

Bi-modal DBM

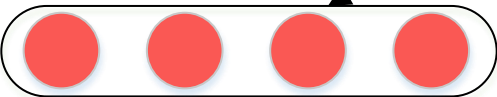
$h^{(3)}$ (Joint Representation)



$W_a^{(3)}$

$W_t^{(3)}$

$h_a^{(2)}$



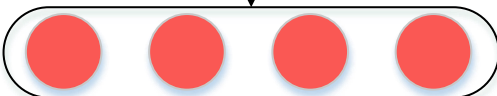
$h_t^{(2)}$



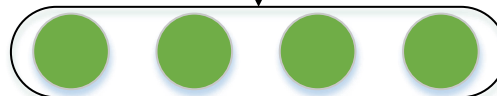
$W_a^{(2)}$

$W_t^{(2)}$

$h_a^{(1)}$



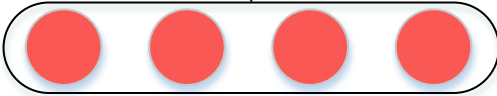
$h_t^{(1)}$



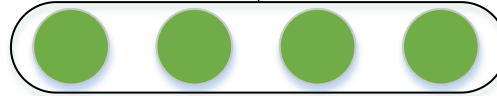
$W_a^{(1)}$

$W_t^{(1)}$

v_a



v_t



Text-modality DBM

Audio-modality DBM