

Inputs (L, V)

"There should be a purple theater banner on your left. Go forward on this street until you come to the first traffic light. Make a right at the light. You should see silver gates on your left."



Linguistic
Embeddings

Visual
Embeddings

Main Model

Transformer Encoder

a)

Classifier

b)

FLPM Framework

Input (tr)



Path Trace
Embeddings

FLPM Framework

PM-VLN

g_{PMTP}

g_{PMF}

g_{Attn}

$e_l \oplus e'_v$

g_{CFns}

$(\widetilde{E}_\eta, U\eta)$

$g_{Class_{max} x_i}$

$\{Forward, Left, Right, Stop\}$

Cross-modal Prioritisation

Cross-modal Attention

Combine Outputs

Action Selection