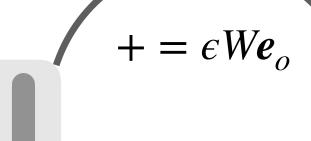
output word embeddings  $e_o$ 



adapted output word embeddings  $e'_{o}$ 



Language Model Hidden Layers

original LM M

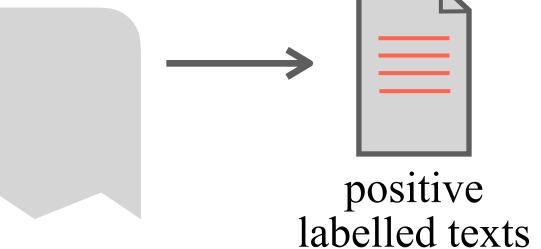
"switched" LM  $M(\epsilon W)$ 

Language Model Hidden Layers

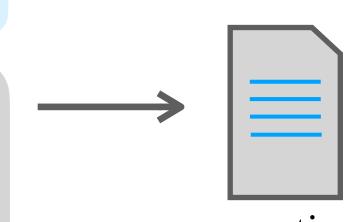
 $M(+\epsilon W)$ 

 $M(-\epsilon W)$ 

objective: maximize likelihood



objective: maximize likelihood

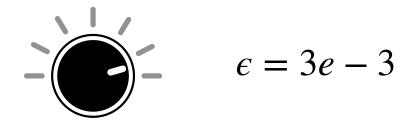


negative labelled texts

(b) Training

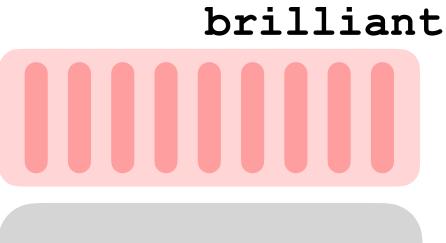
step 1:

setting a "switch" value



step 2: Plugging in and generate

my life is



(c) Generation

(a) LM-Switch overview