

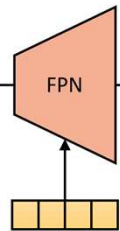
Input

Cross-modal Feature Pyramid

Instance Sequence Segmentation

Inference

- \oplus Element-wise Sum
- C Concatenation
- $*$ Dynamic Convolution



C

$*$

Box Head

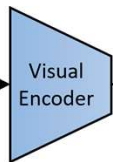
Mask Head

Class Head

Instance
Sequence
Matching

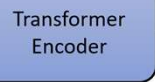


Backbone

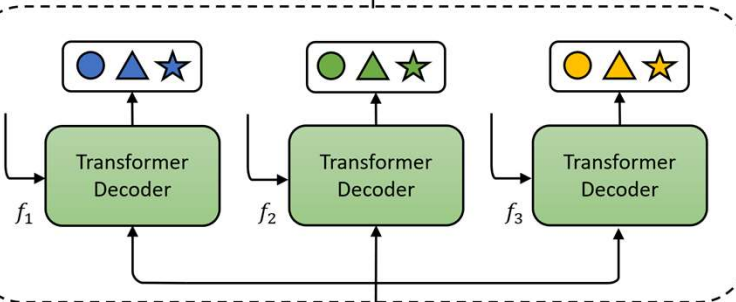


F_v

Stack

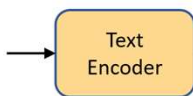


Initial
Conditional
Query



Transformer

a person skateboarding



F_e

$L \times C$

Pool &
Repeat

$N \times C$

