

input size  $m$

normalization

$\ell = m$

relu

$\times d$

$\ell = k$

$\langle \text{value head} \rangle$

relu

$\ell = k$

relu

$\ell = 1$

denorm.

$V_{\theta}(\tilde{b})$

$\langle \text{policy head} \rangle$

relu

$\ell = k$

relu

$\ell = |\mathcal{A}|$

softmax

$P_{\theta}(\tilde{b}, \cdot)$

$\langle \text{failure head} \rangle$

relu

$\ell = k$

relu

$\ell = 1$

sigmoid

$F_{\theta}(\tilde{b})$