

Recent Advances in Adversarial Machine Learning

Nicholas Carlini
Google Research

Recent Advances in Adversarial (Examples in) Machine Learning

Nicholas Carlini
Google Research

The Year is 2014

Someone tells you they have a new algorithm to generate human faces

The Year is 2014

*"more results of how
this helps on real tasks
or real datasets"*



*"the theoretical work
is primitive, and the
experiments are pretty
basic."*

The Year is 2017

Someone tells you they have a new algorithm to generate human faces

The Year is 2017





The Year is 2013

Someone tells you they have discovered
a flaw in the robustness of neural networks

The Year is **2013**

3

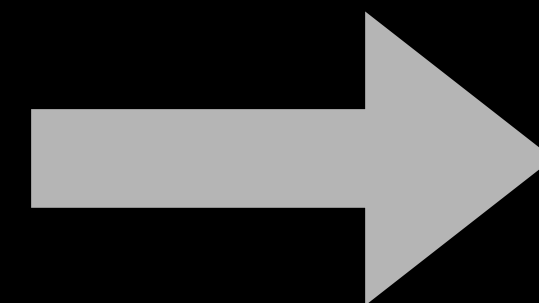
The Year is 2019

Someone tells you they have discovered
a flaw in the robustness of neural networks

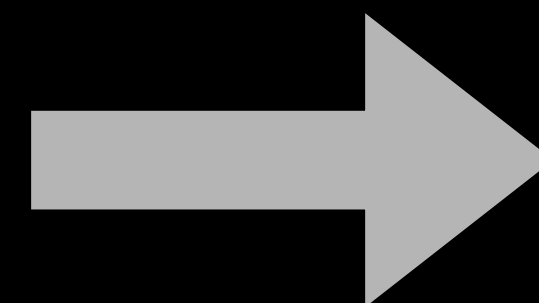
The Year is **2019**

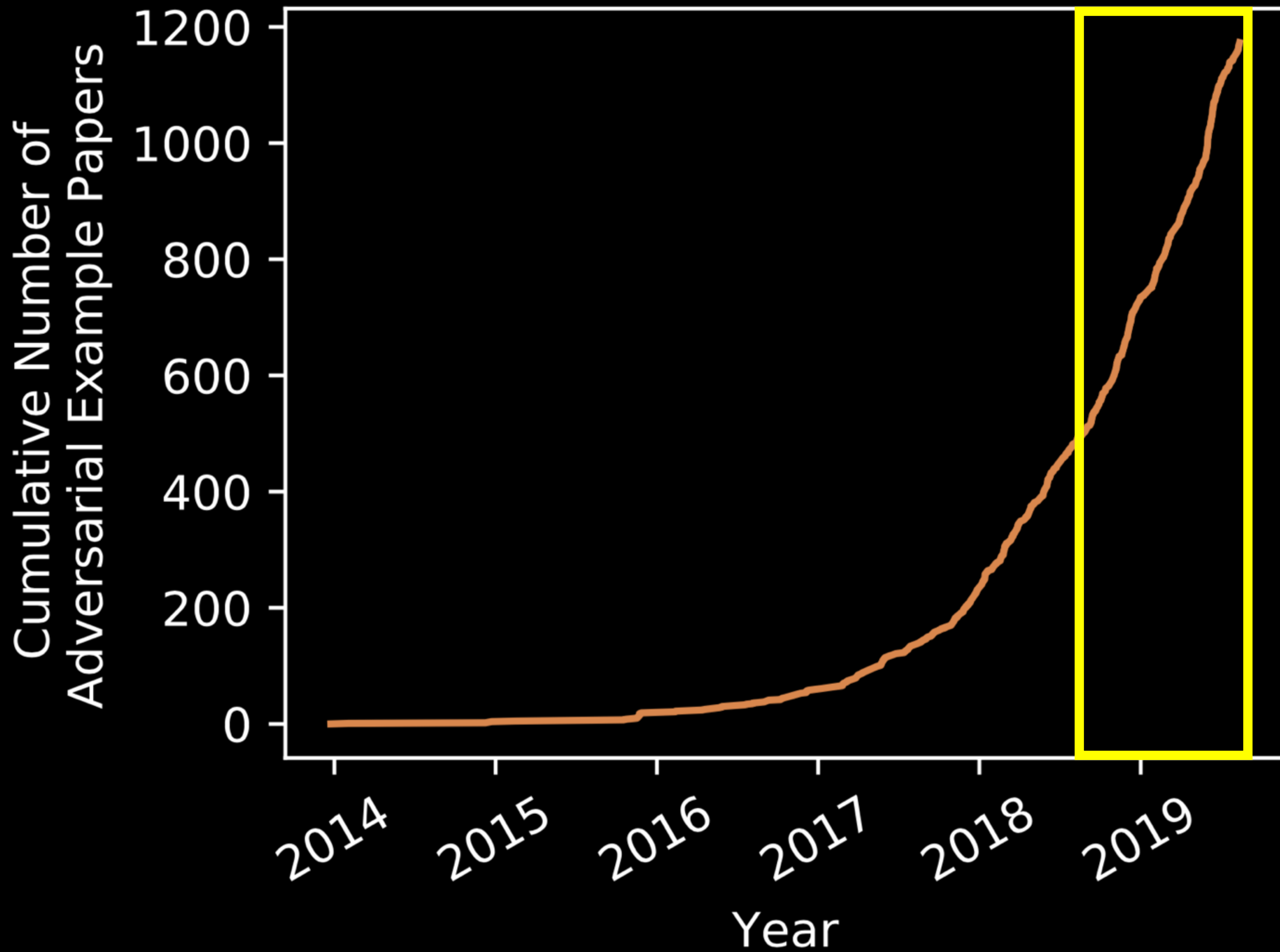
3

3 years:



6 years:





Background: Adversarial Examples



88% **tabby cat**



adversarial
perturbation



88% **tabby cat**



adversarial
perturbation



88% **tabby cat**



adversarial
perturbation

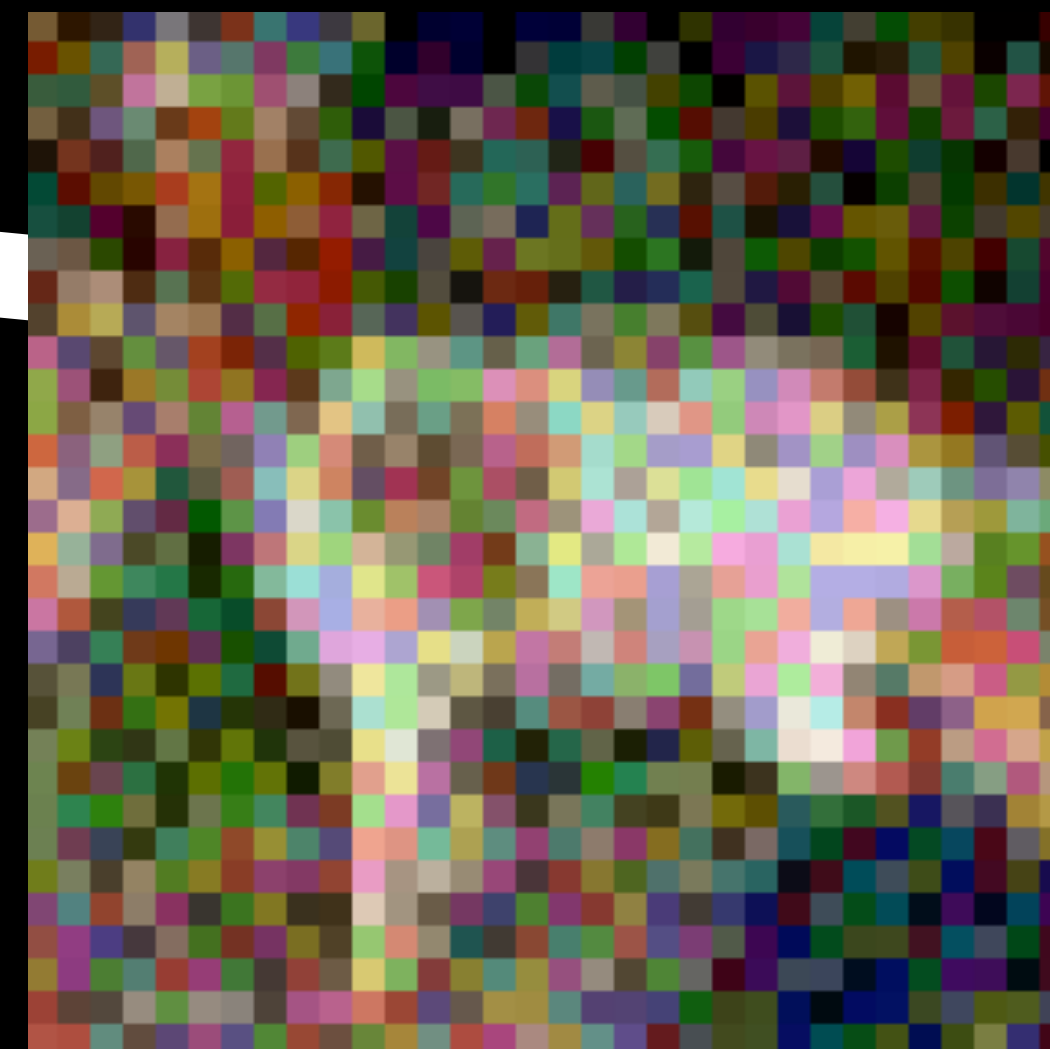
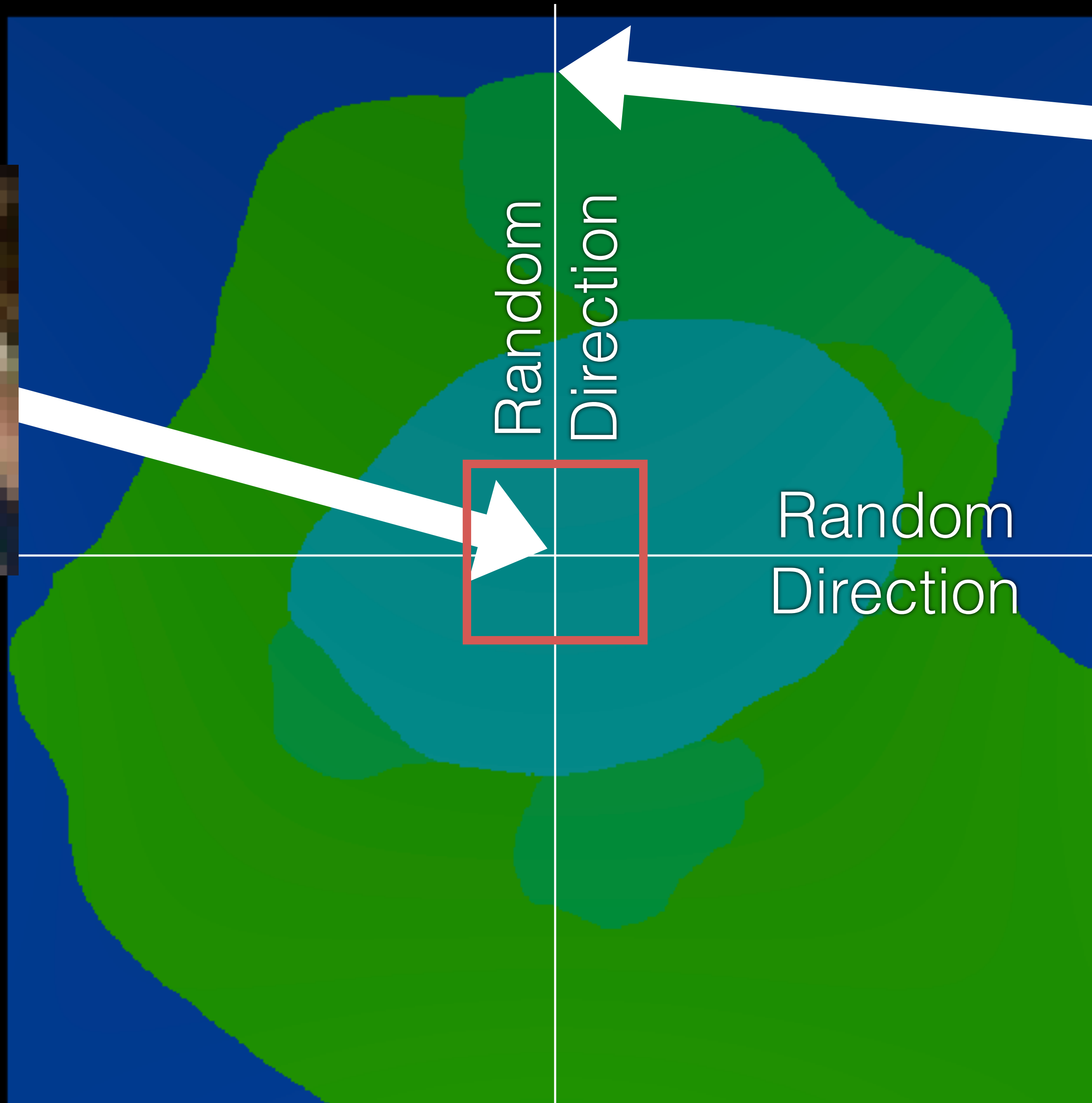


88% **tabby cat**

99% **guacamole**



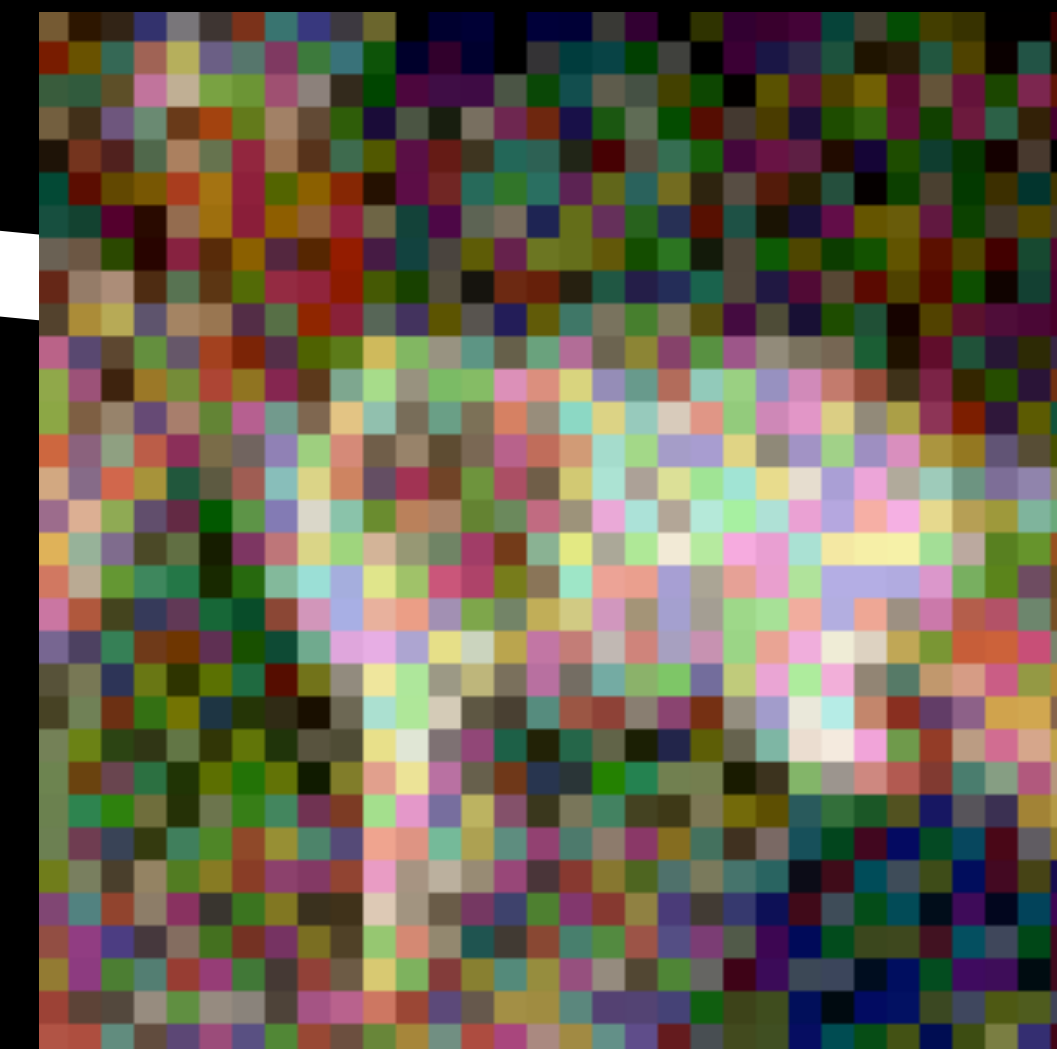
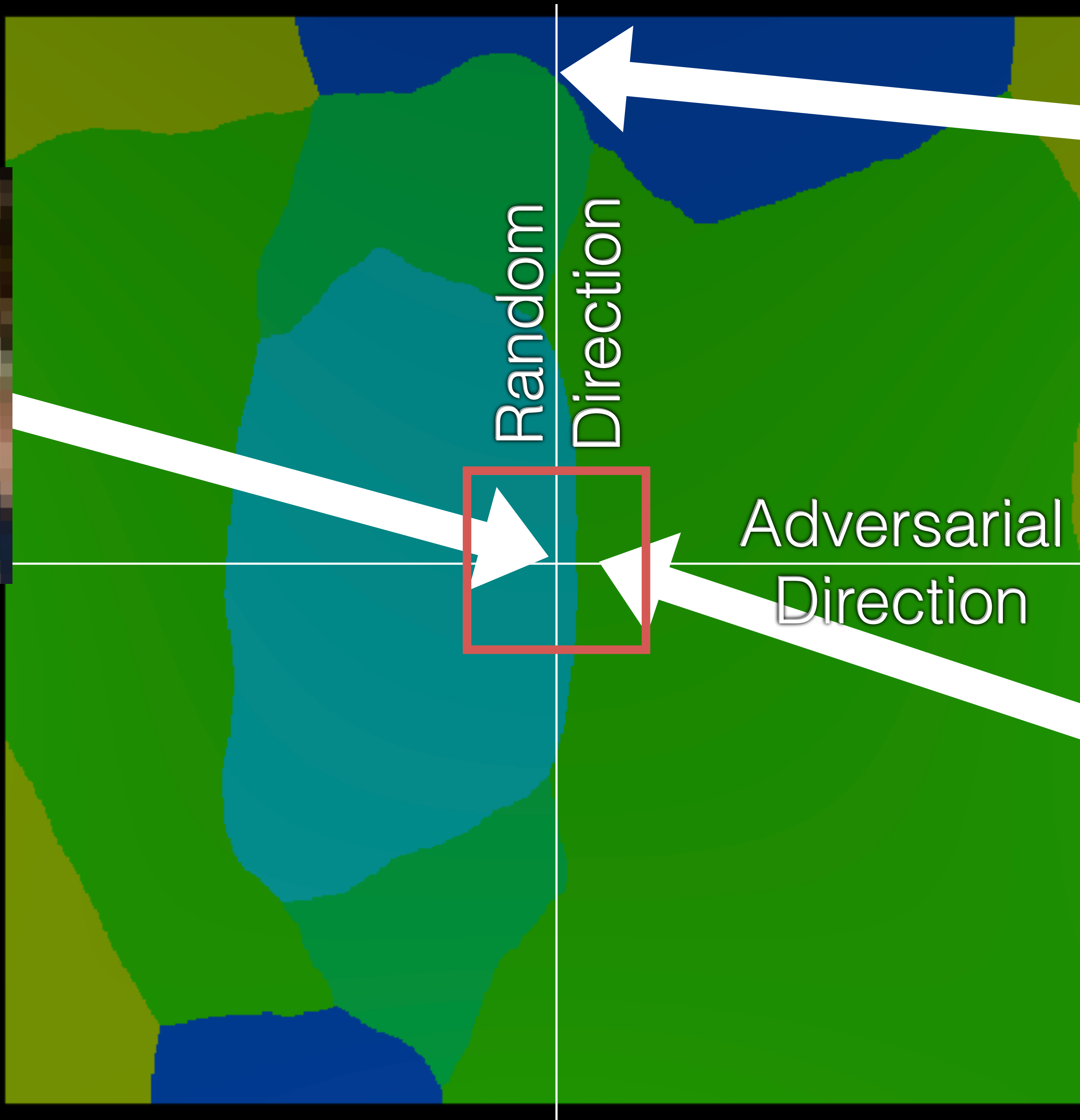
Dog



Truck



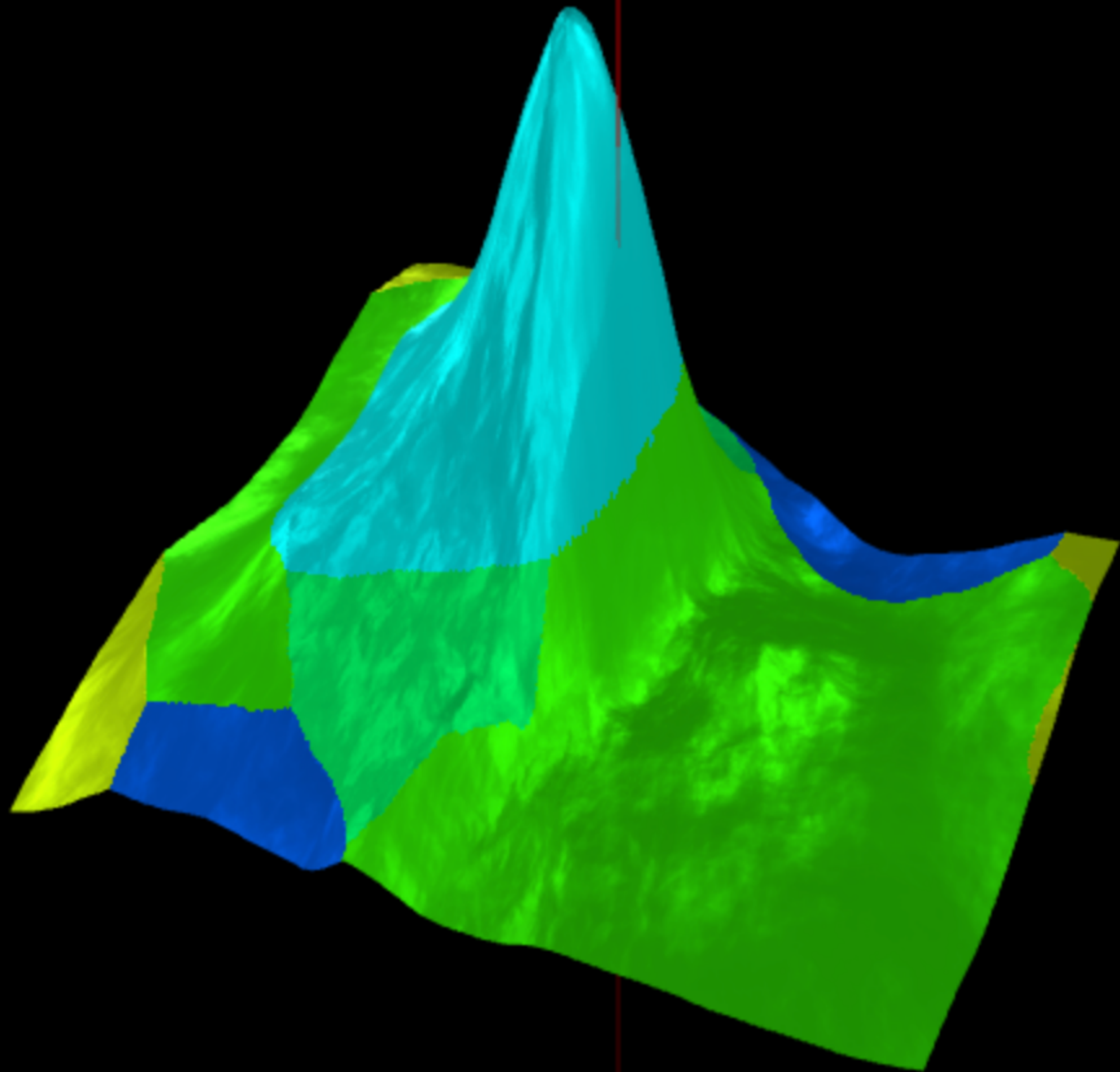
Dog

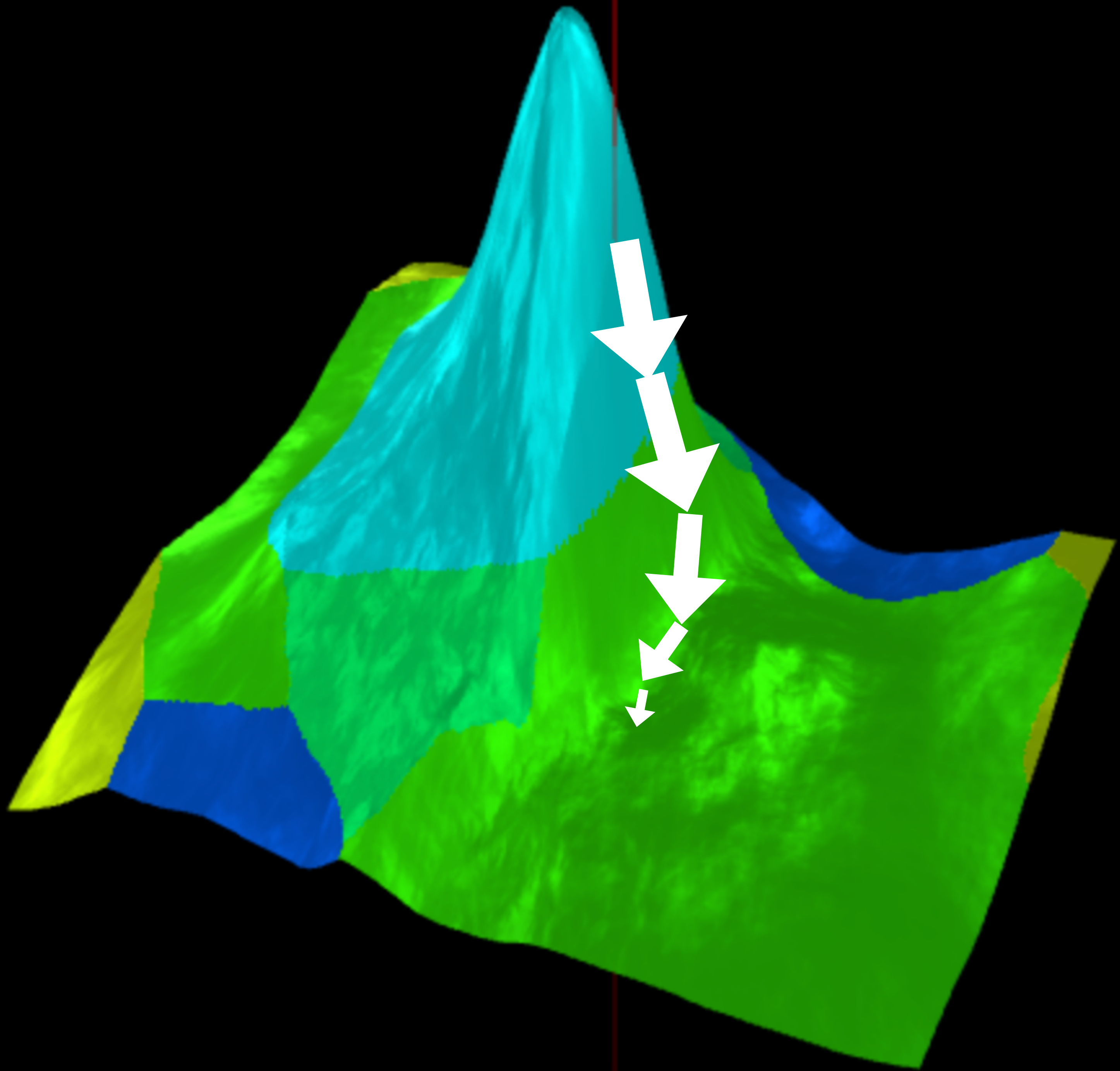


Truck



Airplane





Recent advances in ...

Generating

Adversarial Examples

DECISION-BASED ADVERSARIAL ATTACKS: RELIABLE ATTACKS AGAINST BLACK-BOX MACHINE LEARNING MODELS

Wieland Brendel*, Jonas Rauber* & Matthias Bethge

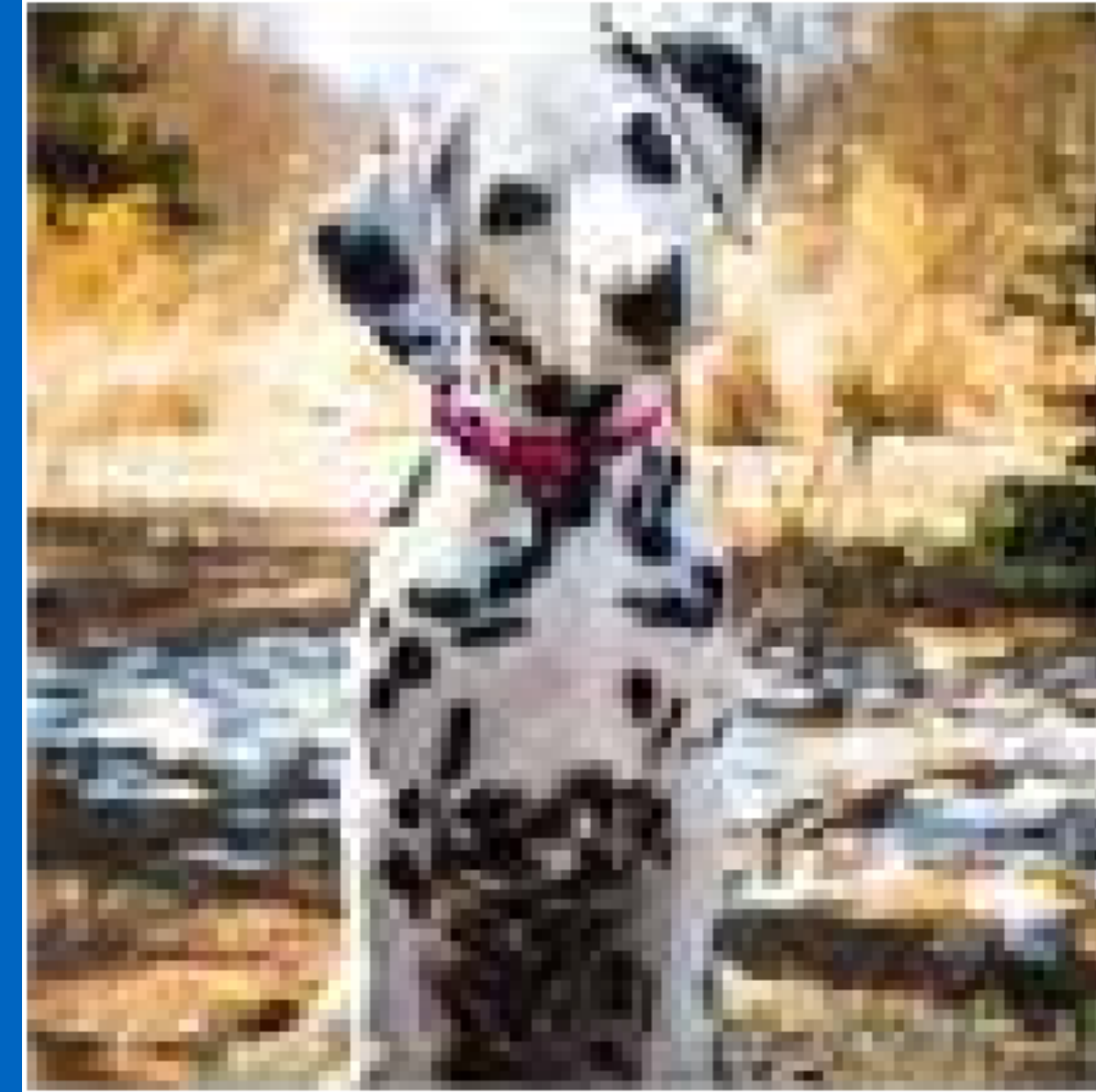
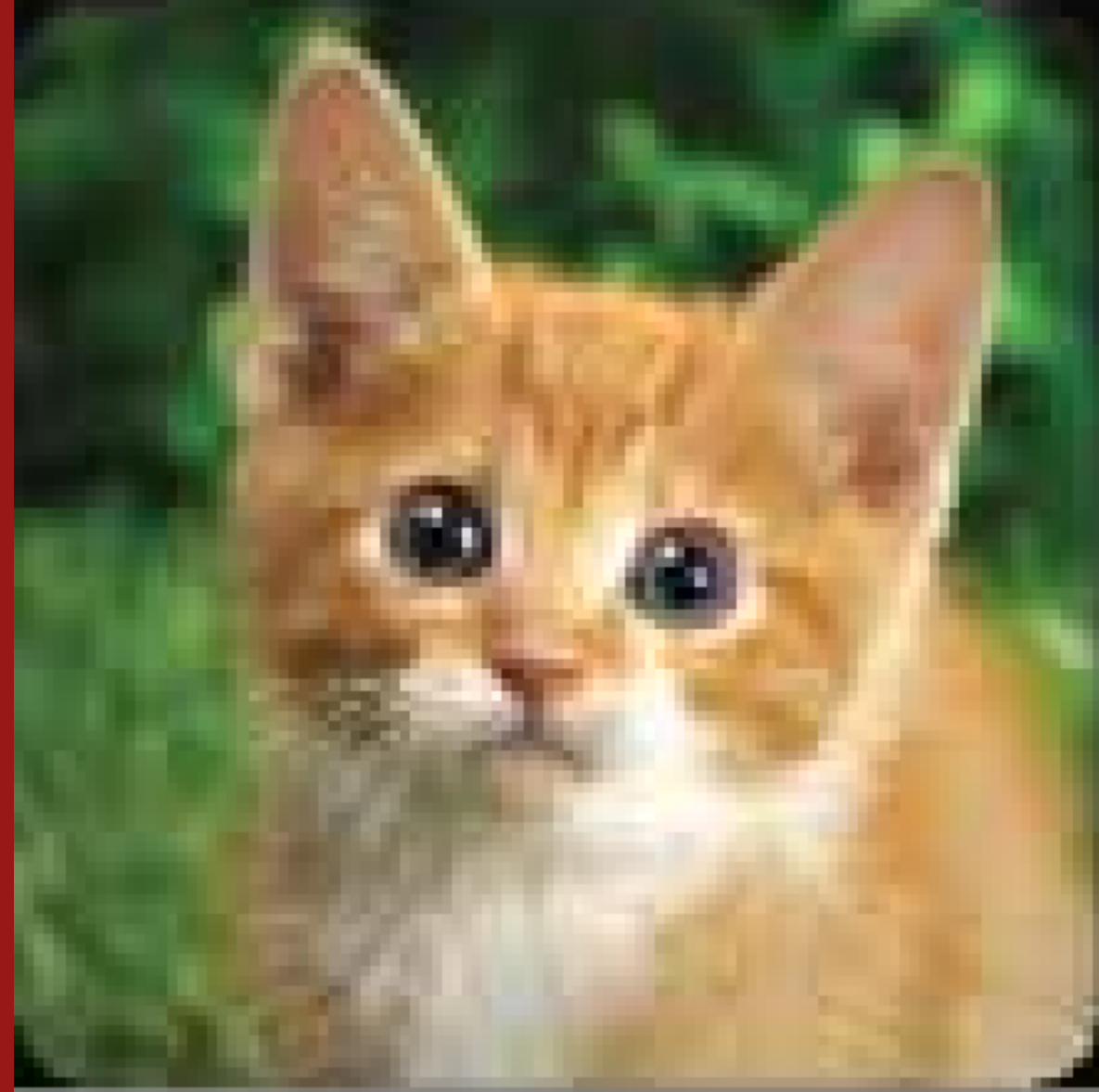
Werner Reichardt Centre for Integrative Neuroscience,

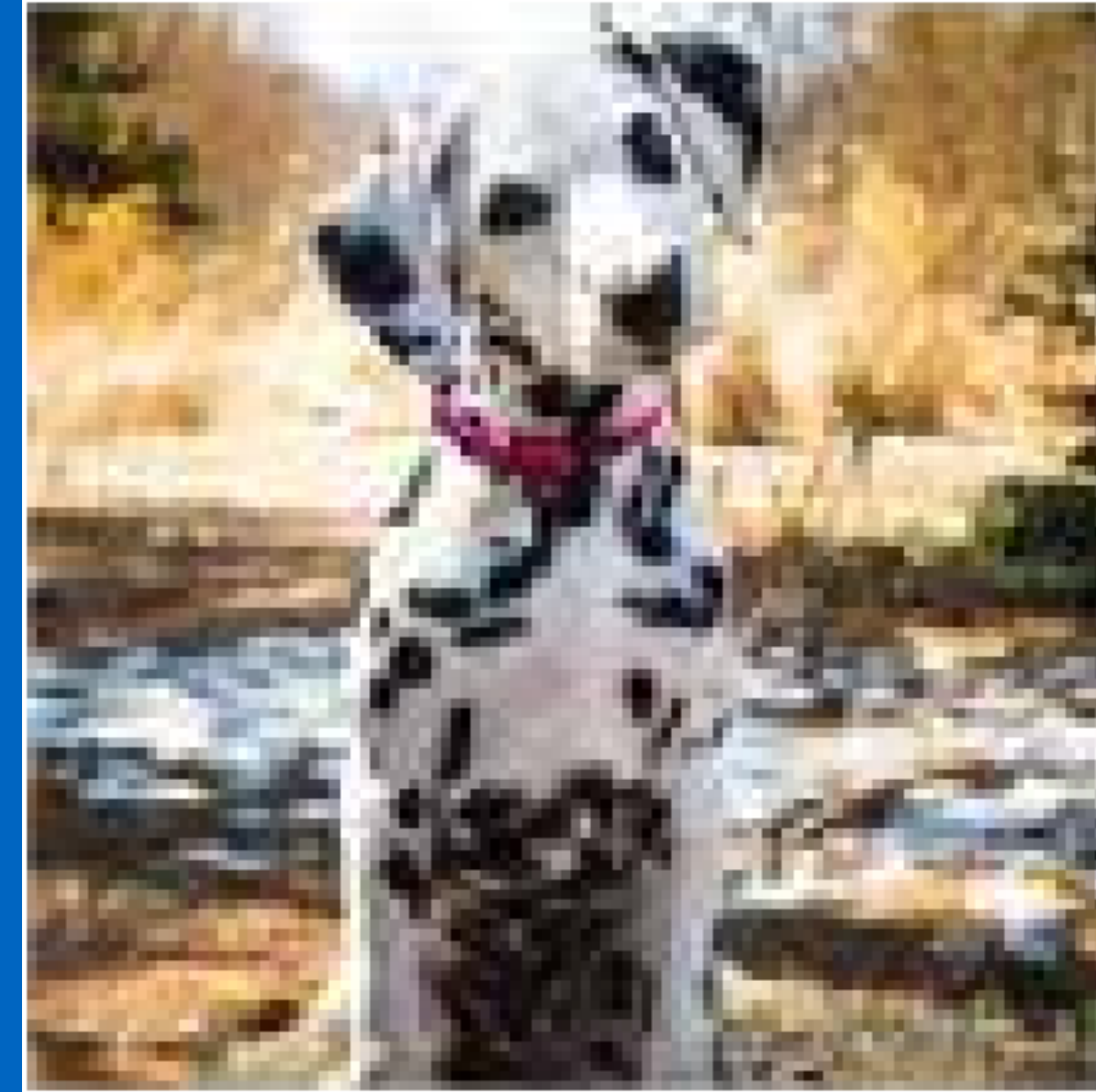
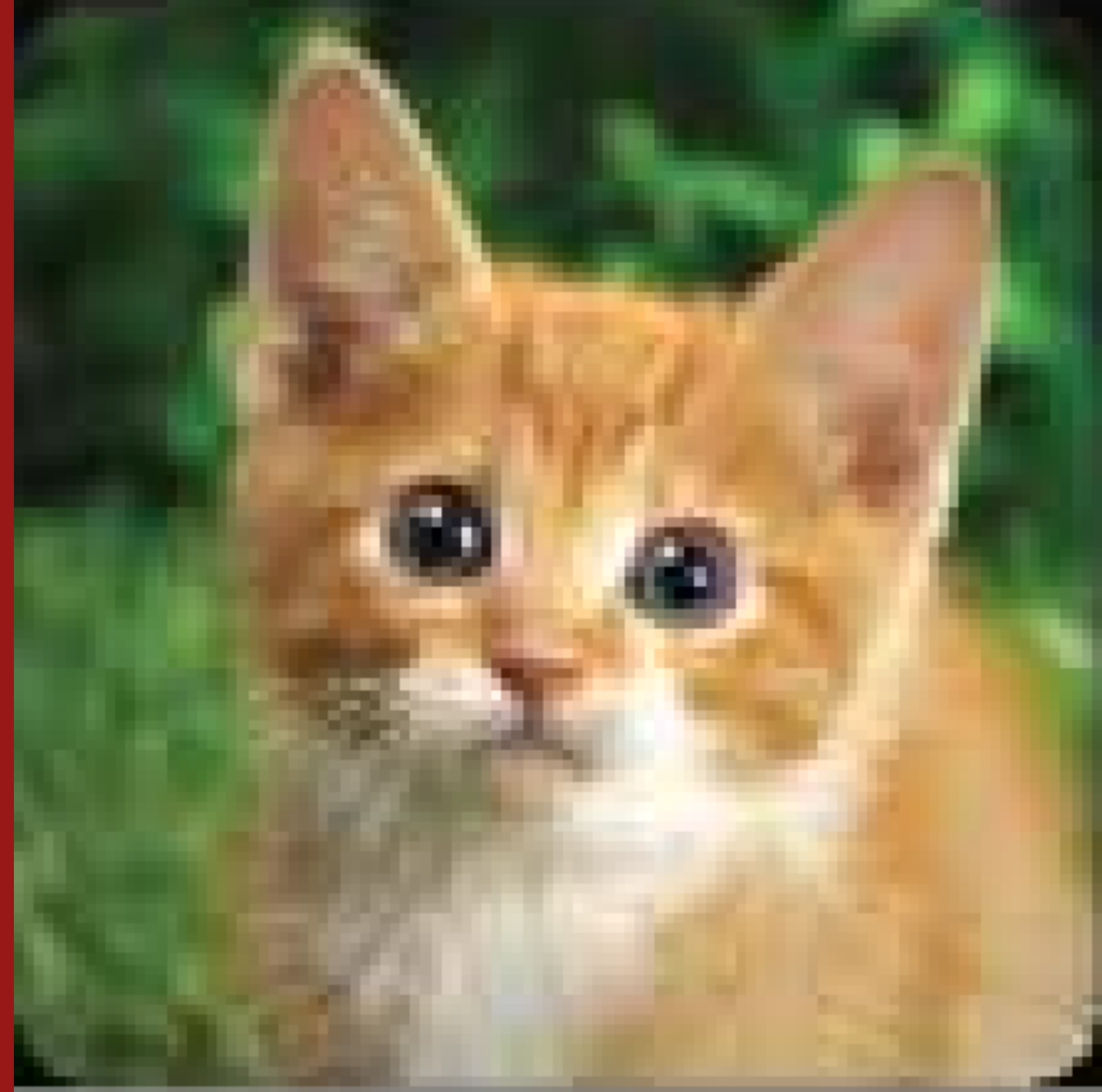
Eberhard Karls University Tübingen, Germany

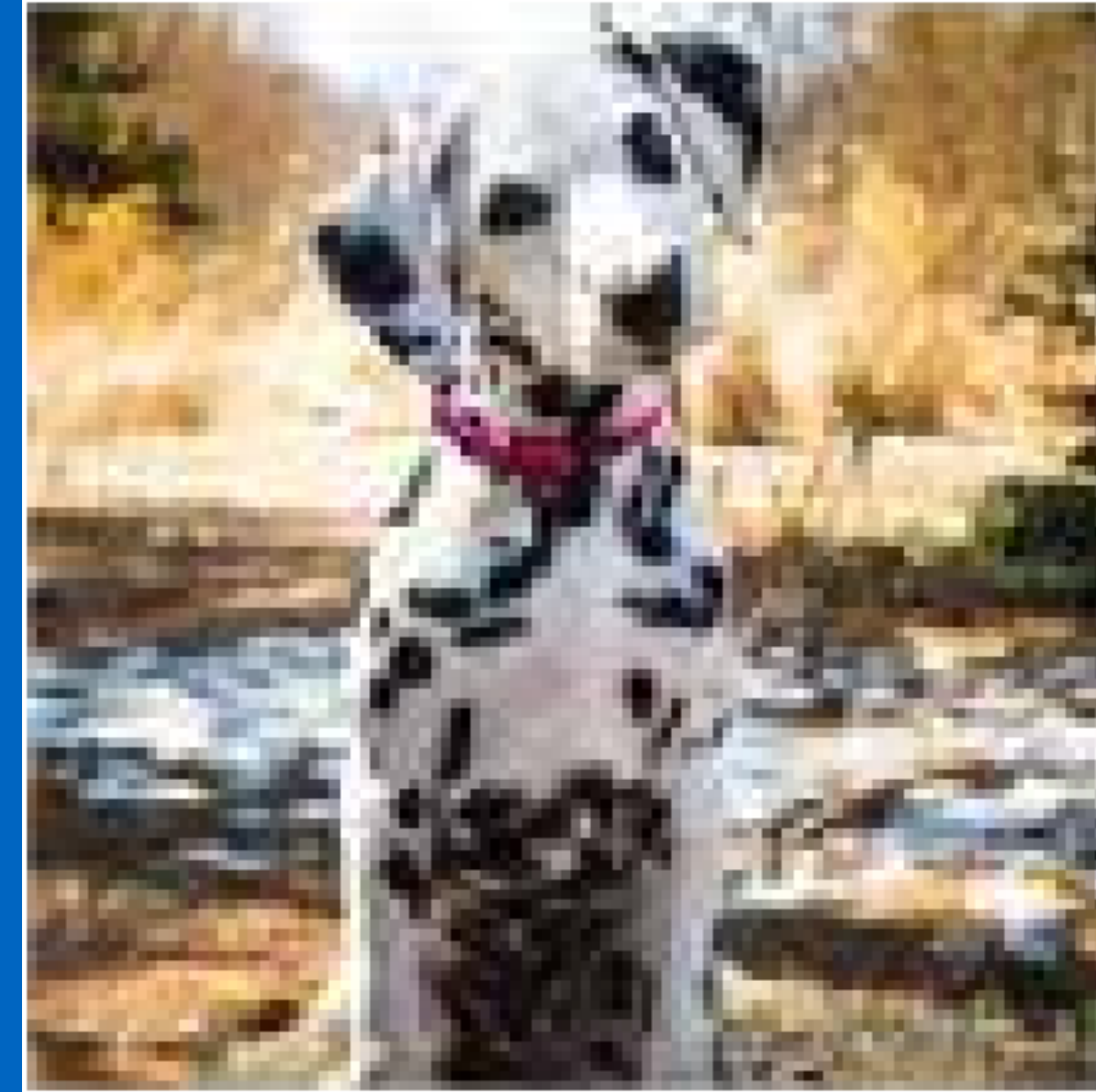
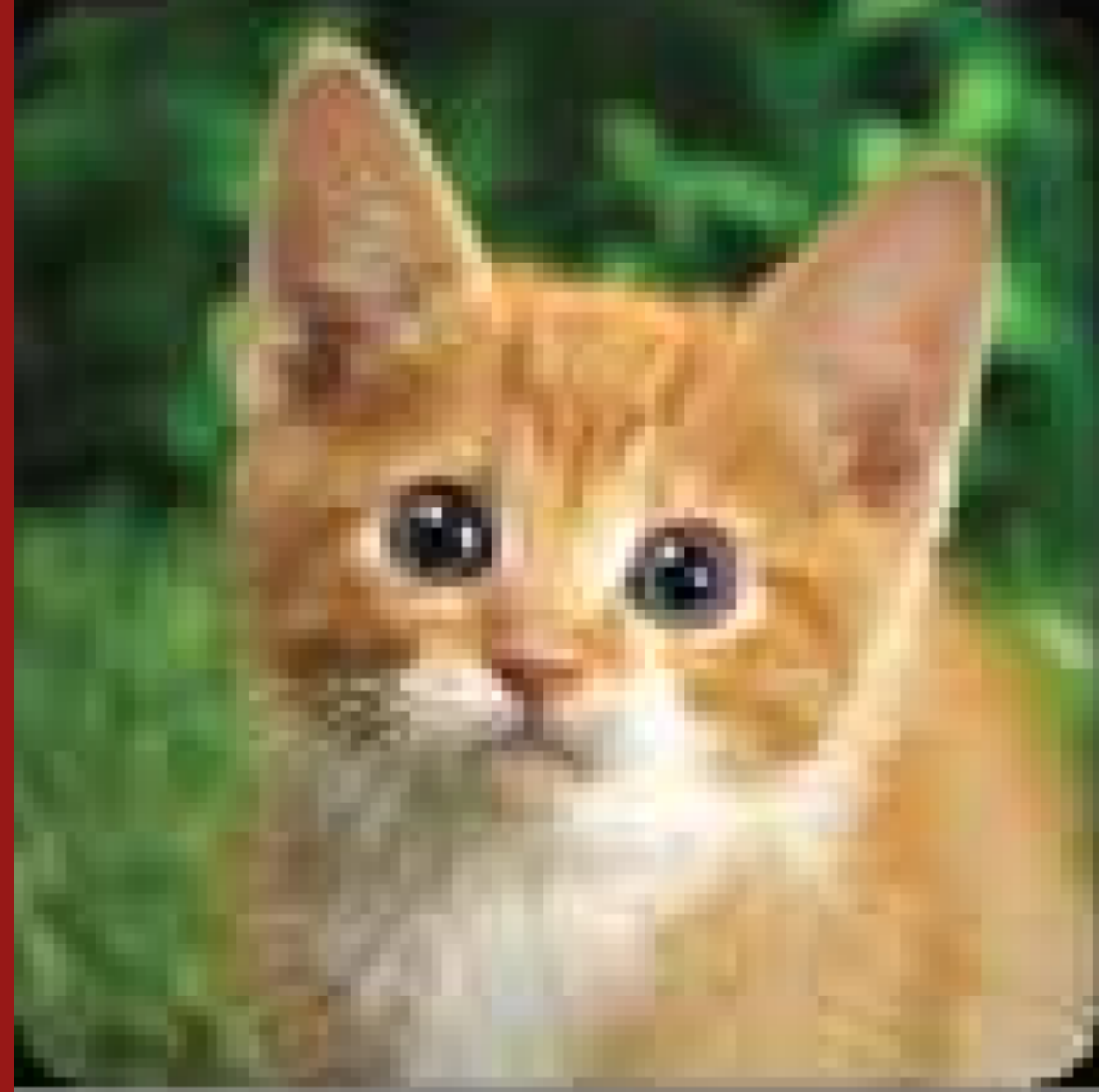
{wieland, jonas, matthias}@bethgelab.org

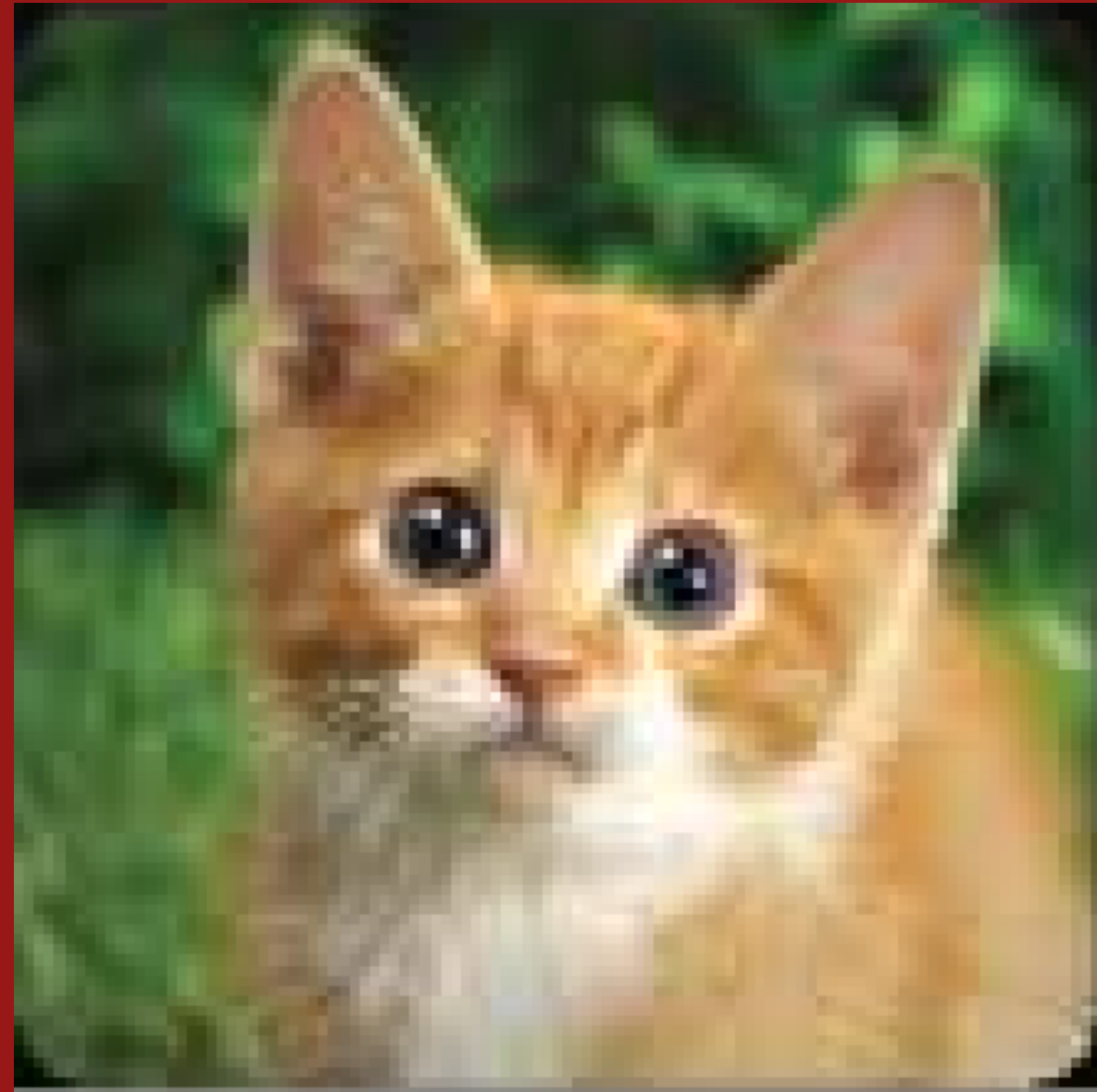
Threat Model:

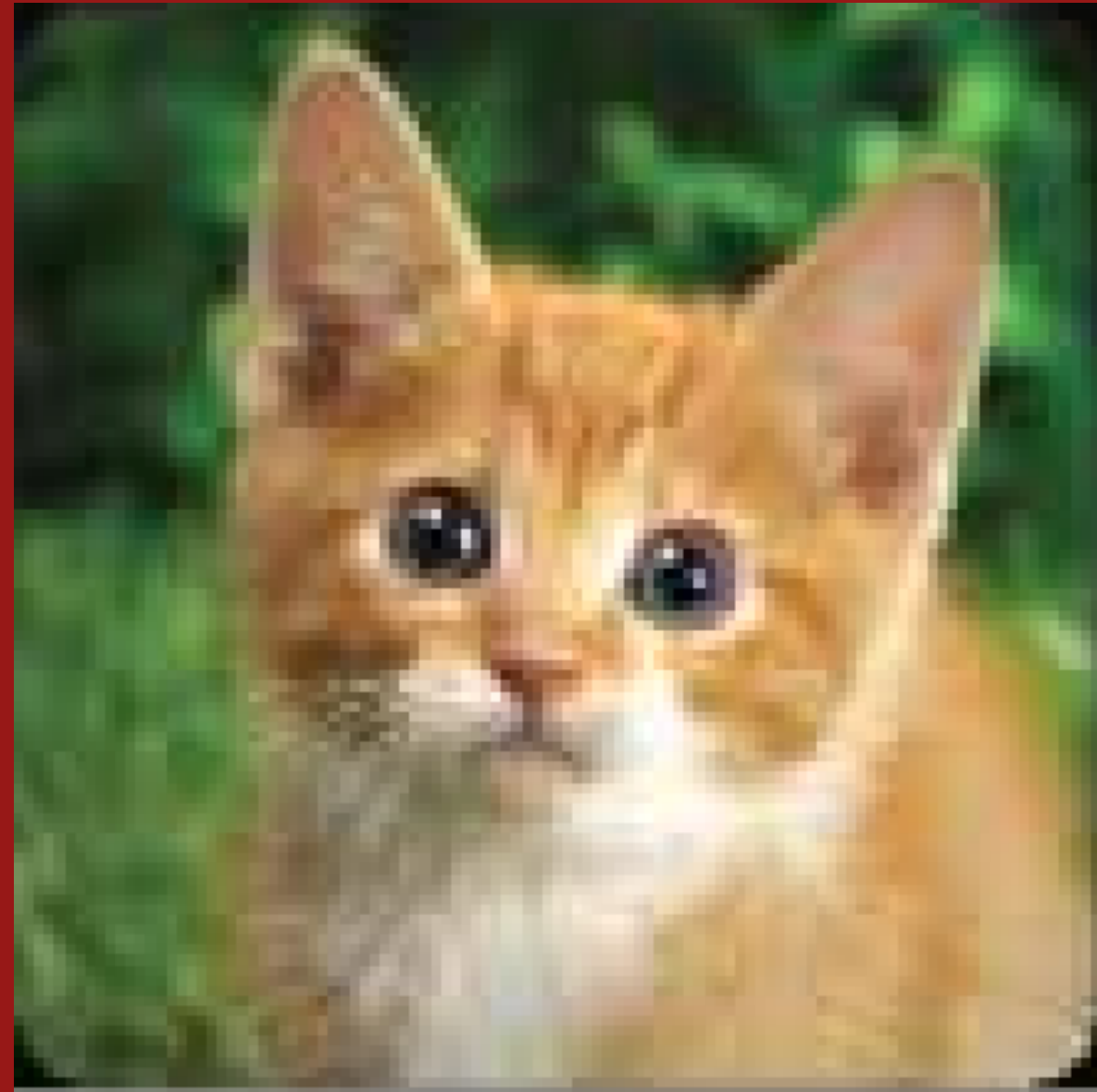
- Black Box
- Hard Label
- Query Access

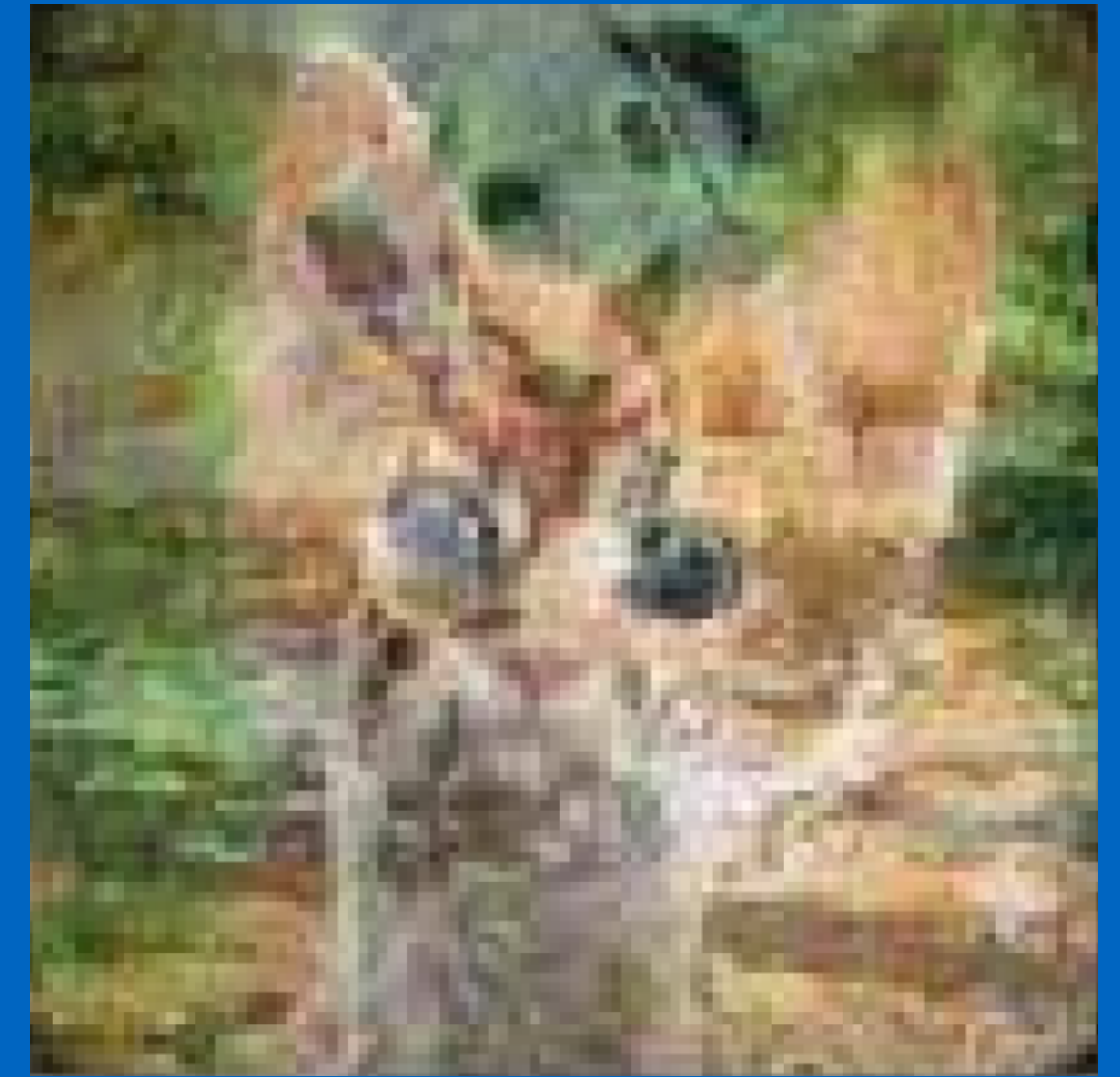
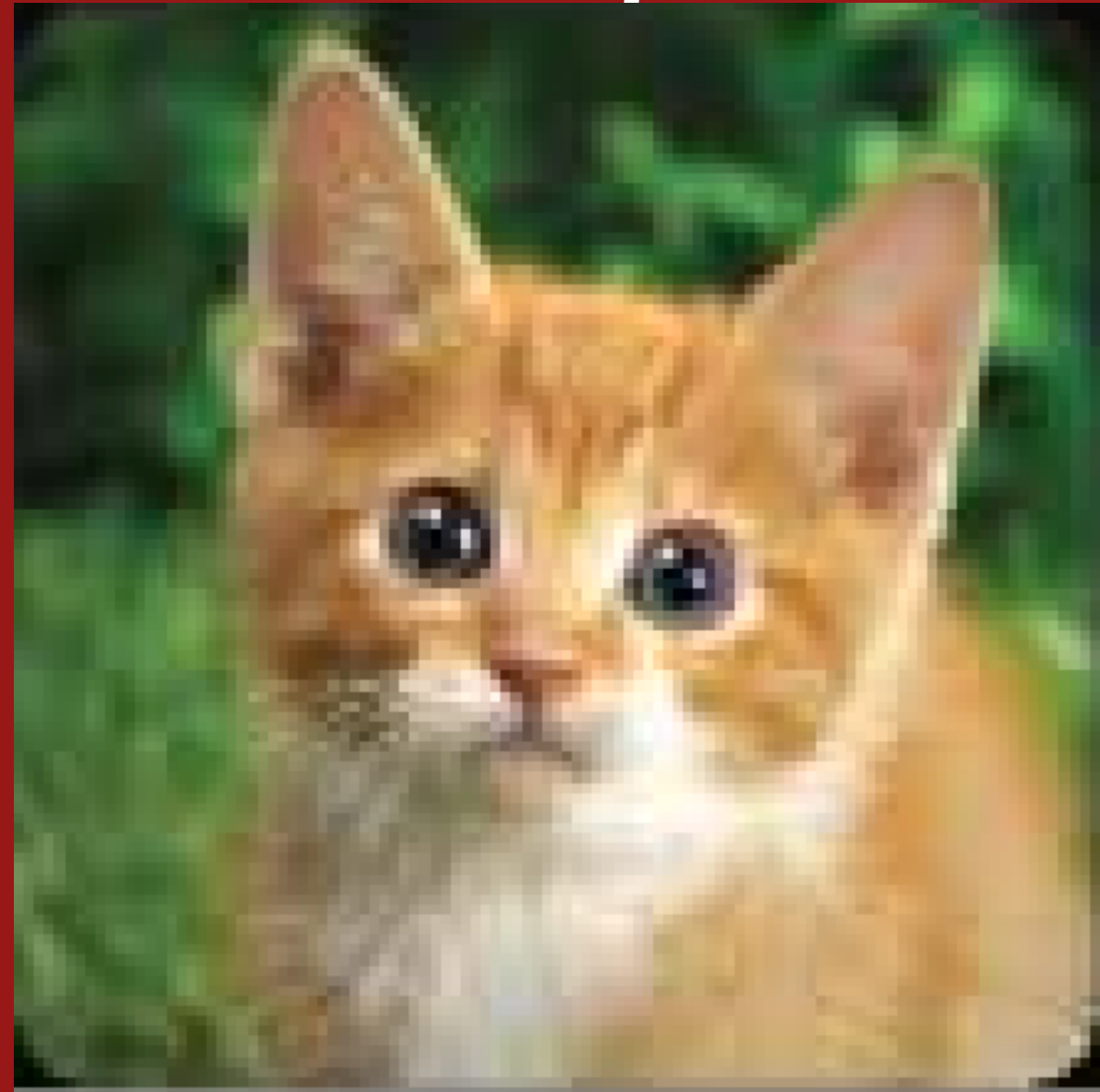


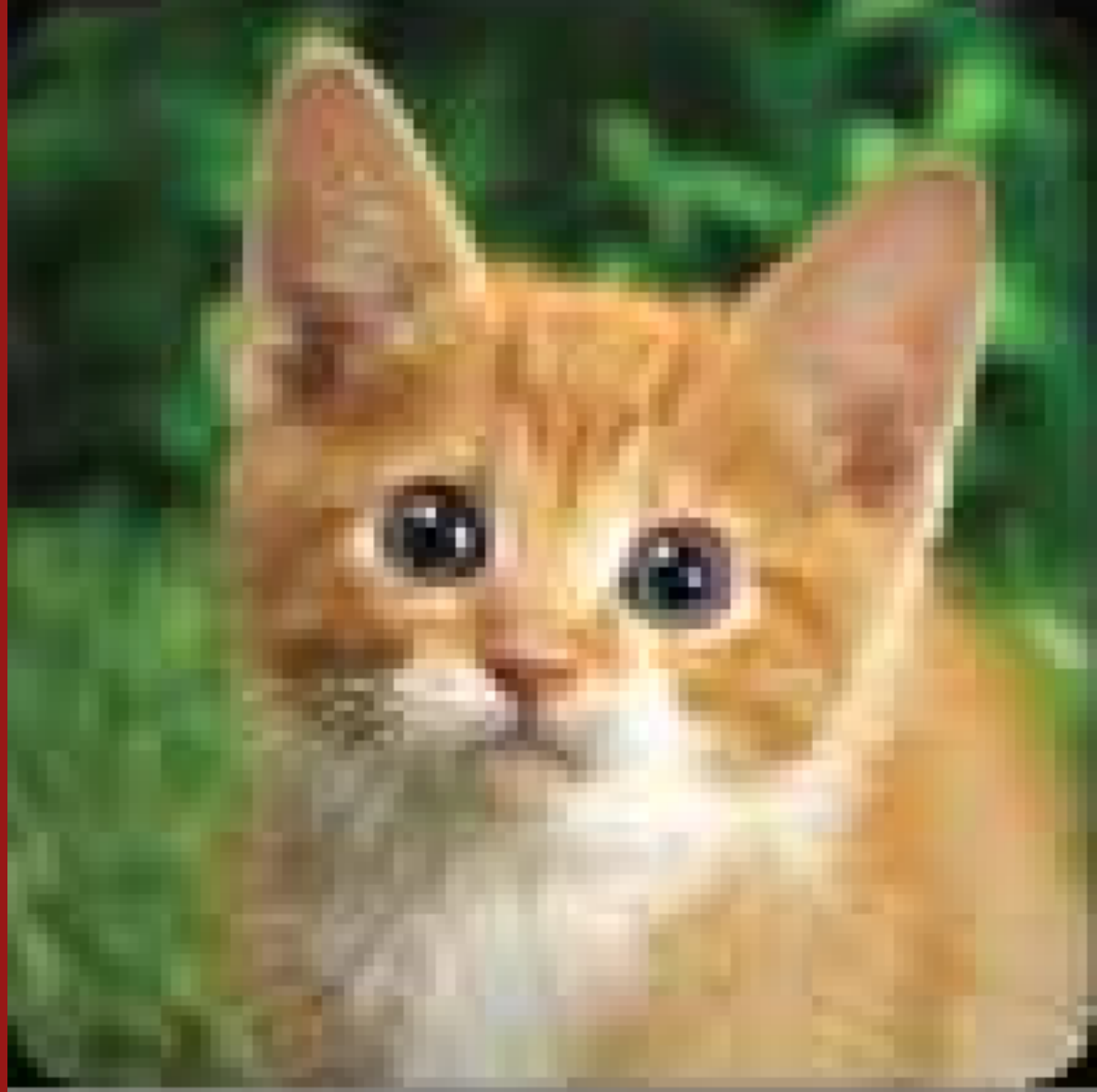


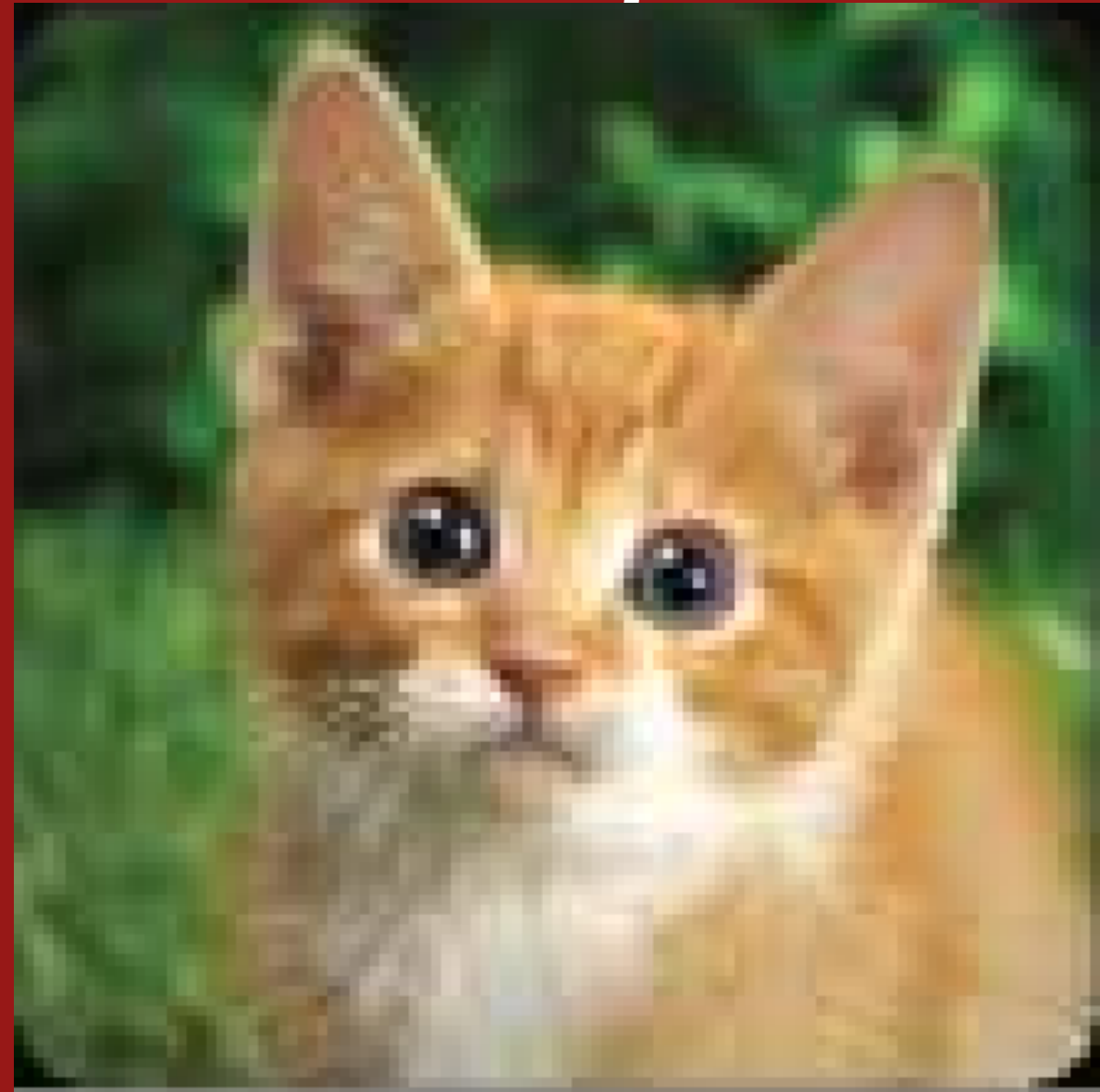


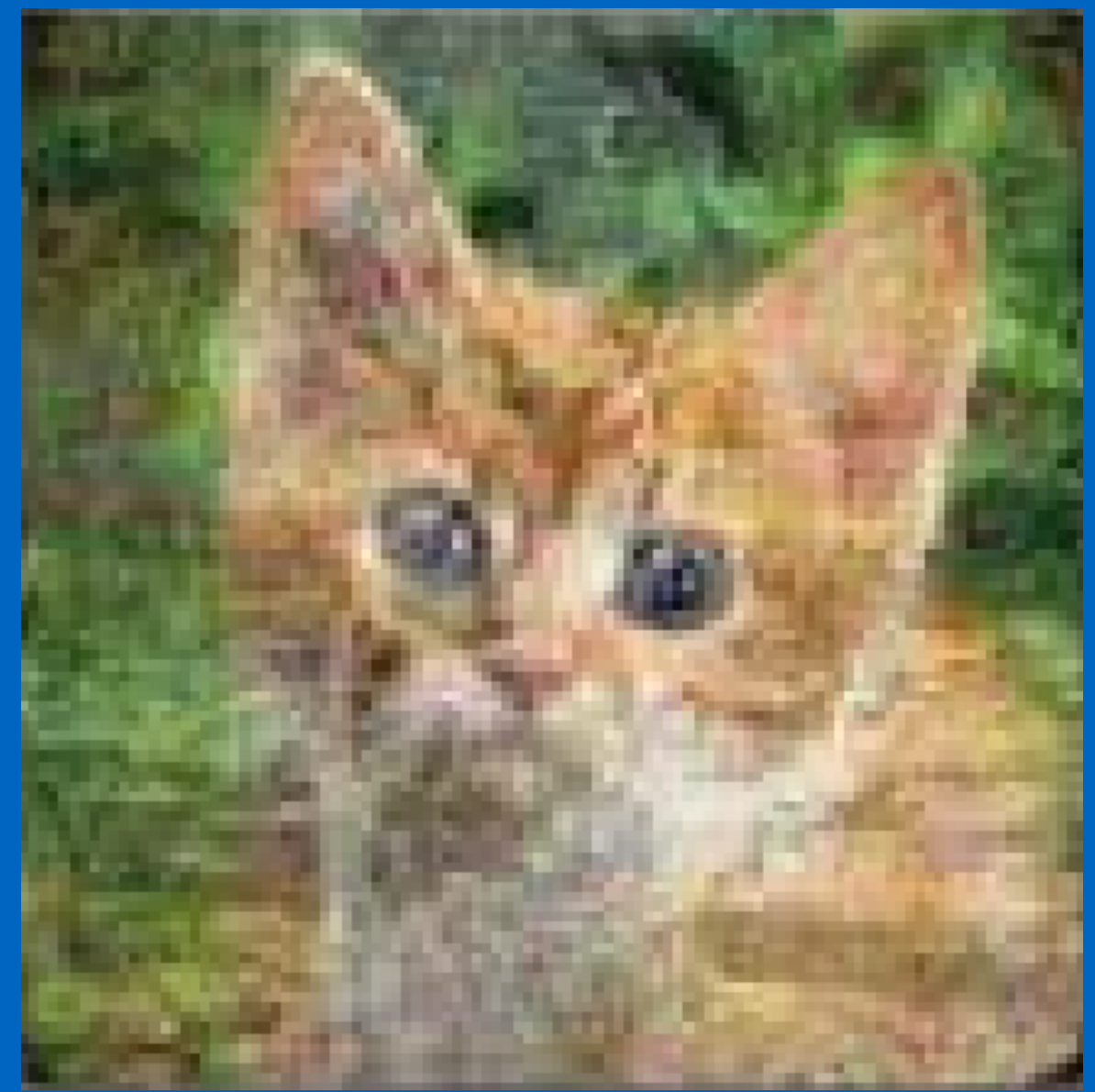
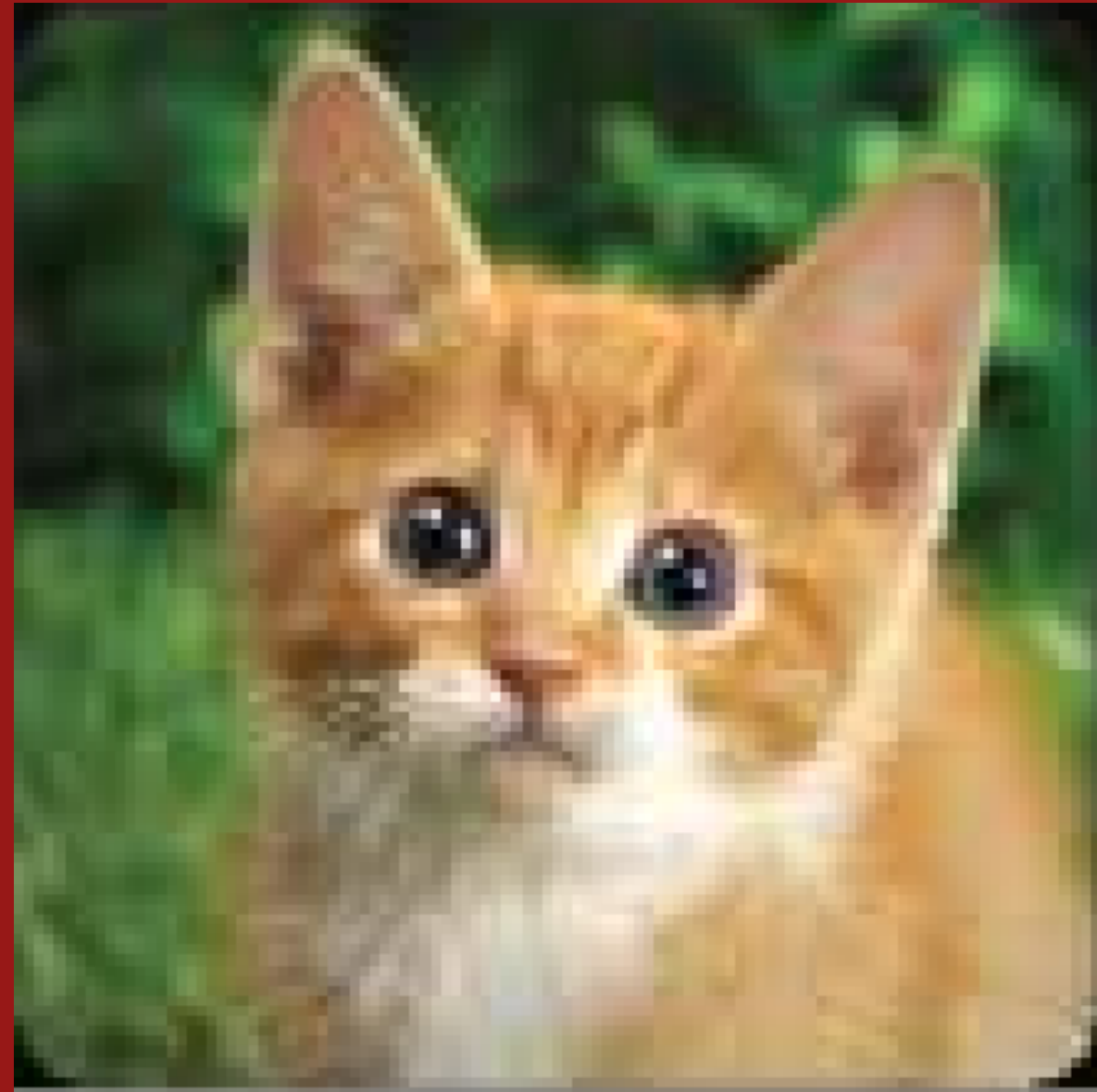


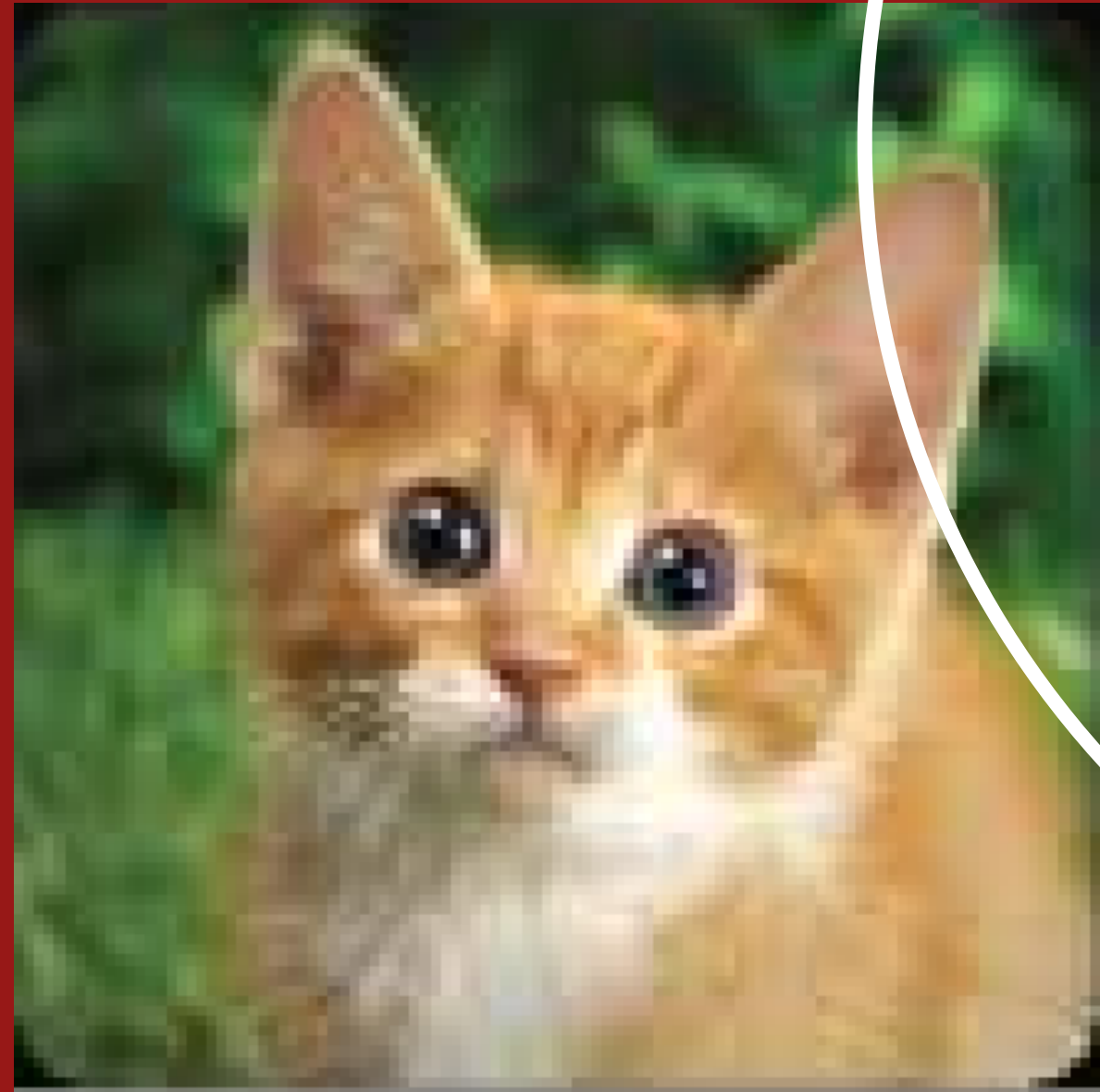


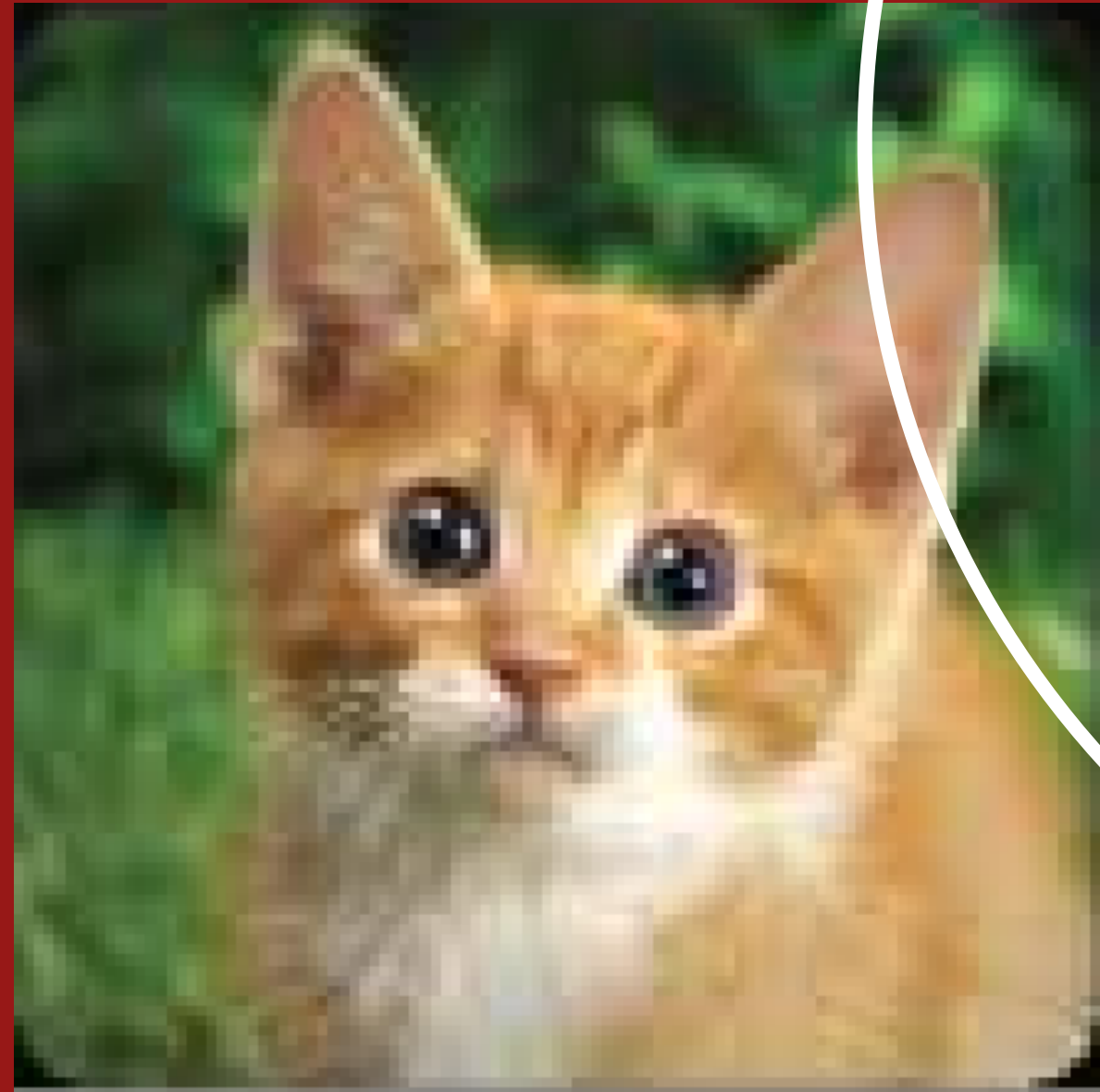


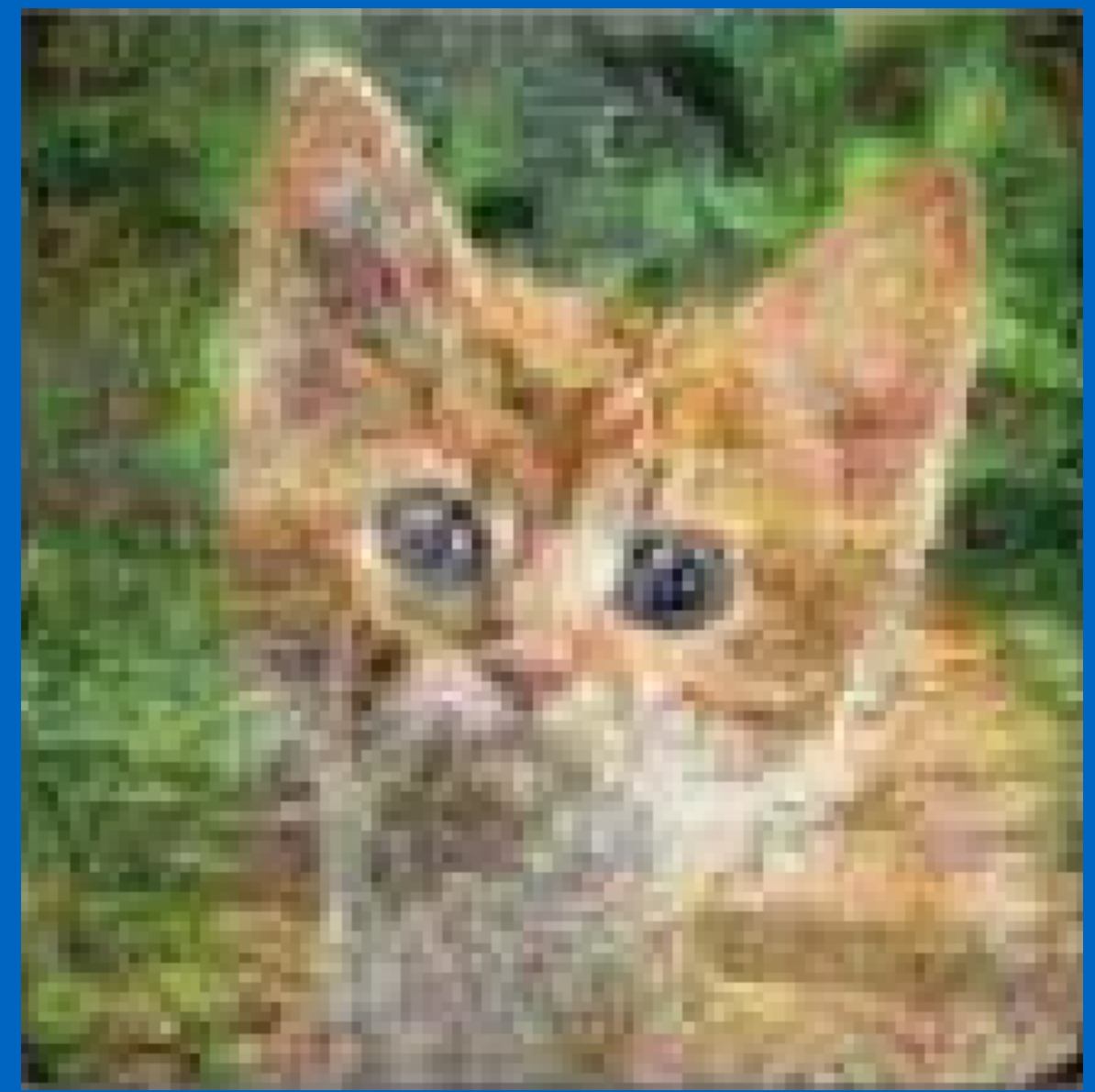
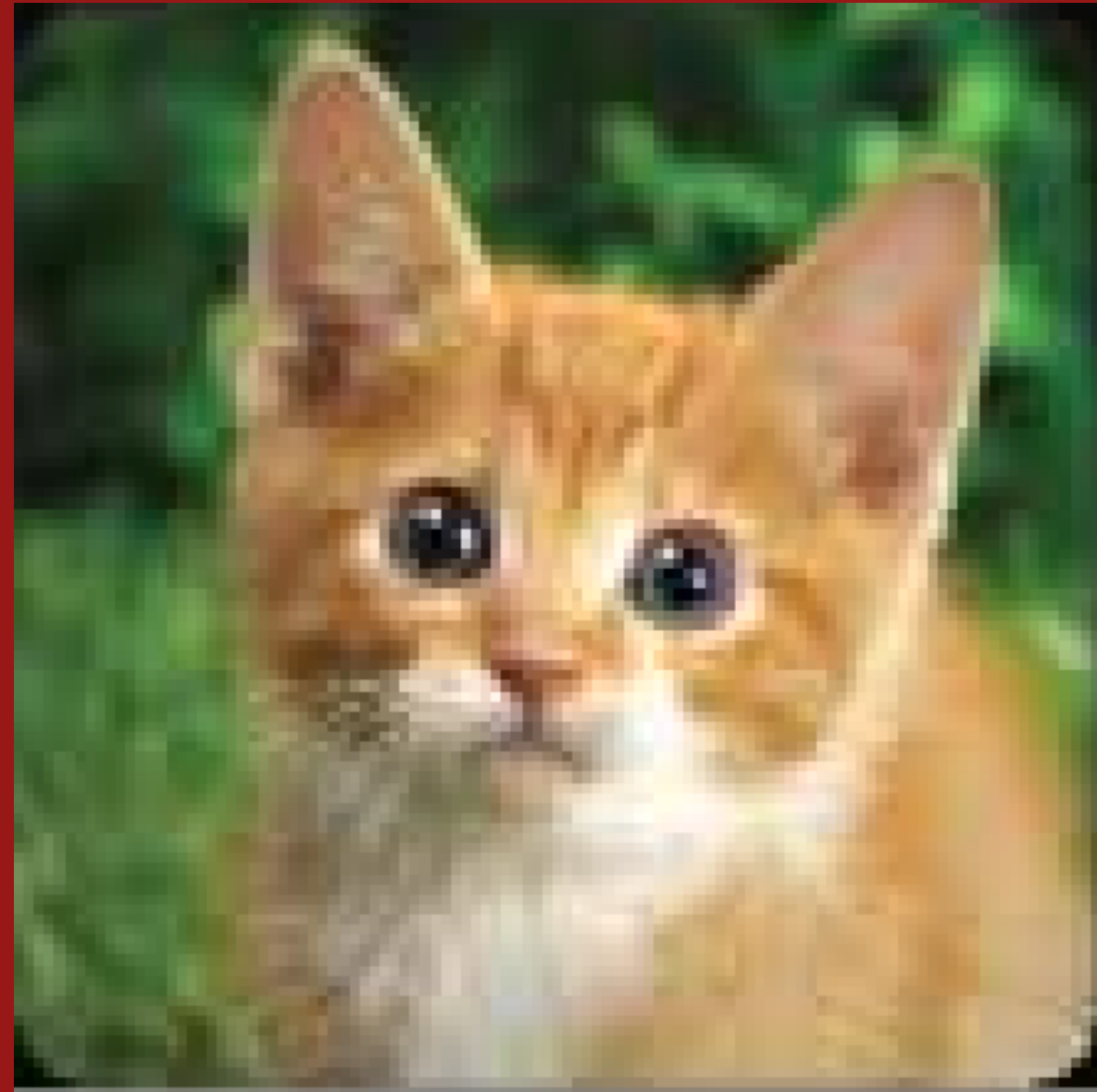


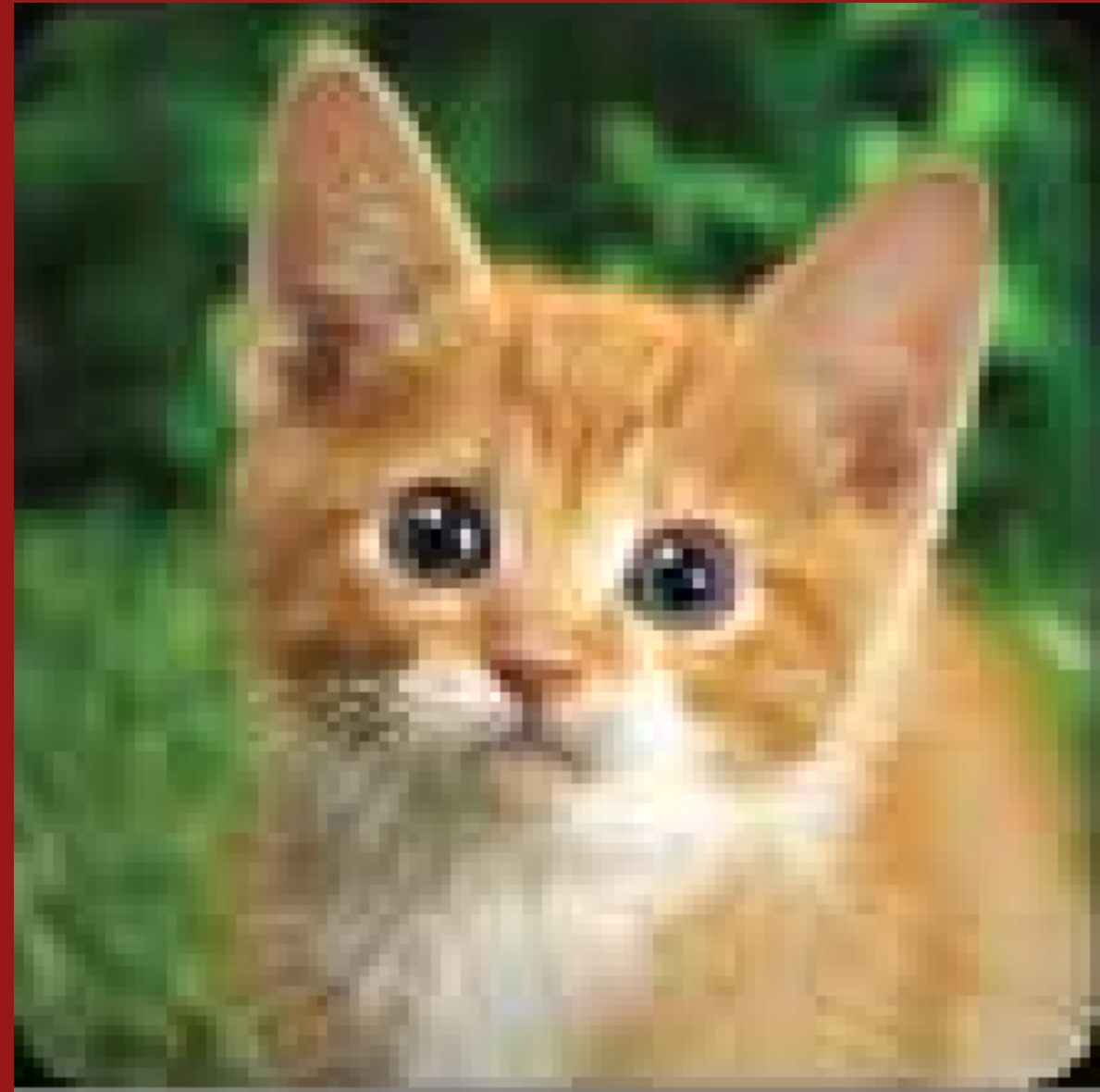


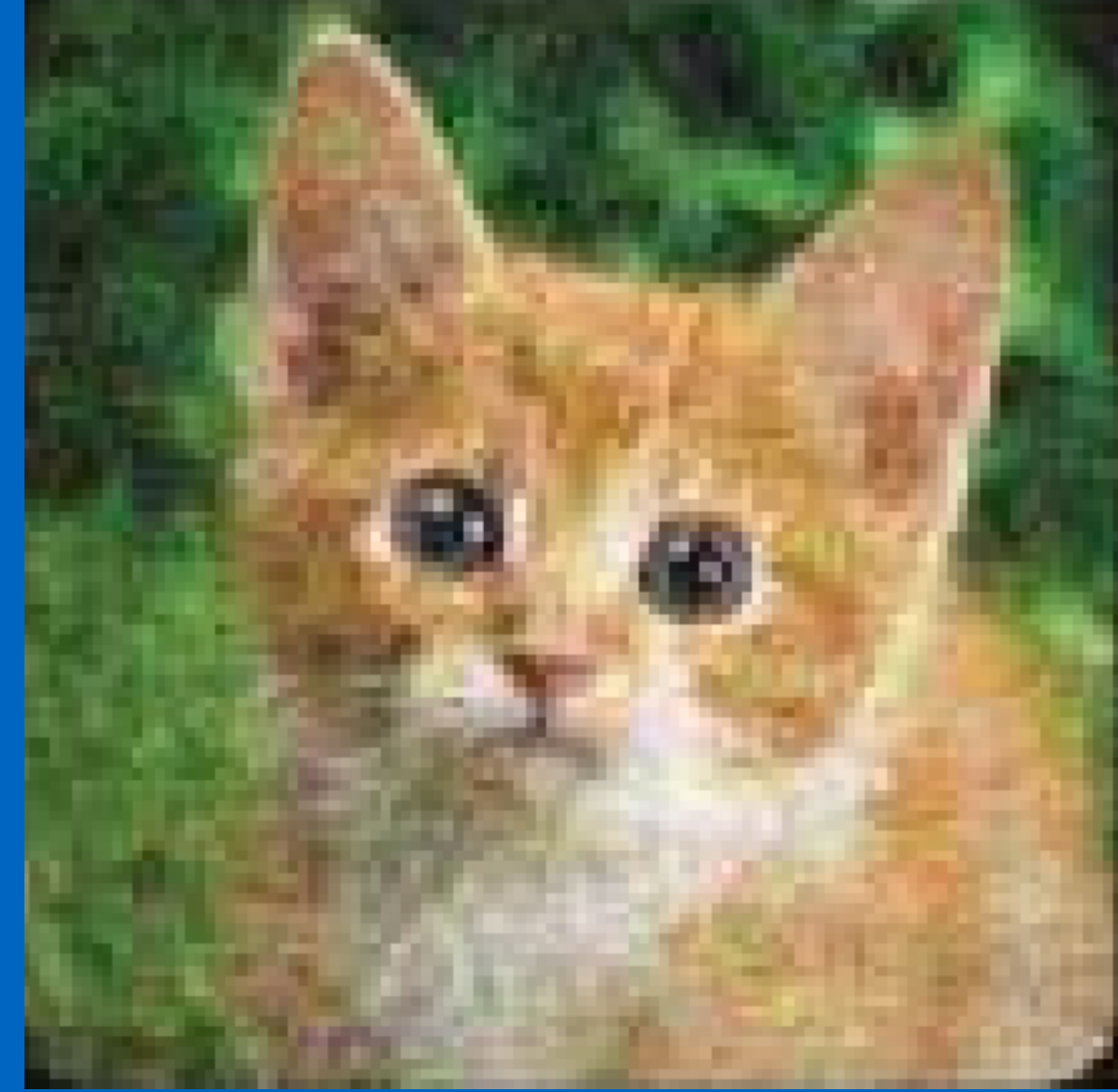
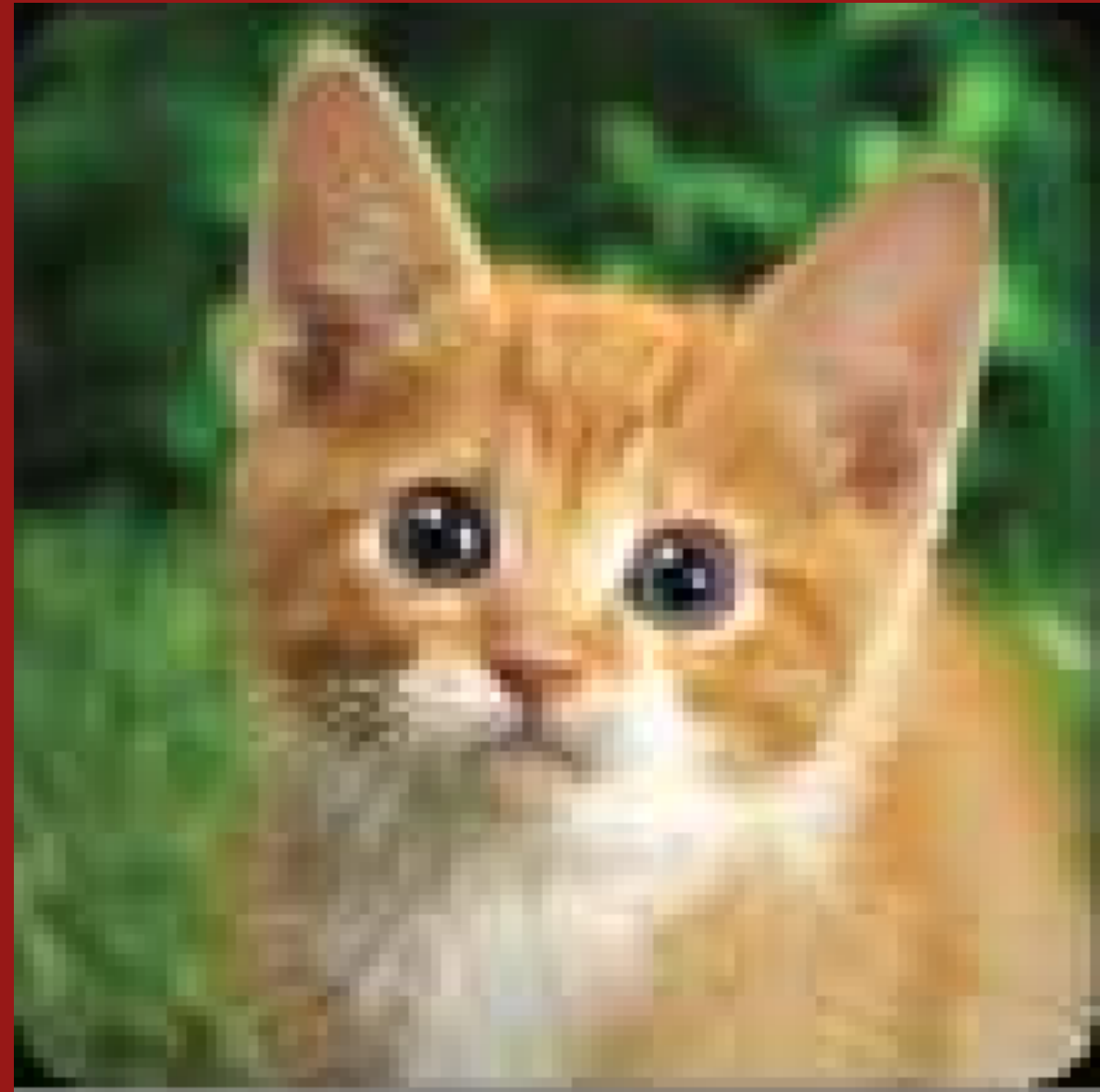


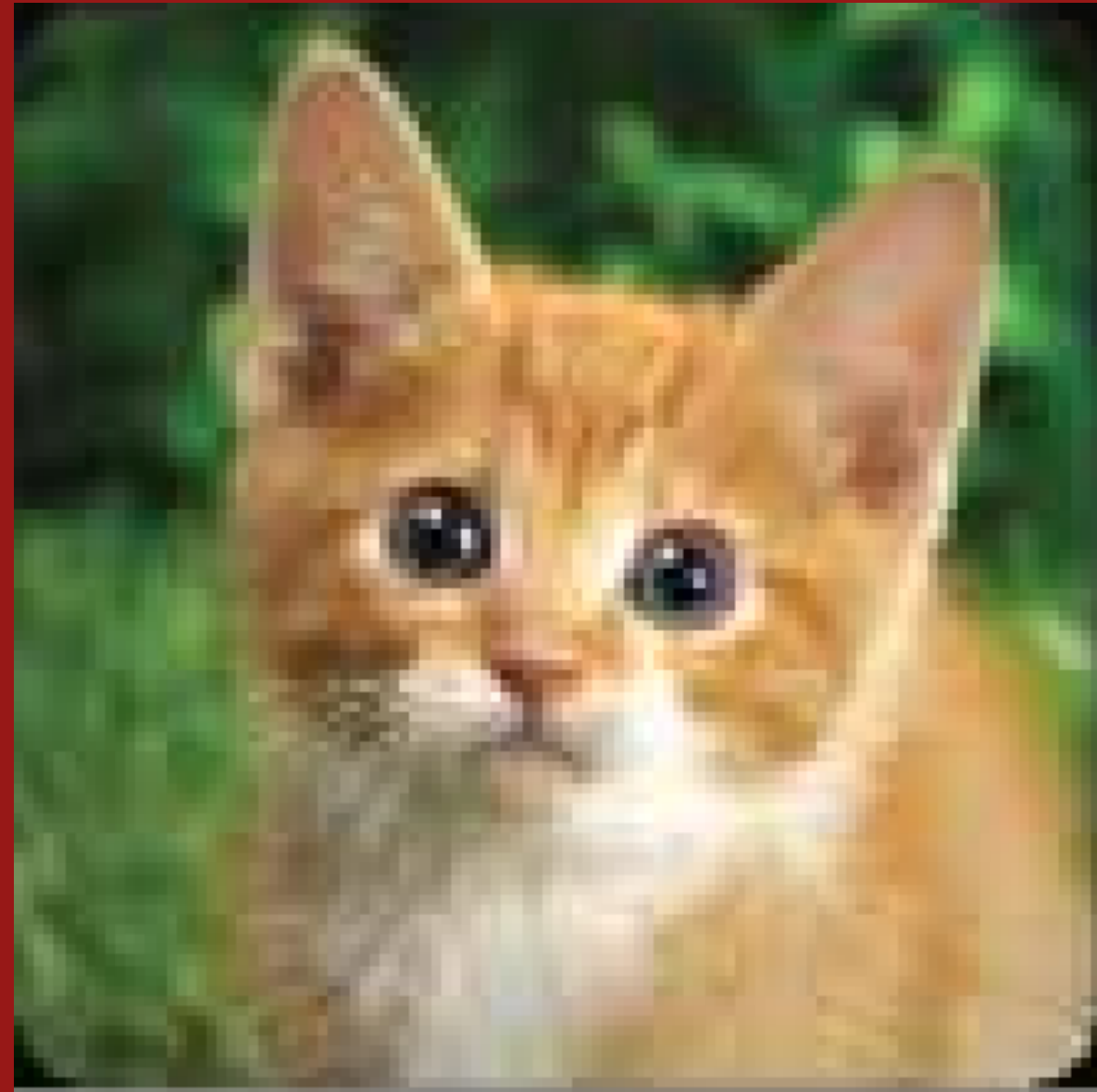


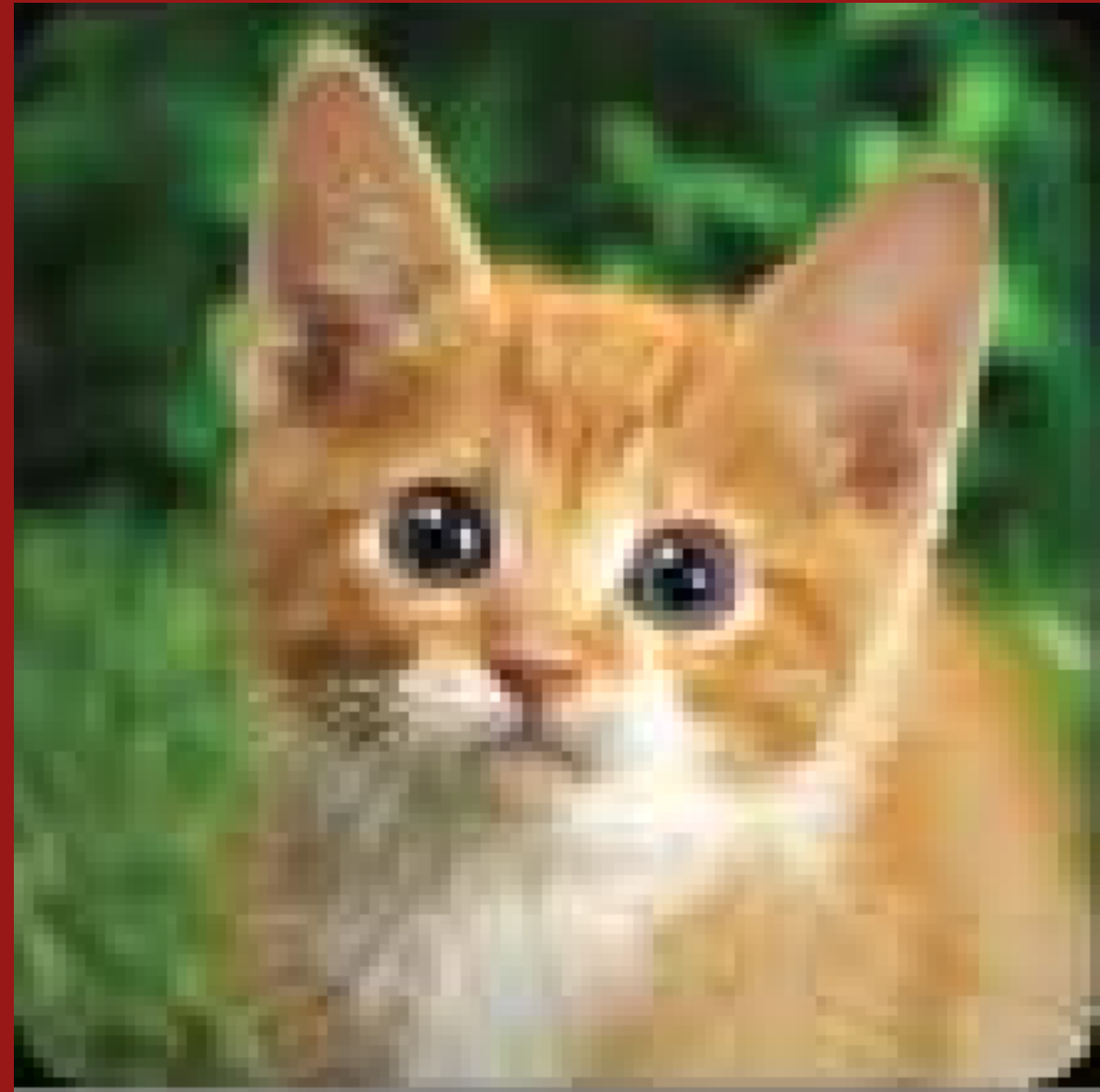












A geometry-inspired decision-based attack

Yujia Liu*

University of Science and Technology of China

Hefei, China

yjcaihon@mail.ustc.edu.cn

Seyed-Mohsen Moosavi-Dezfooli

École Polytechnique Fédérale de Lausanne

Lausanne, Switzerland

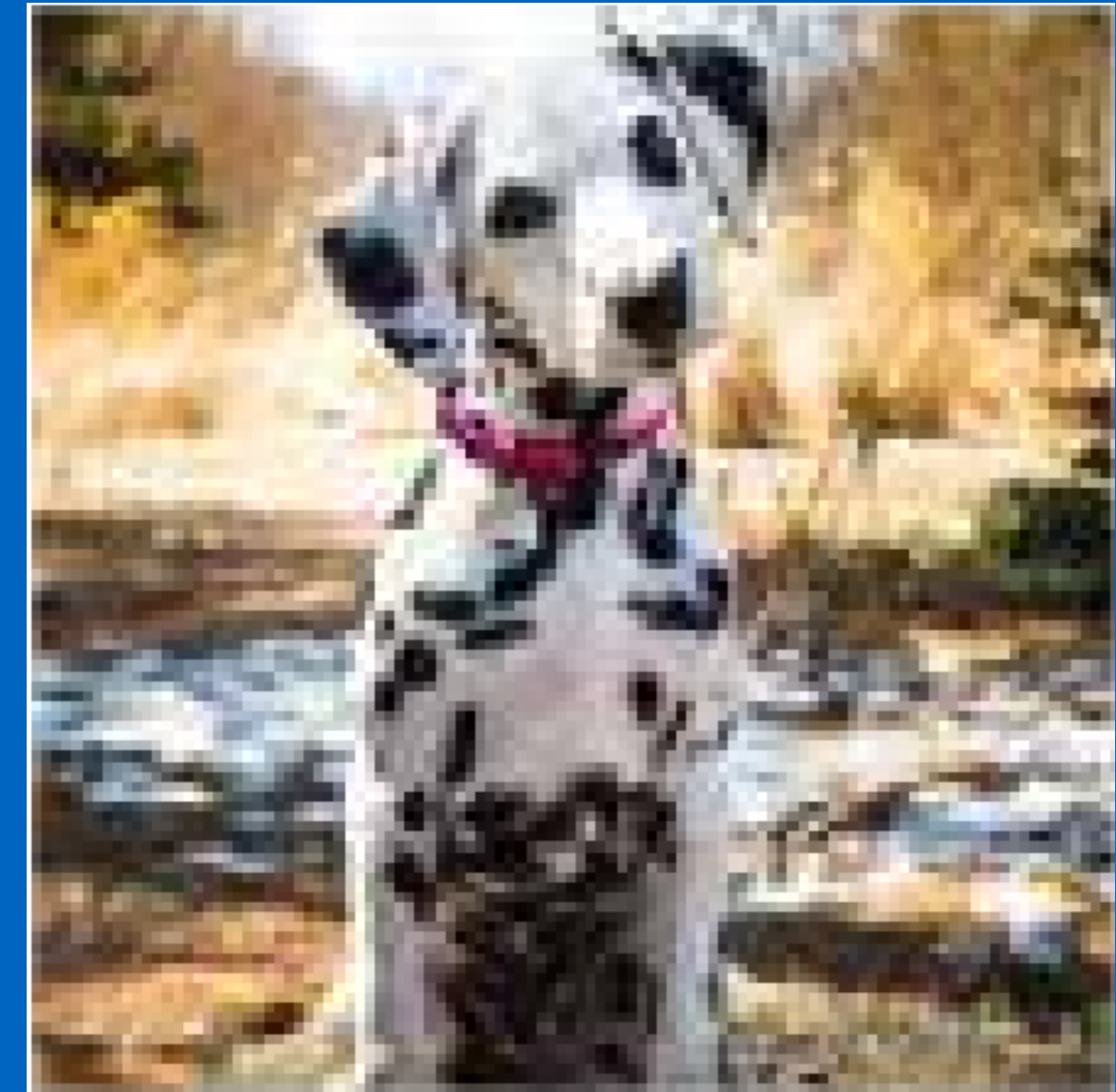
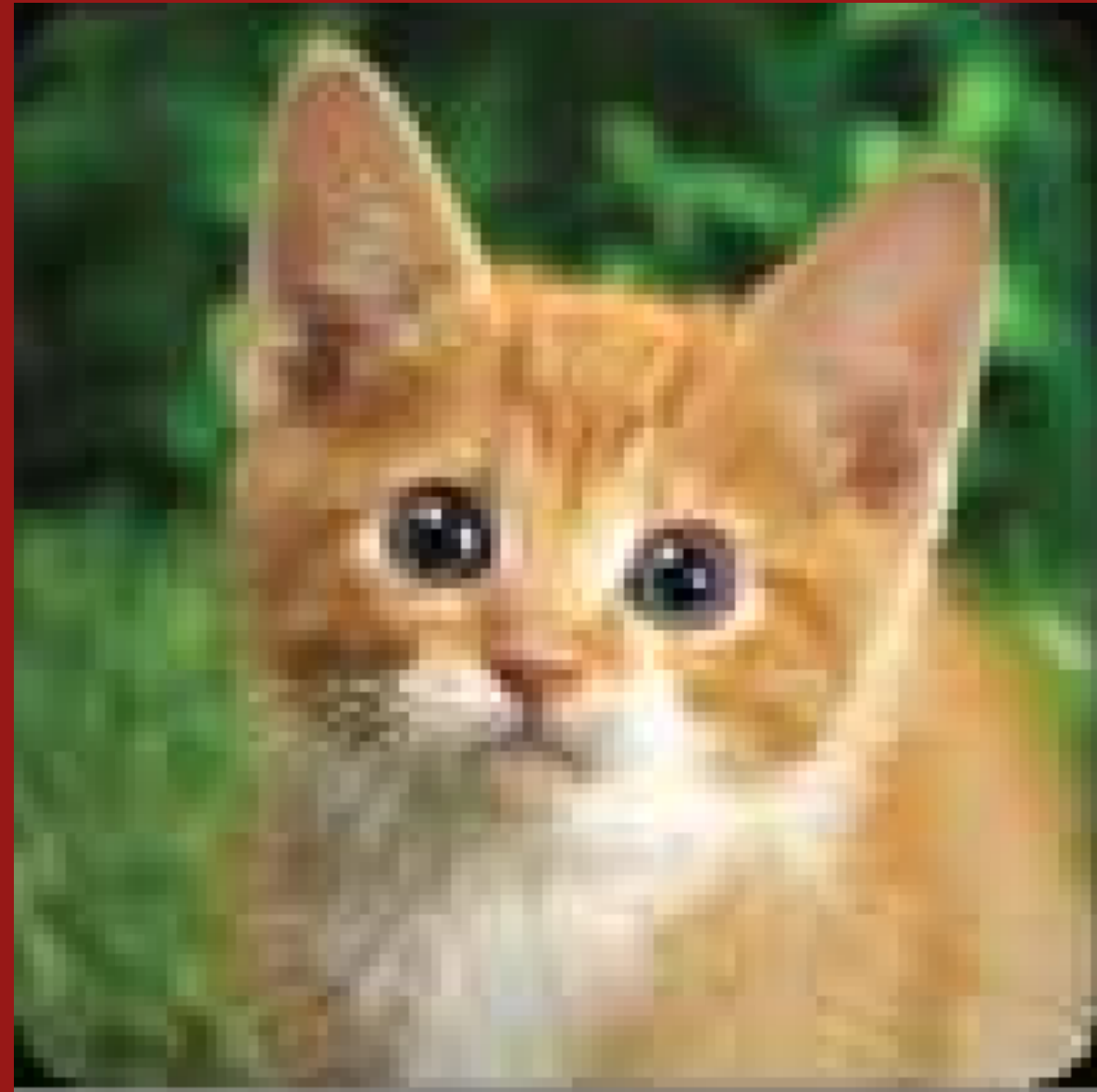
seyed.moosavi@epfl.ch

Pascal Frossard

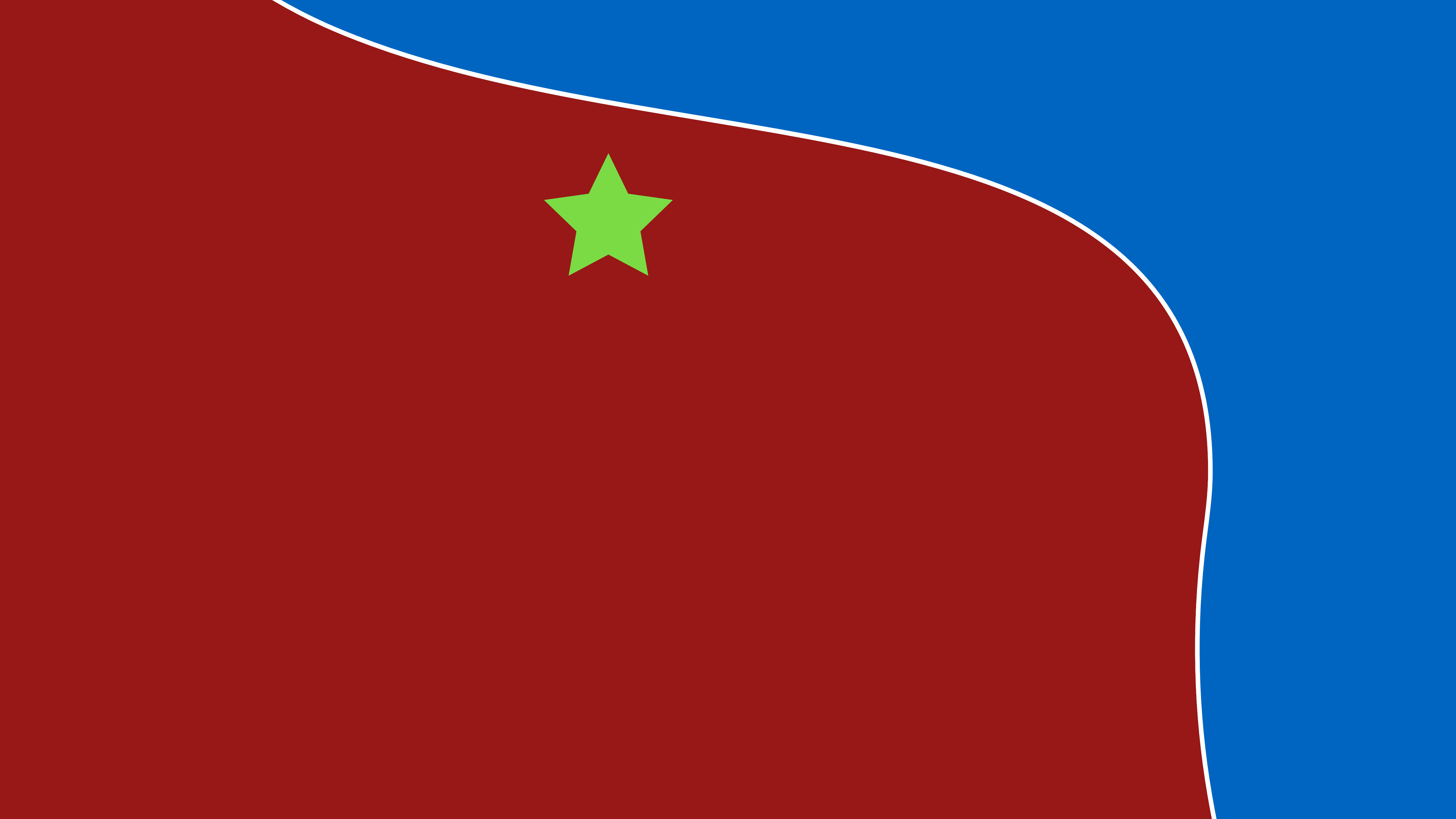
École Polytechnique Fédérale de Lausanne

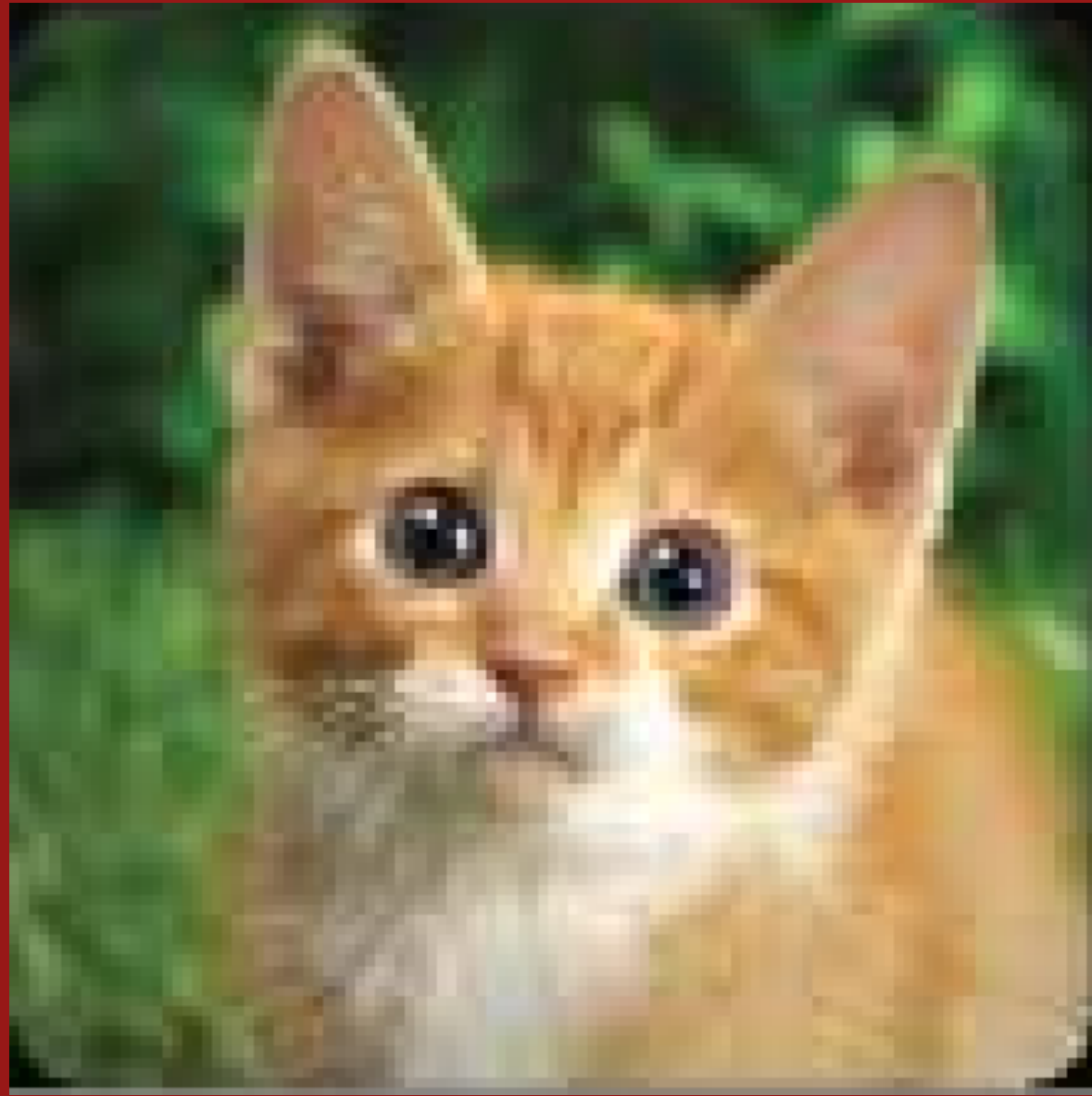
Lausanne, Switzerland

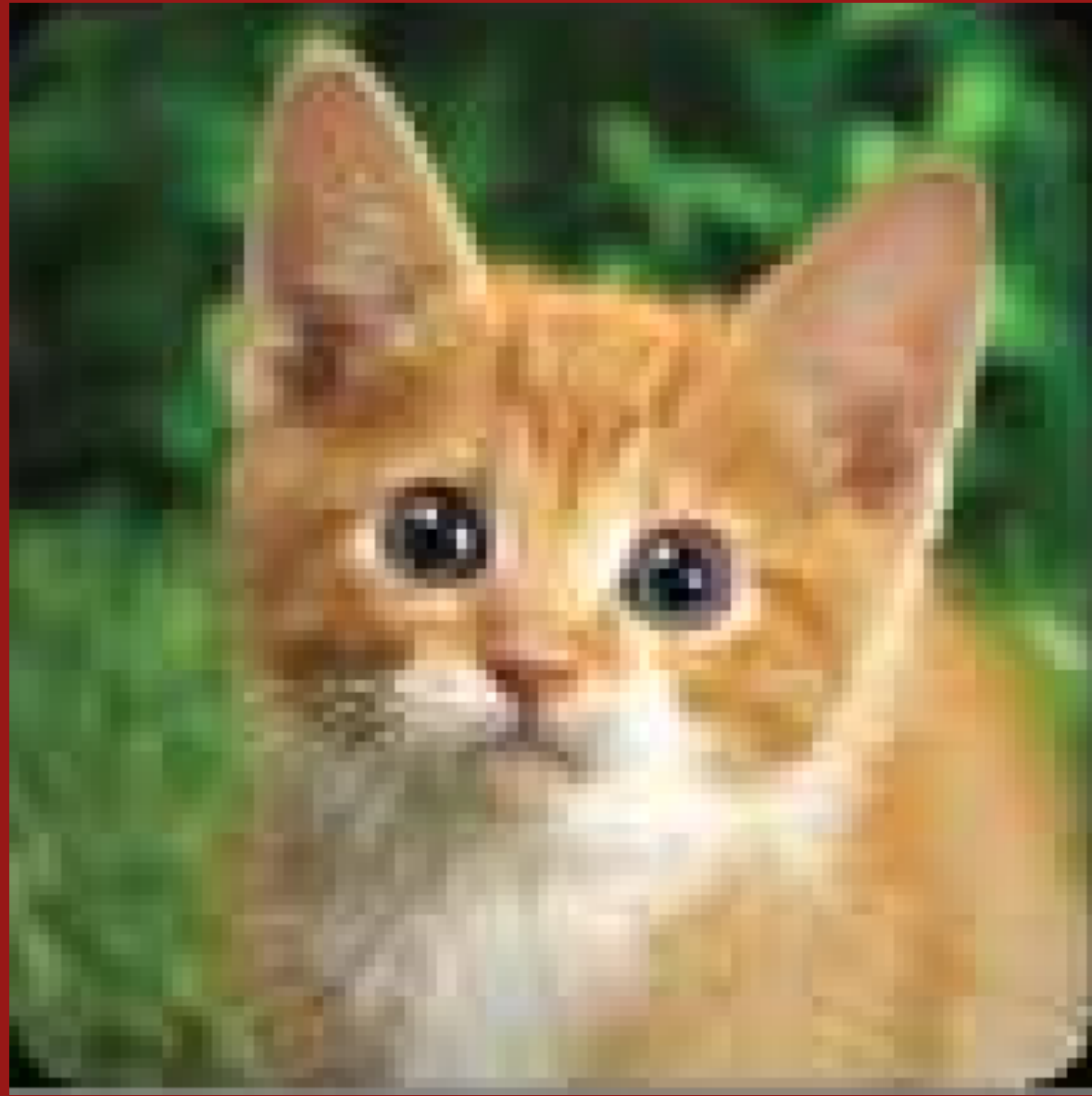
pascal.frossard@epfl.ch

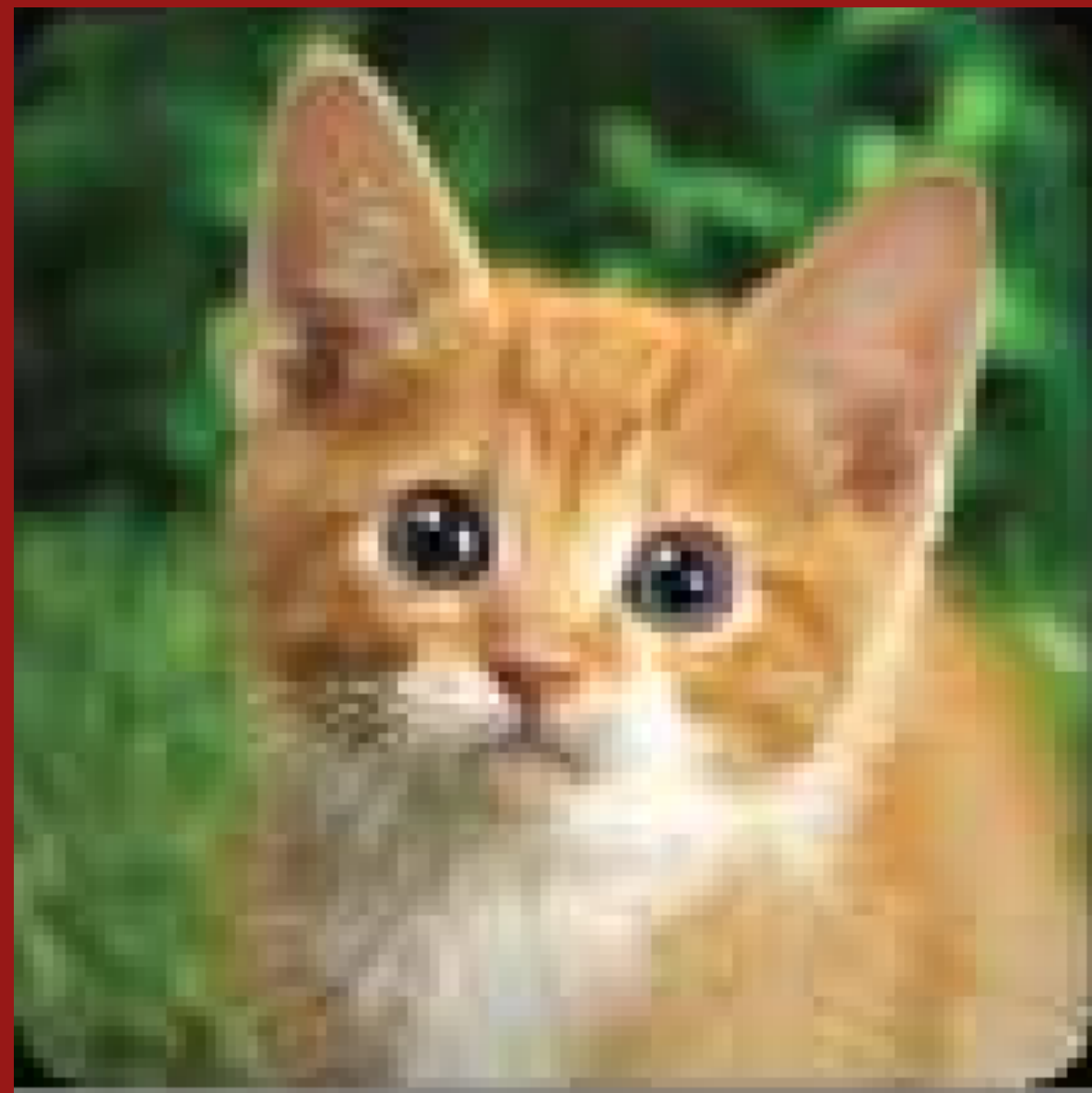


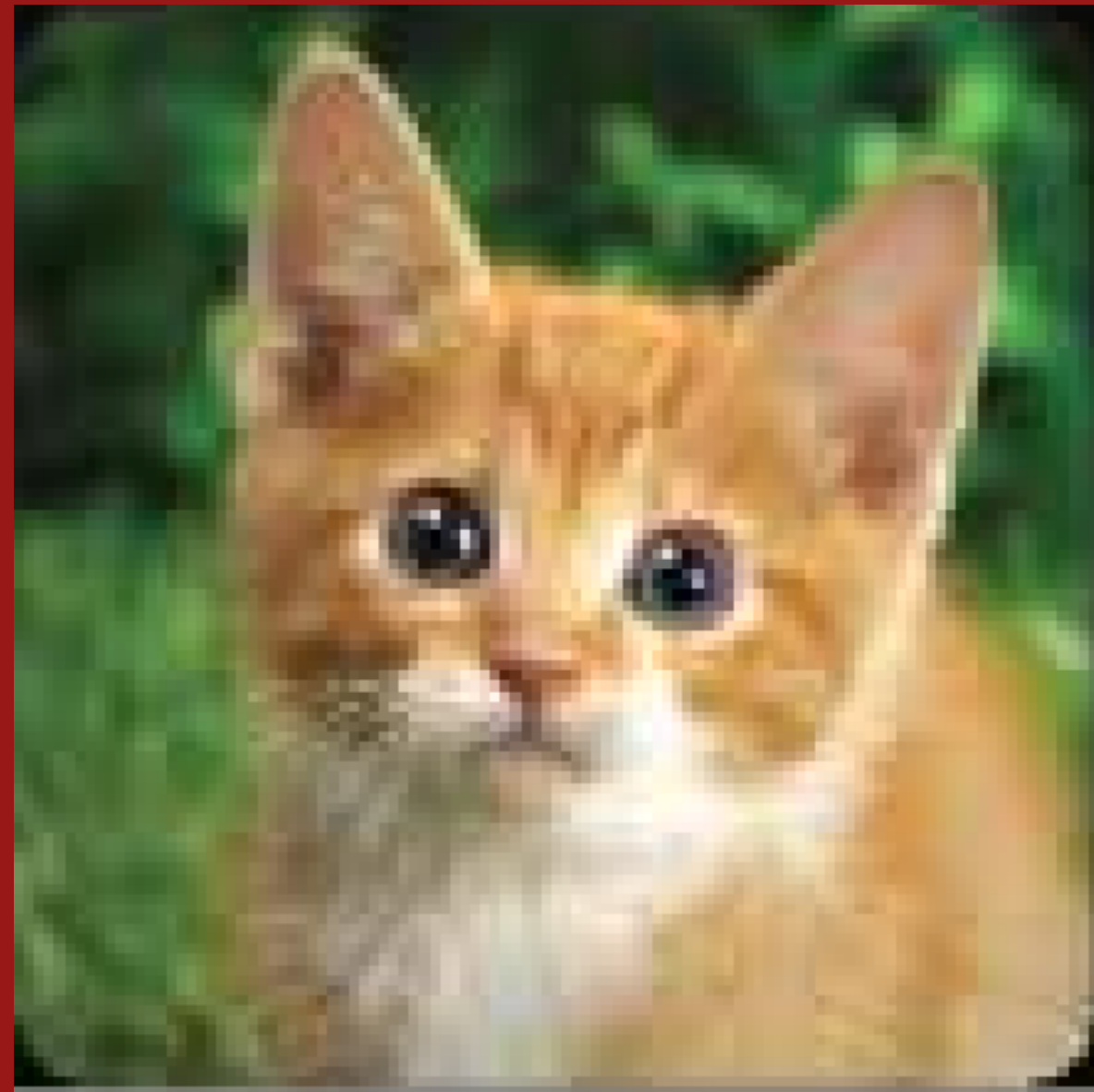
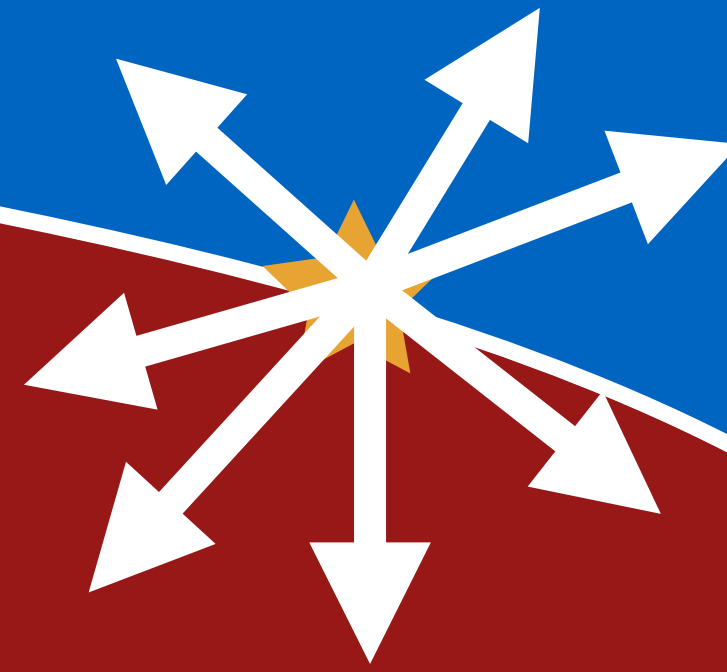


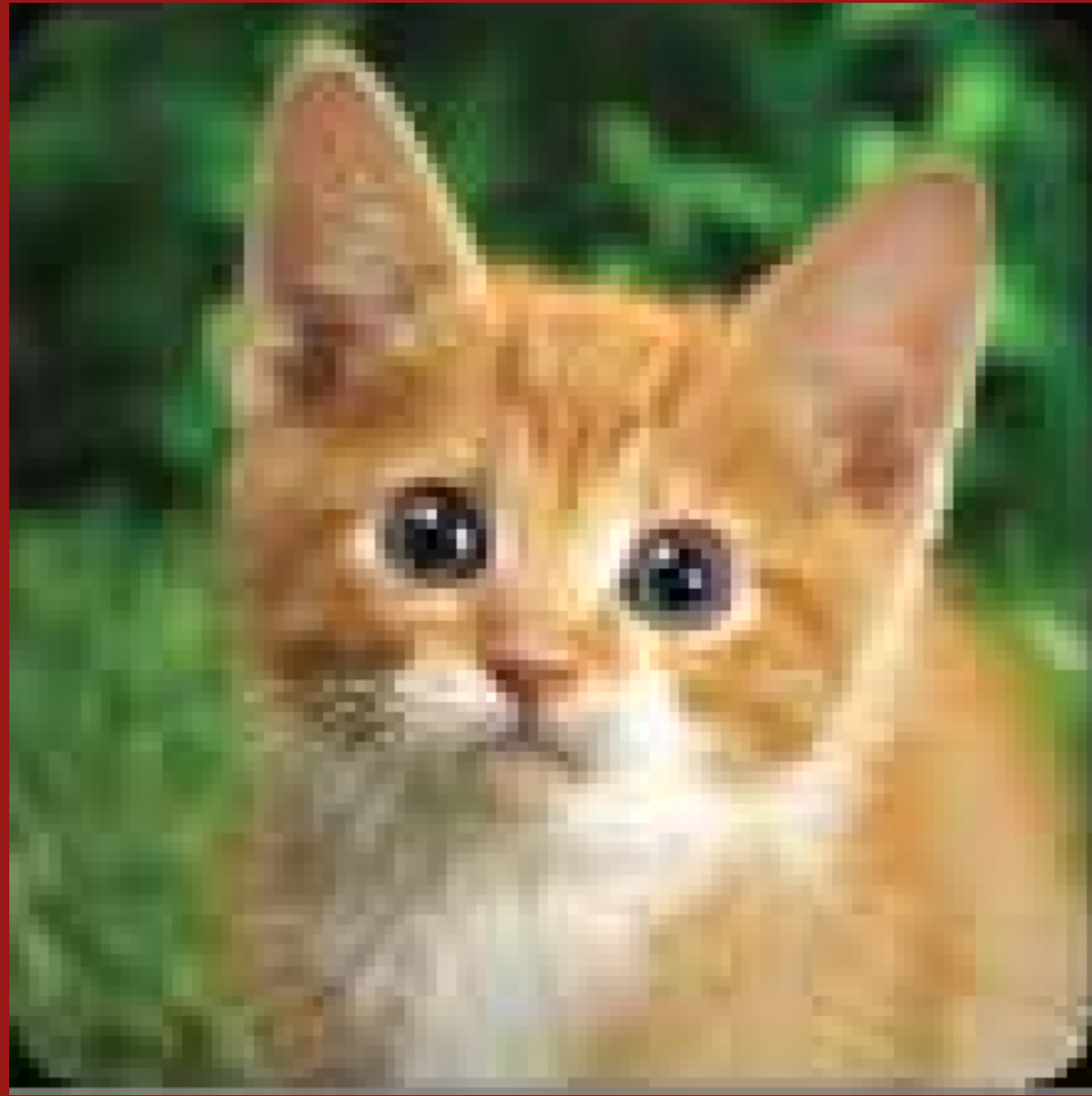


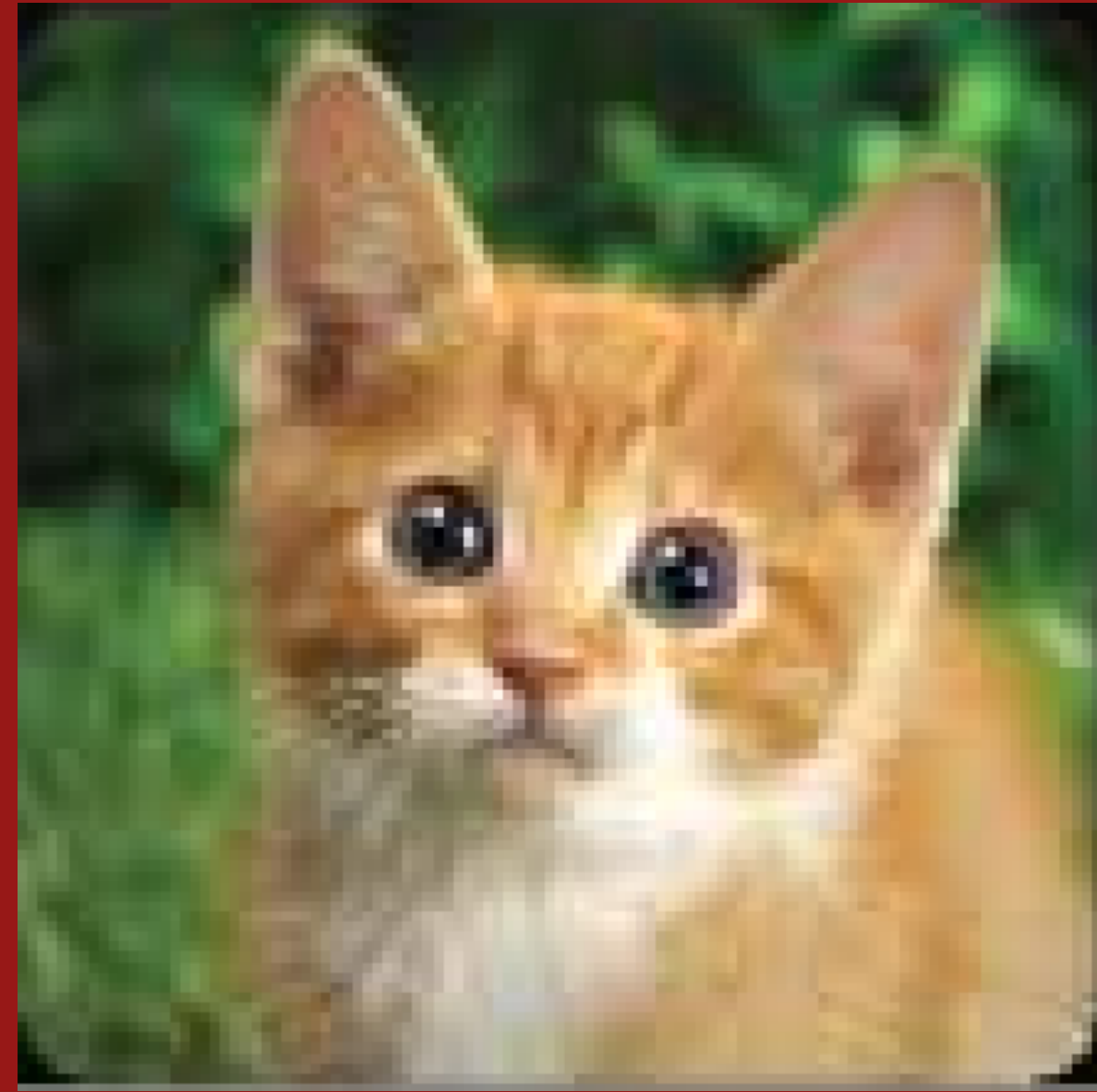


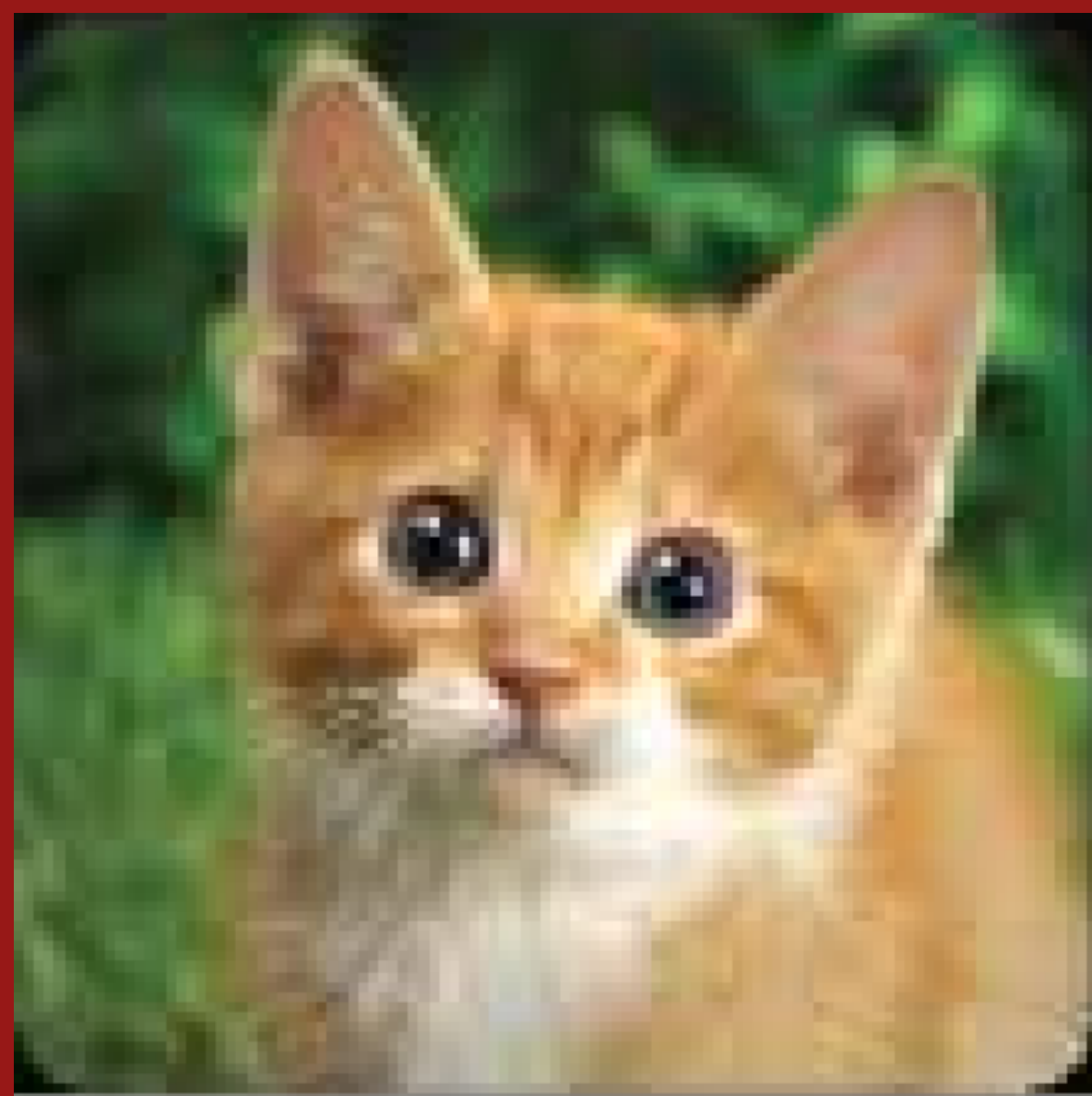
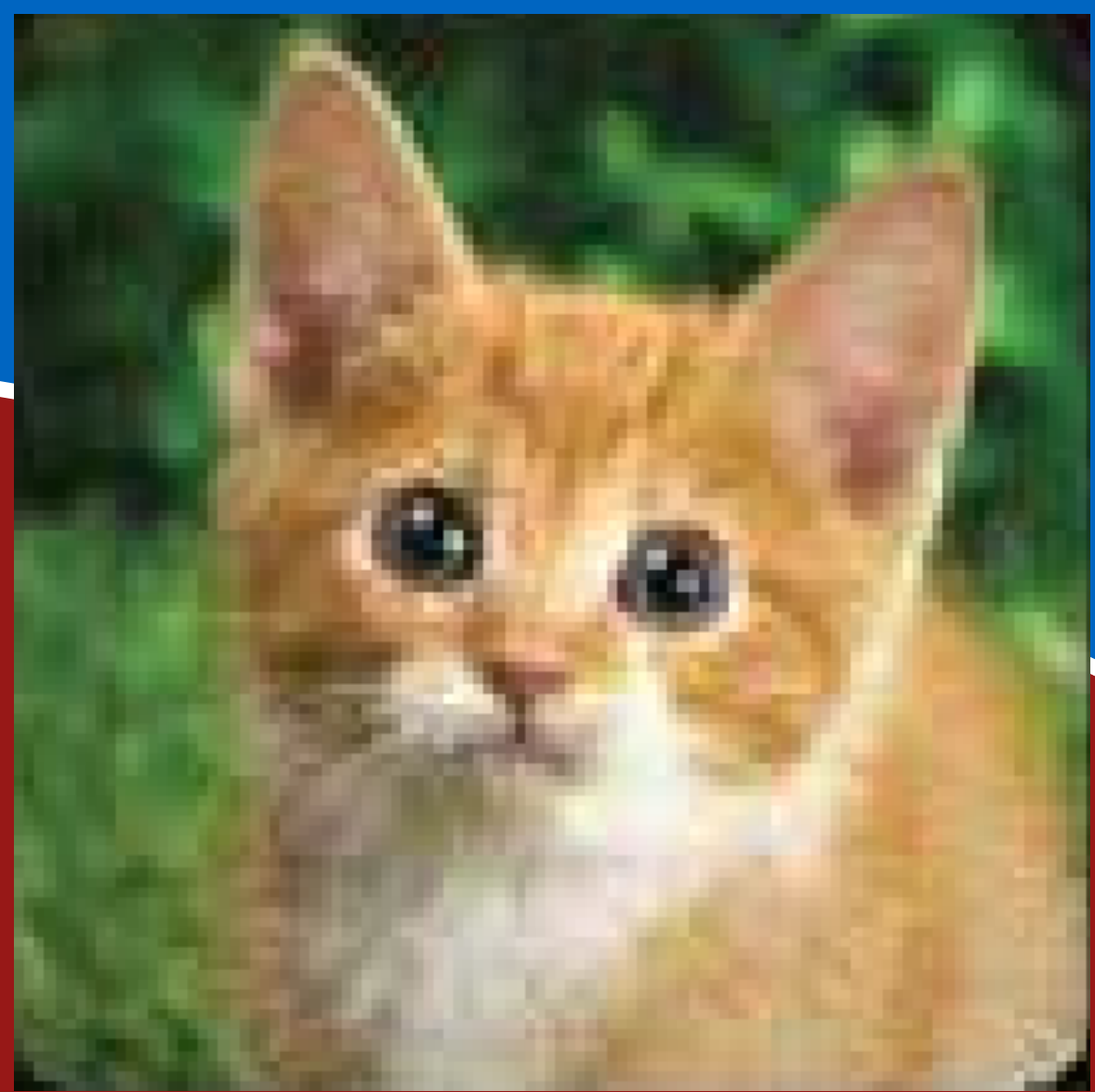












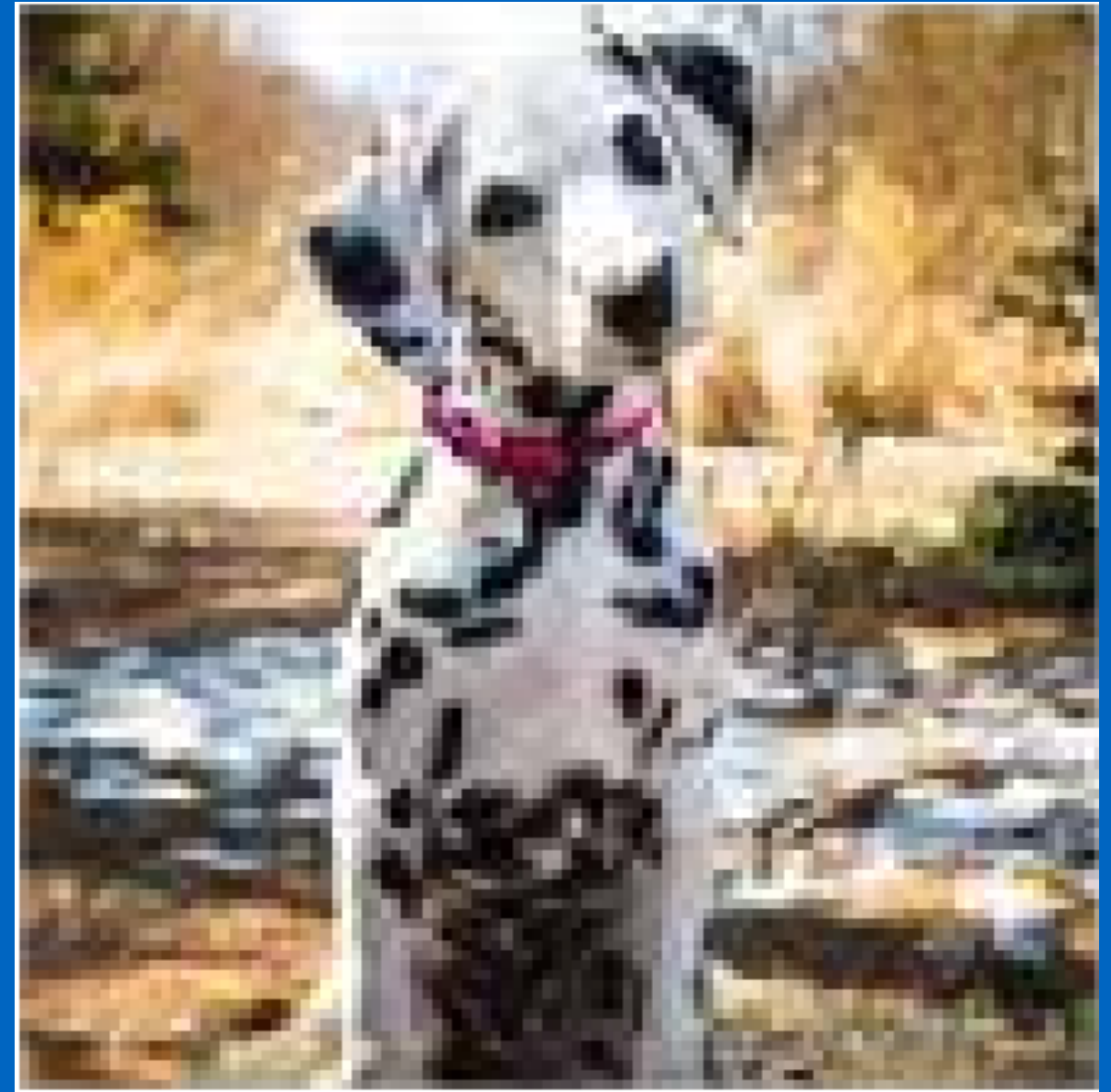
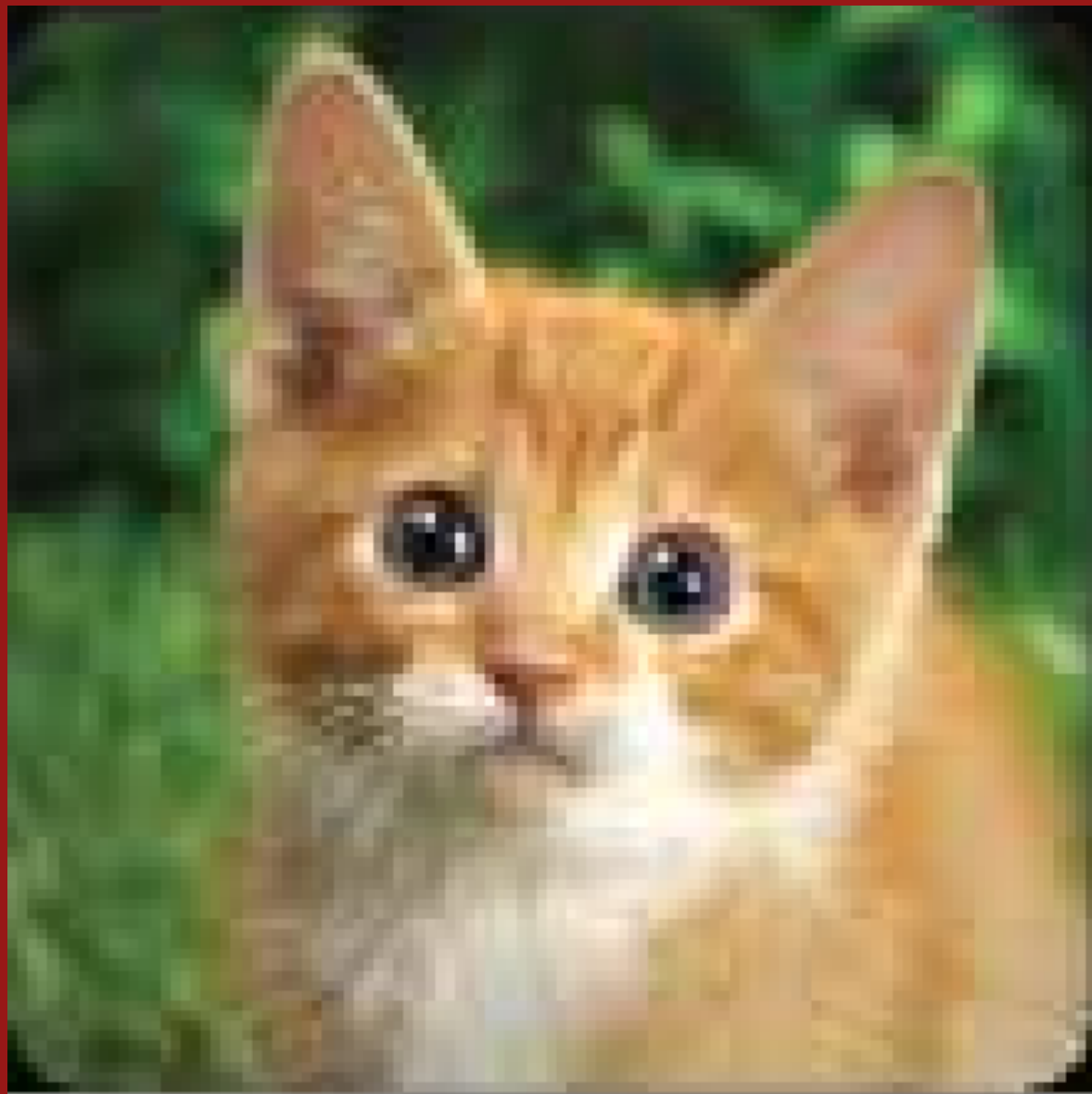
EXCESSIVE INVARIANCE CAUSES ADVERSARIAL VULNERABILITY

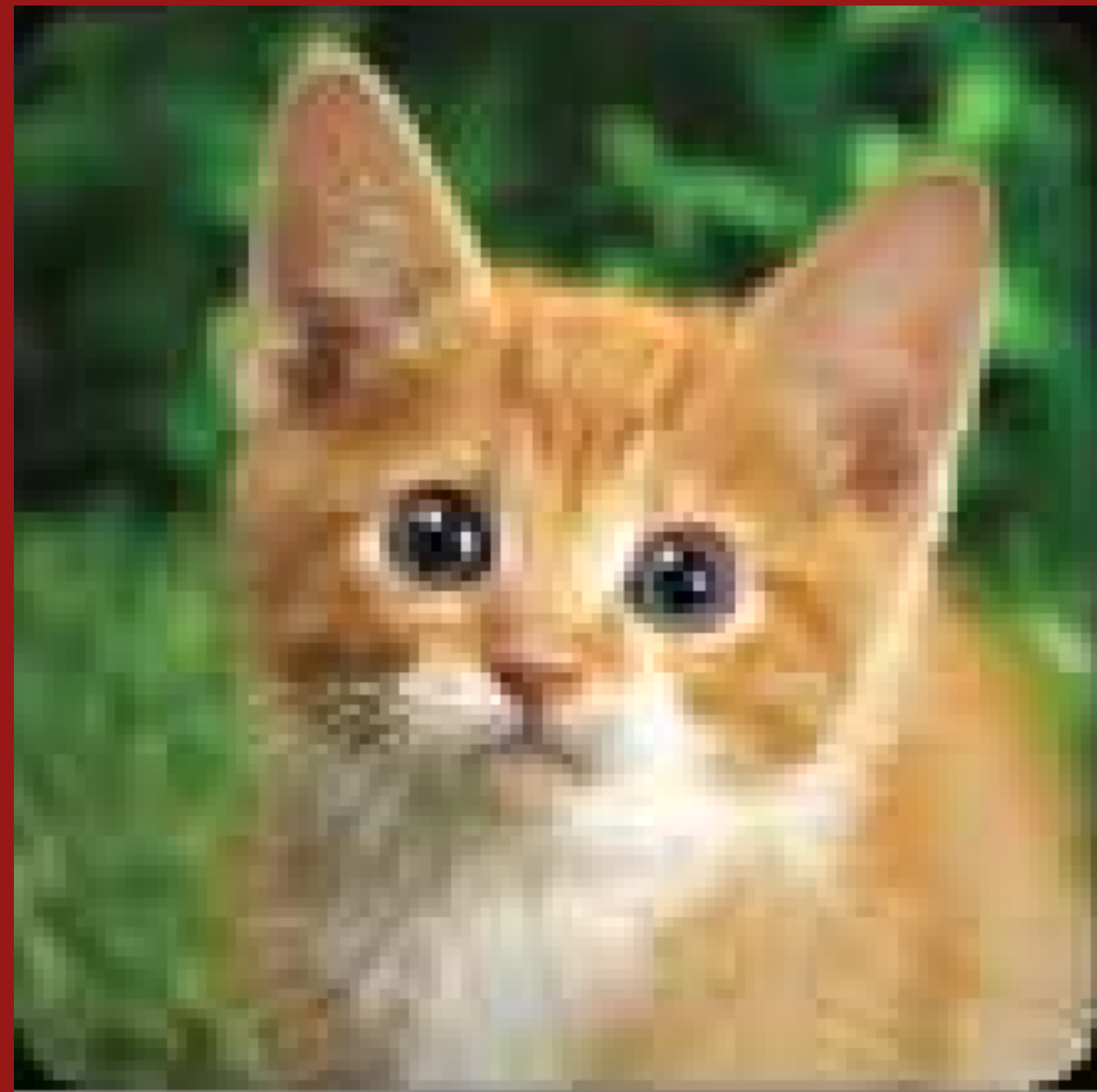
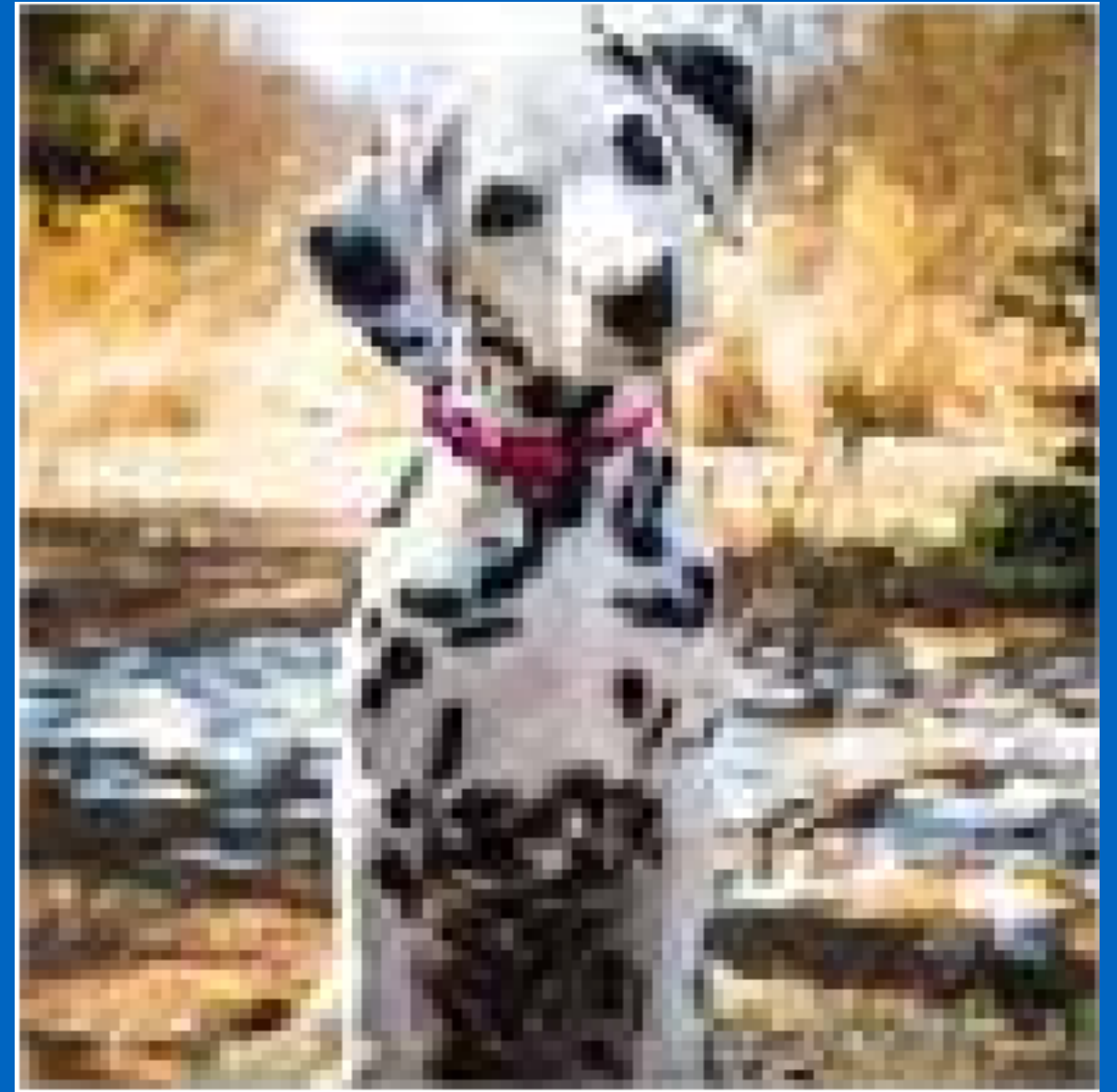
Jörn-Henrik Jacobsen^{1*}, Jens Behrmann^{1,2}, Richard Zemel¹, Matthias Bethge³

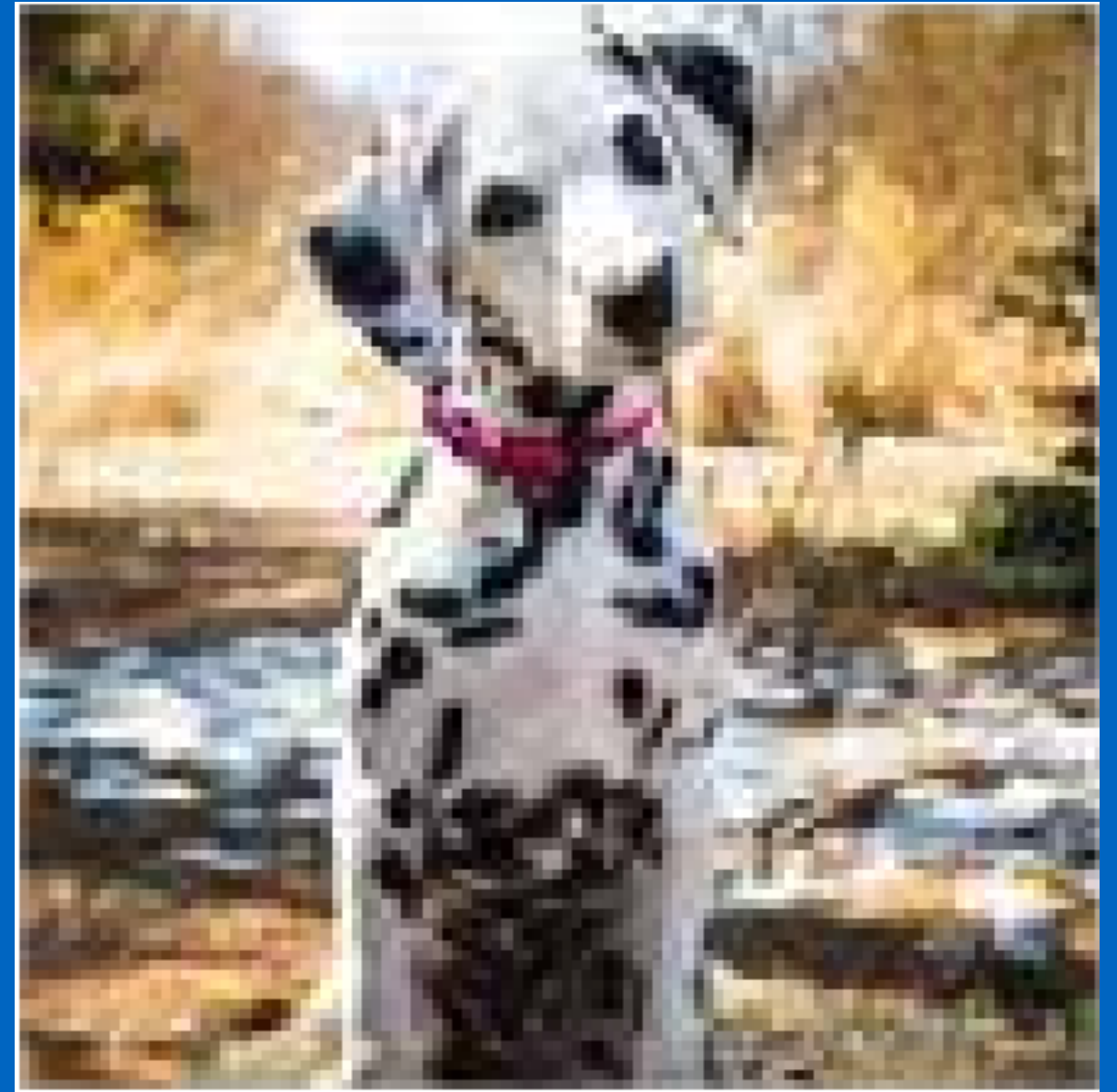
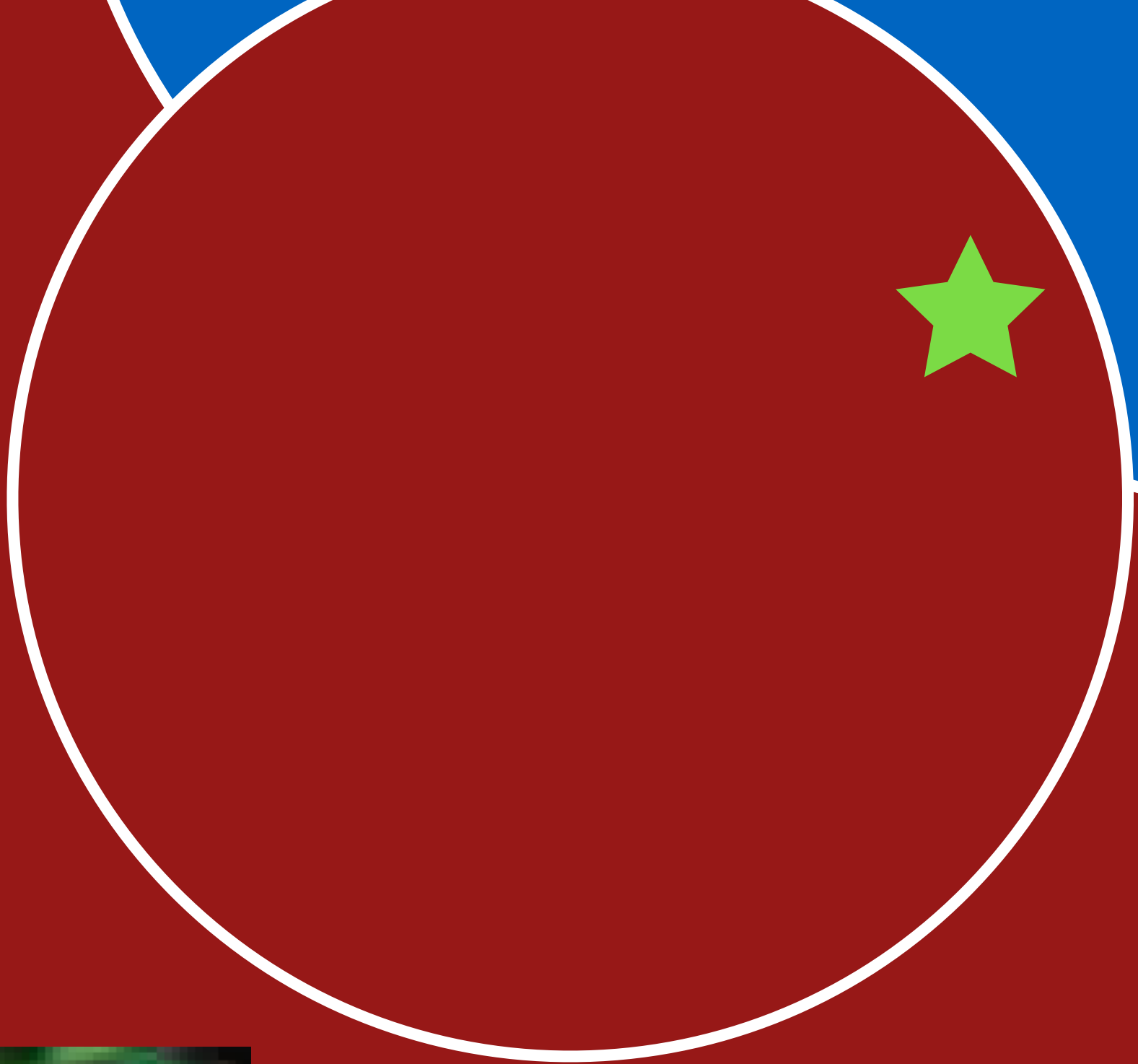
¹Vector Institute and University of Toronto

²University of Bremen, Center for Industrial Mathematics

³University of Tübingen

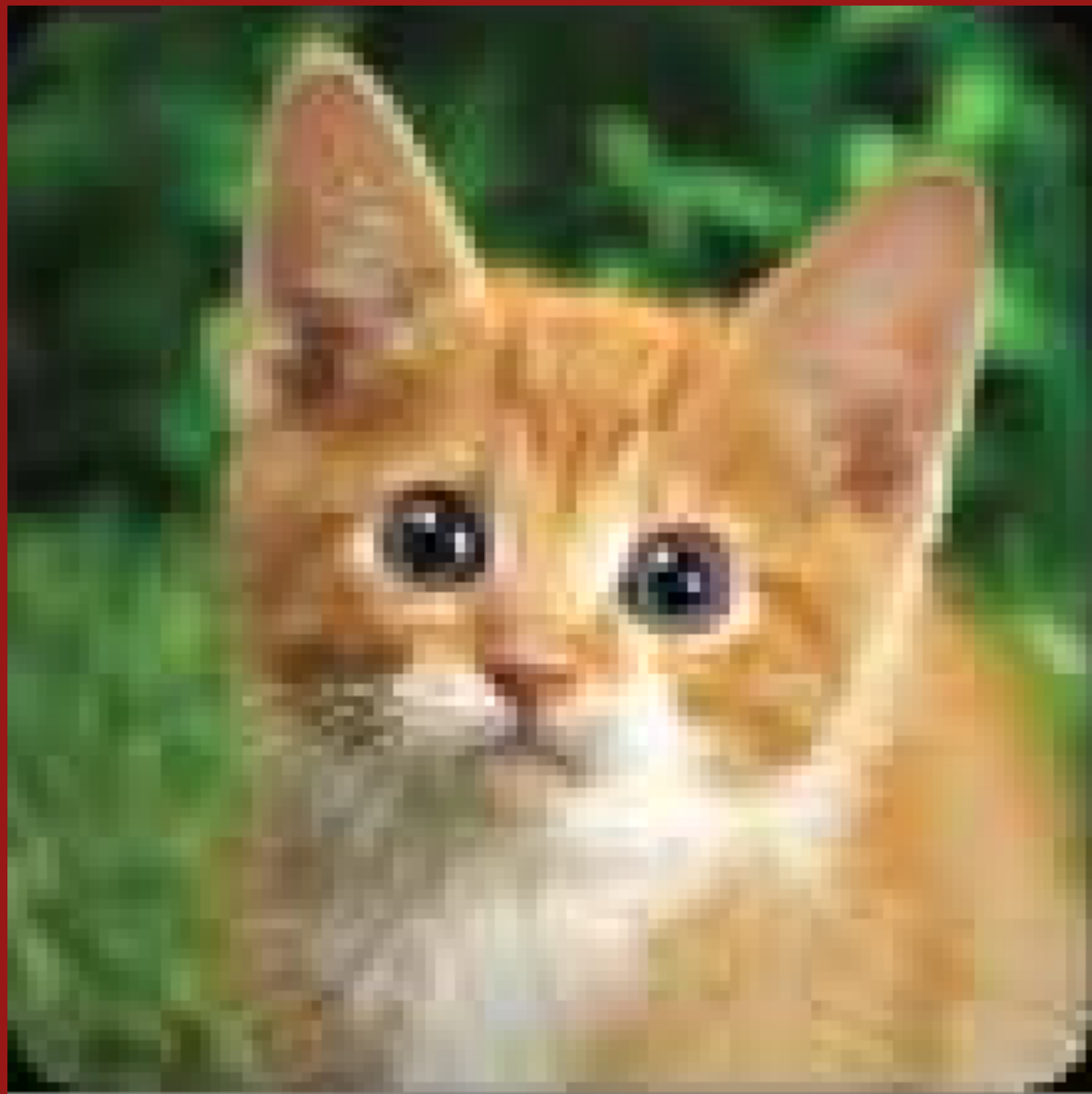
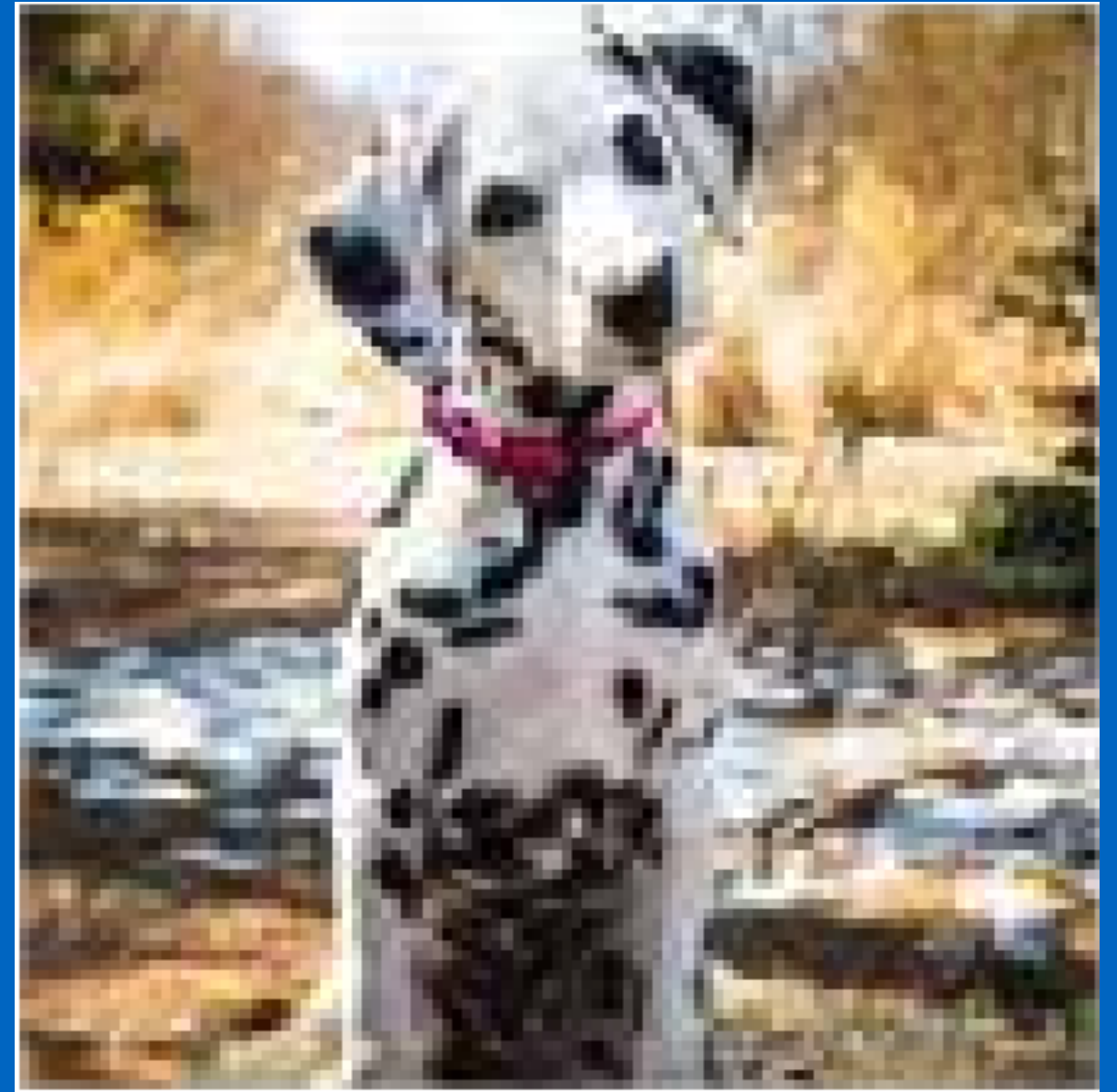


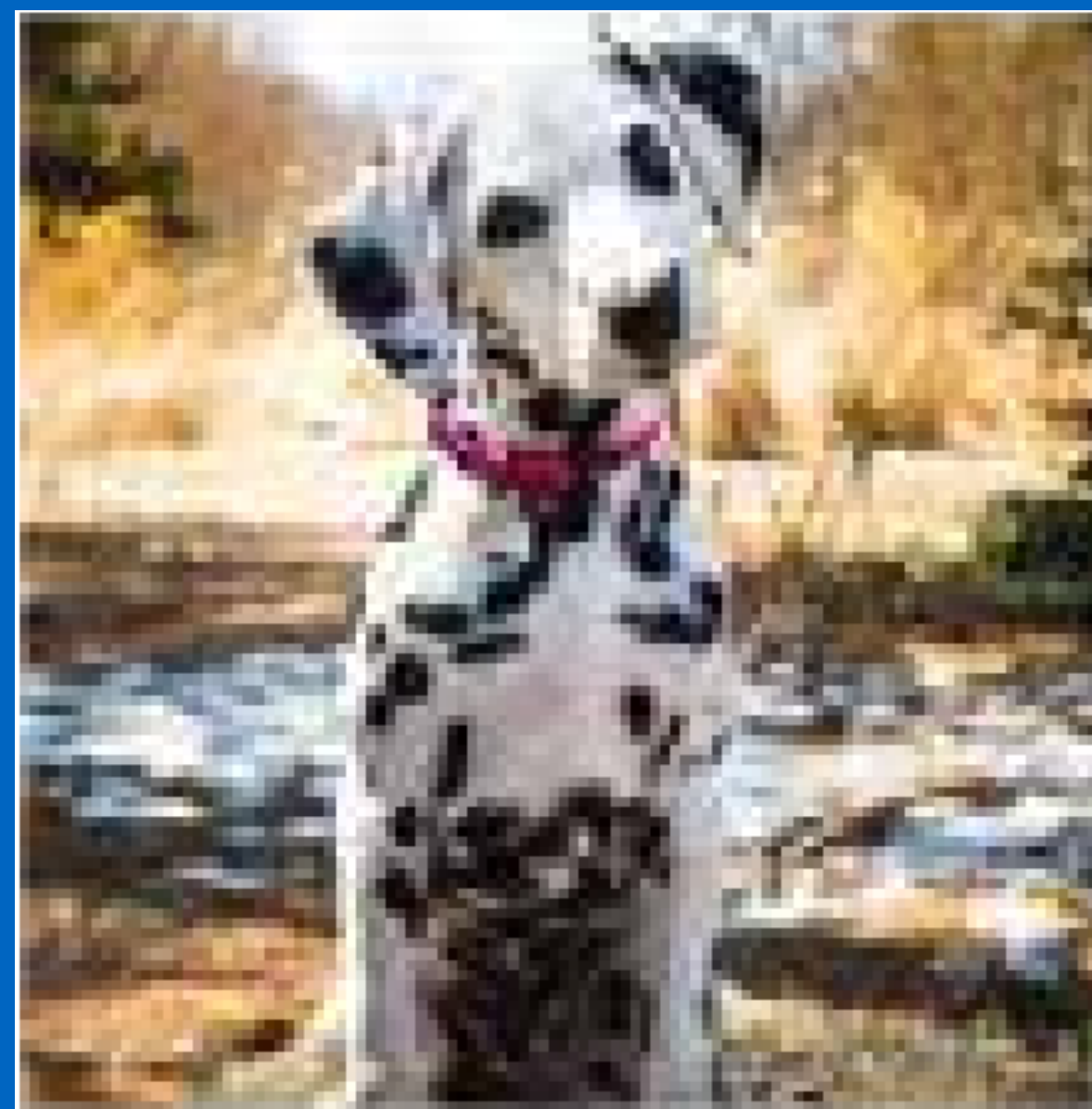
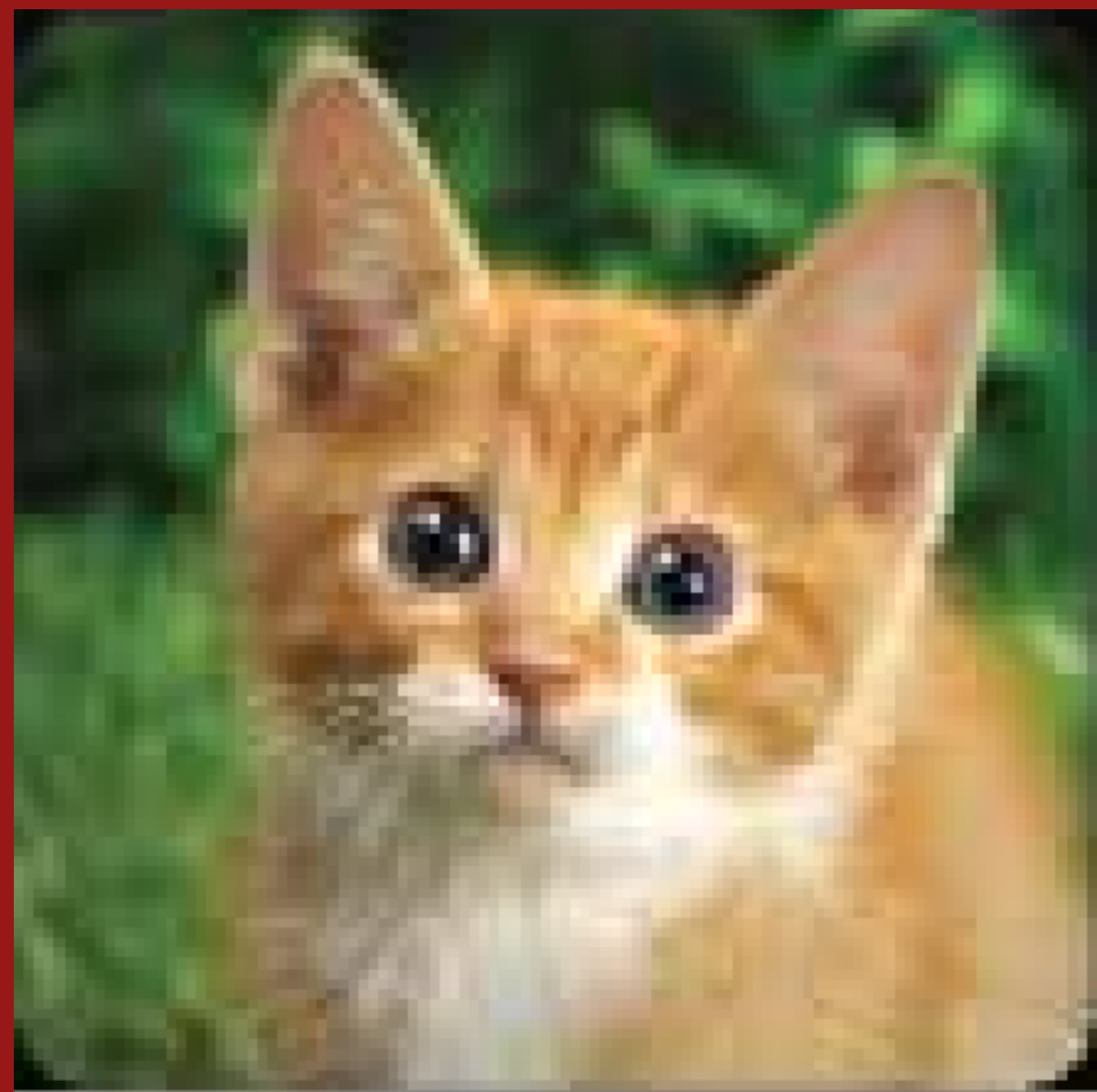








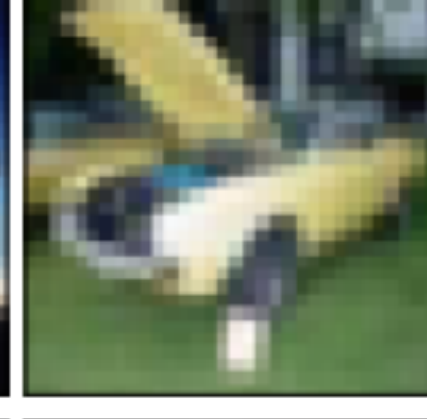
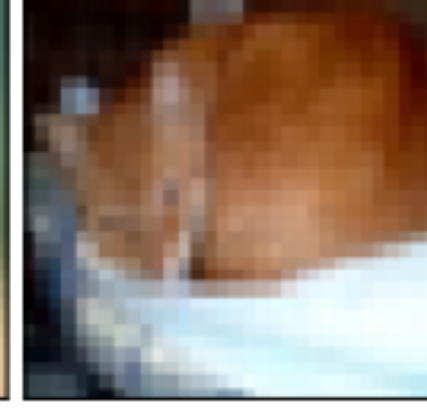
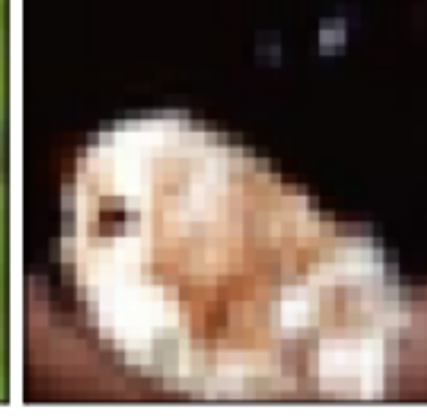
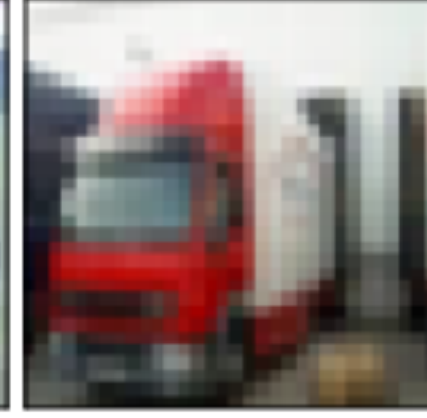
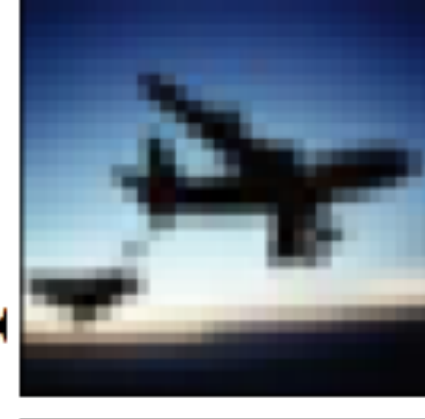
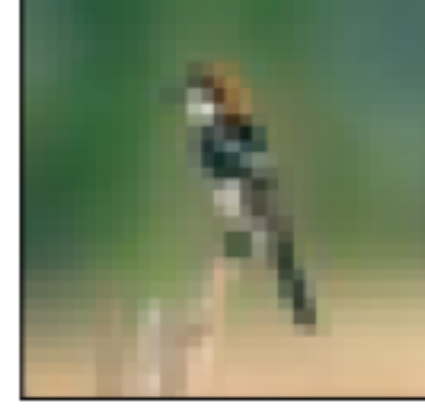
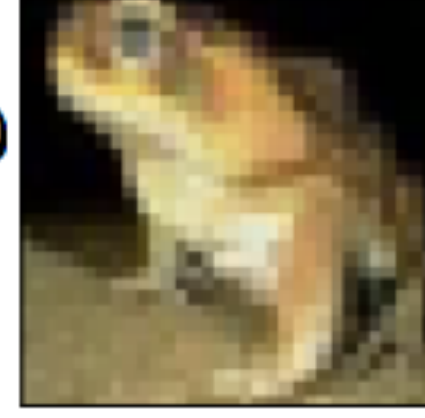




Wasserstein Adversarial Examples via Projected Sinkhorn Iterations

Eric Wong¹ Frank R. Schmidt² J. Zico Kolter^{3,4}

ship frog deer bird plane



truck

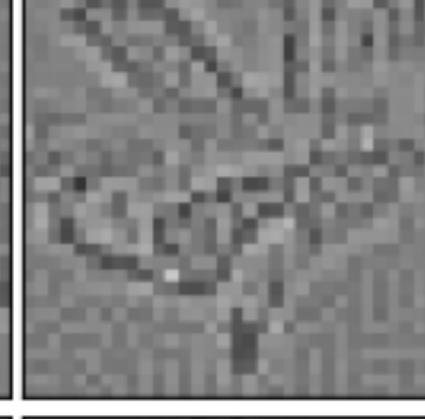
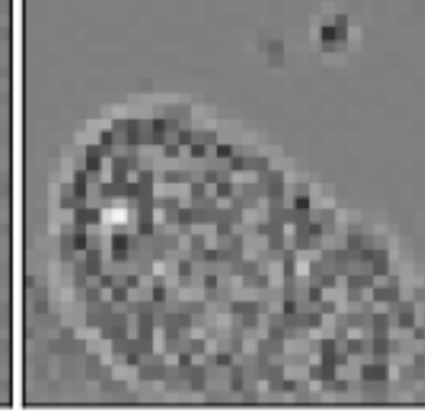
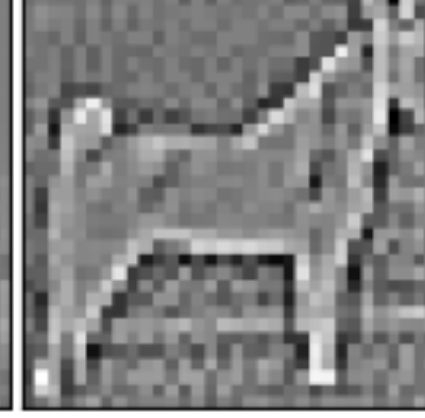
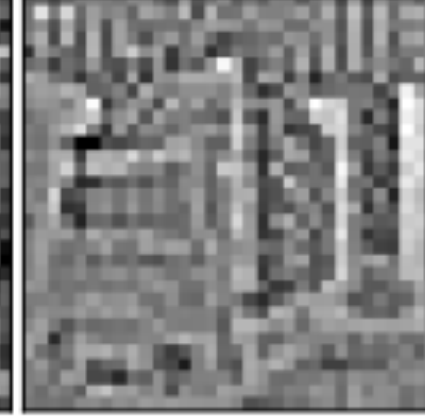
horse

dog

cat

car

+

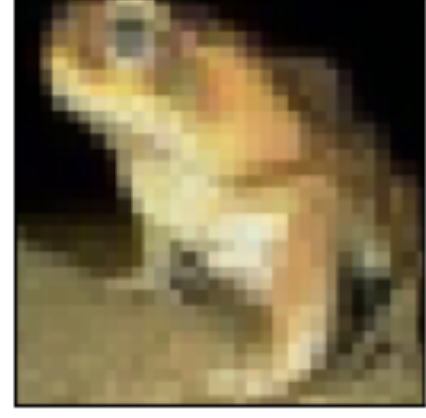


=

plane



bird



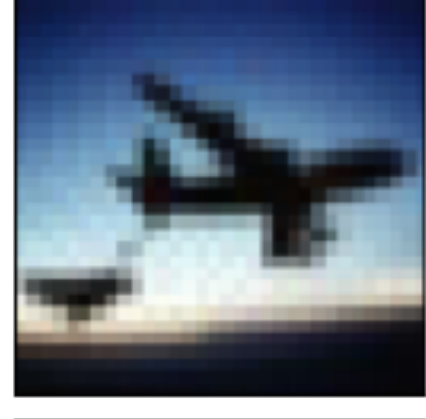
horse



deer



bird



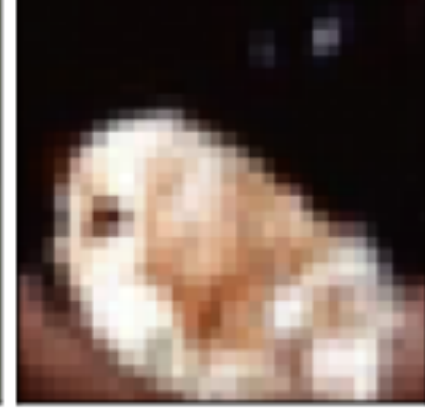
car



dog



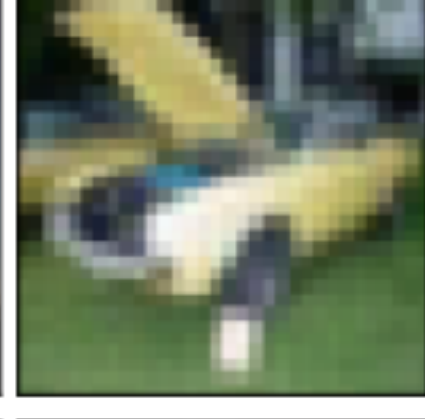
cat



dog



deer



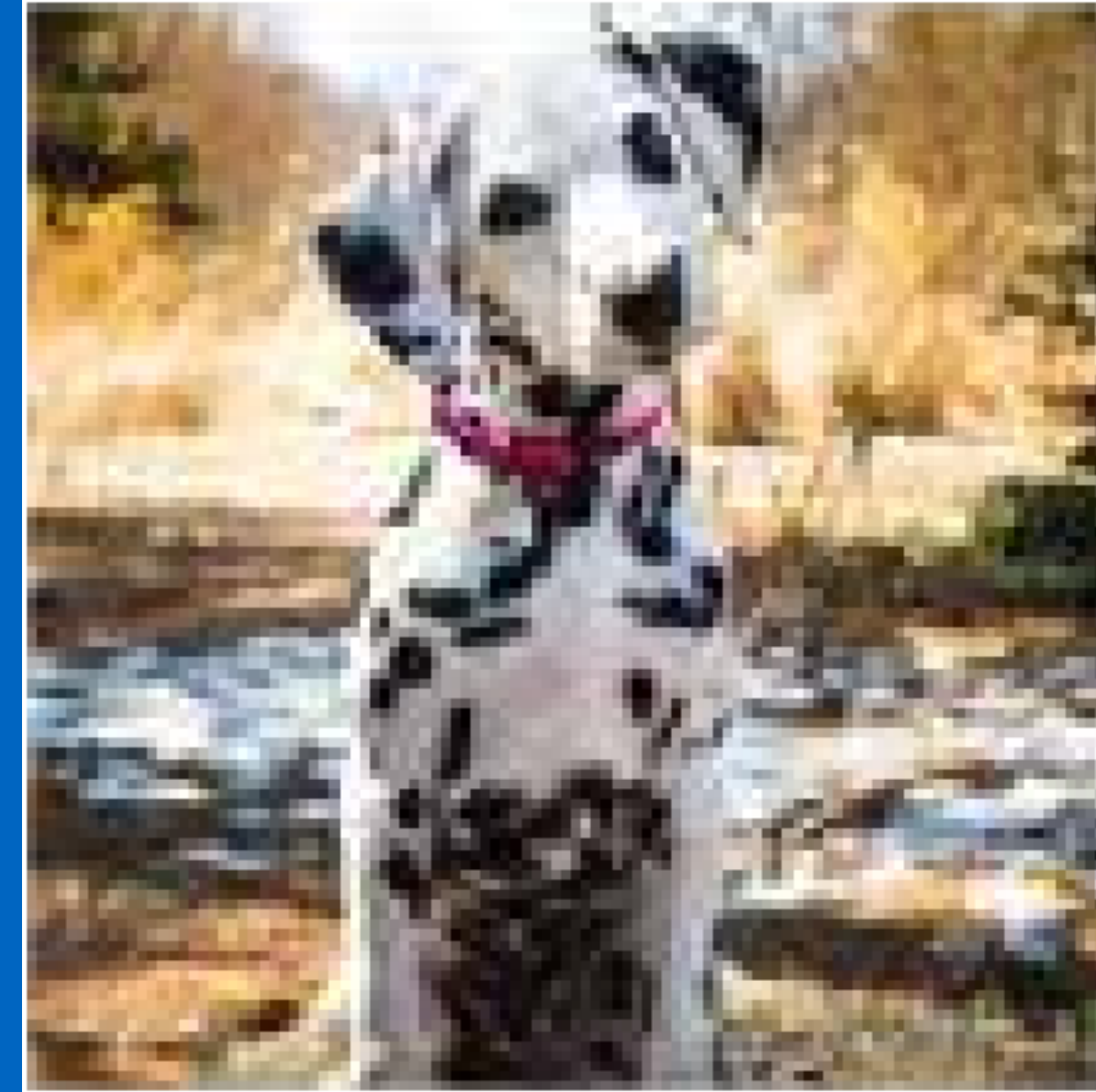
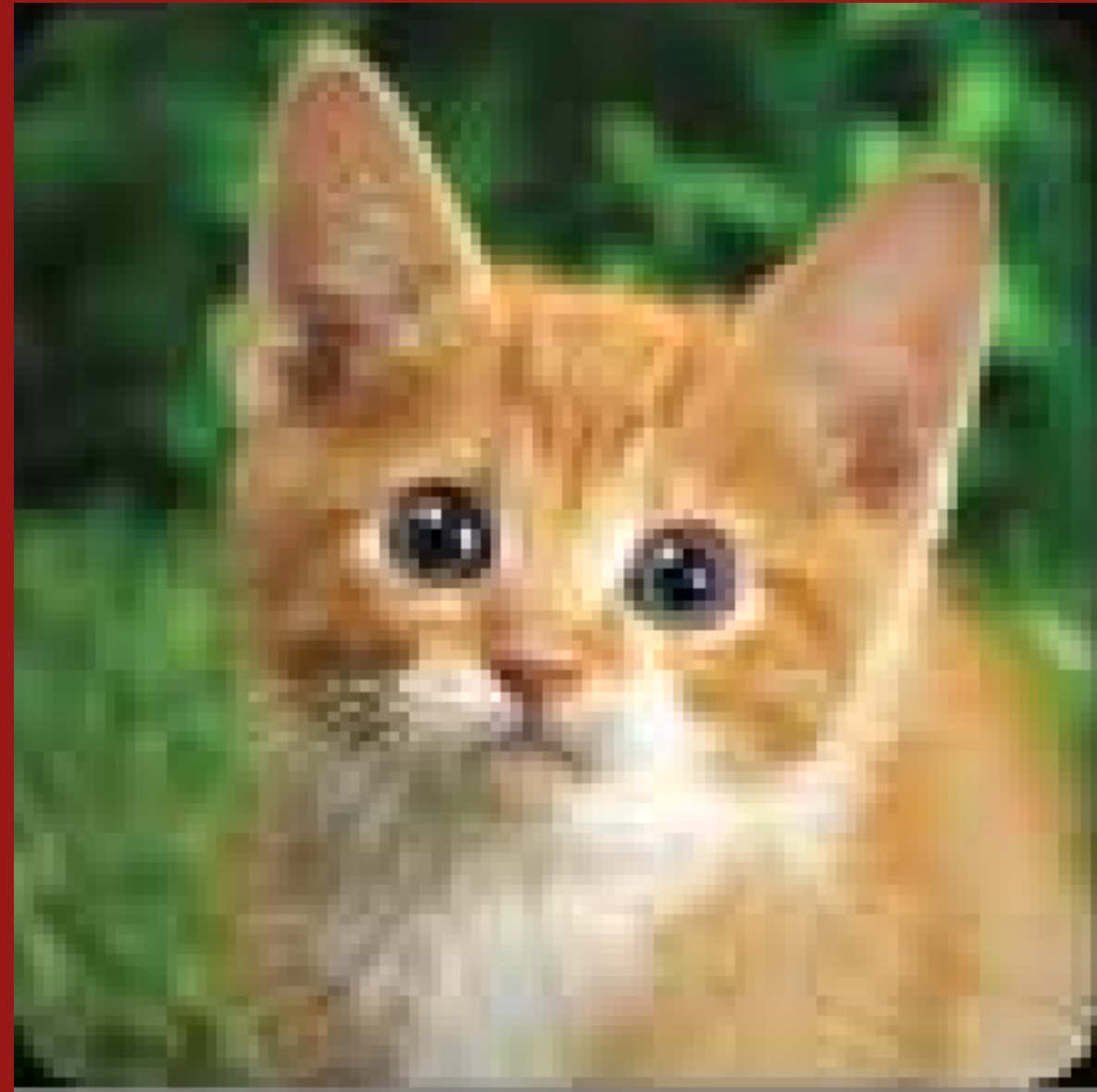
Recent advances in ...

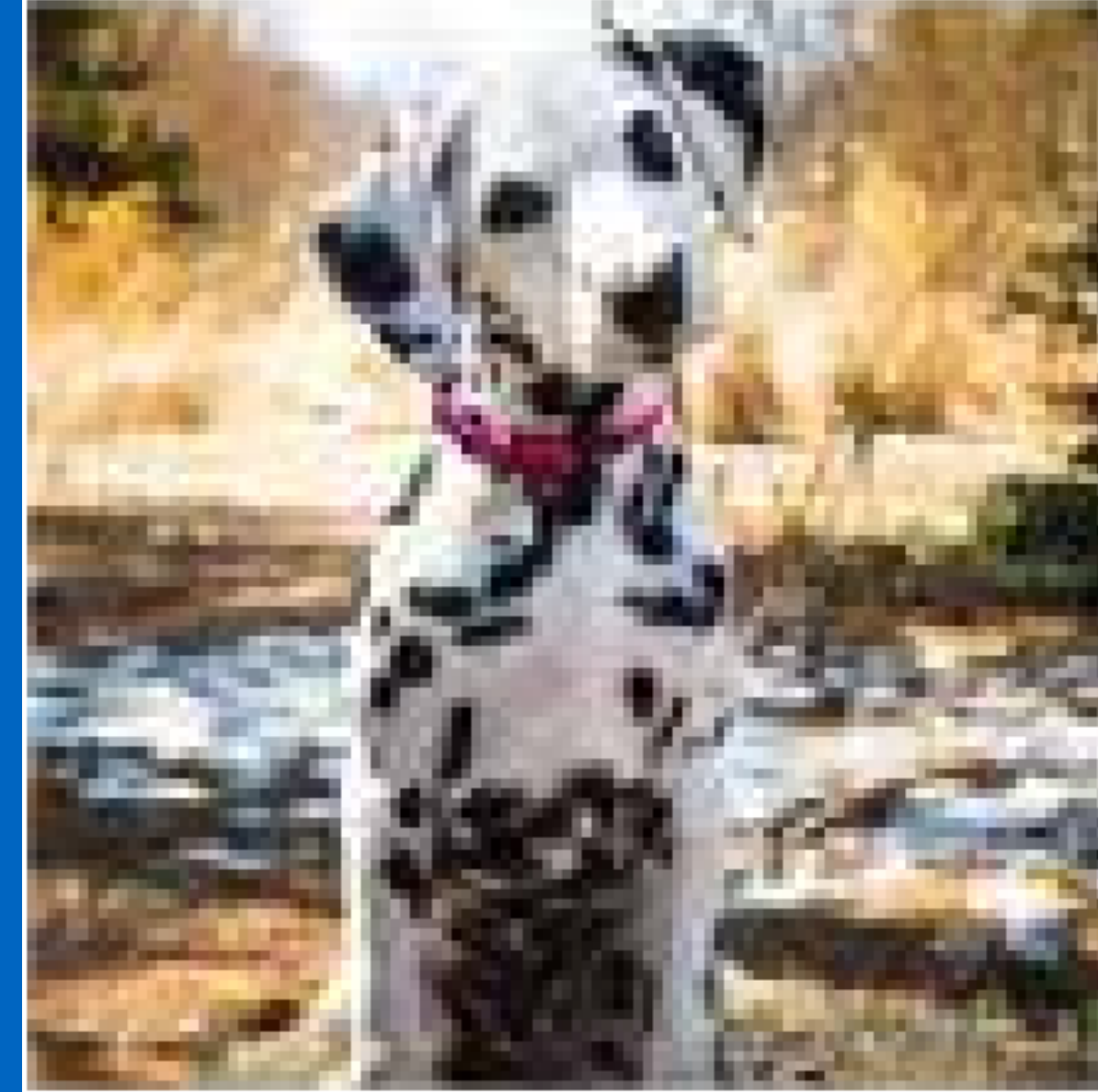
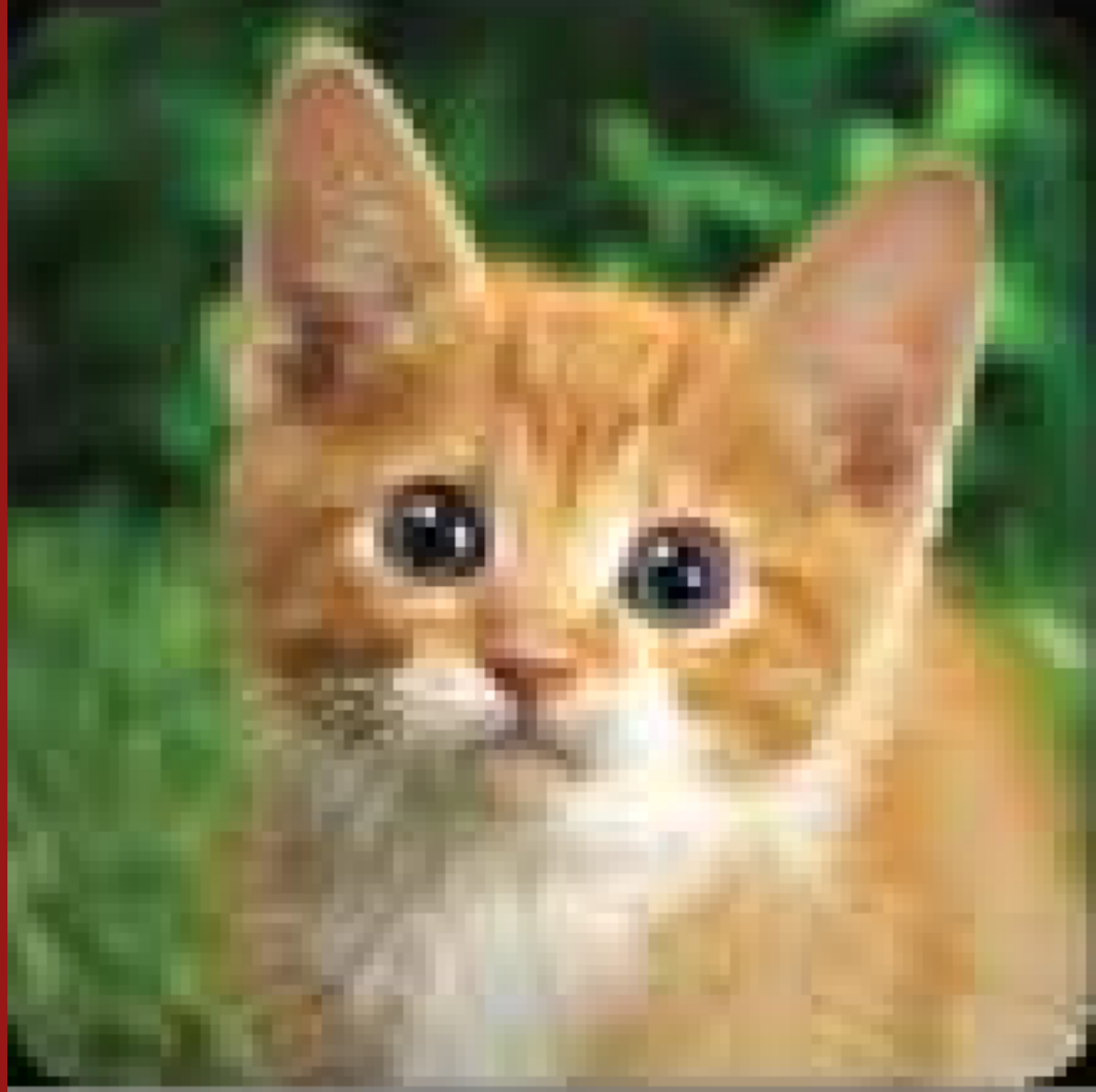
Defending Against
Adversarial Examples

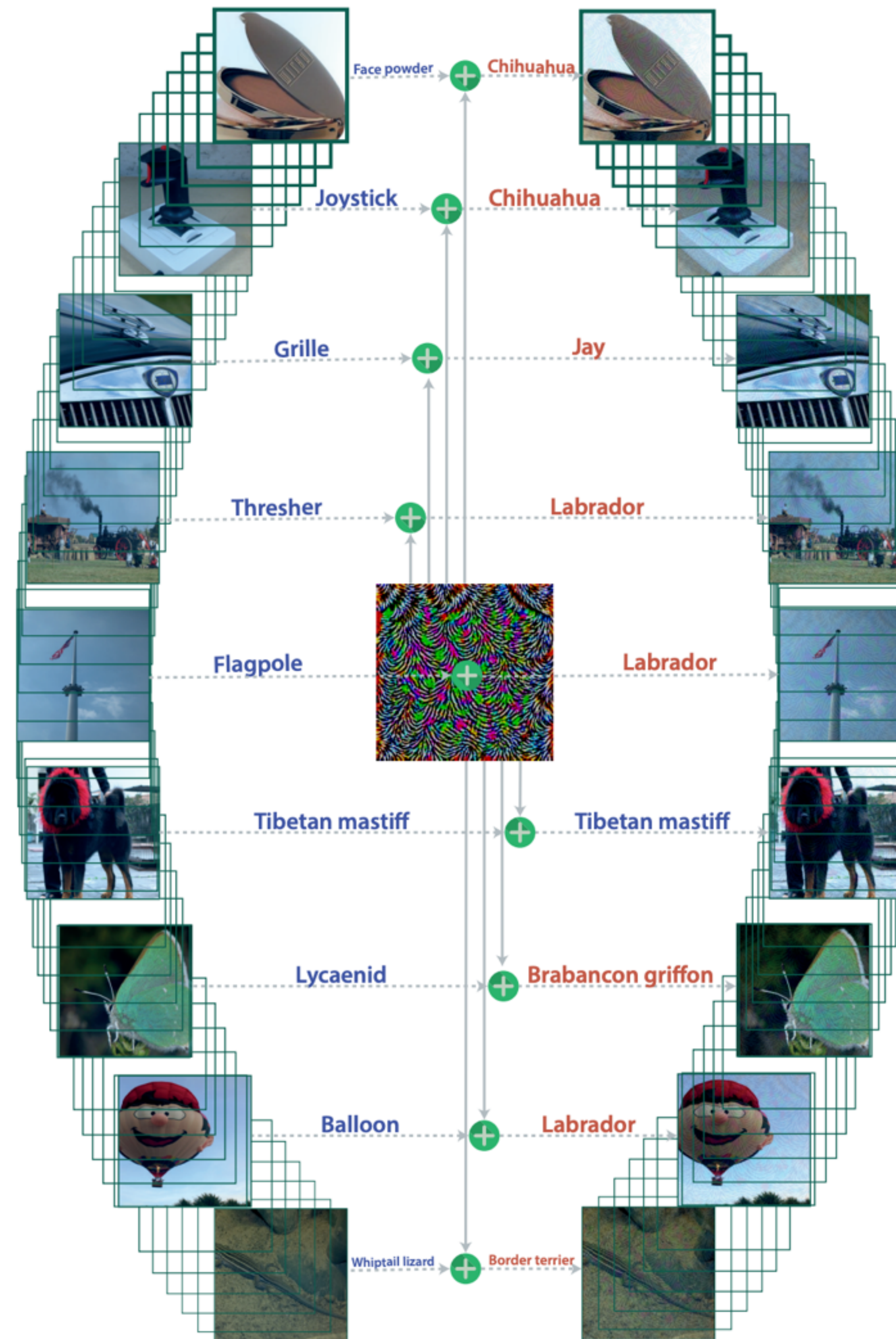
Defenses I *don't*
believe will be effective

... a bit more
background

Transferability







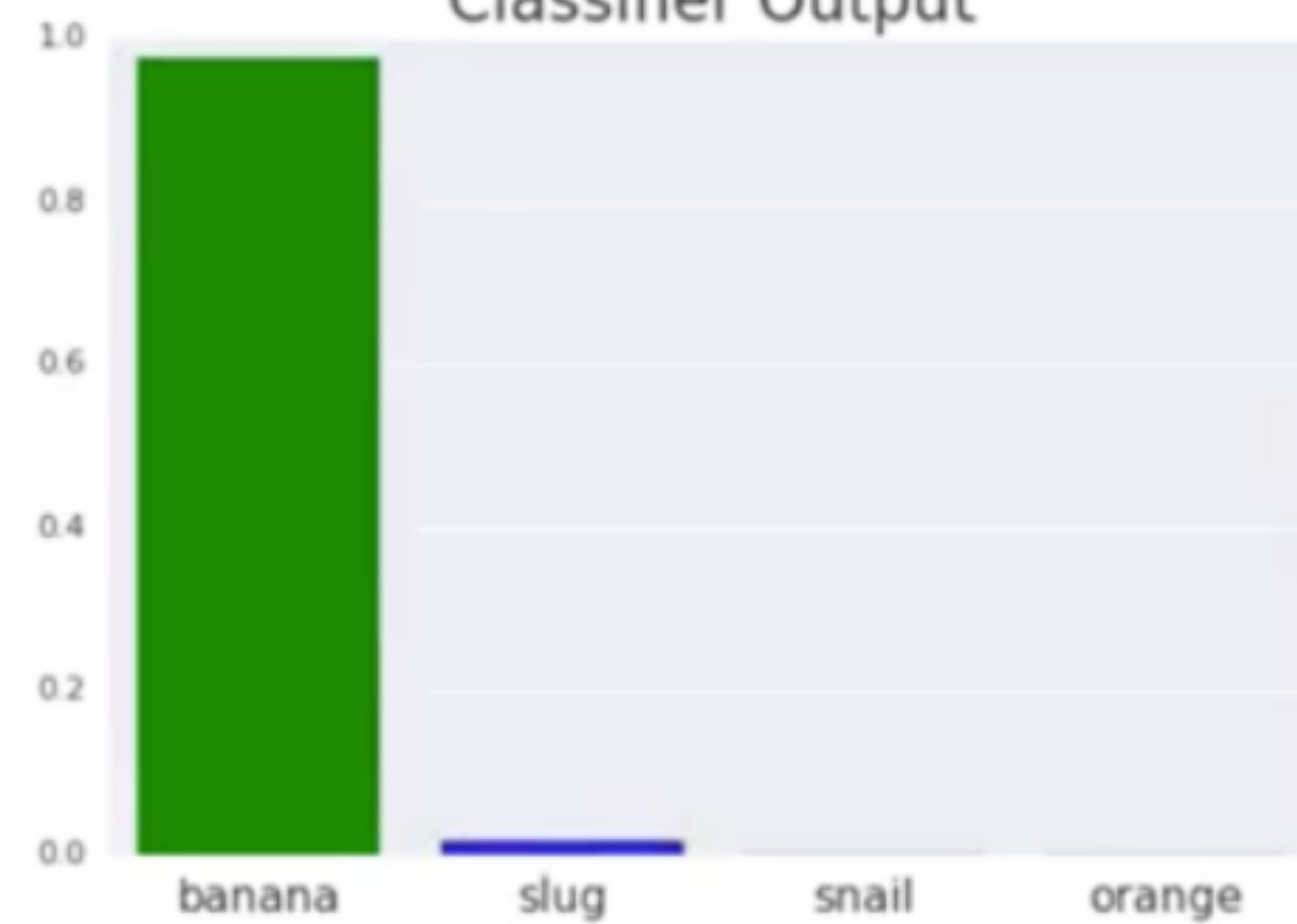
place sticker on table



Classifier Input



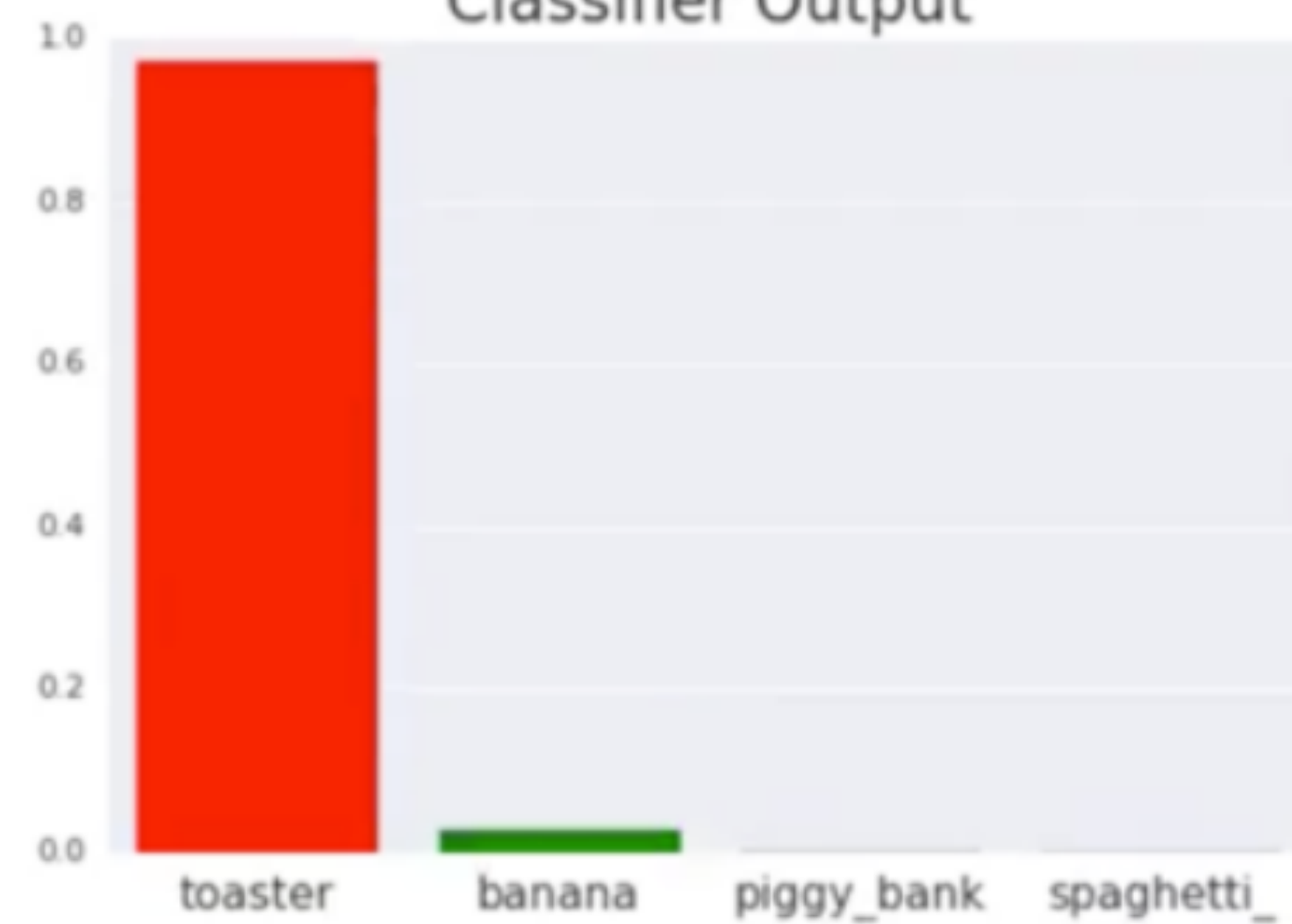
Classifier Output



Classifier Input



Classifier Output



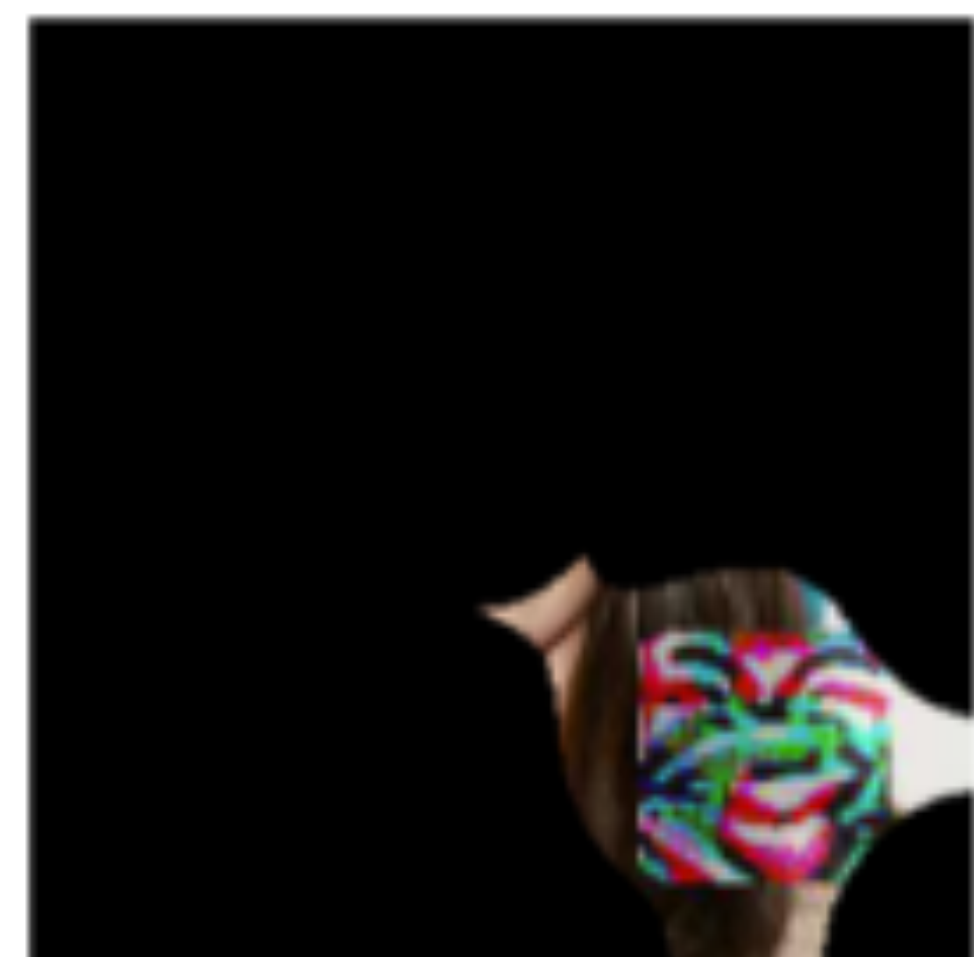
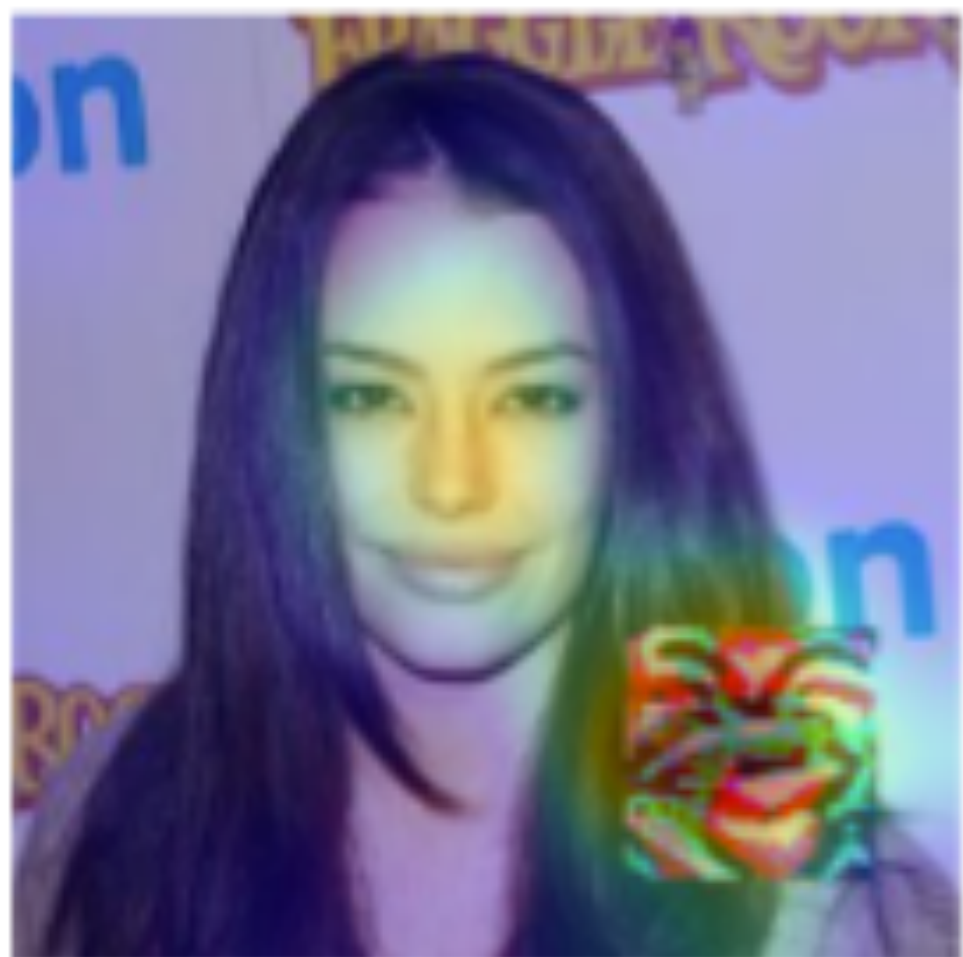
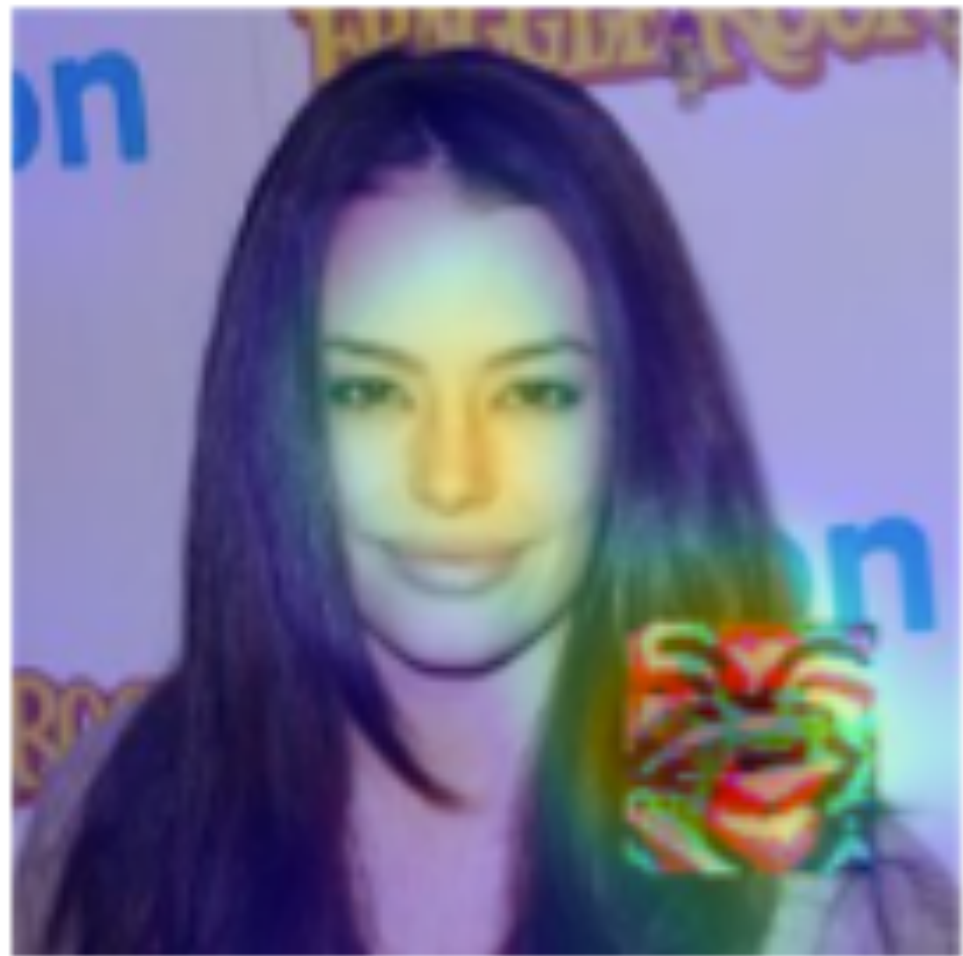
SentiNet: Detecting Physical Attacks Against Deep Learning Systems

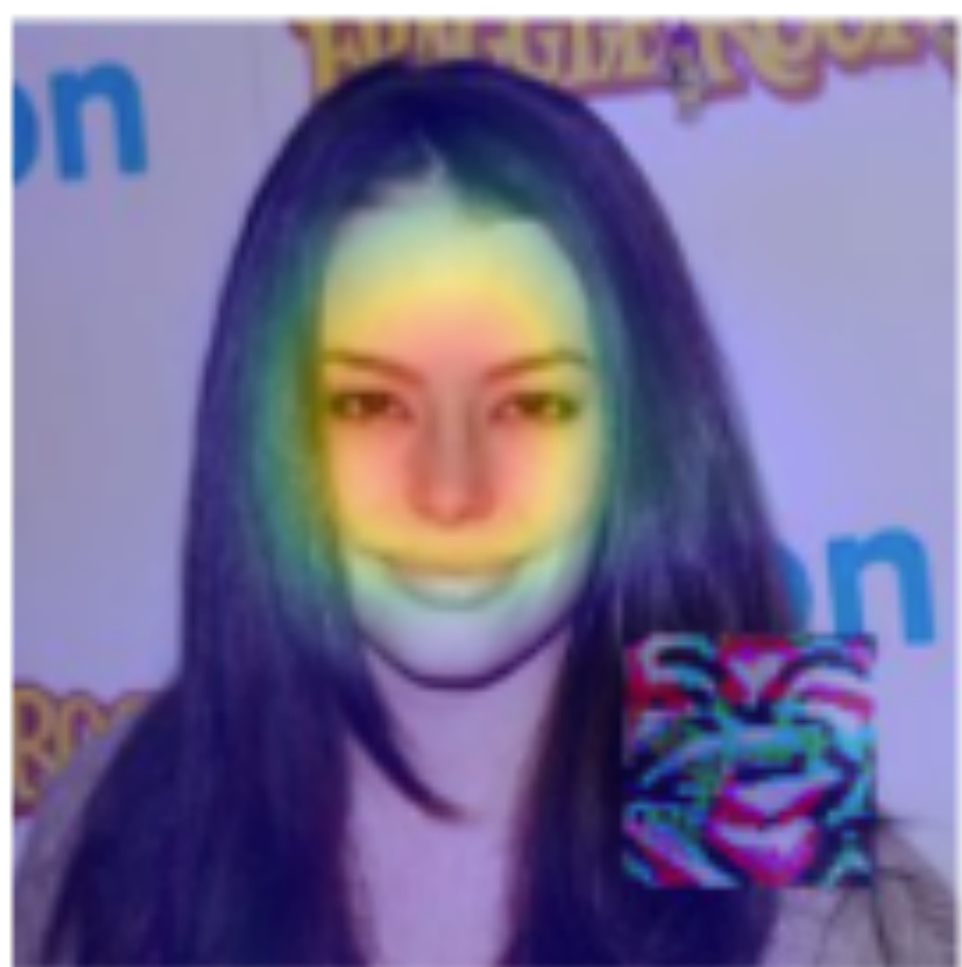
Edward Chou¹

Florian Tramèr¹

Giancarlo Pellegrino^{1,2}

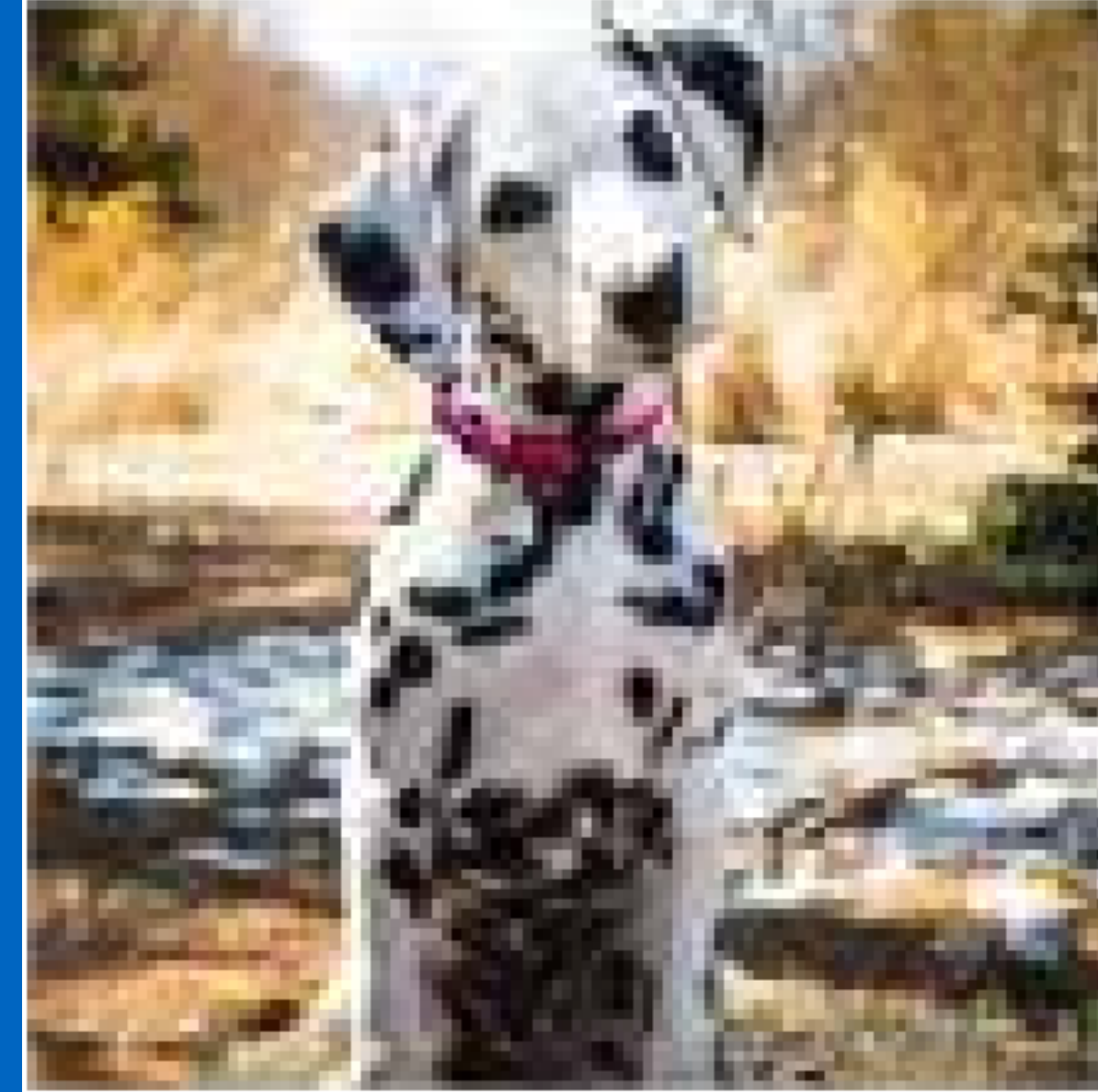
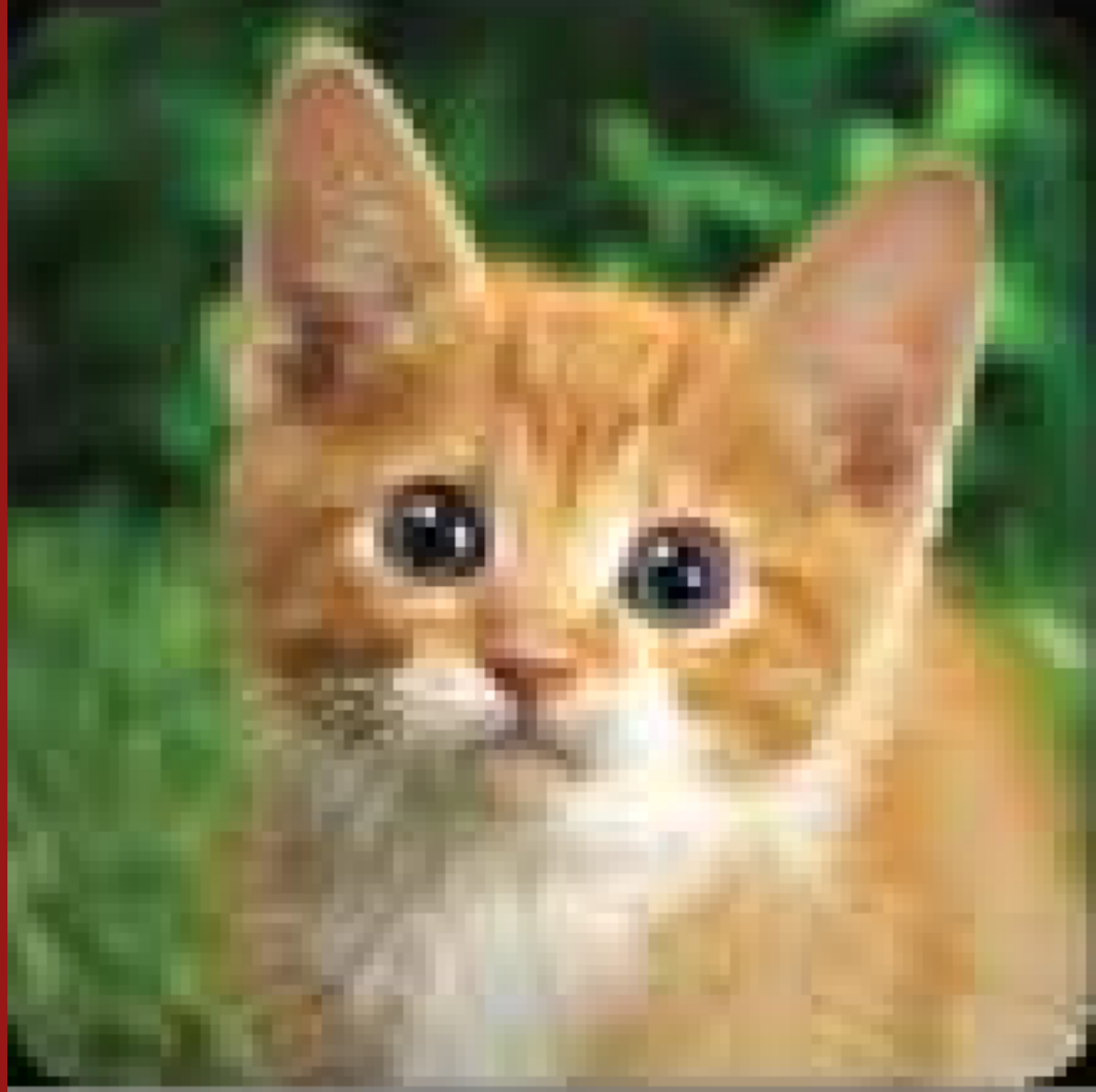
Dan Boneh¹

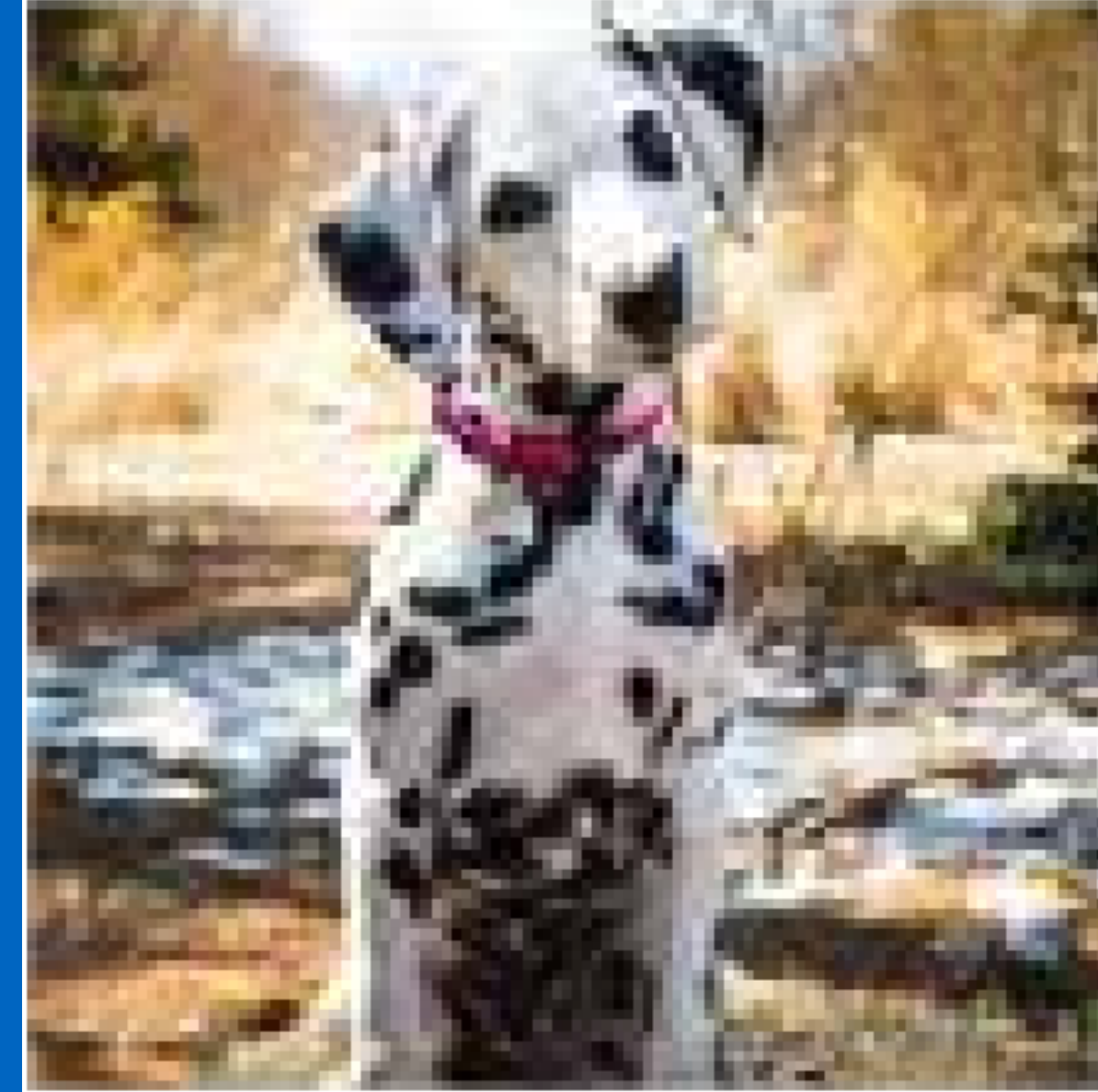
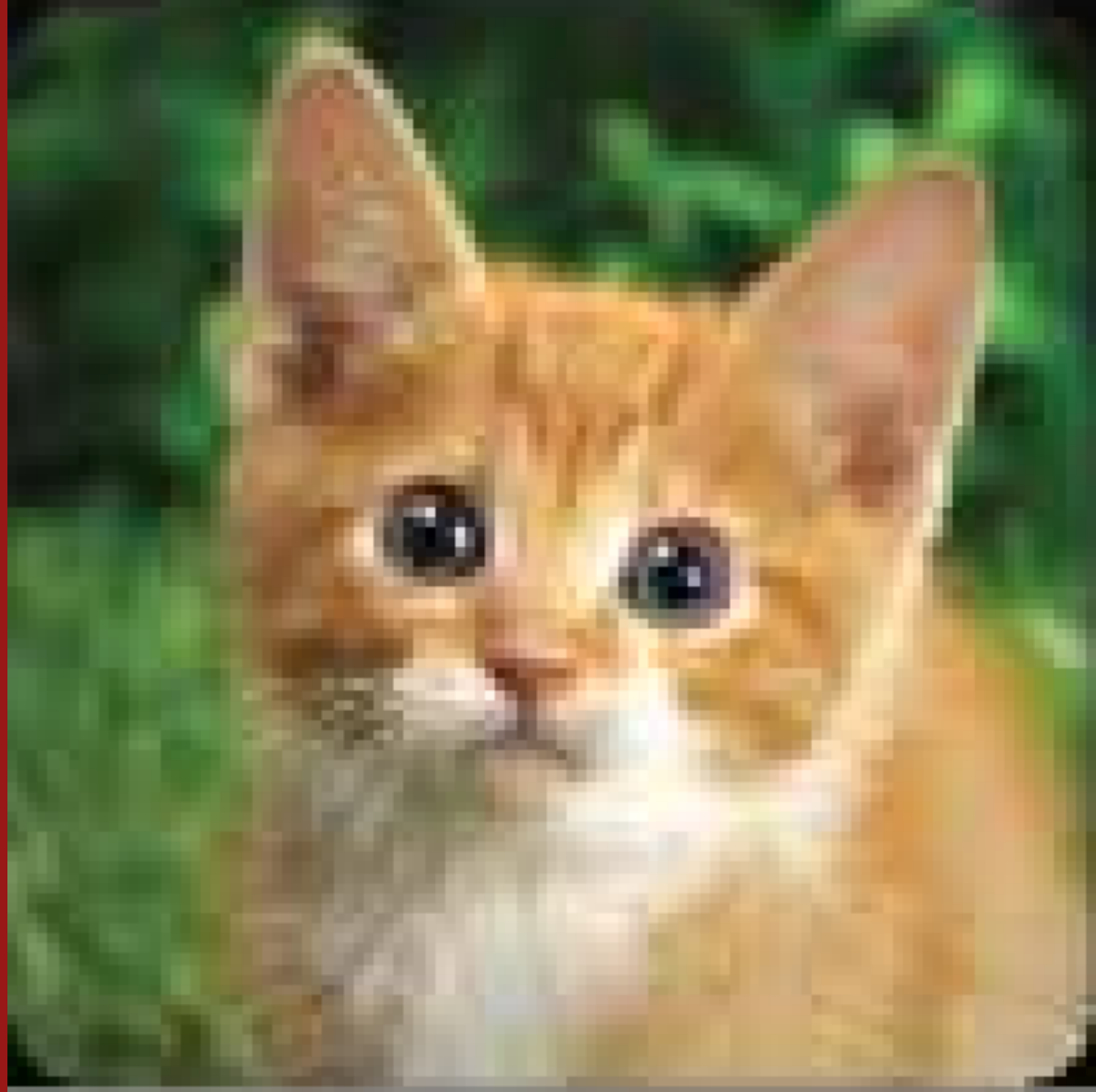




Sitatapatra: Blocking the Transfer of Adversarial Samples

Ilia Shumailov^{*1} Xitong Gao^{*2} Yiren Zhao^{*1} Robert Mullins¹ Ross Anderson¹ Cheng-Zhong Xu²





Stateful Detection of Black-Box Adversarial Attacks

Steven Chen

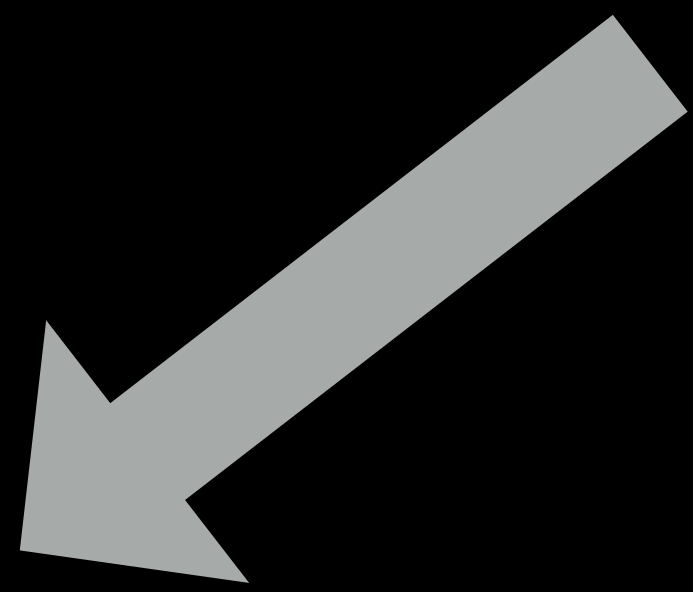
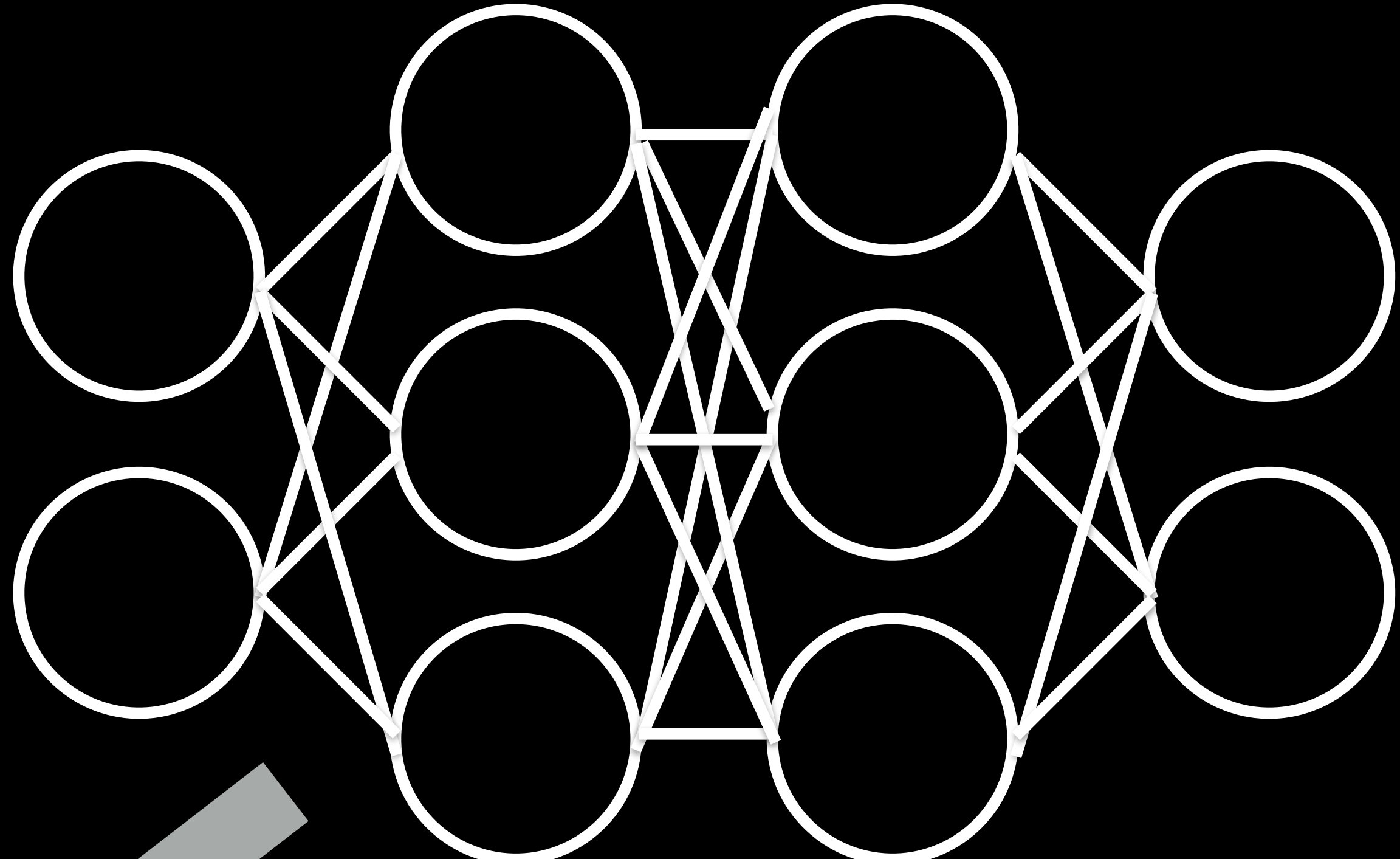
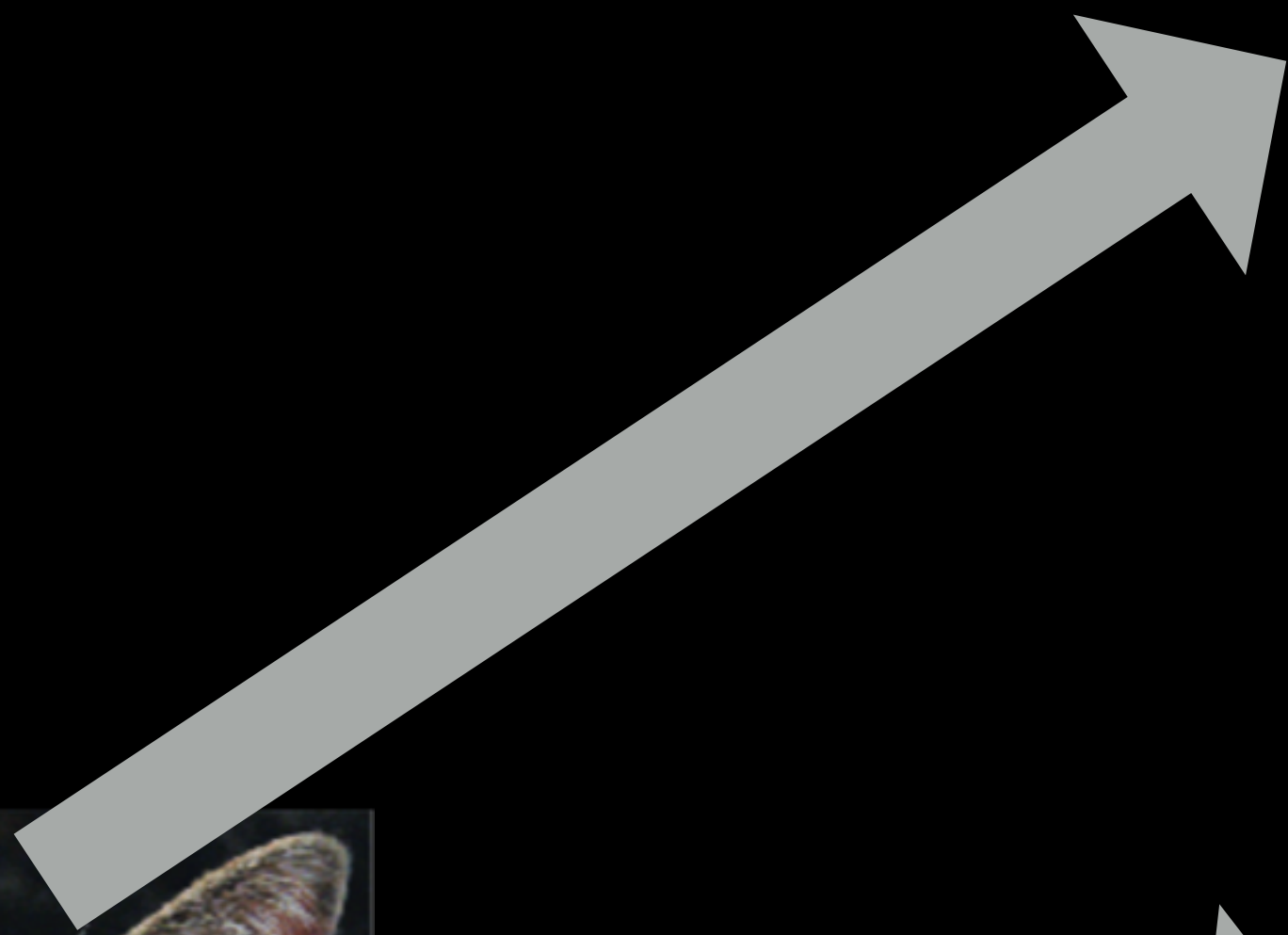
University of California, Berkeley

Nicholas Carlini

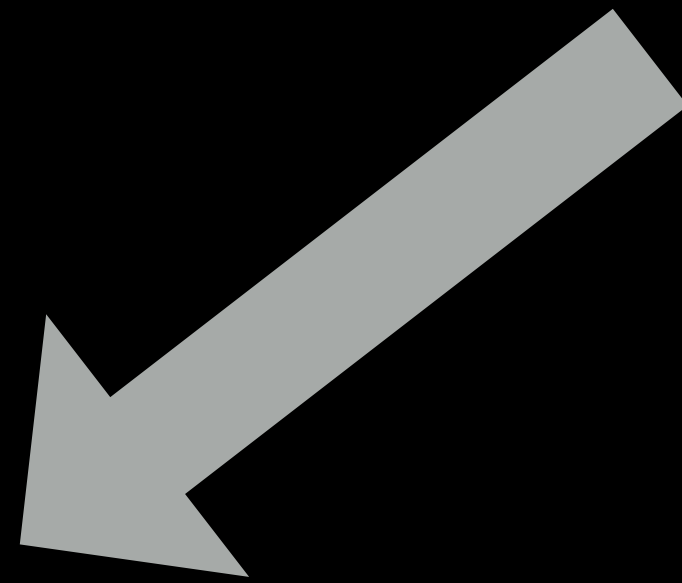
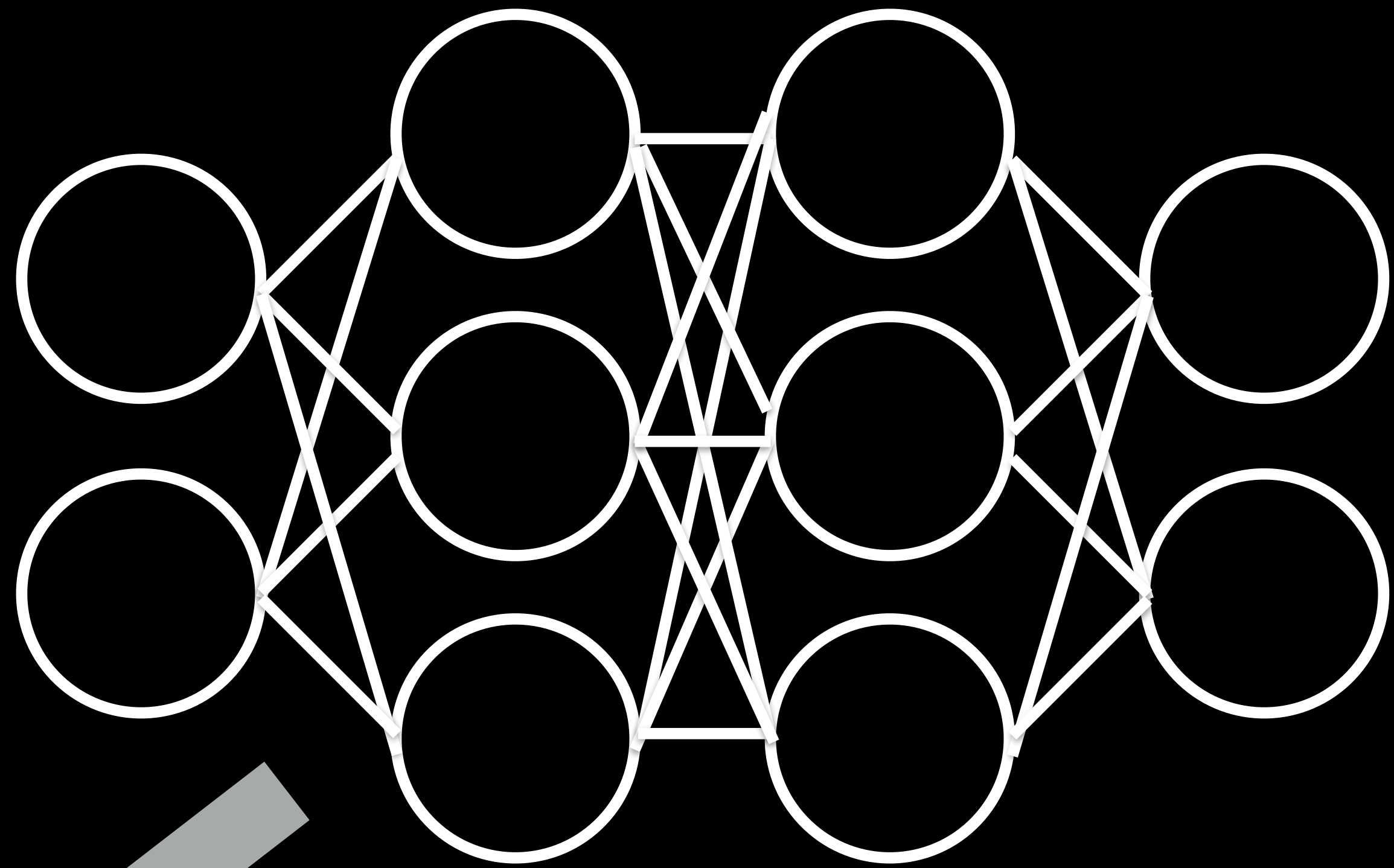
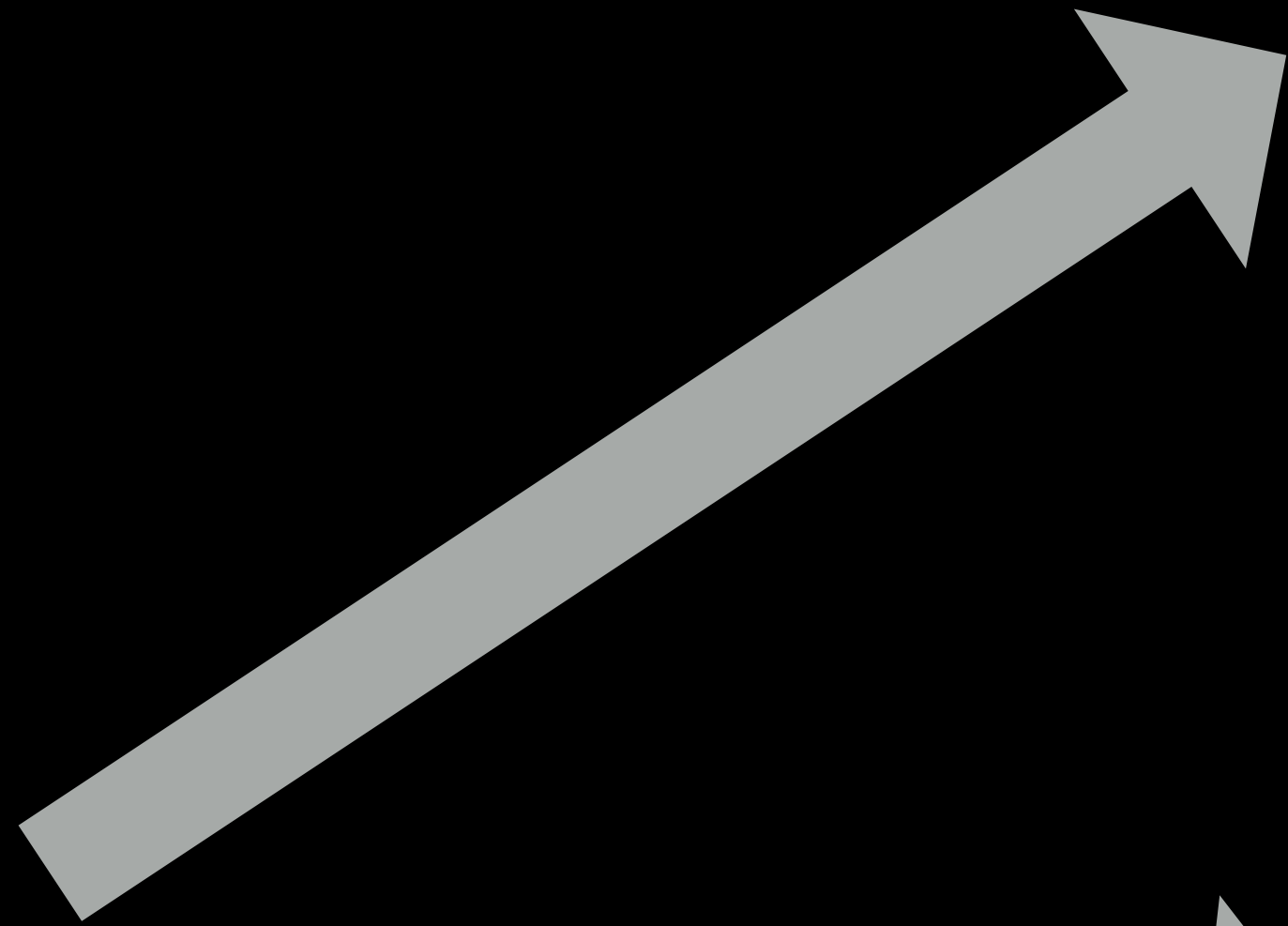
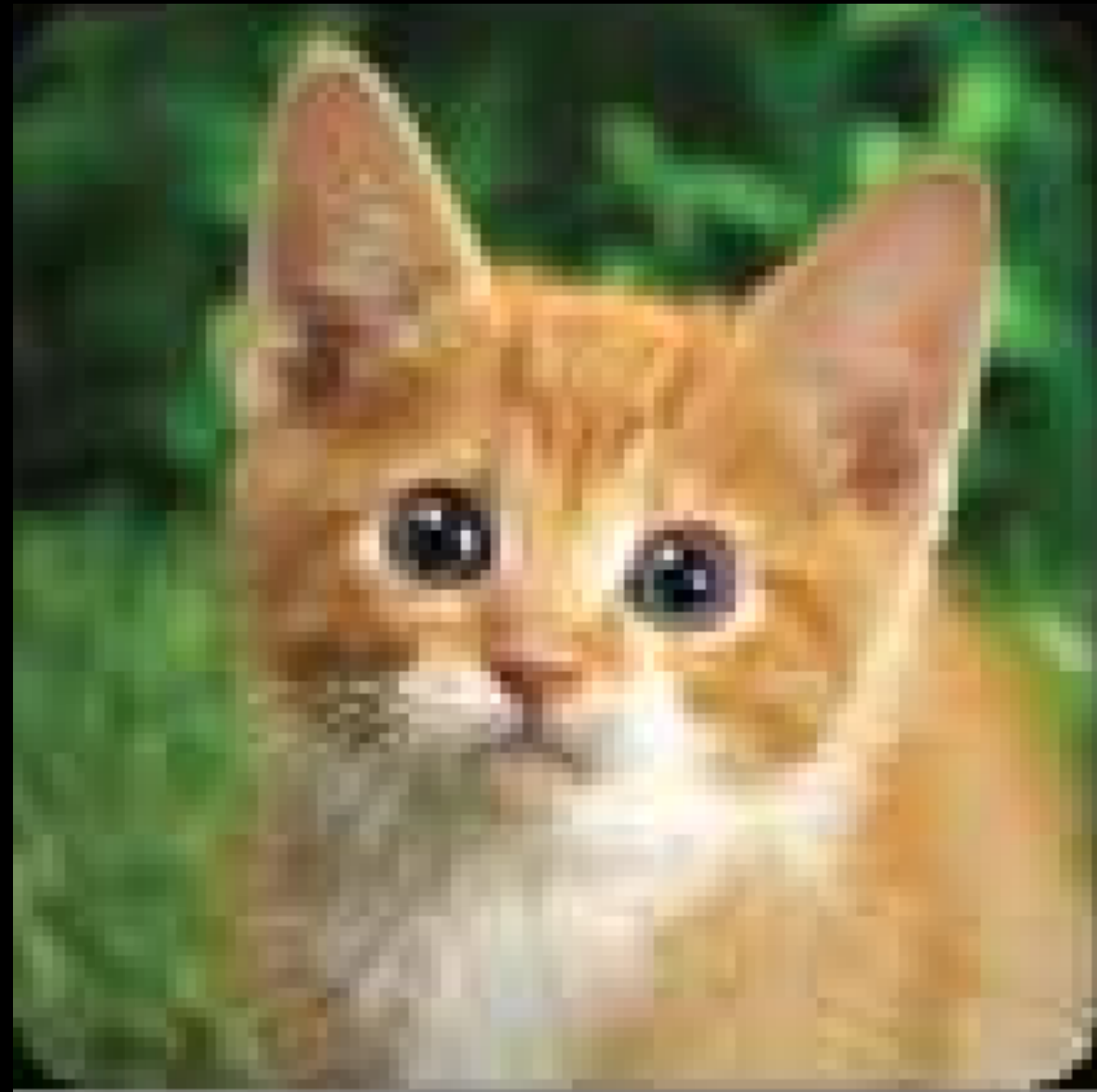
Google Research

David Wagner

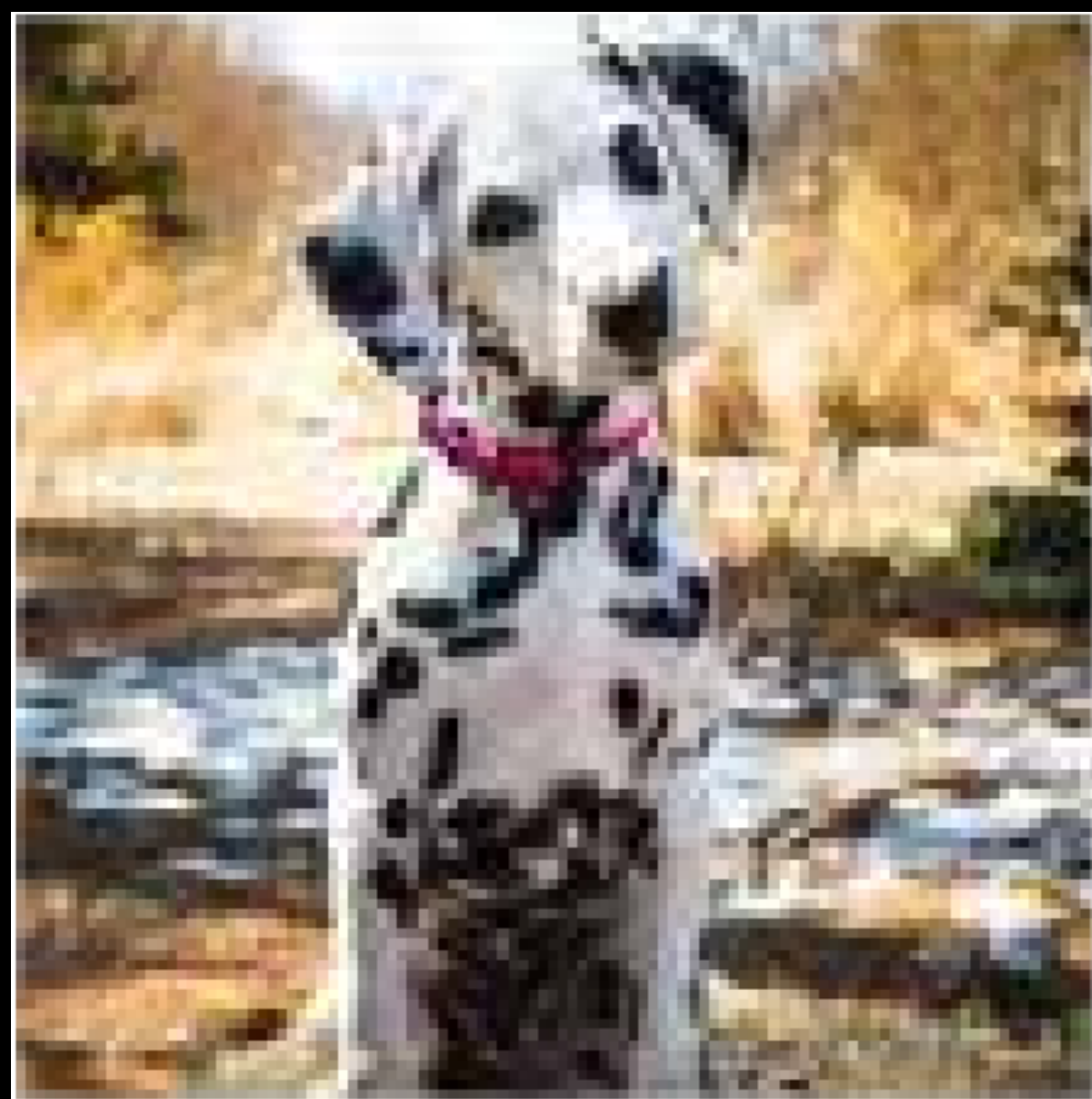
University of California, Berkeley



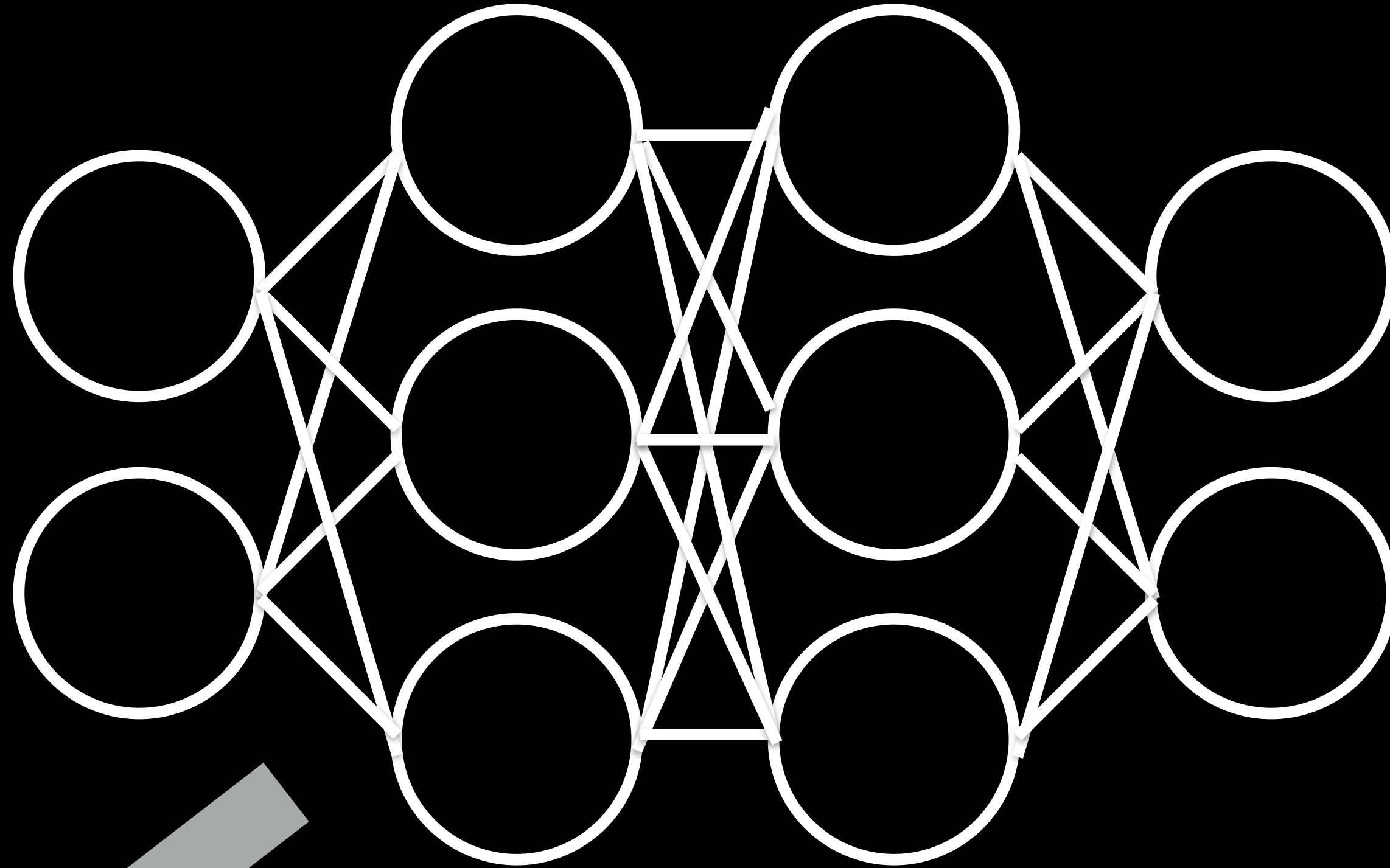
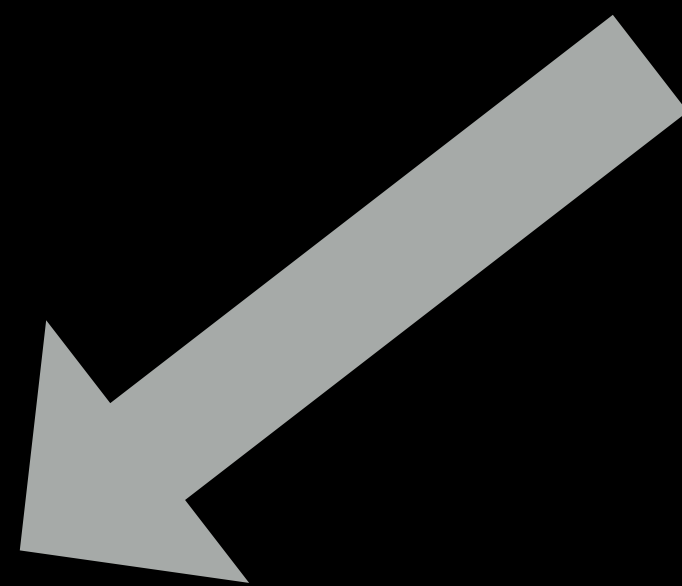
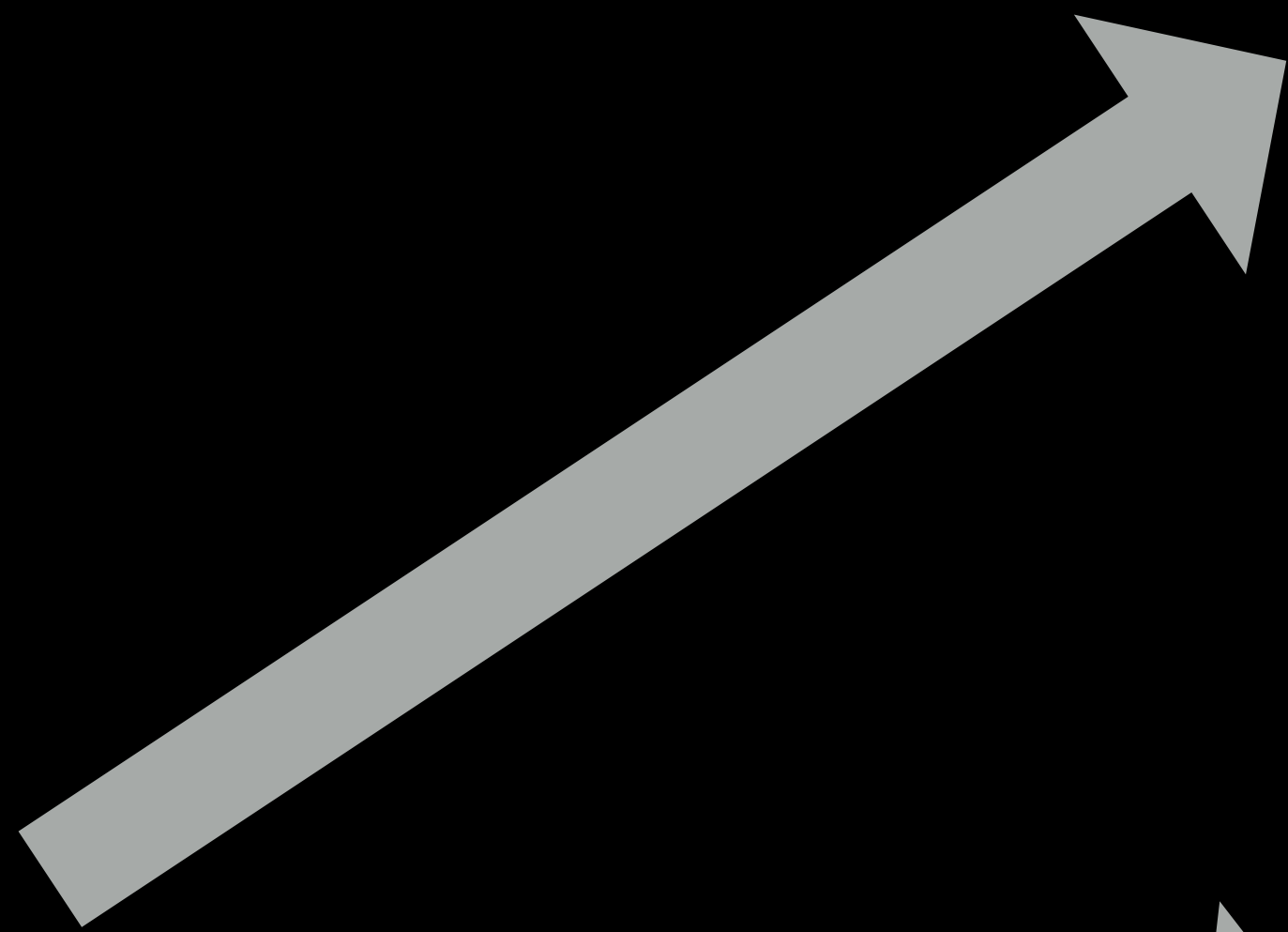
CAT



CAT

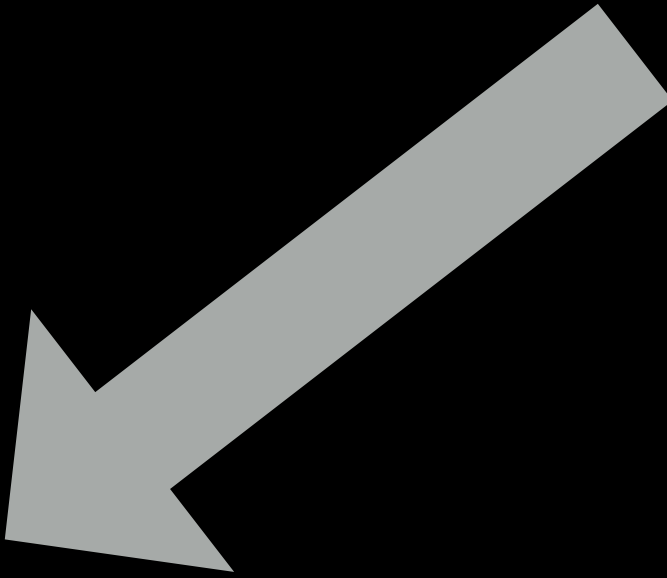
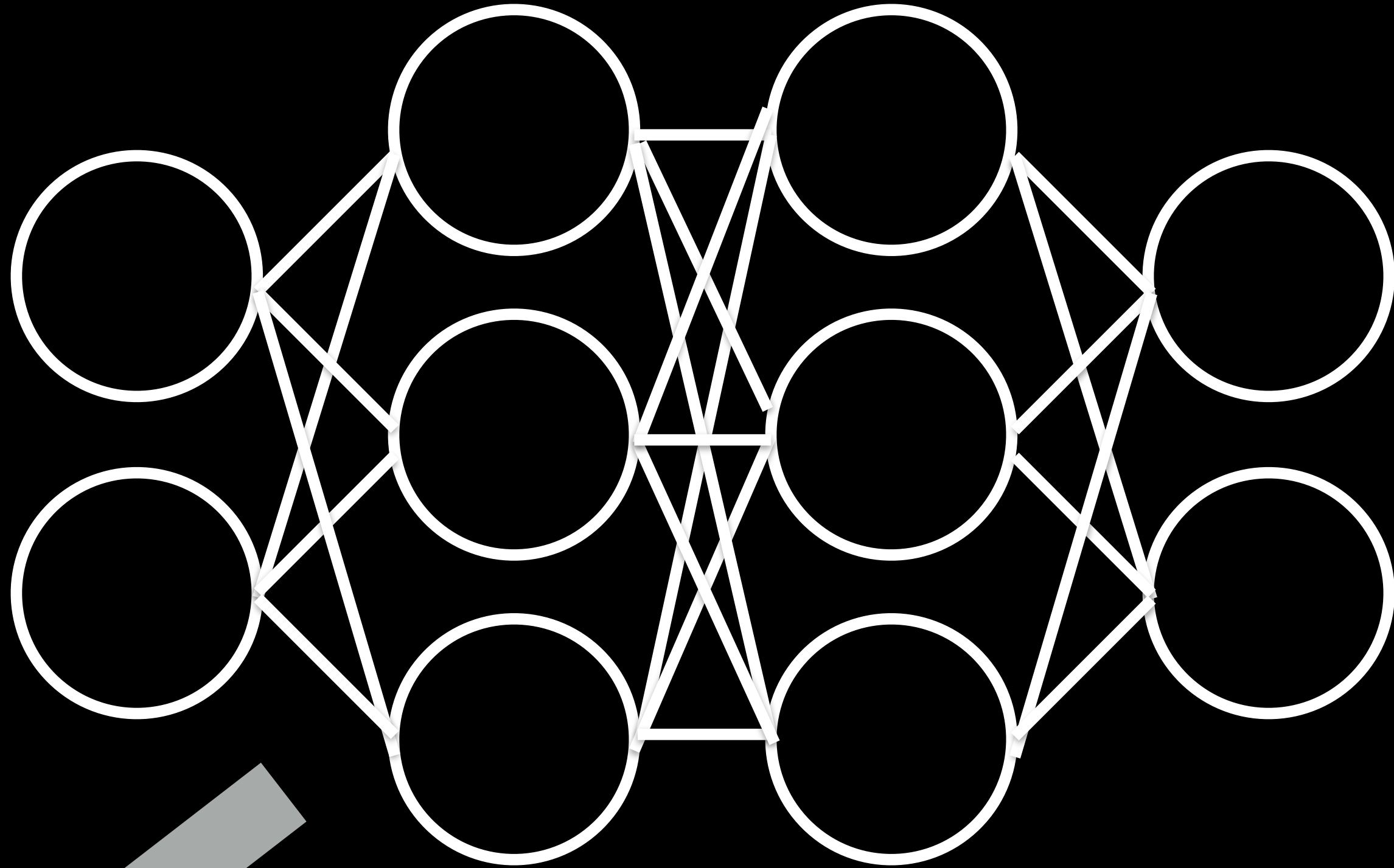
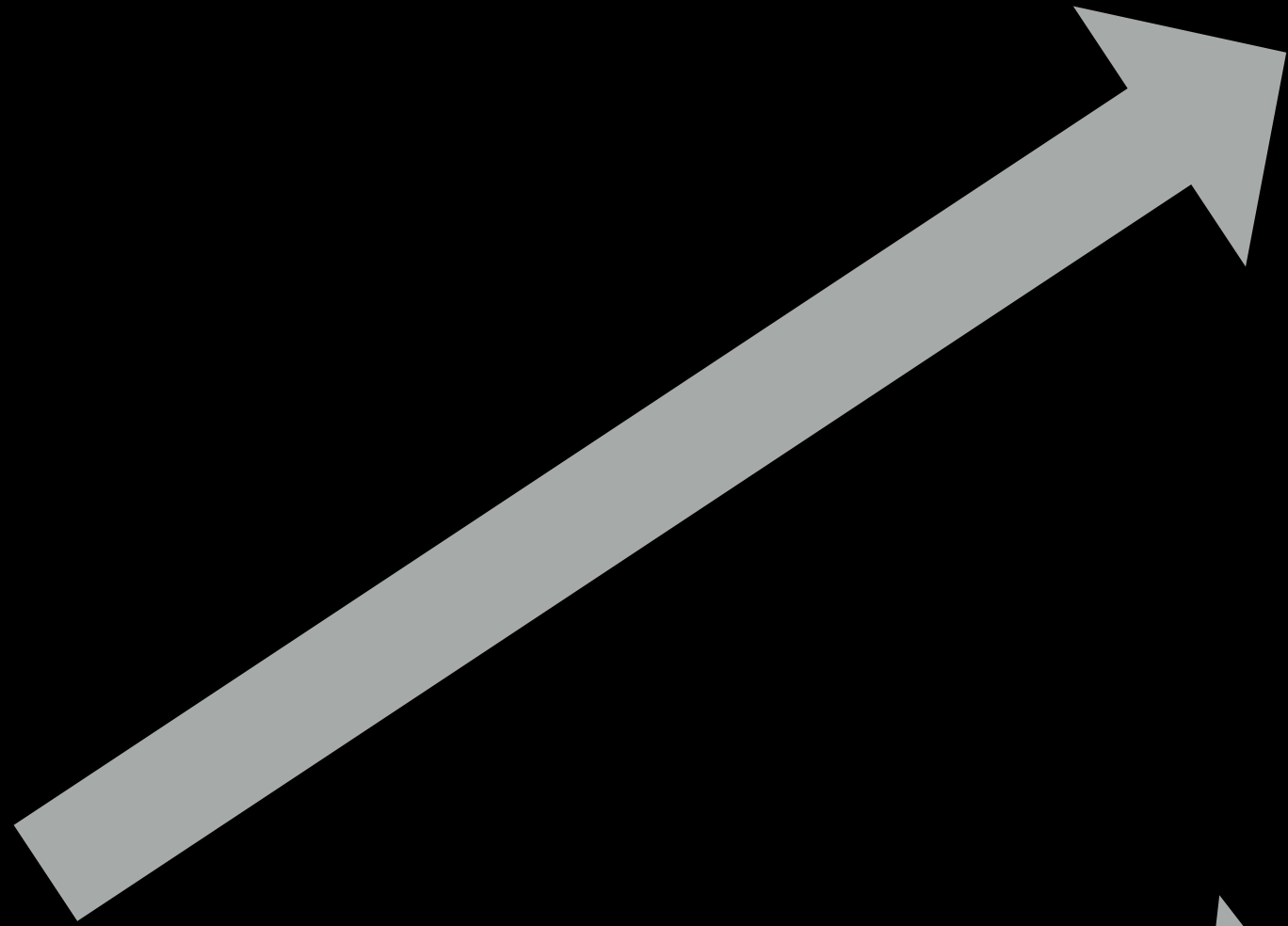


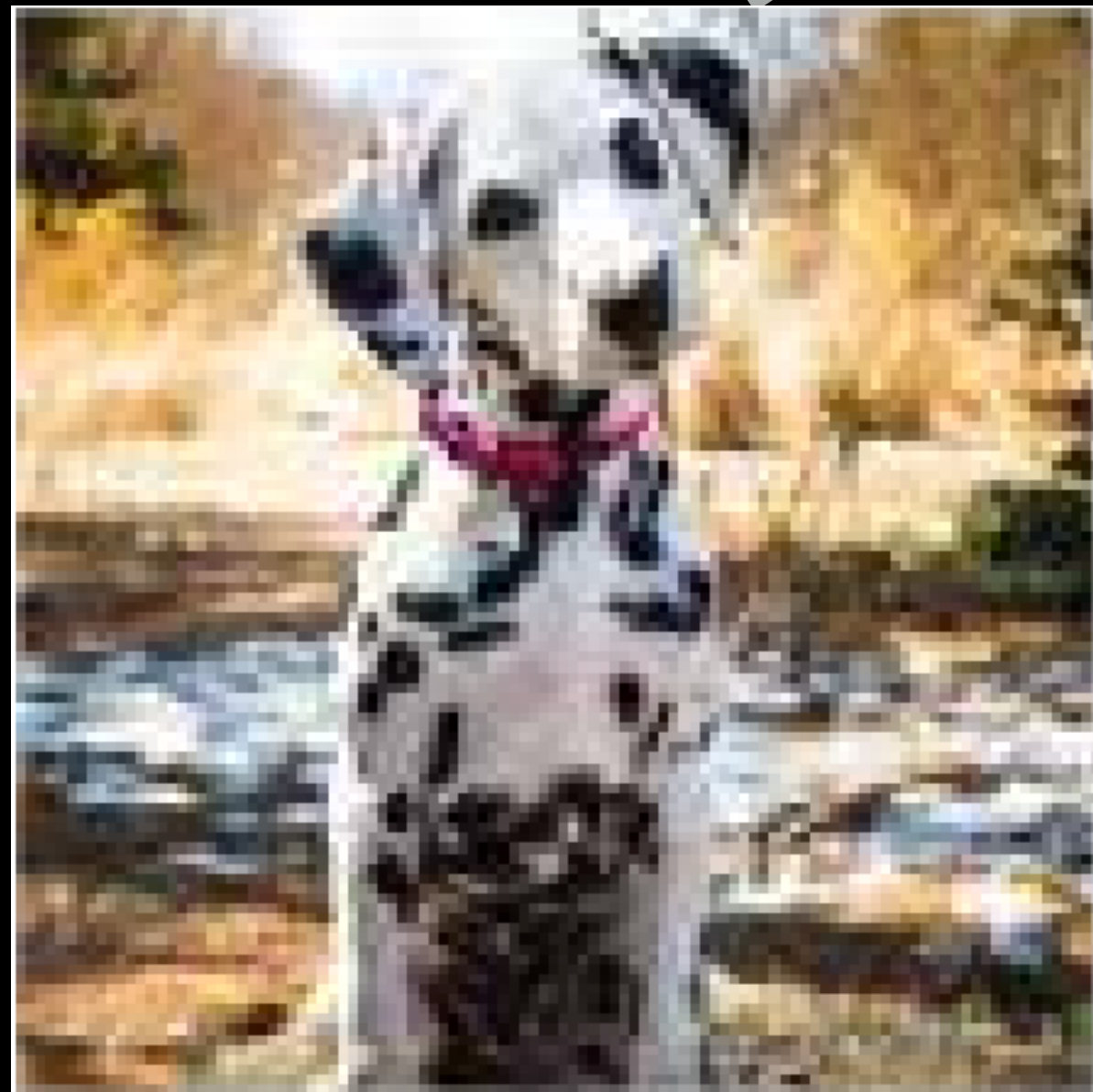
DOG



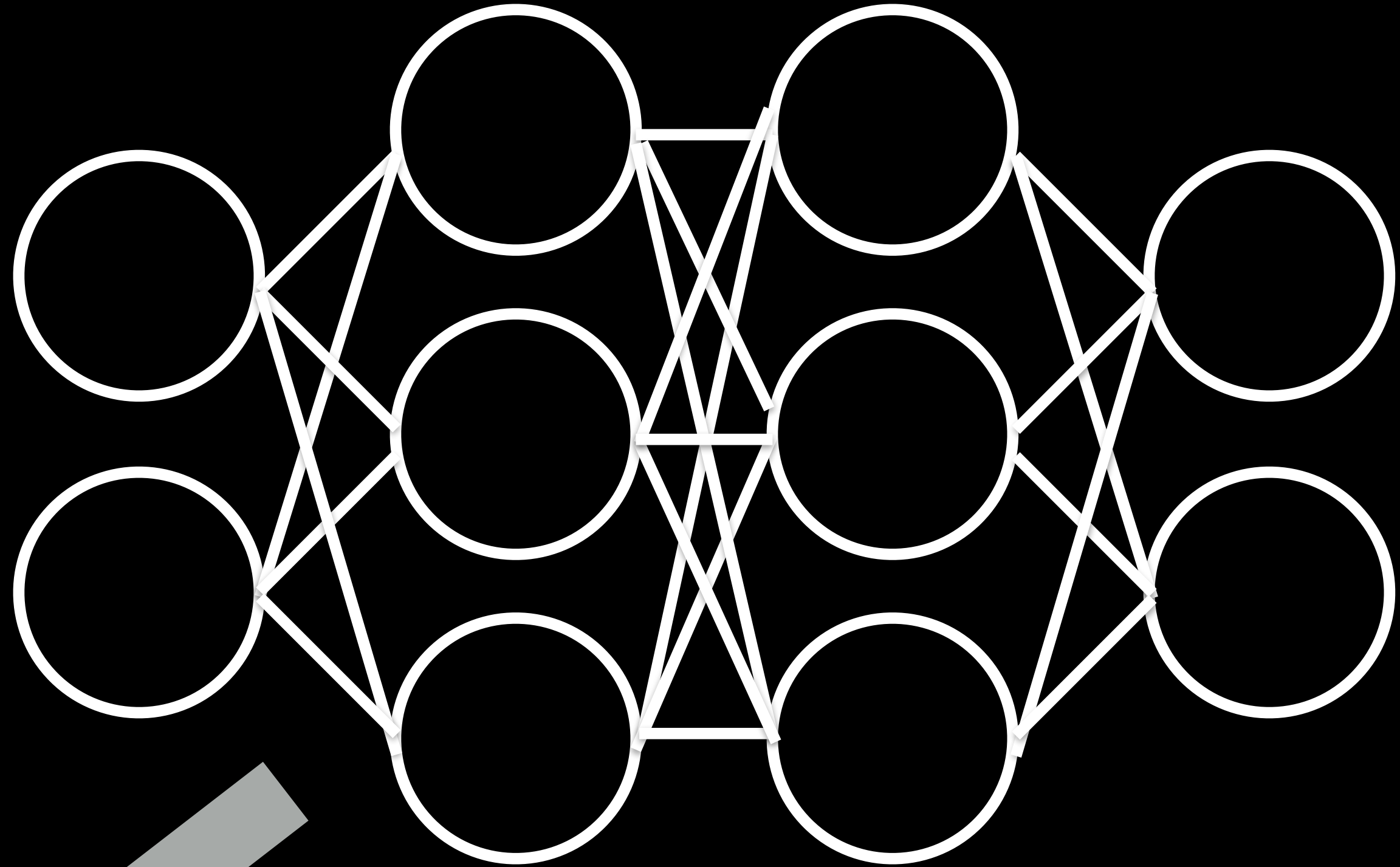
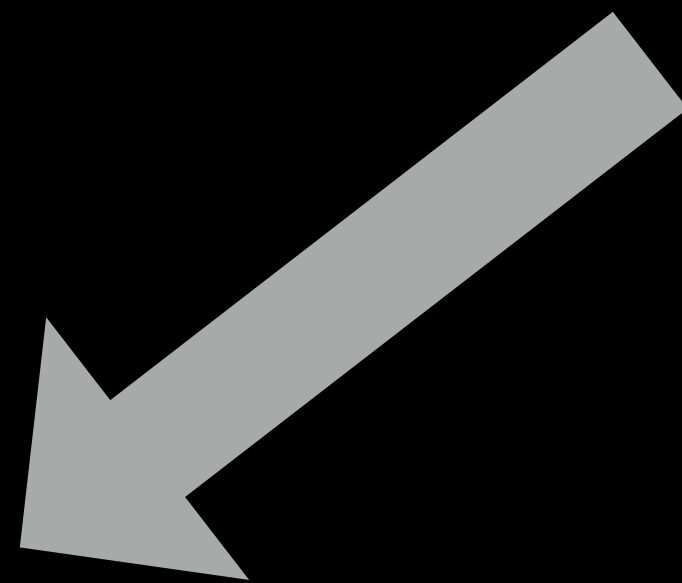
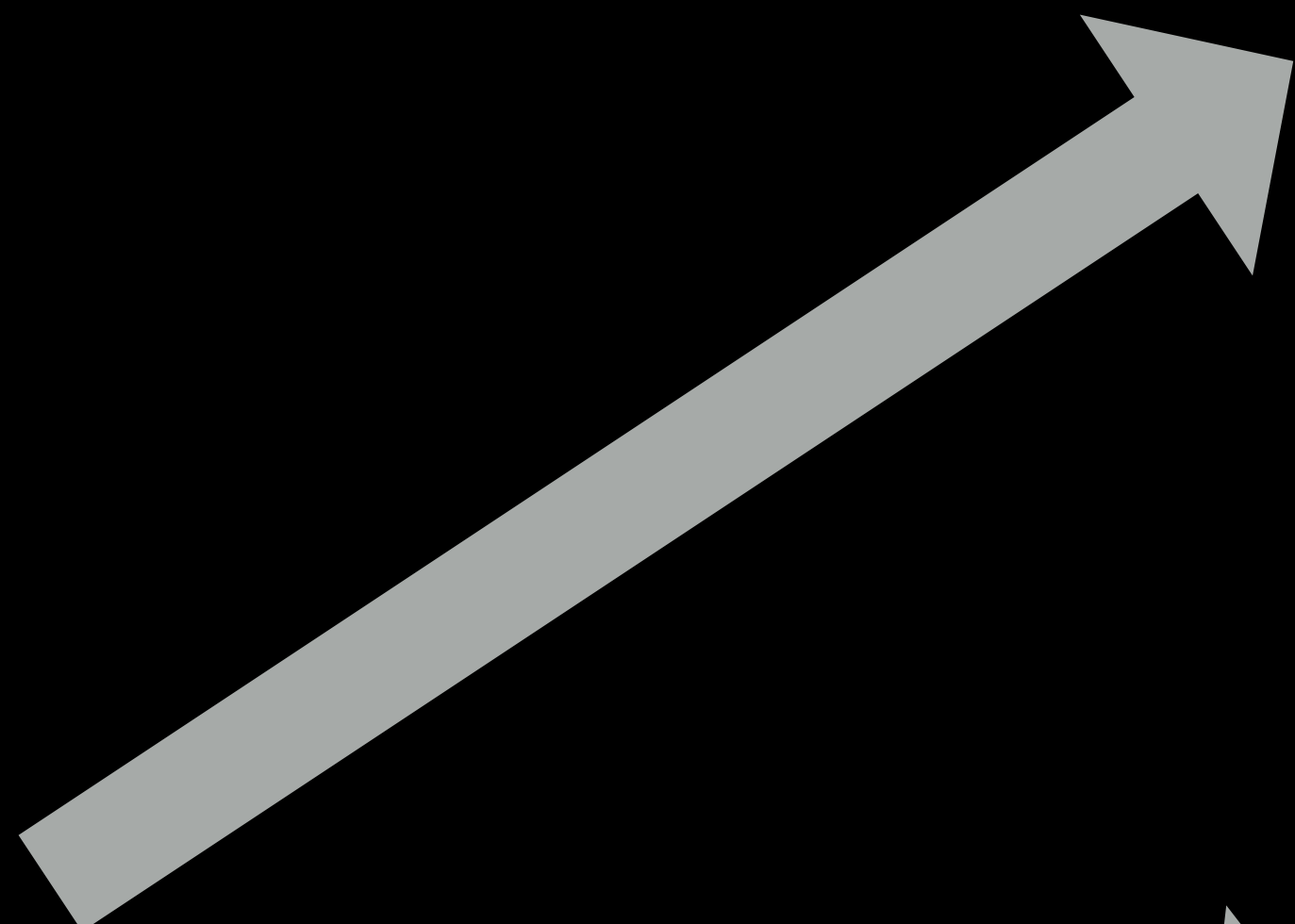


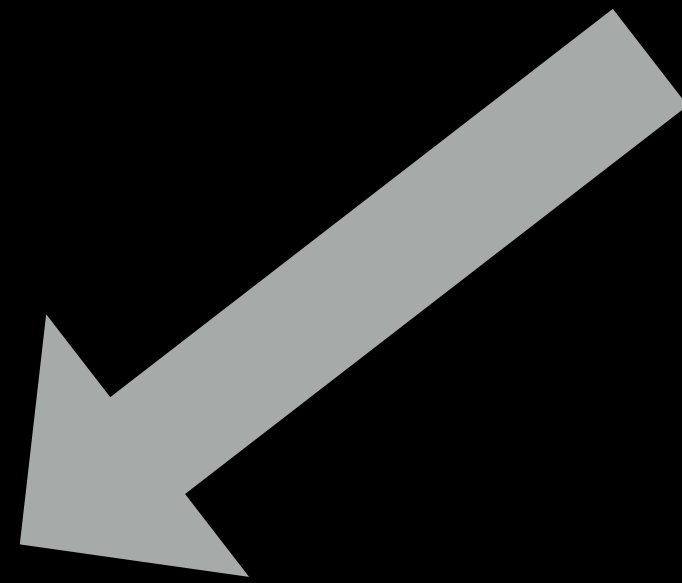
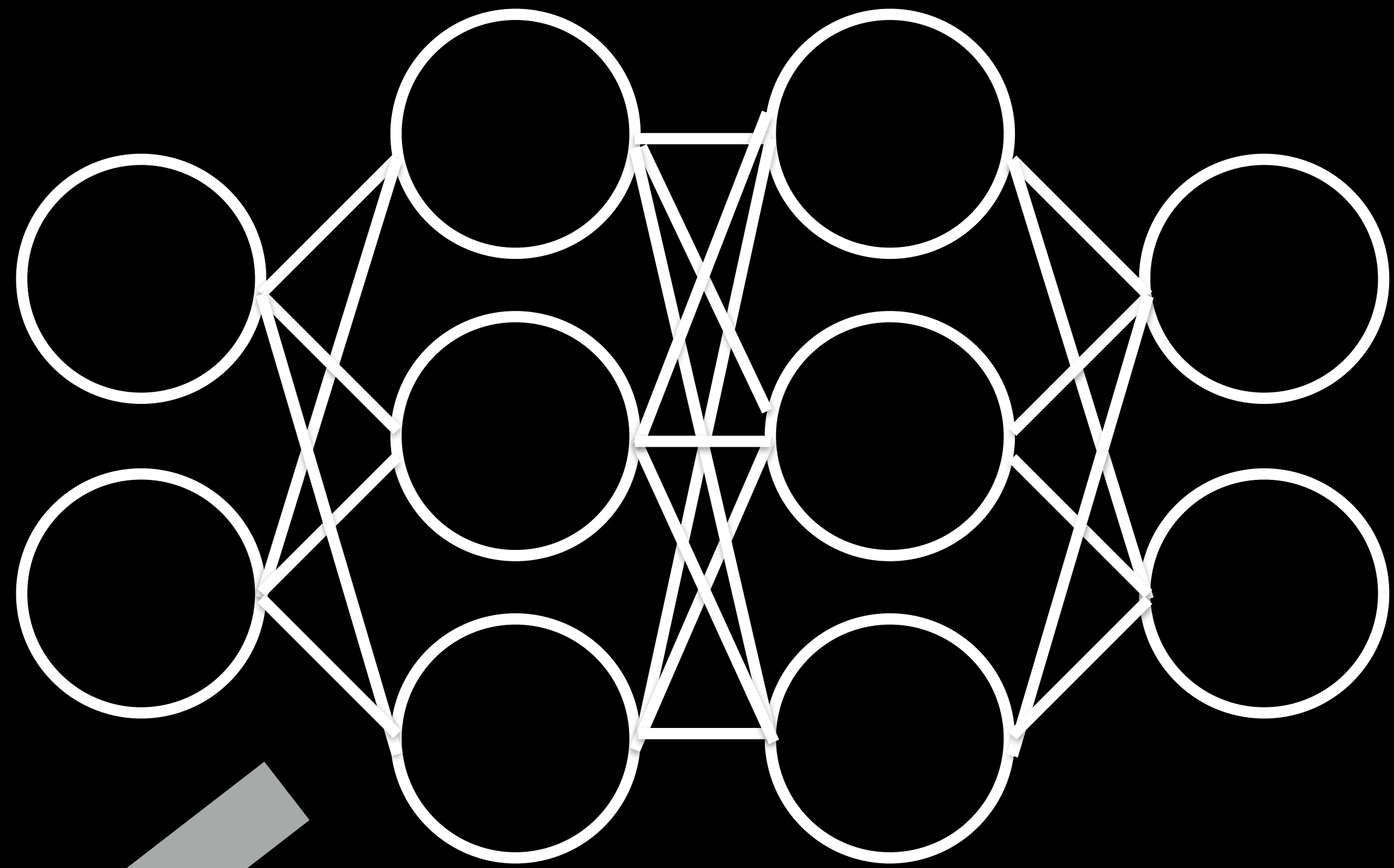
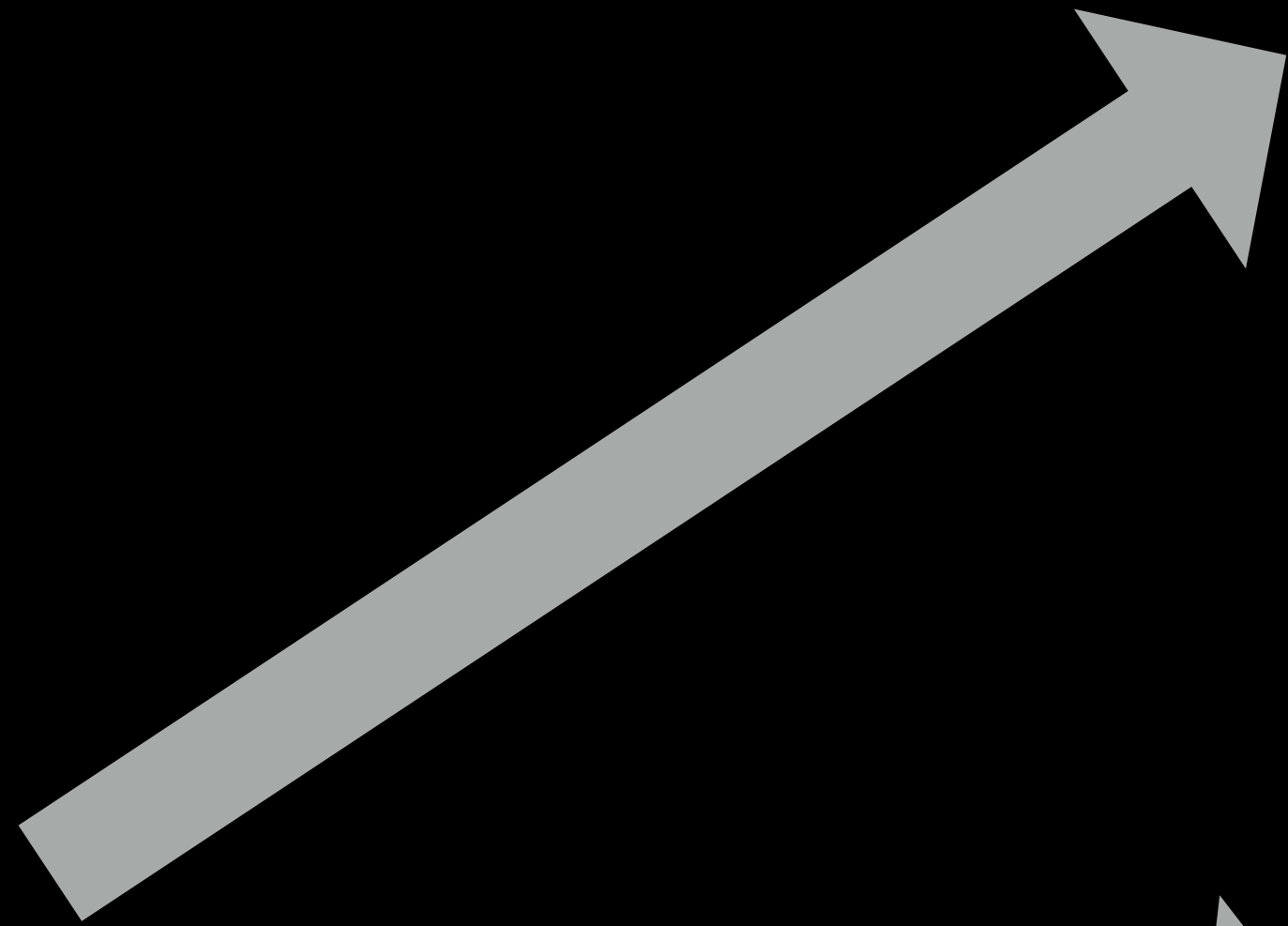
DOG



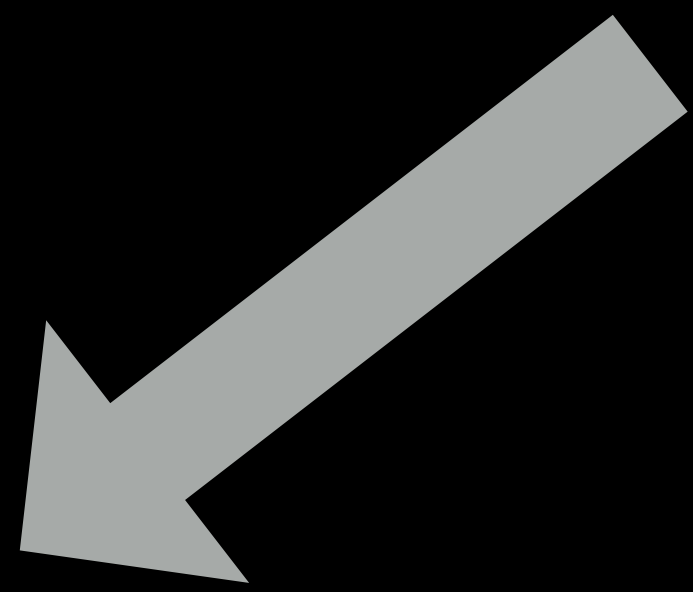
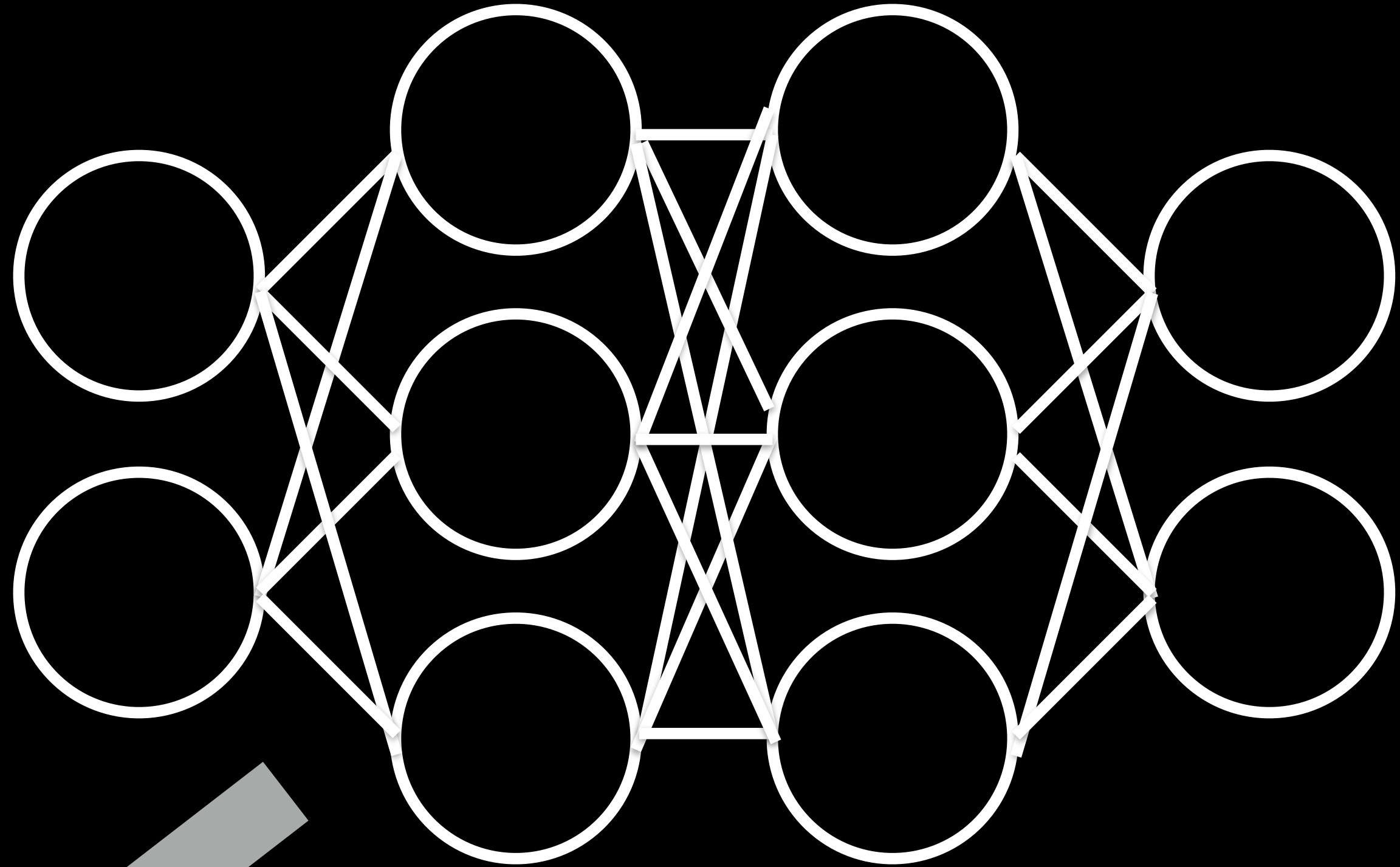
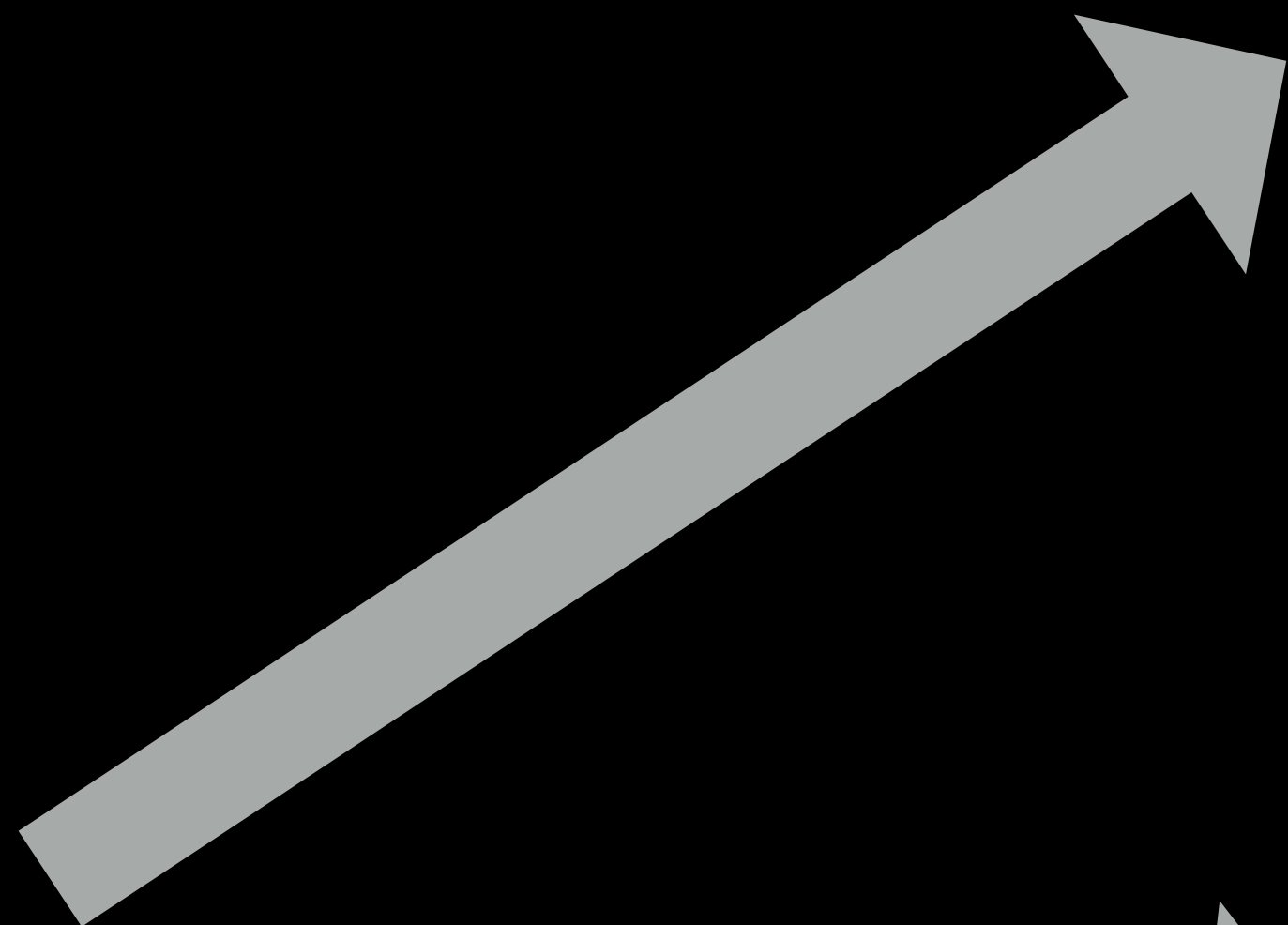
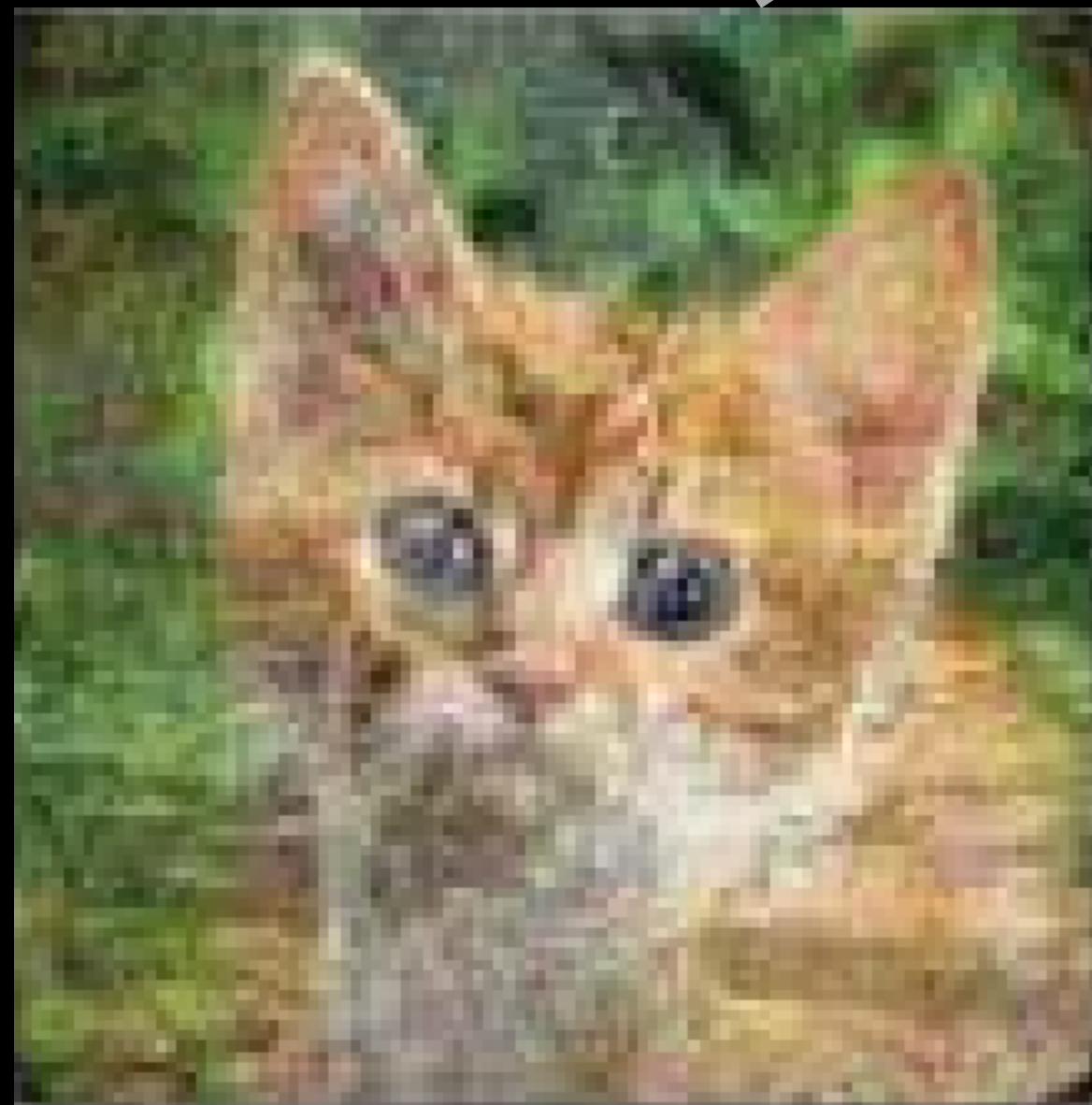


DOG





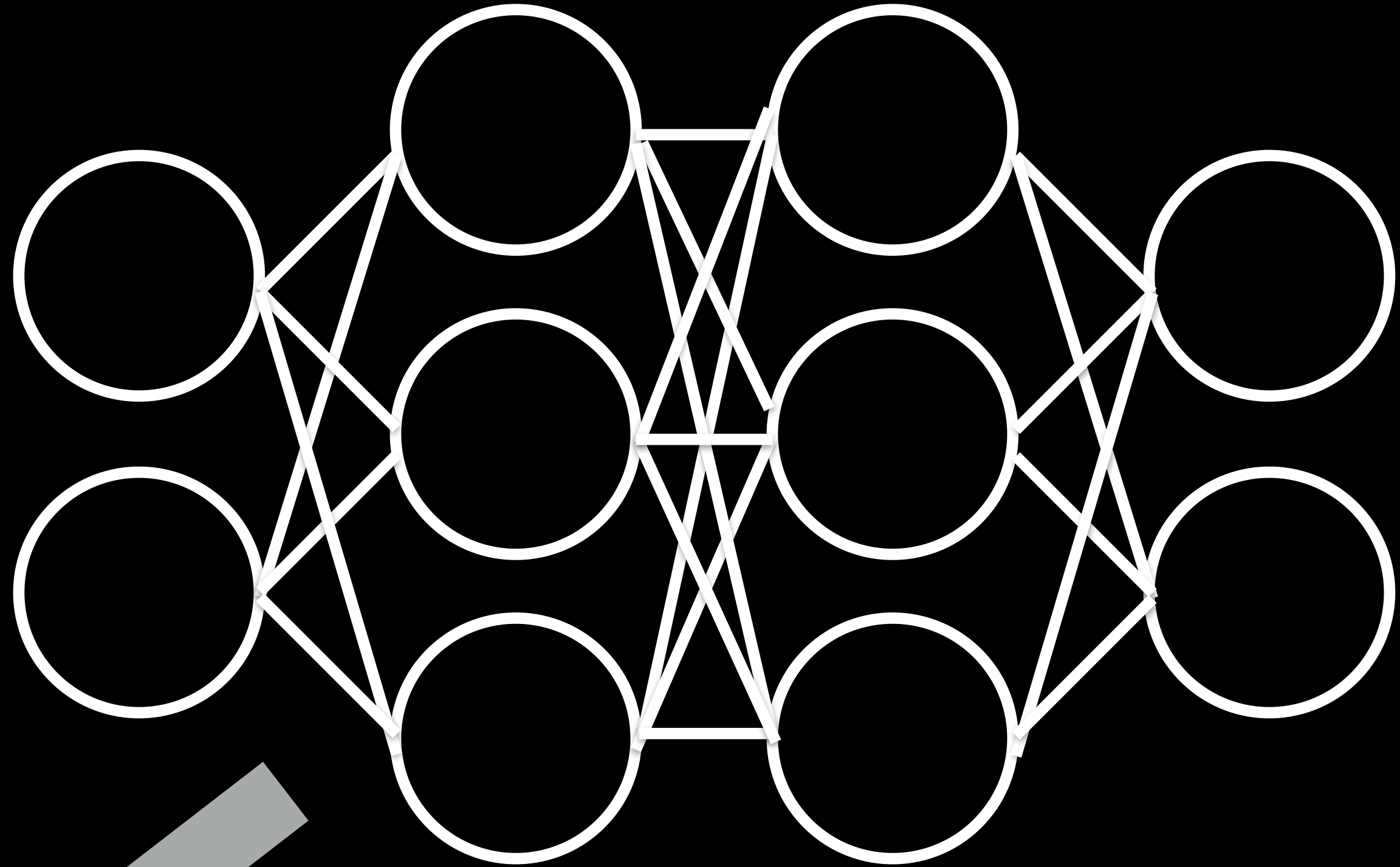
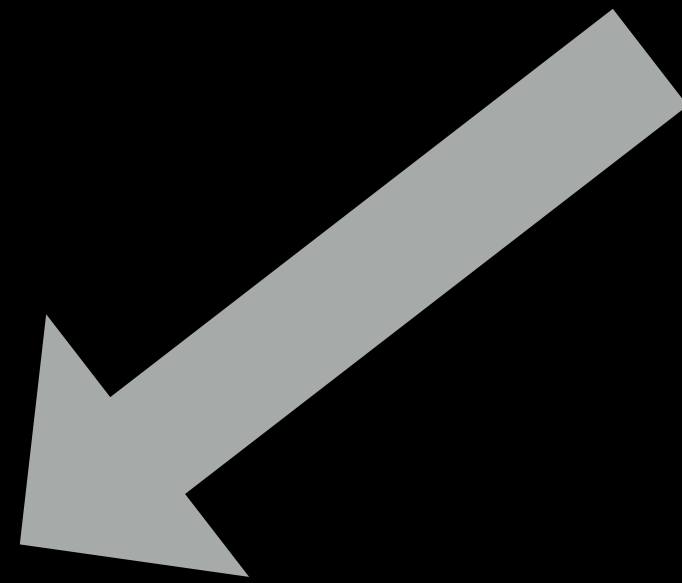
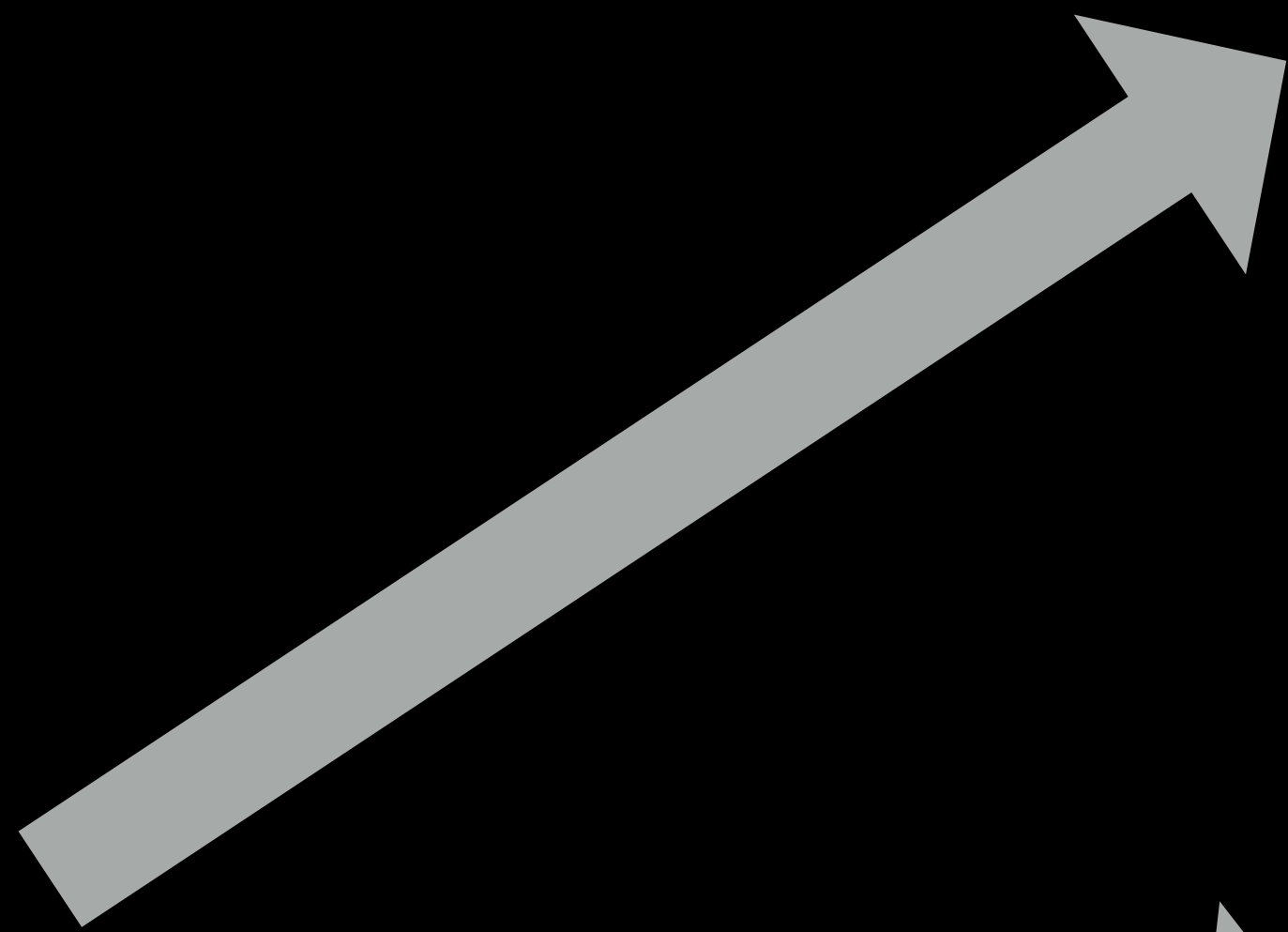
DOG

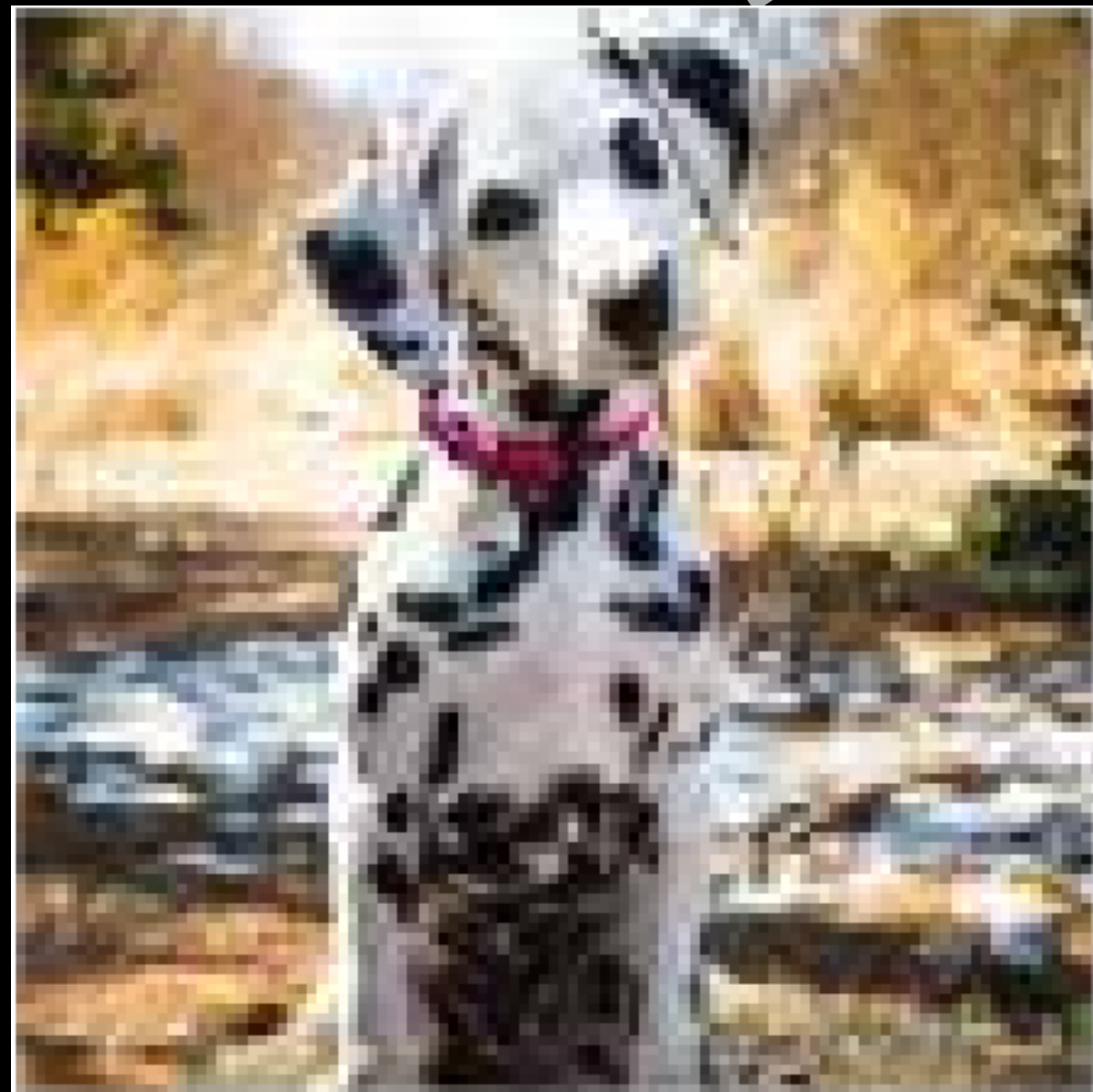


DOG

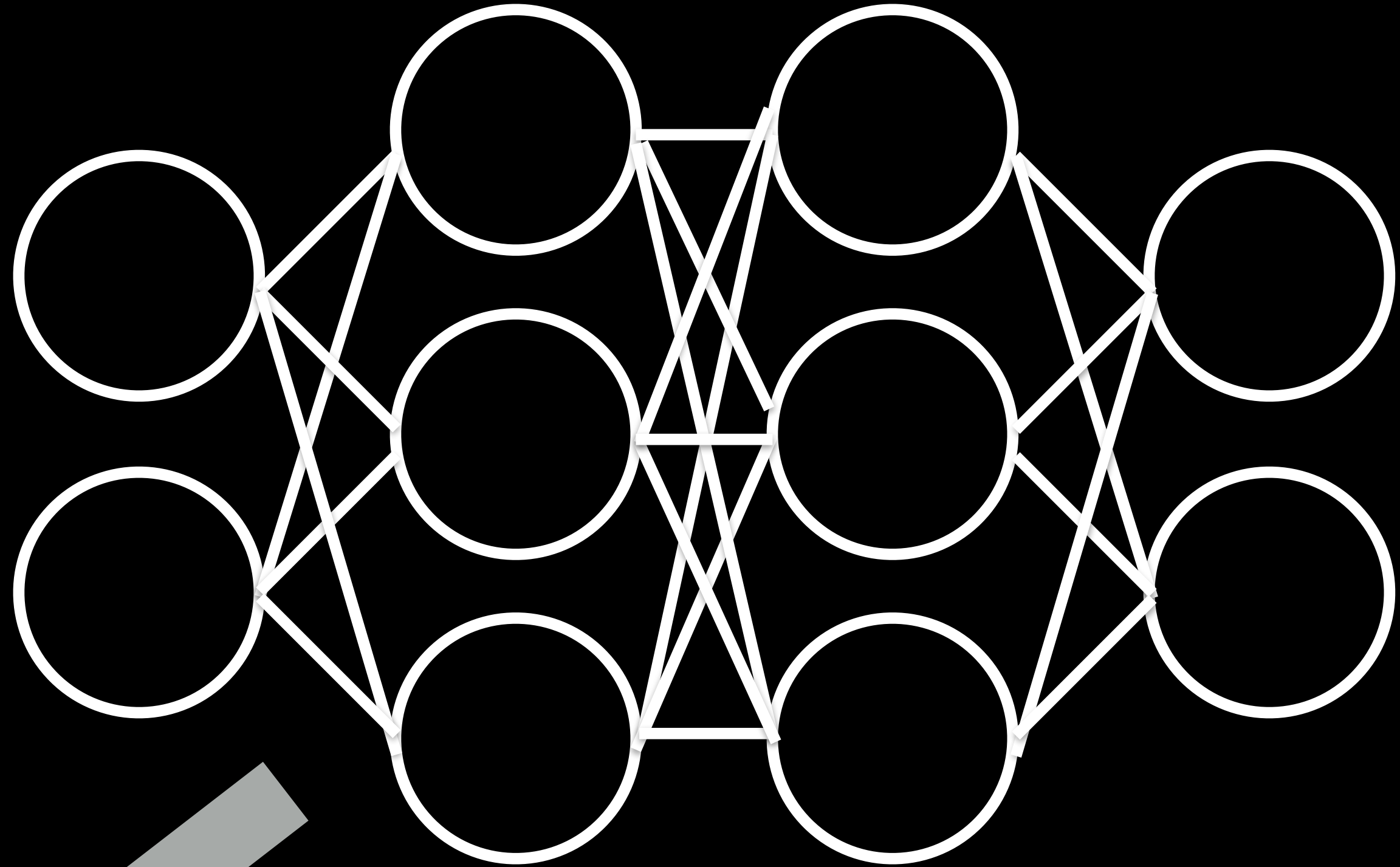
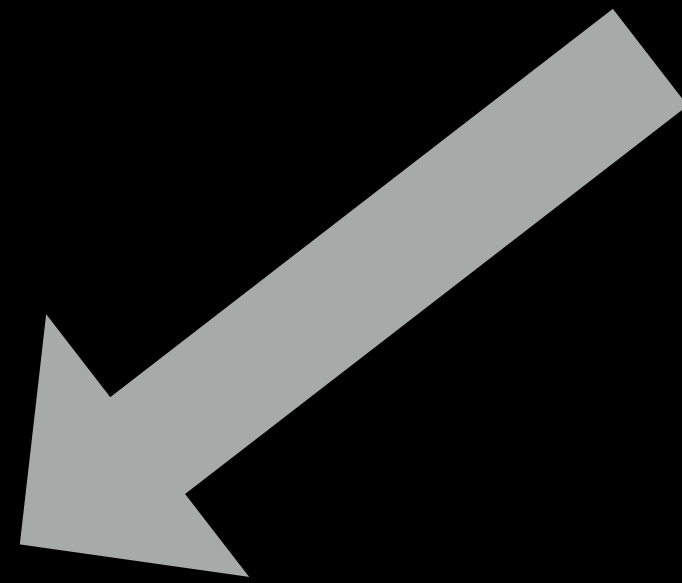
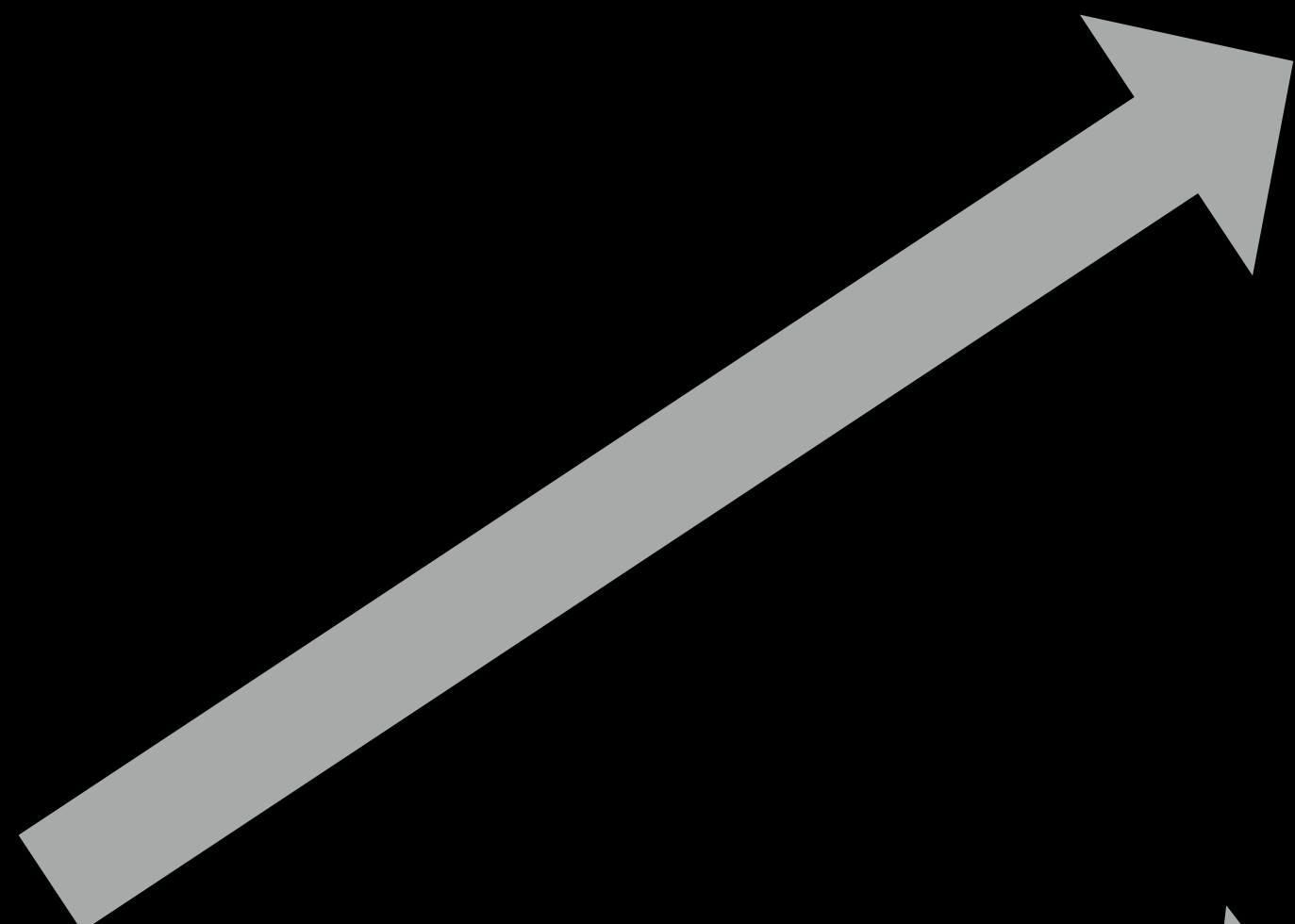


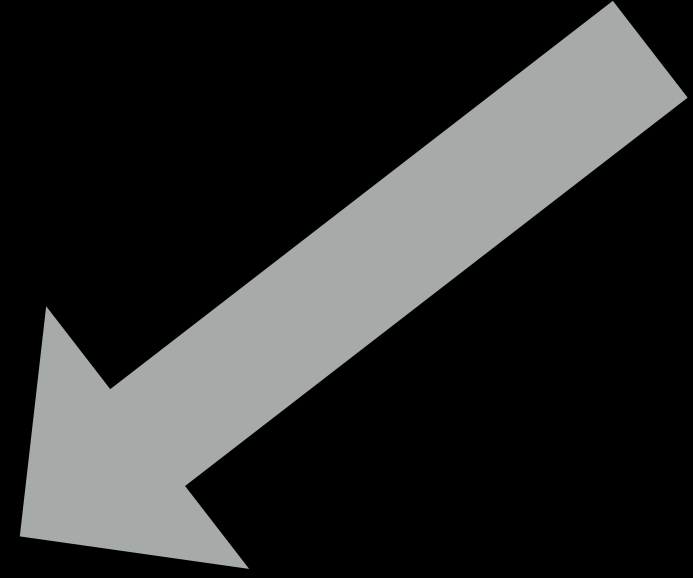
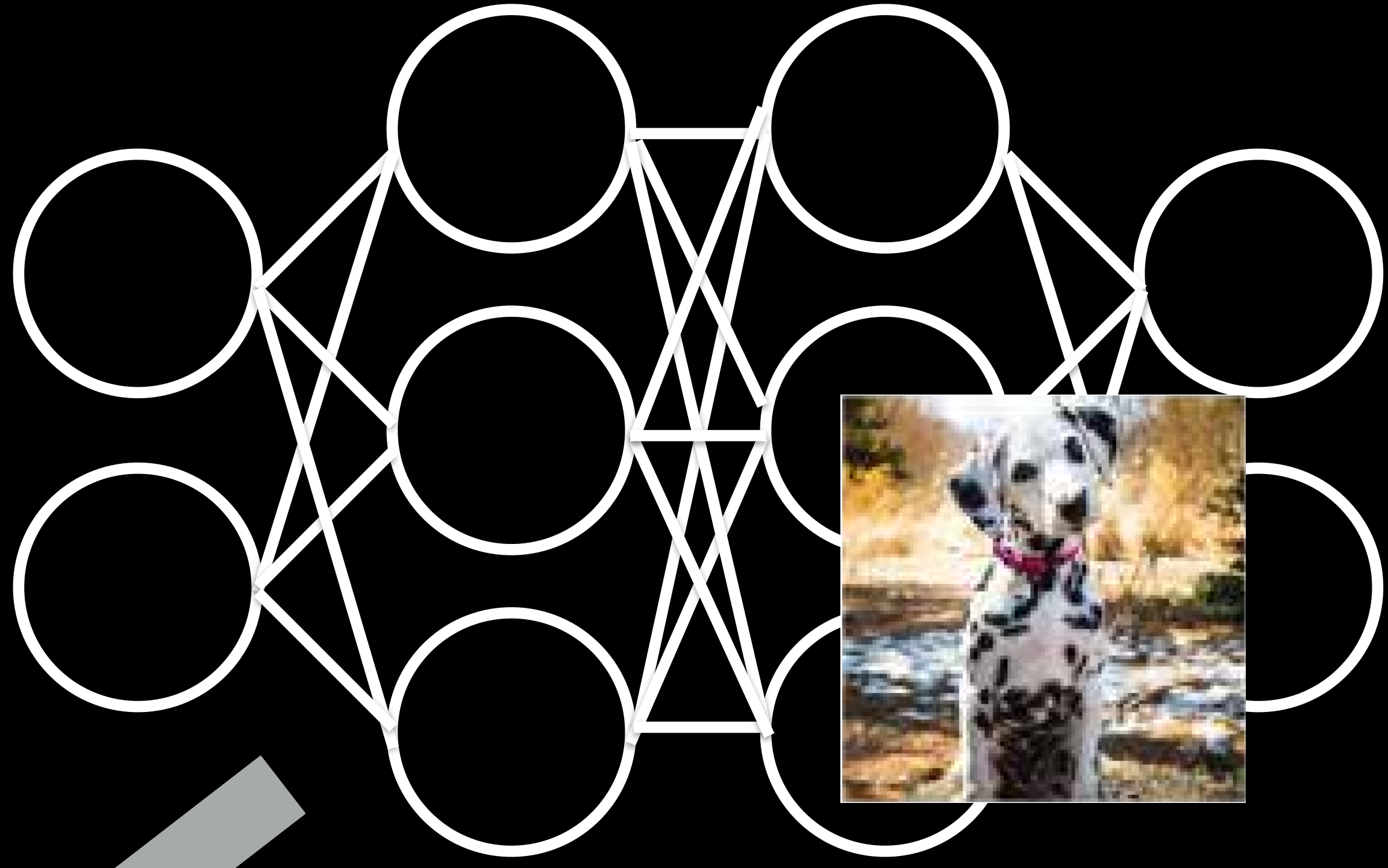
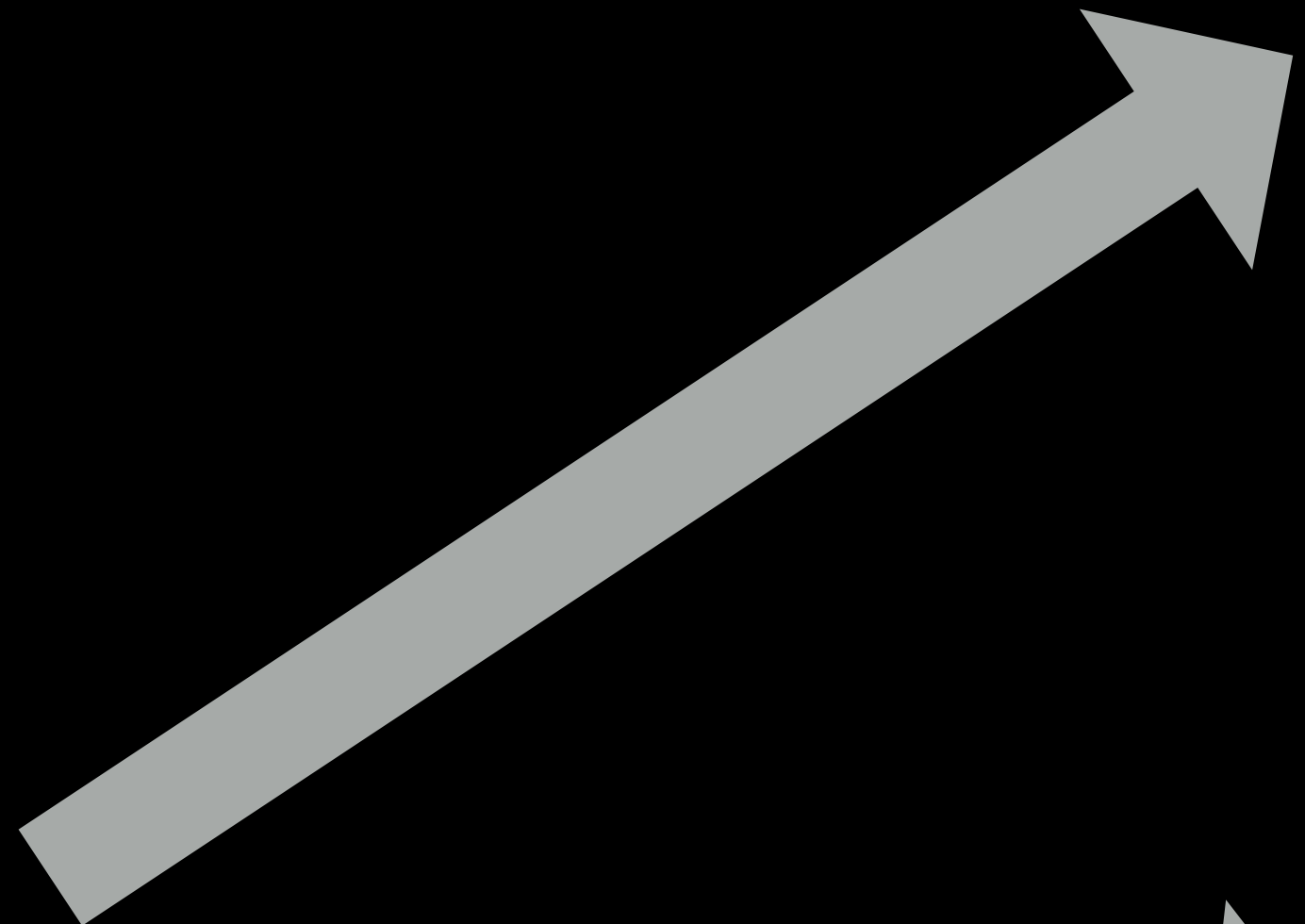
DOG



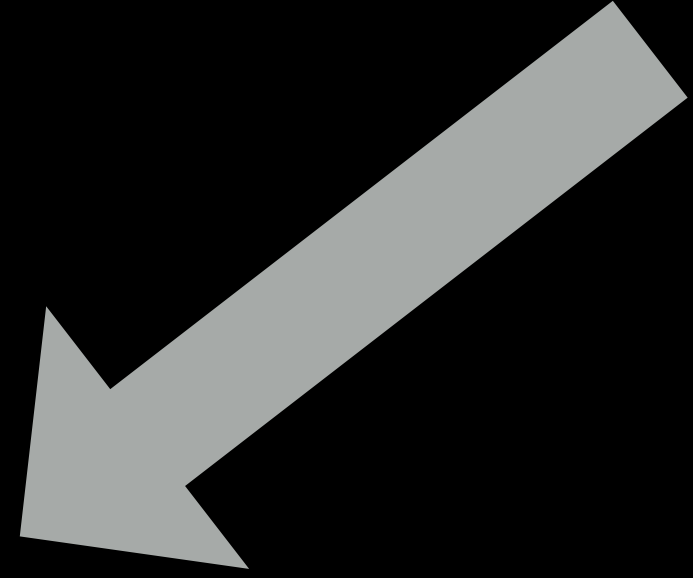
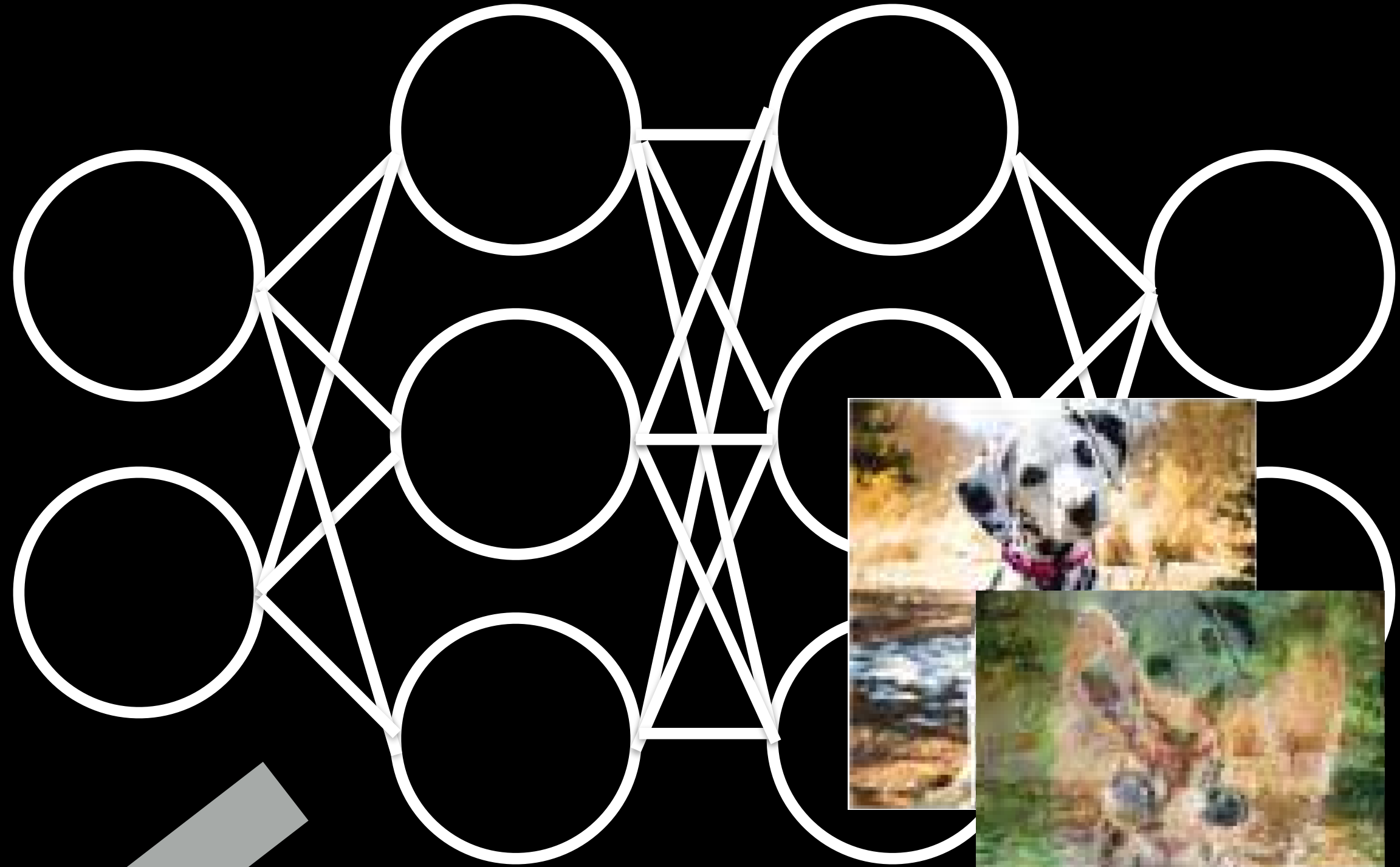
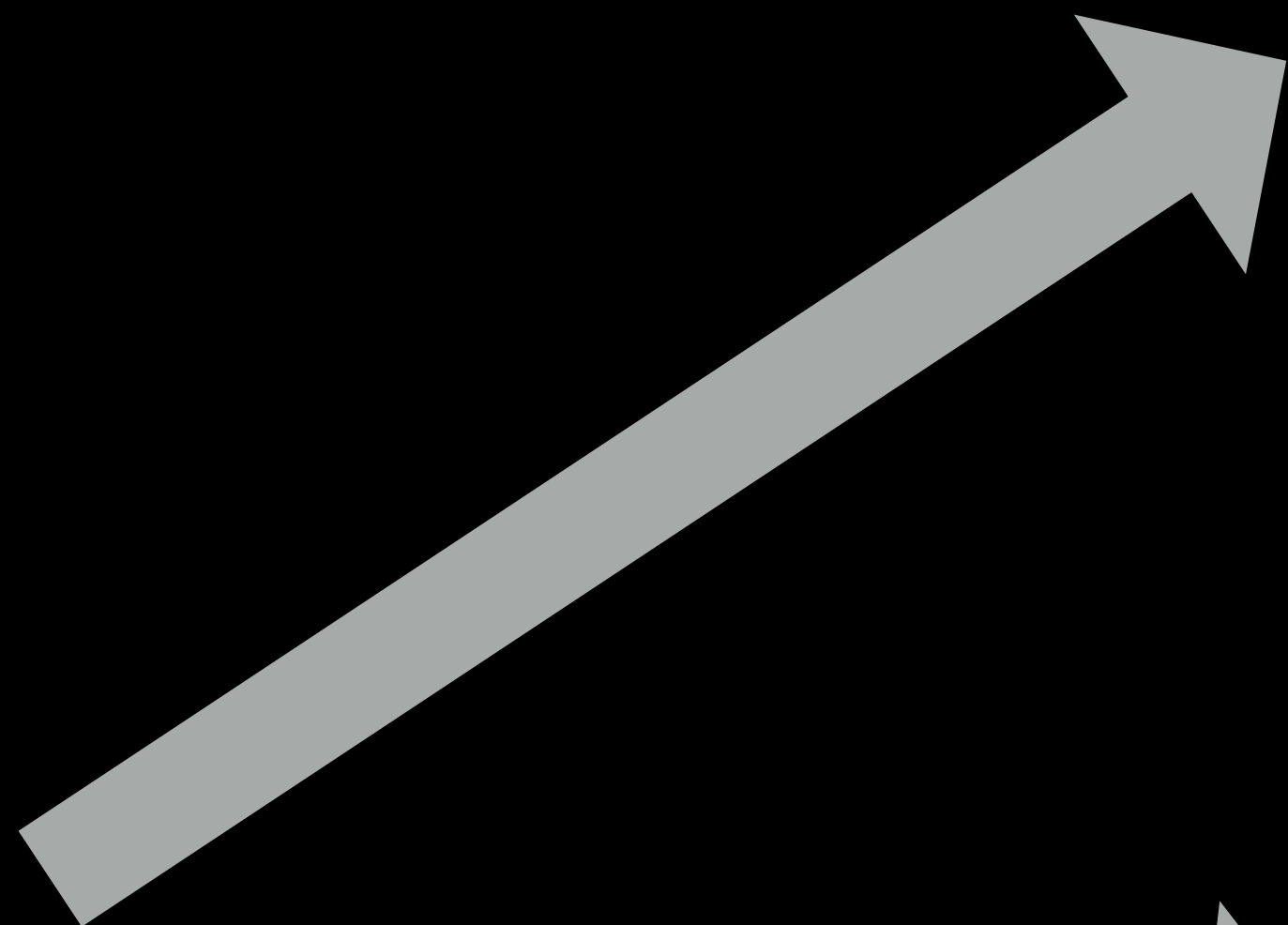
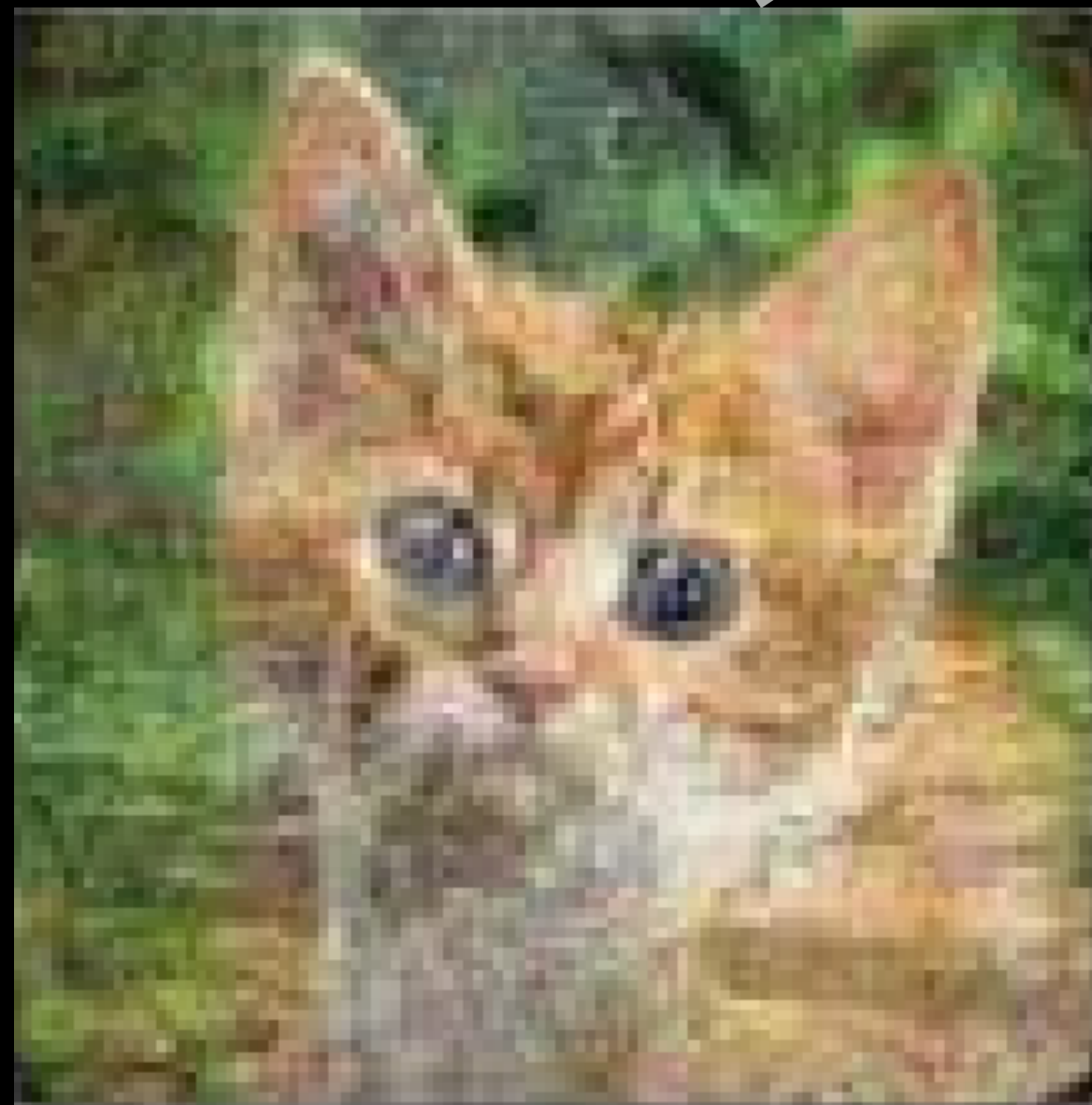


DOG

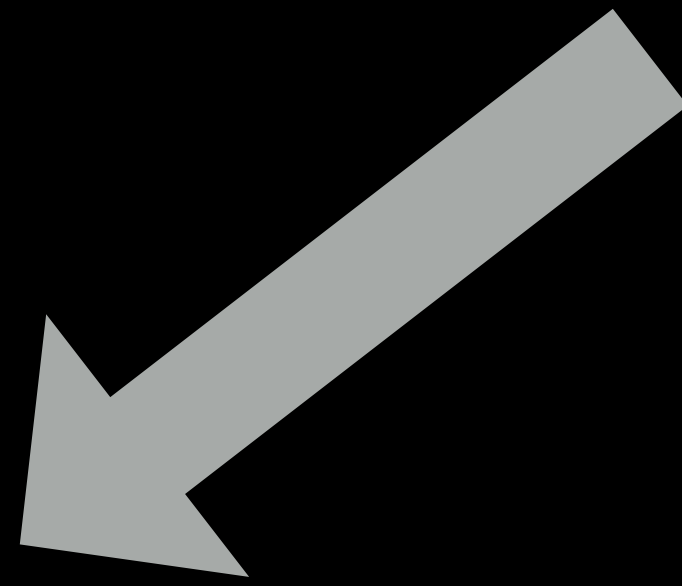
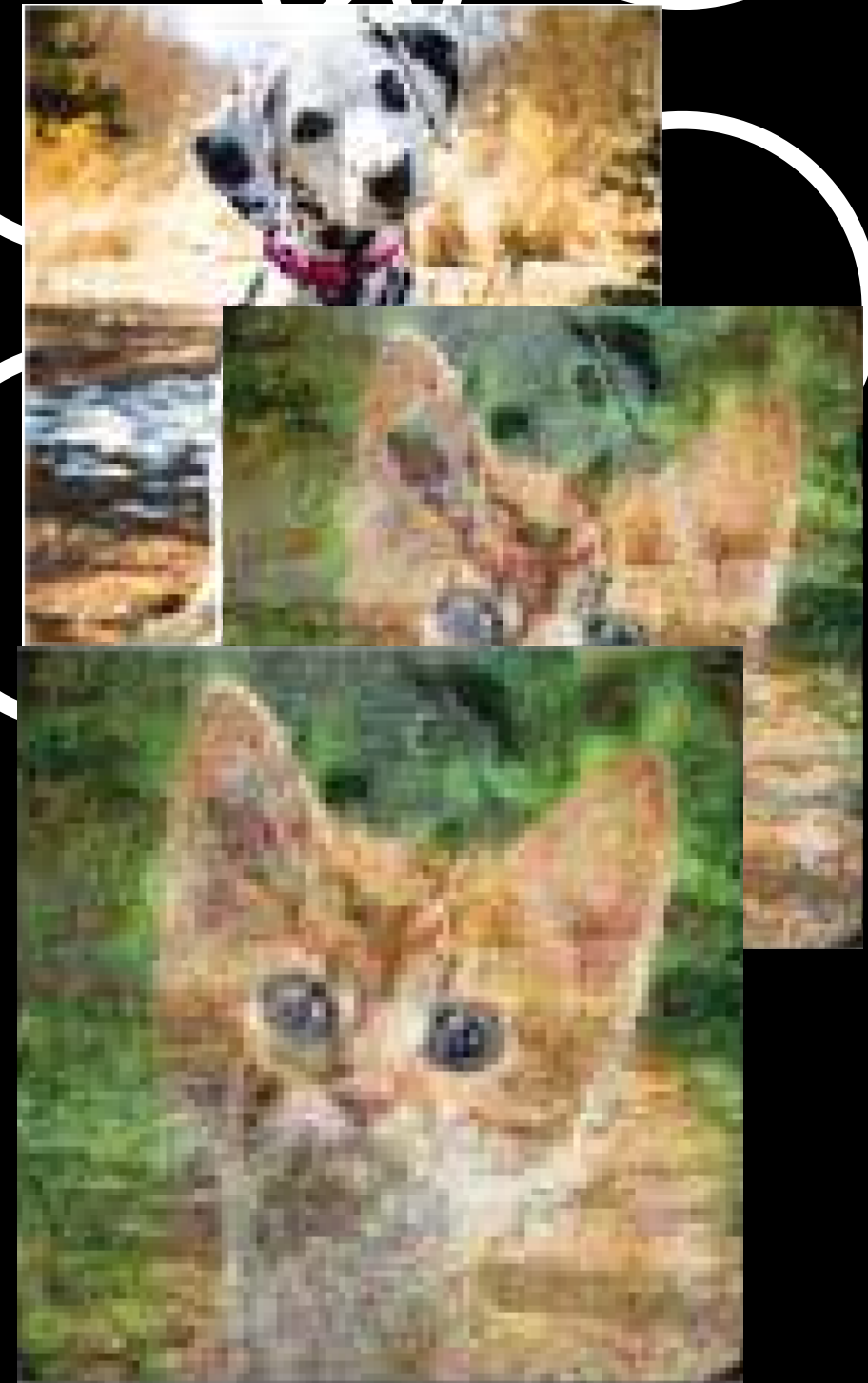
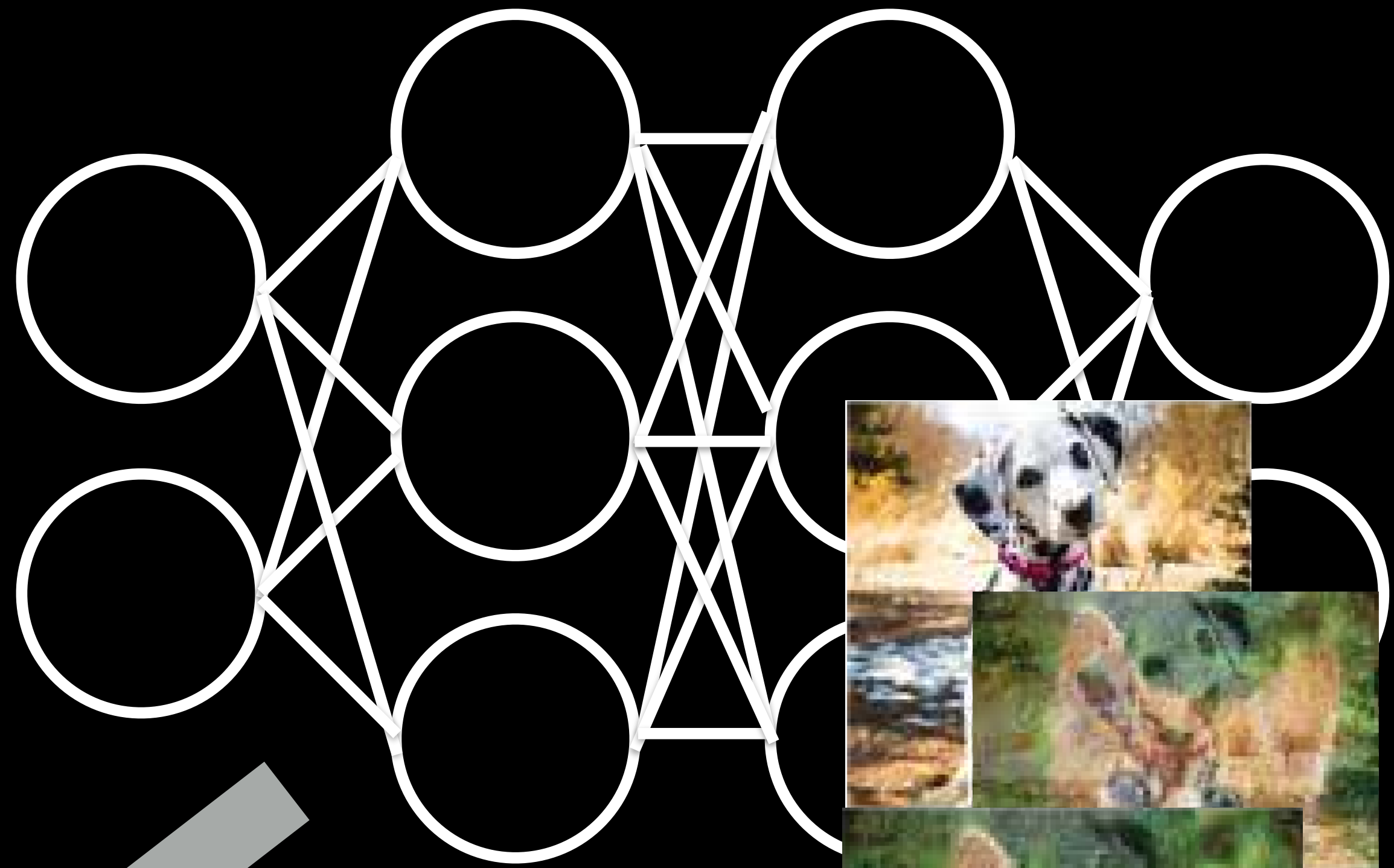
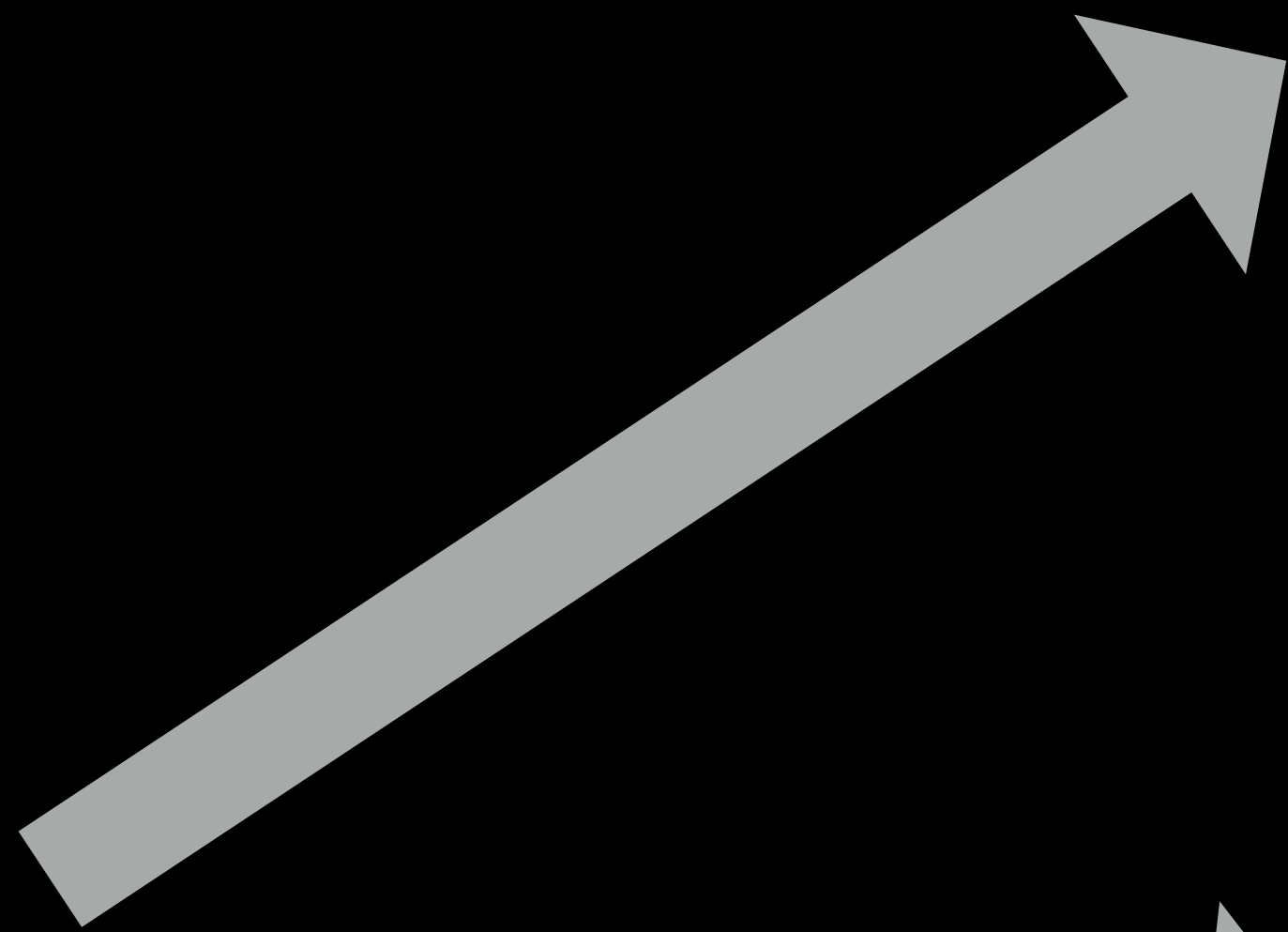




DOG



DOG



You are
being evil

Defenses I *do*
believe will be effective

Adversarially Robust Generalization Requires More Data

Ludwig Schmidt
MIT

Shibani Santurkar
MIT

Dimitris Tsipras
MIT

Kunal Talwar
Google Brain

Aleksander Mądry
MIT

Adversarially Robust Generalization Just Requires More Unlabeled Data

Unlabeled Data Improves Adversarial Robustness

Yair
Stanford
yairc@

Are Labels Required for Improving Adversarial Robustness?

Jonathan Uesato*

Jean-Baptiste Alayrac*

Po-Sen Huang*

Robert Stanforth

Alhussein Fawzi

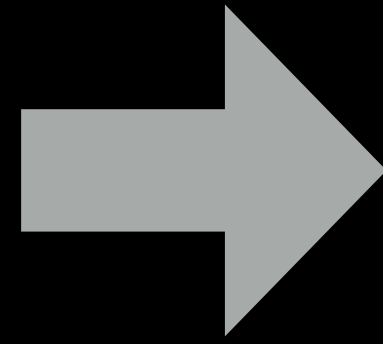
Pushmeet Kohli

Certified Robustness to Adversarial Examples with Differential Privacy

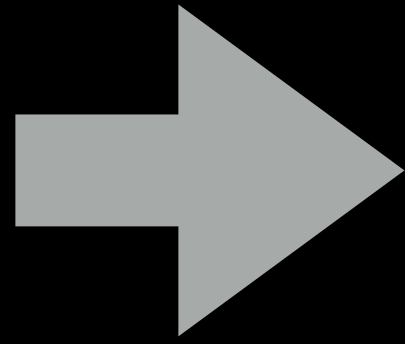
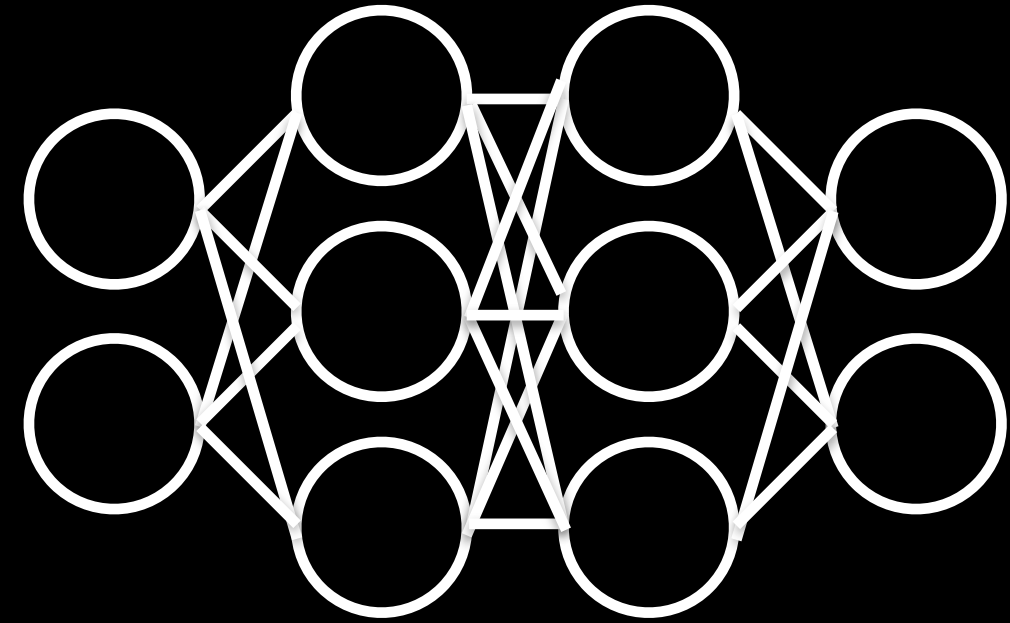
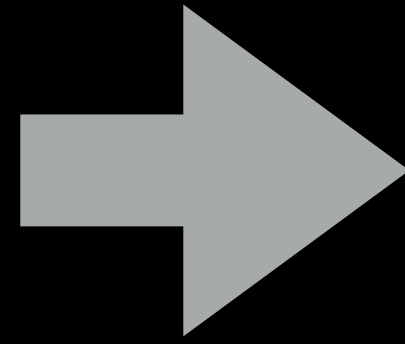
Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana
Columbia University

Certified Adversarial Robustness via Randomized Smoothing

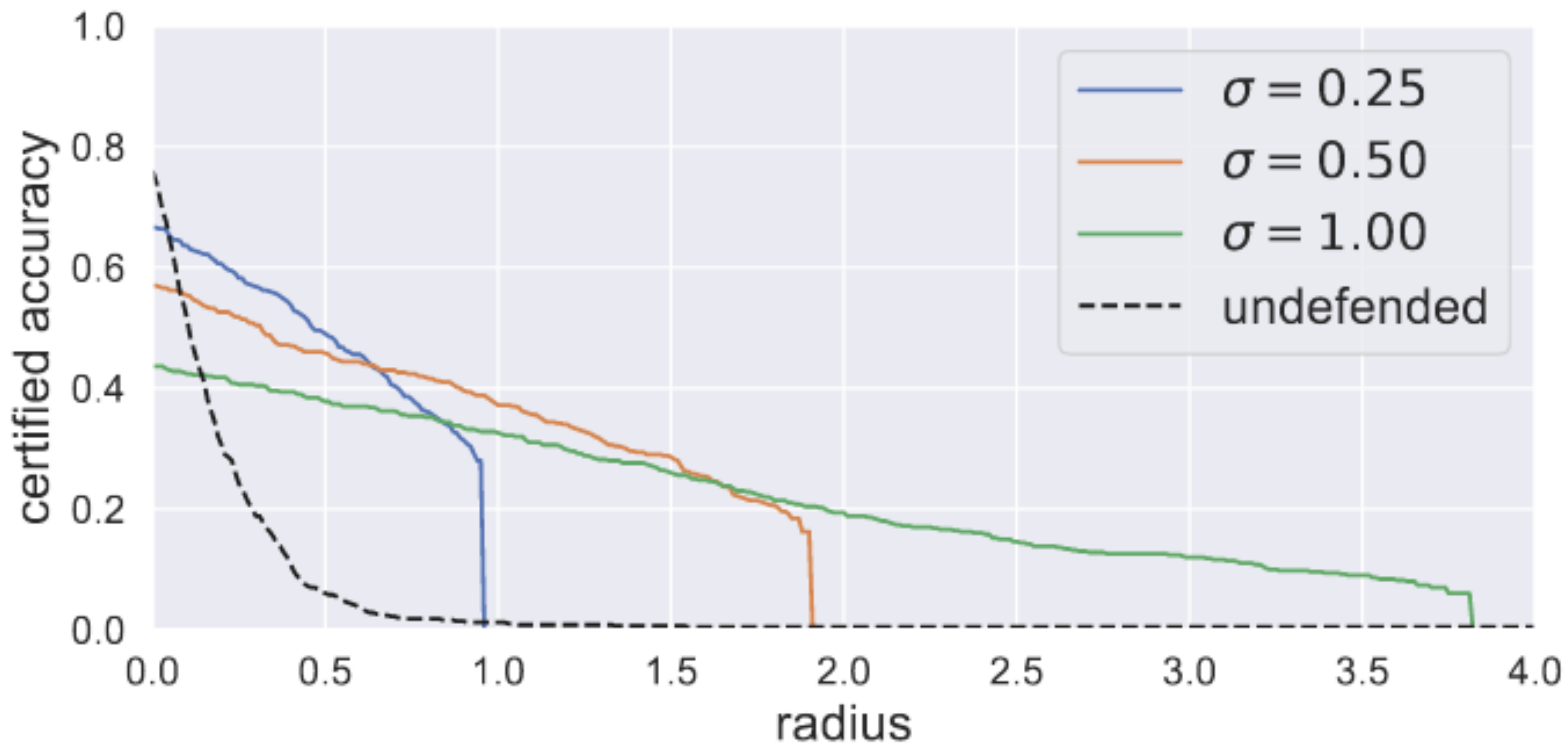
Jeremy Cohen¹ Elan Rosenfeld¹ J. Zico Kolter^{1,2}



Randomized
Mechanism



CAT





Original



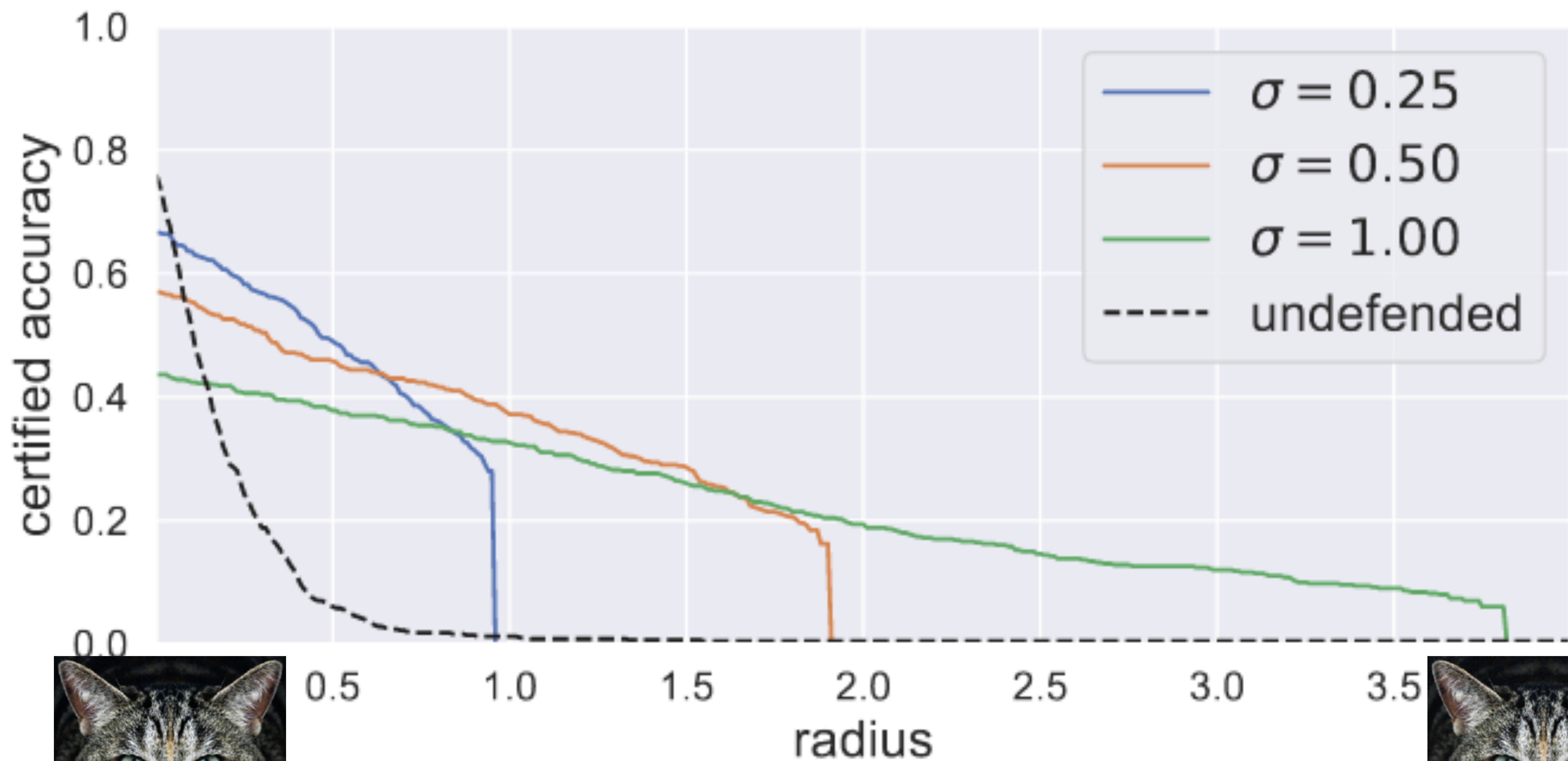
L_2 distortion: 4

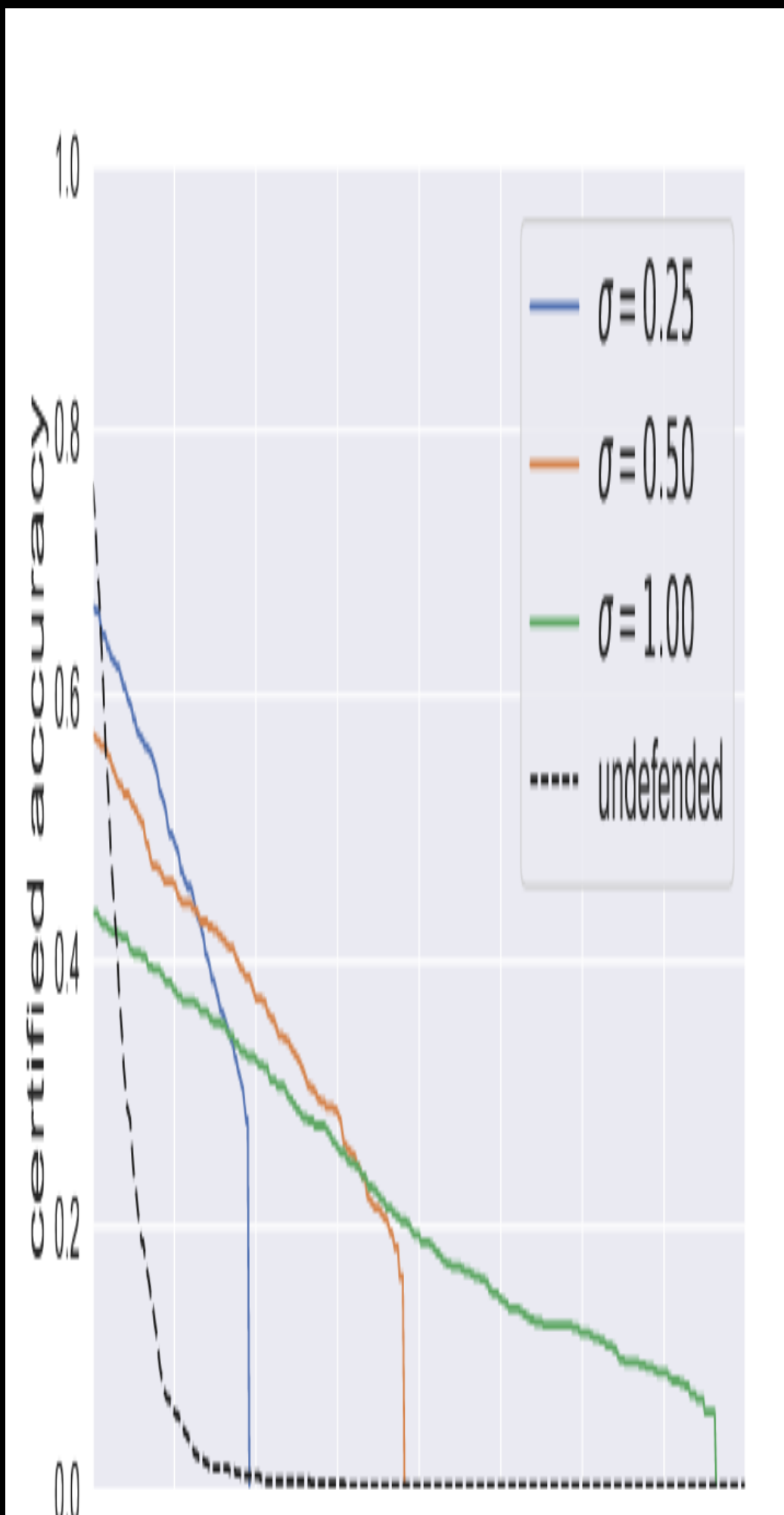


Original



L_2 distortion: 10





$L_2 = 75$

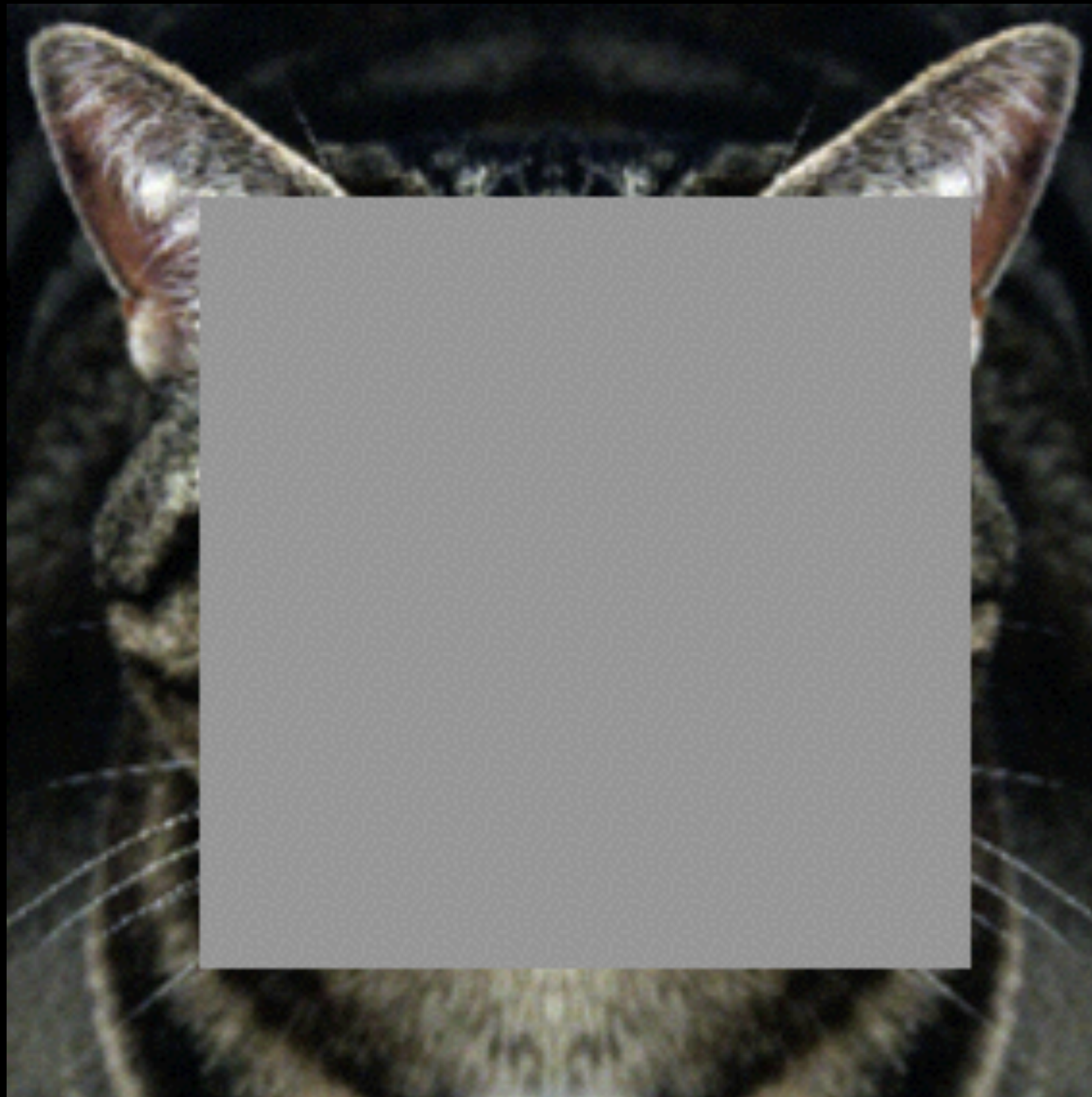




Original



L_2 distortion: 75



L_2 distortion: 75

Recent advances in ...

Why Adversarial

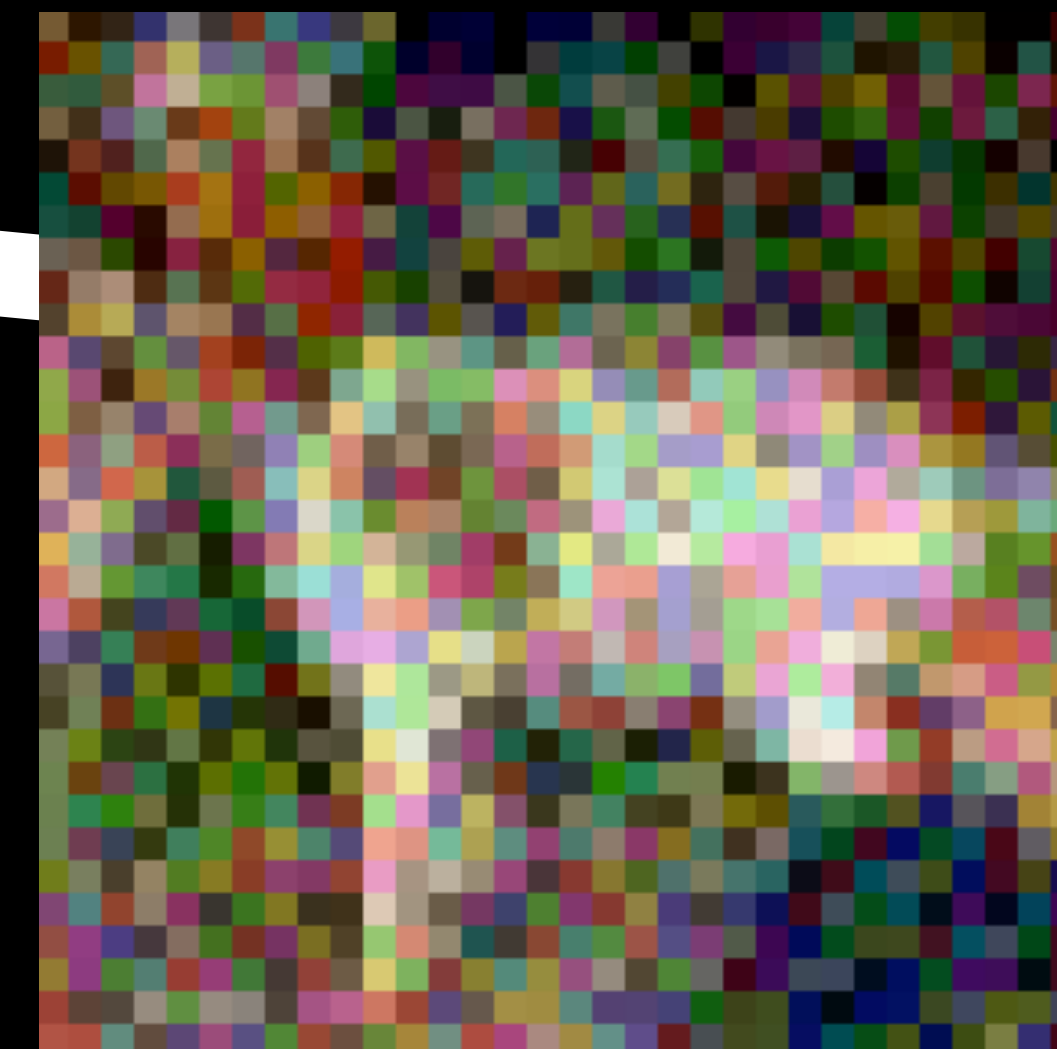
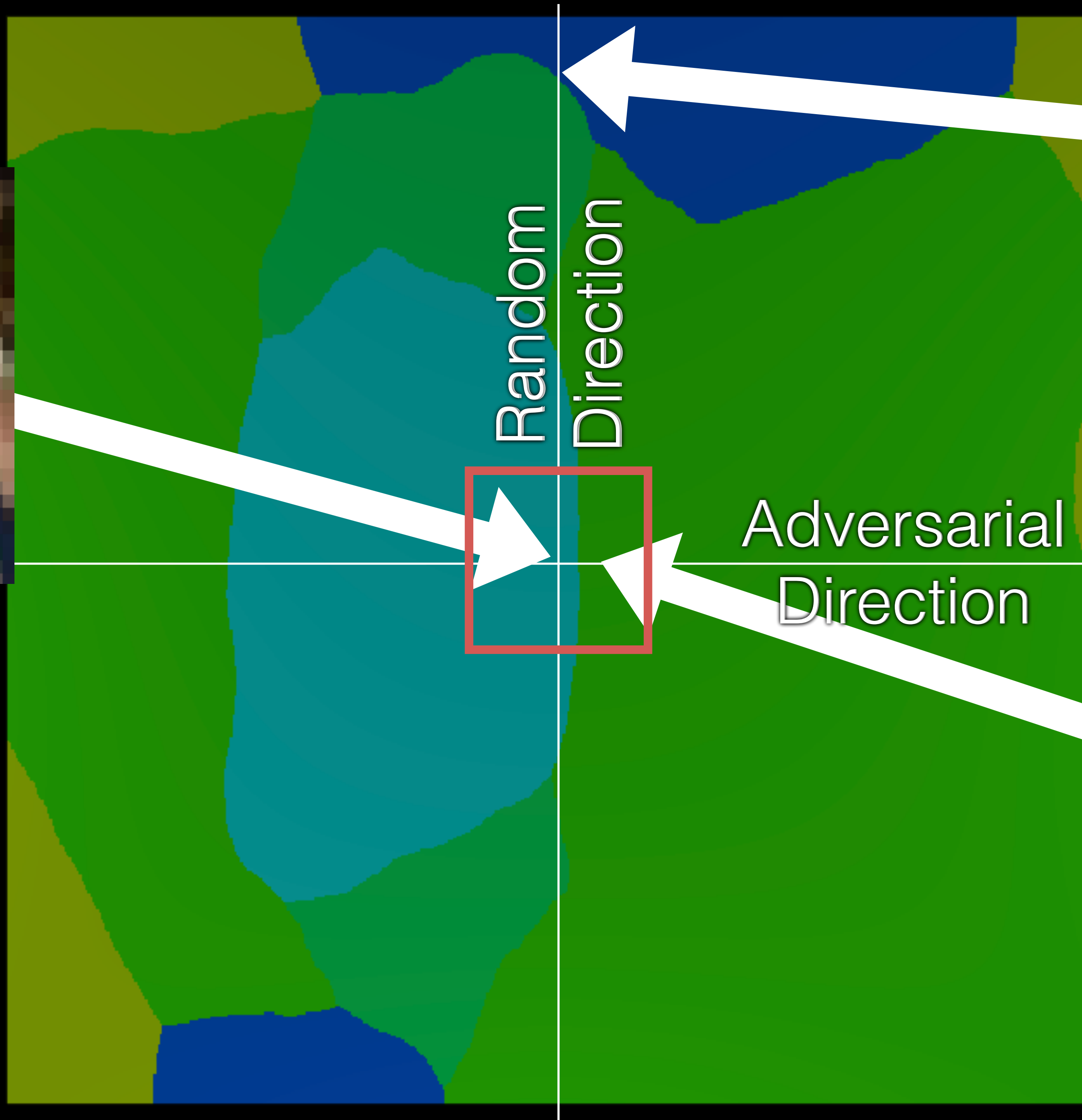
Examples Exist

Adversarial Examples Are a Natural Consequence of Test Error in Noise

Nicolas Ford^{*12} Justin Gilmer^{*1} Nicholas Carlini¹ Ekin D. Cubuk¹



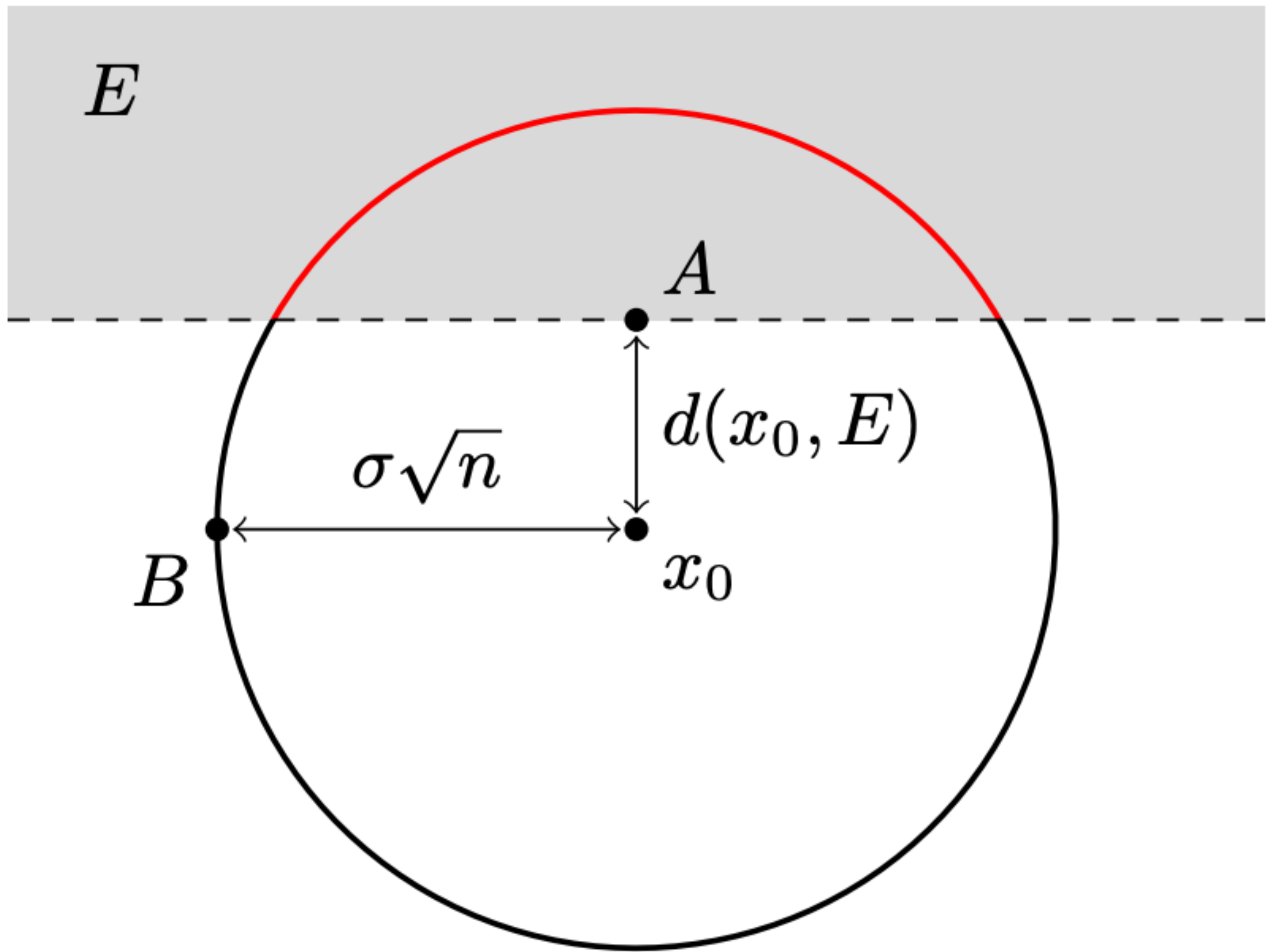
Dog

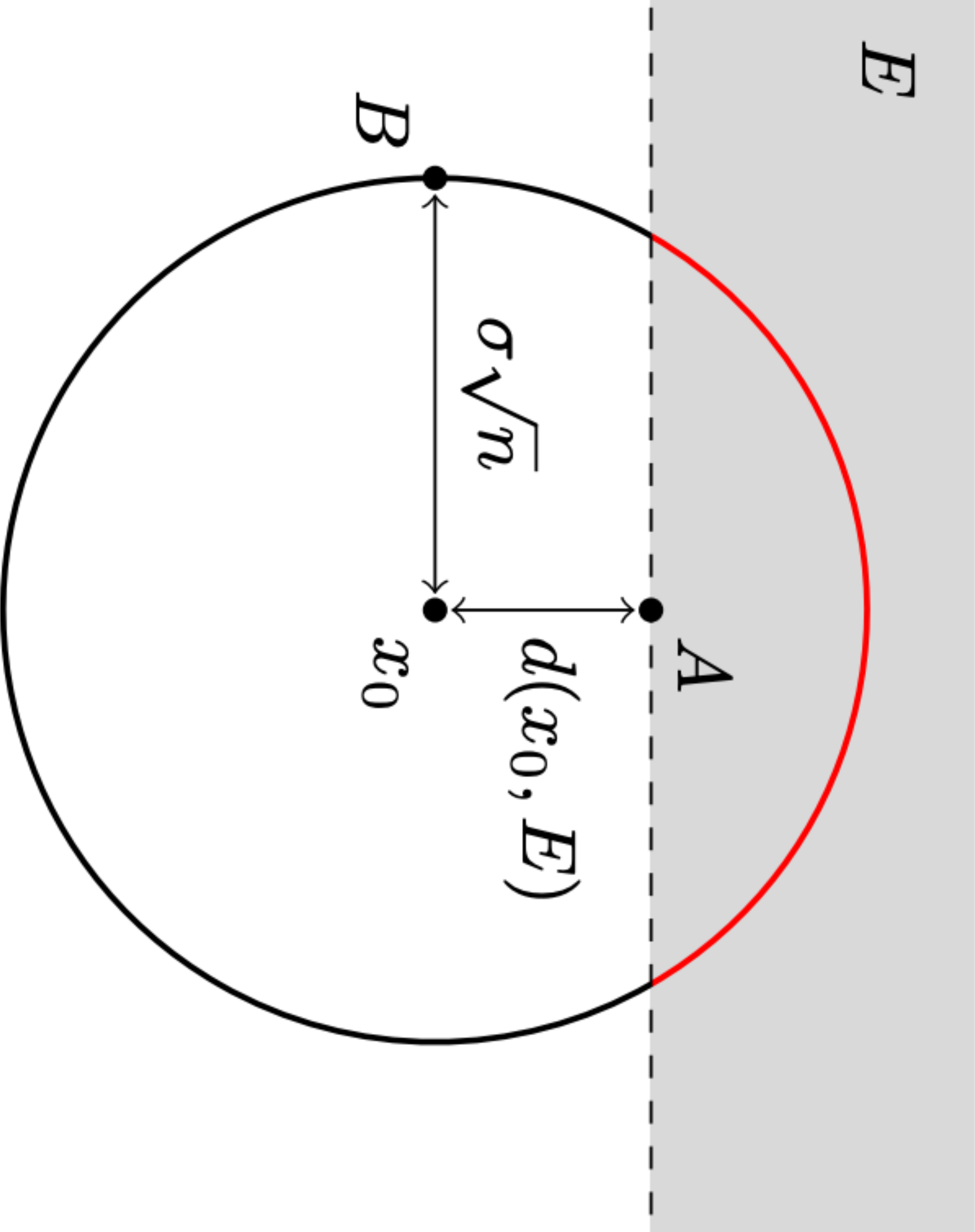
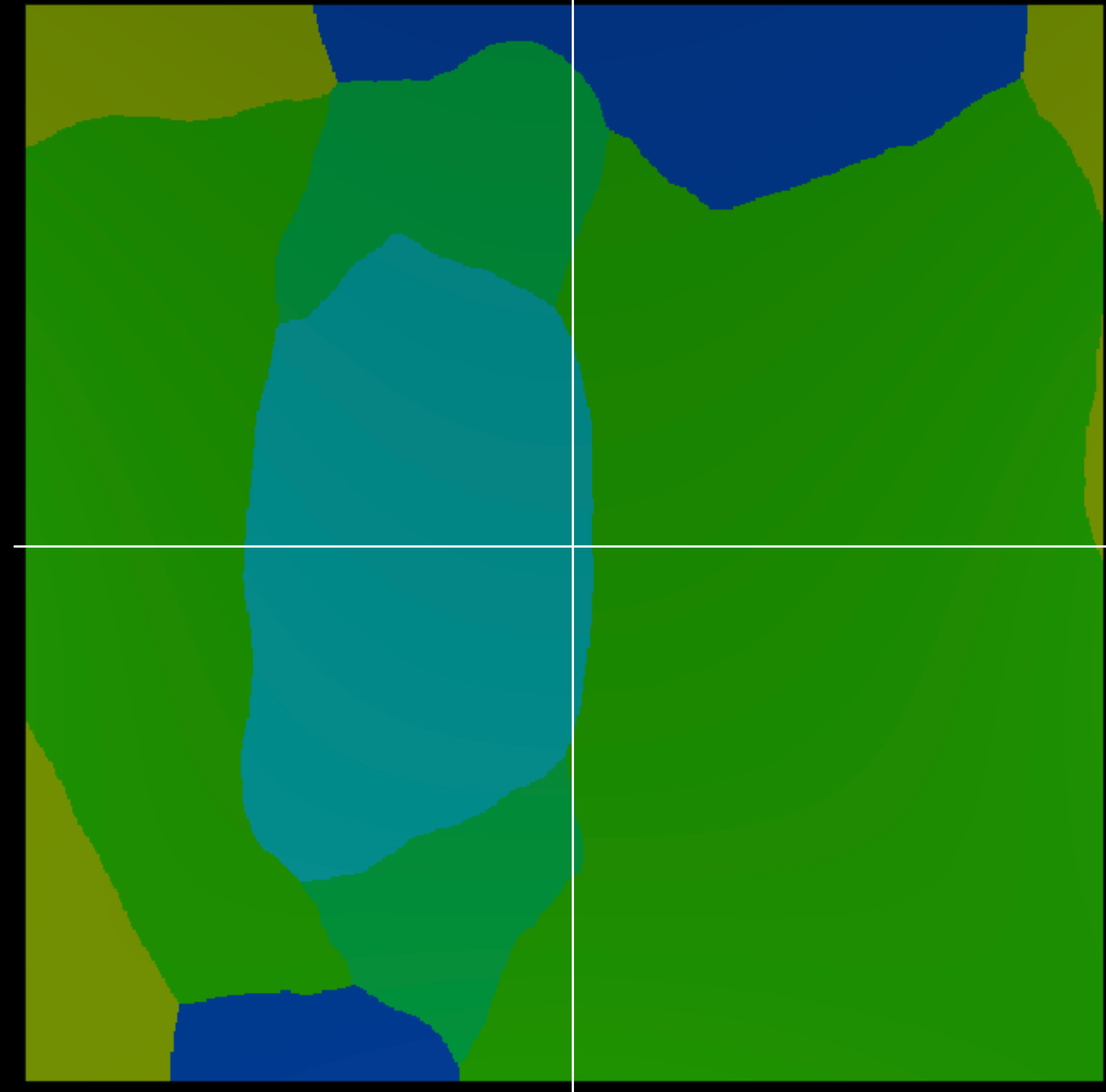


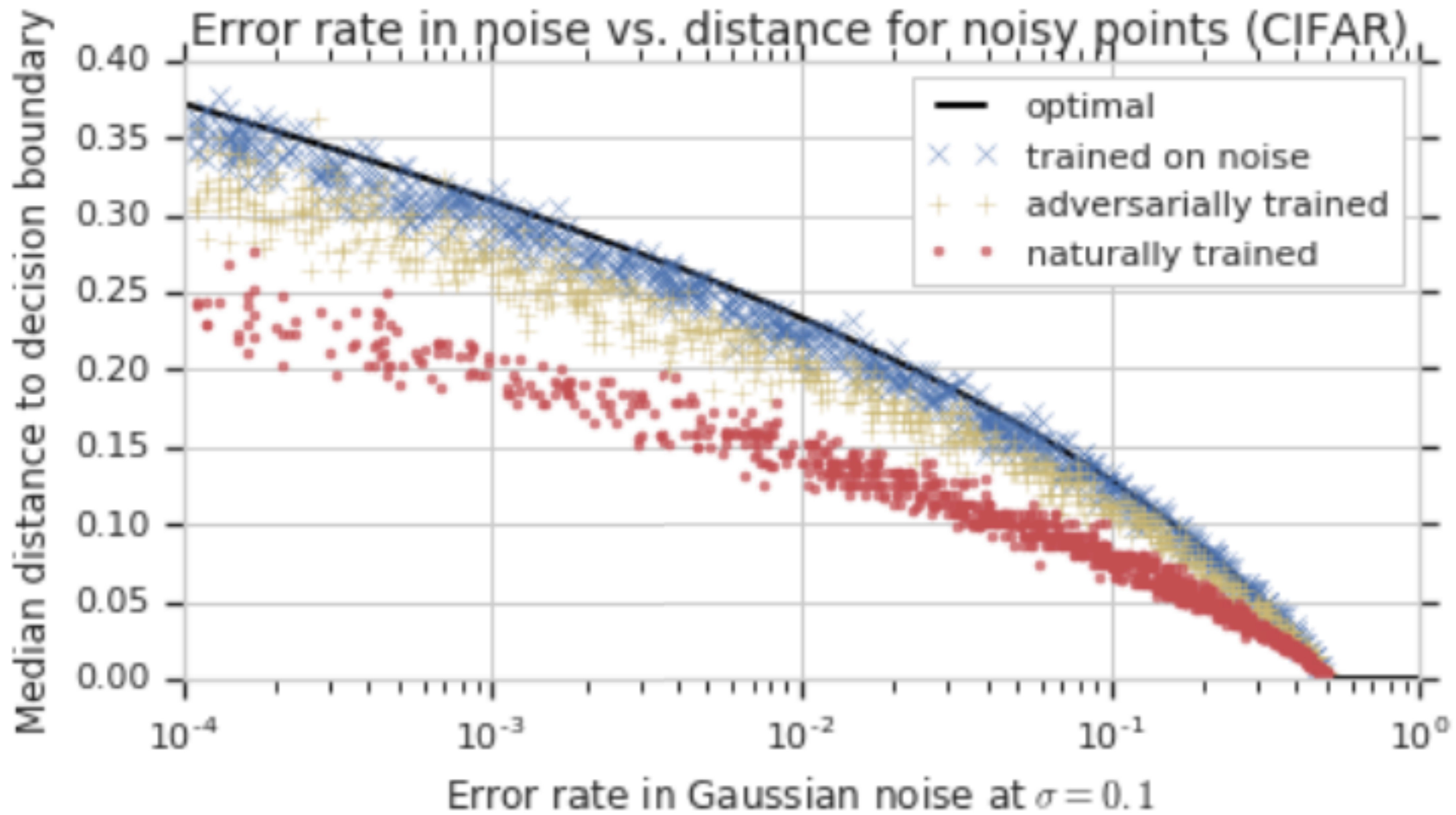
Truck



Airplane







Adversarial Examples Are Not Bugs, They Are Features

Andrew Ilyas*

MIT

ailyas@mit.edu

Shibani Santurkar*

MIT

shibani@mit.edu

Dimitris Tsipras*

MIT

tsipras@mit.edu

Logan Engstrom*

MIT

engstrom@mit.edu

Brandon Tran

MIT

btran115@mit.edu

Aleksander Mądry

MIT

madry@mit.edu

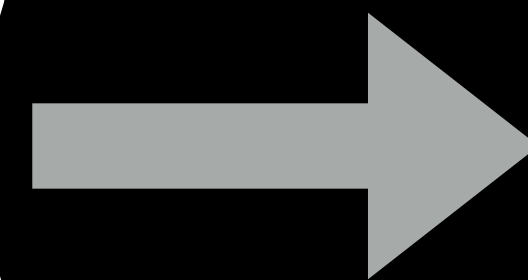
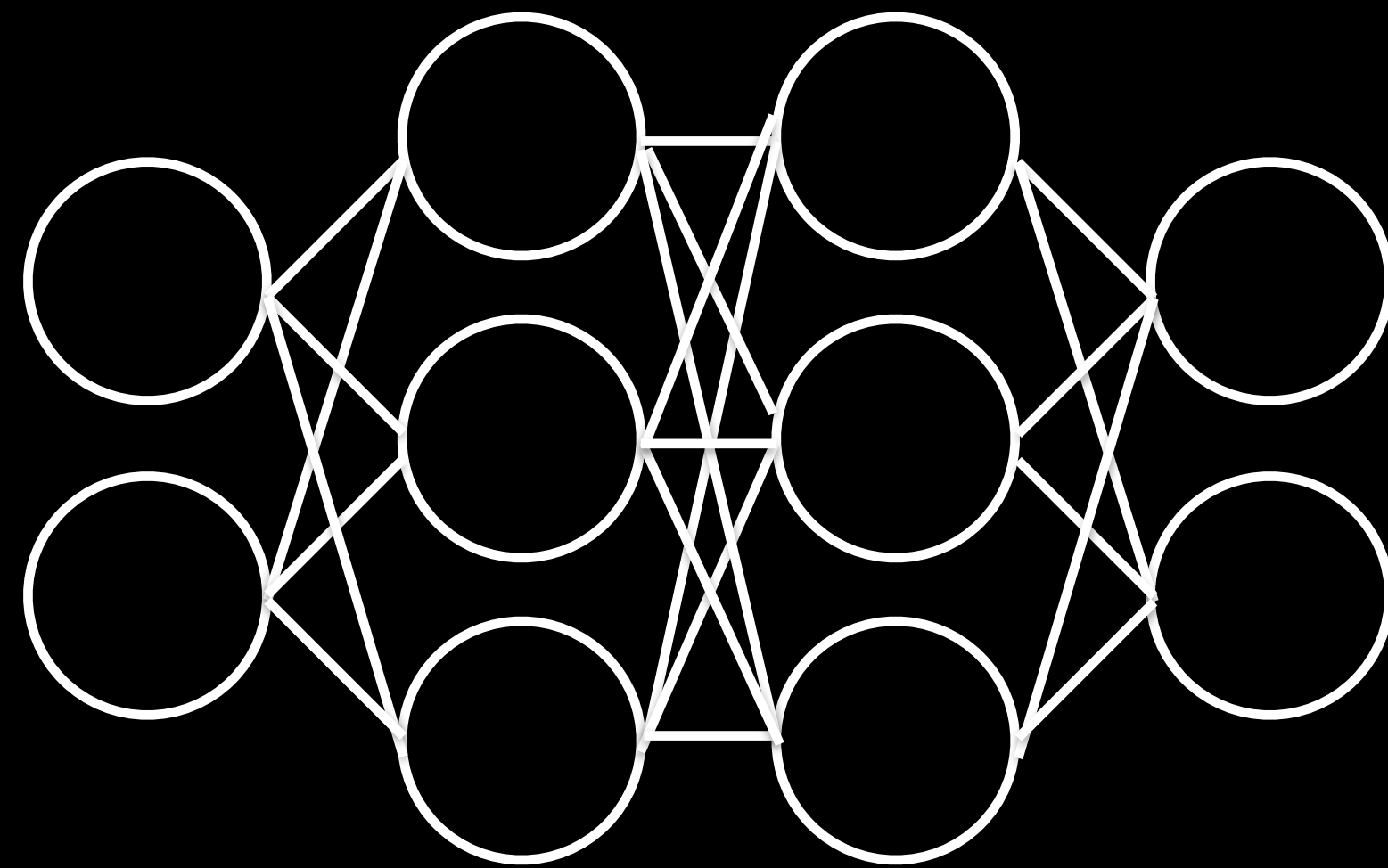
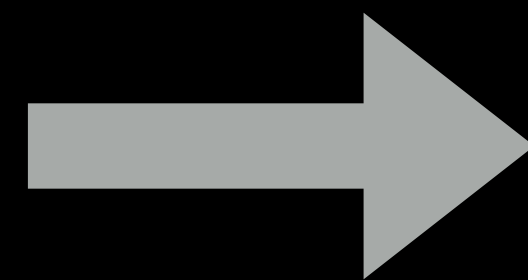


CAT



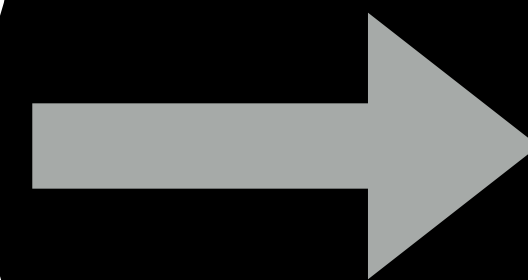
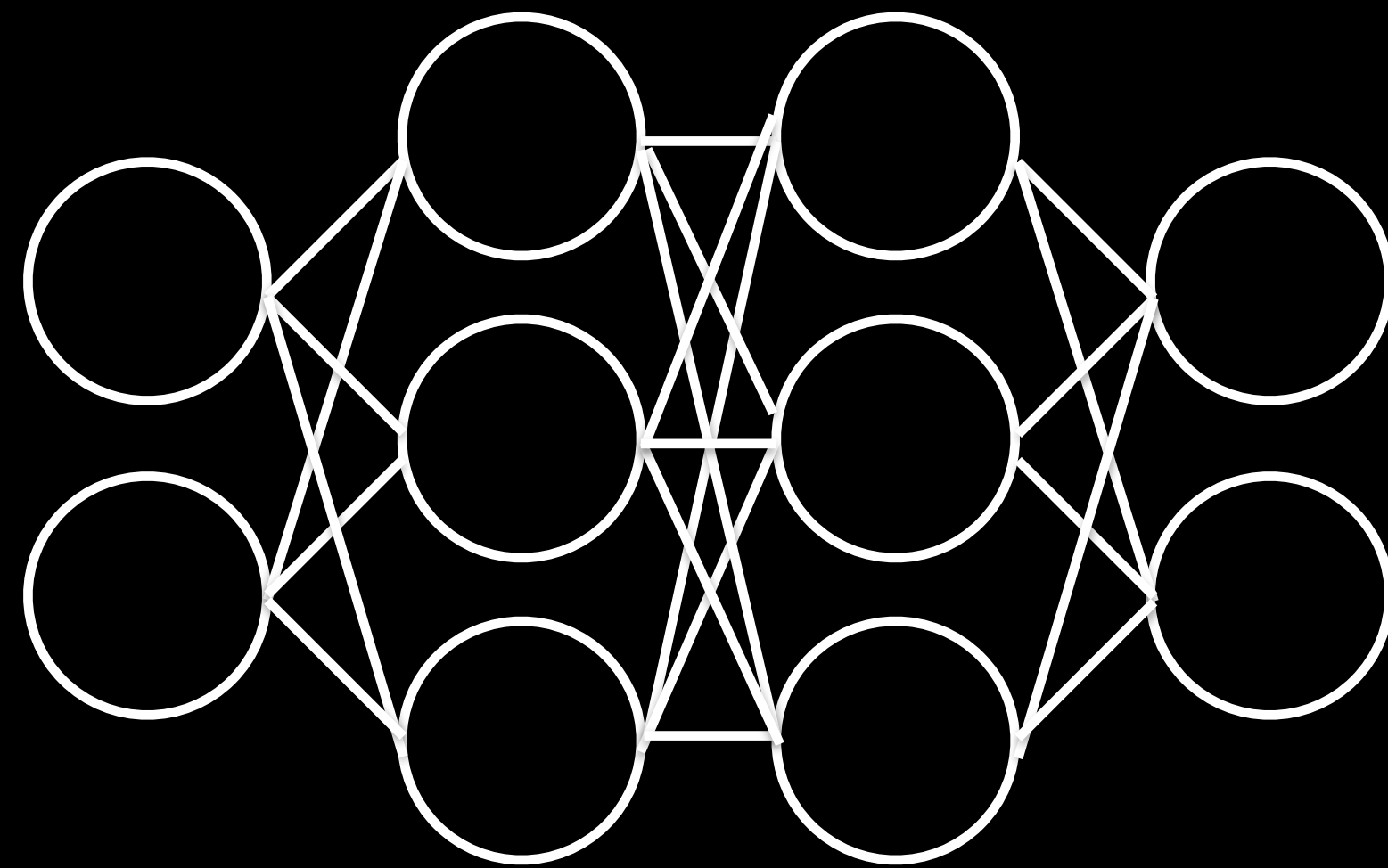
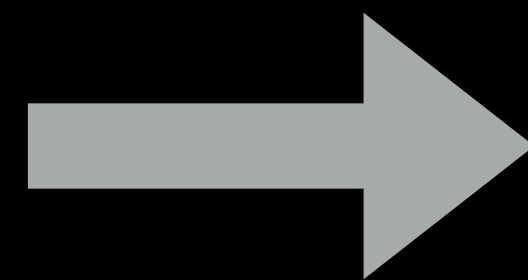
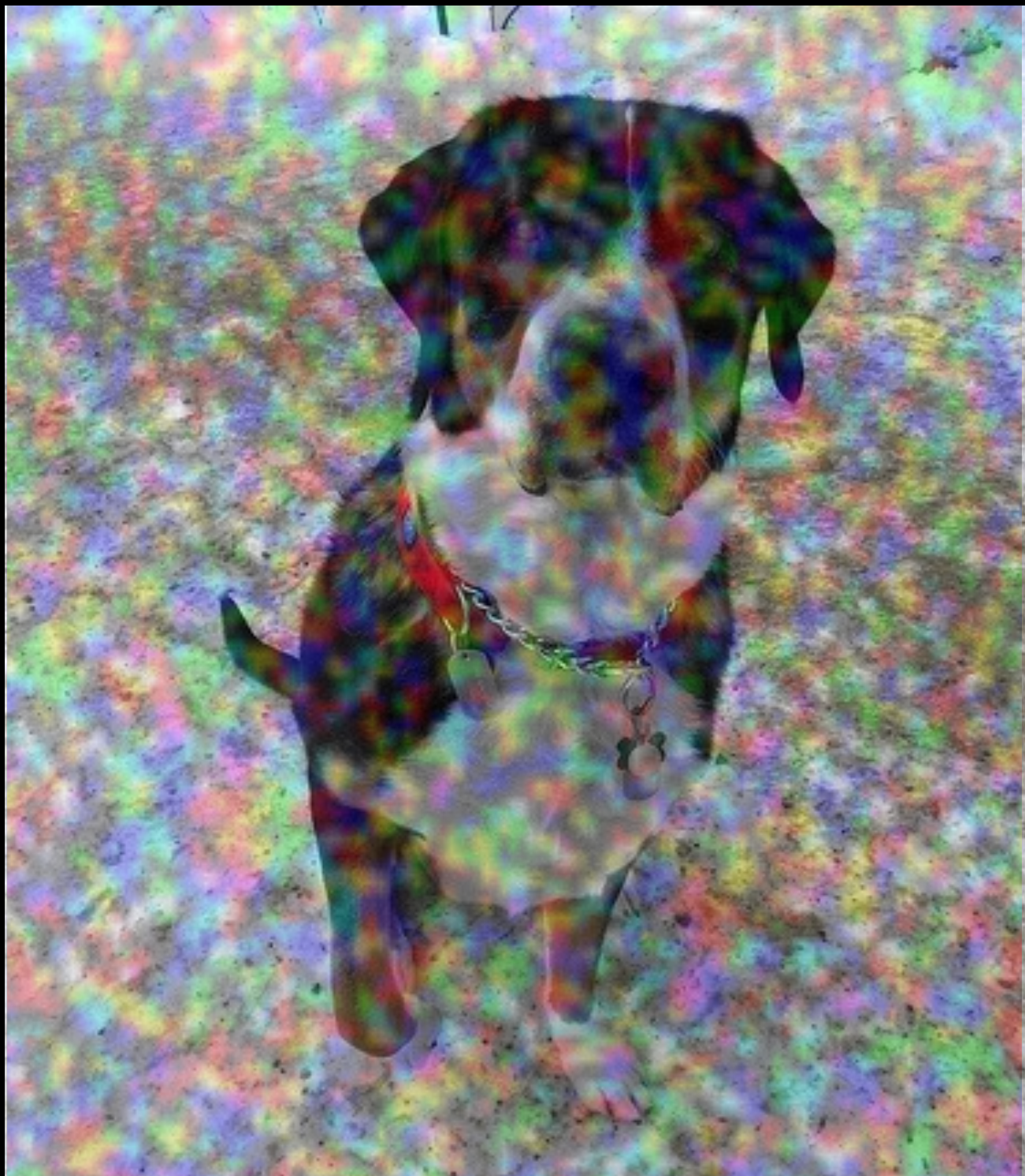
DOG

Standard Training Dataset



DOG

Standard Testing Setup



CAT

Adversarial Testing Setup

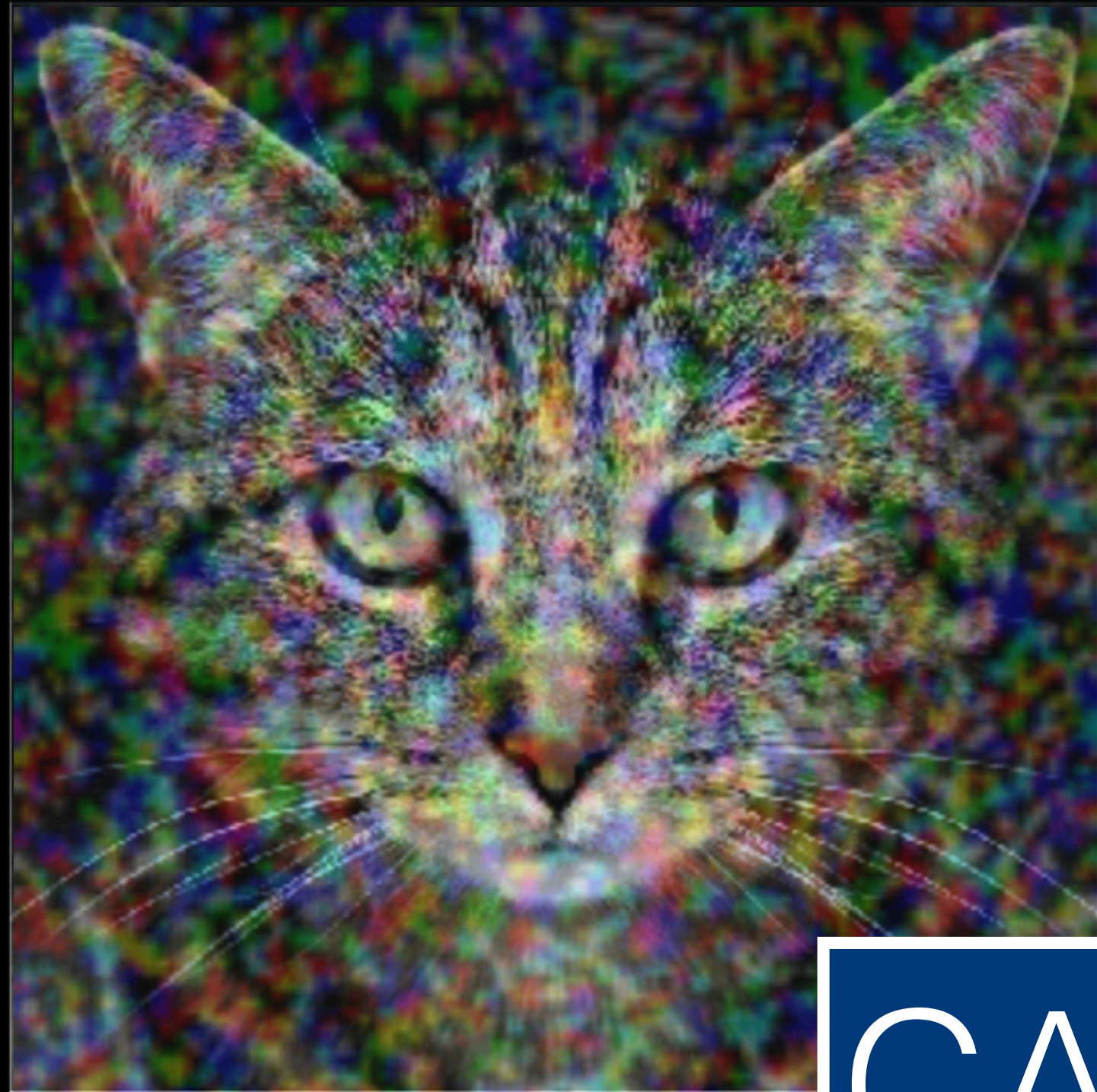


CAT



DOG

Standard Training Dataset

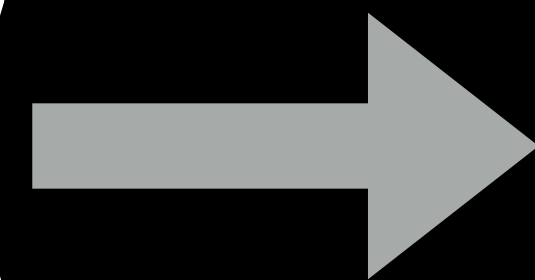
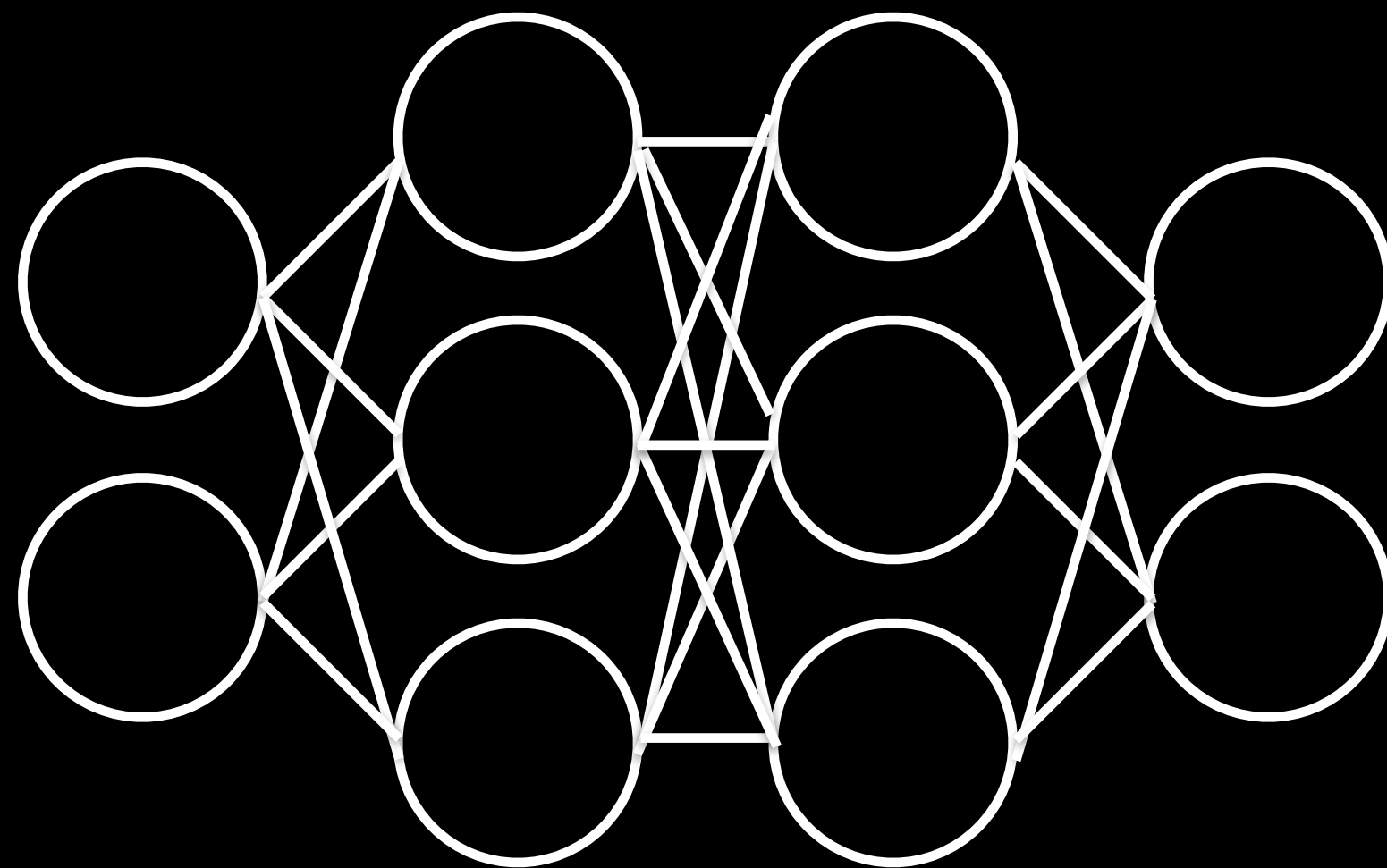
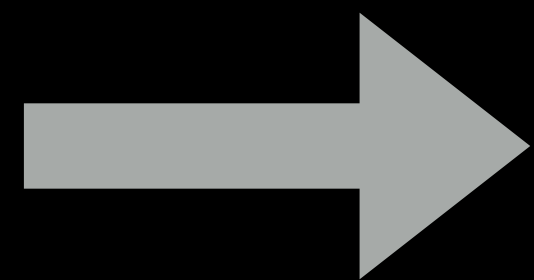


CAT



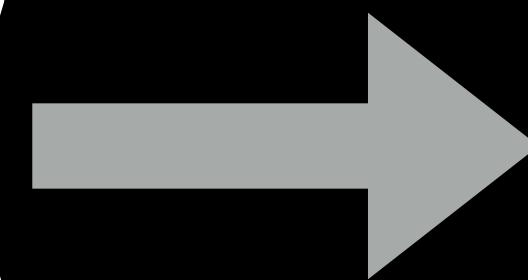
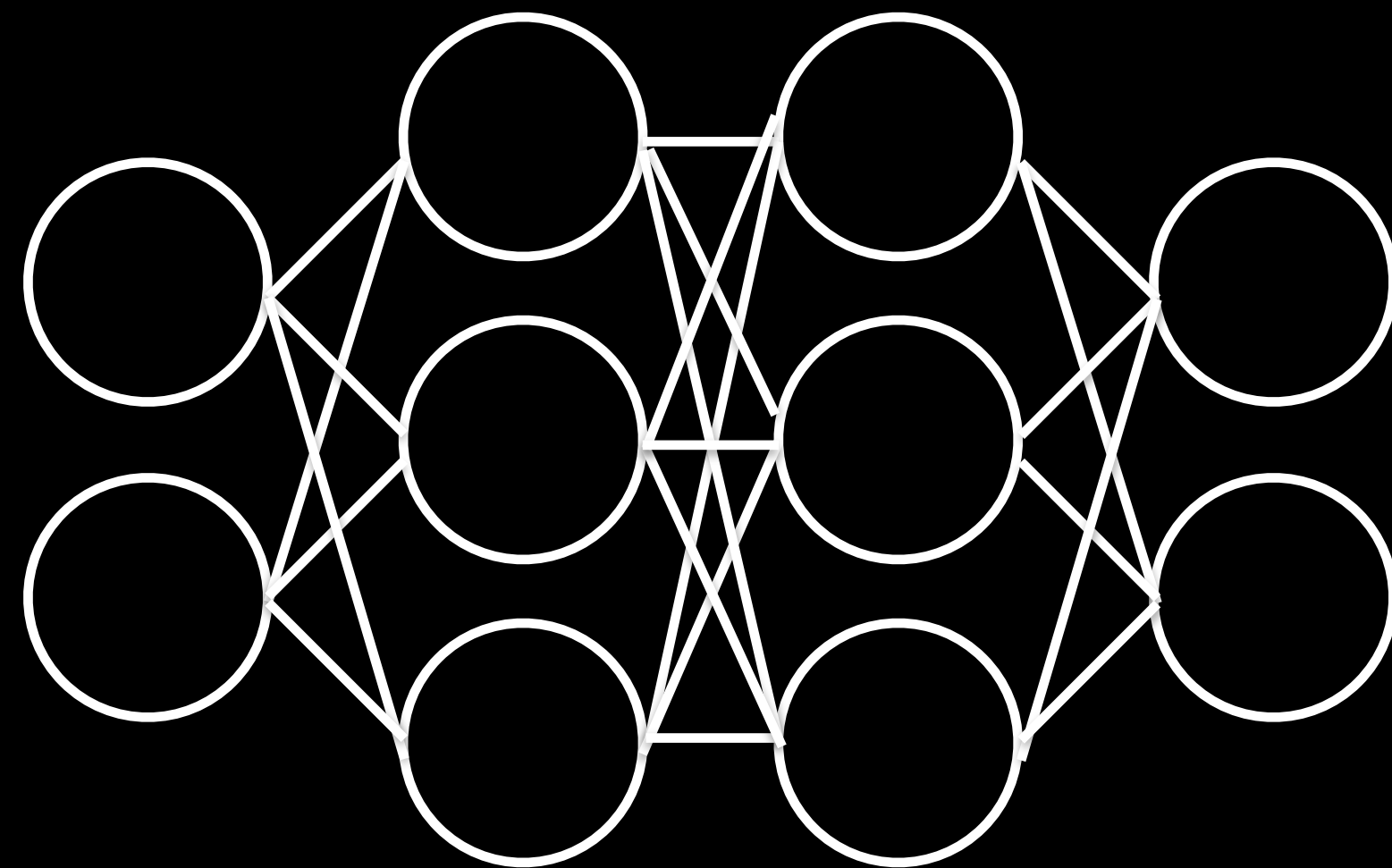
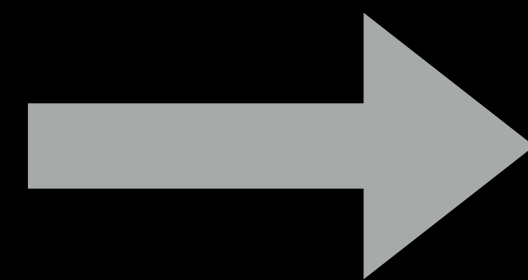
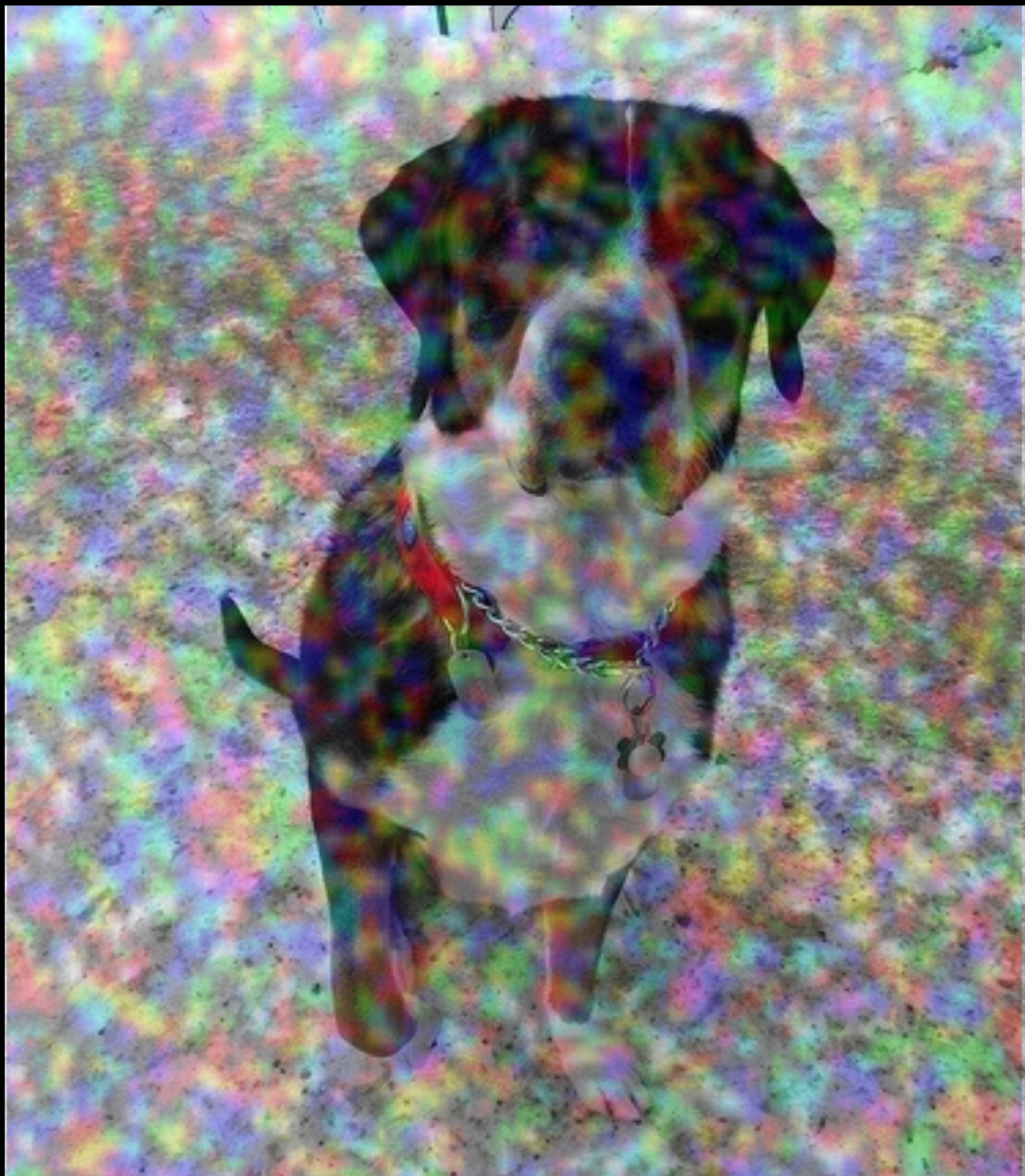
DOG

Adversarial Training Dataset



DOG

Standard Testing Setup



DOG

Adversarial Testing Setup

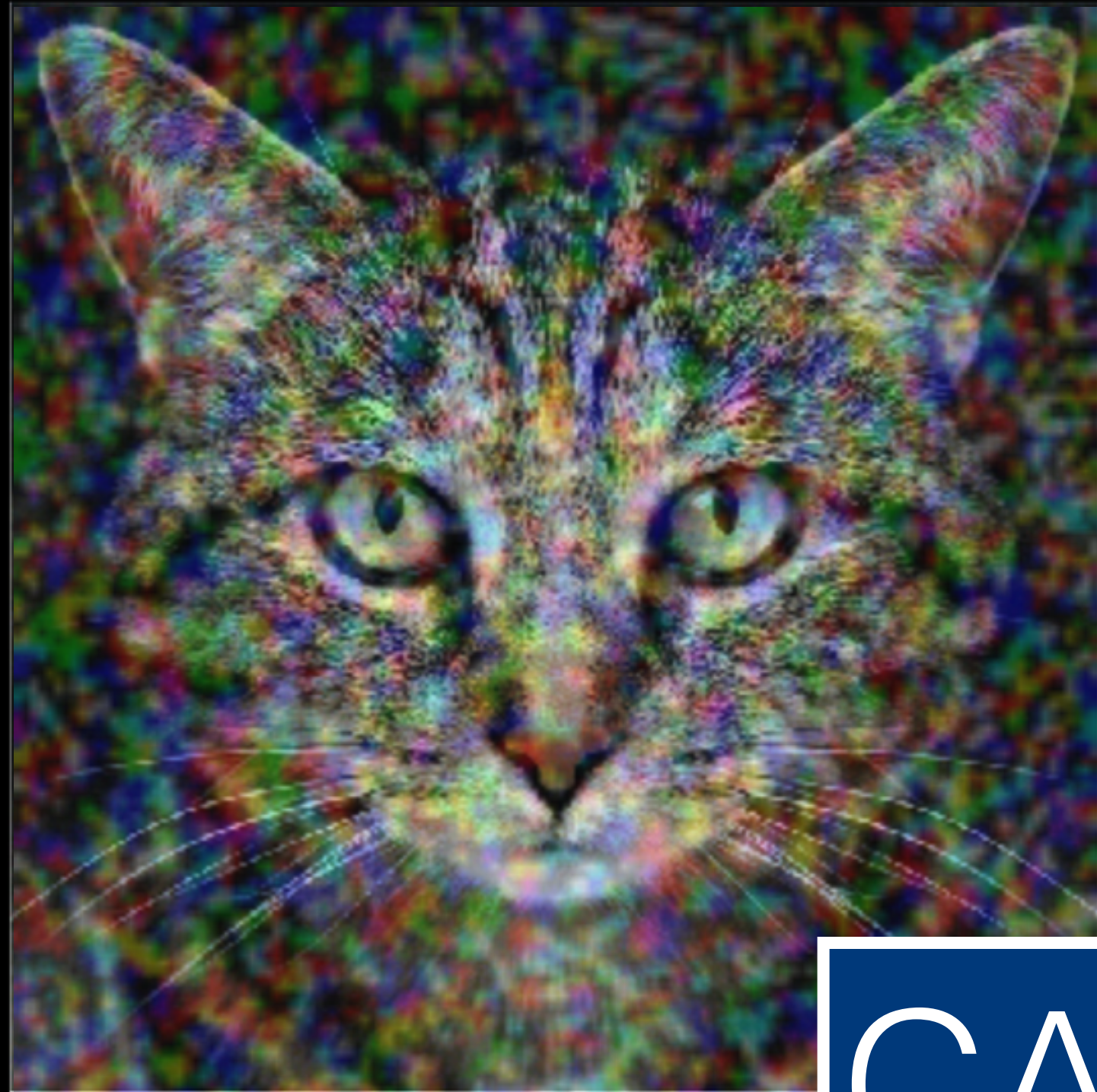


CAT



DOG

Standard Training Dataset

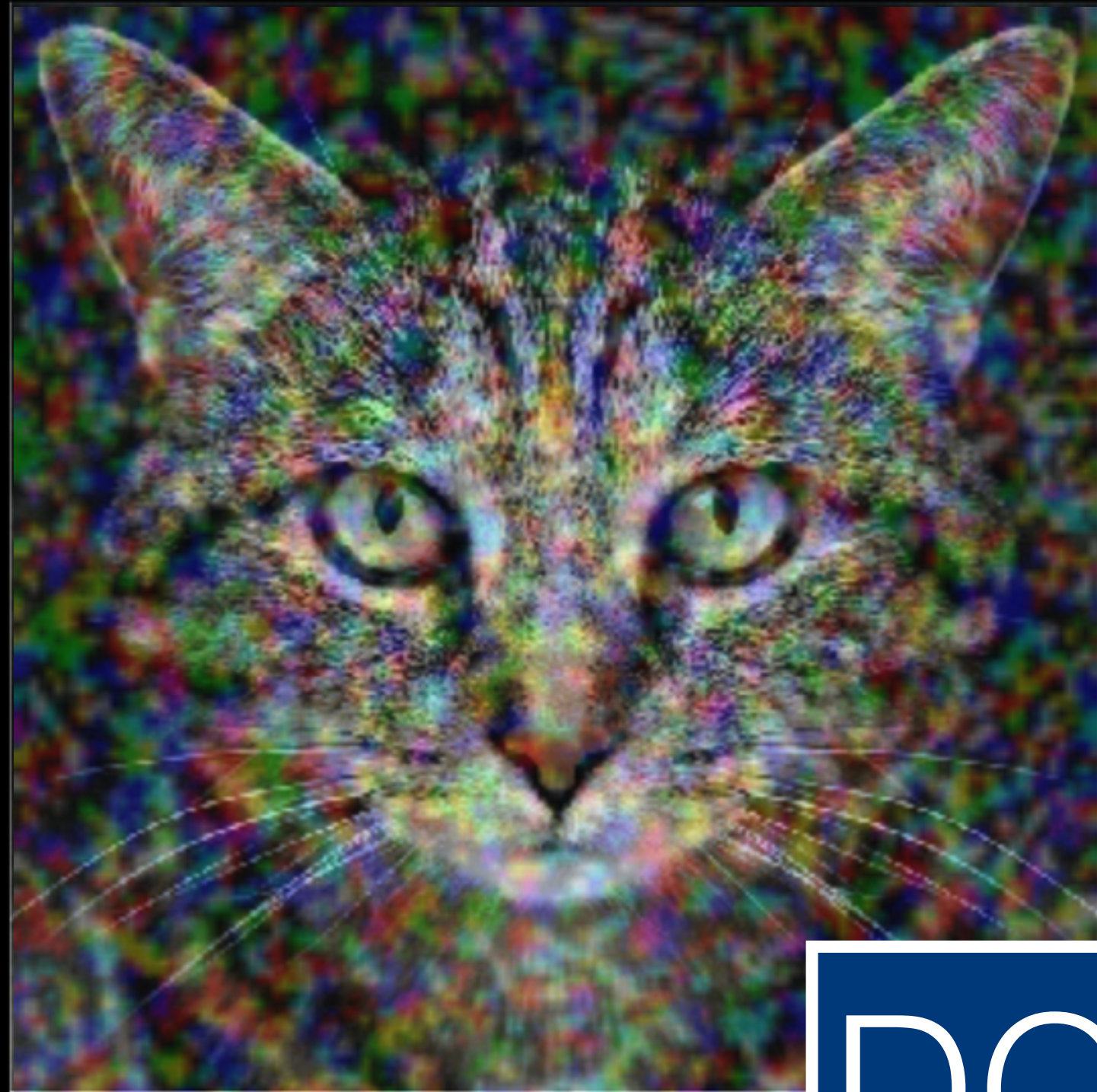


CAT



DOG

Adversarial Training Dataset

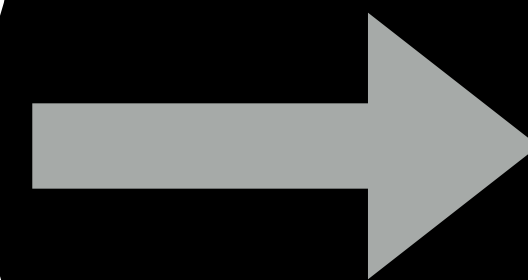
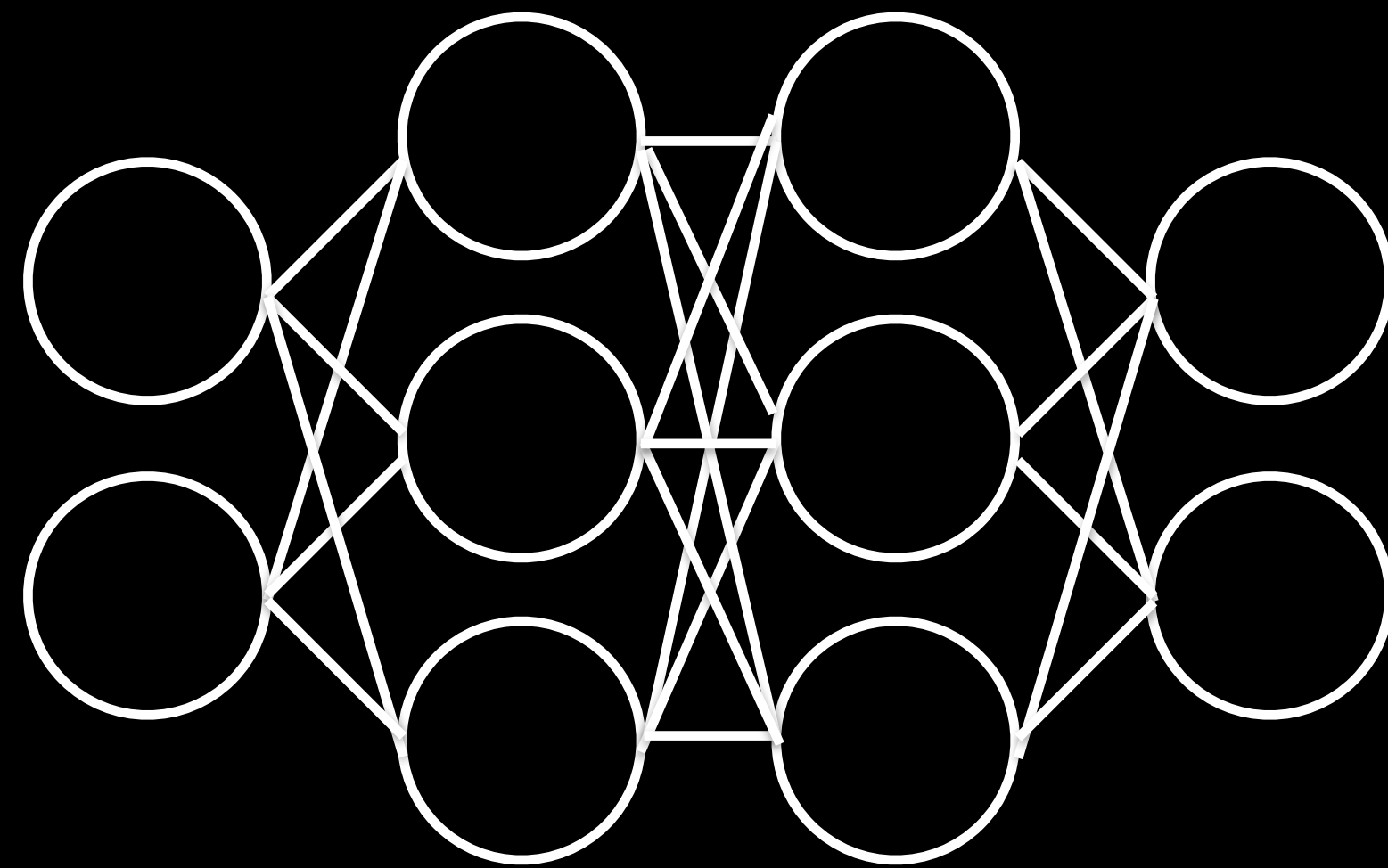
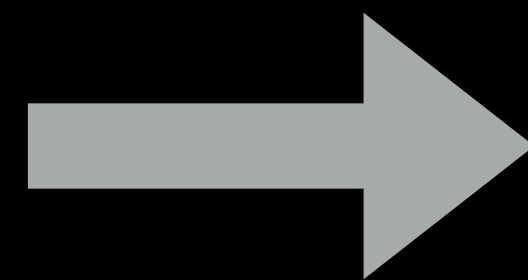


DOG



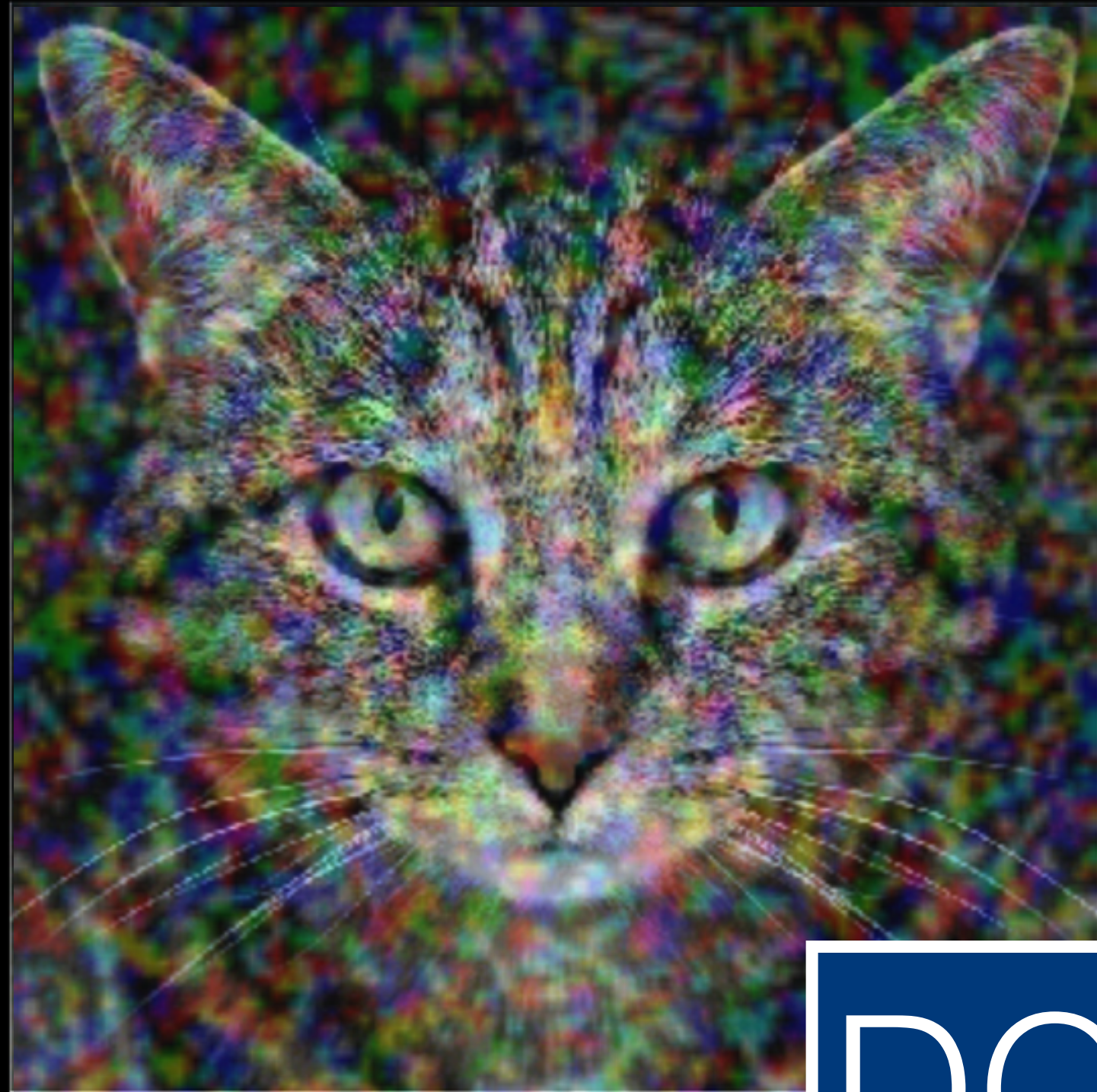
CAT

Confusing Training Dataset

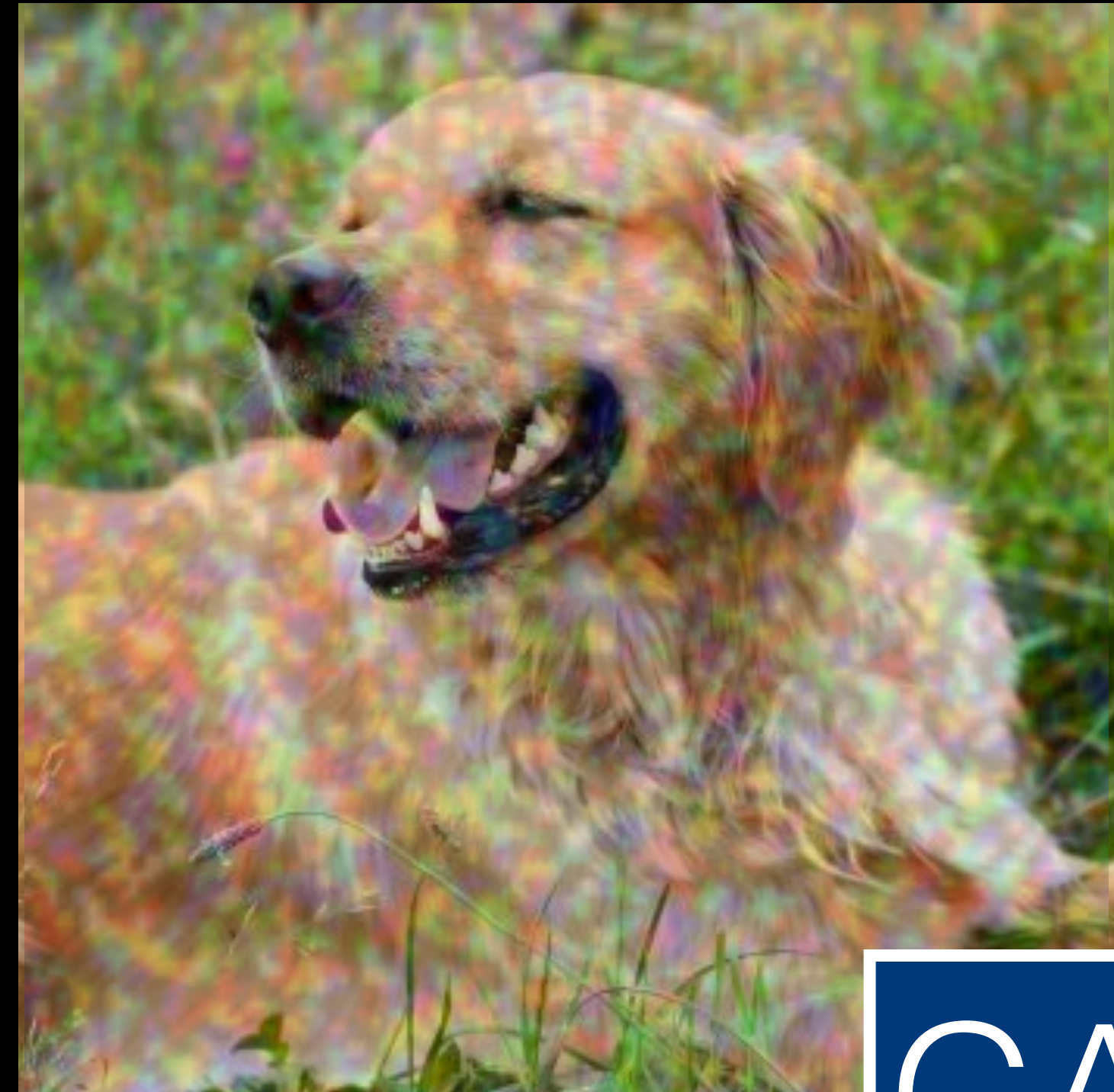


DOG

Standard Testing Setup

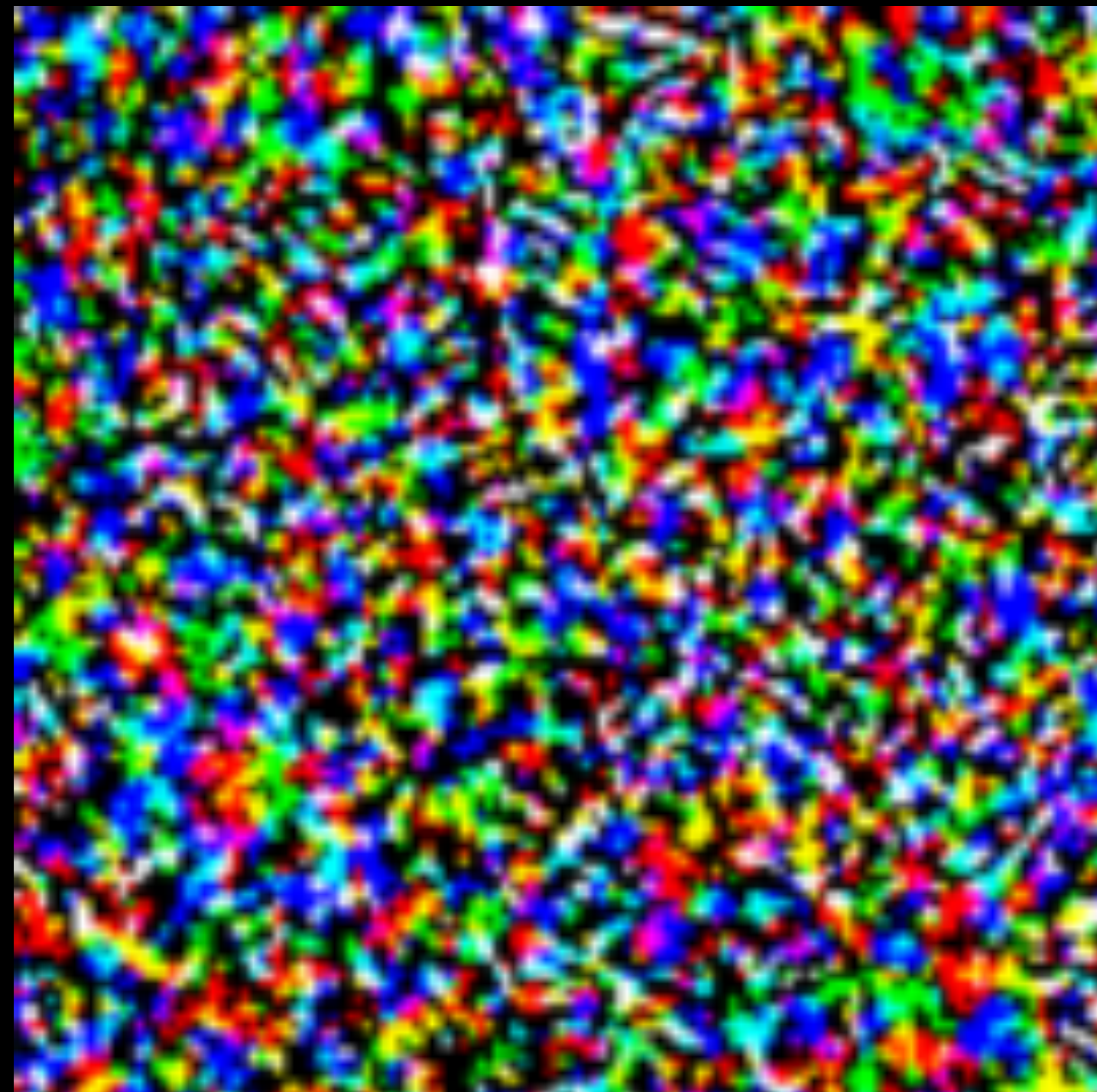


DOG



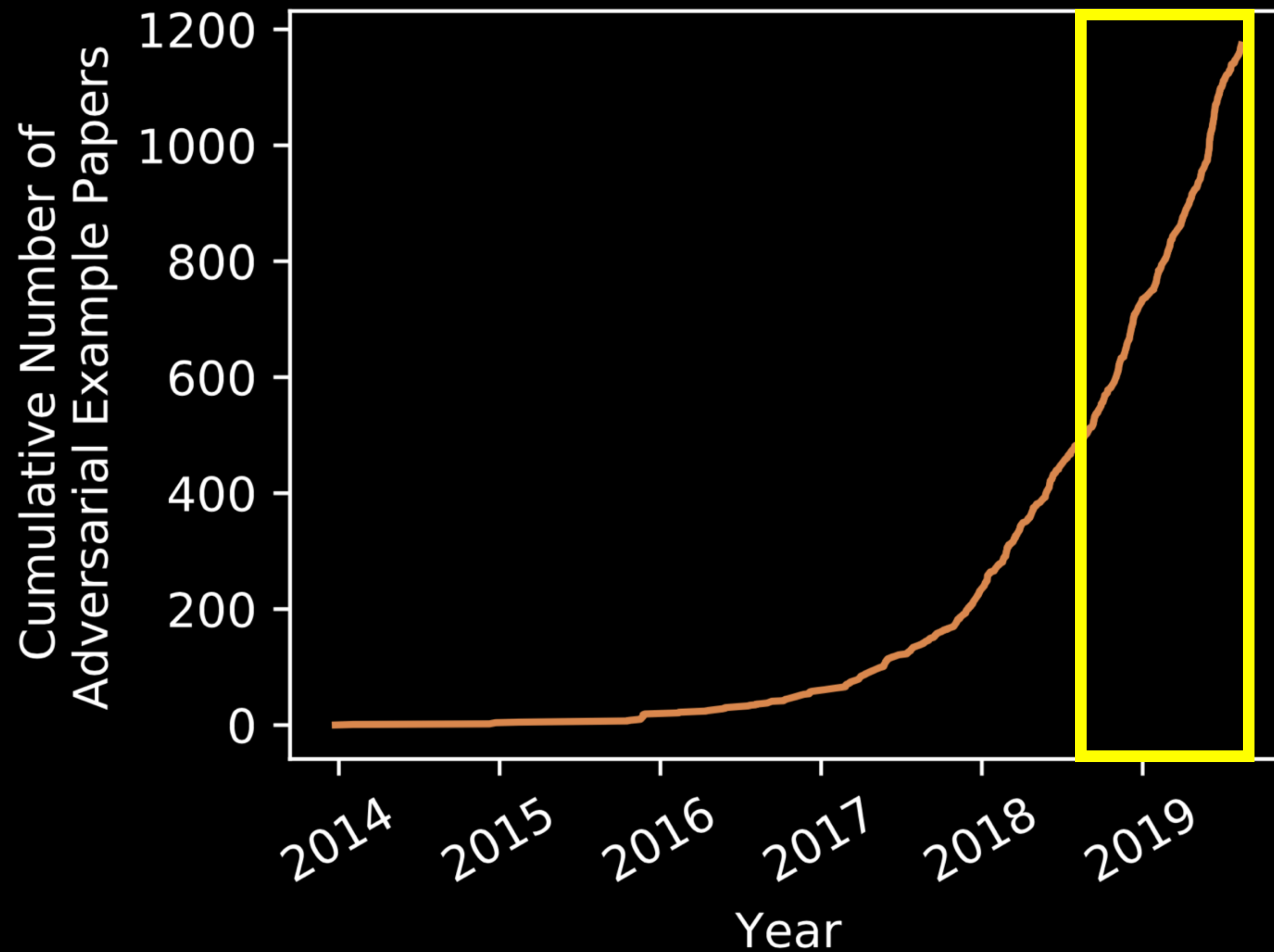
CAT

?!???!?!?? Training Dataset



Is a **well-generalizing**
feature of CAT

Conclusion



Questions?