



Small logit margin  $\Delta$   
 $\Rightarrow$  diffuse attention.

$\Delta$  grows to  $\log(T) \Rightarrow$  attention  
concentrates on needle.

$\nabla_{q_i} \ell \propto \mu - t^* \Rightarrow$  gradient descent moves  $q$  toward  $t^*$  (needle) and away from  $\mu$ .