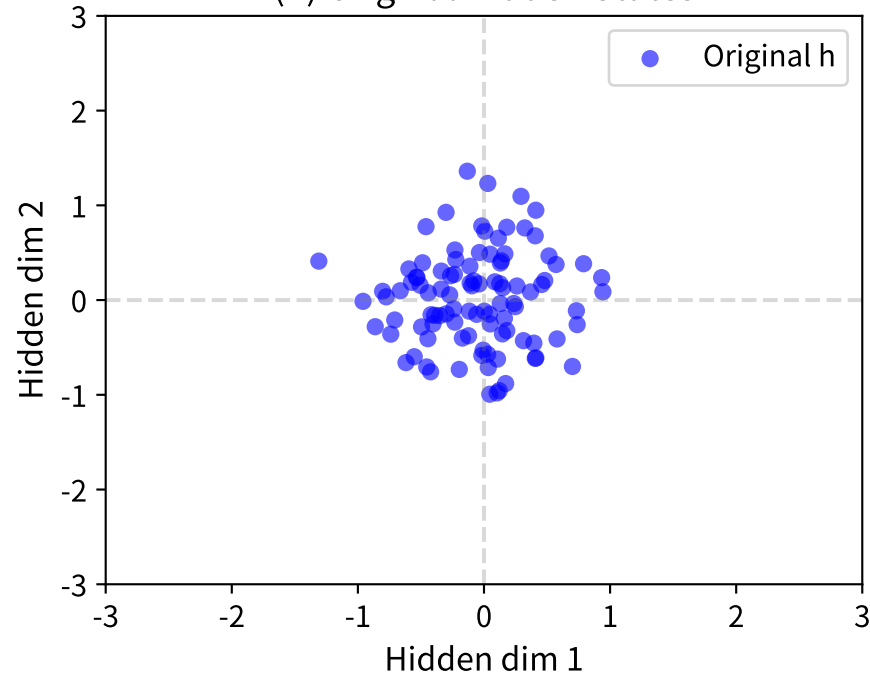
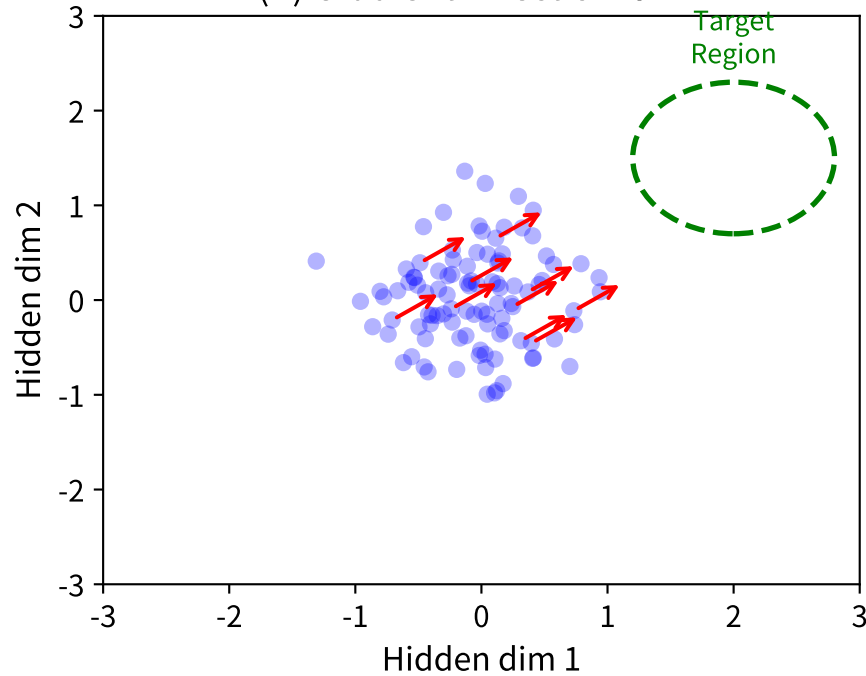


(A) Original Hidden States



(B) Gradient Direction  $\nabla h L$



(C) Perturbed States  $h' = h + \alpha \nabla h$

