# Predicting & Optimizing ML Training Simons "Future of LLMs" Workshop

**April 3, 2025** 



Logan Engstrom



Ben Chen



#### Axel Feldmann



Billy Moses



#### Aleksander Mądry





Logan Engstrom

Ben Chen





#### Axel Feldmann



Billy Moses



#### Aleksander Mądry



Logan Engstrom



Ben Chen



#### Graduating now (from *undergrad*)!



Axel Feldmann



Billy Moses



Aleksander Mądry



#### Logan Engstrom



Ben Chen

# Graduating now!

#### Optimizing ML Training with Metagradient Descent

Logan Engstrom<sup>\*1</sup>, Andrew Ilyas<sup>\*2†</sup>, Benjamin Chen<sup>\*1</sup>, Axel Feldmann<sup>1</sup>, William Moses<sup>3</sup>, Aleksander Mądry<sup>1</sup> <sup>1</sup>MIT, <sup>2</sup>Stanford, <sup>3</sup>UIUC \*Equal contribution

#### Graduating now (from *undergrad*)!



**Axel Feldmann** 



Billy Moses



Aleksander Mądry



Training requires many design choices. Making the right choice is hard!

#### Training requires many design choices. Making the right choice is hard!

Model architecture



#### Training requires many design choices. Making the right choice is hard!

Model architecture

Hyperparameters





#### Training requires many design choices. Making the right choice is hard!

Model architecture

Hyperparameters





Training dataset



#### Training requires many design choices. Making the right choice is hard!

Model architecture

Hyperparameters



#### Fundamental Q: Can we predict and optimize the effect of these choices?

E.g.: Architecture search, scaling laws, data curation

Training dataset



#### The dream

A framework for optimizing and predicting the effect of design choices on model training

Example: Training data

(Sorscher et al. 2022; Gadre et al. 2023; Xie et al. 2024; Xia et al. 2024; Engstrom et al. 2024; many others)

(Sorscher et al. 2022; Gadre et al. 2023; Xie et al. 2024; Xia et al. 2024; Engstrom et al. 2024; many others)

(Sorscher et al. 2022; Gadre et al. 2023; Xie et al. 2024; Xia et al. 2024; Engstrom et al. 2024; many others)



(Sorscher et al. 2022; Gadre et al. 2023; Xie et al. 2024; Xia et al. 2024; Engstrom et al. 2024; many others)





(Sorscher et al. 2022; Gadre et al. 2023; Xie et al. 2024; Xia et al. 2024; Engstrom et al. 2024; many others)







(Sorscher et al. 2022; Gadre et al. 2023; Xie et al. 2024; Xia et al. 2024; Engstrom et al. 2024; many others)

electroniccigarettereviewed.info prestigedentalproducts.com brain-dumps.us





# **Optimizing:** Training

(Sorscher et al. 2022; Gadre et al. 2023; Xie et

electroniccigarettereviewed.info prestigedentalproducts.com brain-dumps.us

Scraped internet data

ac omplishmnl 'o the mi alon and welR I the men. U (SR 600-17b-I) are clear enough and cover everybody. iM-be UeIrU d at home first. IDretewr.. NatM WLWr Of Panama. B Aleol B -w s sin, ambas to V ')... &..... iZ'.". '- .- '."" ""5- ". ^ -\*," ..^ -^:? ^~-. .... ; ,'- i ,? '. " ". 1 the eye to those Europe. like an oriontal: Mrt of the Ch-. te a ter ii wre n capw. to '. %e anti-trust sutt "I made let of money, but . mg C1 Y ay With ne shot of a Lot d Aieles mining and. petrlce Wymore is fac9 aur rtl 1 Dalton. Gr at gag tUdS 't come! Mand Jobaim AlMauccessful- M ." y uw? ie House ar&rkd, "Hold A mtllon yeas? A trillion? t Arrival." Or are they ageless?

#### http://ufdc.ufl.edu/AA00010883/00095

#### iny others)







#### Much of the internet is low-quality: what data should we train on?

Equivalently: can we optimize the training data selection for target performance?

#### http://ufdc.ufl.edu/AA00010883/00095

ac omplishmnl 'o the mi alon and welR I the men. U (SR 600-17b-I) are clear enough and cover everybody. iM-be UelrU d at home first. IDretewr.. NatM WLWr Of Panama. B Aleol B -w s sin,ambas to V '). .. &. ... iZ'." . '- .- '."" "".5- ". ^-\*," ..^ -^:? ^~-. ....; ,'- i ,? '. " ". 1 the eye to those Europe. like an oriontal: Mrt of the Ch-. te a ter ii wre n capw. to '. %e anti-trust sutt "I made let of money, but . mg C1 Y ay With ne shot of a Lot d Aieles mining and. petrice Wymore is fac9 aur rtl 1 Dalton. Gr at gag tUdS 't come! Mand Jobaim AlMauccessful- M ." y uw? ie House ar&rkd, "Hold A mtllon yeas? A trillion? t Arrival." Or are they ageless?

#### iny others)







(Koh & Liang 2017; Feldman 2019; I et al. 2022; Bae et al. 2023, Park et al. 2024; many others)



#### Predict



















(Koh & Liang 2017; Feldman 2019; I et al. 2022; Bae et al. 2023, Park et al. 2024; many others)



"What examples in the training data caused this behavior?"



(Koh & Liang 2017; Feldman 2019; I et al. 2022; Bae et al. 2023, Park et al. 2024; many others)



Equivalently: Can we predict change in this behavior as a function of training data?

"What examples in the training data caused this behavior?"



# Full problem statement

Idea: write model output as a direct function of training setup ("metaparameter")

Idea: write model output as a direct function of training setup ("metaparameter")

Metaparameter  $\mathbf{z} \in \mathbf{Z} \subset \mathbb{R}^d$ 

The aspect(s) of the training setup that we want to study

Metaparameter

 $\mathbf{z} \in \mathbf{Z} \subset \mathbb{R}^d$ 

Idea: write **model output** as a direct function of **training setup** ("metaparameter")



Idea: write **model output** as a direct function of **training setup** ("metaparameter")

Trained model  $\theta = \mathscr{A}(\mathbf{z})$ 

The aspect(s) of the training setup that we want to study

Metaparameter  $\mathbf{z} \in \mathbf{Z} \subset \mathbb{R}^d$ 

Training routine  $\mathscr{A}:\mathbf{Z}\to\Theta$ 

Trained model  $\theta = \mathscr{A}(\mathbf{z})$ 

Map from metaparameter to model, all else fixed

Idea: write **model output** as a direct function of **training setup** ("metaparameter")
#### Idea: write **model output** as a direct function of **training setup** ("metaparameter")

The aspect(s) of the training setup that we want to study

Metaparameter  $\mathbf{z} \in \mathbf{Z} \subset \mathbb{R}^d$ 

Training routine  $\mathscr{A}:\mathbf{Z}\to\Theta$ 

Trained model  $\theta = \mathscr{A}(\mathbf{z})$ 

Map from metaparameter to model, all else fixed

Output function  $\phi: \Theta \to \mathbb{R}$ 

Final output  $\phi(\mathscr{A}(\mathbf{Z}))$ 

#### Idea: write **model output** as a direct function of **training setup** ("metaparameter")

The aspect(s) of the training setup that we want to study

Metaparameter  $\mathbf{z} \in \mathbf{Z} \subset \mathbb{R}^d$ 

Training routine  $\mathscr{A}:\mathbf{Z}\to\Theta$ 

Map from metaparameter to model, all else fixed



#### Idea: write model output as a direct function of training setup ("metaparameter")

The aspect(s) of the training setup that we want to study

Metaparameter

 $\mathbf{z} \in \mathbf{Z} \subset \mathbb{R}^d$ 

Training routine  $\mathscr{A}: \mathbb{Z} \to \Theta$ 

Map from metaparameter to model, all else fixed

**Goals:** 



#### Idea: write **model output** as a direct function of **training setup** ("metaparameter")

The aspect(s) of the training setup that we want to study

Metaparameter

 $\mathbf{z} \in \mathbf{Z} \subset \mathbb{R}^d$ 

Training routine

 $\mathscr{A}: \mathbf{Z} \to \Theta$ 

Map from metaparameter to model, all else fixed

#### **Goals**:

Optimization: Find metaparameter  $\mathbf{z}^* = \arg \min \phi(\mathscr{A}(\mathbf{z}))$  minimizing output Ζ



#### Idea: write **model output** as a direct function of **training setup** ("metaparameter")

The aspect(s) of the training setup that we want to study

Metaparameter  $\mathbf{z} \in \mathbf{Z} \subset \mathbb{R}^d$ 

Training routine  $\mathscr{A}: \mathbf{Z} \to \Theta$ 

Map from metaparameter to model, all else fixed

#### **Goals**:



#### Optimization: Find metaparameter $\mathbf{z}^* = \arg \min \phi(\mathscr{A}(\mathbf{z}))$ minimizing output Z Prediction: Find $\hat{f}(\mathbf{z}) \approx \phi(\mathscr{A}(\mathbf{z}))$ , direct model output estimator (does not invoke $\mathscr{A}$ )



#### Idea: write **model output** as a direct function of **training setup** ("metaparameter")

The aspect(s) of the training setup that we want to study

Metaparameter  $\mathbf{z} \in \mathbf{Z} \subset \mathbb{R}^d$ 

Training routine  $\mathscr{A}: \mathbf{Z} \to \Theta$ 

Map from metaparameter to model, all else fixed

**Goals**:



#### Optimization: Find metaparameter $\mathbf{z}^* = \arg \min \phi(\mathscr{A}(\mathbf{z}))$ minimizing output Z Prediction: Find $\hat{f}(\mathbf{z}) \approx \phi(\mathscr{A}(\mathbf{z}))$ , direct model output estimator (does not invoke $\mathscr{A}$ )

### Captures data curation & attribution but also many established ML challenges



#### Idea: write **model output** as a direct function of **training setup** ("metaparameter")

The aspect(s) of the training setup that we want to study

Metaparameter  $\mathbf{z} \in \mathbf{Z} \subset \mathbb{R}^d$ 

Training routine  $\mathscr{A}:\mathbf{Z}\to\Theta$ 

Map from metaparameter to model, all else fixed



#### Idea: write model output as a direct function of training setup ("metaparameter")



#### Idea: write **model output** as a direct function of **training setup** ("metaparameter")

The aspect(s) of the training setup that we want to study

Metaparameter  $\mathbf{z} \in \mathbf{Z} \subset \mathbb{R}^d$ 

Training routine

 $\mathscr{A}:\mathbf{Z}\to\Theta$ 

Map from metaparameter to model, all else fixed

**Data selection** 



#### Idea: write model output as a direct function of training setup ("metaparameter")

The aspect(s) of the training setup that we want to study

Metaparameter  $\mathbf{z} \in \mathbf{Z} \subset \mathbb{R}^d$ 

**Data selection** 

Training routine  $\mathscr{A}: \mathbb{Z} \to \Theta$ 

Map from metaparameter to model, all else fixed



#### This will be the output of the model trained on **z**

Might be the loss on a specific fixed target example, or performance on a fixed test set

#### Idea: write **model output** as a direct function of **training setup** ("metaparameter")

The aspect(s) of the training setup that we want to study

Metaparameter  $\mathbf{z} \in \mathbf{Z} \subset \mathbb{R}^d$ 

**Data selection** 

Training routine  $\mathscr{A}: \mathbf{Z} \to \Theta$ 

Map from metaparameter to model, all else fixed

#### **Goal (Data curation)**

Find the selection  $\mathbf{z}$  that minimizes the resulting model's test loss





#### Idea: write **model output** as a direct function of **training setup** ("metaparameter")

The aspect(s) of the training setup that we want to study

Metaparameter  $\mathbf{z} \in \mathbf{Z} \subset \mathbb{R}^d$ 

**Data selection** 

Training routine  $\mathscr{A}: \mathbf{Z} \to \Theta$ 

Map from metaparameter to model, all else fixed

#### **Goal (Data curation)**

Find the selection  $\mathbf{z}$  that minimizes the resulting model's test loss



#### **Goal (Data attribution)**

Predict how changing the selection z will change resulting model's output





### Roadmap

Main idea: The metagradient

Prequel: Two challenges

Looking forward

### Roadmap

#### Main idea: The metagradient

#### Prequel: Two challenges

Looking forward



In many cases, metaparameter  $\mathbf{z}$  is a continuous variable (or can be relaxed to one)



In many cases, metaparameter **z** is a continuous variable (or can be relaxed to one) We will show how to **compute** and **use** the *metagradient*:



In many cases, metaparameter **z** is a continuous variable (or can be relaxed to one) We will show how to **compute** and **use** the *metagradient*:



In many cases, metaparameter z is a continuous variable (or can be relaxed to one) We will show how to **compute** and **use** the *metagradient*:



**Cost: In practice, "only" 3-5x the cost of model training!** 

### Metagradients offer a powerful way to optimize and predict (aspects of) model behavior

Given a task:

Given a task:

1.

#### **Rewrite** problem in terms of model outputs (as a function of a metaparameter)

Given a task:

1.

#### **Rewrite** problem in terms of model outputs (as a function of a metaparameter) Decide on appropriate z, $\mathscr{A}$ , $\phi$ so that the problem becomes $\min \phi(\mathscr{A}(z))$ $z \in \mathbb{R}^d$

Given a task:

- 1.
- 2. **Solve** via (meta)gradient descent:

#### **Rewrite** problem in terms of model outputs (as a function of a metaparameter) Decide on appropriate z, $\mathscr{A}$ , $\phi$ so that the problem becomes min $\phi(\mathscr{A}(z))$ $z \in \mathbb{R}^d$

Given a task:

- 1.
- 2. **Solve** via (meta)gradient descent: A. Pick a starting  $\mathbf{z}^{(0)}$  and number of steps T

#### **Rewrite** problem in terms of model outputs (as a function of a metaparameter) Decide on appropriate z, $\mathscr{A}$ , $\phi$ so that the problem becomes min $\phi(\mathscr{A}(z))$ $z \in \mathbb{R}^d$

Given a task:

- 1.
- 2. **Solve** via (meta)gradient descent:
  - A. Pick a starting  $\mathbf{z}^{(0)}$  and number of steps T
  - B. For i = 1...T: Compute  $\mathbf{g} := \nabla_{\mathbf{z}} \phi(\mathscr{A}(\mathbf{z}^{(i-1)}))$  and descend

#### **Rewrite** problem in terms of model outputs (as a function of a metaparameter) Decide on appropriate z, $\mathscr{A}$ , $\phi$ so that the problem becomes $\min \phi(\mathscr{A}(z))$ $z \in \mathbb{R}^d$

Given a task:

- 1.
- 2. **Solve** via (meta)gradient descent:
  - A. Pick a starting  $\mathbf{z}^{(0)}$  and number of steps T
  - B. For i = 1...T: Compute  $\mathbf{g} := \nabla_{\mathbf{z}} \phi(\mathscr{A}(\mathbf{z}^{(i-1)}))$  and descend Requires T total metagradient calculations

#### **Rewrite** problem in terms of model outputs (as a function of a metaparameter) Decide on appropriate **z**, $\mathscr{A}$ , $\phi$ so that the problem becomes min $\phi(\mathscr{A}(\mathbf{z}))$ $Z \in \mathbb{R}^d$

#### Choose train subset



#### Choose train subset



#### Train a model



#### Choose train subset



#### Train a model





# Measure loss on a target/test set:



#### Choose train subset



**Goal:** Select the training data subset minimizing loss on the target set

### Choose train subset



### **Goal:** Select the training data subset minimizing loss on the target set

Equivalent goal: Find the training data weighting minimizing loss on the target set

### Choose train subset



### **Goal:** Select the training data subset minimizing loss on the target set

Equivalent goal: Find the training data weighting minimizing loss on the target set

### Choose train subset



### **Goal:** Select the training data subset minimizing loss on the target set **Equivalent goal:** Find the training data *weighting* minimizing loss on the target set

#### Increasingly recognized as a critical step in large-scale ML pipeline
### **Metaparameter z**: vector $\mathbf{z} \in \mathbb{R}^N$ s.t., $\mathbf{z}_i$ is the weight on *i*-th training example

- **Metaparameter z**: vector  $\mathbf{z} \in \mathbb{R}^N$  s.t.,  $\mathbf{z}_i$  is the weight on *i*-th training example
- Learning algorithm A: standardized CLIP training routine

# $\mathbf{z}_i$ is the weight on i-th training example LIP training routine

- **Metaparameter z**: vector  $\mathbf{z} \in \mathbb{R}^N$  s.t.,  $\mathbf{z}_i$  is the weight on *i*-th training example **Learning algorithm**  $\mathscr{A}$ : standardized CLIP training routine
- **Output function:** Loss function  $L_{targ}$  evaluated on target tasks

**Metaparameter z**: vector  $\mathbf{z} \in \mathbb{R}^N$  s.t.,  $\mathbf{z}_i$  is the weight on *i*-th training example

**Learning algorithm**  $\mathscr{A}$ : standardized CLIP training routine

**Output function:** Loss function  $L_{targ}$  evaluated on target tasks

Repeat, starting from  $z = 1_N$ :

**Metaparameter z**: vector  $\mathbf{z} \in \mathbb{R}^N$  s.t.,  $\mathbf{z}_i$  is the weight on *i*-th training example

**Learning algorithm**  $\mathscr{A}$ : standardized CLIP training routine

**Output function:** Loss function  $L_{targ}$  evaluated on target tasks

Repeat, starting from  $\mathbf{z} = \mathbf{1}_N$ :

1. Train model weighting the training points by **Z** 

**Metaparameter z**: vector  $\mathbf{z} \in \mathbb{R}^N$  s.t.,  $\mathbf{z}_i$  is the weight on *i*-th training example

**Learning algorithm**  $\mathscr{A}$ : standardized CLIP training routine

**Output function:** Loss function  $L_{targ}$  evaluated on target tasks

Repeat, starting from  $\mathbf{z} = \mathbf{1}_N$ :

- Train model weighting the training points by **Z**
- 2. Calculate metagradient  $\mathbf{g} = \nabla_{\mathbf{z}} \phi(\mathscr{A}(\mathbf{z}))$

**Metaparameter z**: vector  $\mathbf{z} \in \mathbb{R}^N$  s.t.,  $\mathbf{z}_i$  is the weight on *i*-th training example

**Learning algorithm**  $\mathscr{A}$ : standardized CLIP training routine

**Output function:** Loss function  $L_{targ}$  evaluated on target tasks

Repeat, starting from  $\mathbf{z} = \mathbf{1}_N$ :

- Train model weighting the training points by **Z**
- 2. Calculate metagradient  $\mathbf{g} = \nabla_{\mathbf{z}} \phi(\mathscr{A}(\mathbf{z}))$
- 3. Update  $\mathbf{z} \leftarrow \mathbf{z} \operatorname{sign}(\mathbf{g})$

**Metaparameter z**: vector  $\mathbf{z} \in \mathbb{R}^N$  s.t.,  $\mathbf{z}_i$  is the weight on *i*-th training example

**Learning algorithm**  $\mathscr{A}$ : standardized CLIP training routine

**Output function:** Loss function  $L_{targ}$  evaluated on target tasks

Repeat, starting from  $\mathbf{z} = \mathbf{1}_N$ :

- Train model weighting the training points by **Z**
- 2. Calculate metagradient  $\mathbf{g} = \nabla_{\mathbf{z}} \phi(\mathscr{A}(\mathbf{z}))$
- 3. Update  $\mathbf{z} \leftarrow \mathbf{z} \operatorname{sign}(\mathbf{g})$

\*simplified

#### **DATACOMP: In search of the next generation of multimodal datasets**

Samir Yitzhak Gadre\*<sup>2</sup>, Gabriel Ilharco\*<sup>1</sup>, Alex Fang\*<sup>1</sup>, Jonathan Hayase<sup>1</sup>, Georgios Smyrnis<sup>5</sup>, Thao Nguyen<sup>1</sup>, Ryan Marten<sup>7,9</sup>, Mitchell Wortsman<sup>1</sup>, Dhruba Ghosh<sup>1</sup>, Jieyu Zhang<sup>1</sup>, Eyal Orgad<sup>3</sup>, Rahim Entezari<sup>10</sup>, Giannis Daras<sup>5</sup>, Sarah Pratt<sup>1</sup>, Vivek Ramanujan<sup>1</sup>, Yonatan Bitton<sup>11</sup>, Kalyani Marathe<sup>1</sup>, Stephen Mussmann<sup>1</sup>, Richard Vencu<sup>6</sup>, Mehdi Cherti<sup>6,8</sup>, Ranjay Krishna<sup>1</sup>, Pang Wei Koh<sup>1,12</sup>, Olga Saukh<sup>10</sup>, Alexander Ratner<sup>1,13</sup>, Shuran Song<sup>2</sup>, Hannaneh Hajishirzi<sup>1,7</sup>, Ali Farhadi<sup>1</sup>, Romain Beaumont<sup>6</sup>, Sewoong Oh<sup>1</sup>, Alex Dimakis<sup>5</sup>, Jenia Jitsev<sup>6,8</sup>, Yair Carmon<sup>3</sup>, Vaishaal Shankar<sup>4</sup>, Ludwig Schmidt<sup>1,6,7</sup>

#### **DATACOMP:** In search of the next generation of multimodal datasets

Samir Yitzhak Gadre\*<sup>2</sup>, Gabriel Ilharco\*<sup>1</sup>, Alex Fang\*<sup>1</sup>, Jonathan Hayase<sup>1</sup>, Georgios Smyrnis<sup>5</sup>, Thao Nguyen<sup>1</sup>, Ryan Marten<sup>7,9</sup>, Mitchell Wortsman<sup>1</sup>, Dhruba Ghosh<sup>1</sup>, Jieyu Zhang<sup>1</sup>, Eyal Orgad<sup>3</sup>, Rahim Entezari<sup>10</sup>, Giannis Daras<sup>5</sup>, Sarah Pratt<sup>1</sup>, Vivek Ramanujan<sup>1</sup>, Yonatan Bitton<sup>11</sup>, Kalyani Marathe<sup>1</sup>, Stephen Mussmann<sup>1</sup>, Richard Vencu<sup>6</sup>, Mehdi Cherti<sup>6,8</sup>, Ranjay Krishna<sup>1</sup>, Pang Wei Koh<sup>1,12</sup>, Olga Saukh<sup>10</sup>, Alexander Ratner<sup>1,13</sup>, Shuran Song<sup>2</sup>, Hannaneh Hajishirzi<sup>1,7</sup>, Ali Farhadi<sup>1</sup>, Romain Beaumont<sup>6</sup>, Sewoong Oh<sup>1</sup>, Alex Dimakis<sup>5</sup>, Jenia Jitsev<sup>6,8</sup>, Yair Carmon<sup>3</sup>, Vaishaal Shankar<sup>4</sup>, Ludwig Schmidt<sup>1,6,7</sup>

Data selection task where the learning algorithm (CLIP training) and evaluation (zero-shot classification) are **fixed** 

#### **DATACOMP:** In search of the next generation of multimodal datasets

Samir Yitzhak Gadre\*<sup>2</sup>, Gabriel Ilharco\*<sup>1</sup>, Alex Fang\*<sup>1</sup>, Jonathan Hayase<sup>1</sup>, Georgios Smyrnis<sup>5</sup>, Thao Nguyen<sup>1</sup>, Ryan Marten<sup>7,9</sup>, Mitchell Wortsman<sup>1</sup>, Dhruba Ghosh<sup>1</sup>, Jieyu Zhang<sup>1</sup>, Eyal Orgad<sup>3</sup>, Rahim Entezari<sup>10</sup>, Giannis Daras<sup>5</sup>, Sarah Pratt<sup>1</sup>, Vivek Ramanujan<sup>1</sup>, Yonatan Bitton<sup>11</sup>, Kalyani Marathe<sup>1</sup>, Stephen Mussmann<sup>1</sup>, Richard Vencu<sup>6</sup>, Mehdi Cherti<sup>6,8</sup>, Ranjay Krishna<sup>1</sup>, Pang Wei Koh<sup>1,12</sup>, Olga Saukh<sup>10</sup>, Alexander Ratner<sup>1,13</sup>, Shuran Song<sup>2</sup>, Hannaneh Hajishirzi<sup>1,7</sup>, Ali Farhadi<sup>1</sup>, Romain Beaumont<sup>6</sup>, Sewoong Oh<sup>1</sup>, Alex Dimakis<sup>5</sup>, Jenia Jitsev<sup>6,8</sup>, Yair Carmon<sup>3</sup>, Vaishaal Shankar<sup>4</sup>, Ludwig Schmidt<sup>1,6,7</sup>

Data selection task where the learning algorithm (CLIP training) and evaluation (zero-shot classification) are **fixed** 

Goal: Curate the best possible dataset

#### **DATACOMP:** In search of the next generation of multimodal datasets

Samir Yitzhak Gadre\*<sup>2</sup>, Gabriel Ilharco\*<sup>1</sup>, Alex Fang\*<sup>1</sup>, Jonathan Hayase<sup>1</sup>, Georgios Smyrnis<sup>5</sup>, Thao Nguyen<sup>1</sup>, Ryan Marten<sup>7,9</sup>, Mitchell Wortsman<sup>1</sup>, Dhruba Ghosh<sup>1</sup>, Jieyu Zhang<sup>1</sup>, Eyal Orgad<sup>3</sup>, Rahim Entezari<sup>10</sup>, Giannis Daras<sup>5</sup>, Sarah Pratt<sup>1</sup>, Vivek Ramanujan<sup>1</sup>, Yonatan Bitton<sup>11</sup>, Kalyani Marathe<sup>1</sup>, Stephen Mussmann<sup>1</sup>, Richard Vencu<sup>6</sup>, Mehdi Cherti<sup>6,8</sup>, Ranjay Krishna<sup>1</sup>, Pang Wei Koh<sup>1,12</sup>, Olga Saukh<sup>10</sup>, Alexander Ratner<sup>1,13</sup>, Shuran Song<sup>2</sup>, Hannaneh Hajishirzi<sup>1,7</sup>, Ali Farhadi<sup>1</sup>, Romain Beaumont<sup>6</sup>, Sewoong Oh<sup>1</sup>, Alex Dimakis<sup>5</sup>, Jenia Jitsev<sup>6,8</sup>, Yair Carmon<sup>3</sup>, Vaishaal Shankar<sup>4</sup>, Ludwig Schmidt<sup>1,6,7</sup>

Data selection task where the learning algorithm (CLIP training) and evaluation (zero-shot classification) are **fixed** 

Goal: Curate the best possible dataset

### Example leaderboard (small):

Rank	Created 🔶	Submission	ImageNet acc.	Average perf.
1	03-14-2025	M-FLYT + SCS (alpha=0.5)	0.080	0.185
2	08-18-2024	EcoDatum	0.053	0.182
3	08-12-2023	> W <	0.056	0.180
4	11-04-2023	Hype Sampler	0.063	0.180
5	01-08-2024	Content Alignment Model	0.064	0.178
6	09-07-2023	Detection   CLIP score	0.053	0.178
7	09-06-2023	BLIP2-COCO-finetuned_similarity_top-35%	0.053	0.177
8	08-03-2023	WS (baselines)	0.059	0.175
9	04-28-2023	Baseline: CLIP score (L/14 30%)	0.051	0.173
10	09-07-2023	OCR and Naive english filtering	0.057	0.173
11	08-05-2024	MoDO	0.051	0.173
12	08-25-2023	Fondant baseline, CLIP score (L/14 1%-30%), ascii > 4	0.053	0.168
13	08-25-2023	text masked CLIP 30%	0.048	0.164
14	08-02-2023	CLIP score (L/14 30%) baseline	0.049	0.164
15	04-28-2023	Baseline: Image-based	0.043	0.159









Method	Score	Δ	
Baseline: No filtering	0.13	_	
Best baseline from [GIF+24]	0.17	+0.04	
Previous SOTA [Eco24]	0.18	+0.05	
	0.22	+0.09	





### MGD improvement (relative to no selection) is **double** existing methods'

Method	Score	Δ	
Baseline: No filtering	0.13	_	
Best baseline from [GIF+24]	0.17	+0.04	
Previous SOTA [Eco24]	0.18	+0.05	
	0.22	+0.09	

Best existing method



### MGD improvement (relative to no selection) is **double** existing methods'

Method	Score	Δ	
<ul> <li>- Baseline: No filtering</li> </ul>	0.13	_	
Best baseline from [GIF+24]	0.17	+0.04	
Previous SOTA [Eco24]	0.18	+0.05	
	0.22	+0.09	

Best existing method

### Now working on DataComp-medium (already #1!)

New goal: Select instruction fine-tuning data for LM benchmarks (2B param model)

New goal: Select instruction fine-tuning data for LM benchmarks (2B param model)





New goal: Select instruction fine-tuning data for LM benchmarks (2B param model)



	<b>BBH</b> [SSS+22]		MMLU	[HBB+20
	Acc.	Δ	Acc.	Δ
All Data	35.2%		41.2%	
LESS	35.2%	-0.0%	41.8%	+0.5%
MGD-DS	36.7%	+1.5%	42.5%	+1.3%

New goal: Select instruction fine-tuning data for LM benchmarks (2B param model)



	BBH [	SSS+22]	MMLU [HBB+20	
	Acc.	Δ	Acc.	Δ
All Data	35.2%		41.2%	
LESS	35.2%	-0.0%	41.8%	+0.5%
<b>MGD-DS</b>	36.7%	+1.5%	42.5%	+1.3%

Previous methods harm performance in BBH

New goal: Select instruction fine-tuning data for LM benchmarks (2B param model)



2x the improvement in MMLU	J
----------------------------	---

	<b>BBH</b> [SSS+22]		MMLU [HBB+20	
	Acc.	Δ	Acc.	Δ
All Data	35.2%		41.2%	
LESS	35.2%	-0.0%	41.8%	+0.5%
MGD-DS	36.7%	+1.5%	42.5%	+1.3%

Previous methods harm performance in BBH

Insert poison data z





Train model using fixed learning algorithm





Train model using fixed learning algorithm

**Goal:** reduce overall test accuracy

CIFAR10: 85%







Train model using fixed learning algorithm

**Goal:** reduce overall test accuracy



CIFAR10: 85%

**Goal:** Insert poison datapoints that *maximize* <u>overall</u> test loss





### **Goal:** Insert poison datapoints that *maximize* <u>overall</u> test loss

[Lu Kamath Yu '23]



**Metaparameter z**: tensor  $\mathbf{z} \in \mathbb{R}^{\varepsilon N \times 32 \times 32 \times 3} \times \mathbb{R}^{\varepsilon N \times 10}$  of poison <u>pixels and labels</u>

**Metaparameter z**: tensor  $\mathbf{z} \in \mathbb{R}^{\varepsilon N \times 32 \times 32 \times 3} \times \mathbb{R}^{\varepsilon N \times 10}$  of poison <u>pixels and labels</u> **Learning algorithm**  $\mathscr{A}$ : standardized CIFAR-10 training (ResNet)

**Learning algorithm**  $\mathscr{A}$ : standardized CIFAR-10 training (ResNet) **Output function:** Loss  $L_{val}$  evaluated on a heldout validation set

**Metaparameter z**: tensor  $\mathbf{z} \in \mathbb{R}^{\varepsilon N \times 32 \times 32 \times 3} \times \mathbb{R}^{\varepsilon N \times 10}$  of poison pixels and labels
**Learning algorithm**  $\mathscr{A}$ : standardized CIFAR-10 training (ResNet) **Output function:** Loss  $L_{val}$  evaluated on a heldout validation set

Repeat, starting from z being the first  $\varepsilon$ -fraction of the training data:

- **Metaparameter z**: tensor  $z \in \mathbb{R}^{\varepsilon N \times 32 \times 32 \times 3} \times \mathbb{R}^{\varepsilon N \times 10}$  of poison pixels and labels

**Learning algorithm**  $\mathscr{A}$ : standardized CIFAR-10 training (ResNet) **Output function:** Loss  $L_{val}$  evaluated on a heldout validation set

Repeat, starting from z being the first  $\varepsilon$ -fraction of the training data:

Train model, replacing the first  $\varepsilon$ -fraction of the training set with z1.

- **Metaparameter z**: tensor  $z \in \mathbb{R}^{\varepsilon N \times 32 \times 32 \times 3} \times \mathbb{R}^{\varepsilon N \times 10}$  of poison pixels and labels

**Learning algorithm**  $\mathscr{A}$ : standardized CIFAR-10 training (ResNet) **Output function:** Loss  $L_{val}$  evaluated on a heldout validation set

Repeat, starting from z being the first  $\varepsilon$ -fraction of the training data:

- Train model, replacing the first  $\varepsilon$ -fraction of the training set with z
- 1. 2. Calculate metagradient  $\mathbf{g} = \nabla_{\mathbf{z}} \phi(\mathscr{A}(\mathbf{z}))$

- **Metaparameter z**: tensor  $z \in \mathbb{R}^{\varepsilon N \times 32 \times 32 \times 3} \times \mathbb{R}^{\varepsilon N \times 10}$  of poison pixels and labels

**Learning algorithm**  $\mathscr{A}$ : standardized CIFAR-10 training (ResNet) **Output function:** Loss  $L_{val}$  evaluated on a heldout validation set

Repeat, starting from z being the first  $\varepsilon$ -fraction of the training data:

- Train model, replacing the first  $\varepsilon$ -fraction of the training set with z
- 1. 2. Calculate metagradient  $\mathbf{g} = \nabla_{\mathbf{z}} \phi(\mathscr{A}(\mathbf{z}))$
- 3. Update  $\mathbf{z} \leftarrow \mathbf{z} + \eta \cdot \mathbf{g}$

- **Metaparameter z**: tensor  $z \in \mathbb{R}^{\varepsilon N \times 32 \times 32 \times 3} \times \mathbb{R}^{\varepsilon N \times 10}$  of poison pixels and labels

**Learning algorithm**  $\mathscr{A}$ : standardized CIFAR-10 training (ResNet) **Output function:** Loss  $L_{val}$  evaluated on a heldout validation set

Repeat, starting from z being the first  $\varepsilon$ -fraction of the training data:

- Train model, replacing the first  $\varepsilon$ -fraction of the training set with z
- 1. 2. Calculate metagradient  $\mathbf{g} = \nabla_{\mathbf{z}} \phi(\mathscr{A}(\mathbf{z}))$
- 3. Update  $\mathbf{z} \leftarrow \mathbf{z} + \eta \cdot \mathbf{g}$

(this is the full algorithm)

- **Metaparameter z**: tensor  $z \in \mathbb{R}^{\varepsilon N \times 32 \times 32 \times 3} \times \mathbb{R}^{\varepsilon N \times 10}$  of poison pixels and labels





Model	Acc.	Δ
<ul> <li> Original model</li> </ul>	92.0%	_
GradCancel [LKY23]	91.2%	-0.80%
— MGD-DP (ours)	78.1%	-13.9%





Model	Acc.	Δ
Original model	92.0%	_
GradCancel [LKY23]	91.2%	-0.80%
— MGD-DP (ours)	78.1%	-13.9%





# We induce 17x the accuracy drop vs. existing methods

Model	Acc.	Δ
Original model	92.0%	
GradCancel [LKY23]	91.2%	-0.80%
— MGD-DP (ours)	78.1%	-13.9%

Given a **prediction** task:

**Rewrite** problem as a function of a metaparameter 1.

Given a **prediction** task:

1. **Rewrite** problem as a function of a metaparameter Decide on appropriate  $\mathbf{z}, \mathscr{A}, \phi$  so that the problem becomes estimating  $\phi(\mathscr{A}(\mathbf{z}))$ 

- 1. **Rewrite** problem as a function of a metaparameter Decide on appropriate  $\mathbf{z}, \mathscr{A}, \phi$  so that the problem becomes estimating  $\phi(\mathscr{A}(\mathbf{z}))$
- 2. For reference metaparameter  $\mathbf{z}_0$ , compute  $\phi(\mathscr{A}(\mathbf{z}_0))$  and  $\nabla_{\mathbf{z}}\phi(\mathscr{A}(\mathbf{z}_0))$

- 1. **Rewrite** problem as a function of a metaparameter Decide on appropriate  $\mathbf{z}, \mathscr{A}, \phi$  so that the problem becomes estimating  $\phi(\mathscr{A}(\mathbf{z}))$
- 2. For reference metaparameter  $\mathbf{z}_0$ , compute  $\phi(\mathscr{A}(\mathbf{z}_0))$  and  $\nabla_{\mathbf{z}}\phi(\mathscr{A}(\mathbf{z}_0))$
- 3. For any metaparameter **z**, predict  $\phi(\mathscr{A}(\mathbf{z}_0)) + (\mathbf{z} \mathbf{z}_0)^\top \nabla_{\mathbf{z}} \phi(\mathscr{A}(\mathbf{z}_0))$

- **Rewrite** problem as a function of a metaparameter 1. Decide on appropriate  $\mathbf{z}, \mathscr{A}, \phi$  so that the problem becomes estimating  $\phi(\mathscr{A}(\mathbf{z}))$
- 2. For reference metaparameter  $\mathbf{z}_0$ , compute  $\phi(\mathscr{A}(\mathbf{z}_0))$  and  $\nabla_z \phi(\mathscr{A}(\mathbf{z}_0))$
- 3. For **any** metaparameter **z**, predict  $\phi$ ( This is just a linear function of z, no "extra computation" beyond  $z_0$

$$(\mathscr{A}(\mathbf{z}_0)) + (\mathbf{z} - \mathbf{z}_0)^\top \nabla_{\mathbf{z}} \phi(\mathscr{A}(\mathbf{z}_0))$$

Choose train subset



Choose train subset



Train a model



Choose train subset

Train a model







# Measure loss on test example *x*:



Choose train subset



**Goal:** Predict the effect of removing some training data on the model's loss on x

Train a model

Measure loss on test example *x*:







#### Choose train subset



Train a model

Measure loss on test example *x*:





- **Goal:** Predict the effect of removing some training data on the model's loss on x
- **Equivalent goal:** Predict the model's loss on x as a function of training data weighting

Choose train subset



State of the art: predictions that are moderately correlated with model behavior

Train a model

Measure loss on test example *x*:





- **Goal:** Predict the effect of removing some training data on the model's loss on x
- **Equivalent goal:** Predict the model's loss on x as a function of training data weighting

Choose train subset



**Goal:** Predict the effect of removing some training data on the model's loss on x

**Equivalent goal:** Predict the model's loss on x as a function of training data weighting

State of the art: predictions that are moderately correlated with model behavior [Park Georgiev Leclerc Madry '24; Grosse Bae Anil et al. '24; many others]

Train a model

Measure loss on test example *x*:





**Metaparameter z**: vector  $\mathbf{z} \in \mathbb{R}^n$  s.t.,  $\mathbf{z}_i$  is the weight on *i*-th training example

- **Metaparameter z**: vector  $\mathbf{z} \in \mathbb{R}^n$  s.t.,  $\mathbf{z}_i$  is the weight on *i*-th training example **Learning algorithm**  $\mathscr{A}$ : standard model training routine

**Metaparameter z**: vector  $\mathbf{z} \in \mathbb{R}^n$  s.t.,  $\mathbf{z}_i$  is the weight on *i*-th training example

**Learning algorithm**  $\mathscr{A}$ : standard model training routine

**Output function:** Loss function  $\ell(\cdot; x)$  evaluated on a specific test example

**Metaparameter z**: vector  $\mathbf{z} \in \mathbb{R}^n$  s.t.,  $\mathbf{z}_i$  is the weight on *i*-th training example

**Learning algorithm**  $\mathscr{A}$ : standard model training routine

**Output function:** Loss function  $\ell(\cdot; x)$  evaluated on a specific test example

Choose reference parameter to be all-ones vector and Taylor expand

**Metaparameter z**: vector  $\mathbf{z} \in \mathbb{R}^n$  s.t.,  $\mathbf{z}_i$  is the weight on *i*-th training example

**Learning algorithm**  $\mathscr{A}$ : standard model training routine

**Output function:** Loss function  $\ell(\cdot; x)$  evaluated on a specific test example

Choose reference parameter to be all-ones vector and Taylor expand

 $\phi(\mathscr{A}(\mathbf{Z})) \approx \phi(\mathscr{A}((\mathbf{1}_n)))$ 

$$+ \frac{d\phi(\mathscr{A}(\mathbf{z}))}{d\mathbf{z}} \Big|_{\mathbf{z}=\mathbf{1}_n} (\mathbf{z}-\mathbf{1}_n)$$

**Metaparameter z**: vector  $\mathbf{z} \in \mathbb{R}^n$  s.t.,  $\mathbf{z}_i$  is the weight on *i*-th training example

**Learning algorithm**  $\mathscr{A}$ : standard model training routine

**Output function:** Loss function  $\ell(\cdot; x)$  evaluated on a specific test example

Choose reference parameter to be all-ones vector and Taylor expand

 $\phi(\mathscr{A}(\mathbf{z})) \approx \phi(\mathscr{A}((\mathbf{1}_n)))$ 



$$+ \frac{d\phi(\mathscr{A}(\mathbf{z}))}{d\mathbf{z}} \Big|_{\mathbf{z}=\mathbf{1}_n} (\mathbf{z}-\mathbf{1}_n)$$

AKA the <u>exact</u> influence function (Hampel '47)

Deleting 1% random subsets from CIFAR-10 training:

#### Deleting 1% random subsets from CIFAR-10 training:





#### Deleting 1% random subsets from CIFAR-10 training:


## **Task: Training data attribution** Results

#### Deleting 1% random subsets from CIFAR-10 training:



## **Task: Training data attribution** Results

### Deleting 1% random subsets from CIFAR-10 training:





#### Also works for LMs!

#### Deleting 1% of data from Gemma-2B fine-tuning



## **Task: Training data attribution** Results

### Deleting 1% random subsets from CIFAR-10 training:



#### (From a different paper; online soon!)



#### Also works for LMs!

#### Deleting 1% of data from Gemma-2B fine-tuning



## Roadmap

Main idea: The metagradient

#### **Prequel: Two challenges**

#### Looking forward

#### Exact metagradients at small scale [Maclaurin Duvenaud Adams '15; Franceschi et al. '17; ...] Optimized LR, weight decay, init, etc. for 100-iteration MNIST training

Largest-scale setting: four-layer network on phone classification dataset

## Exact metagradients at small scale [Maclaurin Duvenaud Adams '15; Franceschi et al. '17; ...]

Optimized LR, weight decay, init, etc. for 100-iteration MNIST training

Largest-scale setting: four-layer network on phone classification dataset

### Approximate (but large-scale) metagradients via implicit differentiation [Lorraine et al. '19]

estimation & cannot optimize optimizer parameters such as learning rate

- Scalable but coarse approximation based on convergence assumption; requires Hessian

#### Exact metagradients at small scale [Maclaurin Duvenaud Adams '15; Franceschi et al. '17; ...] Optimized LR, weight decay, init, etc. for 100-iteration MNIST training

Largest-scale setting: four-layer network on phone classification dataset

#### Approximate (but large-scale) metagradients via implicit differentiation [Lorraine et al. '19]

estimation & cannot optimize optimizer parameters such as learning rate

#### **Note:** These are just three representative works; significant line of work on differentiating through optimization [e.g., Domke '12; Nichol et al. '18; Hara '19; Shaban et al. '19; ...]

- Scalable but coarse approximation based on convergence assumption; requires Hessian

#### Exact metagradients at small scale [Maclaurin Duvenaud Adams '15; Franceschi et al. '17; ...] Optimized LR, weight decay, init, etc. for 100-iteration MNIST training

Largest-scale setting: four-layer network on phone classification dataset

#### Approximate (but large-scale) metagradients via implicit differentiation [Lorraine et al. '19]

estimation & cannot optimize optimizer parameters such as learning rate

**Note:** These are just three representative works; significant line of work on differentiating through optimization [e.g., Domke '12; Nichol et al. '18; Hara '19; Shaban et al. '19; ...]

- Scalable but coarse approximation based on convergence assumption; requires Hessian

#### Broadly, nothing has scaled to "medium"-size training (e.g., >10M params, >1k iters)



## Two barriers to optimizing with the metagradient $\nabla_z \phi(\mathscr{A}(\mathbf{z}))$



#### **#1:** Computational tractability

## Two barriers to optimizing with the metagradient $\nabla_z \phi(\mathscr{A}(\mathbf{z}))$



#### **#1:** Computational tractability

## Two barriers to optimizing with the metagradient $\nabla_z \phi(\mathscr{A}(\mathbf{z}))$



**#2:** Optimization (non)smoothness



#### **#1:** Computational tractability

## Two barriers to optimizing with the metagradient $\nabla_z \phi(\mathscr{A}(\mathbf{z}))$



**#2:** Optimization (non)smoothness



## Two barriers to optimizing with the metagradient $\nabla_z \phi(\mathscr{A}(\mathbf{z}))$

**#2:** Optimization (non)smoothness

# Problem #1: Calculating the gradient

Idea: Direct autodiff (e.g., Domke et al., 2012)

**Direct approach:** Auto-differentiation (AD)

**Direct approach:** Auto-differentiation (AD)

Take derivatives of any function defined by simple operations by saving intermediate products

**Direct approach:** Auto-differentiation (AD)

Take derivatives of any function defined by simple operations by saving intermediate products





d = a \* c

**Direct approach:** Auto-differentiation (AD)

Take derivatives of any function defined by simple operations by saving intermediate products

Key insight: Training an ML model is just a (very, very, very long) sequence of operations! Maybe we can just backprop through it?





- **Direct approach:** Auto-differentiation (AD)
  - Take derivatives of any function defined by simple operations by saving intermediate products
- Key insight: Training an ML model is just a (very, very, very long) sequence of operations! Maybe we can just backprop through it?
- **Problem:** Too many intermediate products!

d = a \* cc = a + bc



**Direct approach:** Auto-differentiation (AD)

Take derivatives of any function defined by simple operations by saving intermediate products

Key insight: Training an ML model is just a (very, very, very long) sequence of operations! Maybe we can just backprop through it?

**Problem:** Too many intermediate products!

We will think of AD as a black box that can compute gradients of a many  $\rightarrow$  one compute graph O(# edges) memory



**Naive approach:** AD on full model training - too many "edges" on compute graph to store!

**Naive approach:** AD on full model training - too many "edges" on compute graph to store!





Naive approach: AD on full model training - too many "edges" on compute graph to store!



**Stepwise AD insight:** step t derivatives depend only on step t + 1 derivatives, step t state









**Stepwise AD insignt:** step t derivatives depend only on step t + 1 derivatives, step t state

df(z)	df(z)	$d\theta_{t+1}$
$\frac{d\theta_t}{d\theta_t}$	$\frac{1}{d\theta_{t+1}}$	$d\theta_t$











A simple modification: REPLAY

#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order

#### A simple modification: REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order

(Maclaurin Duvenaud & Adams, 2015)

#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Deterministic training allows us to "replay" from a previous state to rematerialize

#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Deterministic training allows us to "replay" from a previous state to rematerialize Design a data structure that allows us to trade off time and space:

#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Deterministic training allows us to "replay" from a previous state to rematerialize Design a data structure that allows us to trade off time and space:

**Step-wise AD:** Save every state, traverse states in reverse
#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Deterministic training allows us to "replay" from a previous state to rematerialize Design a data structure that allows us to trade off time and space:





#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Deterministic training allows us to "replay" from a previous state to rematerialize Design a data structure that allows us to trade off time and space:





#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Deterministic training allows us to "replay" from a previous state to rematerialize Design a data structure that allows us to trade off time and space:





#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Deterministic training allows us to "replay" from a previous state to rematerialize Design a data structure that allows us to trade off time and space:





#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Deterministic training allows us to "replay" from a previous state to rematerialize Design a data structure that allows us to trade off time and space:





#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Deterministic training allows us to "replay" from a previous state to rematerialize Design a data structure that allows us to trade off time and space:





#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Deterministic training allows us to "replay" from a previous state to rematerialize Design a data structure that allows us to trade off time and space:





#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Deterministic training allows us to "replay" from a previous state to rematerialize Design a data structure that allows us to trade off time and space:





#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Deterministic training allows us to "replay" from a previous state to rematerialize Design a data structure that allows us to trade off time and space:





#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Deterministic training allows us to "replay" from a previous state to rematerialize Design a data structure that allows us to trade off time and space:





#### **A simple modification:** REPLAY

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Deterministic training allows us to "replay" from a previous state to rematerialize Design a data structure that allows us to trade off time and space:

Idea #1: Save every k-th state, "hop" back to saved states and replay to traverse  $(T/k + k \text{ space} \rightarrow 2\sqrt{T} \text{ space}, \text{ T extra training steps})$ 





#### **A simple modification:** REPLAY (Recursive rematerialization for metagradients)

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Design a data structure that allows us to trade off time and space:



- Deterministic training allows us to "replay" from a previous state to rematerialize



#### **A simple modification:** REPLAY (Recursive rematerialization for metagradients)

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Design a data structure that allows us to trade off time and space:

Idea #2: Recursively apply same idea across k training segments



- Deterministic training allows us to "replay" from a previous state to rematerialize



#### **A simple modification:** REPLAY (Recursive rematerialization for metagradients)

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Design a data structure that allows us to trade off time and space:

Idea #2: Recursively apply same idea across k training segments



- Deterministic training allows us to "replay" from a previous state to rematerialize



#### **A simple modification:** REPLAY (Recursive rematerialization for metagradients)

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Design a data structure that allows us to trade off time and space:

**Idea #2:** Recursively apply same idea across k training segments



- Deterministic training allows us to "replay" from a previous state to rematerialize



#### **A simple modification:** REPLAY (Recursive rematerialization for metagradients)

Key primitive of step-wise AD: **traversing** the states  $\theta_T \dots \theta_1$  in reverse order (Maclaurin Duvenaud & Adams, 2015)

Design a data structure that allows us to trade off time and space:

**Idea #2:** Recursively apply same idea across k training segments



(For any k, uses  $k \log_k(T)$  space,  $T \log_k(T)$  extra training steps)

- Deterministic training allows us to "replay" from a previous state to rematerialize



#### A simple modification: REPLAY (Recursive rematerialization for metagradients)

ø.

Sasha Rush 🗹

Key p

Deter

Desic

 $\theta_1$ 

Idea #2:

I really want to figure out an application for Recursive Halving AD M (picture from Elements of Differentiable Programming) arxiv.org/pdf/2403.14606). It feels like it should be useful for something cool. Maybe like on-device test-time derivatives...



Figure 8.8: Illustration of checkpointing with recursive halving, for a chain of 8 functions. The chain is first fully evaluated while storing some computations as checkpoints in memory. Then, during the backward pass, we recompute some intermediate values from the latest checkpoint available. In contrast, vanilla reversemode autodiff (with full caching of the intermediate computations) would lead to a simple triangle shape

- ing the states  $\theta_T \dots \theta_1$  in reverse order
- ay" from a previous state to rematerialize o trade off time and space:



(For any k, uses  $k \log_k(T)$  space,  $T \log_k(T)$  extra training steps)

# Metagradient roadmap: Two challenges



**#1:** Computational tractability

Two barriers to optimizing with the metagradient  $\nabla_{\tau}\phi(\mathscr{A}(z))$ 



**#2:** Optimization (non)smoothness

# Metagradient roadmap: Two challenges



**#1:** Computational tractability

Two barriers to optimizing with the metagradient  $\nabla_{\tau} \phi(\mathscr{A}(z))$ 



#### A gradient is only "useful" if it locally predicts function behavior In optimization terms, we want $\phi(\mathscr{A}(\mathbf{z}))$ to be "smooth"

A gradient is only "useful" if it locally predicts function behavior In optimization terms, we want  $\phi(\mathscr{A}(\mathbf{z}))$  to be "smooth"

Given the complexity of  $\mathscr{A}(\mathbf{z})$  for large-scale models, cannot take this for granted!

A gradient is only "useful" if it locally predicts function behavior In optimization terms, we want  $\phi(\mathscr{A}(\mathbf{z}))$  to be "smooth"

Given the complexity of  $\mathscr{A}(\mathbf{z})$  for large-scale models, cannot take this for granted!

Indeed:

### A gradient is only "useful" if it locally predicts function behavior In optimization terms, we want $\phi(\mathscr{A}(\mathbf{z}))$ to be "smooth"

Given the complexity of  $\mathscr{A}(\mathbf{z})$  for large-scale models, cannot take this for granted!



Small changes in training data pixels!

### A gradient is only "useful" if it locally predicts function behavior In optimization terms, we want $\phi(\mathscr{A}(\mathbf{z}))$ to be "smooth"

Given the complexity of  $\mathscr{A}(\mathbf{z})$  for large-scale models, cannot take this for granted!



Small changes in training data pixels!

...standard model training is very "non-smooth!"

A gradient is only "useful" if it locally predicts function behavior In optimization terms, we want  $\phi(\mathscr{A}(\mathbf{z}))$  to be "smooth"

Given the complexity of  $\mathscr{A}(\mathbf{z})$  for large-scale models, cannot take this for granted!



Small changes in training data pixels!

...standard model training is very "non-smooth!" **Need to find "smooth" model training routines!** 



**Problem:** need a cheap + accurate measure of "smoothness"

- **Goal**: design model training routines that *do* have "smooth" optimization landscapes

**Goal**: design model training routines that *do* have "smooth" optimization landscapes

**Problem:** need a cheap + accurate measure of "smoothness"

**Desired diagnostic ("Smoothness"):** Approximate Taylor predictiveness  $\approx \max_{\Delta \mathbf{w}} \| \mathscr{A}(\mathbf{w} + \Delta \mathbf{w}) \|$ 

$$\mathbf{w}) - \mathscr{A}(\mathbf{w}) - \frac{d\mathscr{A}(\mathbf{w})}{d\mathbf{w}} \cdot \Delta \mathbf{w}$$

**Goal**: design model training routines that *do* have "smooth" optimization landscapes

**Problem:** need a cheap + accurate measure of "smoothness"

**Desired diagnostic ("Smoothness"):** Approximate Taylor predictiveness  $\approx \max_{\Delta \mathbf{w}} \| \mathscr{A}(\mathbf{w} + \Delta \mathbf{w})$ 

$$\mathbf{w}) - \mathscr{A}(\mathbf{w}) - \frac{d\mathscr{A}(\mathbf{w})}{d\mathbf{w}} \cdot \Delta \mathbf{w} \|$$

**First-order effect of metaparameter changes** 

**Goal**: design model training routines that *do* have "smooth" optimization landscapes

**Problem:** need a cheap + accurate measure of "smoothness"

**Desired diagnostic ("Smoothness"):** Approximate Taylor predictiveness

Can't take a max!

max

 $\mathscr{A}(\mathbf{w} + \Delta \mathbf{w})$ 

$$\mathbf{w}) - \mathscr{A}(\mathbf{w}) - \frac{d\mathscr{A}(\mathbf{w})}{d\mathbf{w}} \cdot \Delta \mathbf{w} \|$$

**First-order effect of metaparameter changes** 

**Goal**: design model training routines that *do* have "smooth" optimization landscapes

**Problem:** need a <u>cheap</u> + accurate measure of "smoothness"

**Desired diagnostic ("Smoothness"):** Approximate Taylor predictiveness

max

Can't take a max!

$$\mathscr{A}(\mathbf{w} + \Delta \mathbf{w}) - \mathscr{A}(\mathbf{w}) - \frac{d\mathscr{A}(\mathbf{w})}{d\mathbf{w}} \cdot \Delta \mathbf{w}$$

First-order effect of metaparameter changes Hard to interpret this quantity!

**Goal**: design model training routines that *do* have "smooth" optimization landscapes

**Problem:** need a <u>cheap</u> + accurate measure of "smoothness"

**Desired diagnostic ("Smoothness"):** Approximate Taylor predictiveness

max

Can't take a max!



Don't have full Jacobian!

$$\mathscr{A}(\mathbf{w} + \Delta \mathbf{w}) - \mathscr{A}(\mathbf{w}) - \frac{d\mathscr{A}(\mathbf{w})}{d\mathbf{w}} \cdot \Delta \mathbf{w}$$

First-order effect of metaparameter changes Hard to interpret this quantity!

**Goal**: design model training routines that *do* have "smooth" optimization landscapes

**Problem:** need a cheap + accurate measure of "smoothness"

**Desired diagnostic ("Smoothness"):** Approximate Taylor predictiveness Can't take a max!  $\mathscr{A}(\mathbf{w} + \Delta \mathbf{w}) - \mathscr{A}(\mathbf{w})$ max Hard to interpret this quantity!



 $d\mathscr{A}(\mathbf{W})$ 

Don't have full Jacobian!

 $d\mathbf{W}$ First-order effect of metaparameter changes

- We adapt this estimator to a new metric Metasmoothness  $(\mathscr{A}) \in [-1,1]$ 
  - **Our final estimator:** feasible, interpretable, and accurate!

Our metric accurately predicts optimization success in practice!
















How to train your smooth model:



How to train your *smooth* model:

Layer composition, width



How to train your *smooth* model:

Layer composition, width Learning rate/weight decay/batch size



How to train your *smooth* model:

Layer composition, width Learning rate/weight decay/batch size Adam  $\epsilon_{root}$ , momentum/warmup



How to train your smooth model:

Layer composition, width Learning rate/weight decay/batch size Adam  $\epsilon_{root}$ , momentum/warmup Normalizer  $\epsilon$ 



How to train your smooth model:

Layer composition, width Learning rate/weight decay/batch size Adam  $\epsilon_{root}$ , momentum/warmup Normalizer  $\epsilon$ Logit scales



## Our metric accurately predicts optimization success in practice!

#### How to train your smooth model:

Layer composition, width Learning rate/weight decay/batch size Adam  $\epsilon_{root}$ , momentum/warmup Normalizer  $\epsilon$ Logit scales LM only:



## Our metric accurately predicts optimization success in practice!

#### How to train your smooth model:

Layer composition, width Learning rate/weight decay/batch size Adam  $\epsilon_{root}$ , momentum/warmup Normalizer  $\epsilon$ Logit scales LM only: QK-layernorm



## Our metric accurately predicts optimization success in practice!

#### How to train your smooth model:

Layer composition, width Learning rate/weight decay/batch size Adam  $\epsilon_{\rm root}$ , momentum/warmup Normalizer  $\epsilon$ Logit scales LM only: QK-layernorm Uncoupled (un)embedding



How to train your smooth model:

Layer composition, width Learning rate/weight decay/batch size Adam  $\epsilon_{root}$ , momentum/warmup Normalizer  $\epsilon$ 

Logit scales

**Key property:** Optimized parameters "transfer back" to non-smooth training!

Uncoupled (un)embedding

## Our metric accurately predicts optimization success in practice!



Smoothness

"Smooth model training"



# Roadmap

Main idea: The metagradient

Prequel: Two challenges

#### Looking forward

# Future work

#### Our goal/hope is for the following recipe to be a general one:

- Pick a design choice you're interested in optimizing/predicting 1.
- 2. Write problem as a metaparameter prediction or optimization problem
- 3. Find a smooth training setup
- 4. Compute the metagradient
- 5. Apply gradient descent/Taylor approximation

#### What do we need to get there?

# Future work

#### What do we need to get there?

- 1. A more systematic way of finding smooth models
- 2. More computational efficiency gains
- 3. More creative ways of applying metagradients
  - Hyperparameter optimization? (Simple example in the paper)
  - Data weighting networks?
  - Synthetic data generators?
  - Gradient pre-conditioners?

#### "Anything we can continuously parameterize, we can (try to) optimize"

# Conclusion

Design space of ML is too large to explore  $\rightarrow$  how do we think about design choices? **Approach**: Take the **gradient** of a model output with respect to metaparameters

#### **Challenges:**

Computing the metagradient at scale Making the metagradient useful for optimization

#### Addressing these challenges unlocks:

Data selection (> 2x more effective) & attribution

Data poisoning (> 10x more effective)

Gradient-based hyperparameter optimization

What next?

See the paper for more!



<u>https://arxiv.org/abs/2503.13751</u>

andrewi@stanford.edu

andrewilyas.com