*A discussion on*

# Eve

## A Gradient Based Optimization Method with Locally and Globally Adaptive Learning Rates

Presented by     David Ryan MacDonald
Toronto AI

Paper:   https://arxiv.org/abs/1611.01505

Eve

authors

**Jayanth Koushik**
PhD student at CMU,          H-index: 4
@jayanth_koushik

**Graham Neubig**
Assistant Professor of CS at CMU,   H-index: 26
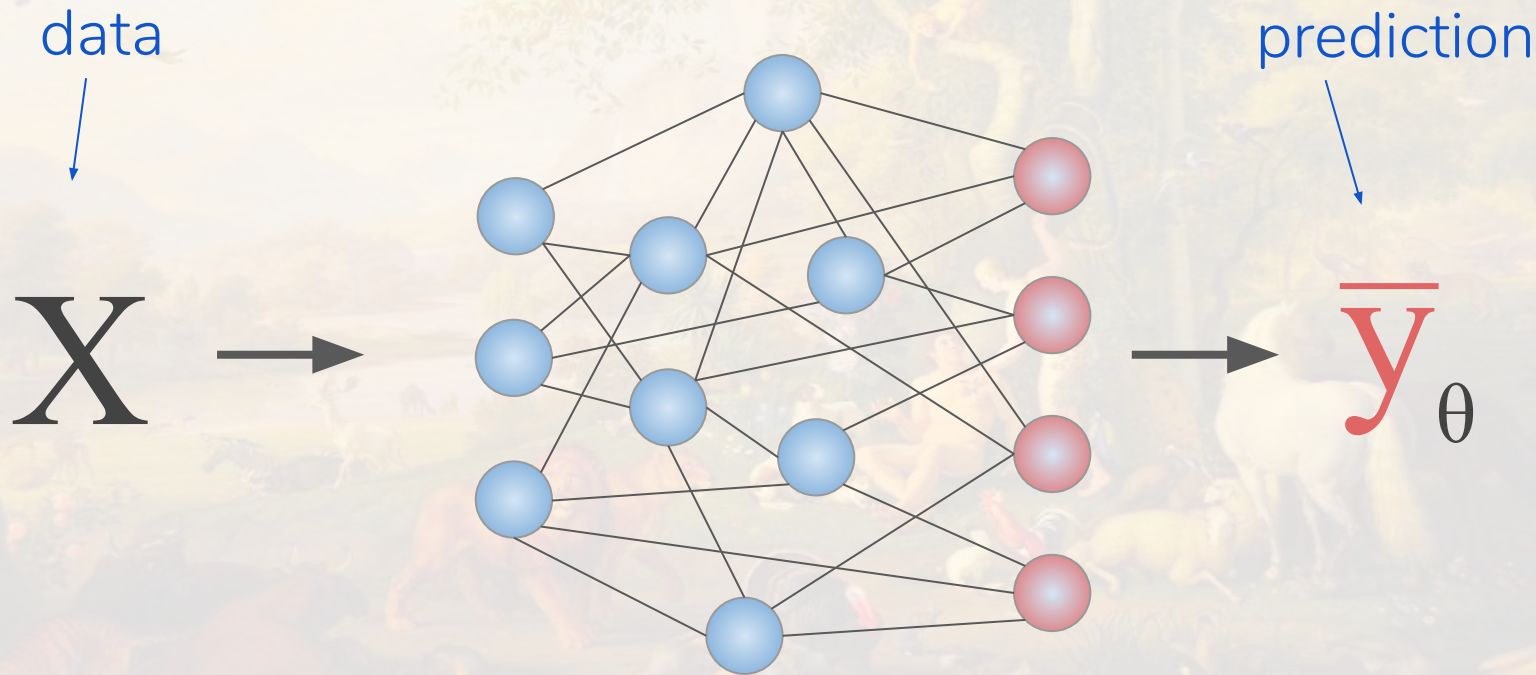@gneubig

**Hiroaki Hayashi**
student at CMU
@hiroberri

Paper:   https://arxiv.org/abs/1611.01505

# Eve



Adam and Eve in the Garden of Eden
Wenzel Peter
(Karlsbad 1745 - Rome 1829)

**Model**

$$f_\theta(X) \longrightarrow \overline{y}_\theta$$

# Goal:  for predictions to match ground truth

$$\overline{y}_\theta \longrightarrow y$$

- we want our output predictions to converge to the ground truth
- The gap between our predictions and our target is the called the **loss**

# Losses

poisson

mean squared error

huber

kullback-leibler divergence

cross entropy

endless variety of others...

cosine proximity

pairwise distance

hinge loss

# Loss

$$\mathcal{L}_\theta(\overline{y}_\theta, y)$$

**Rephrasing our goal**:

- We want to update the parameters of our model, to find a parameterization that minimizes the measured loss (our distance from the target)
- By minimizing the loss, our predicted values approach the ground truth labels.

# Stochastic Gradient Descent

(learning rate alpha)
step size

(i.e. jacobian/gradient of the loss)
step direction

$$\theta_{t+1} = \theta_t - \alpha \nabla_\theta \mathcal{L}_\theta(\bar{y}, y)$$
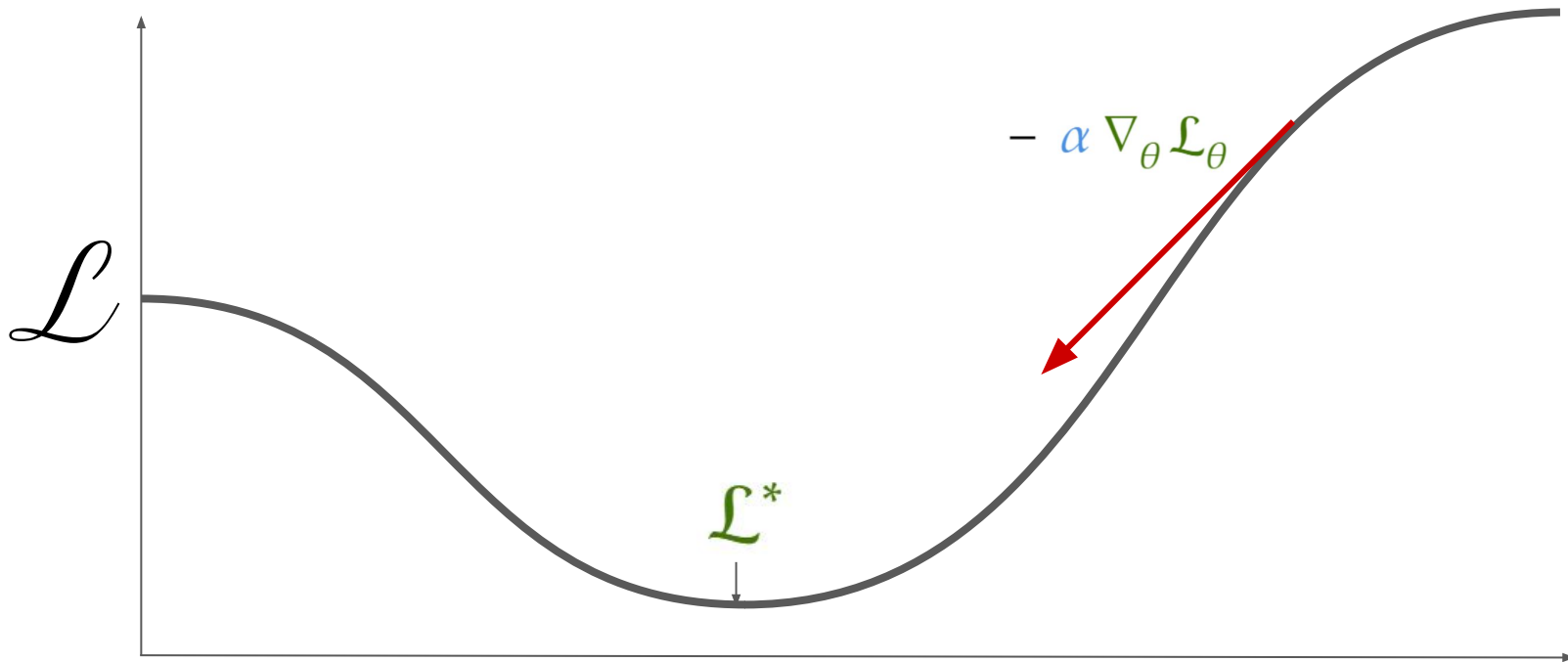
new position in parameter space

current position in parameter space
(i.e. weights & biases)
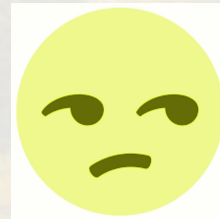
head downhill
(head to regions of lower loss)

# Parameter update

$$\theta_{t+1} = \theta_t - \alpha \nabla_\theta \mathcal{L}_\theta(\bar{y}, y)$$

$$- \alpha \nabla_\theta \mathcal{L}_\theta$$
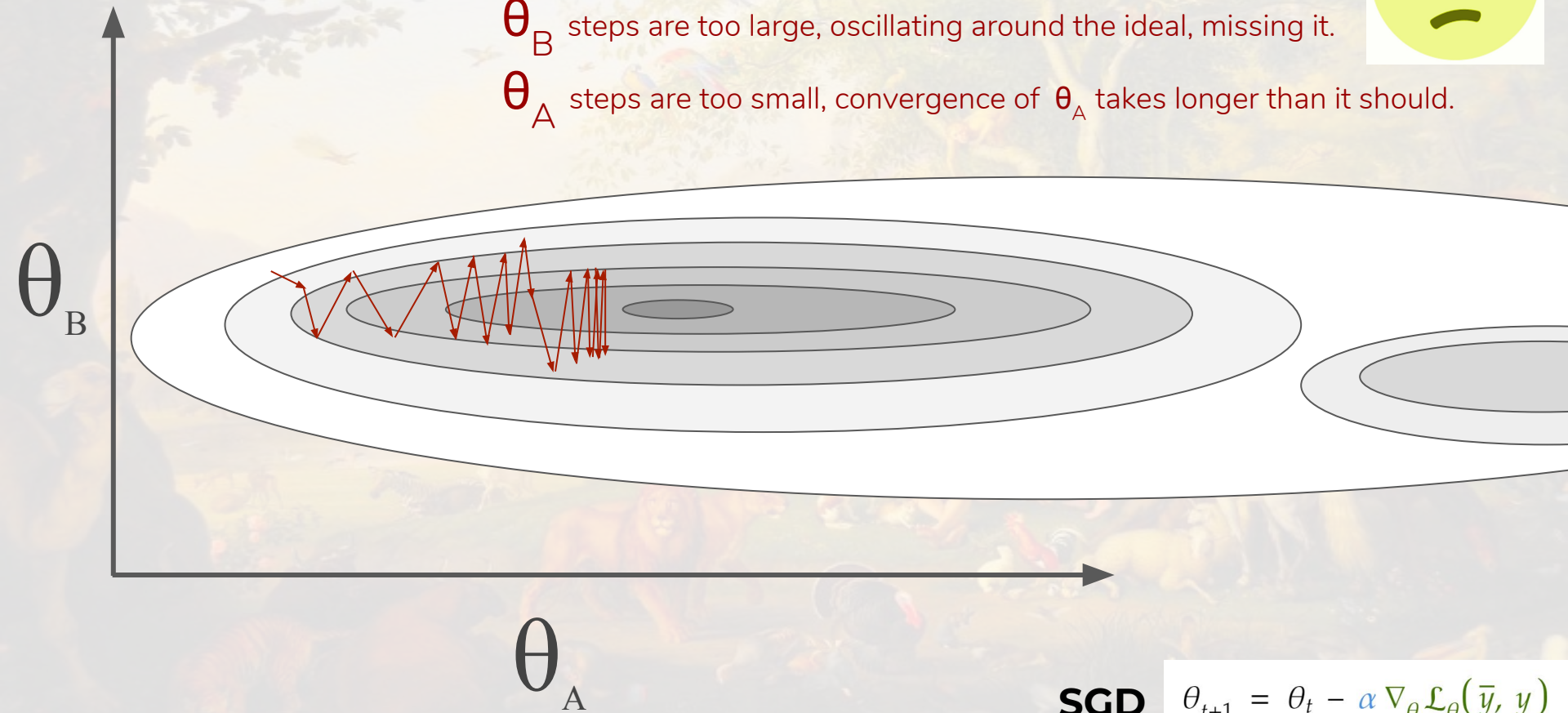
$$\mathcal{L}$$

$$\mathcal{L}^*$$

# Parameter space

$\theta_B$ steps are too large, oscillating around the ideal, missing it.

$\theta_A$ steps are too small, convergence of $\theta_A$ takes longer than it should.
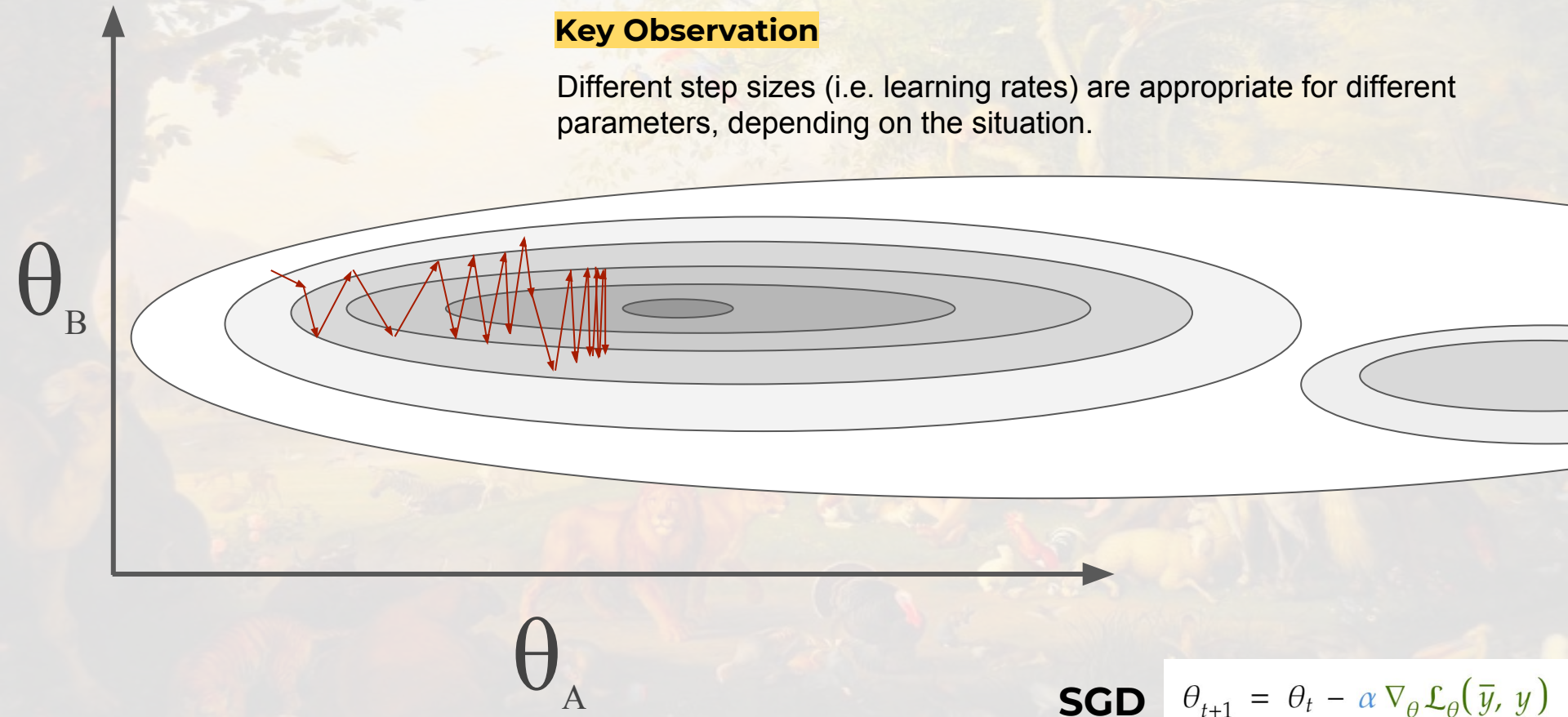


$\theta_B$

$\theta_A$

**SGD** $\quad \theta_{t+1} = \theta_t - \alpha \nabla_\theta \mathcal{L}_\theta(\bar{y}, y)$

# How do we improve this?

**Key Observation**

Different step sizes (i.e. learning rates) are appropriate for different parameters, depending on the situation.



$\theta_B$

$\theta_A$

**SGD** $\quad \theta_{t+1} = \theta_t - \alpha \nabla_\theta \mathcal{L}_\theta(\bar{y}, y)$

# Solution: Adaptive learning rates

- Give each parameter their own learning rate


- ADAGRAD family of algorithms:
  - RMSProp
  - Adadelta
  - Adam


- Adam is a very popular alternative to SGD

**SGD:**

$$\theta_{t+1} = \theta_t - \alpha \nabla_t$$

**Adam:**

$$\theta_{t+1} = \theta_t - \alpha \frac{average\ recent\ gradient}{average\ recent\ deviation\ in\ the\ gradient}$$

# Adam

$$\theta_{t+1} = \theta_t - \alpha \, \frac{average \; recent \; gradient}{average \; recent \; deviation \; in \; the \; gradient}$$

**Unpacking 'average recent gradient':**

- SGD always takes steps in the 'downhill' direction, based on its immediate local position in the loss landscape.

- But the loss landscape is noisy! The loss landscape in parameter space that the optimizer must navigate through totally depends on the input data (loss -> predicted values -> training data). Since our input data changes with every batch we send into the network, so does our loss landscape!

- Adam attempts to cancel this noise by adding momentum to its trajectory through loss space. It does this by heading in the direction based on the running average of recent gradients, instead of in the direction of the current (potentially noisy) gradient

# Adam

$$\theta_{t+1} = \theta_t - \alpha \, \frac{momentum}{average\ recent\ deviation\ in\ the\ gradient}$$

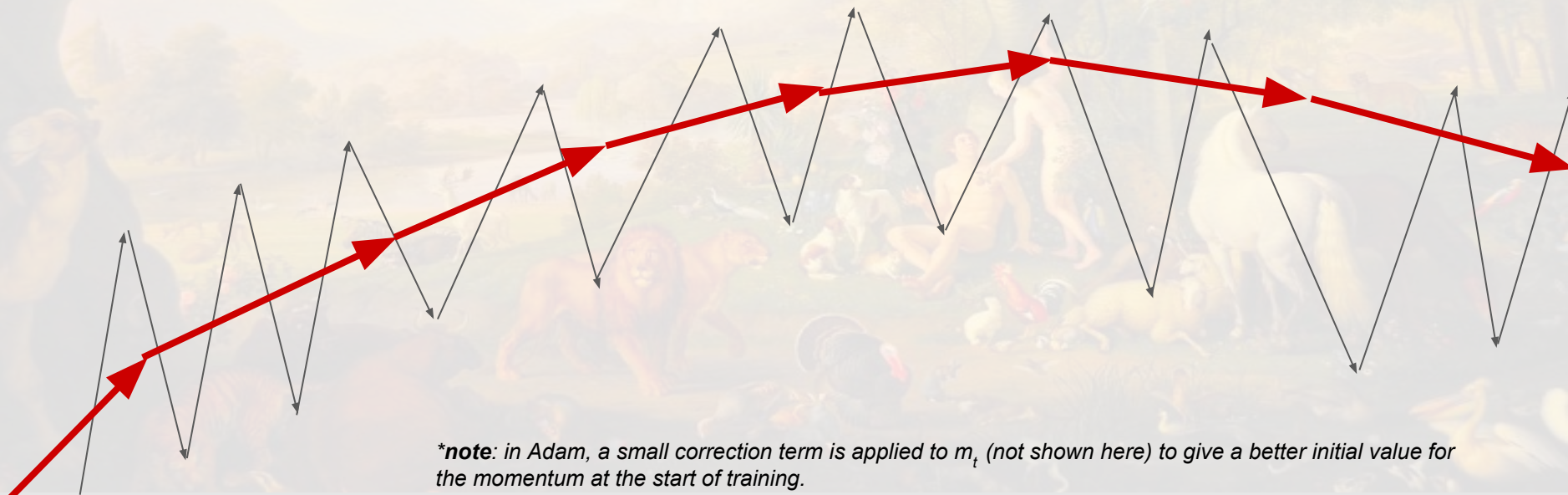Adam's momentum:
smoothes out the impact of noisy gradients

# Adam's Momentum

$$m_t = \beta_1 m_{t-1} + (1-\beta_1)\nabla_t$$

$$suppose\ \beta_1 = 99\%:$$

$$m_t = 99\% \cdot m_{t-1} + 1\% \cdot \nabla_t$$

$$\theta_{t+1} = \theta_t - \alpha \frac{m_t}{average\ recent\ deviation\ in\ the\ gradient}$$

*note: in Adam, a small correction term is applied to $m_t$ (not shown here) to give a better initial value for the momentum at the start of training.

# Adam

$$\theta_{t+1} = \theta_t - \alpha \frac{momentum}{average\ recent\ deviation\ in\ the\ gradient}$$

## The denominator

- The denominator term allows the Adam optimizer to slow training after entering an area of the parameter space where the loss landscape becomes very steep.  Taking big steps is risky in steep terrain.

- Similarly, when a parameter enters a 'valley' in the loss landscape, the denominator allows this parameter to speed up training, allowing for bigger steps to be taken.

# Adam

$$\theta_{t+1} = \theta_t - \alpha \frac{m_t}{\sqrt{v_t} + \epsilon}$$

**Denominator**

$$v_t = \beta_2 v_{t-1} + (1-\beta_2)\nabla_t^2$$

- $v_t$ is used to regulate the learning rate based on a moving average of the squared gradient for the parameter
- The moving average is controlled by our second hyperparameter, **$\beta_2$**
- This moving average is usually an order of magnitude slower to change than the momentum. (i.e. by default **$\beta_1$** = 99% while **$\beta_2$**= 99.9%)

- epislon is used to put an upper limit on the 'speed gain' in shallow regions

***note**: in Adam, a small correction term is applied to $v_t$ (not shown here) to give a better initial value for the moving average of the squared gradients at the start of training.*

# Adam

$$\theta_{t+1} = \theta_t - \alpha \frac{m_t}{\sqrt{v_t} + \epsilon}$$

# Adam

$$\theta_{t+1} = \theta_t - \alpha \frac{m_t}{\sqrt{v_t} + \epsilon}$$

# Eve

$$\theta_{t+1} = \theta_t - \frac{\alpha}{d_t} \frac{m_t}{\sqrt{v_t} + \epsilon}$$

# Eve

$$\theta_{t+1} = \theta_t - \frac{\alpha}{d_t} \frac{m_t}{\sqrt{v_t} + \epsilon}$$

New term

- With Eve, we dynamically adjust the learning rate based on how close we are to the objective.

# Eve - Intuitions

$$\theta_{t+1} = \theta_t - \frac{\alpha}{d_t} \frac{m_t}{\sqrt{v_t} + \epsilon}$$

1. A large variation in the loss between timesteps should be given less weight, and so a smaller step should be taken.

2. If we are far from the minimum, take big steps.

# Eve

$$\theta_{t+1} = \theta_t - \frac{\alpha}{d_t} \frac{m_t}{\sqrt{v_t} + \epsilon}$$

1. A large variation in the loss between timesteps
   → smaller steps

Thus we want $\quad d_t \propto |\mathcal{L}_t - \mathcal{L}_{t-1}|$

# Eve

$$\theta_{t+1} = \theta_t - \frac{\alpha}{d_t} \frac{m_t}{\sqrt{v_t} + \epsilon}$$

2. If the current loss $\mathcal{L}$ is far from the minimum $\mathcal{L}^*$, take big steps.

$$\text{Thus,} \quad d_t \propto \frac{1}{\mathcal{L}_t - \mathcal{L}^*}$$

# Eve   Putting both intuitons together

$$\theta_{t+1} = \theta_t - \frac{\alpha}{d_t} \frac{m_t}{\sqrt{v_t} + \epsilon}$$

$$\frac{\alpha}{d_t} \propto \alpha \frac{\mathcal{L}_t - \mathcal{L}^*}{|\mathcal{L}_t - \mathcal{L}_{t-1}|}$$

A large variation in the loss between timesteps → smaller steps

If the current loss $\mathcal{L}$ is far from the minimum $\mathcal{L}^*$, take big steps.

# Unstable

$\theta_{t+1} = \theta$ **KABOOM**

$$\frac{\alpha}{d_t} \propto \alpha \frac{\mathcal{L}_t - \mathcal{L}^*}{|\mathcal{L}_t - \mathcal{L}_{t-1}|}$$

In the numerator, we want to take big steps when we are far from the minimum.

However, this is at risk of a positive fe___ ___ increases, then ___ ___es, which can ___ ___ the minim___ ___es the step size, wh___ ___n further from the minim___ __…
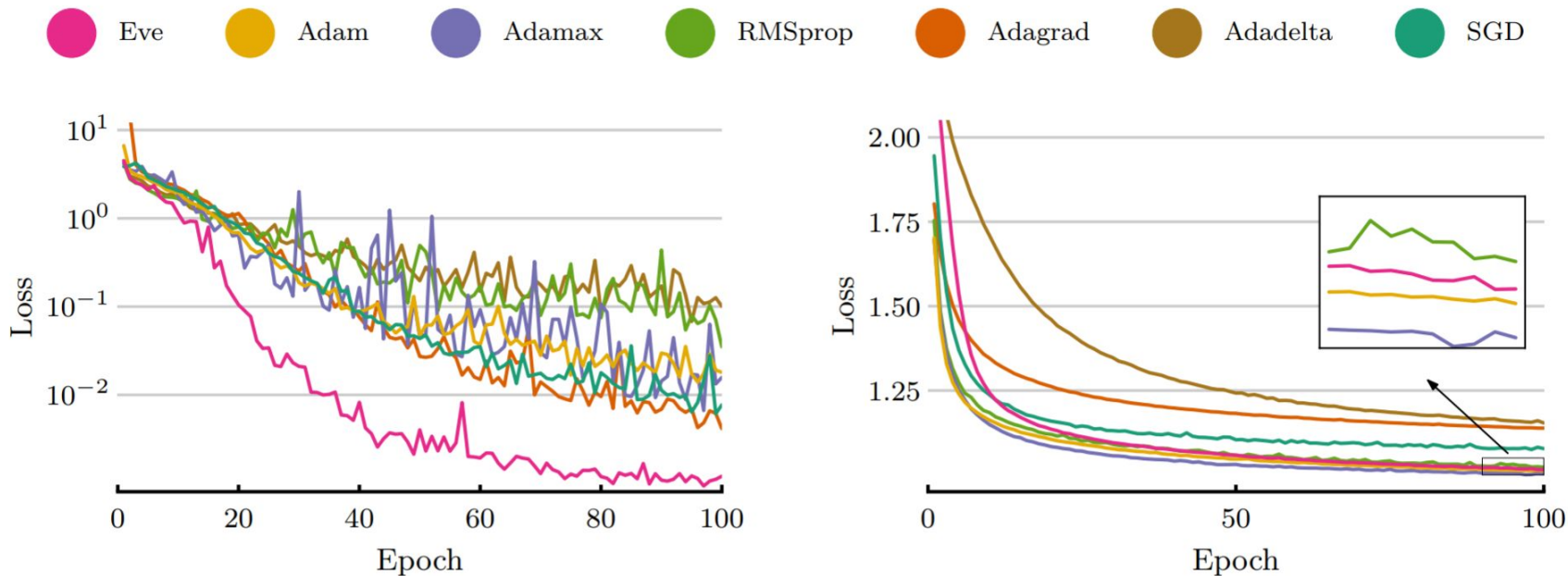
**KABOOM**

**Eve** Experiments

# **Eve**   Optimizer comparison



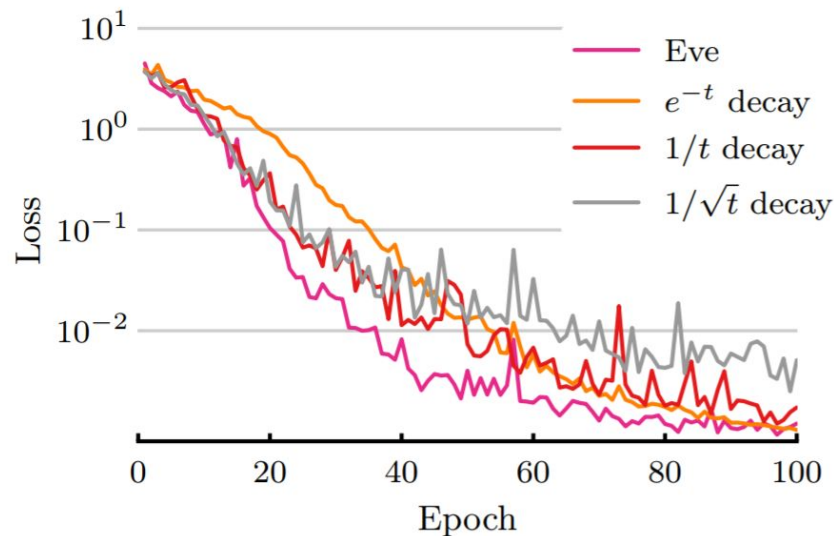Training loss comparison. In both experiments, Eve achieves similar or lower loss than other optimizers.
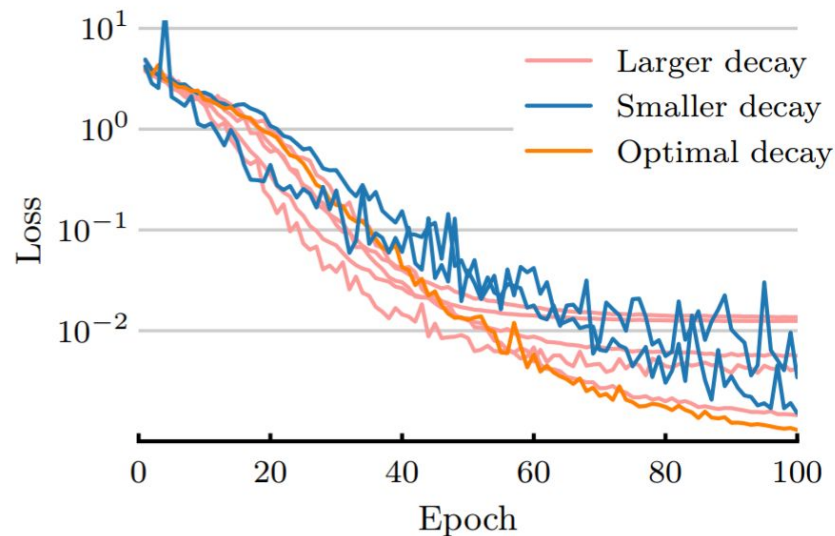
# **Eve**    Comparison with decay strategies



**a.** Eve compared with different decay strategies.

**b.** Adam with different exponential decay strengths.

# Keras implementation

tdeboissiere / **DeepLearningImplementations**

👁 Watch ▾    122          ★ Star    1,397          ⑂ Fork    495

&lt;&gt; **Code**        ⊘ Issues  **8**        ⎇ Pull requests  **3**        ▥ Projects  **0**        📖 Wiki        📊 Insights

Branch: master ▾        **DeepLearningImplementations** / **Eve** /

Create new file    Upload files    Find file    History

Eve

Thank you!

Adam and Eve in the Garden of Eden
Wenzel Peter
(Karlsbad 1745 - Rome 1829)