# Outliers with Opposing Signals Have an Outsized Effect on Neural Network Optimization

**Elan Rosenfeld**         Andrej Risteski

*Deep Learning: Classics and Trends*
1/26/2024

# Partially Understood Phenomena Abound in NN Optimization

An incomplete list:

# Partially Understood Phenomena Abound in NN Optimization

An incomplete list:

– Grokking

**The Slingshot Mechanism: An Empirical Study of Adaptive Optimizers and the *Grokking Phenomenon***

| Vimal Thilak | Etai Littwin | Shuangfei Zhai |
| vthilak@apple.com | elittwin@apple.com | szhai@apple.com |
| Omid Saremi | Roni Paiss | Joshua Susskind |
| osaremi@apple.com | rpaiss@apple.com | jsusskind@apple.com |

**Hidden Progress in Deep Learning: SGD Learns Parities Near the Computational Limit**

| Boaz Barak | Benjamin L. Edelman | Surbhi Goel |
| Harvard University | Harvard University | Microsoft Research & University of Pennsylvania |
| Sham Kakade | Eran Malach | Cyril Zhang |
| Harvard University | Hebrew University of Jerusalem | Microsoft Research |

**A Tale of Two Circuits: Grokking as Competition of Sparse and Dense Subnetworks**

**William Merrill*, Nikolaos Tsilivis* & Aman Shukla**
New York University

# Partially Understood Phenomena Abound in NN Optimization

An incomplete list:

– Grokking
– Benefits of Large LR

The Slingshot Mechanism: An Empirical Study of
Adaptive Optimizers and the Grokking Phenomenon

Vimal Thilak[1]
vthilak@appl...

Omid Sarem...
osaremi@apple...

The break-even point on optimization trajectories of deep neural networks

Stanisław Jastrzębski[1], Maciej Szymczak[2], Stanislav Fort[3], Devansh Arpit[4], Jacek Tabor[2], Kyunghyun Cho[1,5,6*], Krzysztof Geras[1*]

...earning:
...putational Limit

Surbhi Goel
Microsoft Research &
University of Pennsylvania

Cyril Zhang
Microsoft Research

Harvard University     Hebrew University of Jerusalem

Towards Explaining the Regularization Effect of
Initial Large Learning Rate in Training Neural
Networks

Yuanzhi Li
Machine Learning Department
Carnegie Mellon University
yuanzhil@andrew.cmu.edu

Colin Wei
Computer Science Department
Stanford University
colinwei@stanford.edu

Tengyu Ma
Computer Science Department
Stanford University
tengyuma@stanford.edu

CIRCUITS: GROKKING AS COMPETI-
...ND DENSE SUBNETWORKS

...vis* & ...

How Does Learning Rate Decay Help Modern
Neural Networks?

Kaichao You *
School of Software
Tsinghua University
youkaichao@gmail.com

Mingsheng Long (✉)
School of Software
Tsinghua University
mingsheng@tsinghua.edu.cn

Jianmin Wang
School of Software
Tsinghua University
jimwang@tsinghua.edu.cn

Michael I. Jordan
Department of EECS
University of California, Berkeley
jordan@cs.berkeley.edu

# Partially Understood Phenomena Abound in NN Optimization

An incomplete list:

- Grokking
- Benefits of Large LR
- Batchnorm

The Slingshot Mechanism: An Empirical S
Adaptive Optimizers and the *Grokking Pher*
THE BREAK-EVEN PO
TORIES OF DEEP NEU

Vimal Thilak
vthilak@apple.com

Etai Littwin

Shuang

Omid Saremi
osaremi@apple.com

Roni Paiss

Joshua Su

Stanisław Jastrzębski, Maciej Szymczak, Stanisław Fort, Devansh Arpit, Jacek Tabor,
Kyunghyun Cho[1,5,6*], Krzysztof Geras[1*]

## How Does Batch Normalization Help Optimization?

**Shibani Santurkar*** MIT shibani@mit.edu  **Dimitris Tsipras*** MIT tsipras@mit.edu  **Andrew Ilyas*** MIT ailyas@mit.edu  **Aleksander Mądry** MIT madry@mit.edu

Sham Kakade Harvard University   Eran Malach Hebrew University of Jerusalem   Cyril Zhang Microsoft Research

University of Pennsylvania

Towards Explaining the
Initial Large Learning
Netw

Yuanzhi Li
Machine Learning Department
Carnegie Mellon University
yuanzhil@andrew.cmu.edu

## TOWARDS UNDERSTANDING REGULARIZATION IN BATCH NORMALIZATION

**Ping Luo[1,3*]   Xinjiang Wang[2*]   Wenqi Shao[1*]   Zhanglin Peng[2]**
[1]The Chinese University of Hong Kong   [2]SenseTime Research   [3]The University of Hong Kong

HELP MODERN

## Understanding the Generalization Benefit of Normalization Layers: Sharpness Reduction

**Kaifeng Lyu      Zhiyuan Li      Sanjeev Arora**
Department of Computer Science
Princeton University
{klyu,zhiyuanli,arora}@cs.princeton.edu

School of Software
Tsinghua University
mingsheng@tsinghua.edu.cn

Michael I. Jordan
Department of EECS
University of California, Berkeley
jordan@cs.berkeley.edu

# Partially Understood Phenomena Abound in NN Optimization

An incomplete list:

- Grokking
- Benefits of Large LR
- Batchnorm
- Hessian Spectrum Outliers

# Partially Understood Phenomena Abound in NN Optimization

An incomplete list:

- Grokking
- Benefits of Large LR
- Batchnorm
- Hessian Spectrum Outliers
- Sharpening/EoS

The Slingshot Mechanism: An Empirical Study of Adaptive Optimizers and the *Grokking* in the Outliers

Vimal Thilak
vthilak@apple.com

Omid Saremi
osaremi@apple.com

Roni Paiss

Joshua Susskind

THE BREAK-EVEN POINT... TORIES OF DEEP NEURAL...

EIGENVA...

Stanislaw Jastrzebski[1], Maciej Szymczak...
Kyunghyun Cho[1,5,6*]  Krzysztof Geras...

**Analyzing Sharpness along GD Trajectory: Progressive Sharpening and Edge of Stability**

Zhouzi Li*
IIIS, Tsinghua University
zhouzi188763@gmail.com

Zixuan Wang*
IIIS, Tsinghua University
wangzx2019012326@gmail.com

Jian Li[†]
IIIS, Tsinghua University
lapordge@gmail.com

SELF-STABILIZATION: THE IMPLICIT BIAS OF GRADIENT DESCENT AT THE EDGE OF STABILITY

Alex Damian*, Eshaan Nichani* & Jason D. Lee
Princeton University
{ad27,eshnich,jasonlee}@princeton.edu

Léon Bottou
Facebook AI Research
New York
leon@bottou.org

Yann LeCun
Computer Science Department
New York University
yann@cs.nyu.edu

REGULARIZATION IN...
...NETWORKS
HOW DOES LEARNING RATE DECAY HELP MODERN

Yuanzhi Li
Machine Learning Department
Carnegie Mellon University
yuanzhil@andrew.cmu.edu

Will...
Ping Li...
New York...

colinwe...

**Understanding Gradient Descent on the Edge of Stability in Deep Learning**

Sanjeev Arora *[1]   Zhiyuan Li *[1]   Abhishek Panigrahi *[1]

The Anisotropic Noise in Stoch...

from Sharp M...

Understanding the Generalization Benefit of Normalization Layers: Sharpness Reduction

Zhanxing Zhu*[1 2 3]   Jingfeng Wu*[1]   Bing Yu[1]   Lei Wu[1]   Jinwen Ma

School of Softwa...

...tsinghua.edu.cn

Michael I. Jordan
Department of EECS
University of California, Berkeley
jordan@cs.berkeley.edu

Beyond the Quadratic Approximation: The Multiscale Structure of Neural Network Loss Landscapes

Chao Ma *[1],  Daniel Kunin[† 2],  Lei Wu[‡ 3], and  Lexing Ying[§ 1]

oss-Class Structure
earning Spectra

PAPYAN@STANFORD.EDU

Stanford University
Stanford, CA 94305, USA

# Partially Understood Phenomena Abound in NN Optimization

An incomplete list:

- Grokking
- Benefits of Large LR
- Batchnorm
- Hessian Spectrum Outliers
- Sharpening/EoS
- Simplicity Bias

# Partially Understood Phenomena Abound in NN Optimization

An incomplete list:

- Grokking
- Benefits of Large LR
- Batchnorm
- Hessian Spectrum Outliers
- Sharpening/EoS
- Simplicity Bias
- Adaptive Methods

**Toward Understanding Why Adam Converges Faster Than SGD for Transformers**

Yan Pan                                                        YPAN2@ANDREW.CMU.EDU
Yuanzhi Li                                                YUANZHIL@ANDREW.CMU.EDU
*Carnegie Mellon University*

**Why are Adaptive Methods Good for Attention Models?**

Jingzhao Zhang                 Sai Praneeth Karimireddy                 Andreas Veit
MIT                                      EPFL                                      Google Research
jzhzhang@mit.edu        sai.karimireddy@epfl.ch        aveit@google.com

Seungyeon Kim                 Sashank Reddi                 Sanjiv Kumar
Google Research                 Google Research                 Google Research
seungyeonk@google.com        sashank@google.com        sanjivk@google.com

Suvrit Sra
MIT
suvrit@mit.edu

**NOISE IS NOT THE MAIN FACTOR BEHIND THE GAP BETWEEN SGD AND ADAM ON TRANSFORMERS, BUT SIGN DESCENT MIGHT BE**

Frederik Kunstner, Jacques Chen, J. Wilder Lavington & Mark Schmidt[†]
University of British Columbia, Canada CIFAR AI Chair (Amii)[†]

# Partially Understood Phenomena Abound in NN Optimization

An incomplete list:

- Grokking
- Benefits of Large LR
- Batchnorm
- Hessian Spectrum Outliers
- Sharpening/EoS
- Simplicity Bias
- Adaptive Methods
- Unstable Training

# Partially Understood Phenomena Abound in NN Optimization

An incomplete list:

- Grokking
- Benefits of Large LR
- Batchnorm
- Hessian Spectrum Outliers
- Sharpening/EoS
- Simplicity Bias
- Adaptive Methods
- Unstable Training
- Double Descent

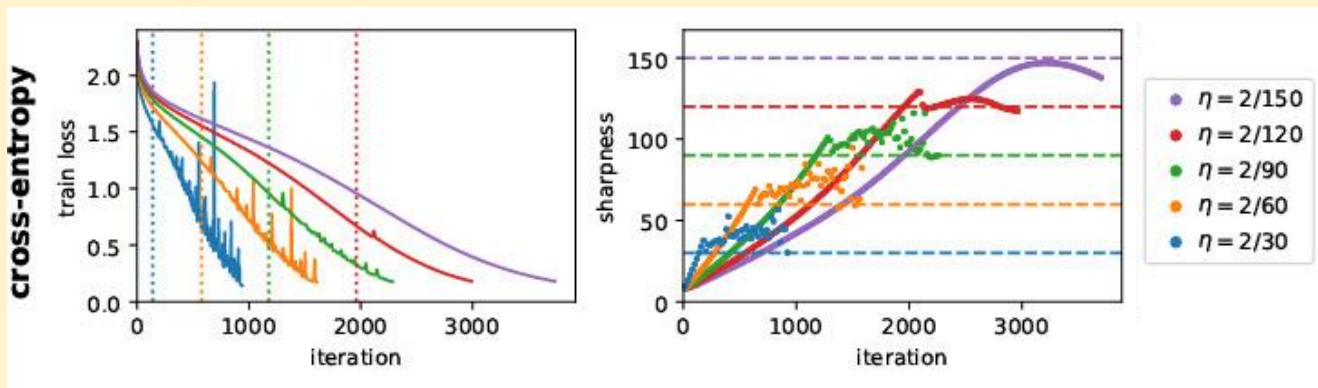Surely, *some* of these results are related… but unclear how.

# Progressive Sharpening + Edge of Stability



Meanwhile, loss decreases non-monotonically, with frequent "spikes".

"Sharpness" = top eigenvalue of loss Hessian
First rises to $2/\eta$...
Then hovers around that value.

[1] Gradient Descent on Neural Networks Typically Occurs at the Edge of Stability. Cohen et al. 2020.

# Progressive Sharpening + Edge of Stability



This is just more evidence that **something more is needed** to understand NN training dynamics…

[1] Gradient Descent on Neural Networks Typically Occurs at the Edge of Stability. Cohen et al. 2020.

# Yet Another Phenomenon

I'm going to present our finding:

*another* interesting phenomenon in neural network optimization.

**But the goal is not just to add to the growing list.**

Instead, we hope it can help explain and unify

these observations via a **shared underlying cause**.*

# Yet Another Phenomenon

(We also look at SGD)

Let's run the following experiment:

1. Train a neural network with full-batch gradient descent on CIFAR-10.
2. Track losses on each training point *individually*.
3. Fix some iteration $T$.
4. Calculate changes in loss on each point from step $T$ to step $T+1$.
5. Visualize the samples with the *most positive* and *most negative* changes.

What should we expect to see?

# Yet Another Phenomenon



VGG–11

ResNet-18

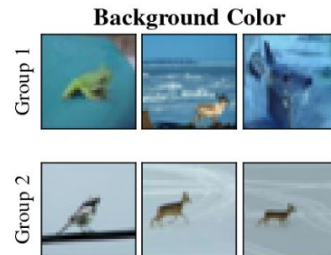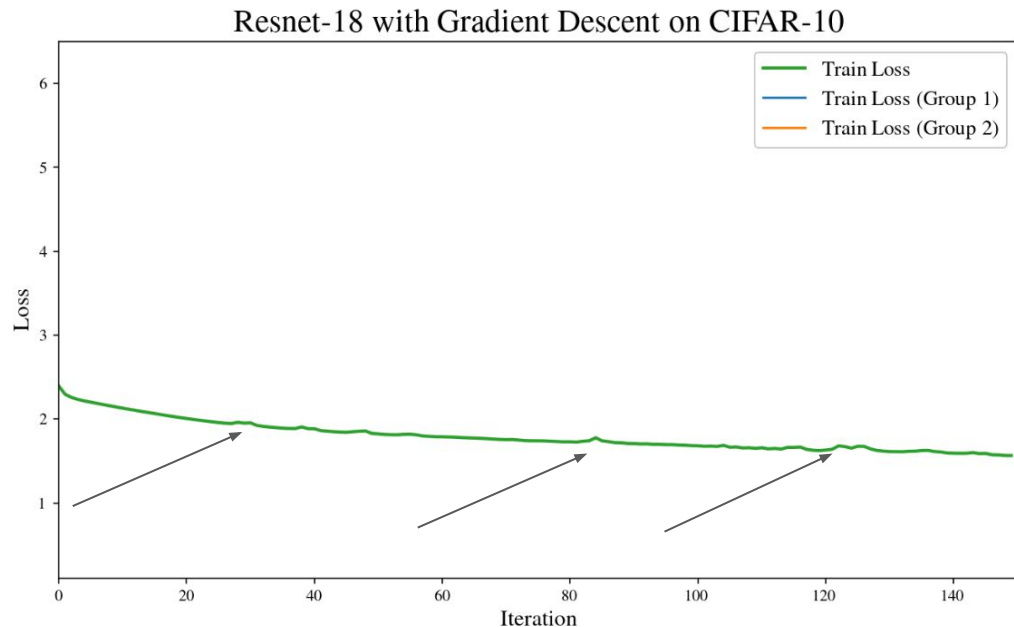**Largest increase in loss**

**Largest decrease in loss**

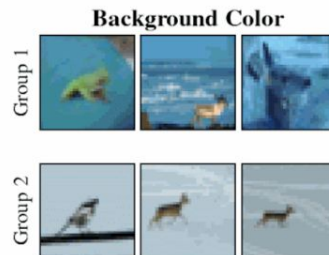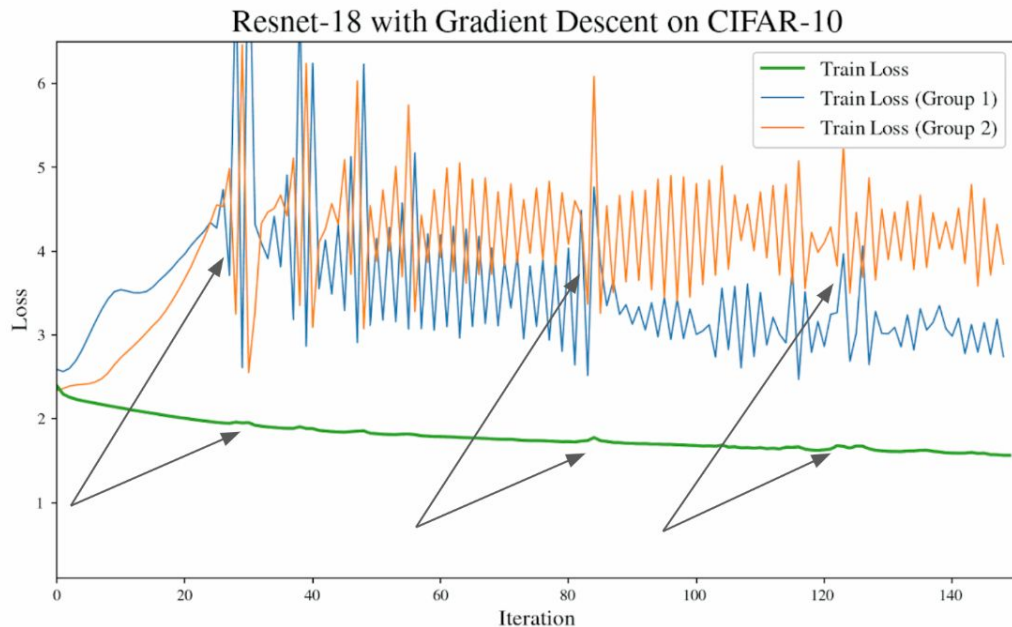The precise patterns change, but this occurs *all throughout training*.

# Visualizing the Group Losses
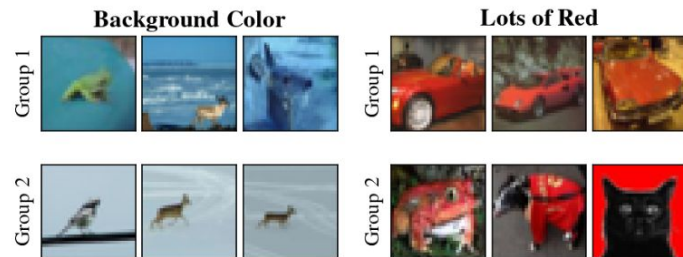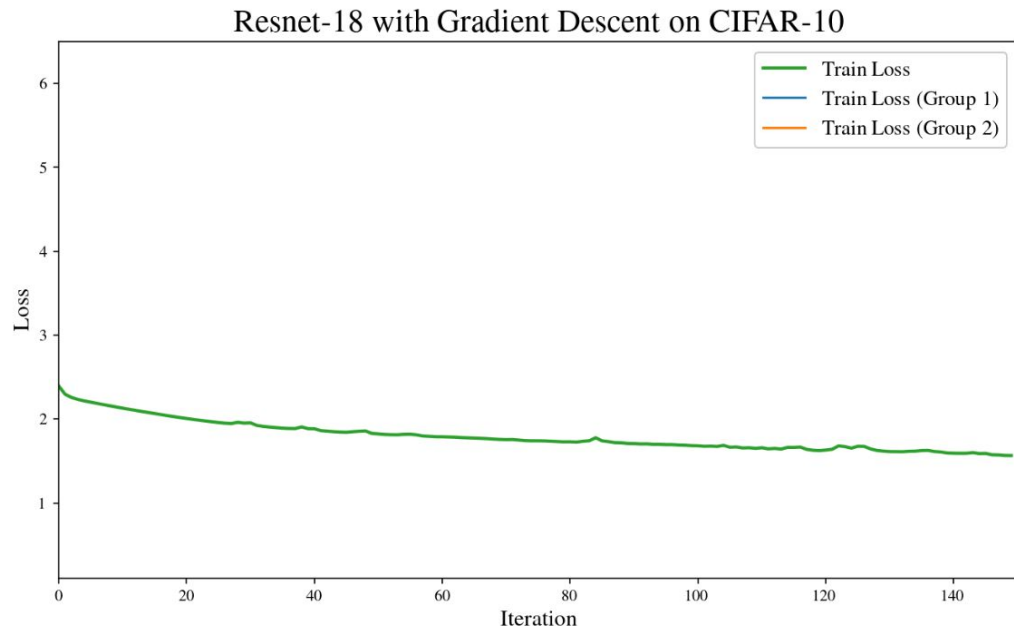
These groups are ~20 samples each.



Samples were selected for largest change in loss, so we expect a "spike" somewhere.
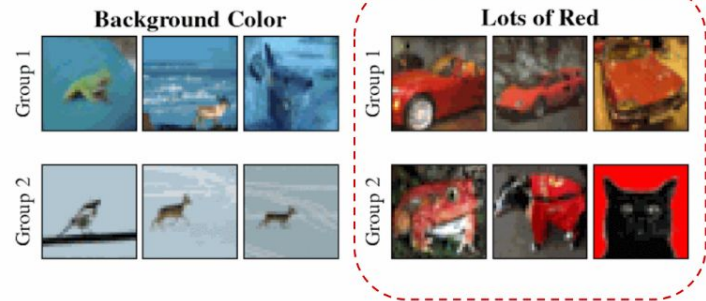
# Visualizing the Group Losses



These opposing groups oscillate with large amplitude *continuously*!

# Visualizing the Group Losses



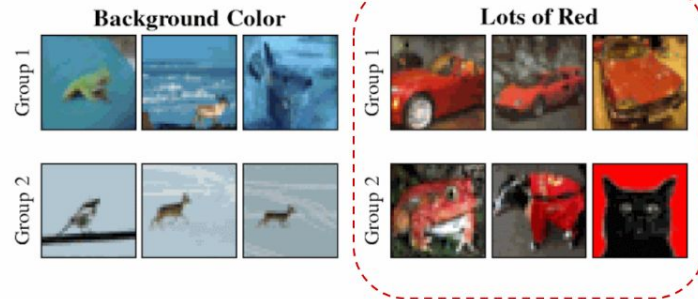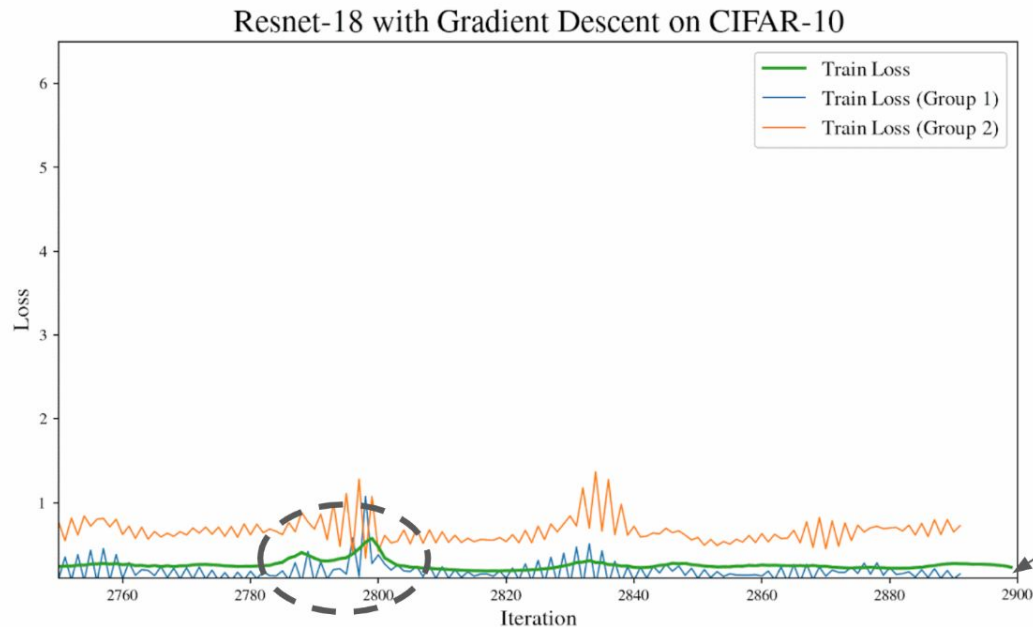Resnet-18 with Gradient Descent on CIFAR-10

What about another group?

# Visualizing the Group Losses



When we're close to interpolating, shouldn't this effect be reduced?

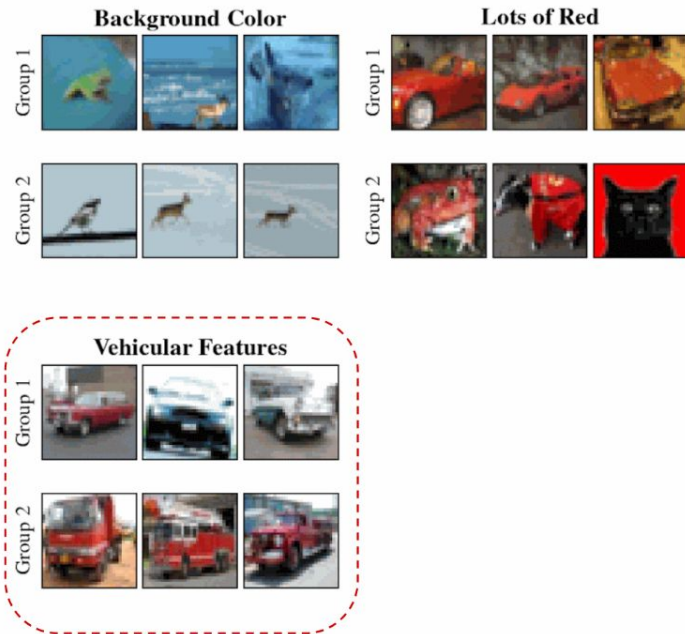# Visualizing the Group Losses



Resnet-18 with Gradient Descent on CIFAR-10
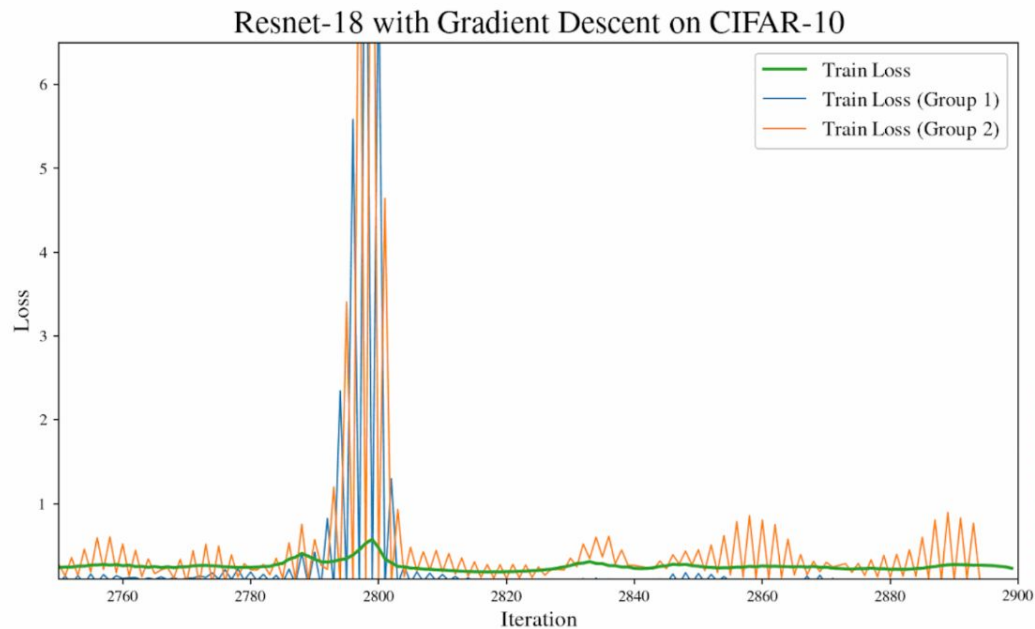
(Note the change in $x$-axis)

Yes, the amplitude is substantially smaller…

So what's causing these loss increases?

# Visualizing the Group Losses



Resnet-18 with Gradient Descent on CIFAR-10

- Train Loss
- Train Loss (Group 1)
- Train Loss (Group 2)

**Background Color**

Group 1

Group 2

**Lots of Red**

Group 1

Group 2

**Vehicular Features**

Group 1

Group 2

# Visualizing the Group Losses



Even at the end stages of training, large loss swings are still occurring.

# What's Going On?

- **Prevalent features**, often with distinct colors.
  - Roughly, "prevalent" ≈ "fills a lot of the image"

- **Begin simple**, become progressively more complex.
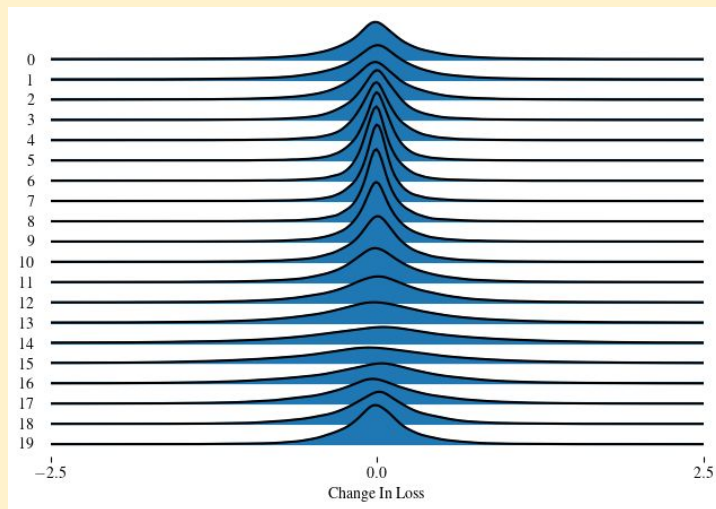
  - "Simple" ≈ "available at random initialization"

- Large gradients pointing in **opposite directions**.
  - Learning "red = car" decreases loss on red cars, increases loss on red *non*-cars

We call these features—or the gradients they induce—*Opposing Signals*.

# What's Going On?

Does this occur for *every* training sample?
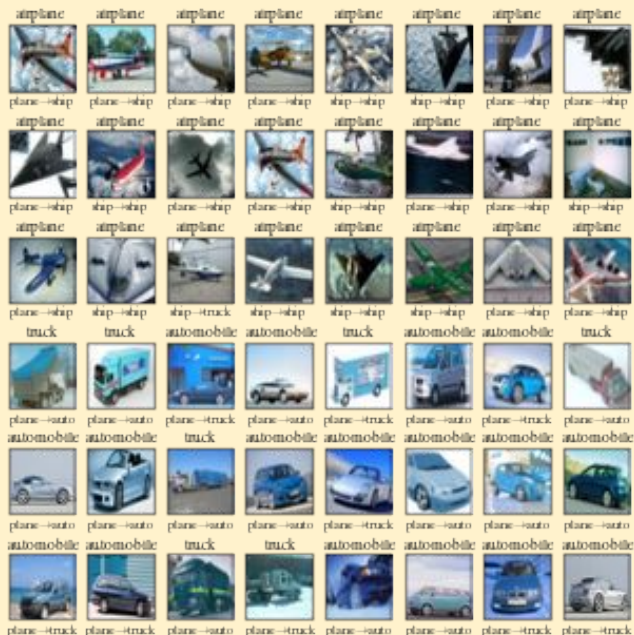
Distribution of changes in loss:



These samples are *significant* outliers.

# What Causes Opposing Signals?

Is this a property of architecture (ConvNet)?

**No.** Same occurs in a Vision Transformer.

# What Causes Opposing Signals?

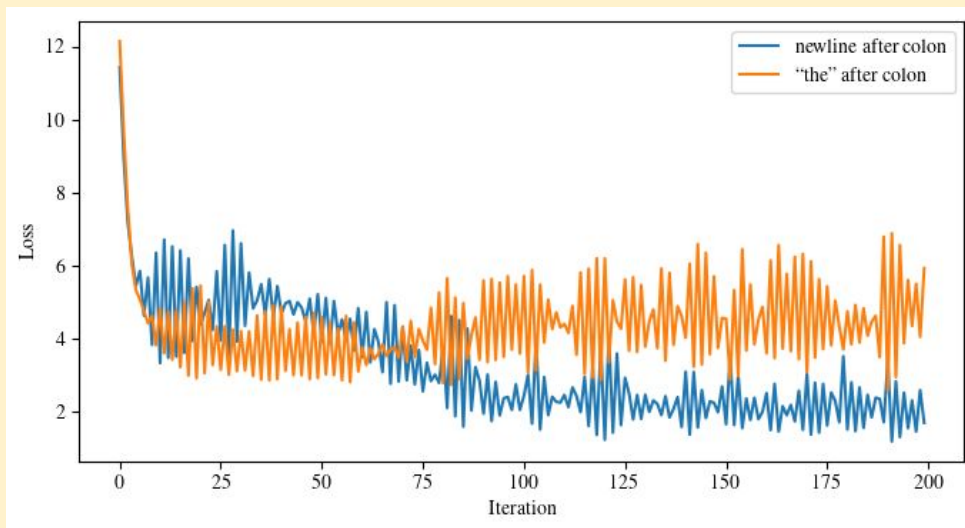Maybe it's a property of the data modality (images)?

**Also no.**

## Group 1

```
Salcedo said of the work:[\n]
Enter your email address:[\n]
According to the CBO update:[\n]
Here's how the Giants can still make the
playoffs:[\n]
in early 2018.\n\nAccording to the CBO update:[\n]
other than me being myself."\n\nWATCH:[\n]
```

## Group 2

```
MPs in Westminster. But to me it is obvious: [the]
The wheelset is the same as that on the model above:
[the]
all other acts of love, both divine and human: [the]
from the Kurds' two main political parties: [the]
title of precisely what makes it so wonderful: [the]
you no doubt noticed something was missing: [the]
```
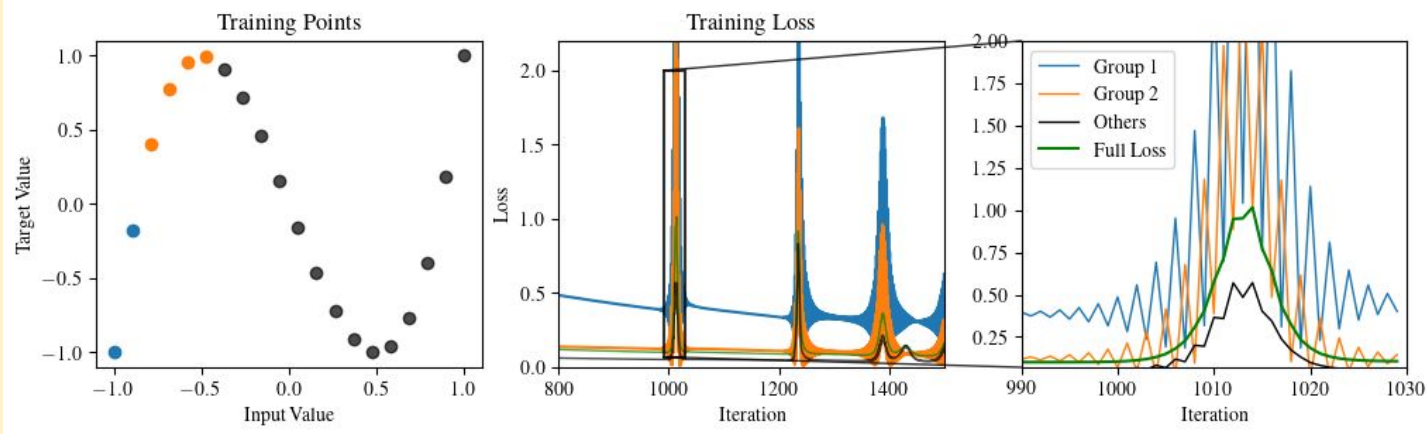
(bracket is next token)



GPT-2 on OpenWebText

# What Causes Opposing Signals?

Maybe it's a property of the data modality (images)?

**Also no.**

What about the loss (cross-entropy)?

# What Causes Opposing Signals?

Remainder of this talk gives our current best understanding, with experiments.
    We believe it a consequence of *depth* and *steepest descent.*

**We don't fully understand the mechanism here.**

- If there are parts you think aren't fully explained, you're right.

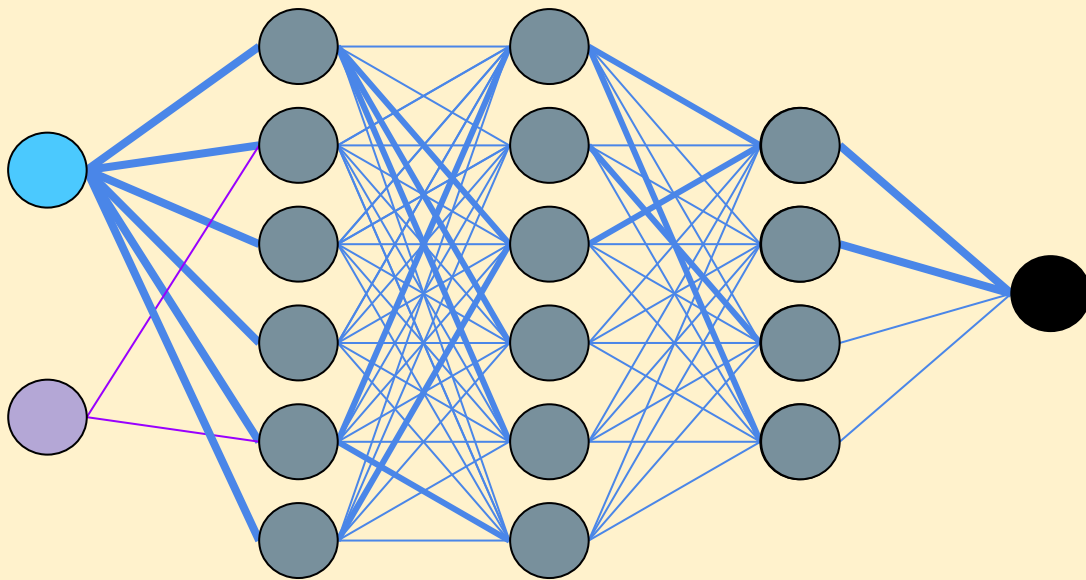- If there are parts you think are *flat out wrong*, you could be right.

However:
- We have a reasonably descriptive high-level story...
    - and we prove this behavior for a simple model on a 2-layer linear net.[*]

- It enables *specific* qualitative predictions which we then verify...
    - and it naturally fits into several existing narratives of other phenomena.

# A Simplified Story of Gradient Descent on Deep Neural Networks

Consider a randomly initialized MLP with two input features:

1.  "Sky": large magnitude + pervasive (propagated to all neurons).
    -   Only sufficient for predicting $p($ class | "sky" $)$.

2.  "Shape": small magnitude, needs to be learned.
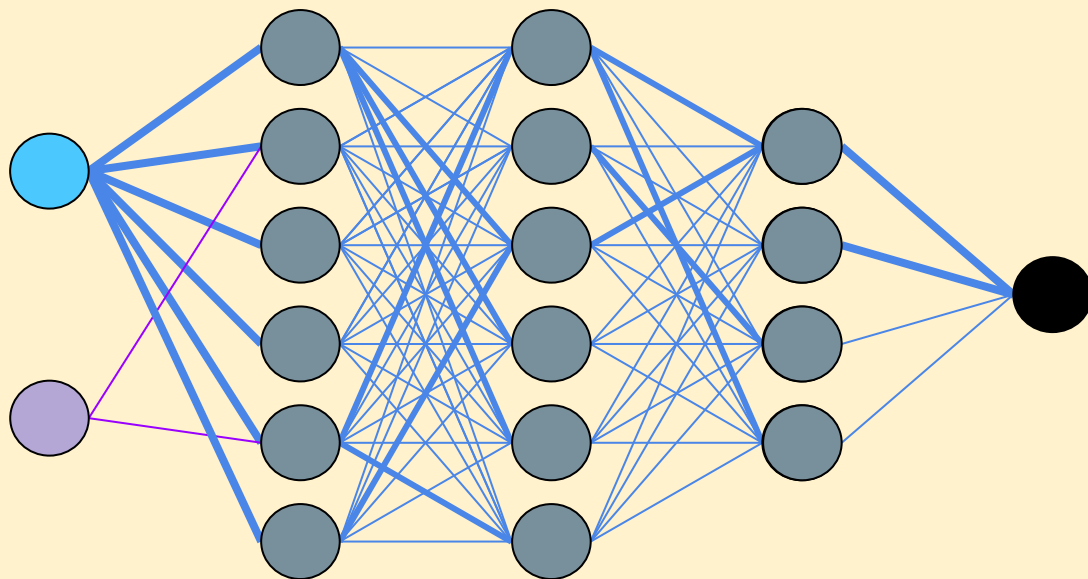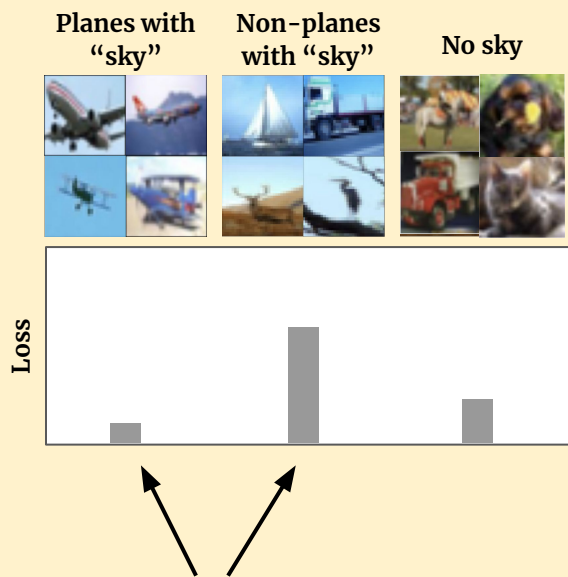    -   But much more useful for loss reduction.

At initialization, network activations are dominated by "sky" on outliers.

- (Suppose network happens to predict "sky = plane")

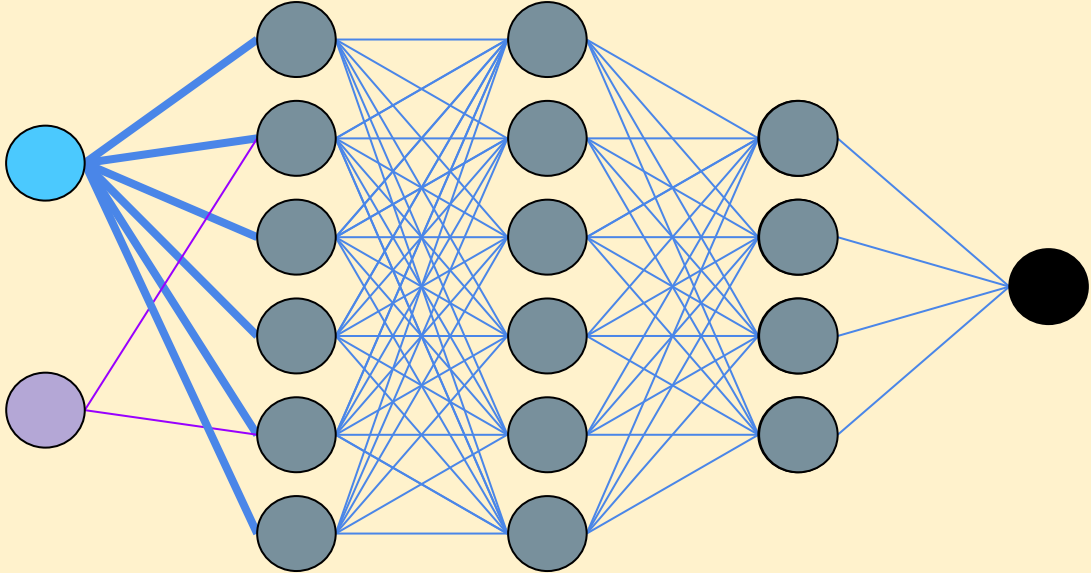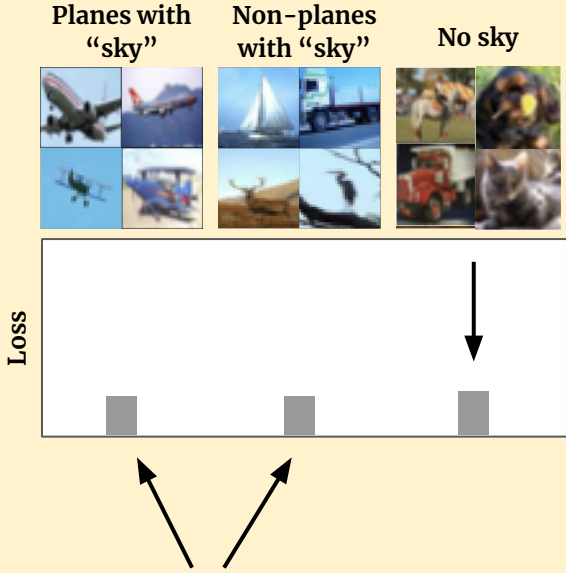High loss → large gradients → rebalance towards predicting $p($ class $|$ "sky" $)$.

- (This "linear first" behavior has been previously observed[1, 2])
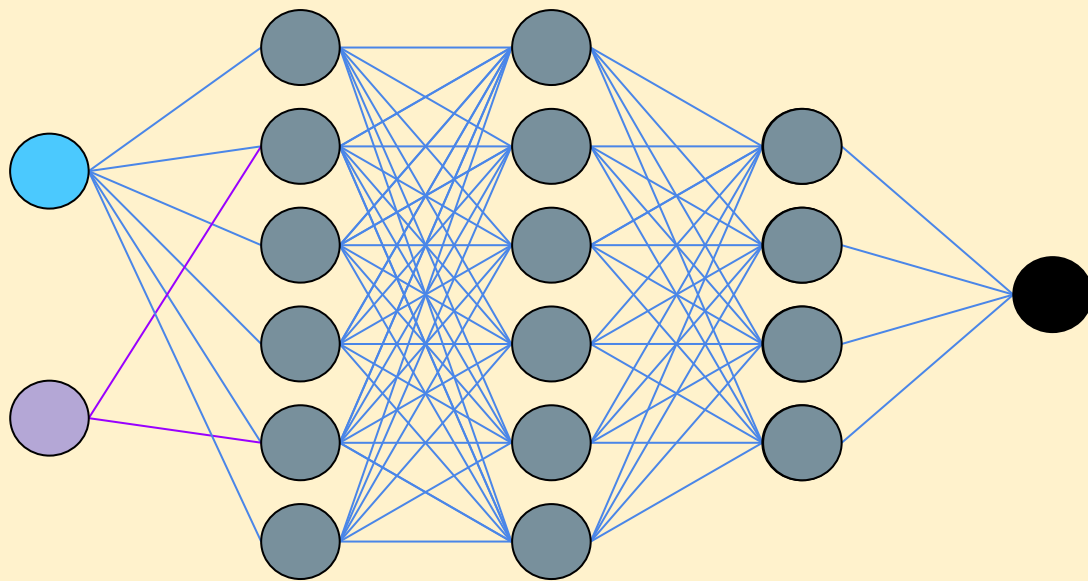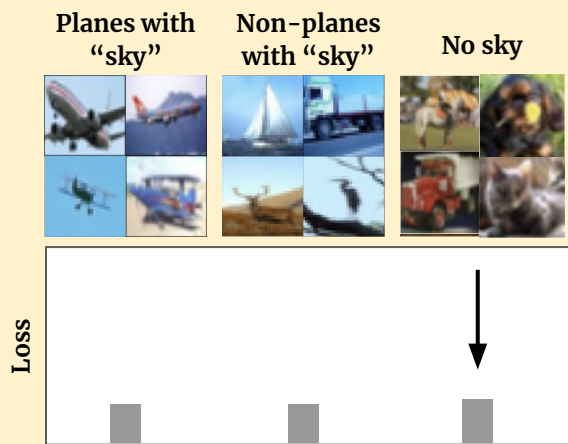
[1] SGD on Neural Networks Learns Functions of Increasing Complexity. Nakkiran et al. 2019.
[2] Do deep neural networks learn shallow learnable examples first? Mangalam and Prabhu 2019.

Once this happens, the network can now upweight the more useful "shape" feature.

Since the outliers' loss no longer dominates the gradient, let's visualize a non-outlier.

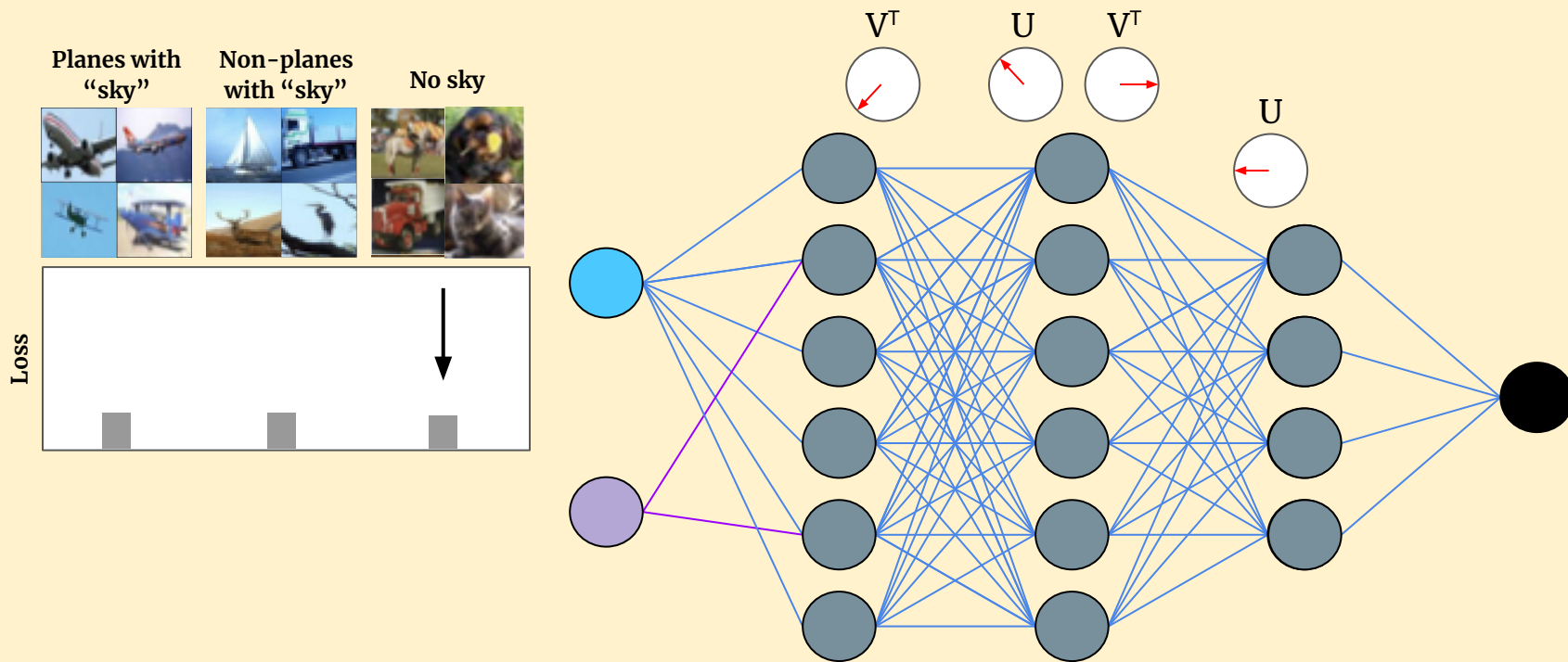Planes with "sky"

Non-planes with "sky"

No sky

Loss

Once this happens, the network can now upweight the more useful "shape" feature.

Since the outliers' loss no longer dominates the gradient, let's visualize a non-outlier.

Planes with "sky"

Non-planes with "sky"

No sky

Loss

As training progresses, the top singular vectors of adjacent layers align to amplify meaningful subspaces. [3, 4]

This is how the "shape" feature gets *upweighted.*



Planes with "sky"   Non-planes with "sky"   No sky

Loss

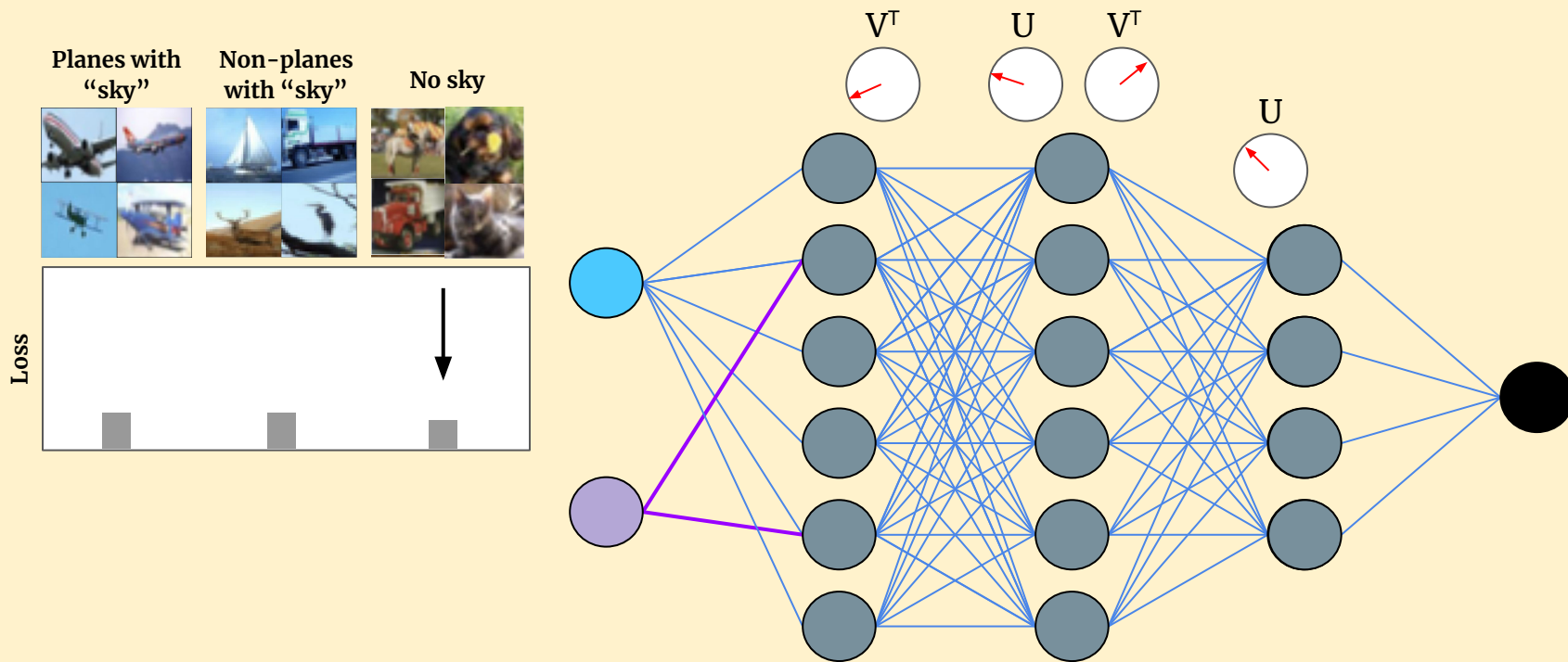V$^\mathsf{T}$   U   V$^\mathsf{T}$   U

[3] Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. Saxe et al. 2013
[4] Unique properties of flat minima in deep networks. Muyaloff and Michaeli, 2020.

As training progresses, the top singular vectors of adjacent layers align to amplify meaningful subspaces. [3, 4]
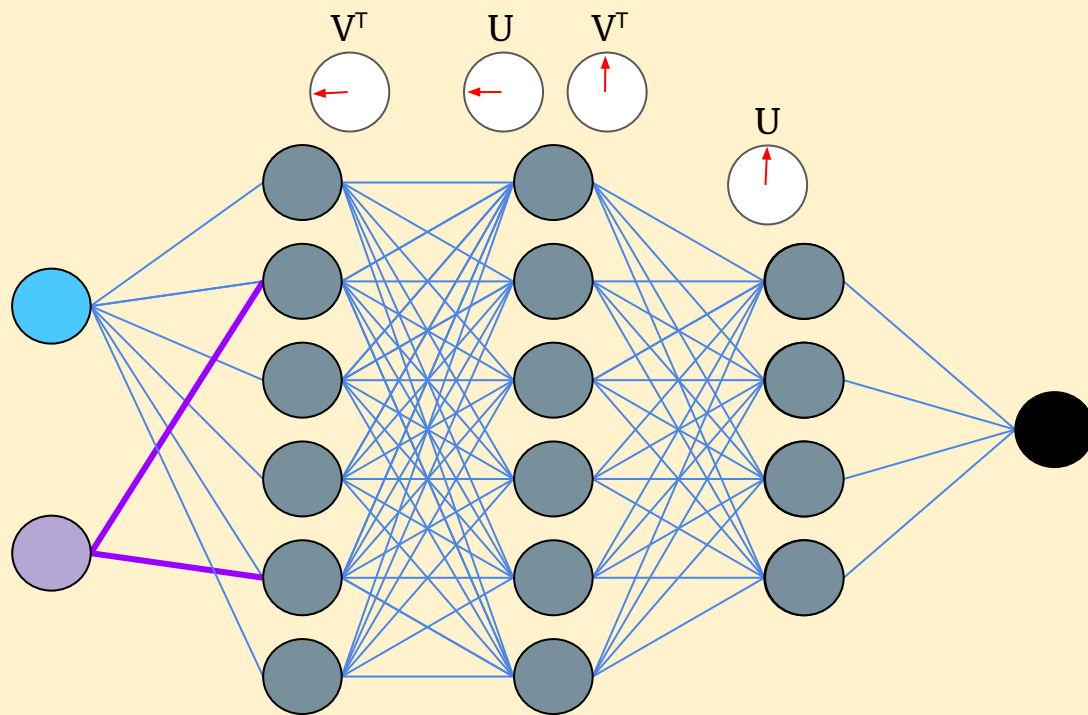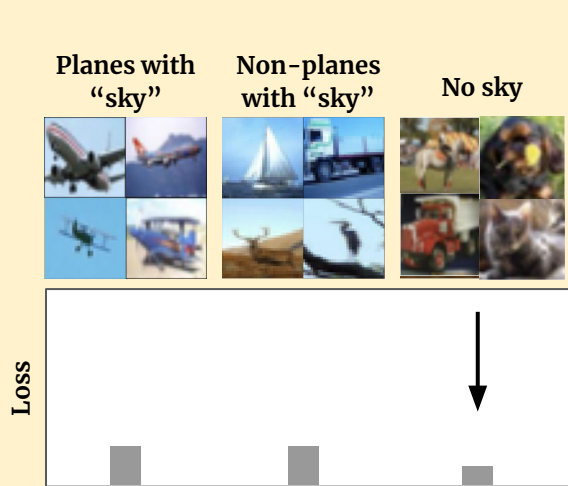
This is how the "shape" feature gets *upweighted.*



Planes with "sky"    Non-planes with "sky"    No sky

Loss

[3] Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. Saxe et al. 2013
[4] Unique properties of flat minima in deep networks. Muyaloff and Michaeli, 2020.

# This alignment has been continuously upweighting the more useful signal.



Planes with "sky"

Non-planes with "sky"

No sky

Loss

$V^{\top}$   $U$   $V^{\top}$

$U$

[3] Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. Saxe et al. 2013
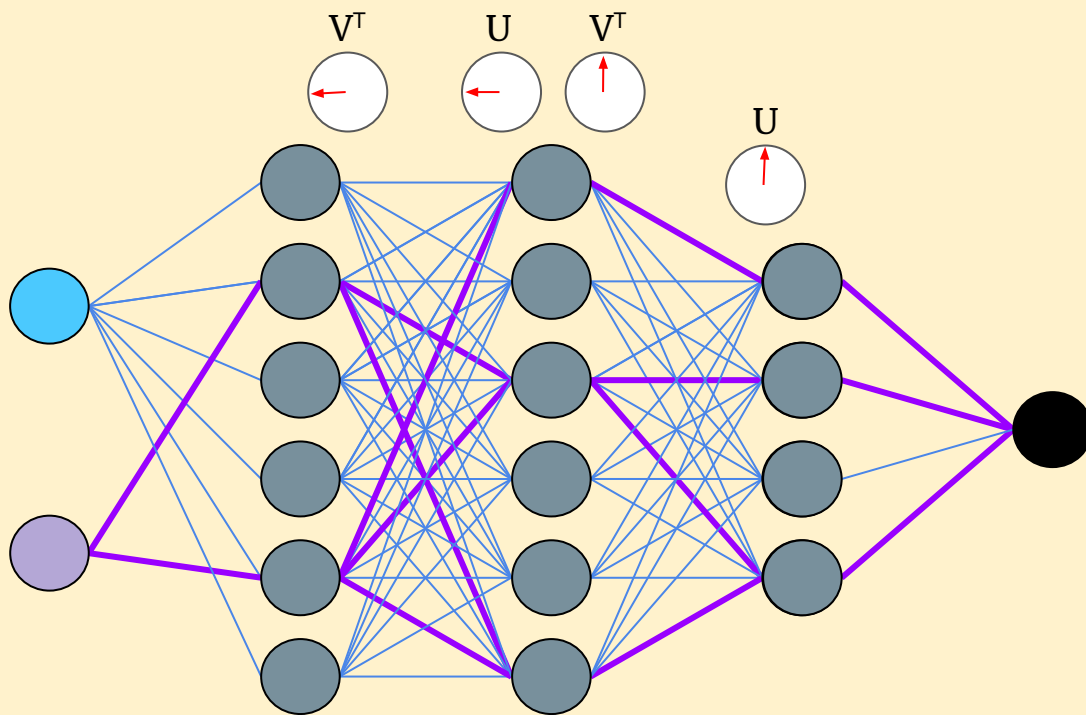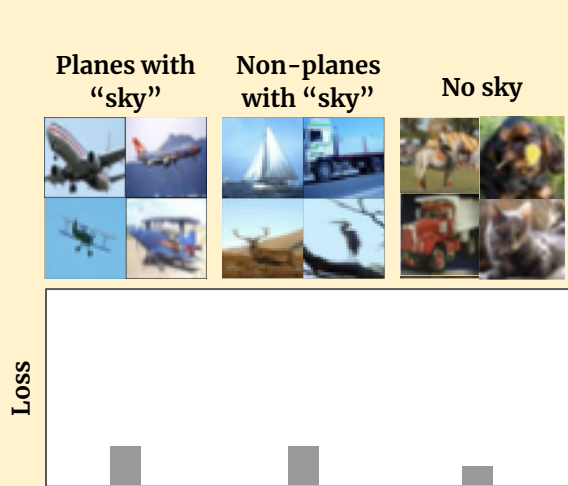[4] Unique properties of flat minima in deep networks. Muyaloff and Michaeli, 2020.

I've left one important part out of this visualization:

When "shape" is amplified, **"sky" is amplified too**.
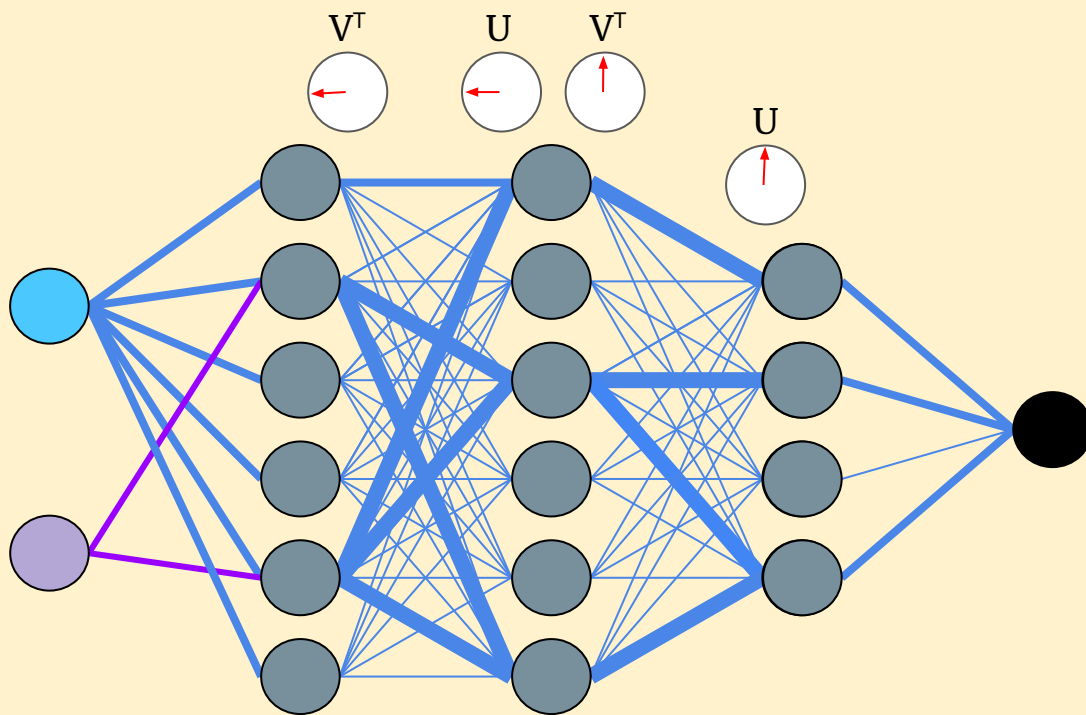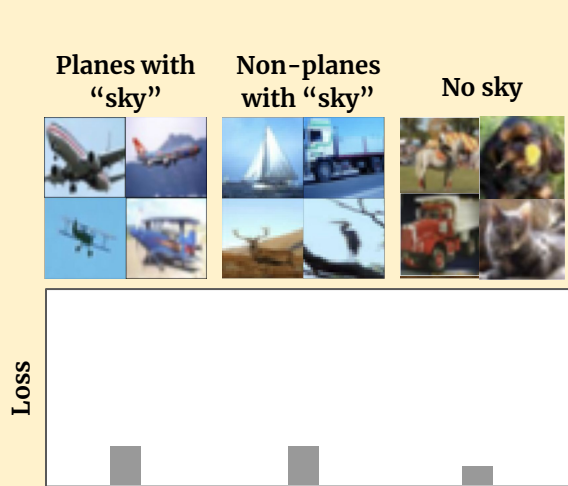
This is the activation pattern for a *non-outlier*.

**What would it look like for an outlier with a sky background?**

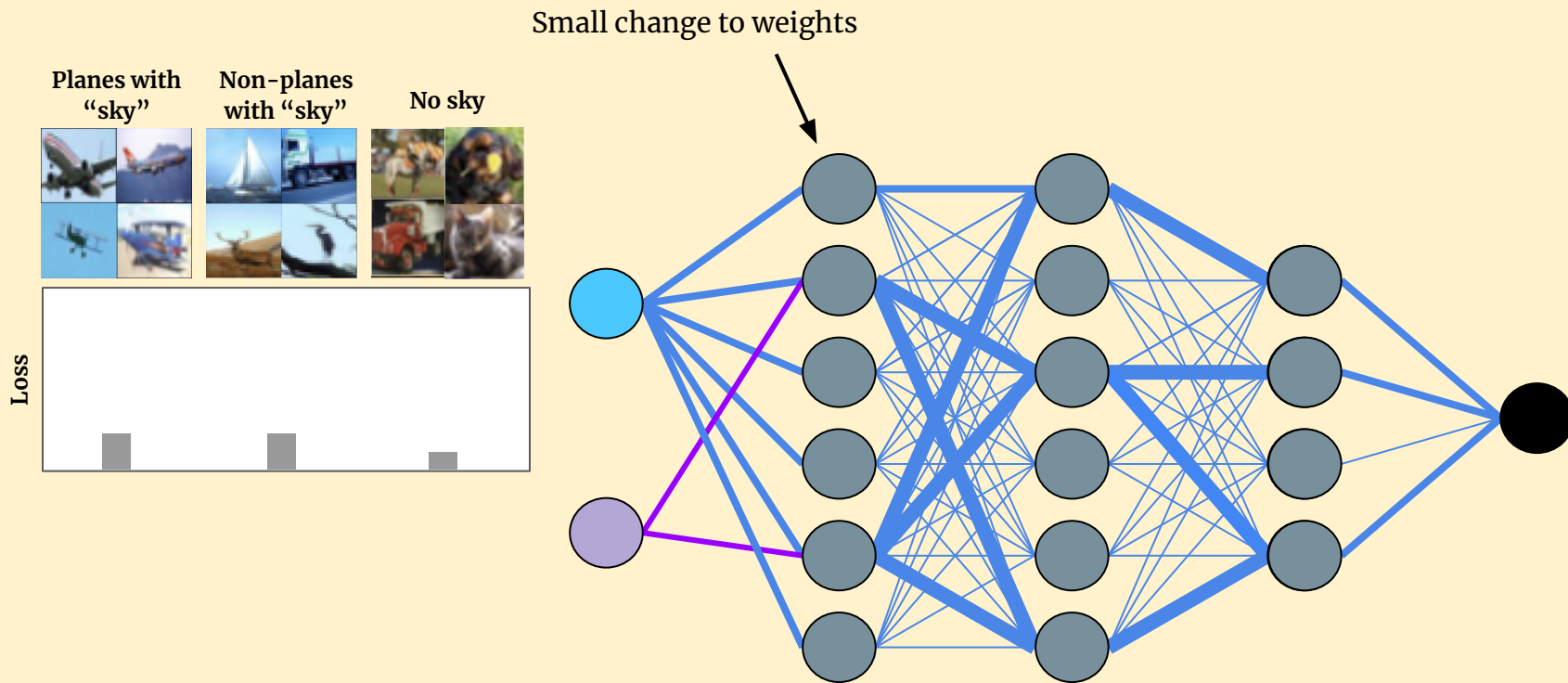I've left one important part out of this visualization:

When "shape" is amplified, **"sky" is amplified too**.

Because it is larger + more pervasive, it still dominates the network's activations.

# This causes large sensitivity to small changes in *how the network uses* "sky".

- Small, targeted change to predict one group massively increases loss on the other.



Small change to weights

Planes with "sky"

Non-planes with "sky"

No sky

Loss

# This causes large sensitivity to small changes in *how the network uses* "sky".

- Small, targeted change to predict one group massively increases loss on the other.

Changes how feature propagates...

**Planes with "sky"**

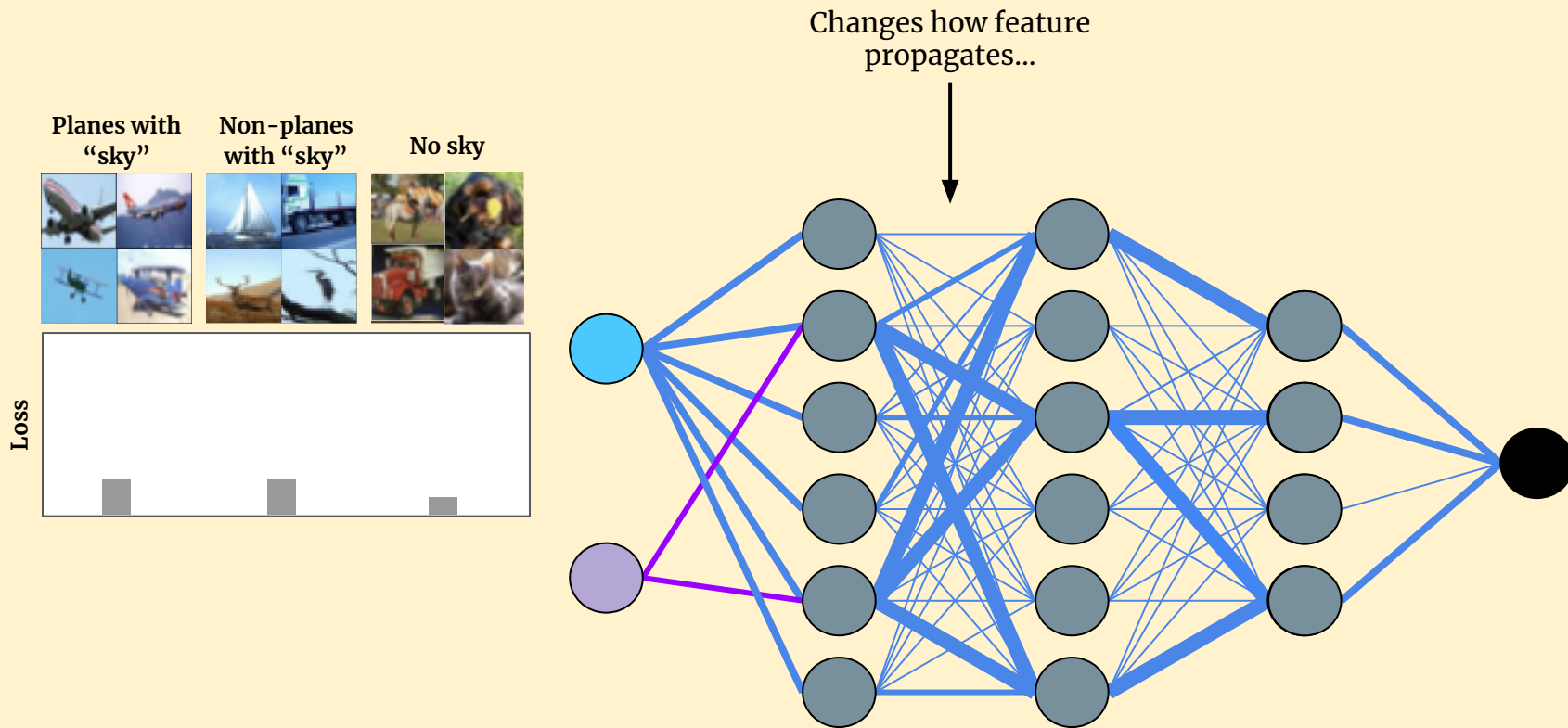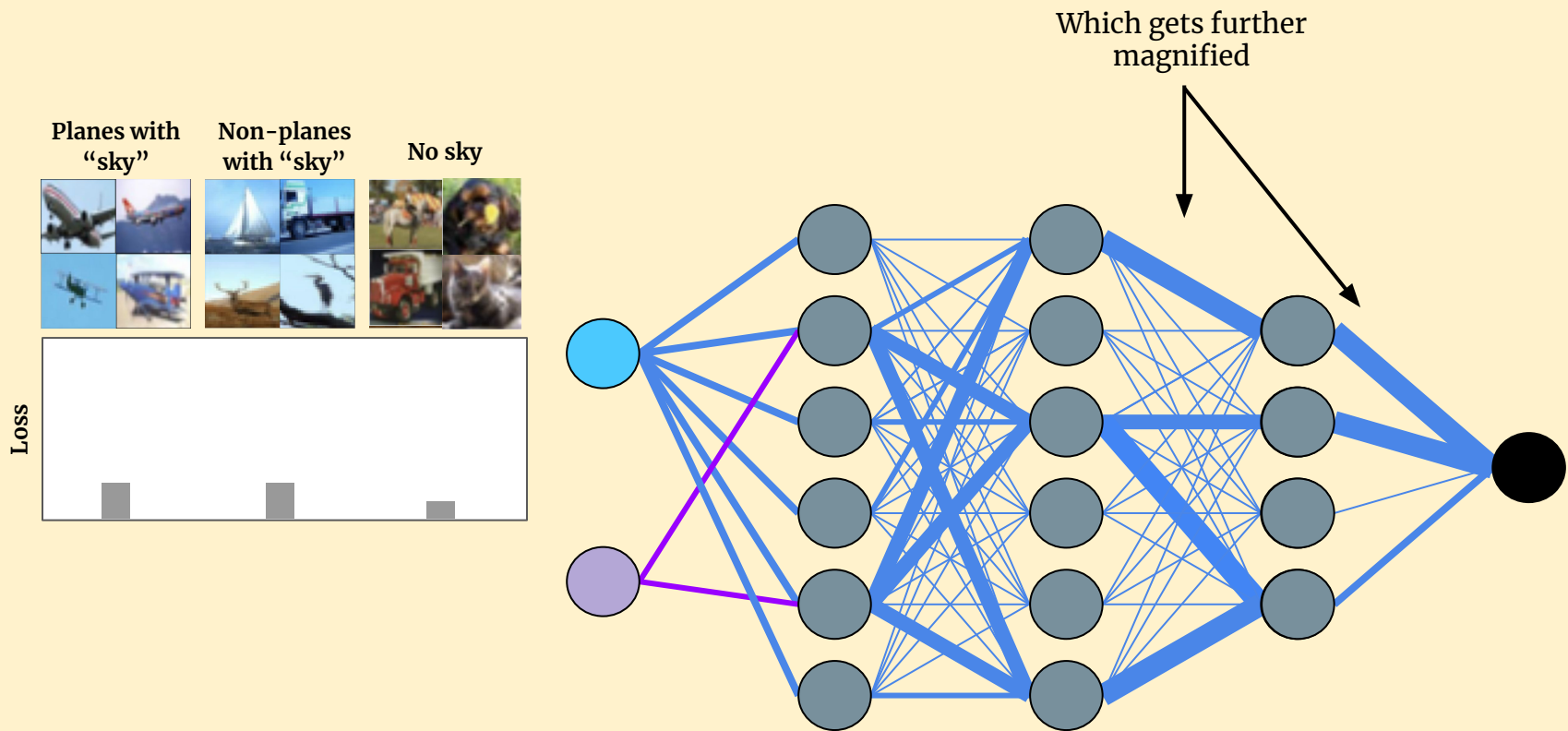**Non-planes with "sky"**

**No sky**

Loss

This causes large sensitivity to small changes in *how the network uses* "sky".
- Small, targeted change to predict one group massively increases loss on the other.

Which gets further magnified

**Planes with "sky"**

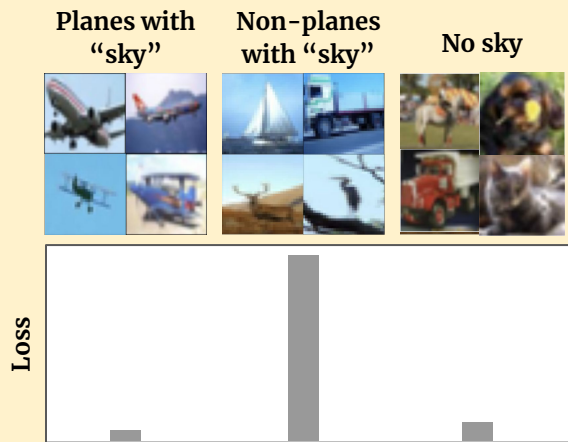**Non-planes with "sky"**

**No sky**

Loss

# This causes large sensitivity to small changes in *how the network uses* "sky".

- Small, targeted change to predict one group massively increases loss on the other.



Planes with "sky"
Non-planes with "sky"
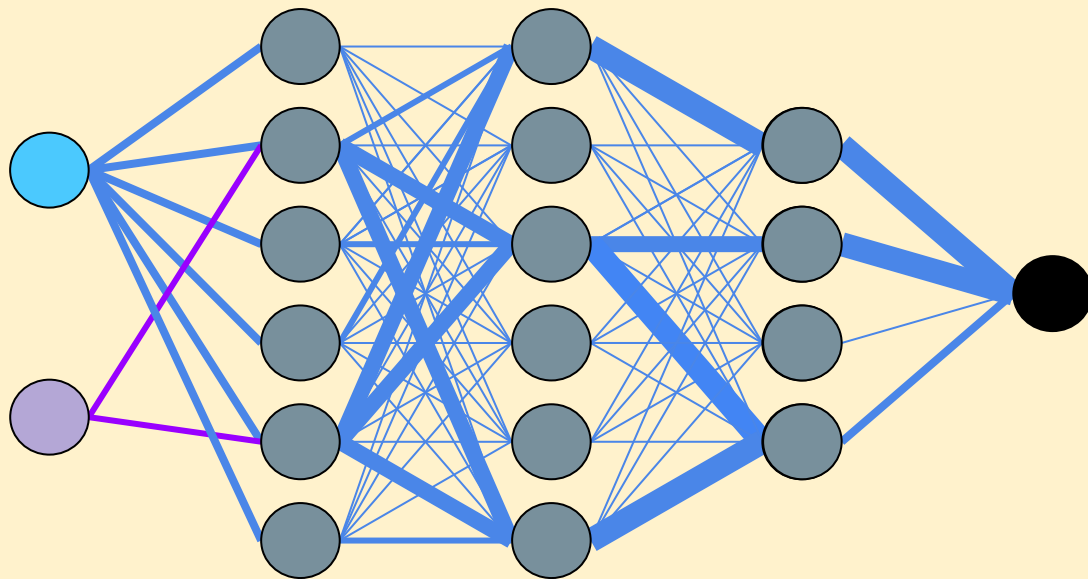No sky

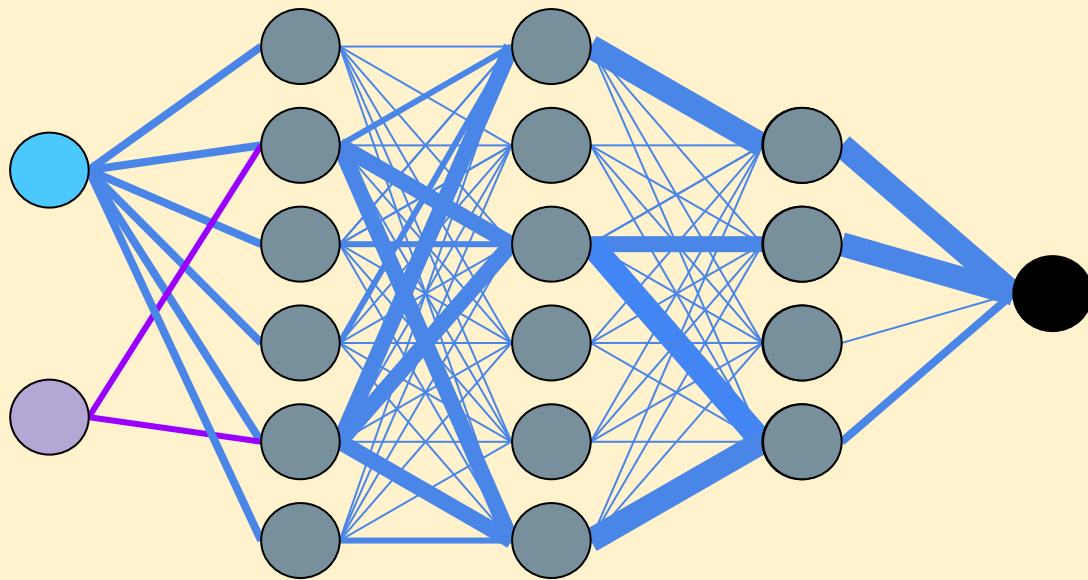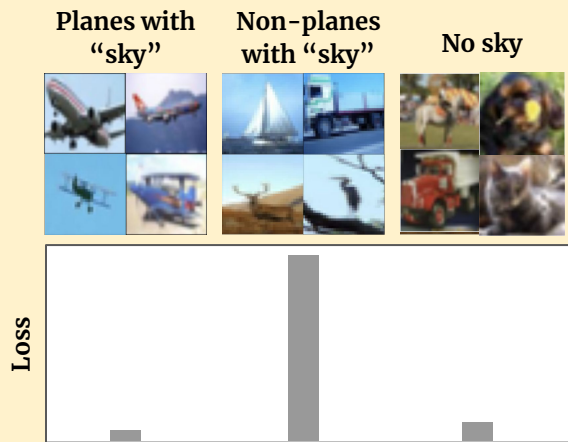Loss

And can have huge effect on loss

This causes large sensitivity to small changes in *how the network uses* "sky".

- Small, targeted change to predict one group massively increases loss on the other.

In other words, **loss on outliers becomes very *sharp* w.r.t. parameters.**

- ("growth in sensitivity" was previously noted, e.g. weight/Jacobian norm[5, 6])



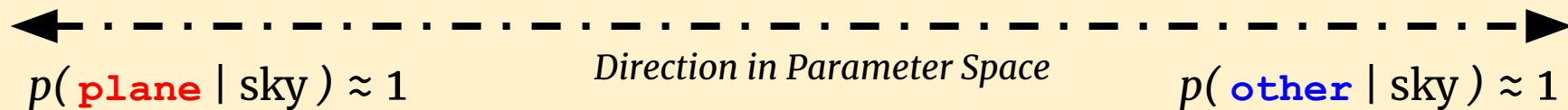**Planes with "sky"**  **Non-planes with "sky"**  **No sky**

Loss

[5] On linear stability of sgd and input-smoothness of neural networks. Ma and Ying, 2021.
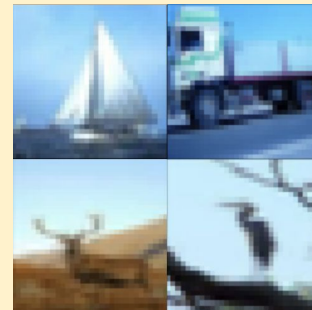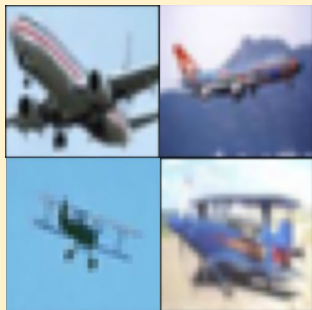[6] On the lipschitz constant of deep networks and double descent. Gamba et al. 2023.

This story is pretty abstract.

Let's visualize something more concrete:

The (hypothetical) loss in a 1D parameter space.

$\longleftarrow$ ·—·—·—·—·—·—·—·—·—·—·—·—·—· $\longrightarrow$

$p(\text{ plane } | \text{ sky }) \approx 1$         *Direction in Parameter Space*         $p(\text{ other } | \text{ sky }) \approx 1$

Optimization continues "through the valley"[1]

How does early optimization move along this axis?

[1] A Walk with SGD. Xing et al. 2018.

Loss on images of **plane** with sky

Loss on images of **non-planes** with sky

$p(\ \text{plane}\ |\ \text{sky}\ ) \approx 1$

*Direction in Parameter Space*

$p(\ \text{other}\ |\ \text{sky}\ ) \approx 1$

What happens when norm of "sky" grows?

Sensitivity to *how we use the sky feature* grows.

Hence, the loss **sharpens** along this direction.

*Direction in Parameter Space*

$p(\text{plane} \mid \text{sky}) \approx 1$

$p(\text{other} \mid \text{sky}) \approx 1$

Eventually, sharpness crosses step size threshold, and iterates begin to diverge!

Loss on images with little/no sky

Let's also visualize the loss on the non-outliers.

$p(\texttt{plane} \mid \text{sky}) \approx 1$

*Direction in Parameter Space*
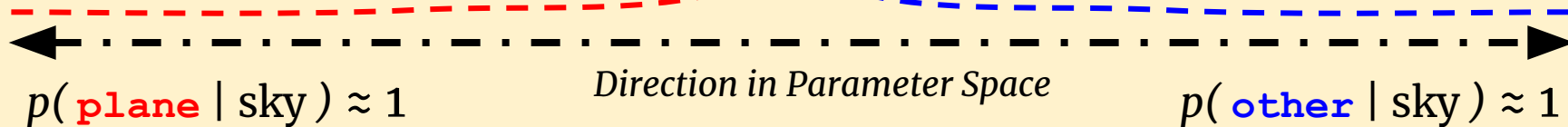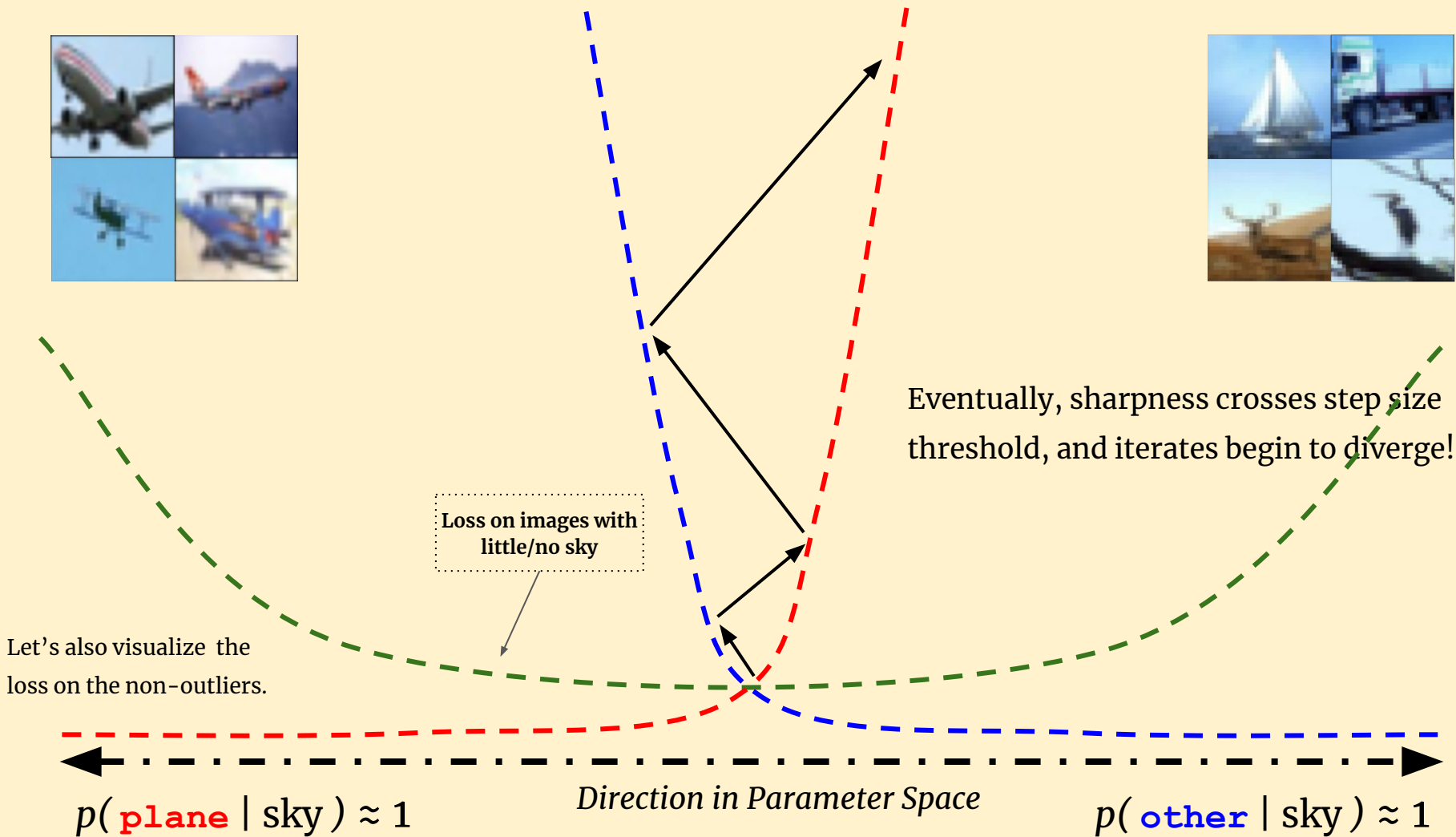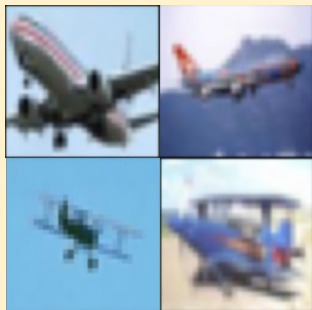
$p(\texttt{other} \mid \text{sky}) \approx 1$

What does the loss look like on each group?

"catapult" / "slingshot"

How far does this continue?

Why should it go back down?

*Direction in Parameter Space*

$p(\,\texttt{plane}\,|\,\text{sky}\,) \approx 1$

$p(\,\texttt{other}\,|\,\text{sky}\,) \approx 1$

Here, losses are balanced.

So are opposing *gradients*.

Feature growth continues.
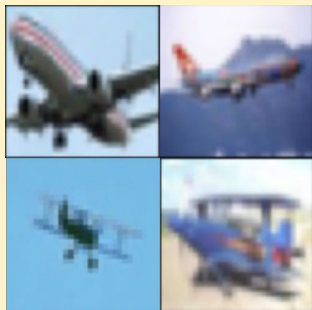
$p(\,\texttt{plane} \mid \text{sky}\,) \approx 1$

*Direction in Parameter Space*

$p(\,\texttt{other} \mid \text{sky}\,) \approx 1$

Here, losses are imbalanced.

But outliers still have small
influence on overall gradient.

*Direction in Parameter Space*

$p(\texttt{plane} \mid \text{sky}) \approx 1$

$p(\texttt{other} \mid \text{sky}) \approx 1$

Here, gradient on `plane` dominates.

Two ways to decrease loss:

1. Use feature differently.
2. *Downweight* feature.

Valley flattens, we descend again.

$p(\,\texttt{plane}\mid\text{sky}\,) \approx 1$

*Direction in Parameter Space*

$p(\,\texttt{other}\mid\text{sky}\,) \approx 1$

# Experimental Verification

The value of a theory (even a non-rigorous one) is in its ability to make predictions.
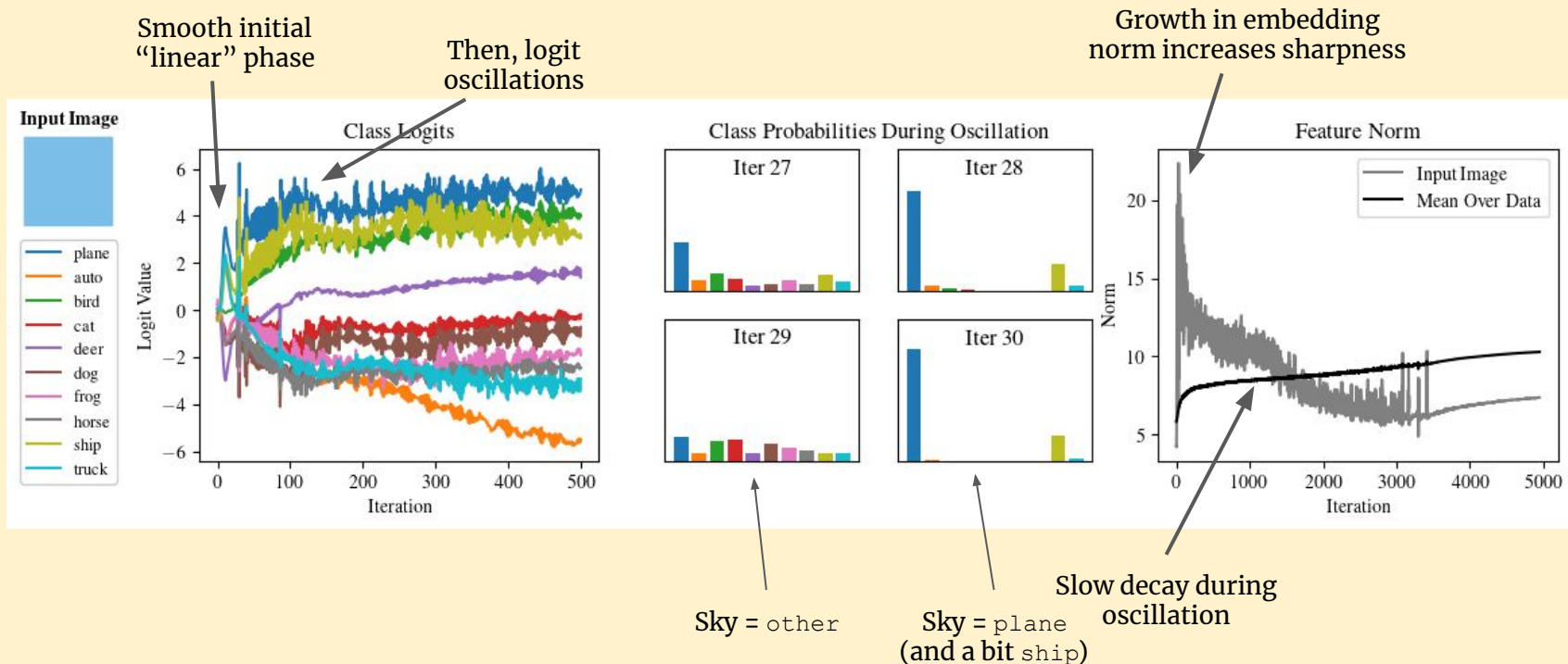
So far we've described:

1. Initial phase of fitting a "linear" model.  ← (previously observed)
2. Growth in activation magnitude among images with this feature.  ← (least well understood)
3. Upon reaching Edge of Stability, predictions *oscillate* between "sky = `plane`" and "sky = `other`".
4. Oscillation results in shrinking of activation magnitude.

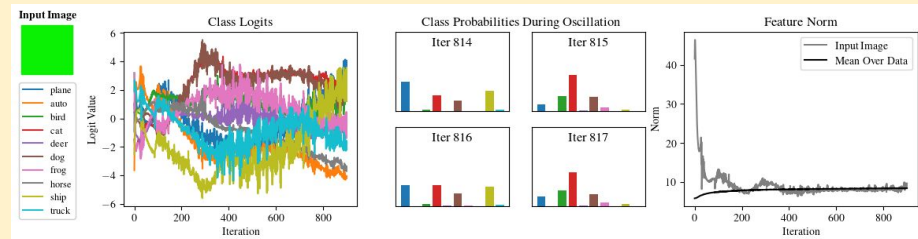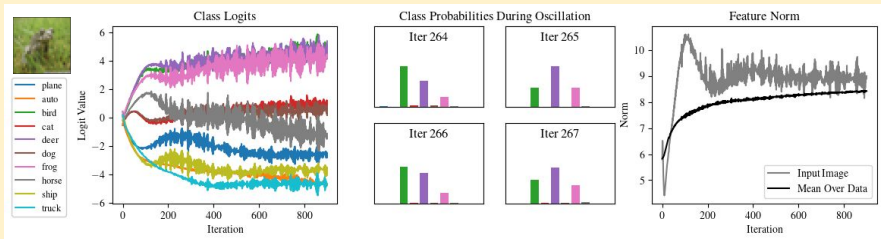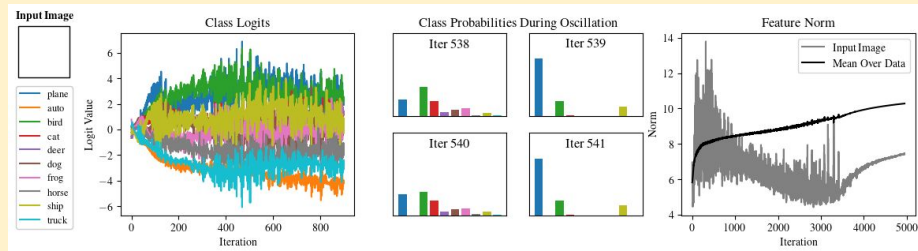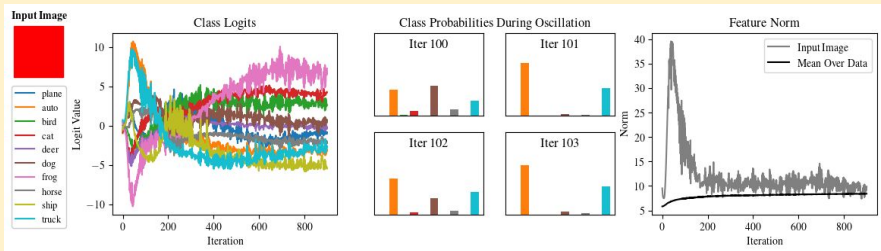What does this story imply, *behaviorally?* Can we test it more directly?

# Experimental Verification

To avoid confounders, we'll pass a pure "sky" image through a ResNet-18.



Smooth initial "linear" phase

Then, logit oscillations

Growth in embedding norm increases sharpness

Slow decay during oscillation

Sky = `other`

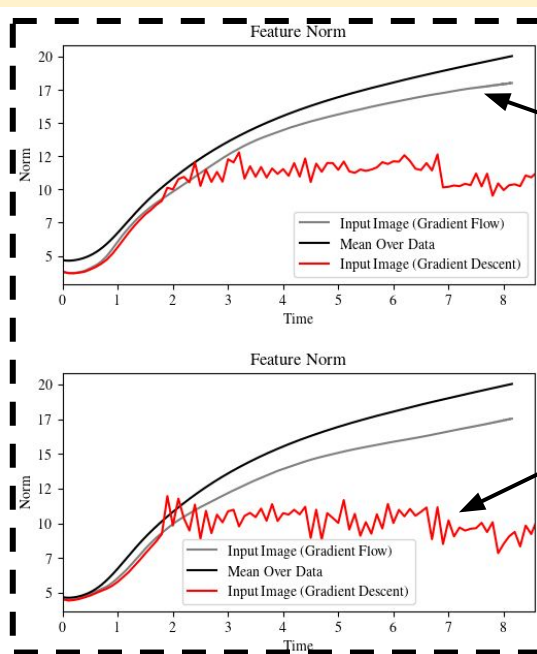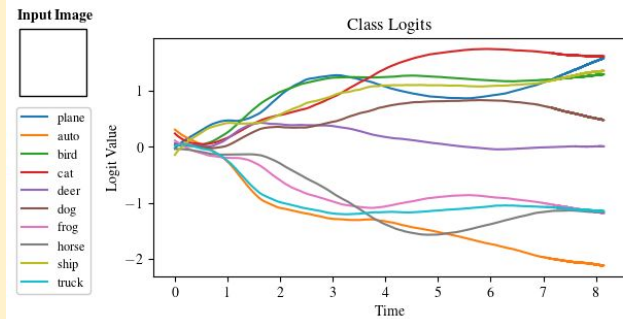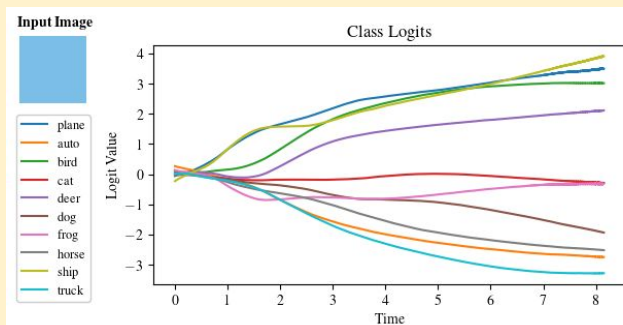Sky = `plane` (and a bit `ship`)

# Experimental Verification

(Doesn't happen as cleanly for all archs/colors, but it's pretty consistent.)

# Experimental Verification

Oscillation seems valuable for downweighting the "simple" but "incomplete" features.

- *Gradient Flow doesn't oscillate.* Maybe that's part of why it generalizes poorly?



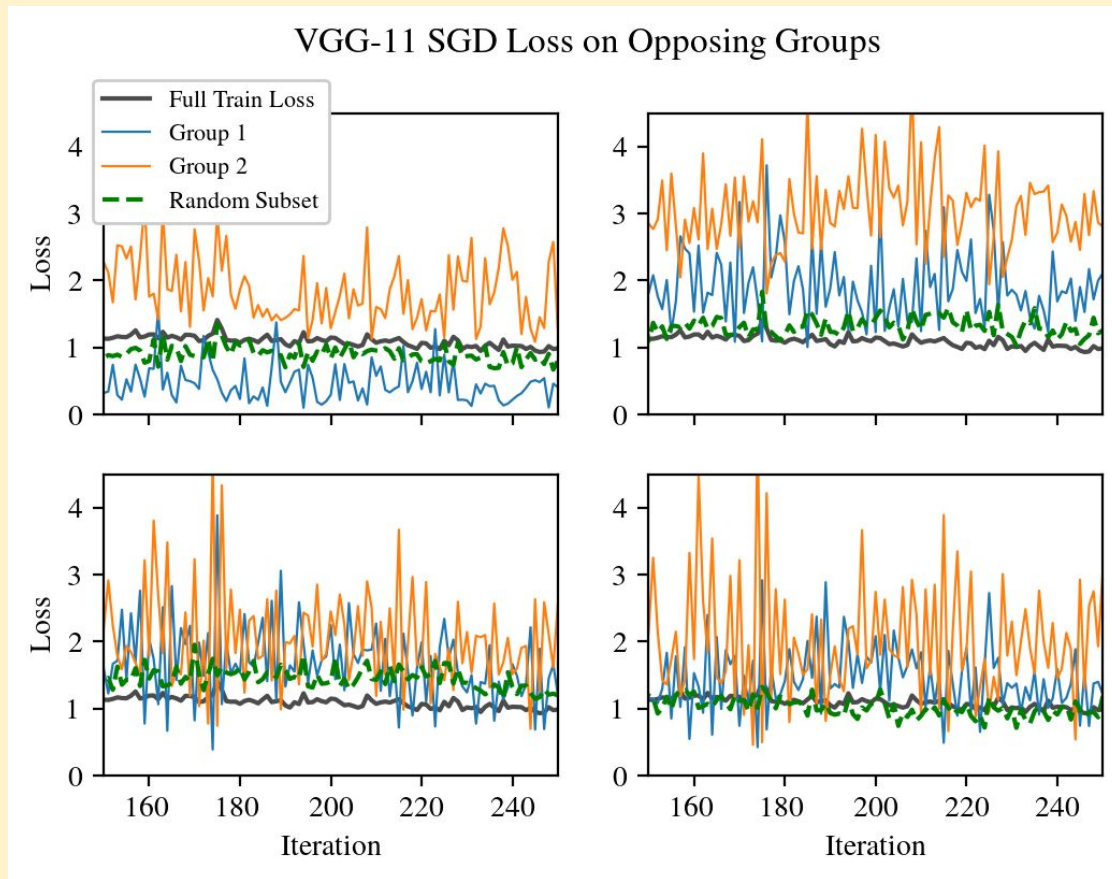Under gradient flow, feature norm grows continuously.

Gradient descent matches flow initially, but *norm starts decreasing* once oscillation begins

# Does this Occur for SGD?

Long story short, **Yes**.

Alternations are not every step.

Groups are not always opposite.



VGG-11 SGD Loss on Opposing Groups

Legend:
- Full Train Loss
- Group 1
- Group 2
- Random Subset

Opposing Signals have clear *potential* connections to existing tools in stochastic optimization, for both training speed and generalization:

- *Batch Normalization*
- *Adaptive Gradient Methods*
- *Sharpness-Aware Minimization*
- *Large Initial Learning Rate*

Maybe these methods work because of how they handle Opposing Signals?

- **Could this help us design new improvements to SGD?**

**Lots** of unanswered questions.
Very happy to discuss further.

*Outliers with Opposing Signals Have an Outsized Effect on Neural Network Optimization*
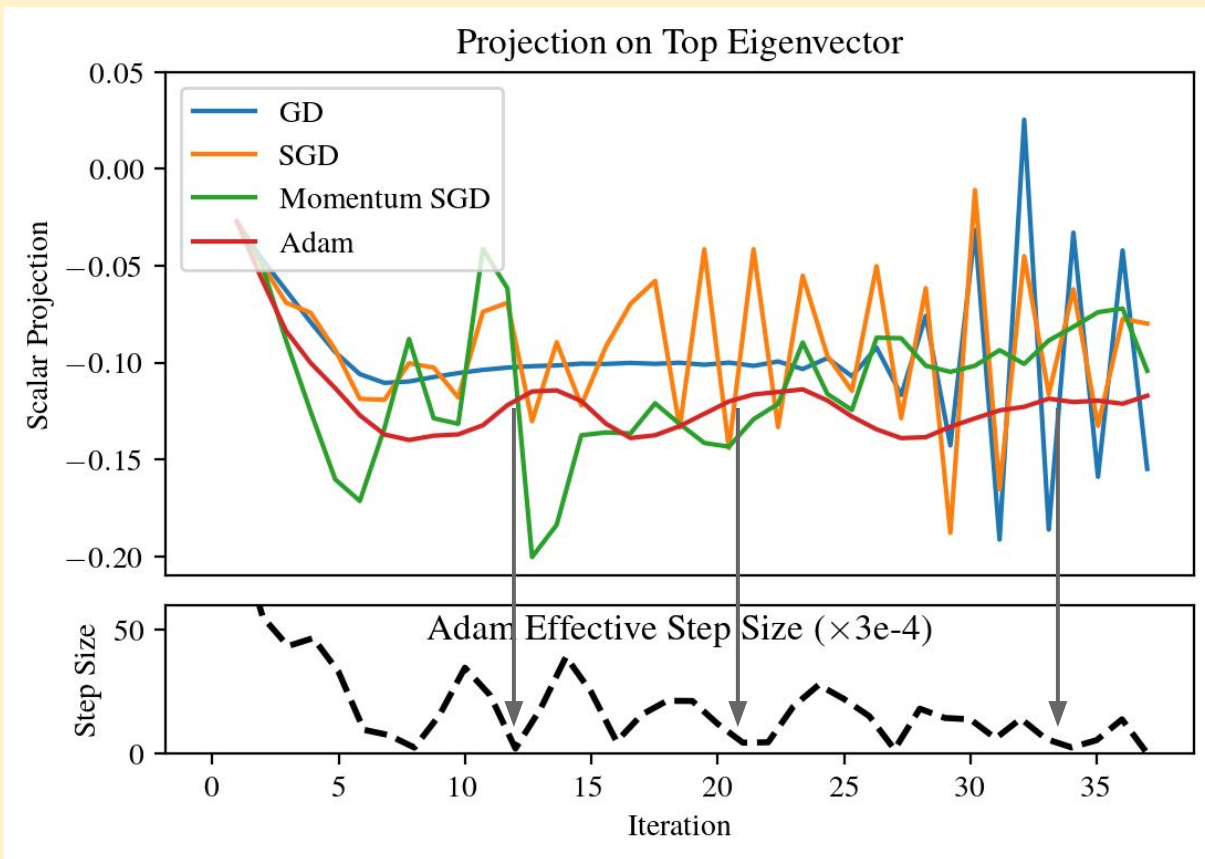
Elan Rosenfeld & Andrej Risteski
https://arxiv.org/abs/2311.04163

# Implications for Stochastic Optimization

# A Case Study of Adam vs. SGD

Adam looks markedly different!

Prevents steps that would *approach the local minimum*

Effective step size **drops sharply** when approaching valley floor.



Projection on Top Eigenvector

Remember this?