# Transferability Attack

Target model with unknown weights, machine learning algorithm, training set; maybe non-differentiable

Substitute model mimicking target model with known, differentiable function

Train your own model

Adversarial crafting against substitute

Adversarial examples

Deploy adversarial examples against the target; transferability property results in them succeeding