

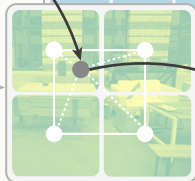
(a) Feature Query



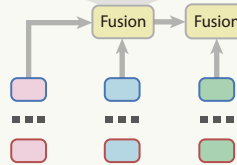
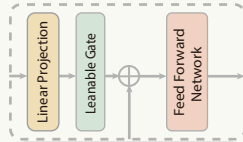
Input Image

ViT Encoder

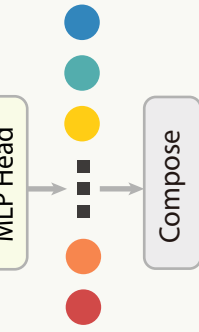
Reassemble



Bilinear Feature Interpolation

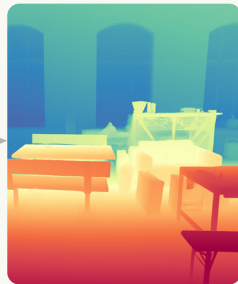


Queried Feature Tokens



Queried Depths

(b) Depth Decoding



Output Depth Map