

1 We thank the reviewers for their valuable comments. In the following, we will address each comment in detail.

2 **R1, R2: Novelty.** To our knowledge, learning-based density models of point cloud sweeps [spatial+intensity] for
 3 compression are not widely explored, and modeling temporal relationships between point clouds is non-trivial. Our
 4 novel tree-structured density models combine existing techniques (deep sets, cont. convs) to significantly advance
 5 beyond the scope of OctSqueeze – fully exploiting temporal redundancies and modeling intensity attributes.

6 **R2: Additional baselines.** We benchmarked MPEG TMC13 on UrbanCity and SemanticKITTI (see Fig. 1). On
 7 average, our approach achieves a better bitrate vs. reconstruction trade-off, especially at lower bitrates. We could not
 8 benchmark MPEG TMC2 due to segmentation fault errors in their reference implementation¹. We have emailed
 9 the authors and hope to include TMC2 in our final paper. That said, TMC2 is related to our MPEG Range baseline:
 10 both compress 2D projections of the point cloud using video codecs. We found that tree-based codecs perform better.

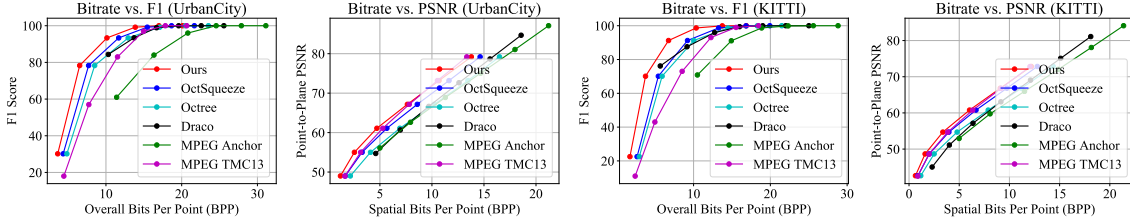


Figure 1: Bitrate vs. reconstruction quality on UrbanCity and SemanticKITTI.

11 **R2, R3: Runtime/GPU passes.** We benchmarked our approach on a workstation with an Nvidia GTX 1080Ti GPU
 12 and Intel Xeon E5 CPU. Our octree data structure and range coder is written in C++ for CPU execution, and our
 13 spatial/intensity entropy models run on the GPU. During encoding, we predict distributions over occupancy symbols
 14 with one forward pass of the entropy models. During decoding, we require $O(D)$ forward passes (D is the octree’s
 15 depth), each interleaved with range decoding to reconstruct the octree level by level. We also benchmarked OctSqueeze,
 16 Draco, MPEG Anchor, and MPEG TMC13. Our encoding speed is fast, and our total runtime is on par with TMC13.
 17 Moreover, there are many opportunities to speed up our research code significantly, especially in CPU/GPU I/O.

Depth	MuSACLE (Spatial)		MuSACLE (Intensity)		OctSqueeze (Spatial)		Draco (Spatial)		MPEG Achor (Spatial + Intensity)		MPEG TMC13 (Spatial + Intensity)	
	Encoding (ms)	Decoding (ms)	Encoding (ms)	Decoding (ms)	Encoding (ms)	Decoding (ms)	Encoding (ms)	Decoding (ms)	Encoding (ms)	Decoding (ms)	Encoding (ms)	Decoding (ms)
11	65.6	279.5	20.8	49.0	16.5	93.0	25.8	15.3	59.6	33.6	1219.0	272.0
12	94.4	477.5	31.4	89.9	26.2	165.7	29.0	16.2	84.2	50.3	1906.9	461.4
13	135.6	764.0	42.7	137.9	40.1	299.4	30.3	17.2	106.5	67.2	2539.0	674.4
14	177.1	1084.2	49.2	162.6	53.4	486.4	30.7	17.5	125.5	84.1	3038.9	839.1
15	202.4	1352.1	50.7	166.3	63.6	699.0	30.9	17.6	143.0	98.9	3406.3	925.8
16	208.4	1569.4	55.5	165.6	65.4	902.3	31.1	17.6	160.2	112.8	3695.1	985.0

Table 1: Encoding/decoding time (ms) of MuSACLE spatial intensity model, OctSqueeze, Draco, MPEG Anchor, and MPEG TMC13.

18 **R2, R3: Experiment details.** The average number of points per LiDAR point cloud is 80,156 in UrbanCity and
 19 120,402 in SemanticKITTI. During training, our full model requires 6-10GB of GPU memory. 16 GPUs was our
 20 starting point, but we retrained models with 4 GPUs and found roughly equal convergence / test metrics.

21 **R1: Clarification of line 138 and bitrates.** The $\sigma: \mathbb{R}^3 \rightarrow \mathbb{R}^d$ in line 138 is an MLP that predicts d -dimensional
 22 weights for continuous convolution. Bitrate is the total number of bits used to store the LiDAR dataset divided by the
 23 total number of points it contains. We do not count the one-time transmission of network weights (standard practice in
 24 learned compression) since this is negligible compared to the size of long LiDAR streams, e.g., 1 hour. We emphasize
 25 that the middle and right columns of Fig. 2 depict spatial-only bitrates since they compare spatial reconstruction quality.

26 **R4: Typos and clarification of line 71.** Thanks, we will fix the typos. Regarding line 71, our quantization scheme
 27 may produce a point cloud with fewer points than the original point cloud. Thus, we determine the intensity value of
 28 each point in the quantized point cloud by taking that of its nearest neighbour in the original point cloud.

29 **R4: Simpler entropy models.** Our tree-based baselines use a spectrum of statistical models with varying complexity.
 30 For example, Octree and MPEG Anchor use an empirical histogram-based entropy model, MPEG TMC13 uses a
 31 cascade of adaptive frequency tables, and OctSqueeze uses a deep learning model.

32 **R4: Invariance to scale, density, and tree depth.** Thanks for the exciting insight! We agree that ideally our density
 33 estimation model based on tree-structured message passing and continuous convolutions should be invariant to scale,
 34 density, and tree depth, echoing R4’s thoughts on potential applications to other perception tasks. However, in our
 35 current implementation, several mechanisms make this not the case: 1) our input feature include information on the
 36 depth; 2) we use separate continuous convolution modules per depth; and 3) we do not use data augmentation (e.g.,
 37 random scaling) and so our network learns the inductive bias of the dataset (e.g., the real-world scale). That said, by
 38 modifying these mechanisms, we should get a density model that is invariant to scale, density, and tree depth.

¹MPEG TMC13: github.com/MPEGGroup/mpeg-pcc-tmc13, MPEG TMC2: github.com/MPEGGroup/mpeg-pcc-tmc2.