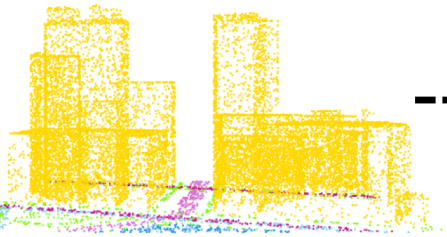


Data Input



Feature Encoding

Map Tile
Encoder

Point Cloud
Encoder

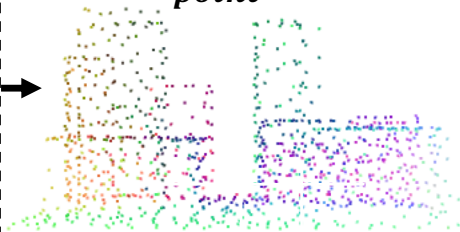
Panorama
Encoder

Late Stage with Feature Fusion and Aggregation

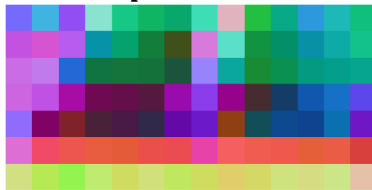
F_{tile}



F_{point}



F_{pano}



Pixel-to-Point Fusion

Max Pooling

F are feature maps
 f are global feature vectors
 L are InfoNCE loss

f_{map}

SAFA

f_{pano}

Neural Feature Alignment

$$L = L_{pano} + L_{map} + L_{cross}$$