

METAPOISON: LEARNING TO CRAFT POISON

**W. Ronny Huang,* Jonas Geiping,*
Liam Fowl,^ Tom Goldstein**

*Equal Contribution
^Speaker

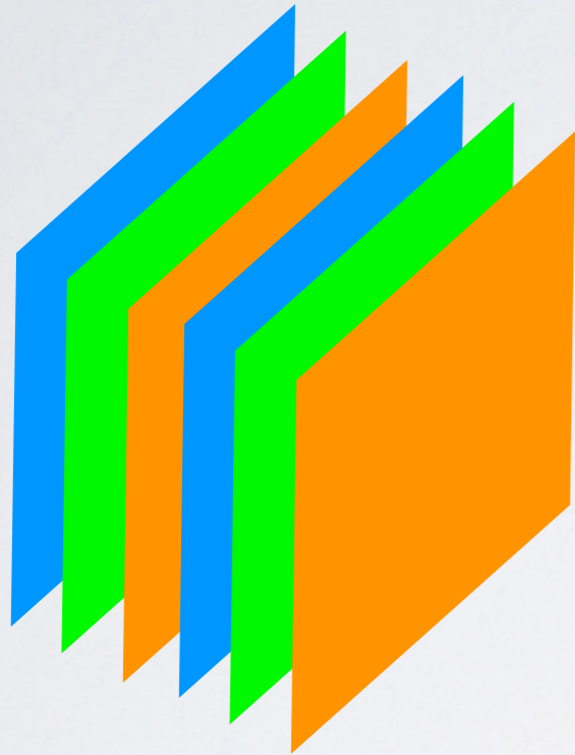
University of Maryland



NeurIPS MetaLearn 2019

DATA POISONING

Training data



Testing example

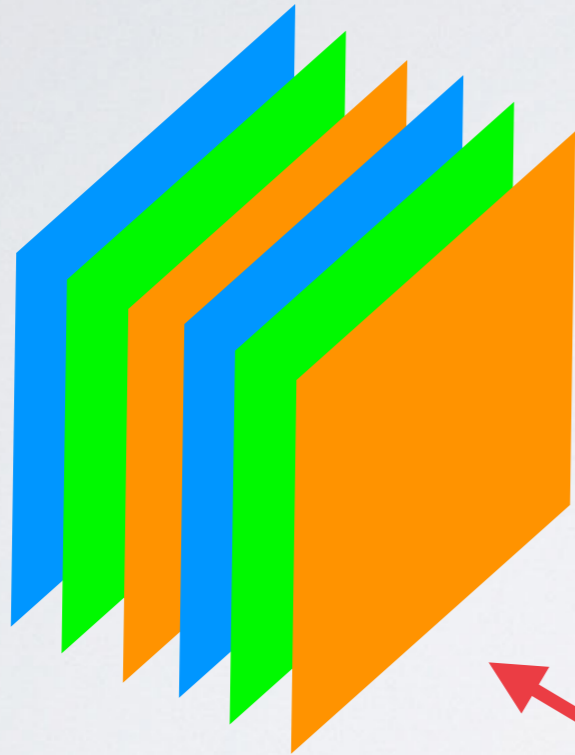


Base



DATA POISONING

Training data



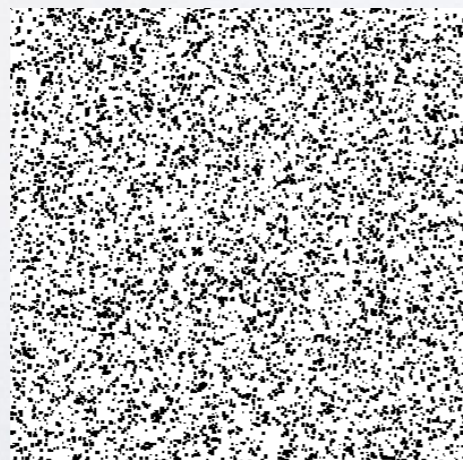
Testing example



Base



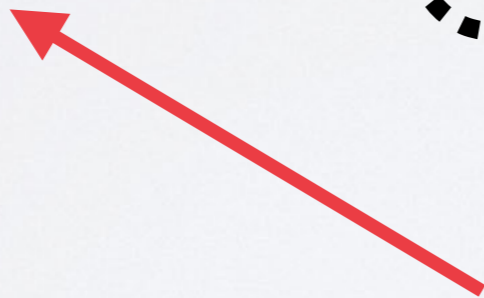
+



=

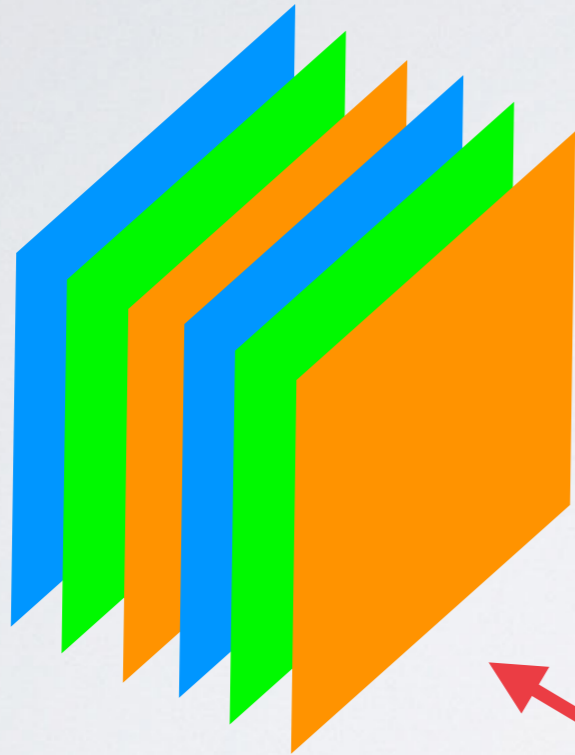


Poison!



DATA POISONING

Training data



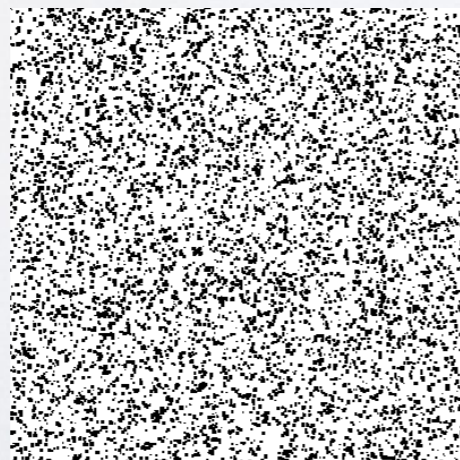
Testing example



Base



+



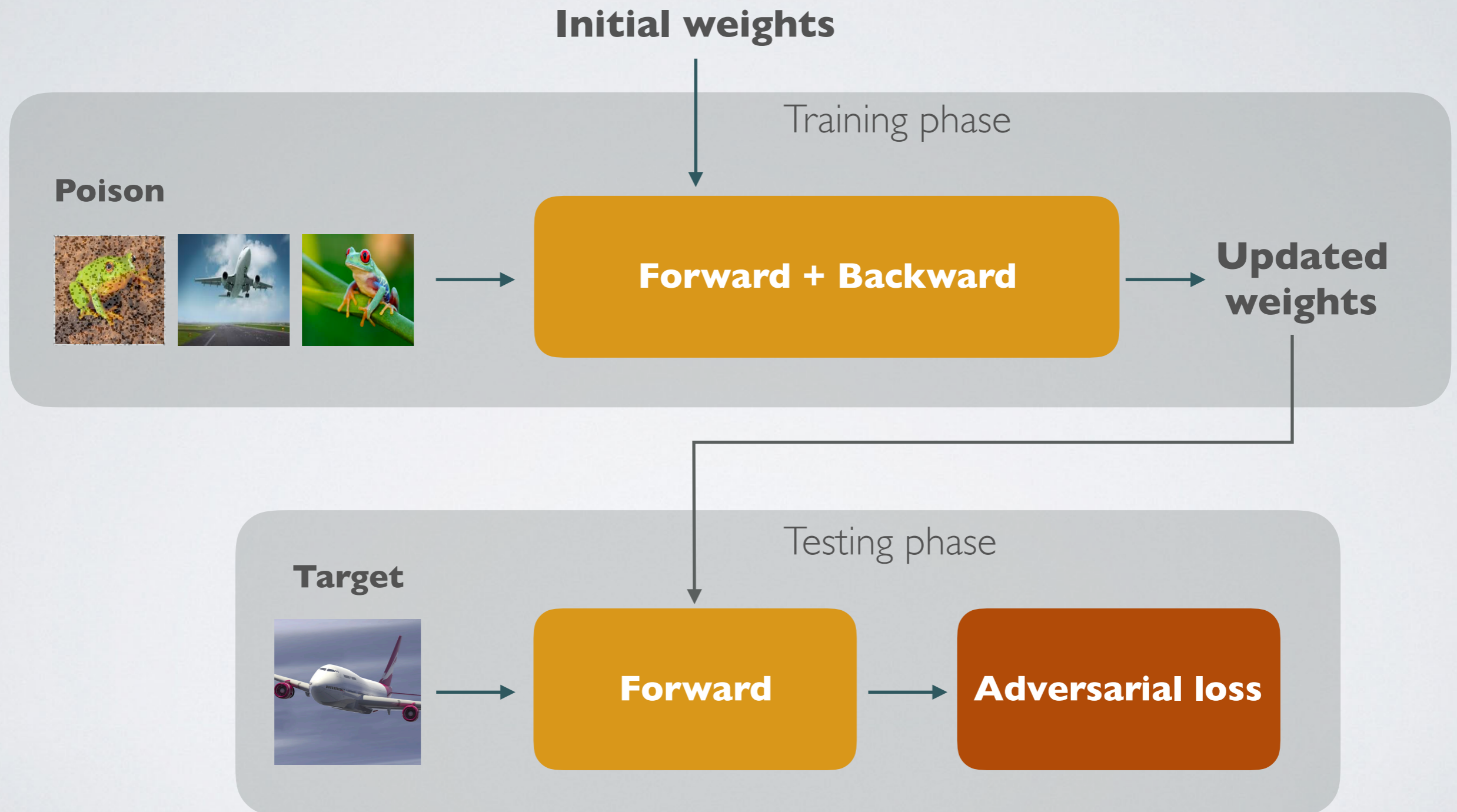
=



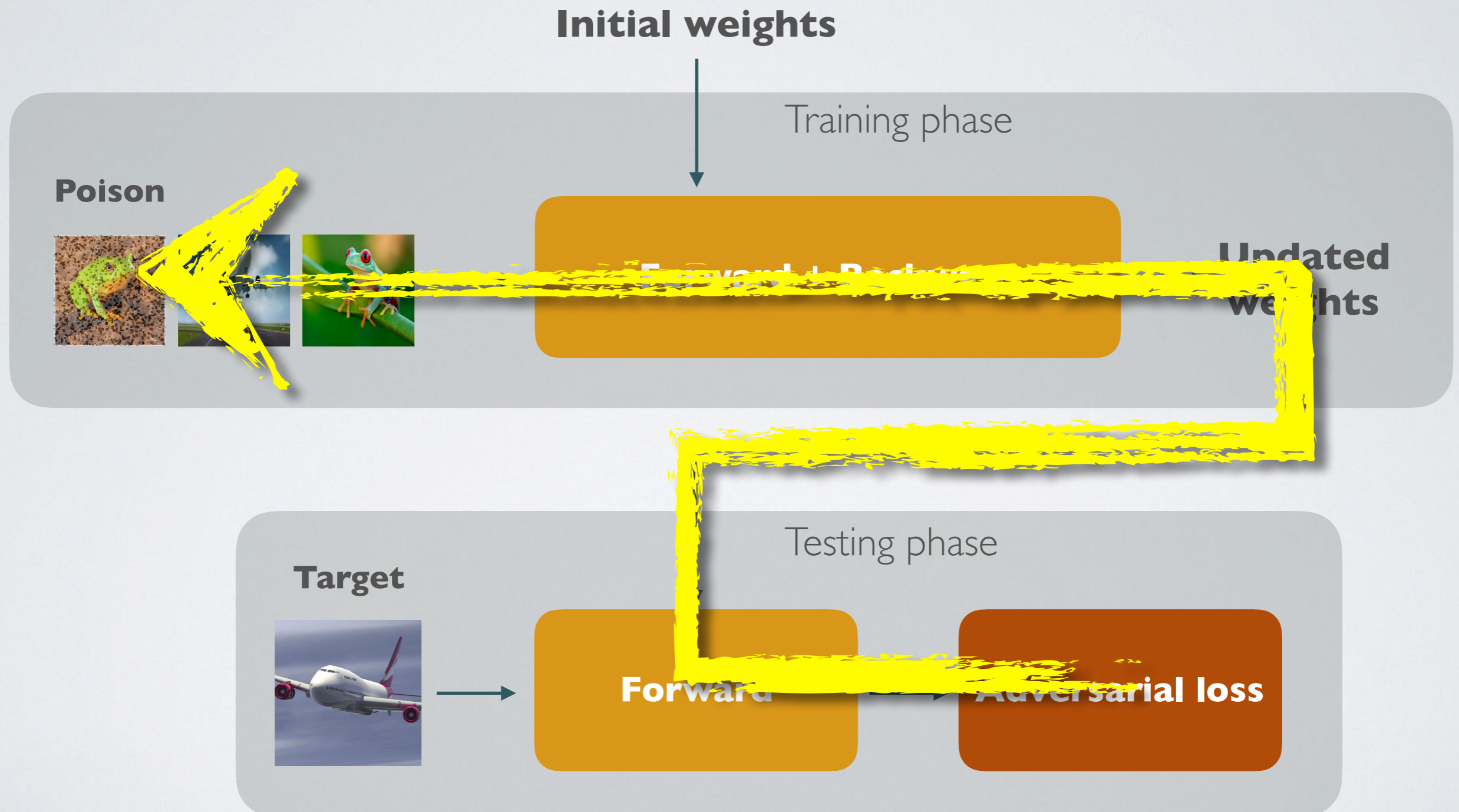
Poison!



LEARNING TO CRAFT

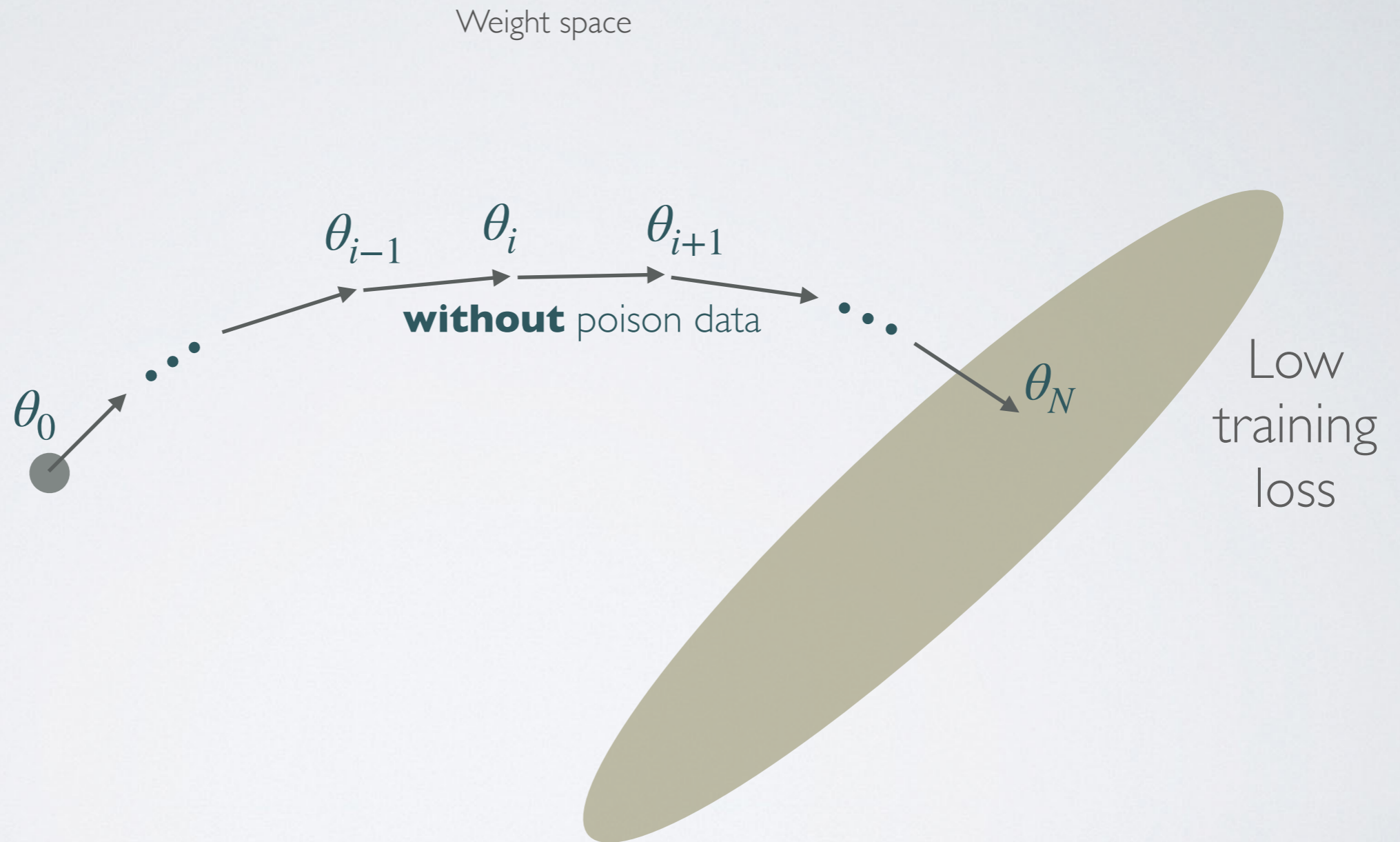


LEARNING TO CRAFT

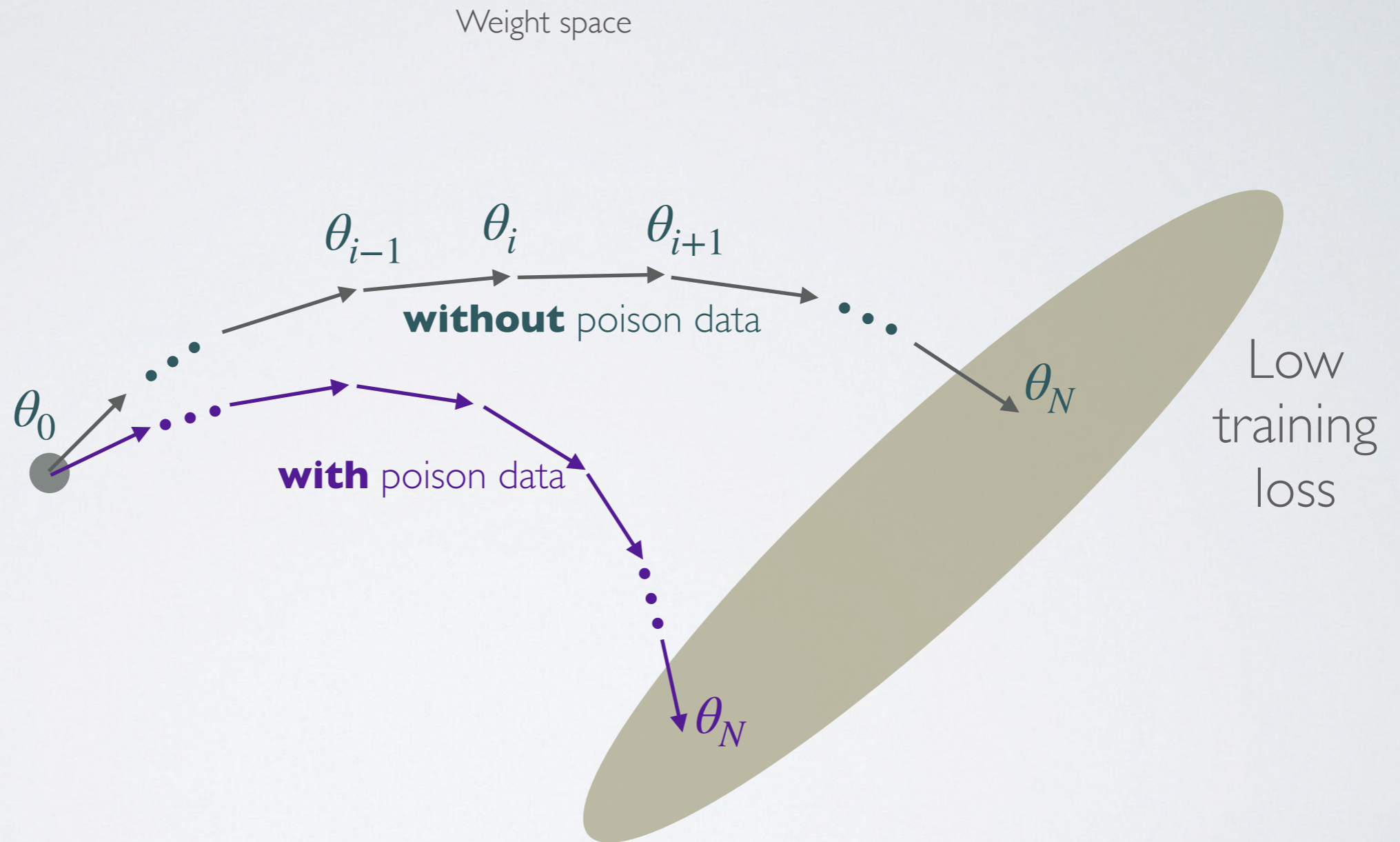


Backprop to the poison!

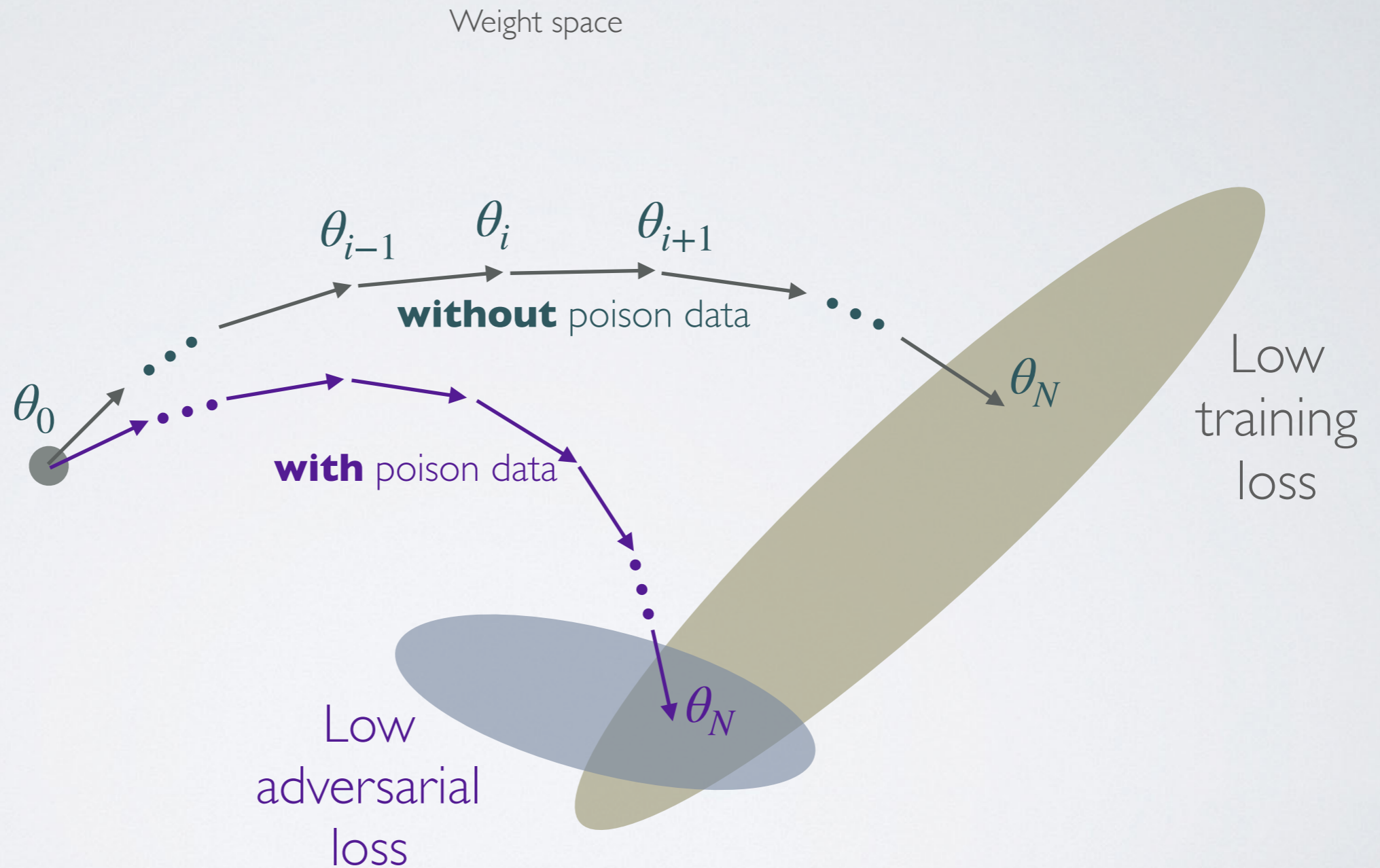
POISONED TRAINING DYNAMICS

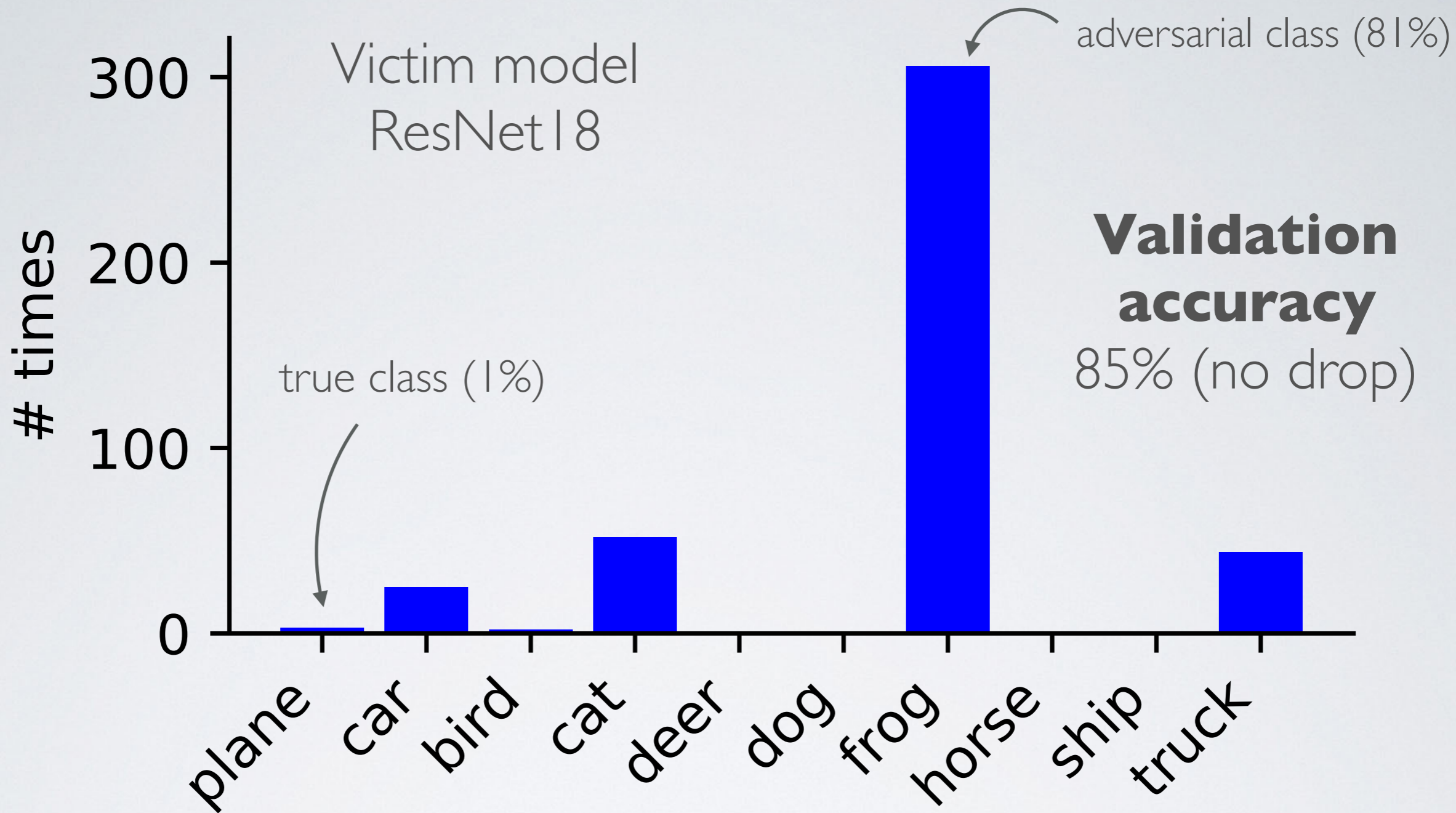


POISONED TRAINING DYNAMICS



POISONED TRAINING DYNAMICS





5000 poisons (10%)



cause



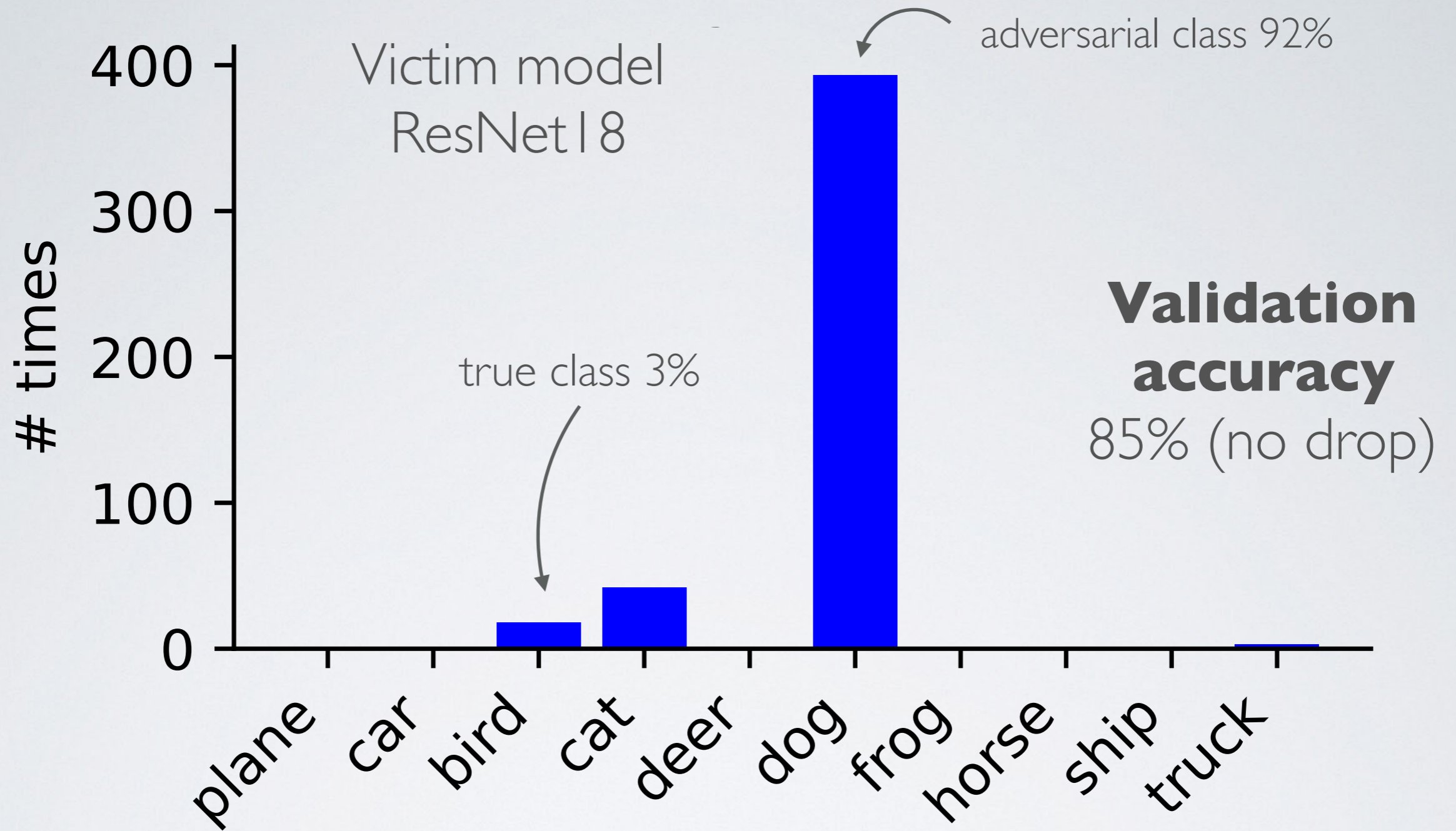
Target



classified

as





5000 poisons (10%)



cause



Target



classified

as

