



Entire Scene
 I

Global Feature
Extractor
(GFE)

F_G

δ_G

\mathbf{f}_G

Local Feature
Extractor
(LFE)

F_L

δ_L

\mathbf{f}_L

Intrinsic Feature
Extractor

F_I

δ_I

\mathbf{f}_I

Complete Features
 $\mathbf{f}_T \in \mathbb{R}^l$

Attributes vector
 $\hat{\mathbf{a}} \in \mathbb{R}^{n_c}$

Interpreter
 \mathcal{I}

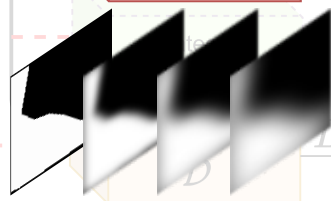
Cropped Image
 I_C

Mask
 M



Category Embedding
(e.g. floor)

$\hat{\mathbf{c}}$



D

A_G

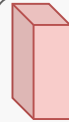
\mathcal{G}_G

A_L

\mathcal{G}_L

A_I

\mathcal{G}_I



Using Informed
Convolution
Layers



Concatenation



Feature-level
Multiplication &
Spatial Reduction



Attention Gate

$\hat{\mathbf{c}}$

Estimated from LFE by Category Estimator