# GRADIENT DESCENT GAN OPTIMIZATION IS LOCALLY STABLE.

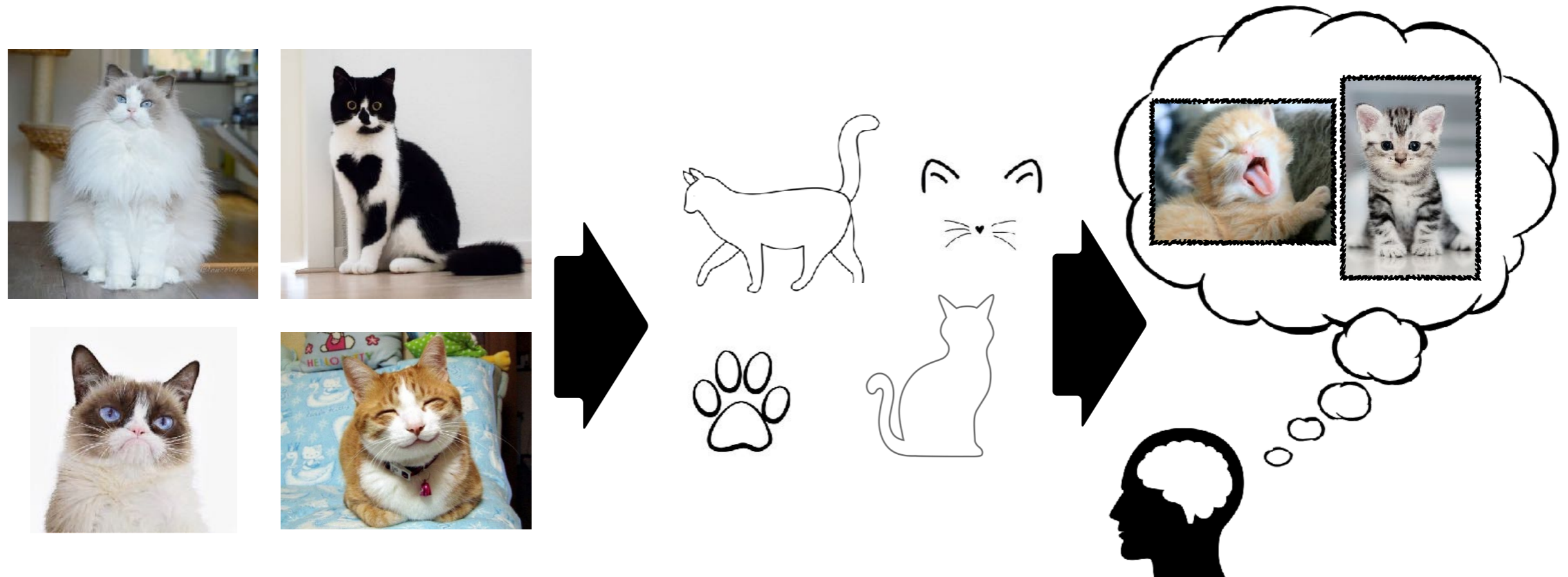**Vaishnavh Nagarajan** | Zico Kolter

**(Based on NIPS '17 Oral paper)**

These slides are adapted from a 1hr talk I presented at CMU for a general CS audience.

1

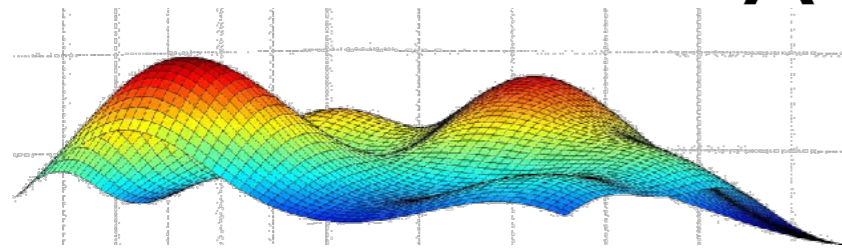# GENERATIVE ADVERSARIAL NETWORKS (GANs)

## A goal of AI: "Understand" data



Build an agent that generates new data (which it does by learning an abstract representation of training data)
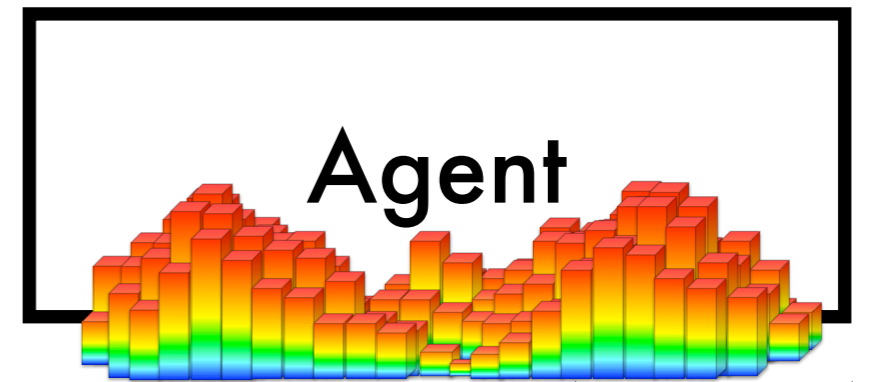
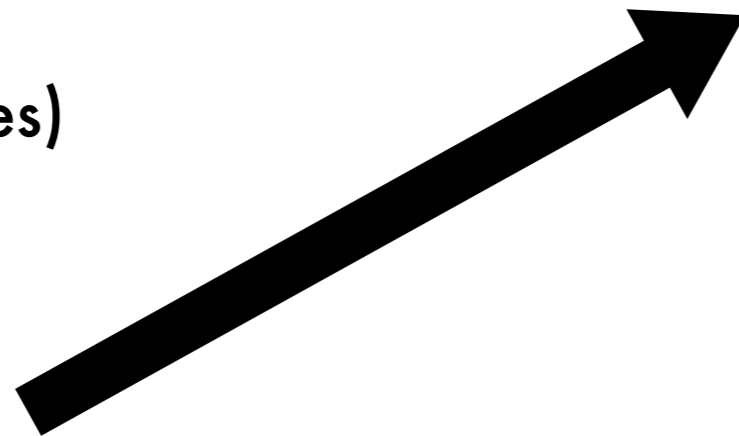# GENERATIVE ADVERSARIAL NETWORKS (GANs)

A generative model



TRUE DISTRIBUTION
(over e.g., cat images)

TRAINING DATA

Agent

(Implicitly)
LEARNED DISTRIBUTION
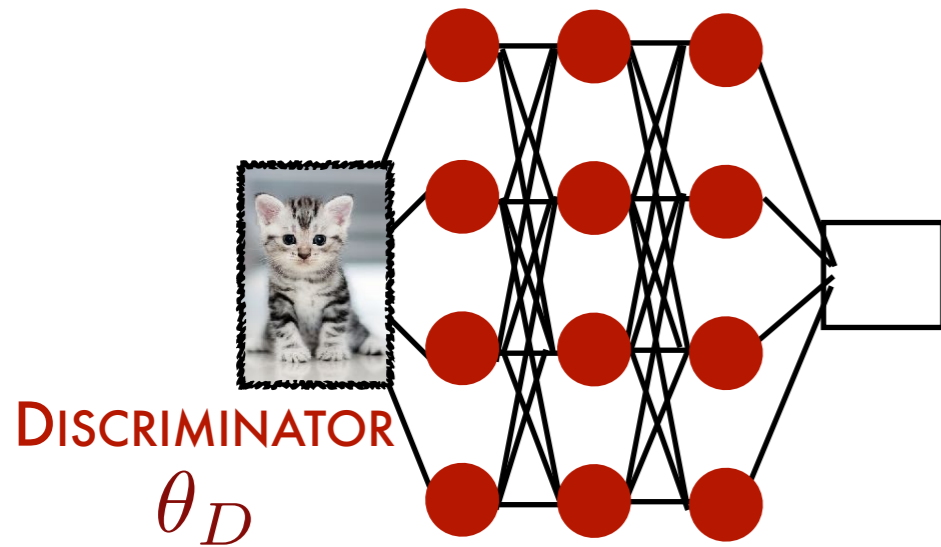
NEW DATA

# PAST WORK

- **GANs were introduced by Goodfellow et al., '14**

- **Many, many variants**: Improved GAN, WGAN, Improved WGAN, Unrolled GAN , InfoGAN MMD-GAN, McGAN,  f-GAN, Fisher GAN, EBGAN, …

- **Wide-ranging applications:** image generation (DCGAN), text-to-image generation (StackGAN), super-resolution (SRGAN) …
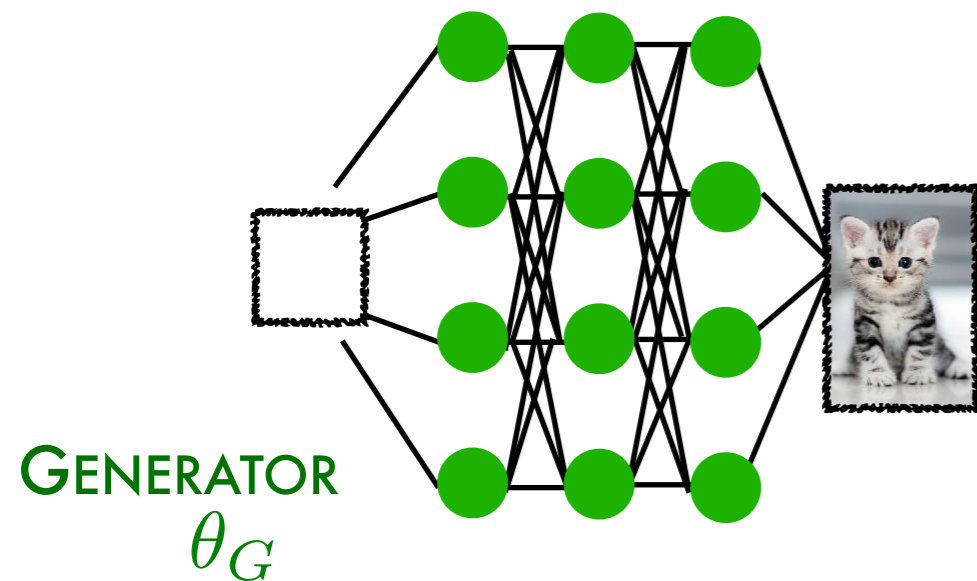
# PAST WORK



"One hour of imaginary celebrities" [Karras et al., '17]

# GENERATIVE ADVERSARIAL NETWORKS (GANS)

DISCRIMINATOR
$\theta_D$

tries its best to tell apart generated images from real images

like a game

$$\min_{\theta_G} \max_{\theta_D} V(\theta_G, \theta_D)$$

GENERATOR
$\theta_G$

tries its best to generate images that discriminator finds real

$\theta_G$

$\theta_D$

GAN OPTIMIZATION: Parameters of two models are iteratively updated (in a standard way) to find "**equilibrium**" of a "min-max objective".

We study dynamics of standard GAN optimization:

Is the equilibrium "locally stable"?
When it is not, how do we make it stable?

# OUTLINE

- **GAN Formulation**

- Toolbox: *Non-linear systems*

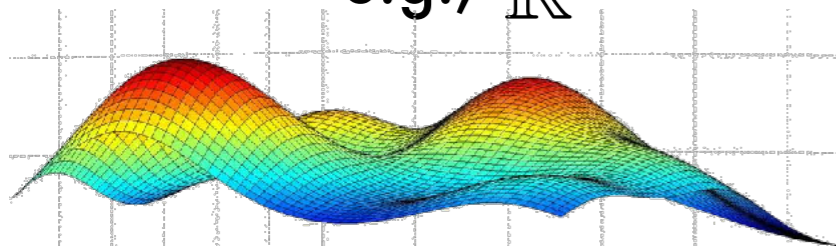- Challenge: *Why is proving stability hard?*

- Main result

- Stabilizing WGANs

# GAN FORMULATION

**Unknown TRUE P.D.F** $p_{\mathrm{data}}(\cdot)$ **over INPUT DOMAIN** $\mathcal{X}$

**e.g.,** $\mathbb{R}^{32 \times 32}$



DISCRIMINATOR $\theta_D$

$$D: \mathcal{X} \to \mathbb{R}$$



GENERATOR $\theta_G$

$$G: \mathcal{Z} \to \mathcal{X}$$

Random input



**Known distribution over latent space** $\mathcal{Z}$ **with P.D.F** $p_{\mathrm{latent}}(\cdot)$



**Generated distribution of** $G(z)$ **over** $\mathcal{X}$ **with P.D.F** $p_{\theta_G}(\cdot)$

# GAN FORMULATION

**Unknown TRUE P.D.F** $p_{\text{data}}(\cdot)$ **over INPUT DOMAIN** $\mathcal{X}$

e.g., $\mathbb{R}^{32 \times 32}$

DISCRIMINATOR $\theta_D$

$$D : \mathcal{X} \to \mathbb{R}$$

GENERATOR $\theta_G$

$$G : \mathcal{Z} \to \mathcal{X}$$

Random input

inducing P.D.F $p_{\theta_G}(\cdot)$

Discriminator's objective: Tell real and generated data apart

$D$ thinks $x$ is:

$D(x) > 0$    real

$D(x) < 0$    generated

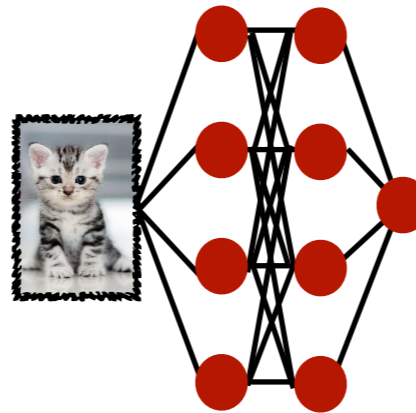$D(x) = 0$    equally both

# GAN FORMULATION

Unknown TRUE P.D.F $p_{\text{data}}(\cdot)$ over INPUT DOMAIN $\mathcal{X}$
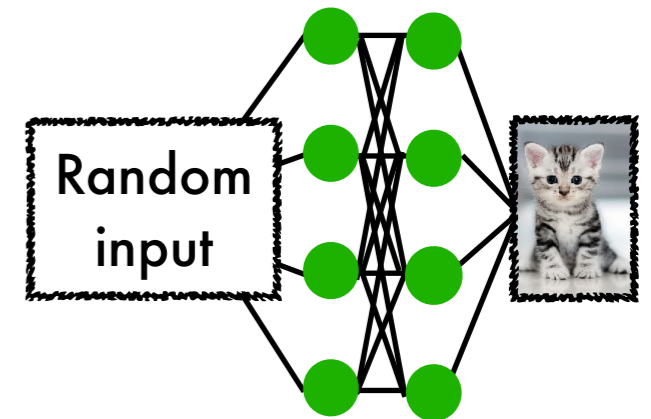
e.g., $\mathbb{R}^{32 \times 32}$

DISCRIMINATOR $\theta_D$
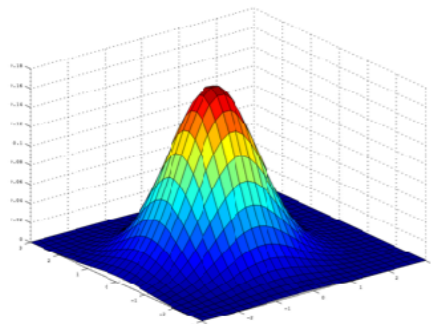
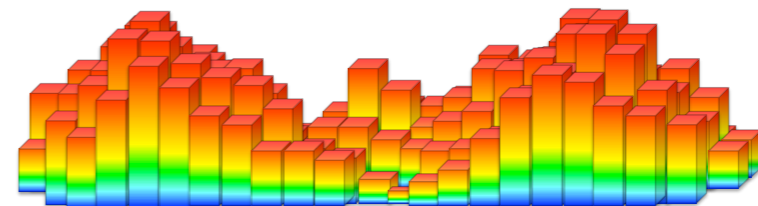$$D: \mathcal{X} \rightarrow \mathbb{R}$$

GENERATOR $\theta_G$

$$G: \mathcal{Z} \rightarrow \mathcal{X}$$

Random input

inducing P.D.F $p_{\theta_G}(\cdot)$



Discriminator's objective: Tell real and generated data apart

$$\max_{\theta_D} V(\theta_G, \theta_D)$$

$$= \mathbb{E}_{x \sim p_{\text{data}}} \left[ f(D(x)) \right] + \mathbb{E}_{z \sim p_{\text{latent}}} \left[ f(-D(G(z))) \right]$$

How real x is, according to the discriminator

How "generated" $G(z)$ looks according to the discriminator
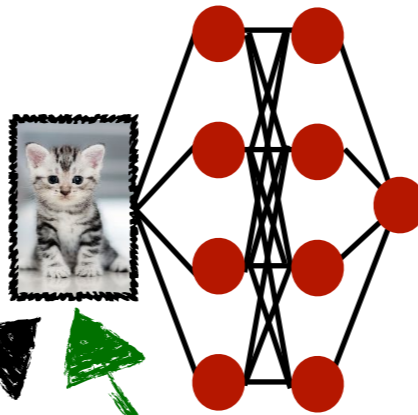
# GAN FORMULATION

Unknown
TRUE P.D.F  $p_{\mathrm{data}}(\cdot)$
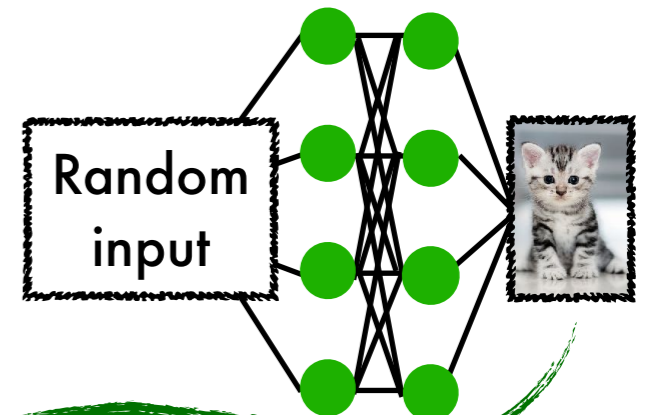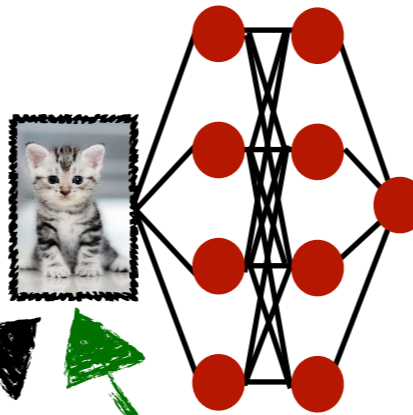over INPUT DOMAIN $\mathcal{X}$

e.g., $\mathbb{R}^{32 \times 32}$

DISCRIMINATOR  $\theta_D$
$$D: \mathcal{X} \to \mathbb{R}$$

??

GENERATOR  $\theta_G$
$$G: \mathcal{Z} \to \mathcal{X}$$

Random input

inducing P.D.F  $p_{\theta_G}(\cdot)$

Generator's objective: Generate data that even the best discriminator can't tell apart from real data

$$\min_{\theta_G} \left[ \max_{\theta_D} V(\theta_G, \theta_D) \right]$$

$$= \mathbb{E}_{x \sim p_{\mathrm{data}}} [f(D(x))] + \mathbb{E}_{z \sim p_{\mathrm{latent}}} [f(-D(G(z)))]$$

How real x is,
according to the discriminator

How "generated" $G(z)$ looks
according to the discriminator

# GAN Formulation

Unknown True P.D.F $p_{\text{data}}(\cdot)$ over Input domain $\mathcal{X}$
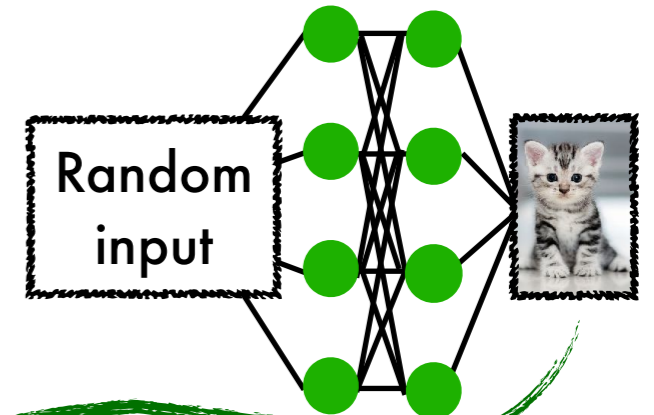
e.g., $\mathbb{R}^{32 \times 32}$

Discriminator $\theta_D$

$$D : \mathcal{X} \to \mathbb{R}$$

??

Generator $\theta_G$

$$G : \mathcal{Z} \to \mathcal{X}$$

Random input

inducing P.D.F $p_{\theta_G}(\cdot)$

| Traditional GAN | Wasserstein GAN (WGAN) |
|---|---|
| $f(t) = \log\left(\dfrac{1}{1 + \exp(-t)}\right)$ | $f(t) = t$ |

$$= \mathbb{E}_{x \sim p_{\text{data}}}\left[f(D(x))\right] + \mathbb{E}_{z \sim p_{\text{latent}}}\left[f(-D(G(z)))\right]$$
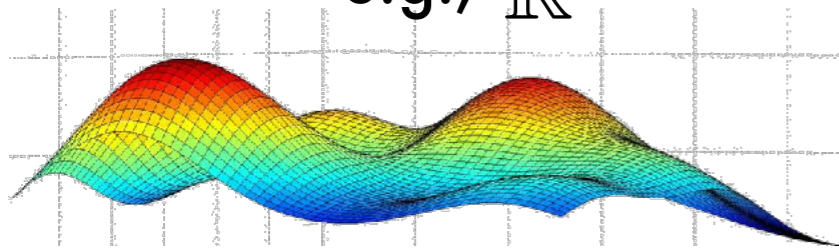
How real x is, according to the discriminator

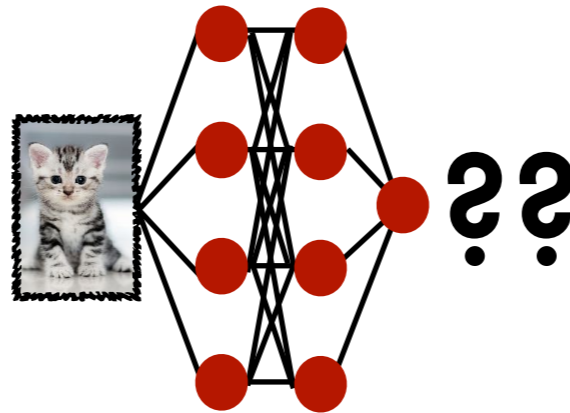How "generated" $G(z)$ looks according to the discriminator

# GAN FORMULATION

Unknown
TRUE P.D.F $p_{\text{data}}(\cdot)$
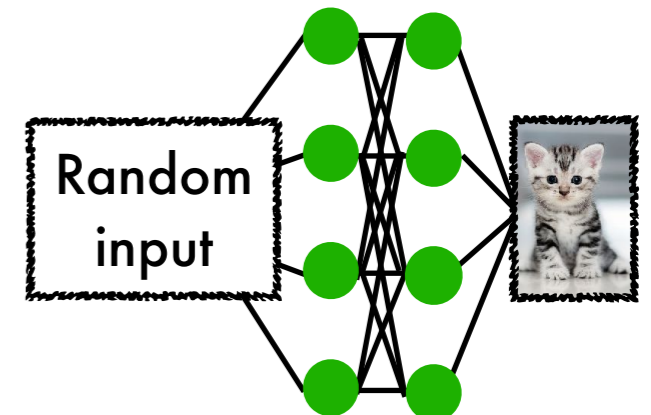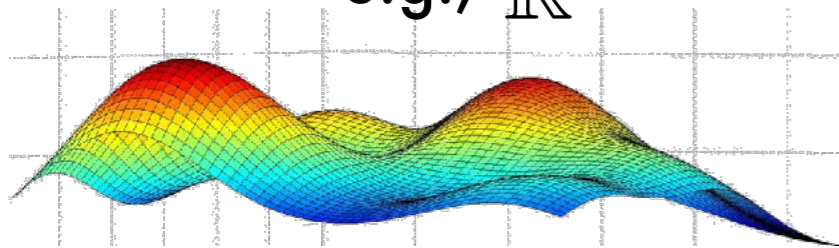over INPUT DOMAIN $\mathcal{X}$

e.g., $\mathbb{R}^{32 \times 32}$

DISCRIMINATOR $\theta_D$
$$D: \mathcal{X} \to \mathbb{R}$$

??

GENERATOR $\theta_G$
$$G: \mathcal{Z} \to \mathcal{X}$$

Random input

inducing P.D.F $p_{\theta_G}(\cdot)$

SOLUTION: Generator matches true distribution and discriminator cannot tell apart data from either. How do we find this solution?

$$\min_{\theta_G} \left[ \max_{\theta_D} V(\theta_G, \theta_D) \right]$$

$$= \mathbb{E}_{x \sim p_{\text{data}}} \left[ f(D(x)) \right] + \mathbb{E}_{z \sim p_{\text{latent}}} \left[ f(-D(G(z))) \right]$$

How real x is,
according to the discriminator

How "generated" $G(z)$ looks
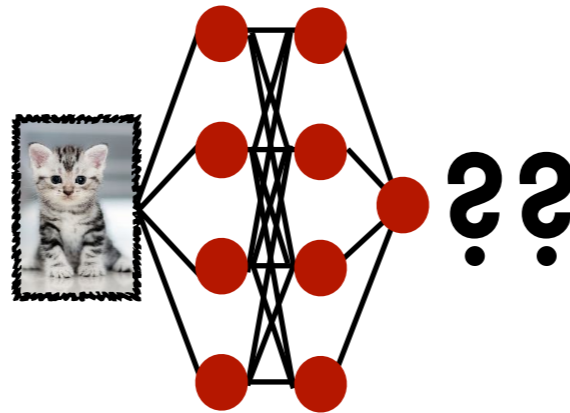according to the discriminator

# GAN Optimization

We consider: **infinitesimal, simultaneous** gradient ascent/descent updates

$$\min_{\theta_G} \left[ \max_{\theta_D} V(\theta_G, \theta_D) \right]$$

Repeat simultaneously:

time derivative

$$\dot{\theta}_D = \nabla_{\theta_D} V(\theta_G, \theta_D)$$

$$\dot{\theta}_G = -\nabla_{\theta_G} V(\theta_G, \theta_D)$$

until **equilibrium**:

$$\dot{\theta}_D = 0$$

$$\dot{\theta}_G = 0$$

# OUTLINE

- GAN Formulation

- **Toolbox: *Non-linear systems***

- Challenge: *Why is proving stability hard?*

- Main result

- Stabilizing WGANs

# LOCALLY EXPONENTIALLY STABLE

Consider a dynamical system $\dot{\theta} = h(\theta)$ for which $\theta^\star$ is an equilibrium point i.e., $h(\theta^\star) = 0$

INFORMAL DEFINITION: The equilibrium point is **locally exponentially stable** if **any** initialization of the system sufficiently close to the equilibrium, converges to the equilibrium point "very quickly" (distance to equilibrium decays at the rate $\propto e^{-O(t)}$)

# PROVING STABILITY

Consider a dynamical system $\dot{\theta} = h(\theta)$ for which $\theta^\star$ is an equilibrium point i.e., $h(\theta^\star) = 0$

LINEARIZATION THEOREM: The equilibrium of this (non-linear) system is locally exponentially stable if and only if its Jacobian at equilibrium HAS EIGENVALUES WITH **STRICTLY NEGATIVE REAL PARTS:**

$$J = \left.\frac{\partial h(\theta)}{\partial \theta}\right|_{\theta^\star} = \left[\begin{array}{ccc} \frac{\partial h_1(\theta)}{\partial \theta_1} & \frac{\partial h_1(\theta)}{\partial \theta_2} & \cdots \\ \frac{\partial h_2(\theta)}{\partial \theta_1} & \frac{\partial h_2(\theta)}{\partial \theta_2} & \cdots \\ \frac{\partial h_3(\theta)}{\partial \theta_1} & \vdots & \vdots \cdots \end{array}\right]_{\theta = \theta^\star}$$

(asymmetric, real square matrix with possibly complex eigenvalues)

$$Jv = \lambda v \implies Re(\lambda) < 0$$

# PROVING STABILITY

Consider a dynamic[al system with] equilibrium point i.e.[ ]

LINEARIZATION THEOR[EM]

locally exponentiall[y ]

HAS EIGENVALUES WIT[H]

1D Sanity Check/Intuition:

$$\dot{\theta} = -\theta$$

Origin

$$J = -1 < 0$$

$$\theta$$

$$J = \left. \frac{\partial h(\theta)}{\partial \theta} \right|_{\theta^{\star}}$$

(asymmetric, real square matrix with possibly complex eigenvalues)

$$Jv = \lambda v \implies Re(\lambda) < 0$$

# OUTLINE

- GAN Formulation

- Toolbox: Non-linear systems

- **Challenge: *Why is proving stability hard?***
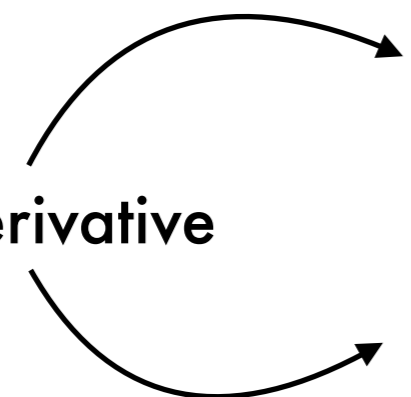
- Main result

- Stabilizing GANs

# RECALL: GAN OPTIMIZATION

We consider: **infinitesimal, simultaneous** gradient descent updates

$$\min_{\theta_G} \left[ \max_{\theta_D} V(\theta_G, \theta_D) \right]$$

Repeat simultaneously:

time derivative

$$\dot{\theta}_D = \nabla_{\theta_D} V(\theta_G, \theta_D)$$

$$\dot{\theta}_G = -\nabla_{\theta_G} V(\theta_G, \theta_D)$$

until equilibrium:
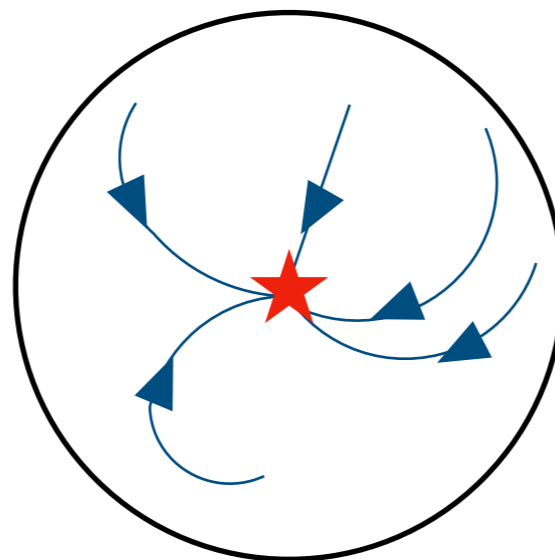
$$\dot{\theta}_D = 0$$

$$\dot{\theta}_G = 0$$

# WHY IS PROVING **GAN** STABILITY HARD?

GAN involves **concave-minimization**–**concave-maximization**, even for a linear **discriminator** and a **generator**.

$$D(x) = \theta_D x \qquad G(z) = \theta_G z$$

$$\min_{\theta_G} \max_{\theta_D} \mathbb{E}_{x \sim p_{\text{data}}}[f(\theta_D x)] + \mathbb{E}_{z \sim p_{\text{latent}}}[f(-\theta_D \theta_G z)]$$

Given $f$ is concave for GANs ,
objective is **concave w.r.t** $\theta_G$.

# WHY IS PROVING **GAN** STABILITY HARD?

If objective were convex-concave, would've been easy!

but for GANs, it is concave-concave. Sad!



The descent-ascent updates individually point in the direction of equilibrium.

The generator's descent updates take us away from equilibrium!

# WHY IS PROVING GAN STABILITY HARD?

SOME CONCURRENT WORK:

Mescheder et al., '17: GANs may **not** be stable.

Heusel et al., '17, Li et al., '17: Stable provided
<span style="color:darkred">discriminator updates</span> "dominate" <span style="color:green">generator</span>
<span style="color:green">updates</span> in some way. e.g.,

$$\dot{\theta}_D = \nabla_{\theta_D} V(\theta_G, \theta_D) \text{ x } 100$$

$$\dot{\theta}_G = -\nabla_{\theta_G} V(\theta_G, \theta_D)$$

But GANs  in practice: updated with
"equal weights"...

**Despite a <span style="color:green">concave</span>-<span style="color:red">concave</span> objective,**
simultaneous gradient descent GAN equilibrium
*is*
"locally exponentially stable"
under suitable conditions
on the representational powers of
the discriminator & generator.

# OUTLINE

- GAN Formulation

- Toolbox: *Non-linear systems*

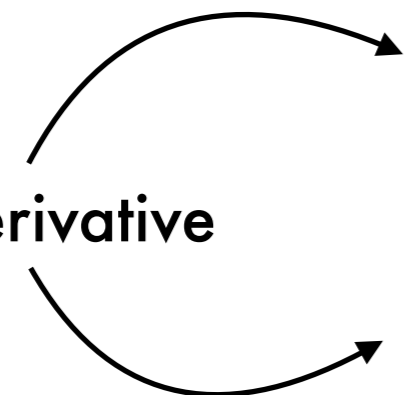- Challenge: *Why is proving stability hard?*

- **Main result: *GANs are stable***

- Stabilizing WGANs

# ASSUMPTION 1

Consider an equilibrium point $(\theta_D^\star, \theta_G^\star)$ such that generated distribution matches true distribution:

$$p_{\theta_G^\star}(\cdot) = p_{\mathrm{data}}(\cdot)$$

and discriminator cannot tell real and generated data apart:

$$D_{\theta_D^\star}(x) = 0 \quad \text{for all } x$$

NOTE:
1. This *is* an equilibrium point (updates are 0 here).
2. Other kinds of equilibria may exist.
3. More relaxations in the paper, but at the cost of other restrictions

# ASSUMPTION 2

Consider the objective at the equilibrium generator,
as a function of the discriminator.

$$V(\theta_D, \theta_G^\star)$$



AT EQUILIBRIUM DISCRIMINATOR, this is already a concave function.

---

We assume stronger curvature.
the corresponding Hessian $\nabla_{\theta_D}^2 V(\theta_D, \theta_G^\star)$ evaluated at equilibrium discriminator is **NEGATIVE DEFINITE.**

# ASSUMPTION 3

Consider

"the magnitude of the objective's gradient w.r.t equilibrium discriminator",

as a function of the generator.

$$\left\|\nabla_{\theta_D} V(\theta_D, \theta_G)\right\|^2\Big|_{\theta_D = \theta_D^\star}$$



$$\theta_G^\star$$

AT EQUILIBRIUM GENERATOR, this is already
a convex function.

---

We assume stronger curvature.

the Hessian $\nabla_{\theta_G}^2 \left\|\nabla_{\theta_D} V(\theta_D, \theta_G)\right\|^2\Big|_{\theta_D = \theta_D^\star}$

evaluated at equilibrium generator is POSITIVE DEFINITE.

These strong curvature assumptions
imply a locally unique equilibrium.
We also consider a specific relaxation allowing
a subspace of equilibria.



flat
direction

strong curvature

# RECALL: GAN OPTIMIZATION

We consider: **infinitesimal, simultaneous** gradient descent updates

$$\min_{\theta_G} \left[ \max_{\theta_D} V(\theta_G, \theta_D) \right]$$

Repeat simultaneously:

time derivative

$$\dot{\theta}_D = \nabla_{\theta_D} V(\theta_G, \theta_D)$$

$$\dot{\theta}_G = -\nabla_{\theta_G} V(\theta_G, \theta_D)$$

until equilibrium:

$$\dot{\theta}_D = 0$$

$$\dot{\theta}_G = 0$$

# MAIN RESULT

THEOREM: Under assumptions 1-3, the equilibrium of the simultaneous gradient descent GAN system is locally exponentially stable.

# Main Result

Theorem: Under assumptions 1-3, the equilibrium of the simultaneous gradient descent GAN system is locally exponentially stable.

Specifically, the Jacobian at equilibrium has eigenvalues with strictly negative real parts.

$$J = \left. \frac{\partial h(\theta)}{\partial \theta} \right|_{\theta^\star} = \begin{bmatrix} \frac{\partial h_1(\theta)}{\partial \theta_1} & \frac{\partial h_1(\theta)}{\partial \theta_2} & \cdots \\ \frac{\partial h_2(\theta)}{\partial \theta_1} & \frac{\partial h_2(\theta)}{\partial \theta_2} & \cdots \\ \frac{\partial h_3(\theta)}{\partial \theta_1} & \vdots & \vdots \ddots \end{bmatrix}_{\theta = \theta^\star}$$

(asymmetric, real square matrix with possibly complex eigenvalues)

$$Jv = \lambda v \implies Re(\lambda) < 0$$

# PROOF OUTLINE

Jacobian at equilibrium:

$$\begin{bmatrix} \color{darkred}{\partial \dot{\theta}_D}/\color{darkred}{\partial \theta_D} & \color{darkred}{\partial \dot{\theta}_D}/\color{green}{\partial \theta_G} \\ \color{green}{\partial \dot{\theta}_G}/\color{darkred}{\partial \theta_D} & \color{green}{\partial \dot{\theta}_G}/\color{green}{\partial \theta_G} \end{bmatrix}$$

# PROOF OUTLINE

Jacobian at equilibrium:

$$\begin{bmatrix} \boxed{\nabla^2_{\theta_D} V(\theta_D, \theta_G)} & \partial\dot{\theta}_D / \partial\theta_G \\ \partial\dot{\theta}_G / \partial\theta_D & \partial\dot{\theta}_G / \partial\theta_G \end{bmatrix} = \begin{bmatrix} \text{negative definite} & \\ & \end{bmatrix}$$



A negative definite diagonal matrix makes it more likely that the whole matrix has eigenvalues with negative real parts.

# PROOF OUTLINE

Jacobian at equilibrium:

$$\begin{bmatrix} \nabla^2_{\theta_D} V(\theta_D, \theta_G) & \nabla_{\theta_G} \nabla_{\theta_D} V(\theta_D, \theta_G) \\ -\left( \nabla_{\theta_G} \nabla_{\theta_D} V(\theta_D, \theta_G) \right)^T \frac{\partial \dot{\theta}_G}{\partial \theta_G} & \Big/ \frac{\partial \dot{\theta}_G}{\partial \theta_G} \end{bmatrix} = \begin{bmatrix} \text{negative definite} & \text{full column rank} \\ \text{negative transpose} & \end{bmatrix}$$

$$\cdot \left( \nabla_{\theta_G} \nabla_{\theta_D} V(\theta_D, \theta_G) \right)^T \nabla_{\theta_G} \nabla_{\theta_D} V(\theta_D, \theta_G) = \nabla^2_{\theta_G} \left\| \nabla_{\theta_D} V(\theta_D, \theta_G) \right\|^2 \Big|_{\theta_D = \theta_D^\star}$$

Assumption 3: positive definite

A negative definite diagonal matrix makes it more likely that the whole matrix has eigenvalues with negative real parts.

# PROOF OUTLINE

Jacobian at equilibrium:

$$\begin{bmatrix} \nabla^2_{\theta_D} V(\theta_D, \theta_G) & \nabla_{\theta_G} \nabla_{\theta_D} V(\theta_D, \theta_G) \\ -\left(\nabla_{\theta_G} \nabla_{\theta_D} V(\theta_D, \theta_G)\right)^T & \boxed{-\nabla^2_{\theta_G} V(\theta_D, \theta_G)} \end{bmatrix}$$



but it is concave-concave!

$\theta_D$

$\theta_G$

The generator's descent updates take us away from equilibrium!

25

A negative definite diagonal n... whole matrix has eigenv...

# PROOF OUTLINE

Jacobian at equilibrium:

$$\begin{bmatrix} \nabla^2_{\theta_D} V(\theta_D, \theta_G) & \nabla_{\theta_G} \nabla_{\theta_D} V(\theta_D, \theta_G) \\ -\left(\nabla_{\theta_G} \nabla_{\theta_D} V(\theta_D, \theta_G)\right)^T & -\nabla^2_{\theta_G} V(\theta_D, \theta_G) \end{bmatrix} = \begin{bmatrix} \text{negative definite} & \text{full column rank} \\ \text{negative transpose} & 0 \end{bmatrix}$$
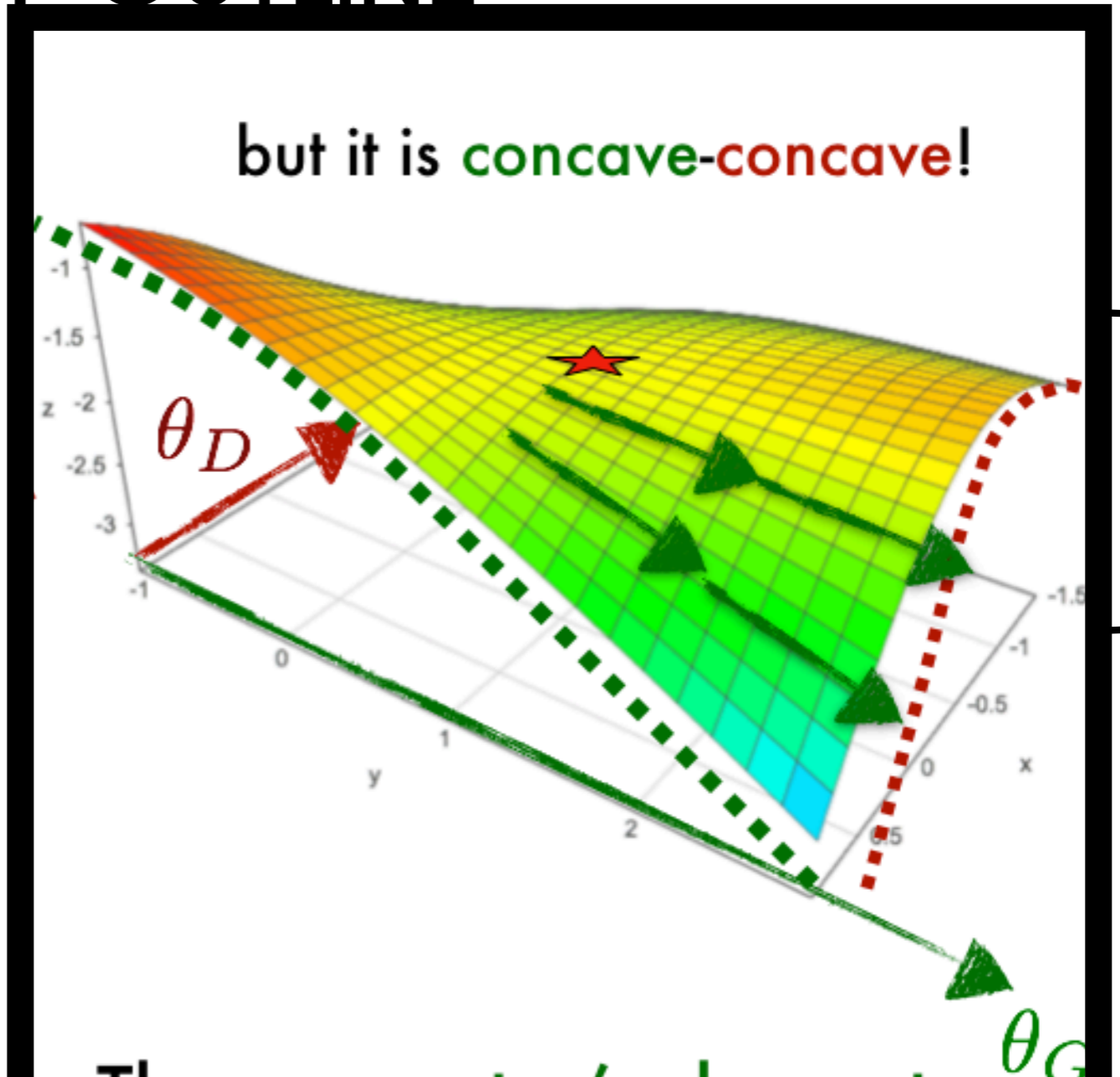
could be (− negative definite)
i.e., positive definite!

A negative definite diagonal matrix makes it more likely that the whole matrix has eigenvalues with negative real parts.

# PROOF OUTLINE

Jacobian at equilibrium:

$$\begin{bmatrix} \nabla^2_{\theta_D} V(\theta_D, \theta_G) & \nabla_{\theta_G}\nabla_{\theta_D} V(\theta_D, \theta_G) \\ -\left(\nabla_{\theta_G}\nabla_{\theta_D} V(\theta_D, \theta_G)\right)^T & -\nabla^2_{\theta_G} V(\theta_D, \theta_G) \end{bmatrix} = \begin{bmatrix} \text{negative definite} & \text{full column rank} \\ \text{negative transpose} & 0 \end{bmatrix}$$

fix discriminator as all-zero equilibrium discriminator, objective is a constant:

$$\mathbb{E}_{p_{\text{data}}}[f(0)] + \mathbb{E}_{p_{\theta_G}}[f(0)] = 2f(0)$$

A negative definite diagonal matrix makes it more likely that the whole matrix has eigenvalues with negative real parts.

# PROOF OUTLINE

Jacobian at equilibrium:

$$\begin{bmatrix} \nabla^2_{\theta_D} V(\theta_D, \theta_G) & \nabla_{\theta_G} \nabla_{\theta_D} V(\theta_D, \theta_G) \\ -\left(\nabla_{\theta_G} \nabla_{\theta_D} V(\theta_D, \theta_G)\right)^T & -\nabla^2_{\theta_G} V(\theta_D, \theta_G) \end{bmatrix} = \begin{bmatrix} \text{negative definite} & \text{full column rank} \\ \text{negative transpose} & 0 \end{bmatrix}$$

MAIN LEMMA: Matrices J of this form have eigenvalues with **strictly** negative real parts:
$$J v = \lambda v \implies Re(\lambda) < 0$$

THUS, THE **GAN** EQUILIBRIUM IS LOCALLY EXPONENTIALLY STABLE.

# OUTLINE

- GAN Formulation

- Toolbox: *Non-linear systems*

- Challenge: *Why is proving stability hard?*

- Main result: *GANs are stable*

- **Stabilizing WGANs**

# WGAN

Jacobian at equilibrium:

$$\begin{bmatrix} \boxed{\nabla^2_{\theta_D} V(\theta_D, \theta_G)} & \nabla_{\theta_G} \nabla_{\theta_D} V(\theta_D, \theta_G) \\ -\left(\nabla_{\theta_G} \nabla_{\theta_D} V(\theta_D, \theta_G)\right)^T & -\nabla^2_{\theta_G} V(\theta_D, \theta_G) \end{bmatrix} = \begin{bmatrix} 0 & \text{full column rank} \\ \text{negative transpose} & 0 \end{bmatrix}$$
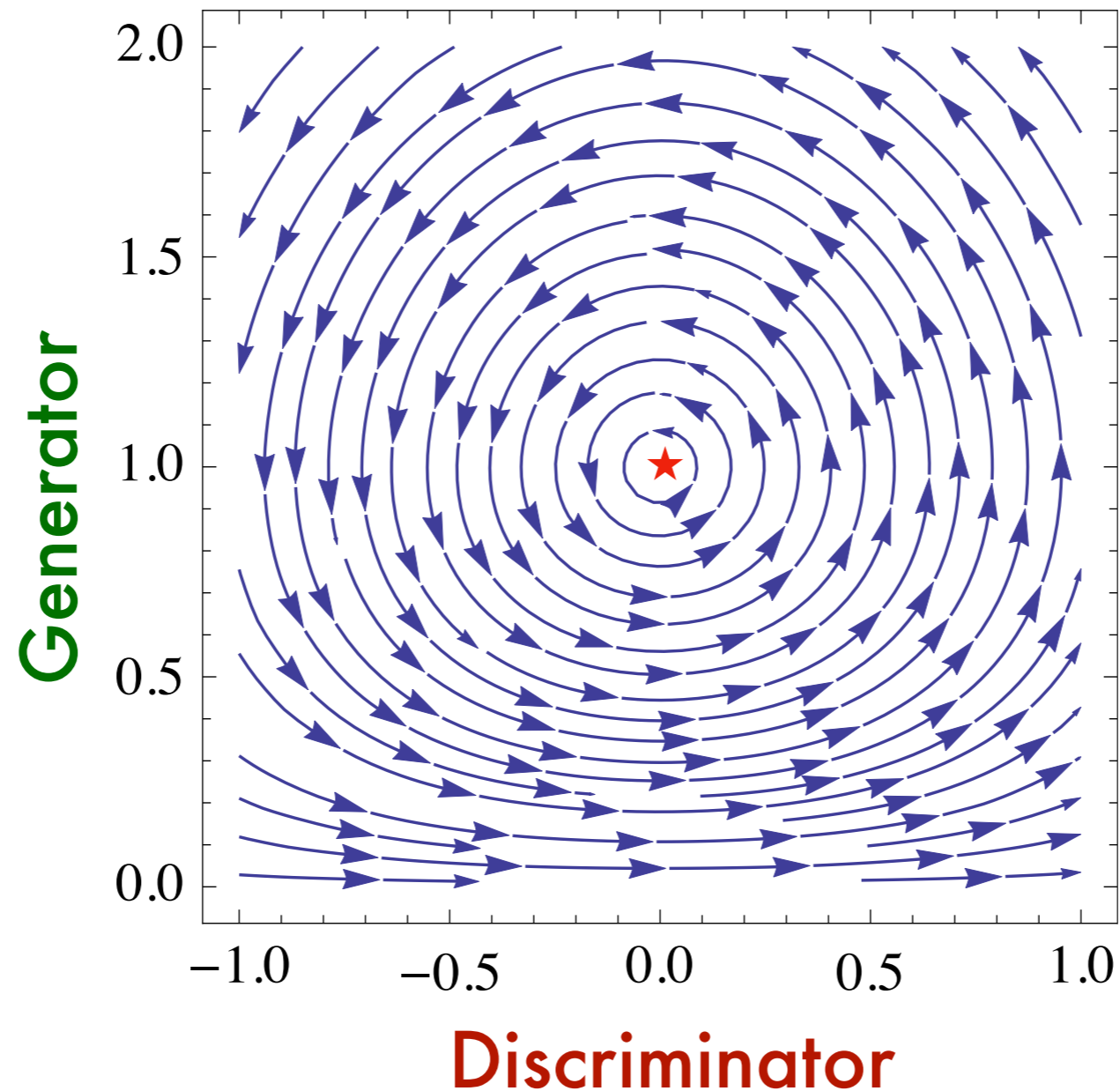
$$f(t) = t$$

fix generator to be at equilibrium:

$$\mathbb{E}_{p_{\text{data}}}[D(x)] + \mathbb{E}_{p_{\theta_G^\star}}[-D(x)] = 0$$

THEOREM: There exists an equilibrium for simultaneous gradient descent WGAN that does not converge locally.
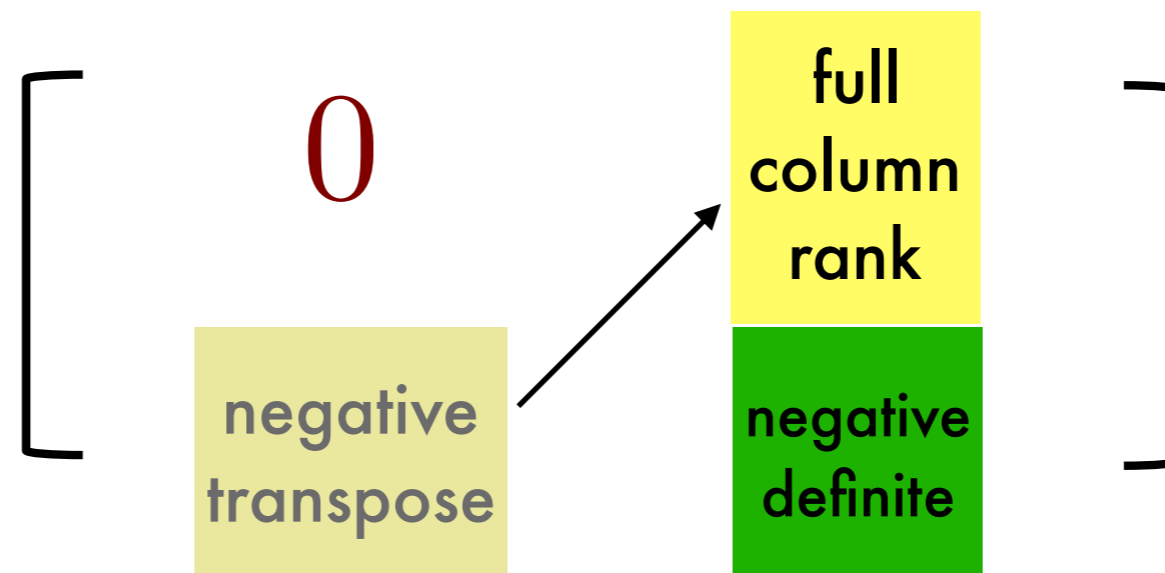
# WGAN



A system learning
a uniform
distribution.

THEOREM: There exists an equilibrium for simultaneous gradient descent WGAN that does not converge locally.

# GRADIENT-NORM BASED REGULARIZATION

$$\dot{\theta}_D = \nabla_{\theta_D} V(\theta_D, \theta_G)$$

$$\dot{\theta}_G = -\nabla_{\theta_G} V(\theta_D, \theta_G) - \eta \nabla_{\theta_G} \, \|\nabla_{\theta_D} V(\theta_D, \theta_G)\|^2$$

Generator minimizes (the objective + the norm of the discriminator's gradient).

$$
\begin{bmatrix}
0 & \text{full column rank} \\
\text{negative transpose} & \text{negative definite}
\end{bmatrix}
$$
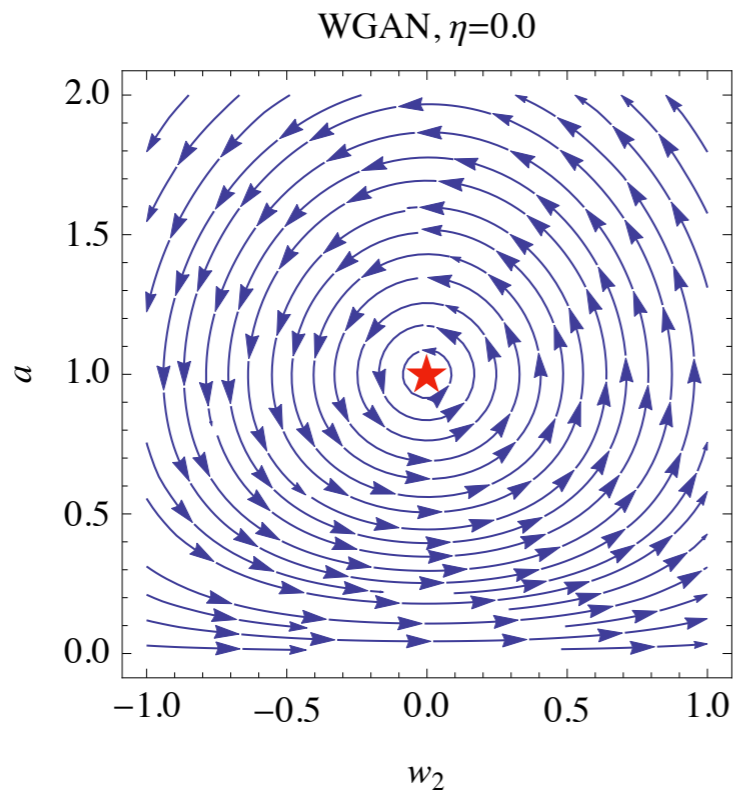
THEOREM: Under similar assumptions, the equilibrium of the regularized simultaneous gradient descent (W)GAN system is locally exponentially stable when $\eta$ not too large.
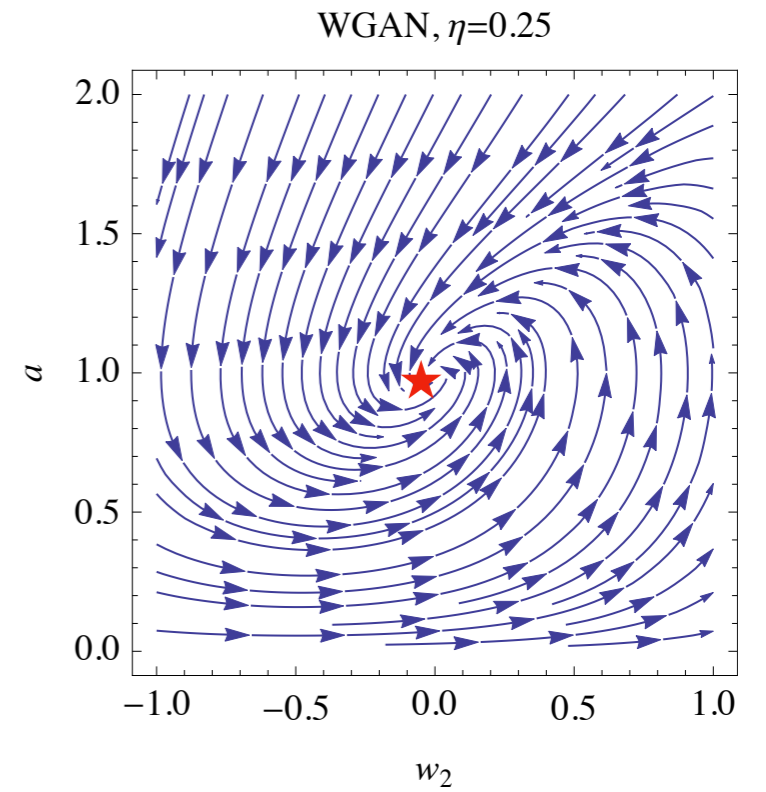
# REGULARIZED WGAN
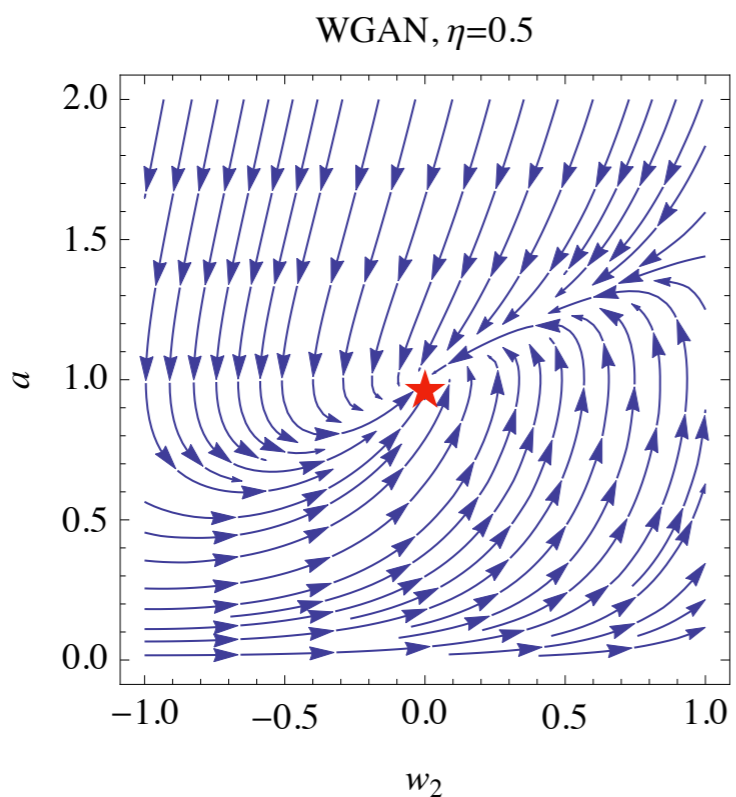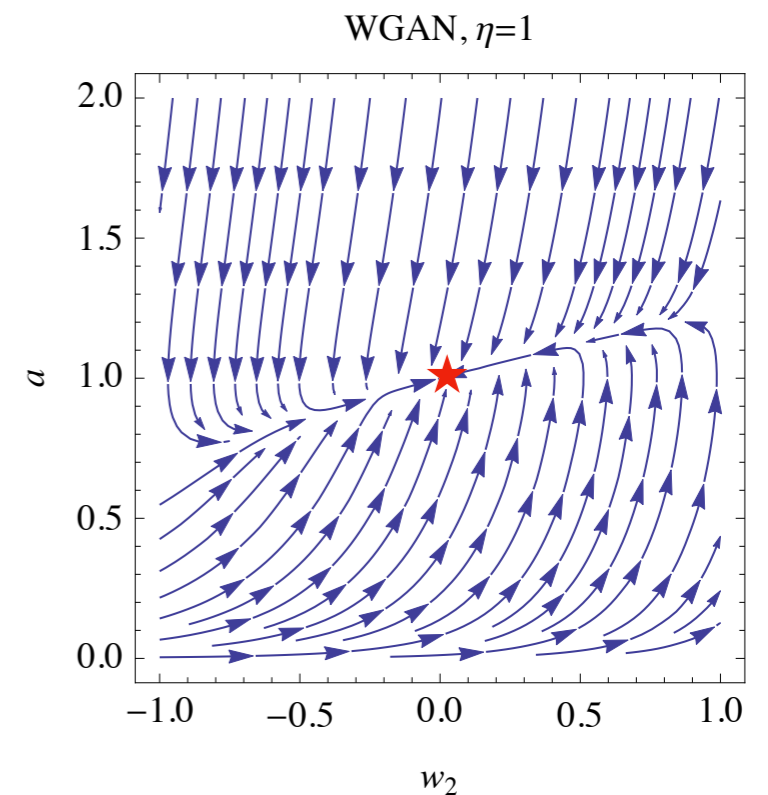## (learning a uniform distribution)



$$\eta = 0$$

$$\eta = 0.25$$

$$\eta = 0.5$$

$$\eta = 1.0$$

# FORESIGHTED GENERATOR

GAN training: a game where discriminator and generator try to outdo each other until neither can outdo the other.



Traditional GAN Generator is "greedy"

(Initialization)

Generated points

Real points

Greedy generator strategy: Generate only one data point: the one to which discriminator has assigned highest value ("most real" according to discriminator).
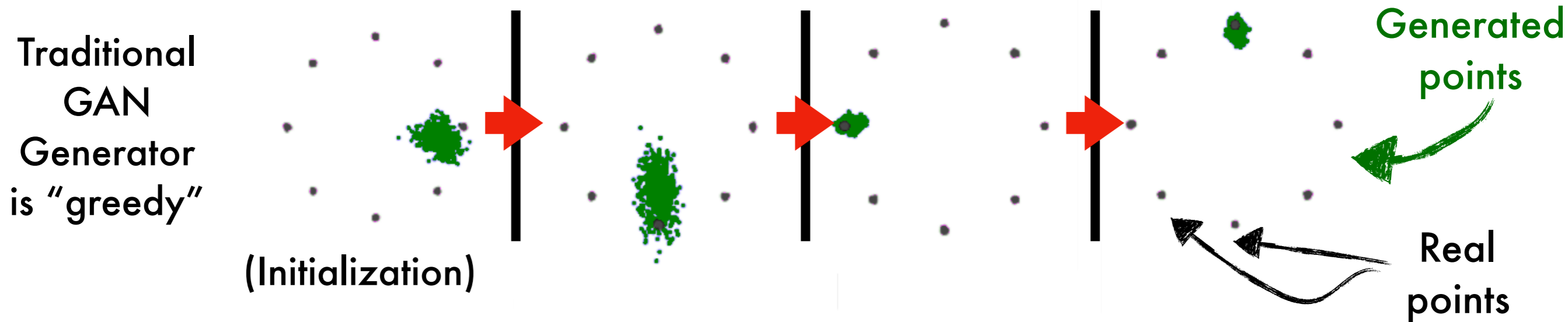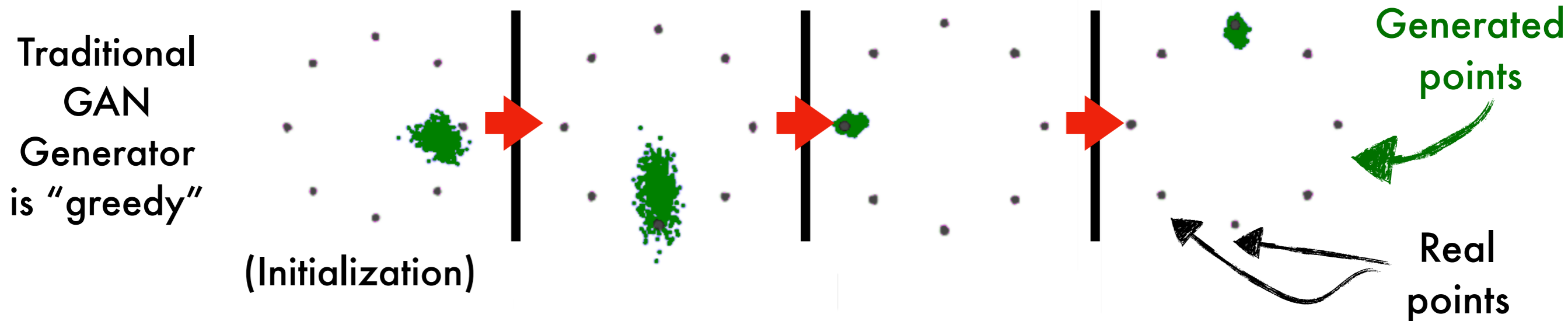
# FORESIGHTED GENERATOR

GAN training: a game where discriminator and generator try to outdo each other until neither can outdo the other.



Traditional GAN Generator is "greedy"

(Initialization)

Generated points

Real points

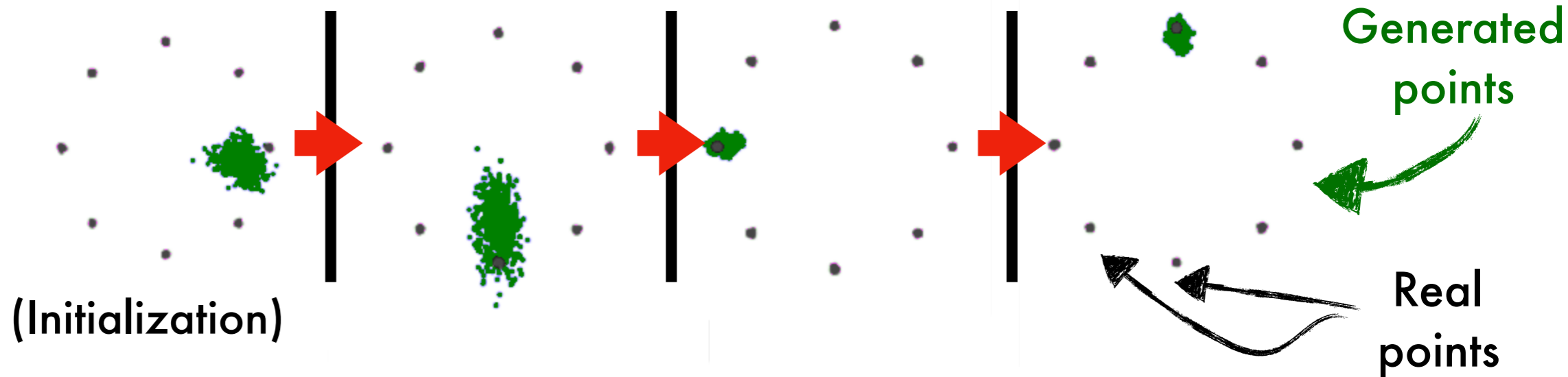OBSERVATION: Generator keeps updating to state where objective $V(\theta_G, \theta_D)$ is small but discriminator update $\|\nabla_{\theta_D} V(\theta_D, \theta_G)\|^2$ is large.

SOLUTION: Generator explicitly seeks state where objective $V(\theta_G, \theta_D)$ is small AND discriminator update $\|\nabla_{\theta_D} V(\theta_D, \theta_G)\|^2$ is small.
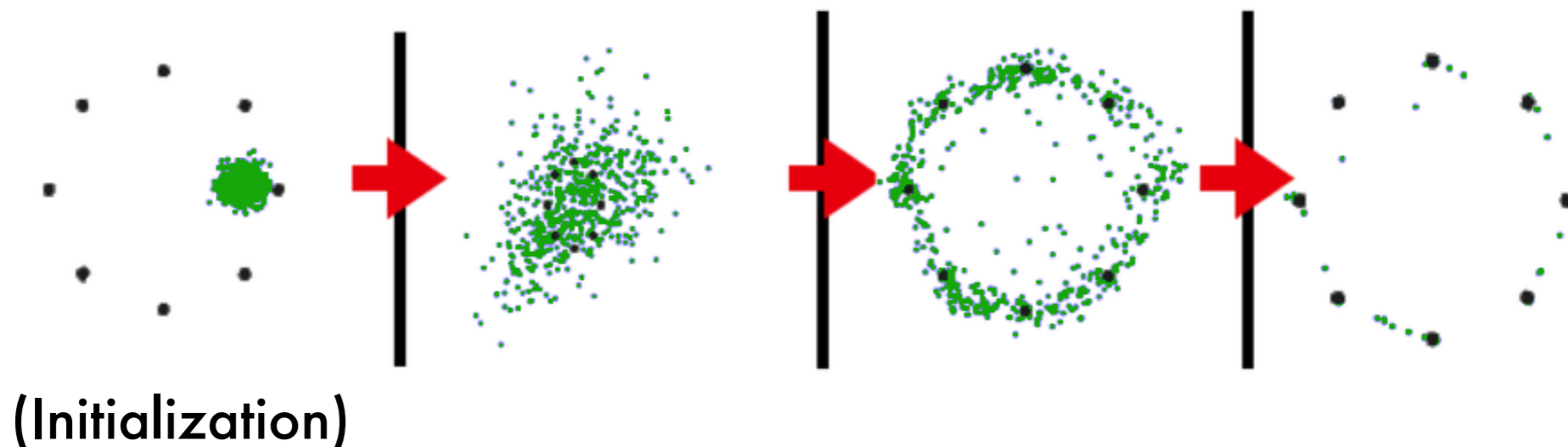
# FORESIGHTED GENERATOR

GAN training: a game where discriminator and generator try to outdo each other until neither can outdo the other.



Traditional GAN Generator is "greedy"

(Initialization)

Generated points

Real points

$$\dot{\theta}_G = -\nabla_{\theta_G} V(\theta_D, \theta_G) - \eta \nabla_{\theta_G} \|\nabla_{\theta_D} V(\theta_D, \theta_G)\|^2$$

Traditional GAN but with Regularized Generator

(Initialization)

# CONCLUSION

- Theoretical analysis of local convergence/stability of simultaneous gradient descent GANs using non-linear systems.

- GAN objective is <span style="color:red">concave</span>-<span style="color:green">concave</span>, yet simultaneous gradient descent is locally stable — perhaps why GANs have worked well in practice.

- Our analysis yields a regularization term that provides more stability.

# OPEN QUESTIONS

- Prove local stability for a more general case

- Global convergence?

- Many more theoretical questions in GANs: when do equilibria exist? Do they generalize?

# THANK YOU.
# QUESTIONS?