$I_{t-k}, ..., I_{t-1}, I_t$

$\hat{I}_{t+1}$

$\mathcal{D}_a$

Relation Learning module

$\mathrm{E}_s$     R

$\mathcal{E}$

$\mathcal{RL}$   $\psi_r$

$\mathrm{w}_1$

$\mathrm{w}_2$   $\mathrm{e}_o$

$\mathcal{M}$

$\mathcal{D}_m$

$\psi_r$

$\tilde{I}_{t+1}$

$\otimes$: Element-wise multiplication    $\odot$: Dot production    $\mathcal{S}_{\mathrm{rl}}$: Relation score

$\tilde{I}$: Warped future frame          $\mathcal{E}$: Encoder          $\mathcal{M}$: Region mask set

$\hat{I}$: Predicted future frame     $\mathcal{D}$: Decoder        $\mathrm{E}_s$: Scene Embedding

$I$: Ground-truth frame          $\mathrm{w}$: Conv filter     $\psi_r$: Relation score map

$\mathcal{RL}$: Relation Learning module    $\mathrm{e}_o$: Object concensus embedding