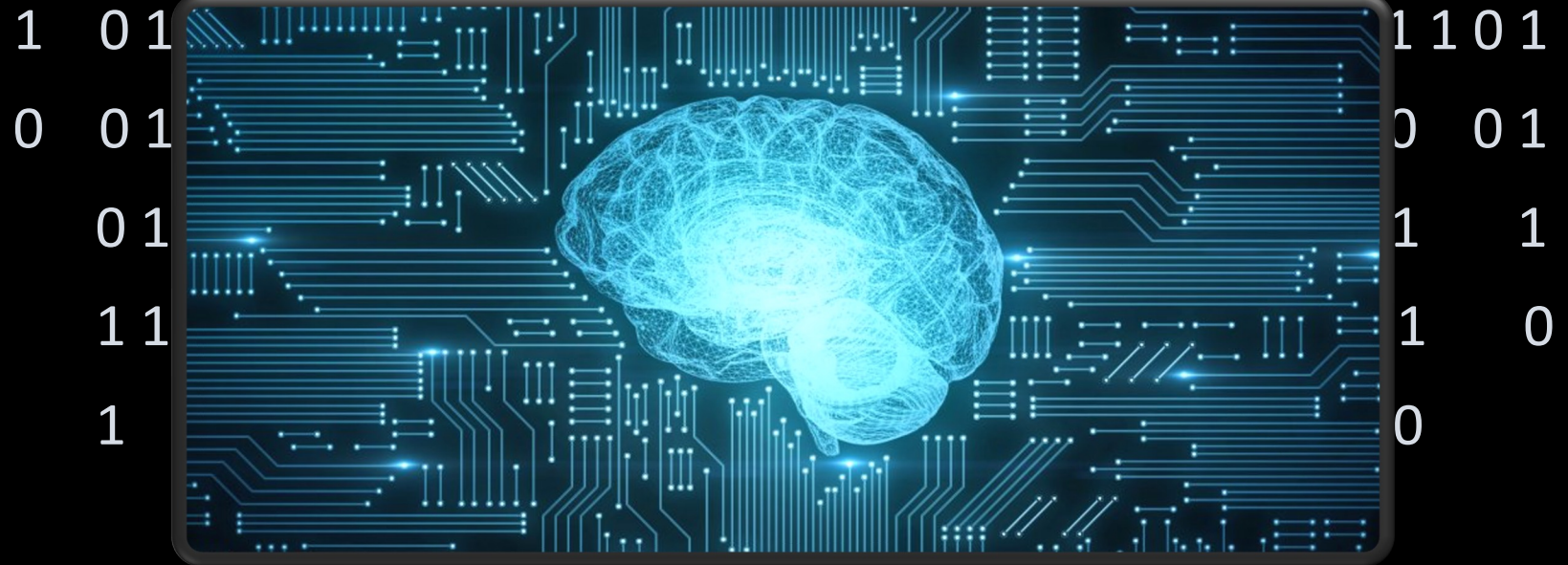# Memory as a lens to understand efficient learning and optimization



Op timization

Memory

Vatsal Sharan (USC)

01000111101011010010110110101011010
11010100010010010110111101011101010
01110110101101110011001101010010010
10101001001001101111100 10010101101
10010 00100111101101111 01101010001
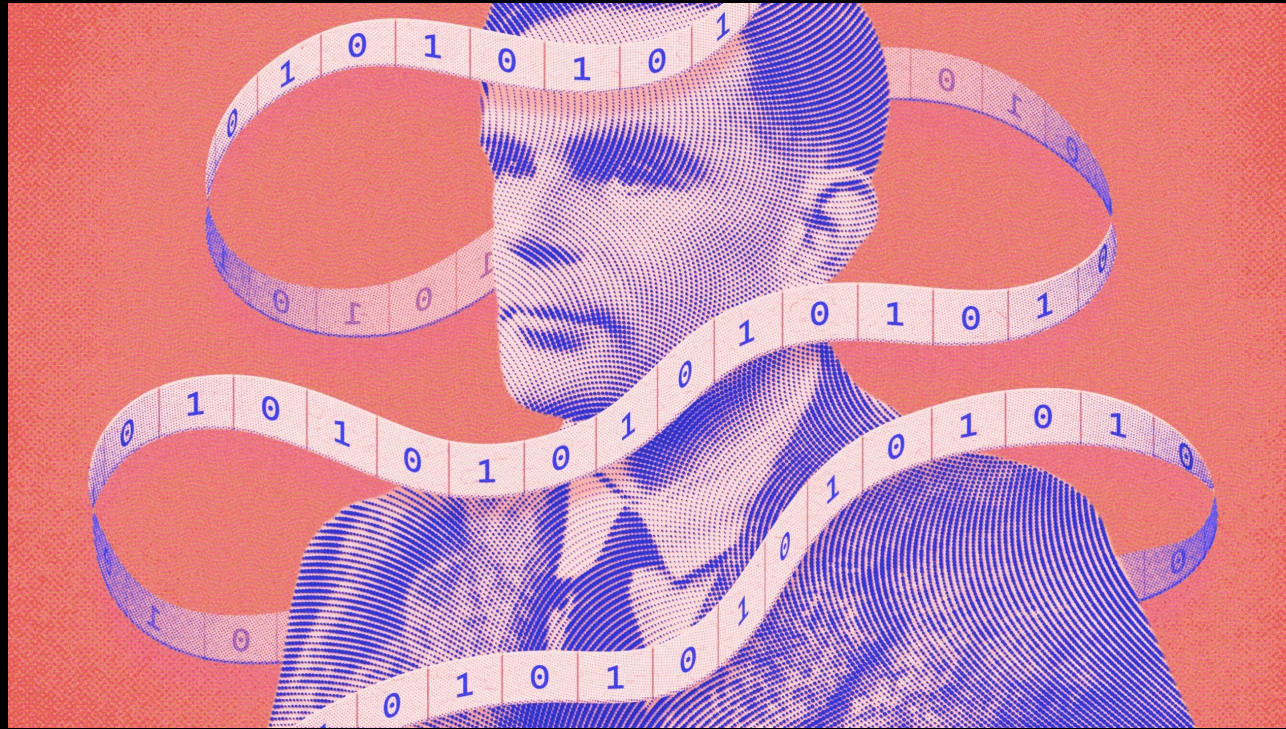01101 010101110110010 010111100100
01001 10111

1 01 1101

0 01 0 01

01 1 1

11 1 0

1 0

Machine Learning

# How do **information** and **computation** interact for **optimization?**

Optimization algorithm

# Memory as the Computational Resource



Traditionally in TCS, Memory has been a fundamental computational resource

Pic: Quanta Magazine

# Memory is a Constraint in Many Modern Practical Settings
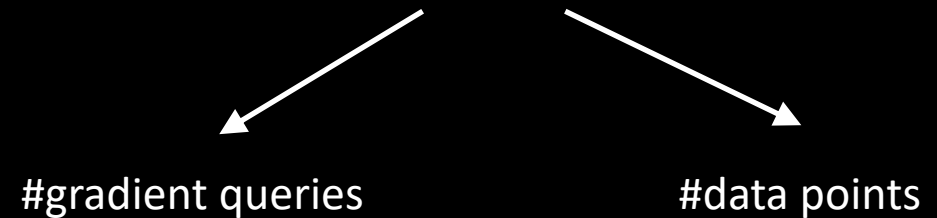


Small memory

Large models

Huge datasets

*``Memory is the dominant performance and energy bottleneck in modern computing systems; data movement is much more expensive than computation, both in latency and energy.”* [Falcao and Ferreira, CACM, 2023]

Memory is a fundamental computation resource, is crucial in practice.

What is the role of memory in learning and **optimization**?
Are there tradeoffs between available **memory** and required **information**?

#gradient queries          #data points

[This talk] Memory Dichotomy Hypothesis: It is not possible to significantly improve on the convergence rate of known memory efficient techniques without using significantly more memory.

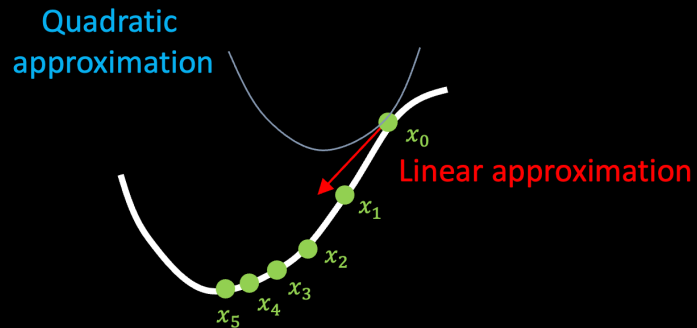**Lower bounds: Convex optimization with first-order oracle**

(with Annie Marsden, Aaron Sidford & Greg Valiant)

**Lower bounds: Convex optimization with stochastic gradient oracle**

(with Aaron Sidford & Greg Valiant)

**Upper bounds: Better convergence with small memory**

(with Jon Kelner, Annie Marsden, Aaron Sidford, Greg Valiant, Honglin Yuan)

# Lower bounds: Convex optimization with first-order oracle

Annie Marsden

Aaron Sidford

Greg Valiant

*Efficient Convex Optimization Requires Superlinear Memory,*
Annie Marsden, Vatsal Sharan, Aaron Sidford, Gregory Valiant, 2022

# A canonical optimization problem

Consider minimizing convex,
1- Lipschitz functions:

$$\text{min. } F(x)$$
$$x \in R^d : \|x\| \leq 1$$

# A canonical optimization problem

Consider minimizing convex,
1- Lipschitz functions:

$$\text{min. } F(x)$$

$$x \in R^d : \|x\| \leq 1$$

Given access to a first-order oracle:

- Algorithm queries some point $x$
- Oracle responds with $(F(x), \nabla F(x))$



$\nabla F(x)$

$F(x)$

Query point $x$

# Algorithms we know

Gradient Descent

Initialize $x_0$ . At time $t$,

Query point $x_t$

Receive gradient $\nabla F(x_t)$ at $x_t$

Update $x_{t+1} \rightarrow x_t - \eta \cdot \nabla F(x_t)$

$-\nabla F(x)$

$x_t$     $x_{t+1} \rightarrow x_t - \eta \cdot \nabla F(x_t)$

# Algorithms we know

## Gradient Descent

Initialize $x_0$. At time $t$,

Query point $x_t$

Receive gradient $\nabla F(x_t)$ at $x_t$

Update $x_{t+1} \to x_t - \eta \cdot \nabla F(x_t)$

- $O(d)$ computation time per query
- $O(d)$ memory per query
- Query complexity large with respect to desired error $\epsilon$: need $\epsilon^{-2}$ queries to find $\epsilon$ optimal answer
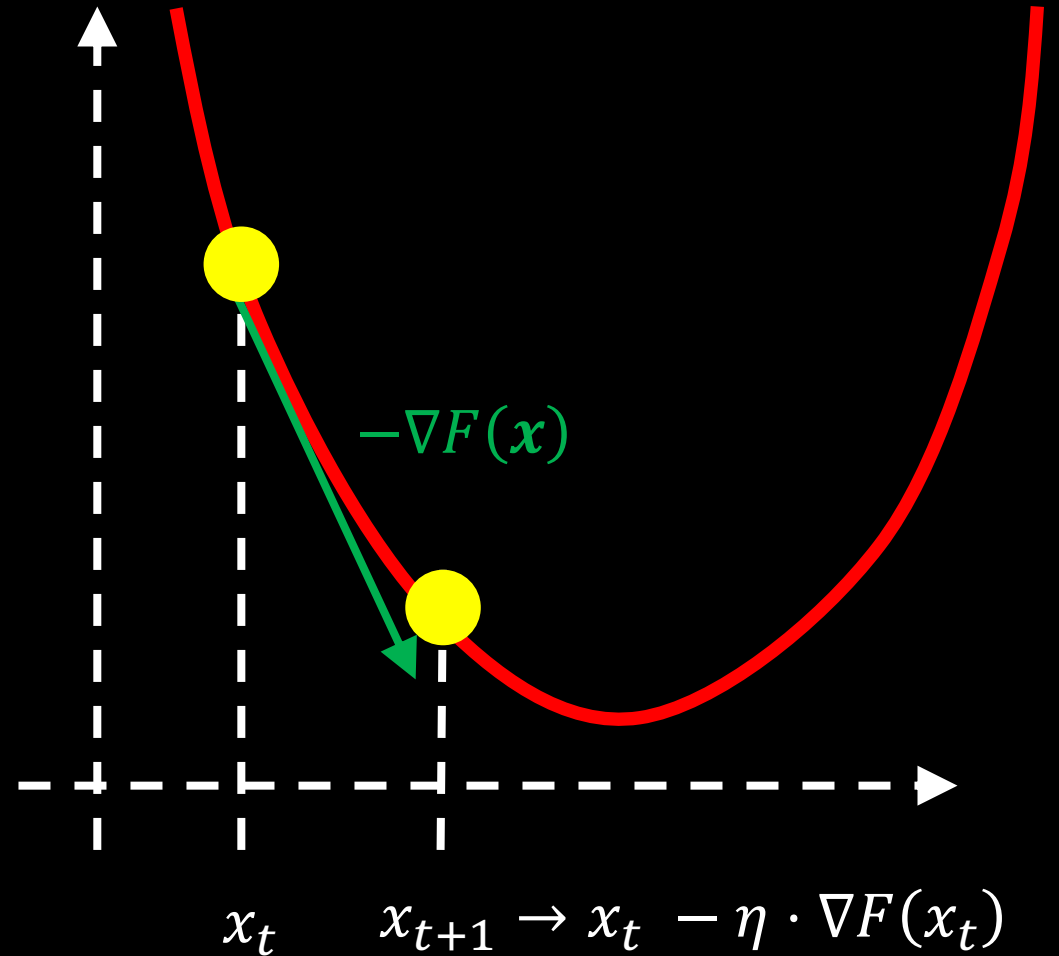
# Algorithms we know

## Gradient Descent

Initialize $x_0$. At time $t$,
Query point $x_t$
Receive gradient $\nabla F(x_t)$ at $x_t$
Update $x_{t+1} \to x_t - \eta \cdot \nabla F(x_t)$

## Suite of other techniques

- Based on the ellipsoid algorithm
- Does something like high-dimensional binary search

- $O(d)$ computation time per query
- $O(d)$ memory per query
- Query complexity large with respect to desired error $\epsilon$: need $\epsilon^{-2}$ queries to find $\epsilon$ optimal answer

- $> d^2$ computation time per query
- $> d^2$ memory per query
- Query complexity small with respect to desired error $\epsilon$: need $d \log \left( \frac{1}{\epsilon} \right)$ queries to find $\epsilon$ optimal answer

# Algorithms we know

**Gradient Descent**

**Suite of other techniques**

Are complex algorithms necessary?

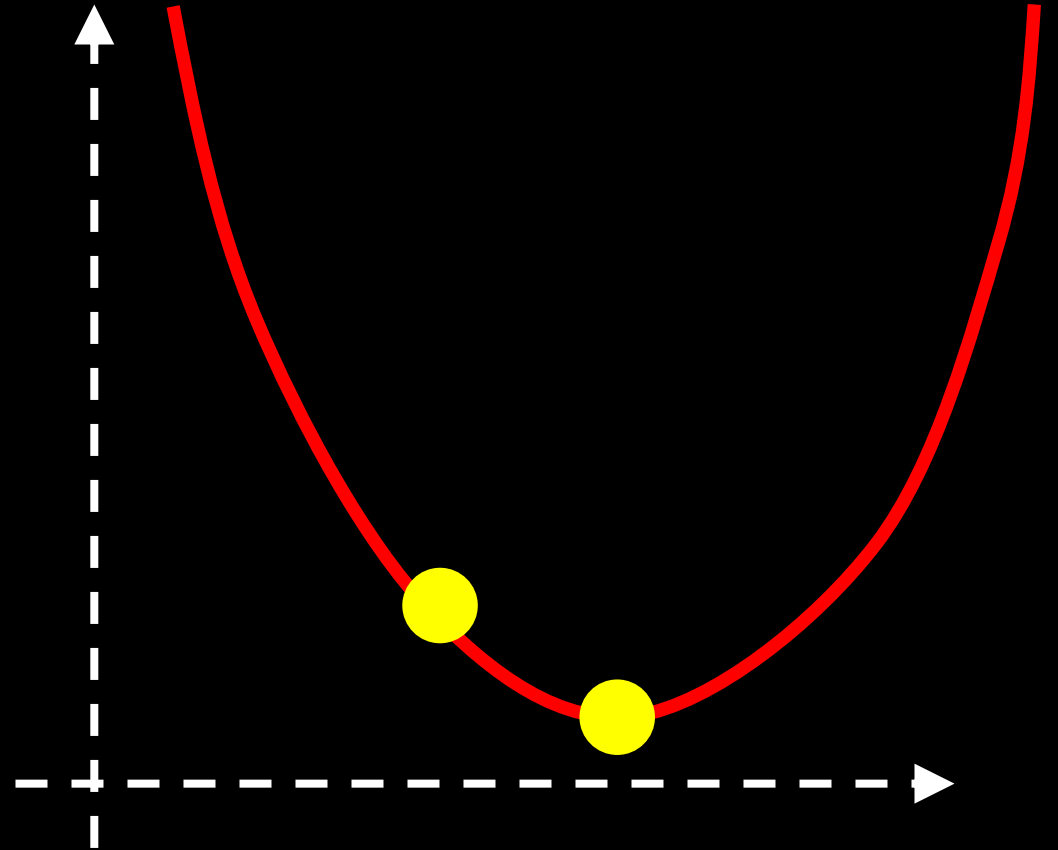# Algorithms we know

## Gradient Descent

Initialize $x_0$ . At time $t$,
Query point $x_t$
Receive gradient $\nabla F(x_t)$ at $x_t$
Update $x_{t+1} \rightarrow x_t - \eta \cdot \nabla F(x_t)$

- $O(d)$ computation time per query
- $O(d)$ memory per query
- Query complexity large with respect to desired error $\epsilon$: need $\epsilon^{-2}$ queries to find $\epsilon$ optimal answer
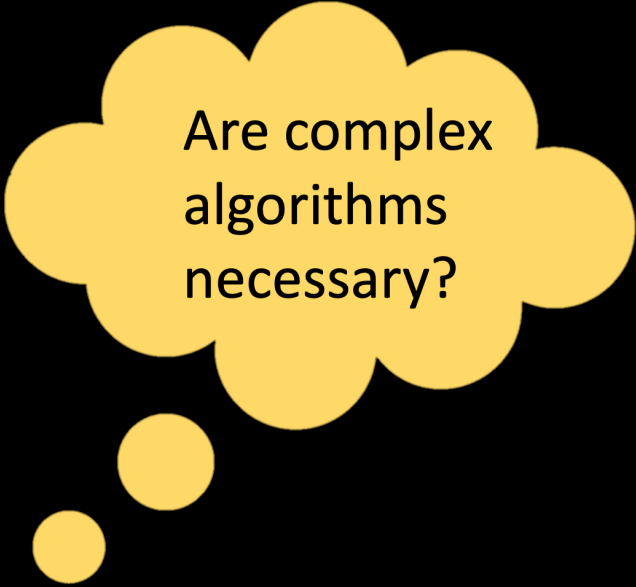
## Suite of other techniques

- Based on the ellipsoid algorithm
- Does something like high-dimensional binary search

- $> d^2$ computation time per query
- $> d^2$ memory per query
- Query complexity small with respect to desired error $\epsilon$: need $d \log\left(\frac{1}{\epsilon}\right)$ queries to find $\epsilon$ optimal answer

# Algorithms we know

## Gradient Descent

Initialize $x_0$. At time $t$,
Query point $x_t$
Receive gradient $\nabla F(x_t)$ at $x_t$
Update $x_{t+1} \to x_t - \eta \cdot \nabla F(x_t)$

- $O(d)$ computation time per query
- $O(d)$ memory per query
- Query complexity large with respect to desired error $\epsilon$: need $\epsilon^{-2}$ queries to find $\epsilon$ optimal answer

## Suite of other techniques

- Based on the ellipsoid algorithm
- Does something like high-dimensional binary search

- $> d^2$ computation time per query
- $> d^2$ memory per query
- Query complexity small with respect to desired error $\epsilon$: need $d \log\left(\frac{1}{\epsilon}\right)$ queries to find $\epsilon$ optimal answer

# Algorithms we know

## Gradient Descent

Initialize $x_0$. At time $t$,
Query point $x_t$
Receive gradient $\nabla F(x_t)$ at $x_t$
Update $x_{t+1} \to x_t - \eta \cdot \nabla F(x_t)$

- $O(d)$ computation time per query
- $O(d)$ memory per query
- Query complexity large with respect to desired error $\epsilon$: need $\epsilon^{-2}$ queries to find $\epsilon$ optimal answer

## Suite of other techniques

- Based on the ellipsoid algorithm
- Does something like high-dimensional binary search

- $> d^2$ computation time per query
- $> d^2$ memory per query
- Query complexity small with respect to desired error $\epsilon$: need $d \log \left( \frac{1}{\epsilon} \right)$ queries to find $\epsilon$ optimal answer

# Algorithms we know

## Gradient Descent

Initialize $x_0$. At time $t$,
Query point $x_t$
Receive gradient $\nabla F(x_t)$ at $x_t$
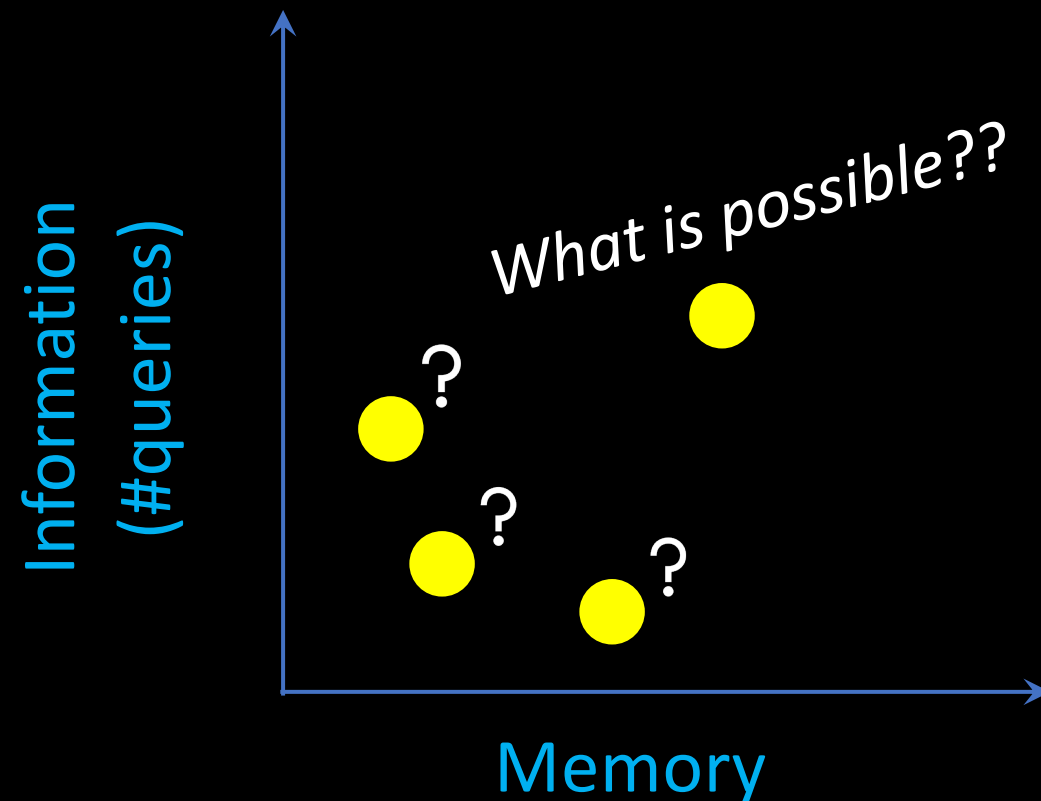Update $x_{t+1} \rightarrow x_t - \eta \cdot \nabla F(x_t)$

- $O(d)$ computation time per query
- $O(d)$ memory per query
- Query complexity large with respect to desired error $\epsilon$: need $\epsilon^{-2}$ queries to find $\epsilon$ optimal answer
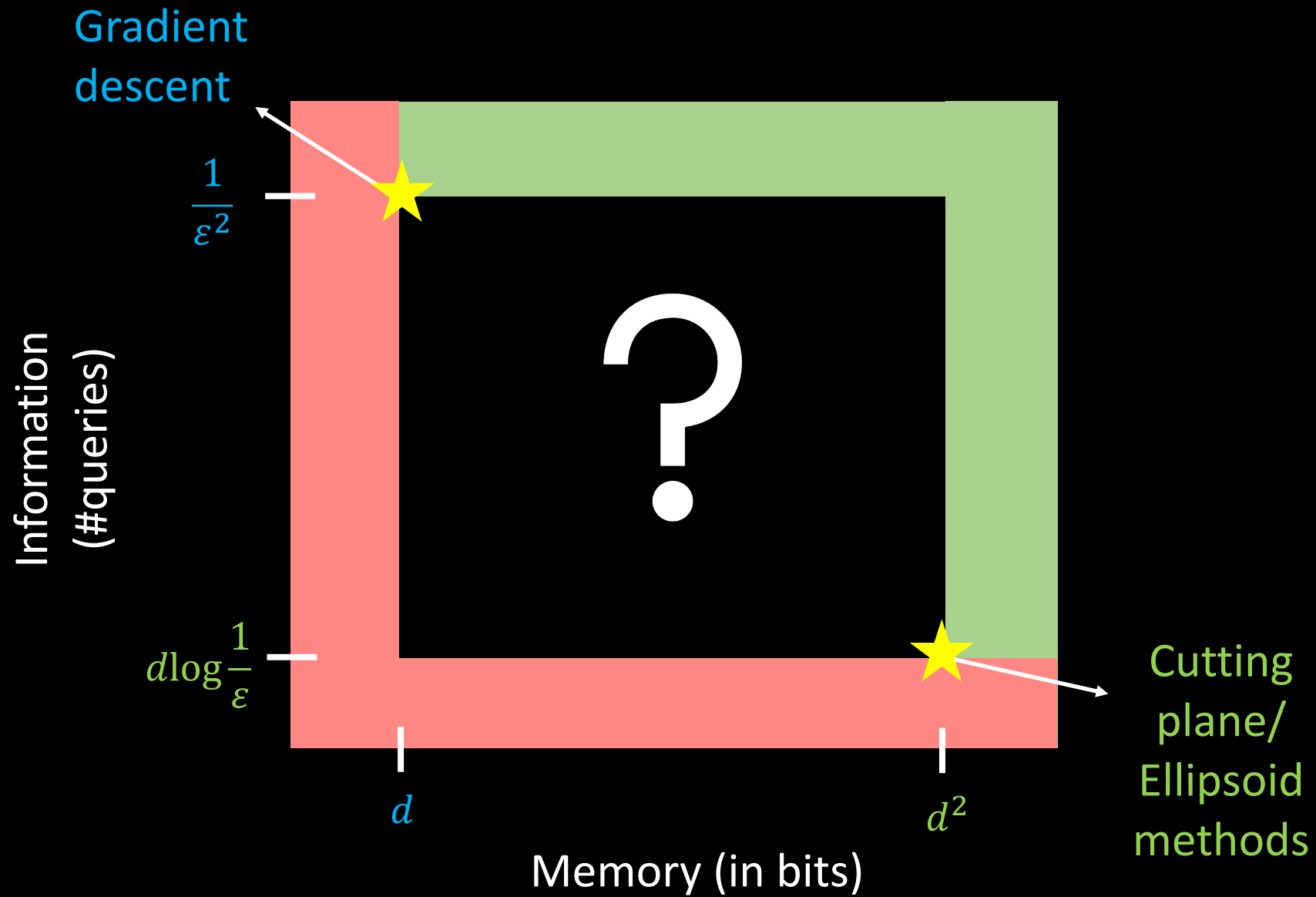
## Suite of other techniques

- Based on the ellipsoid algorithm
- Does something like high-dimensional binary search

- $> d^2$ computation time per query
- $> d^2$ memory per query
- Query complexity small with respect to desired error $\epsilon$: need $d \log \left( \frac{1}{\epsilon} \right)$ queries to find $\epsilon$ optimal answer

# Are there inherent **tradeoffs** between available **memory** and **information** requirement?

Gradient descent

$\dfrac{1}{\varepsilon^2}$

Information (#queries)

$d\log\dfrac{1}{\varepsilon}$

$d$

$d^2$

Memory (in bits)

Cutting plane/ Ellipsoid methods

# What is known?

**Info-theoretic bounds for optimization algorithms**

Nemirovski-Yudin'83,
Shamir'13,
Nesterov'14,
Bubeck'15,
Duchi-Jordan-
Wainwright-Wibisono'15,
Woodworth-Srebro'16,
Carmon-Duchi-Hinder-
Sidford'17ab,
Arjevani-Shamir'17,
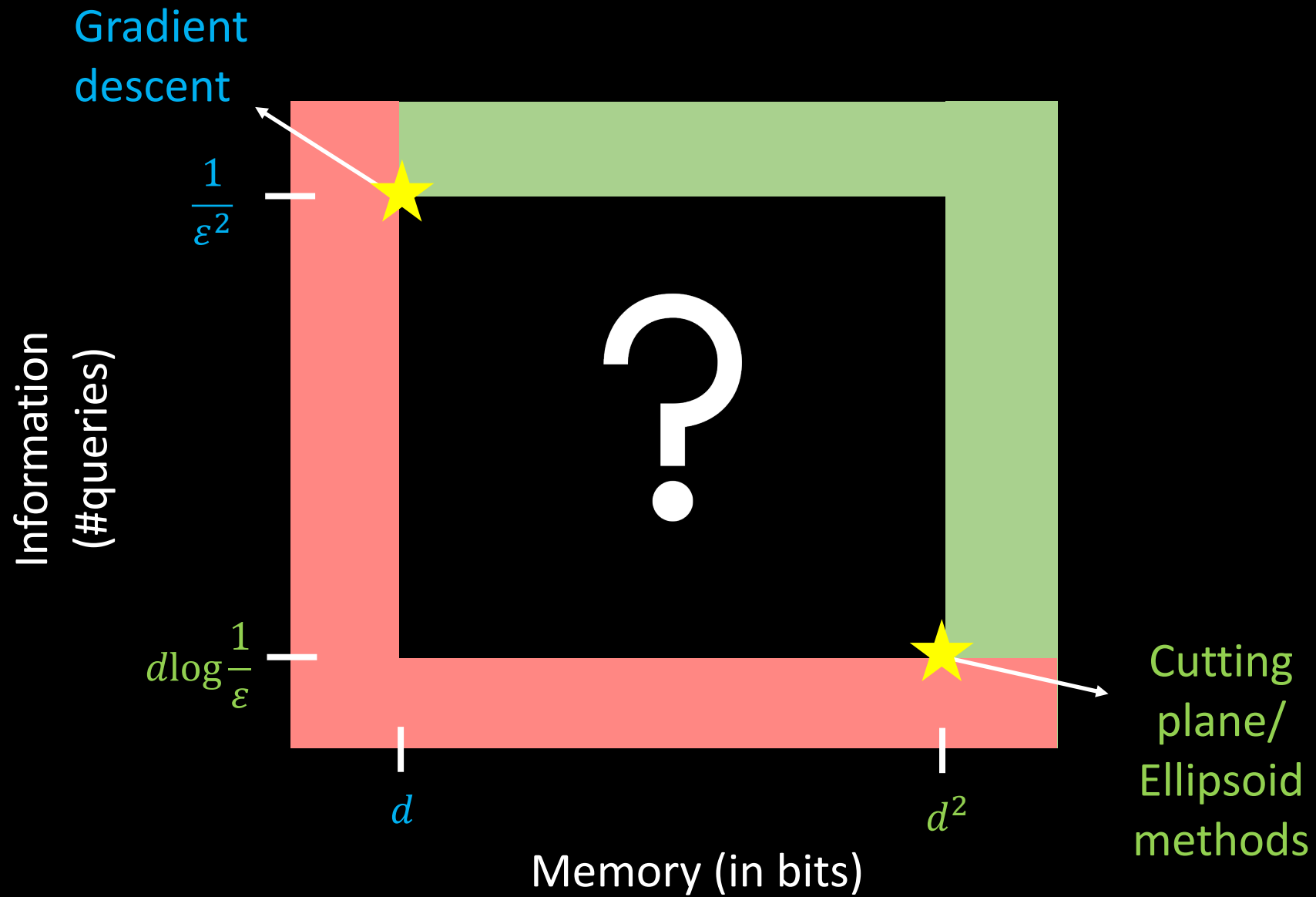Agarwal-Hazan'18,
Diakonikolas-Guzman'19

**Memory bounds for streaming data**

Alon-Matias-Szegedy'99,
Indyk-Woodruff'03
Bar-Yossef-Jayaram-Kumar-
Sivakumar'04,
Nelson-Le Huy'13,
Steinhardt-Duchi'15,
Braverman-Garg-Ma-Nguyen-
Woodruff'16,
Kapralov-Nelson-Pachocki-
Wang-Woodruff-Yahyazadeh'17,
Nelson-Yu'19,
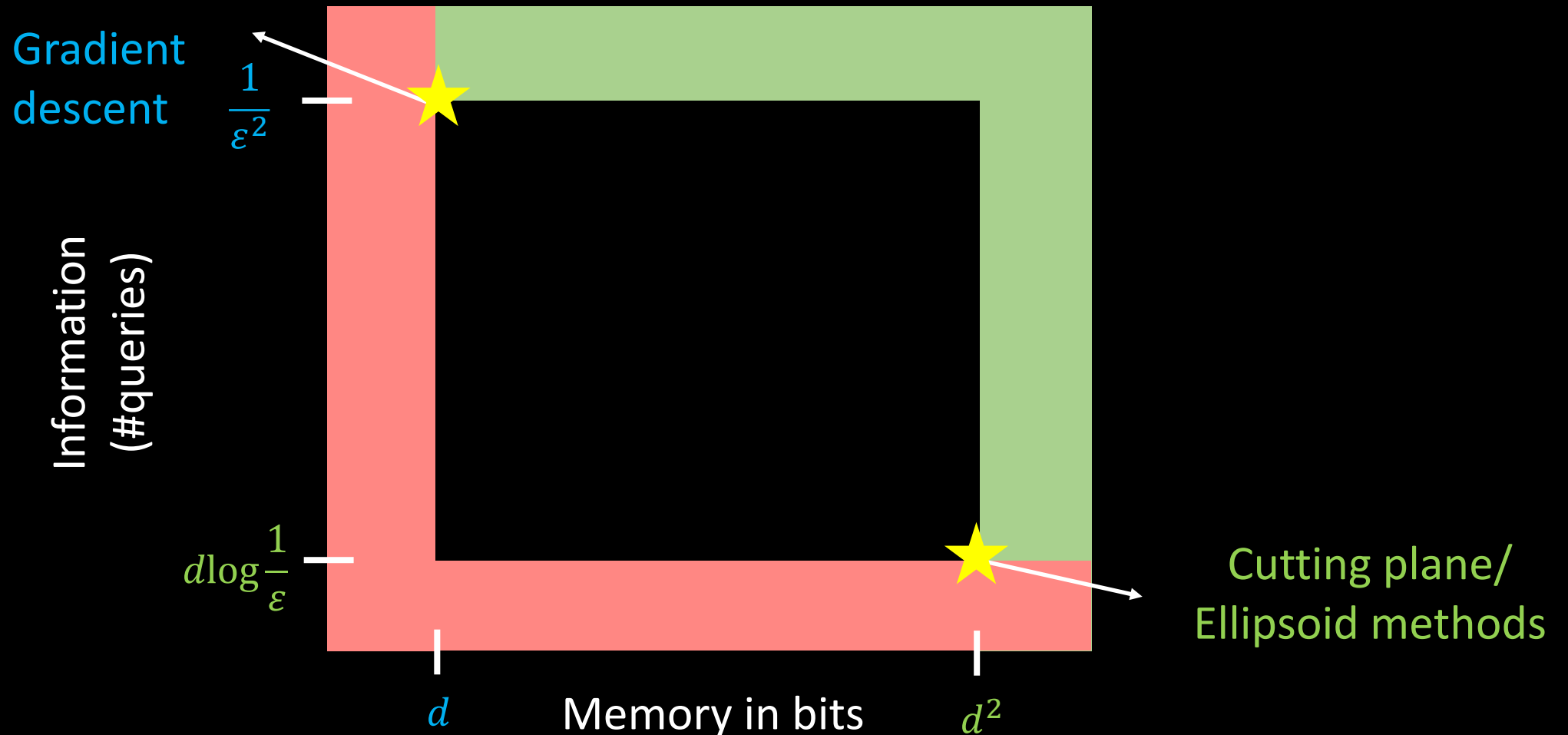Dagan-Kur-Shamir'19

**Memory bounds over finite fields**

Shamir'14,
Steinhardt-Valiant-Wager'16,
Raz'17,
Moshkovitz-Moshkovitz'17
Kol-Raz-Tal'17,
Moshkovitz-Moshkovitz'18,
Garg-Raz-Tal'18,
Beame-Oveis Gharan-Yang'18,
Garg-Raz-Tal'19,
Raz-Zhan'20,
Gonen-Lovett-Moshkovitz'20,
Garg-Kothari-Raz'20

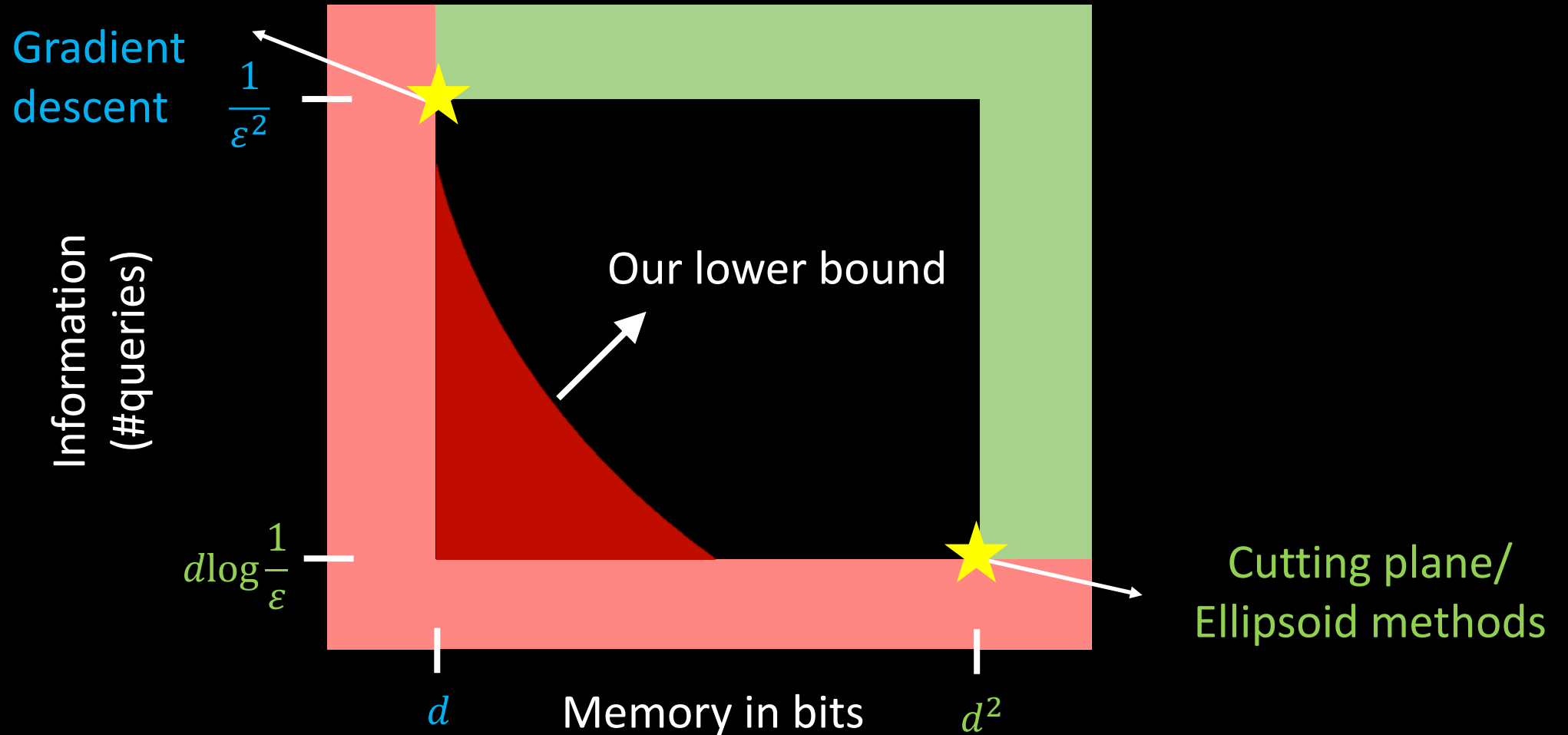**Memory bounds for continuous optimization**

Theorem [Marsden, Sharan, Sidford, Valiant]:
For $\epsilon \geq \frac{1}{\text{poly}(d)}$ and $\delta \in [0, 0.25]$, any (randomized) algorithm with memory $\boldsymbol{d^{1.25-\delta}}$ requires at least $\boldsymbol{d^{1+1.33\delta}}$ first-order queries to find $\epsilon$-optimal point.

Gradient descent

$\frac{1}{\varepsilon^2}$

Information (#queries)

$d\log\frac{1}{\varepsilon}$

Cutting plane/ Ellipsoid methods

$d$

Memory in bits

$d^2$

Theorem [Marsden, Sharan, Sidford, Valiant]:
For $\epsilon \geq \frac{1}{\text{poly}(d)}$ and $\delta \in [0, 0.25]$, any (randomized) algorithm with memory $\boldsymbol{d^{1.25-\delta}}$
requires at least $\boldsymbol{d^{1+1.33\delta}}$ first-order queries to find $\epsilon$-optimal point.



Gradient descent $\frac{1}{\varepsilon^2}$

Information (#queries)

Our lower bound

$d \log \frac{1}{\varepsilon}$

Cutting plane/ Ellipsoid methods

$d$ Memory in bits $d^2$

# High-level proof

## Step one

Construct a distribution over functions that seems hard to optimize with limited memory

## Step two

Relate optimizing these functions to winning a communication game

## Step three

For the communication game, prove a memory/query tradeoff

# Hard distribution over functions

Construct a distribution over functions that seems hard to optimize with limited memory

$$F_{h,A,\eta,\rho}(x) = \max\{\eta\|Ax\|_\infty - \rho, h(x)\}$$

$$A \sim \text{Unif}(\{\pm 1\}^{\frac{d}{2}\times d}) \qquad h(x) = \max_{i\in[N]} v_i^T x - i\gamma \quad \text{(variant of Nemirovski function)}$$

To receive first order information about $h$, must make query which is reasonably orthogonal to $A$

Nemirovski property: To continue receiving new or informative subgradients, queries must be robustly linearly independent

# From optimization to winning a game

Relate optimizing these functions to winning a communication game

$$F_{h,A,\eta,\rho}(x) = \max\{\eta\|Ax\|_{\infty} - \rho, h(x)\}$$

$$h(x) = \max_{i\in[N]} v_i^T x - i\gamma \quad \text{(variant of the Nemirovski function)}$$

Relating optimizing $F_{h,A}(x)$ to winning an **Orthogonal Vector Game**

# The Orthogonal Vector Game



Random matrix
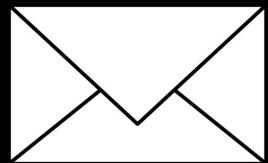
$A \sim \mathrm{Unif}(\{\pm 1\}^{\frac{d}{2} \times d})$
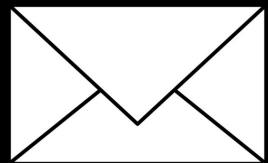
$M$-bit message

Player

Game oracle

To win, find $y_1, y_2, \dots, y_k$ which are roughly orthogonal*  to $A$

$d$: dimension        $k$: #vectors to be returned        $m$: #oracle queries      $M$: size of message (in bits)

# The Orthogonal Vector Game

Query $x_1$

$g_1 \in \partial \|Ax_1\|_\infty$

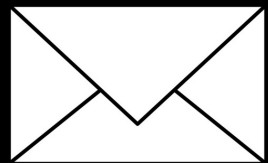$M$-bit message

Player

Game oracle

To win, find $y_1, y_2, \ldots, y_k$ which are roughly orthogonal*  to $A$

$d$: dimension      $k$: #vectors to be returned    $m$: #oracle queries   $M$: size of message (in bits)

# The Orthogonal Vector Game

$(x_1, g_1)$

$M$-bit message

Player

Game oracle

To win, find $y_1, y_2, \ldots, y_k$ which are roughly orthogonal*  to $A$

$d$: dimension        $k$: #vectors to be returned      $m$: #oracle queries    $M$: size of message (in bits)

# The Orthogonal Vector Game

Query $x_2$

$g_2 \in \partial \|Ax_2\|_\infty$

$(x_1, g_1)$

$M$-bit message

Player

Game oracle

To win, find $y_1, y_2, \ldots, y_k$ which are roughly orthogonal*  to $A$

$d$: dimension    $k$: #vectors to be returned    $m$: #oracle queries   $M$: size of message (in bits)

# The Orthogonal Vector Game



$(x_1, g_1)$
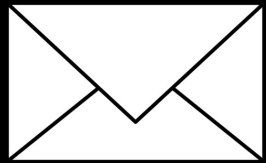$(x_2, g_2)$

$M$-bit message

Player

Game oracle

To win, find $y_1, y_2, \ldots, y_k$ which are roughly orthogonal*  to $A$

$d$: dimension      $k$: #vectors to be returned    $m$: #oracle queries   $M$: size of message (in bits)

# The Orthogonal Vector Game

# The Orthogonal Vector Game

$(x_1, g_1)$

$(x_2, g_2)$

.
.
.

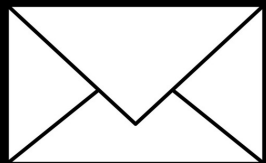$(x_m, g_m)$

$M$-bit message
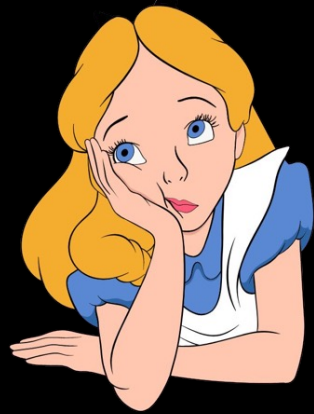
Player

Game oracle

To win, find $y_1, y_2, \ldots, y_k$ which are roughly orthogonal*  to $A$

$d$ : dimension        $k$ : #vectors to be returned      $m$ : #oracle queries    $M$ : size of message (in bits)

# The Orthogonal Vector Game

$(x_1, g_1)$

$(x_2, g_2)$

.
.
.

$(x_m, g_m)$

$M$-bit message

Player

To win, find $y_1, y_2, \ldots, y_k$ which are roughly orthogonal*  to $A$

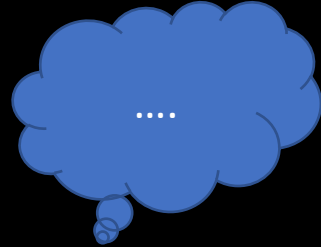$d$: dimension        $k$: #vectors to be returned      $m$: #oracle queries    $M$: size of message (in bits)

# The Orthogonal Vector Game

$(x_1, g_1)$
$(x_2, g_2)$
.
.
.
.
$(x_m, g_m)$

$M$-bit message

Player

To win, find $y_1, y_2, \ldots, y_k$ which are roughly orthogonal*  to $A$

$d$: dimension      $k$: #vectors to be returned    $m$: #oracle queries    $M$: size of message (in bits)

# The Orthogonal Vector Game

$$y_1, y_2, \ldots, y_k$$

Player

To win, $y_1, y_2, \ldots, y_k$ must be:
- Roughly orthogonal to $A$: $\|Ay_i\|_\infty \leq 1/d^4$
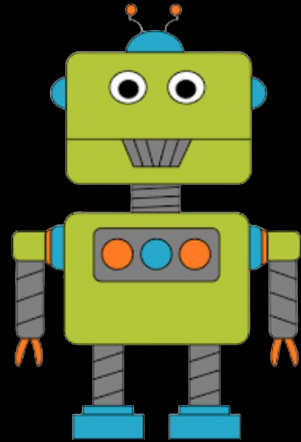- Robustly linearly independent
  $\left\|\mathrm{Proj}_{\mathrm{span}(y_1, \ldots, y_{i-1})}(y_i)\right\| \leq 1 - 1/d^2$

To win, find $y_1, y_2, \ldots, y_k$ which are roughly orthogonal*  to $A$

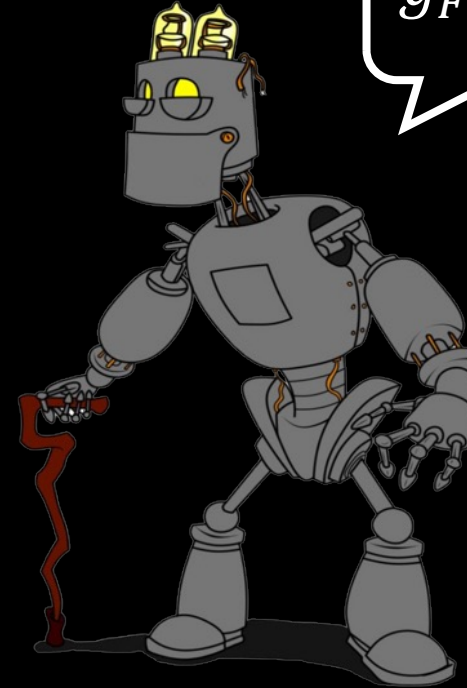$d$: dimension       $k$: #vectors to be returned     $m$: #oracle queries    $M$: size of message (in bits)

# From Optimization to winning the Game

Query $x$

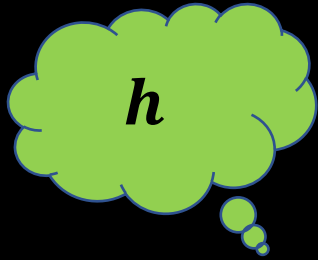$g_{F_{h,A}}(x)$

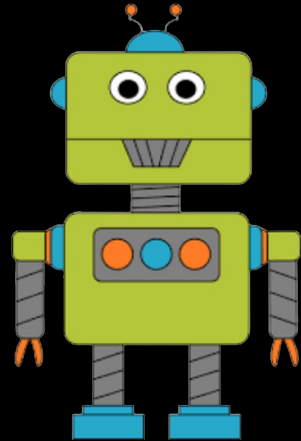Optimization algorithm

Optimization oracle

$M$-bit memory state

$$F_{h,A}(x) = \max\{\eta\|Ax\|_\infty - \rho, h(x)\}$$

# From Optimization to winning the Game

$h$

Generates Nemirovski function $h$
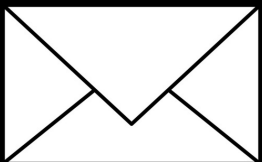Wants to optimize $F_{h,A}$

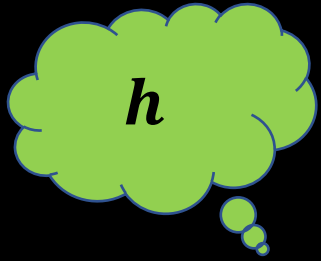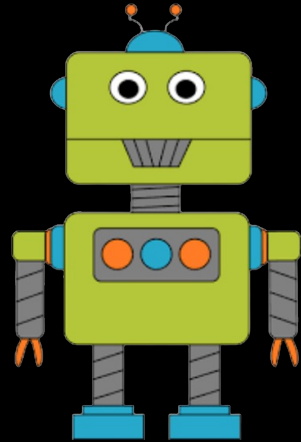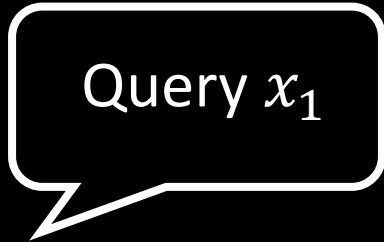Optimization algorithm

$A$

Random matrix

Game oracle

$M$-bit memory state

$$F_{h,A}(x) = \max\{\eta\|Ax\|_\infty - \rho, h(x)\}$$

# From Optimization to winning the Game

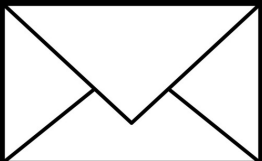# From Optimization to winning the Game

# From Optimization to winning the Game

# Memory/Query tradeoffs for the Game

Step three

For the communication game, prove a memory/query tradeoff

If available memory $< kd$, then Player must make $\approx d$ queries

To win, find $y_1, y_2, \ldots, y_k$ which are roughly orthogonal*  to $A$

$d$ : dimension      $k$ : #vectors to be returned     $m$ : #oracle queries    $M$ : size of message (in bits)

# Open Questions

- Randomized algorithms, for poly-small $\epsilon$?
- What happens for **smooth functions**?
- Can you improve on the poly$(1/\epsilon)$ rate of gradient descent for **super-poly small $\epsilon$**?

Conjecture: Cannot improve gradient descent's convergence rate without using quadratic memory.

[This talk] Memory Dichotomy Hypothesis: It is not possible to significantly improve on the convergence rate of known memory efficient techniques without using significantly more memory.

# Lower bounds: Convex optimization with stochastic gradient oracle

(with Aaron Sidford & Greg Valiant)

Upper bounds: Better convergence with small memory

*Memory-Sample Tradeoffs for Linear Regression with Small Error*

Vatsal Sharan, Aaron Sidford, Gregory Valiant, 2019

Quadratic approximation

Linear approximation

$x_0$
$x_1$
$x_2$
$x_3$
$x_4$
$x_5$

1st order vs. 2nd order methods

Information (#queries)

$\frac{1}{\varepsilon^2}$

$d\log\frac{1}{\varepsilon}$

Computation (memory in bits)

$d$     $d^2$

# Stochastic optimization

In many modern ML settings, we work with stochastic gradients $g(x)$:

$$\text{min. } F(x)$$

$$x \in R^d : \|x\| \leq 1$$

$$\mathrm{E}[g(x)] = \nabla F(x)$$

If $F(x)$ is expected loss with respect to data points sampled from some distribution, we can find stochastic gradient using a randomly sampled labelled datapoint.

What is the tradeoff between available memory and number of samples needed to optimize?

# **Linear model**: Data vs. Memory?

$$\begin{bmatrix} A \end{bmatrix} \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} b \end{bmatrix}$$

# Find $x$

$$x, a_i \in R^d$$
$$b_i \in R$$

$$\begin{bmatrix} 3 & 5 \end{bmatrix} \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} 4 \end{bmatrix}$$

$a_1$       $x$    $b_1$

$$<a_1, x> = b_1$$

# Find $x$

$$\begin{bmatrix} -2 & 2 \end{bmatrix} \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} 8 \end{bmatrix}$$

$a_2 \qquad\qquad\qquad b_2$

$$x, a_i \in R^d$$
$$b_i \in R$$

$$<a_2, x> = b_2$$

# Find $x$

$$x, a_i \in R^d$$
$$b_i \in R$$

$$\begin{bmatrix} -7 & -2 \end{bmatrix} \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} 10 \end{bmatrix}$$

$a_3$ $\qquad b_3$

$$<a_3, x> = b_3$$

Find $x$

$$\begin{bmatrix} x, a_i & \in & R^d \\ b_i & \in & R \end{bmatrix}$$

$$\begin{bmatrix} 4 & -1 \end{bmatrix} \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} -10 \end{bmatrix}$$

$$a_4 \qquad\qquad\qquad b_4$$

$$< a_4, x > = b_4$$

# Find $x$

$$x, a_i \in R^d$$
$$b_i \in R$$

$$\begin{bmatrix} -8 & -6 \end{bmatrix} \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} 4 \end{bmatrix}$$

$a_5$ $\qquad$ $b_5$

$$<a_5, x> = b_5$$

# Find $x$

$$x, a_i \in R^d$$
$$b_i \in R$$

$$\begin{bmatrix} -8 & -6 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} 4 \end{bmatrix}$$

$$a_5 \qquad\qquad b_5$$

1. Memory = #bits
2. Samples drawn from Gaussian

$$<a_5, x> = b_5$$

# What can you do?

Gaussian Elimination

$$\begin{bmatrix} 4 & -1 \\ -8 & -6 \end{bmatrix} \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} -10 \\ 4 \end{bmatrix}$$

# What can you do?

Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ -8 & -6 \end{bmatrix} \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} -5/2 \\ 4 \end{bmatrix}$$

# What can you do?

Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & -8 \end{bmatrix} \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} -5/2 \\ -16 \end{bmatrix}$$

# What can you do?

Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

# What can you do?

Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx d^2$ memory

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx d^2$ memory

## Gradient Descent

Initialize $x_0$. At time $i$,

Get $(a_i, b_i)$. Update $x_i \rightarrow x_{i+1}$.

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx d^2$ memory

## Gradient Descent

Initialize $x_0$. At time $i$,

Get $(a_i, b_i)$. Update $x_i \rightarrow x_{i+1}$.

$$x_0 = (-0.25, 0.98)$$

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx$ d$^2$ memory

## Gradient Descent

Initialize $x_0$. At time $i$,

Get $(a_i, b_i)$. Update $x_i \rightarrow x_{i+1}$.

$$x_0 = (-0.25, 0.98)$$

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx$ d$^2$ memory

## Gradient Descent

Initialize $x_0$ . At time $i$,

Get $(a_i, b_i)$ . Update $x_i \rightarrow x_{i+1}$ .

$$x_1 = (-0.45, 0.74\ )$$

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx d^2$ memory

## Gradient Descent

Initialize $x_0$. At time $i$,

Get $(a_i, b_i)$. Update $x_i \rightarrow x_{i+1}$.

$$x_2 = (-0.74, 2.24\,)$$

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx$ d$^2$ memory

## Gradient Descent

Initialize $x_0$. At time $i$,

Get $(a_i, b_i)$. Update $x_i \rightarrow x_{i+1}$.

$$x_3 = (-1.64, 2.70)$$

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx d^2$ memory

## Gradient Descent

Initialize $x_0$. At time $i$,

Get $(a_i, b_i)$. Update $x_i \rightarrow x_{i+1}$.

$$x_4 = (-1.85, 2.74)$$

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx$ d$^2$ memory

## Gradient Descent

Initialize $x_0$. At time $i$,

Get $(a_i, b_i)$. Update $x_i \rightarrow x_{i+1}$.

$$x_5 = (-2.27, 2.53)$$

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx$ d$^2$ memory

## Gradient Descent

Initialize $x_0$. At time $i$,

Get $(a_i, b_i)$. Update $x_i \rightarrow x_{i+1}$.

$$x_6 = (-1.99, 2.52)$$

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx$ d$^2$ memory

## Gradient Descent

Initialize $x_0$ . At time $i$,

Get $(a_i, b_i)$ . Update $x_i \rightarrow x_{i+1}$ .

$$x_7 = (-1.83, 2.47)$$

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx$ d$^2$ memory

## Gradient Descent

Initialize $x_0$. At time $i$,

Get $(a_i, b_i)$. Update $x_i \rightarrow x_{i+1}$.

$$x_8 = (-1.92, 2.48)$$

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx d^2$ memory

## Gradient Descent

Initialize $x_0$. At time $i$,

Get $(a_i, b_i)$. Update $x_i \rightarrow x_{i+1}$.

$$x_9 = (-2.20, 2.17)$$

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx$ d$^2$ memory

## Gradient Descent

Initialize $x_0$ . At time $i$,

Get $(a_i, b_i)$ . Update $x_i \rightarrow x_{i+1}$ .

$$x_{10} = (-1.97, 2.08 )$$

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx d^2$ memory

## Gradient Descent

Initialize $x_0$. At time $i$,

Get $(a_i, b_i)$. Update $x_i \rightarrow x_{i+1}$.

$$x_{11} = (-2.02, 2.01)$$

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx$ d$^2$ memory

## Gradient Descent

Initialize $x_0$ . At time $i$,

Get $(a_i, b_i)$ . Update $x_i \rightarrow x_{i+1}$ .

$$x_{12} = (-2.01, 2.00\ )$$

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx d^2$ memory

## Gradient Descent

Initialize $x_0$. At time $i$,

Get $(a_i, b_i)$. Update $x_i \rightarrow x_{i+1}$.

$$x_{12} = (-2.01, 2.00)$$

$O\left(d \log \frac{1}{\varepsilon}\right)$ examples

$\approx d$ memory

# What can you do?

## Gaussian Elimination

$$\begin{bmatrix} 1 & -1/4 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 2 \end{bmatrix} = \begin{bmatrix} -5/2 \\ 2 \end{bmatrix}$$

d examples

$\approx d^2$ memory

## Gradient Descent

Initialize $x_0$. At time $i$,

Get $(a_i, b_i)$. Update $x_i \rightarrow x_{i+1}$.

$$x_{12} = (-2.01, 2.00)$$

> d examples

$\approx d$ memory

# Gaussian Elimination
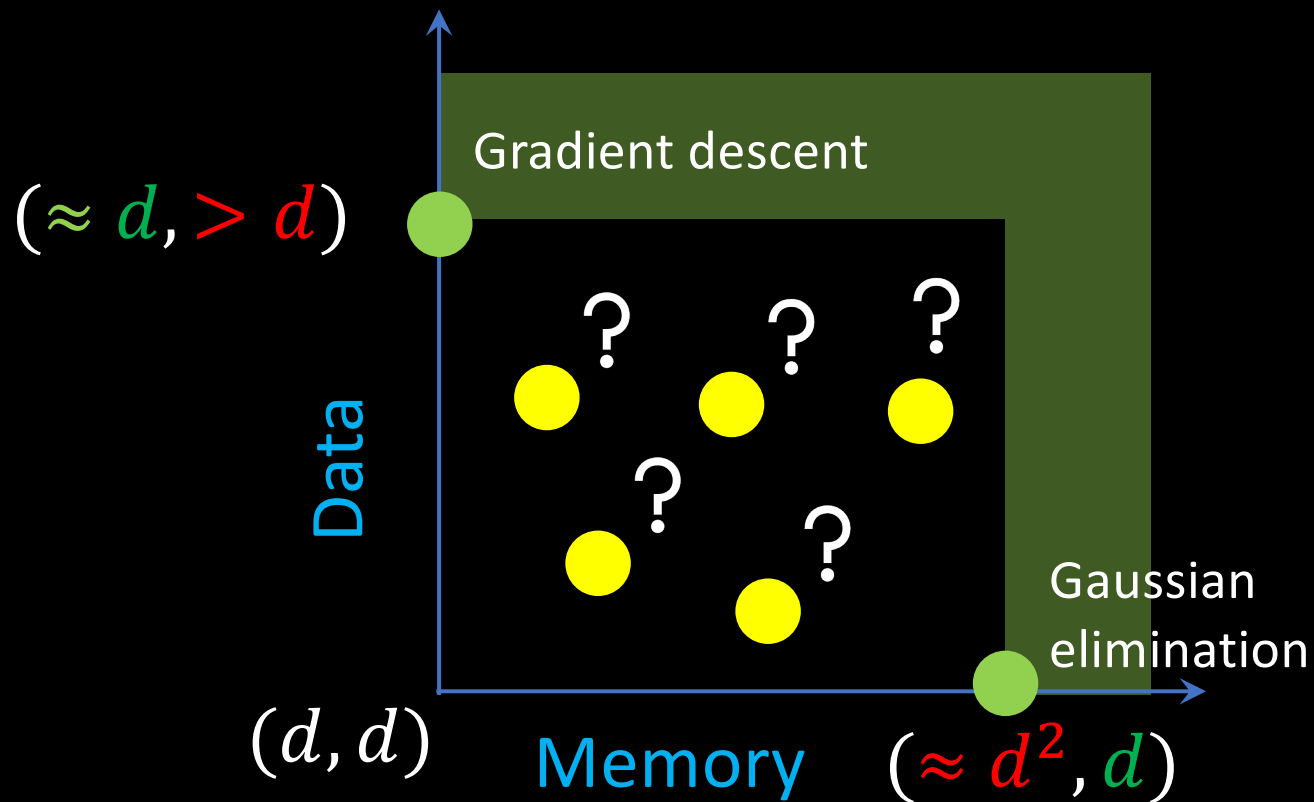
d examples

$\approx d^2$ memory

# Gradient Descent

> d examples

$\approx$ d memory

# Gaussian Elimination
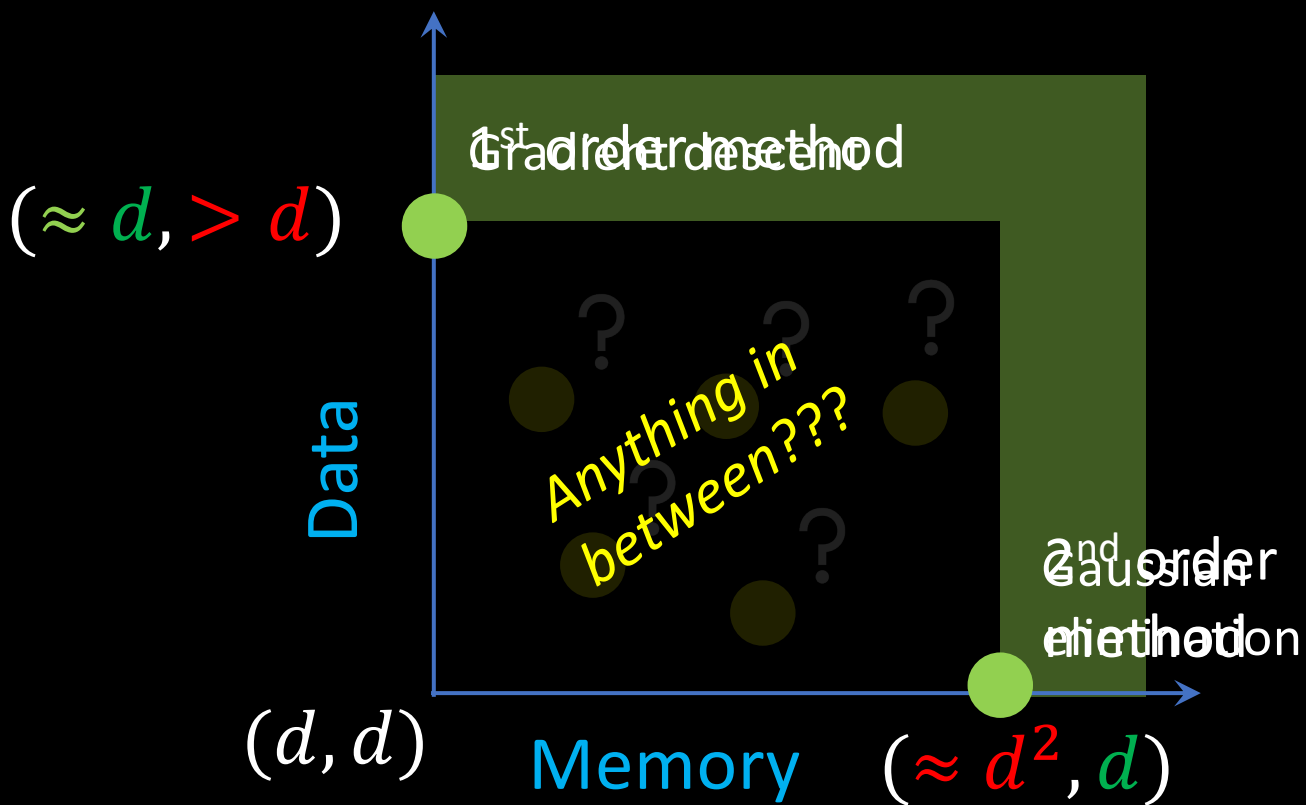
d examples

$\approx d^2$ memory

# Gradient Descent

$> d$ examples

$\approx d$ memory

# Gaussian Elimination

d examples

$\approx d^2$ memory
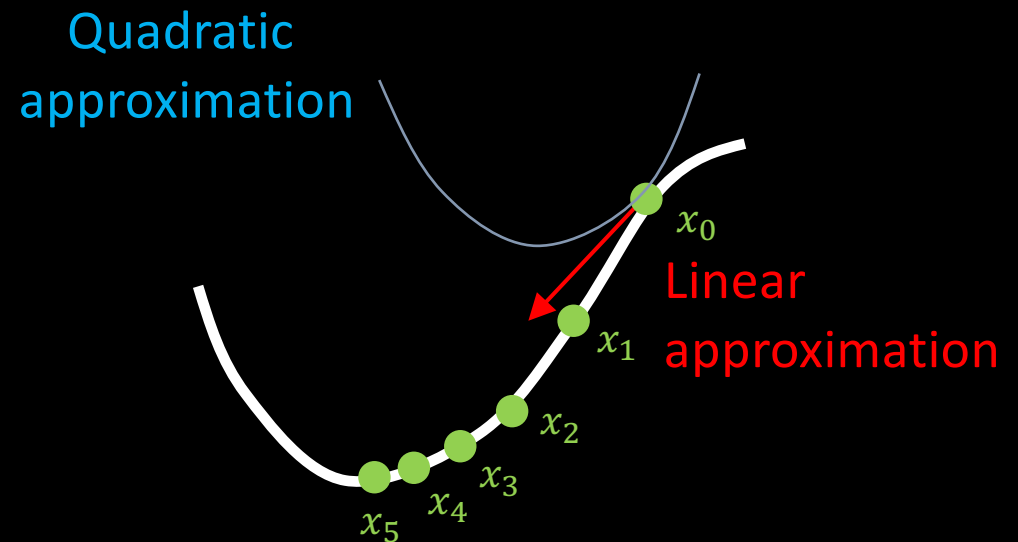
# Gradient Descent

> d examples

$\approx d$ memory



$(\approx d, > d)$

1st order method
Gradient descent

Anything in between???

2nd order
Gaussian elimination

$(d, d)$

Data

Memory

$(\approx d^2, d)$

Quadratic approximation
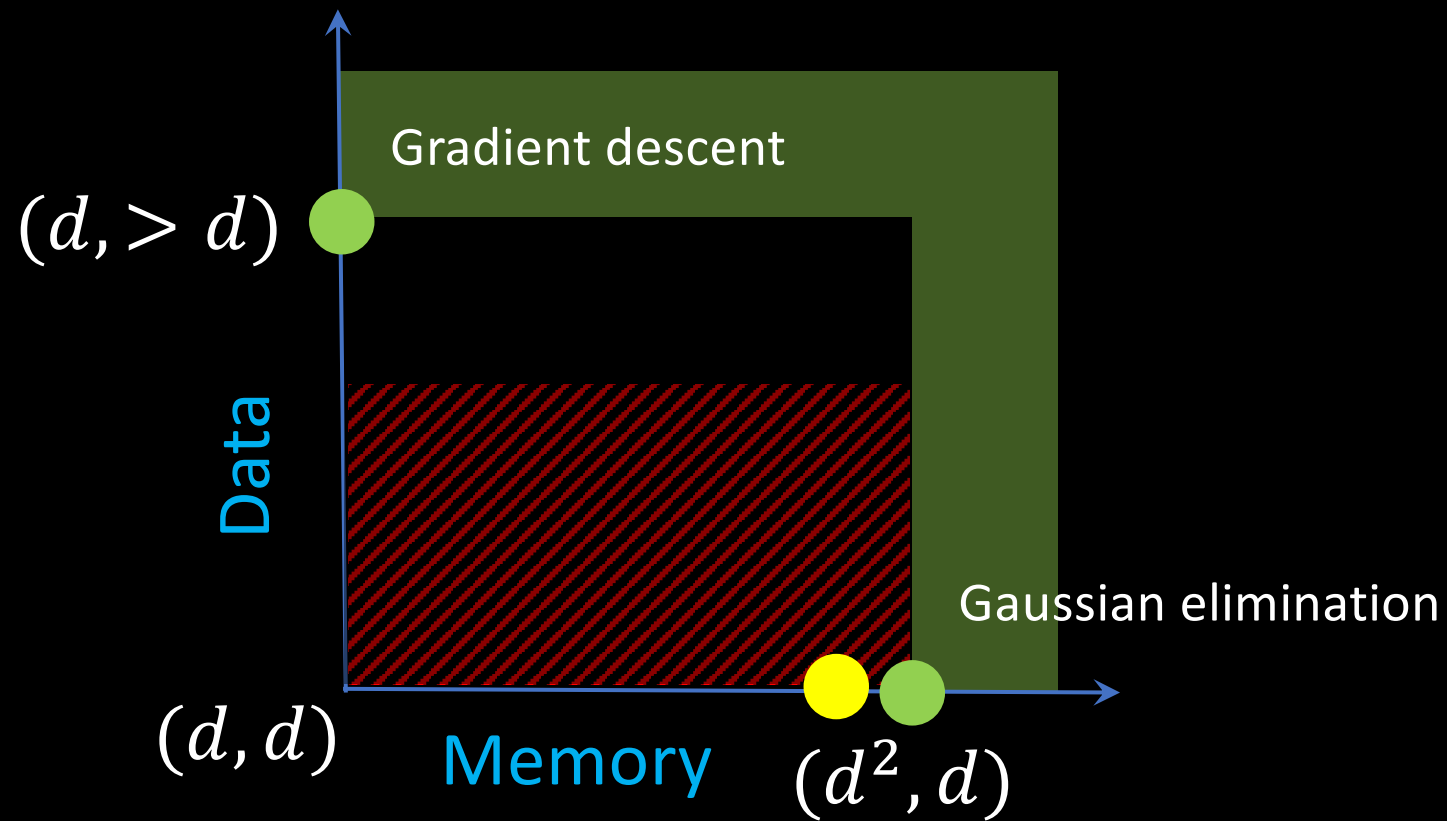
$x_0$

Linear approximation

$x_1$

$x_2$

$x_3$

$x_4$

$x_5$

1st order vs. 2nd order methods

Informal Theorem[Sharan, Sidford, Valiant]:
**Any sub-quadratic memory algorithm** requires more data.

Gradient descent

$(d, > d)$

Data

Gaussian elimination

$(d, d)$

Memory

$(d^2, d)$

Informal Theorem[Sharan, Sidford, Valiant]:
**Any sub-quadratic memory algorithm** requires more data.

Gradient descent

$(d, > d)$

Data

Gaussian elimination

$(d, d)$

Memory

$(d^2, d)$

Informal Theorem[Sharan, Sidford, Valiant]:
**Any sub-quadratic memory algorithm** requires more data.

Gradient descent

$(d, d \log \frac{1}{\varepsilon})$

$d \log \log \frac{1}{\varepsilon}$ samples

Data

Gaussian elimination

$(d, d)$

Memory

$(d^2, d)$

DISCUSSION

# 1.5<sup>th</sup> order method?

Quadratic approximation
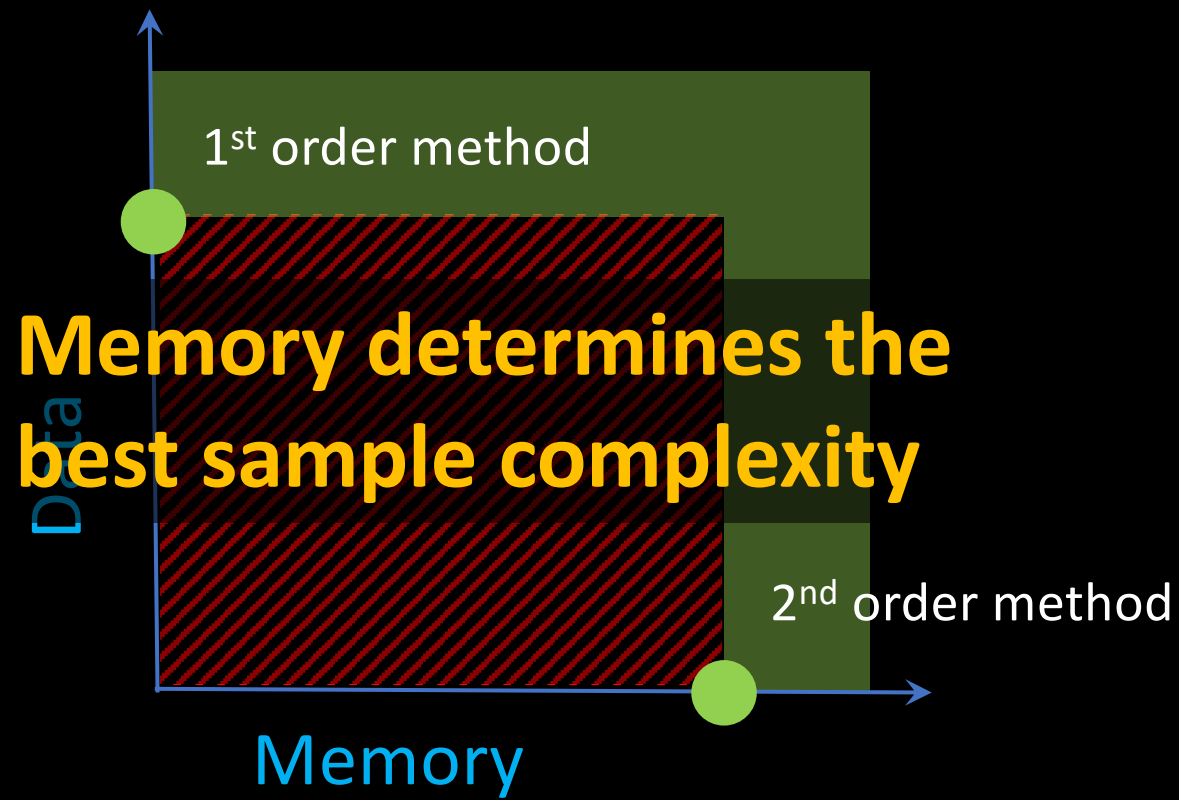
Linear approximation

$x_0$
$x_1$
$x_2$
$x_3$
$x_4$
$x_5$

1st order vs. 2nd order methods

Our Conjecture:
Any algorithm that improves on convergence rate of best known "first-order" methods, requires quadratic memory.

# 1.5<sup>th</sup> order method?

Our Conjecture:
Any algorithm that improves on convergence rate of best known "first-order" methods, requires quadratic memory.

Ill-conditioned distribution:

First order methods need $poly\,(\kappa)$ samples

(e.g. Needell-Srebro-Ward'16, Moritz-Nishihara-Jordan'16, Agarwal-Bullins-Hazan'17 etc. etc.)

Conjecture: There is a class of linear systems with condition number $\kappa$, such that any algorithm either requires $\Omega(d^2)$ memory or $d\,\text{poly}(\kappa)$ examples.

QUADRATIC MEMORY OR
CONDITION NUMBER SAMPLES?

# 1.5<sup>th</sup> order method?

Our Conjecture:
Any algorithm that improves on convergence rate of best known "first-order" methods, requires quadratic memory.

Ill-conditioned distribution:



QUADRATIC MEMORY OR...

CONDITION NUMBER SAMPLES?

[This talk] Memory Dichotomy Hypothesis: It is not possible to significantly improve on the convergence rate of known memory efficient techniques without using significantly more memory.

# 1.5th order method?

Our Conjecture:
Any algorithm that improves on convergence rate of best known "first-order" methods, requires quadratic memory.

Ill-conditioned distribution:

QUADRATIC MEMORY OR...

CONDITION NUMBER SAMPLES?

[This talk] Memory Dichotomy Hypothesis: It is not possible to significantly improve on the convergence rate of known memory efficient techniques without using significantly more memory.

# Using memory considerations to develop
## more efficient optimization algorithms

# Memory-efficient Algorithms for Optimization

Our Conjecture: There is a class of linear systems with condition number $\kappa$, such that any algorithm either requires $\Omega(d^2)$ memory or $d \, \text{poly}(\kappa)$ examples.

With more structure, can get best of both worlds!

Result (Informal):
For some structured linear systems, can get $d \, \textbf{polylog}(\kappa)$ examples with $O(d)$ memory!

This is true more broadly beyond linear systems,
and holds for any "multiscale" optimization problem.

# Memory-efficient Algorithms for Optimization

Result (Informal):
For some structured linear systems, can get
$d$ **polylog($\kappa$)** examples with $O(d)$ memory!

Linear system has
small number of
unique eigenvalues:

$$\begin{bmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & 1 & & & \\ & & & 1/\kappa & & \\ & & & & \ddots & \\ & & & & & 1/\kappa \end{bmatrix}$$

# Memory-efficient Algorithms for Optimization

Linear system has two unique eigenvalues

Too large for larger eigendirections!

$$
\begin{bmatrix}
1 & & & & & \\
& \ldots & & & & \\
& & 1 & & & \\
& & & 1/\kappa & & \\
& & & & \ldots & \\
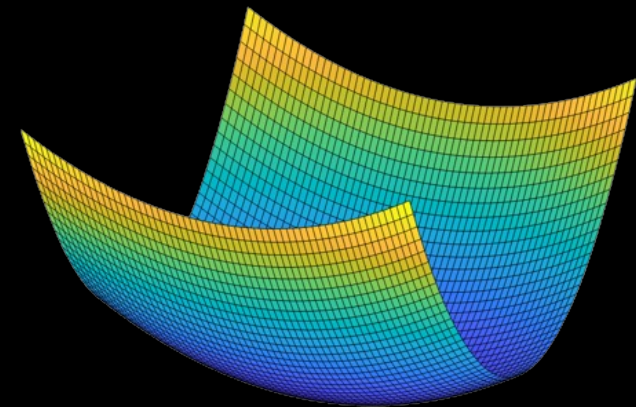& & & & & 1/\kappa
\end{bmatrix}
$$

Need about $\approx \kappa$ steps because of small eigendirections

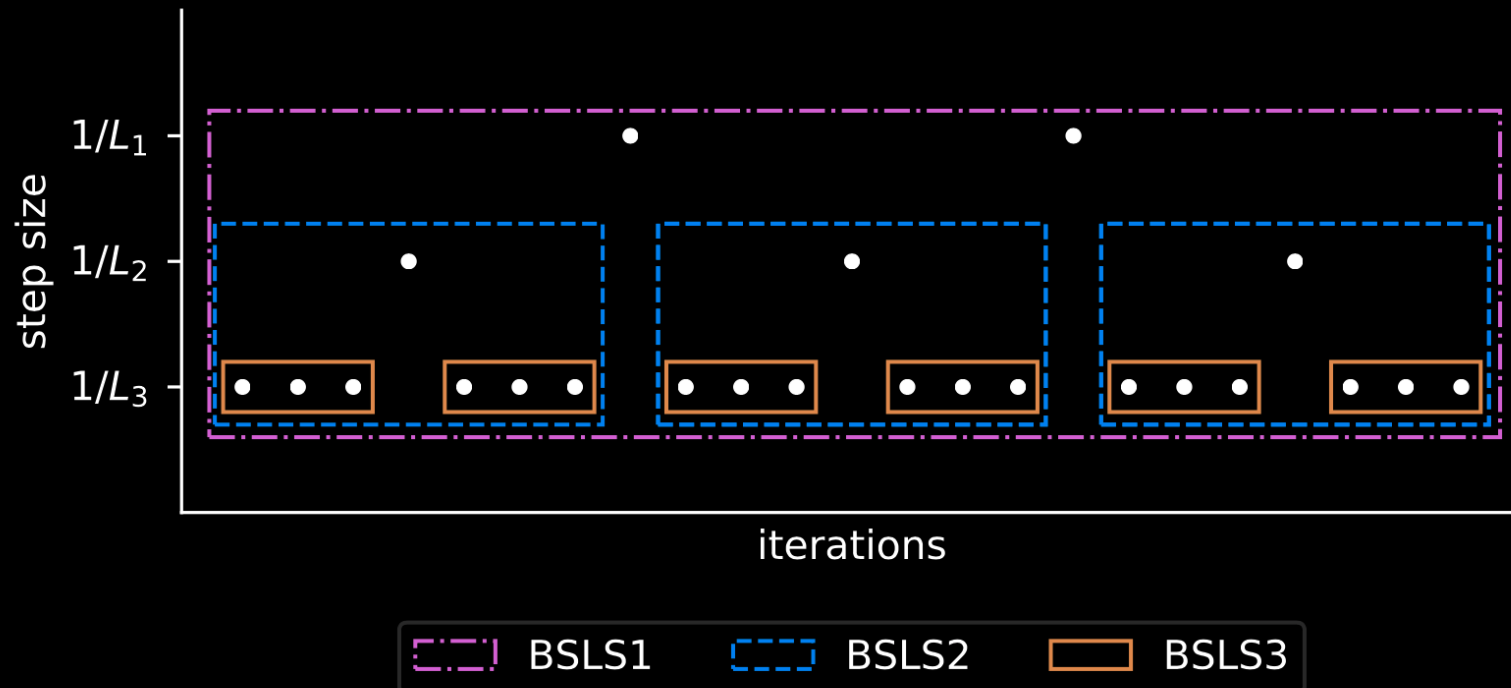Safest choice: Take step size $\approx 1$

Aggressive choice: Take step size $\approx \kappa$

Solution: Follow large step with small steps to fix error along larger eigendirections
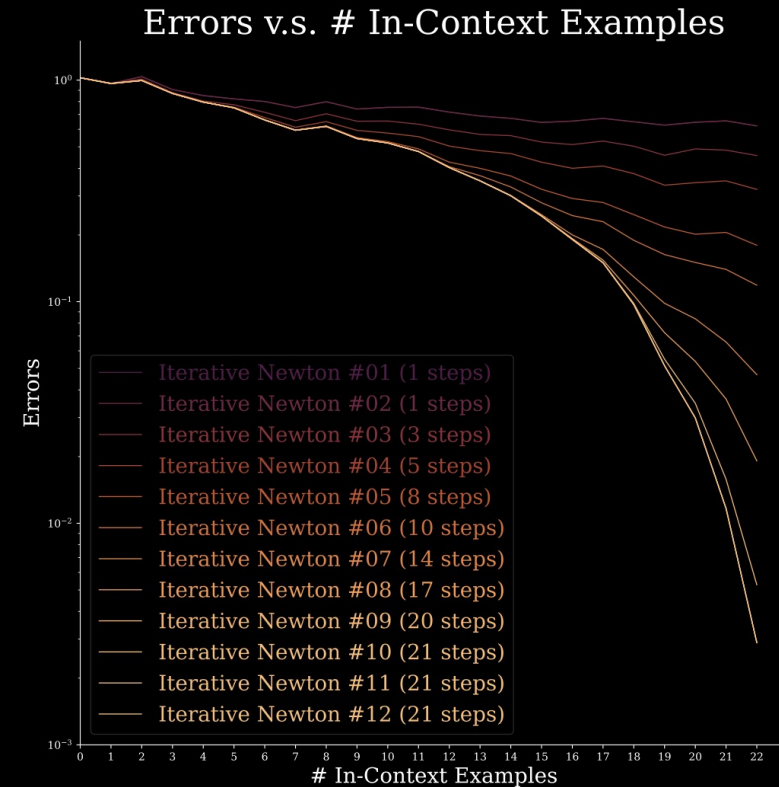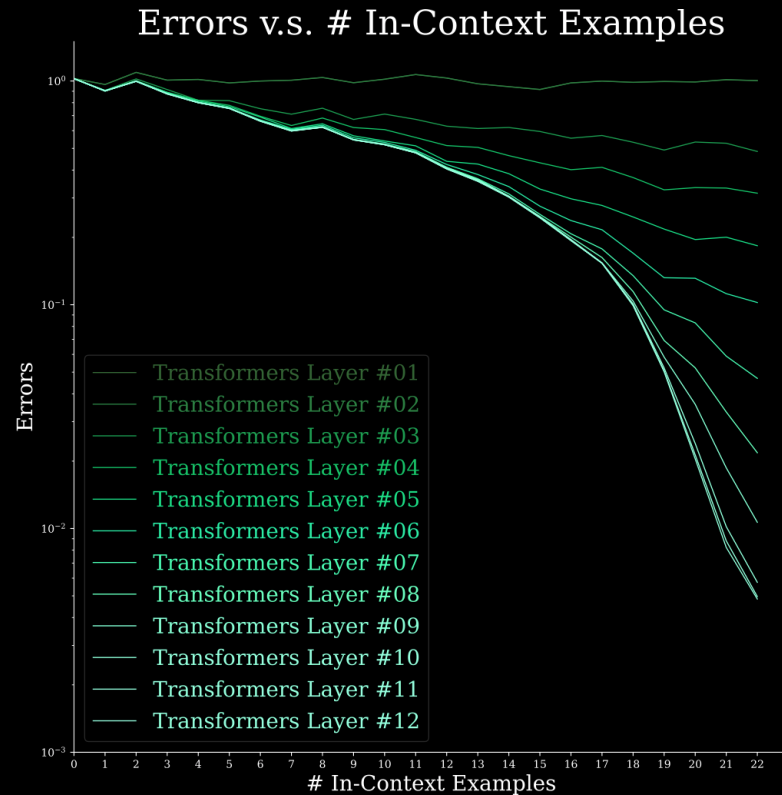
# Memory-efficient Algorithms for Optimization

Theorem (Kelner, Marsden, Sharan, Sidford, Yuan, Valiant):
For some structured linear systems, recursive sequence of large and small steps solves the problem with $d \, \text{polylog}(\kappa)$ **examples/gradient queries** and $O(d)$ **memory**.

# Using theory to understand what deep learning models learn?



*Transformers Learn Higher-Order Optimization Methods for In-Context Learning: A Study with Linear Models*
Deqing Fu, Tian-Qi Chen, Robin Jia, Vatsal Sharan, 2023

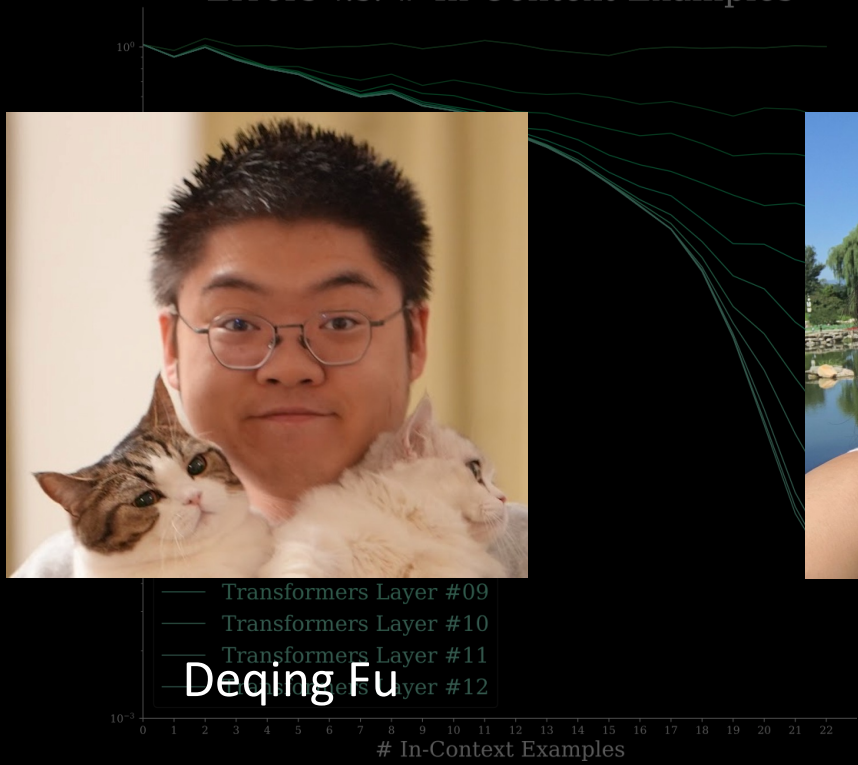# Using theory to understand what deep learning models learn?



*Transformers Learn Higher-Order Optimization Methods for In-Context Learning: A Study with Linear Models*
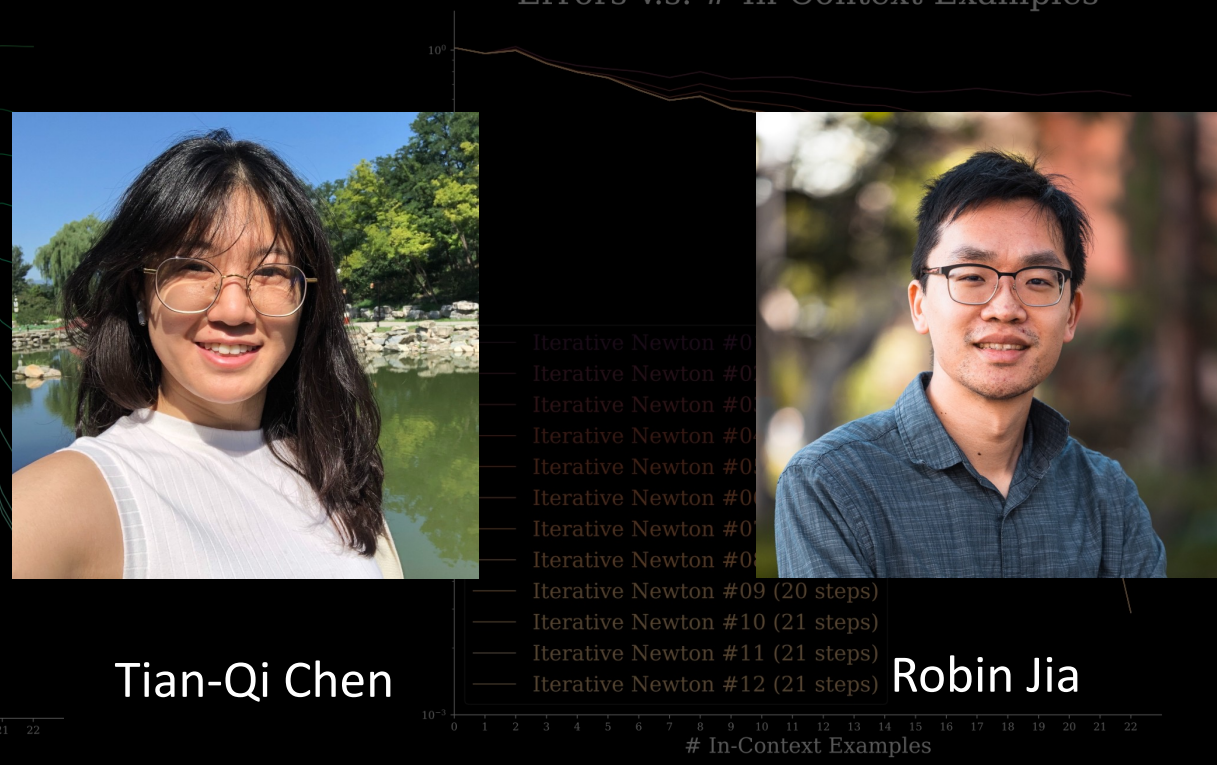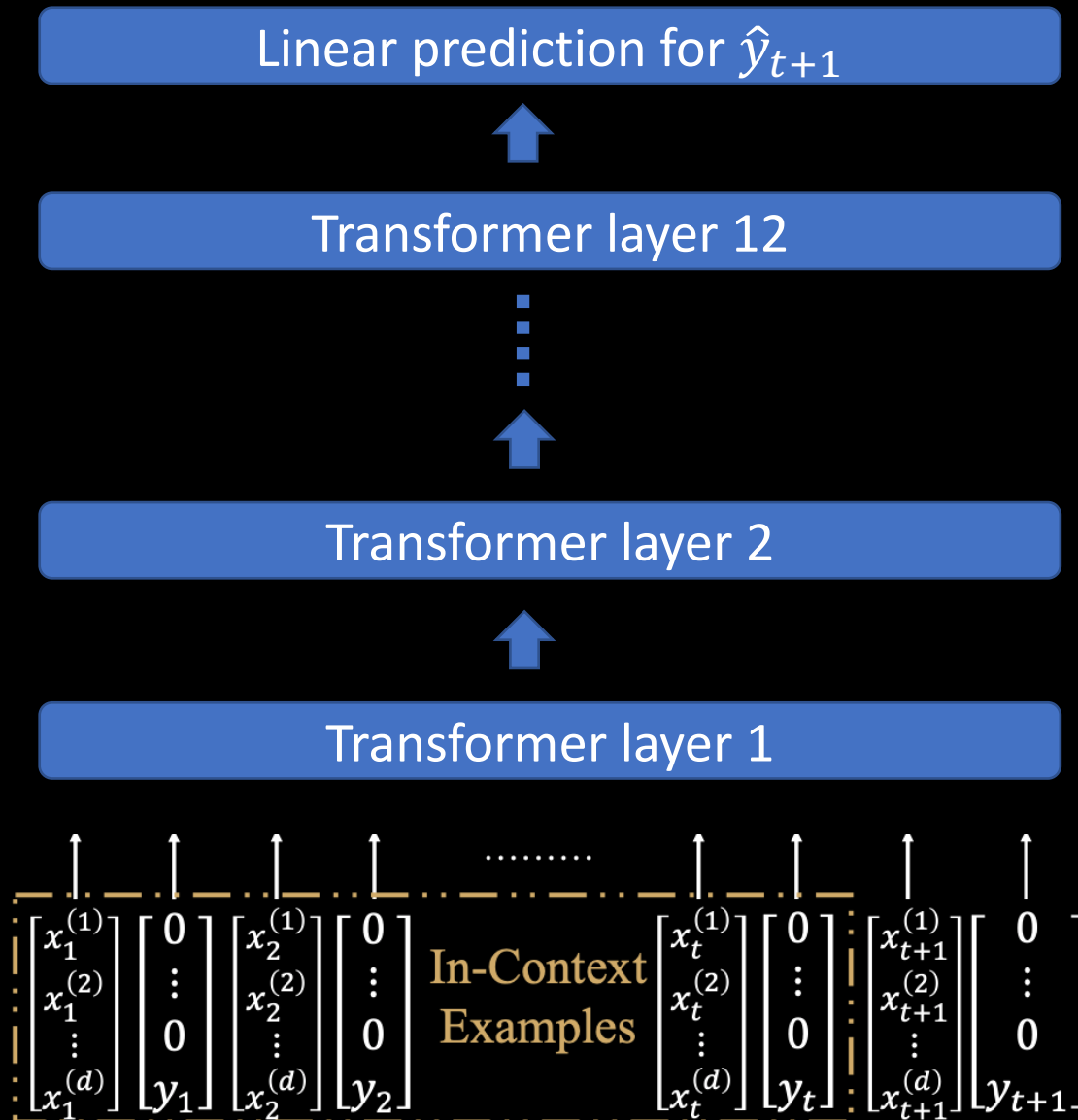Deqing Fu, Tian-Qi Chen, Robin Jia, Vatsal Sharan, 2023

# "Applied theory"?

Claim:

1. We can use understanding of statistical and computational gaps to understand mechanisms of models
2. Available memory may explain differences in behavior between different architectures

# 1. We can use understanding of statistical and computational gaps to understand mechanisms of models



Errors v.s. # In-Context Examples

Based on rate of convergence, argue that Transformers cannot be doing any 1st order method

Test on ill-conditioned settings where gap between 1st and 2nd order methods is largest

# 2. Available memory may explain differences in behavior between different architectures

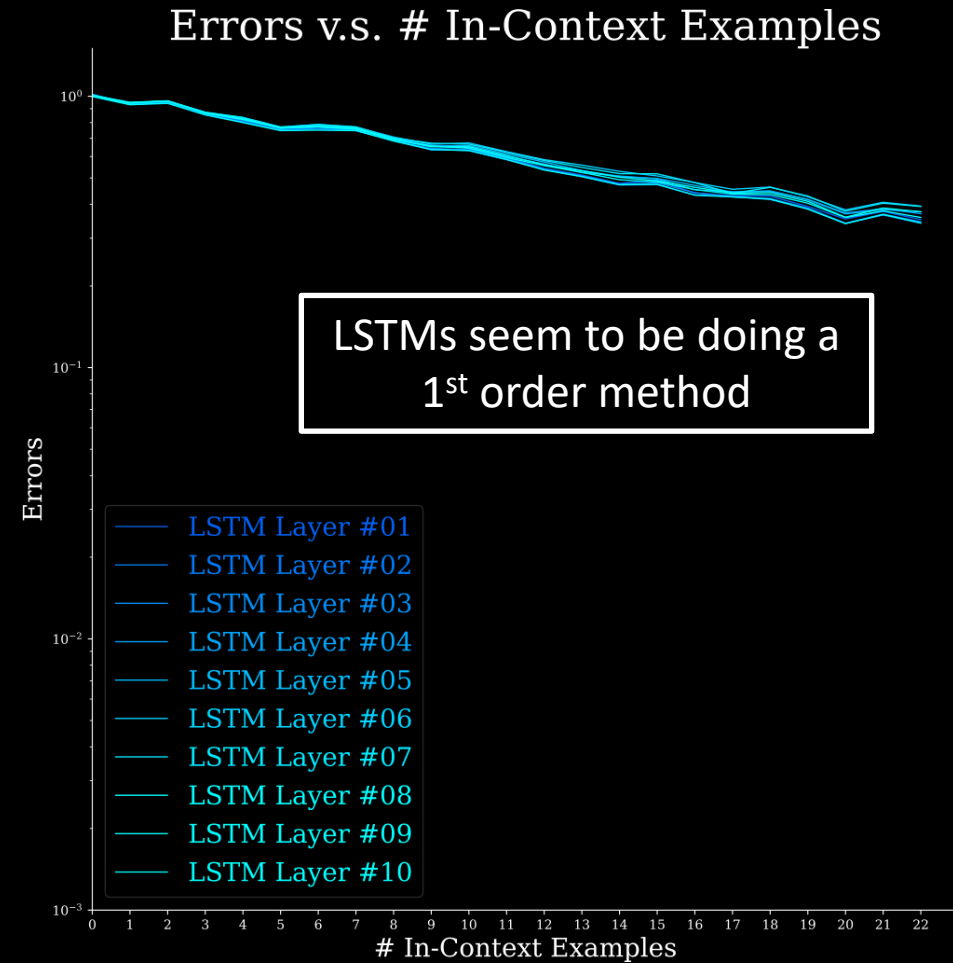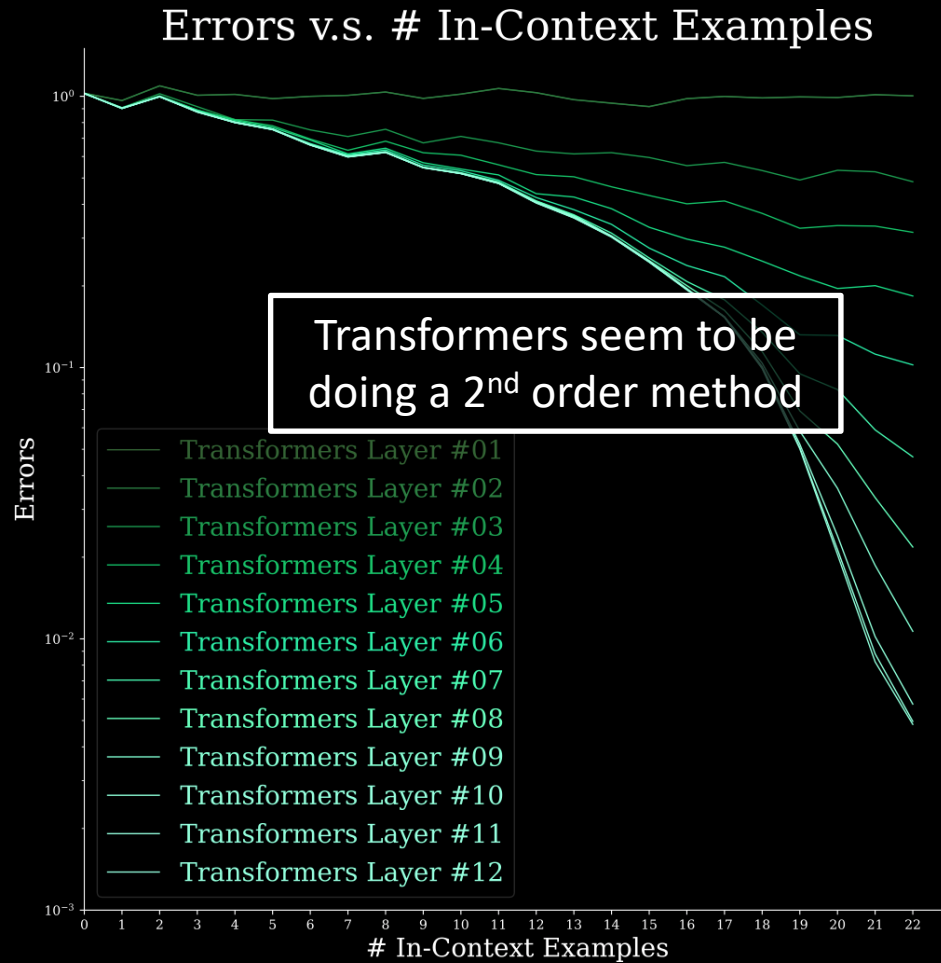# 2. Available memory may explain differences in behavior between different architectures



Errors v.s. # In-Context Examples

Transformers seem to be doing a 2nd order method

- Transformers Layer #01
- Transformers Layer #02
- Transformers Layer #03
- Transformers Layer #04
- Transformers Layer #05
- Transformers Layer #06
- Transformers Layer #07
- Transformers Layer #08
- Transformers Layer #09
- Transformers Layer #10
- Transformers Layer #11
- Transformers Layer #12

Errors v.s. # In-Context Examples

LSTMs seem to be doing a 1st order method

- LSTM Layer #01
- LSTM Layer #02
- LSTM Layer #03
- LSTM Layer #04
- LSTM Layer #05
- LSTM Layer #06
- LSTM Layer #07
- LSTM Layer #08
- LSTM Layer #09
- LSTM Layer #10

Memory is a fundamental computation resource.
Memory considerations are crucial in practice.

What is the role of memory in learning and optimization?
Are there tradeoffs between available memory and
information requirement?

**Op**timization

Memory

Memory Dichotomy Hypothesis: It is not possible to significantly improve on the convergence rate of known memory efficient techniques without using significantly more memory.

- Memory determines the best available convergence rate
- Memory provides a separation between simple and complex techniques
- New problem structures where we can circumvent lower bounds, new variants of GD