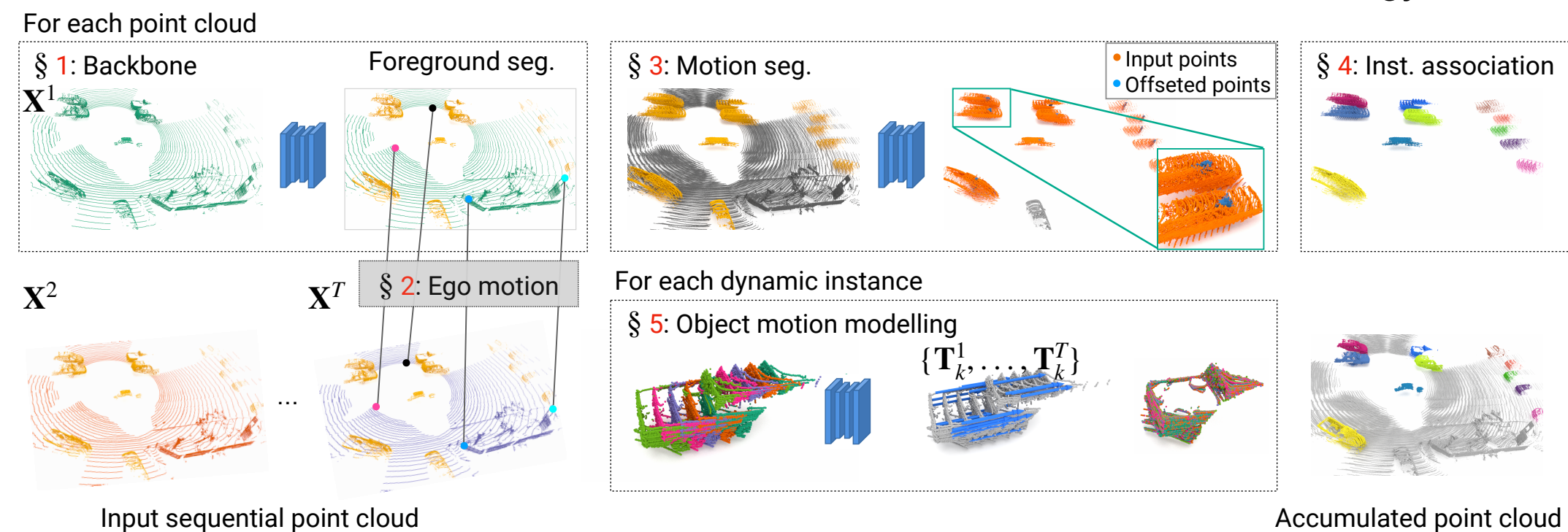


## Problem statement

- Given an *ordered* sequence of  $T$  point cloud frames (a), the goal is to estimate the flow vectors that align each of the *source* frames to the *target* frame, and hence accumulate the point clouds (b). Such accumulation can serve as a pre-processing step to benefit downstream tasks like surface reconstruction (c).



## Methodology



- We first estimate the ego motion (§2) using predicted *background* points. Then we segment the foreground points into *static* and *dynamic* parts through motion segmentation (§3). The instance IDs of the dynamic parts are obtained from spatio-temporal association (§4). Built upon the scene *rigidity*, we explain the flow vectors of the static parts by estimated ego-motion, and that of the dynamic instances by regressed per-instance rigid motion (§5).

$$\text{Loss} = L_{\text{ego}} + L_{\text{FG}} + L_{\text{motion}} + L_{\text{offset}} + L_{\text{obj}}$$

## Experimental results

- SoTA results on multi-frame scene flow evaluation.

Dataset	Method	Static part				Dynamic foreground				
		EPE avg.↓	AccS↑	AccR↑	ROutlier↓	EPE avg. ↓	EPE med.↓	AccS↑	AccR↑	ROutliers ↓
Waymo	PPWC-Net [63]	0.414	17.6	40.2	12.1	0.475	0.201	9.0	29.3	22.4
	FLOT [44]	0.129	65.2	85.0	2.8	0.625	0.231	9.8	27.4	33.8
	WsRSF [19]	0.063	87.3	96.6	0.6	0.381	0.094	31.3	64.0	10.1
	NSFPrior [34]	0.187	79.8	89.1	4.7	0.237	0.077	44.7	68.6	11.5
	Ours	<b>0.028</b>	<b>97.5</b>	<b>99.5</b>	<b>0.1</b>	<b>0.197</b>	<b>0.062</b>	<b>53.3</b>	<b>77.5</b>	<b>5.9</b>
	Ours+	<b>0.018</b>	<b>99.0</b>	<b>99.7</b>	<b>0.1</b>	<b>0.173</b>	<b>0.043</b>	<b>69.1</b>	<b>86.9</b>	<b>5.1</b>
nuScenes	Ours (w. ground)	0.042	91.9	98.8	0.1	0.219	0.071	47.1	72.8	8.5
	PPWC-Net [63]	0.316	16.1	37.0	8.7	0.661	0.307	7.6	24.2	31.9
	FLOT [44]	0.153	51.7	78.3	4.3	1.216	0.710	3.0	10.3	63.9
	WsRSF [19]	0.195	57.4	82.6	4.8	0.539	0.204	17.9	37.4	22.9
	NSFPrior [34]	0.584	38.9	56.7	26.9	0.707	0.222	19.3	37.8	32.0
	Ours	<b>0.111</b>	<b>65.4</b>	<b>88.6</b>	<b>1.1</b>	<b>0.301</b>	<b>0.146</b>	<b>26.6</b>	<b>53.4</b>	<b>12.1</b>
	Ours+	<b>0.091</b>	<b>72.8</b>	<b>91.9</b>	<b>0.9</b>	<b>0.301</b>	<b>0.135</b>	<b>32.7</b>	<b>56.7</b>	<b>13.7</b>
	Ours (w. ground)	0.134	55.3	83.8	1.9	0.37	0.182	18.2	43.8	17.5

- Generalisation across different numbers of input frames on Waymo dataset.

	3	4	5	6	7	8	9	10
static EPE avg.	<b>0.022</b>	0.025	0.028	0.032	0.037	0.044	0.054	0.066
dynamic EPE avg.	0.199	<b>0.168</b>	0.190	0.218	0.250	0.294	0.348	0.412

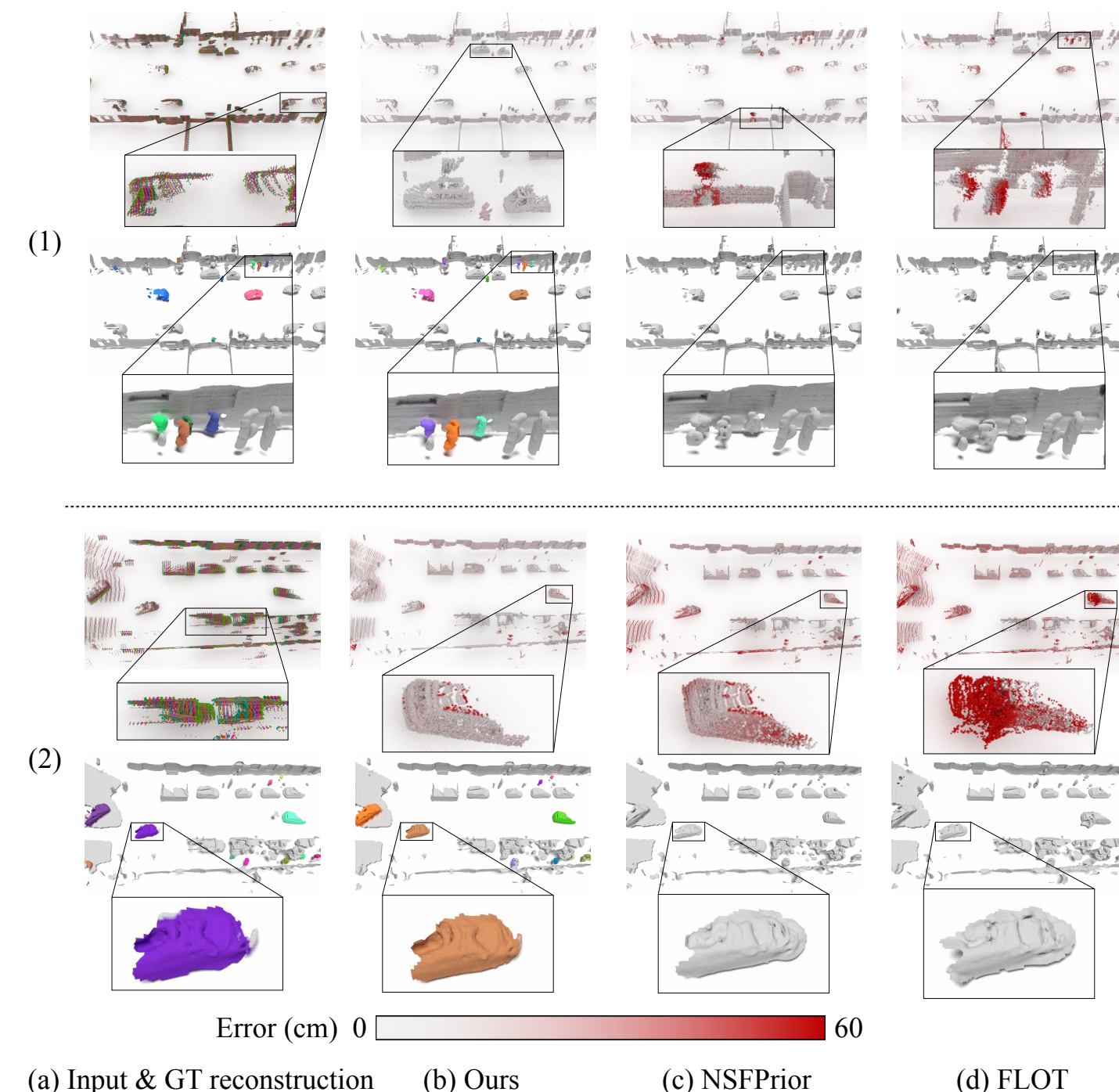
- Superior runtime per 5-frame (Waymo) and 11-frame (nuScenes) sample.

	Waymo	nuScenes
PPWC-Net [63]	0.608	0.990
FLOT [44]	1.028	2.010
WsRSF [19]	1.252	1.460
NSFPrior [34]	212.256	63.460
Ours	<b>0.174</b>	<b>0.250</b>

	Waymo	nuScenes
ego-motion estimation	0.100	0.188
motion segmentation	0.024	0.040
instance association	0.036	0.009
TubeNet	0.014	0.013

- Surface reconstruction on Waymo dataset.



- Surface reconstruction on nuScenes dataset.

