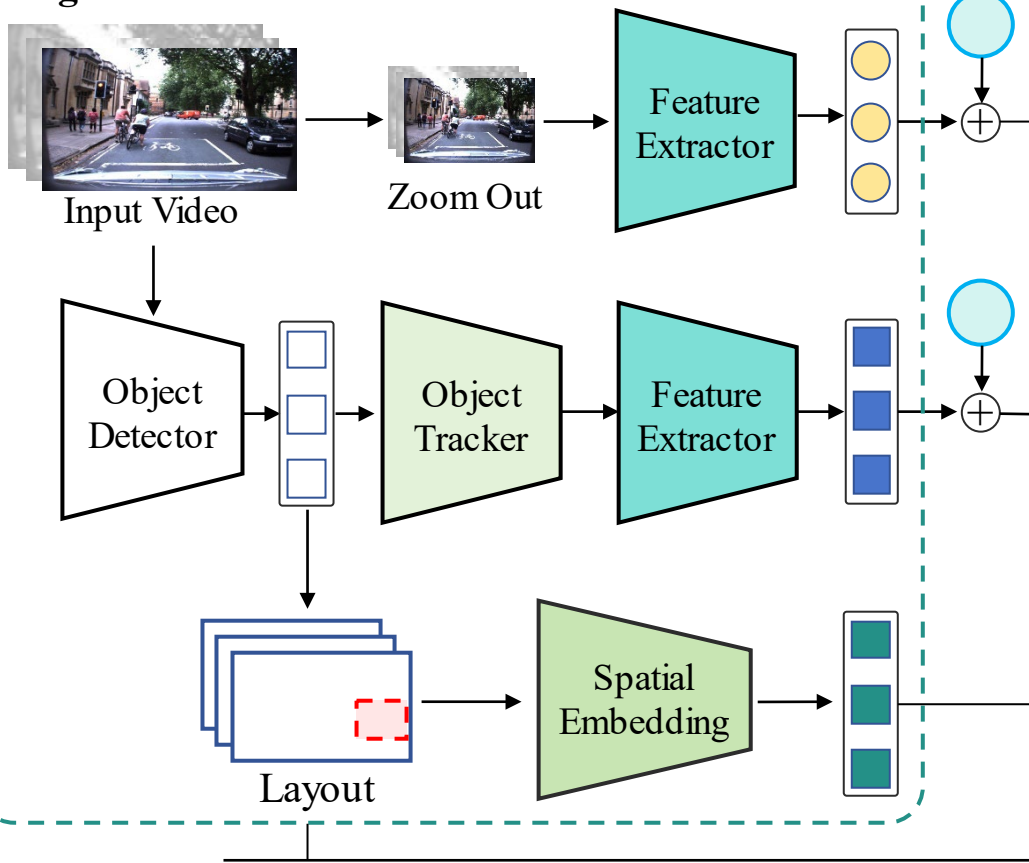
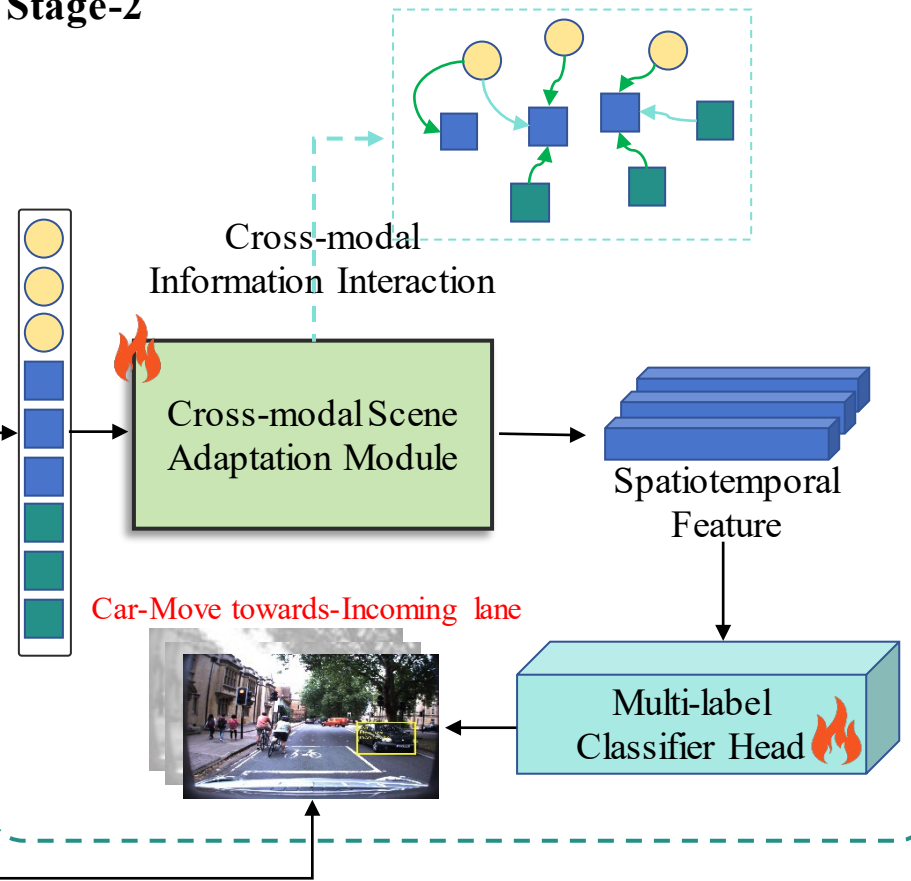


Stage-1



Stage-2



Coarse-grained
global context

Fine-grained
local information

Layout
Embedding

Positional
Embedding