# *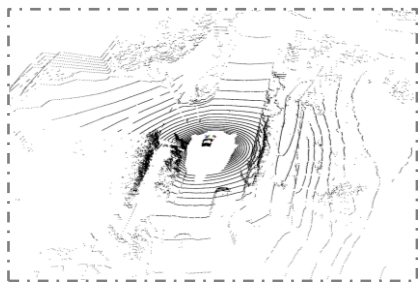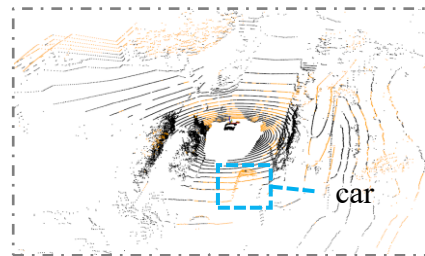Semantic-driven Cross-modal Contrastive Learning*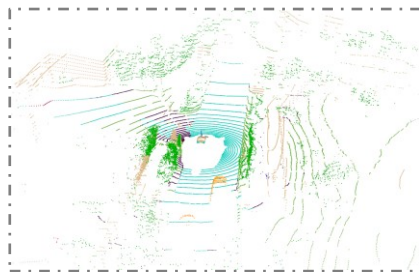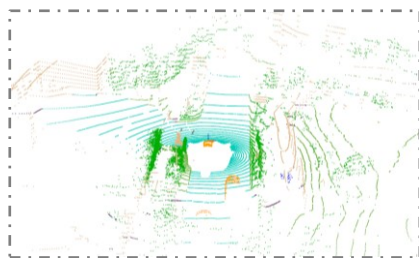