

input size =  $2k + 2$

normalization

$\ell = 2k + 2$

fc  $\Rightarrow$  batch norm

$\ell = 128$

relu

$\ell = 128$

dropout

$\times 3$

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

fc  $\Rightarrow$  batch norm

$\ell = 128$

relu

$\ell = 128$

dropout

$\ell = 128$

relu

$\ell = 1$

denorm.

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

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

fc  $\Rightarrow$  batch norm

$\ell = 128$

relu

$\ell = 128$

dropout

$\ell = 128$

relu

$\ell = |\mathcal{A}| \in \{20, 25\}$

softmax

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