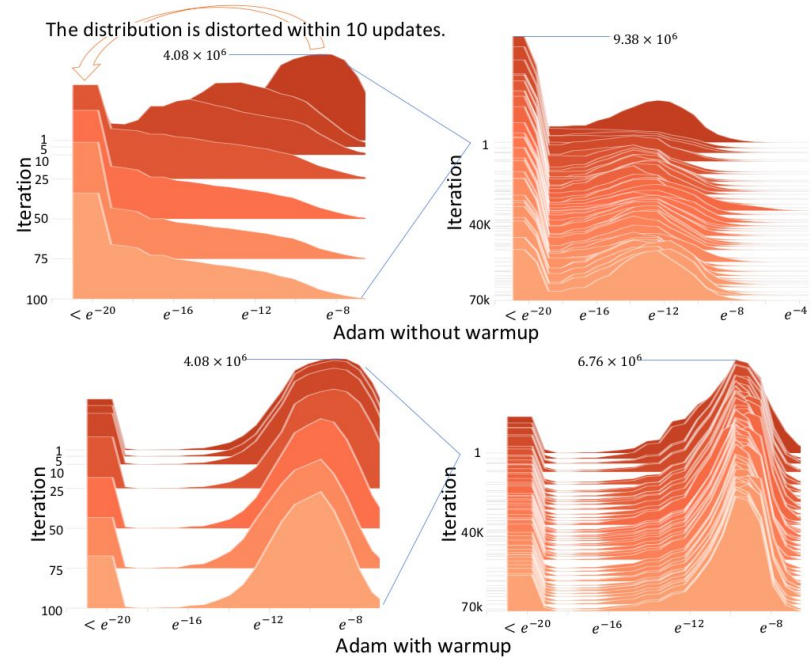# RAdam - Rectified Adam

- learning rate warmup stabilizes training, accelerates convergence and improves generalization

- problem of adaptive learning rates:
  - problematically large variance in the early stage of training

    → *suggests warmup works as a variance reduction technique!*

- RAdam introduces a term to rectify the variance of the adaptive learning rate!

**Absolute gradient histogram:**

Vienna
Deep Learning
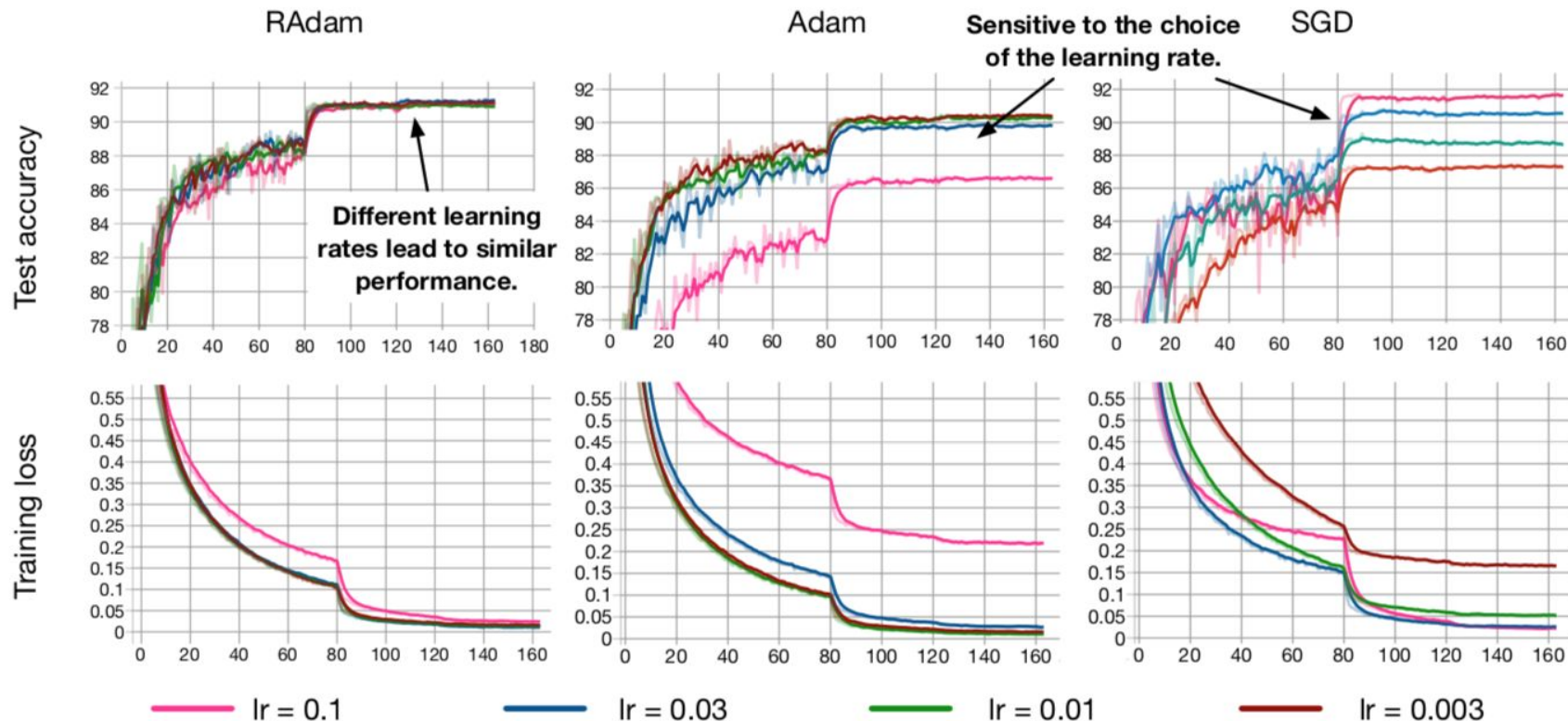Meetup

# RAdam - Update

**Algorithm 2:** Rectified Adam. All operations are element-wise.

---

**Input:** $\{\alpha_t\}_{t=1}^T$: step size, $\{\beta_1, \beta_2\}$: decay rate to calculate moving average and moving 2nd moment, $\theta_0$: initial parameter, $f_t(\theta)$: stochastic objective function.

**Output:** $\theta_t$: resulting parameters

1   $m_0, v_0 \leftarrow 0, 0$ (Initialize moving 1st and 2nd moment)

2   $\rho_\infty \leftarrow 2/(1 - \beta_2) - 1$ (Compute the maximum length of the approximated SMA)

3   **while** $t = \{1, \cdots, T\}$ **do**

4      $g_t \leftarrow \Delta_\theta f_t(\theta_{t-1})$ (Calculate gradients w.r.t. stochastic objective at timestep t)

5      $v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2)g_t^2$ (Update exponential moving 2nd moment)

6      $m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1)g_t$ (Update exponential moving 1st moment)

7      $\widehat{m_t} \leftarrow m_t/(1 - \beta_1^t)$ (Compute bias-corrected moving average)

8      $\rho_t \leftarrow \rho_\infty - 2t\beta_2^t/(1 - \beta_2^t)$ (Compute the length of the approximated SMA)

9      **if** *the variance is tractable, i.e., $\rho_t > 4$* **then**

10         $\widehat{v_t} \leftarrow \sqrt{v_t/(1 - \beta_2^t)}$ (Compute bias-corrected moving 2nd moment)

11         $r_t \leftarrow \sqrt{\frac{(\rho_t - 4)(\rho_t - 2)\rho_\infty}{(\rho_\infty - 4)(\rho_\infty - 2)\rho_t}}$ (Compute the variance rectification term)

12         $\theta_t \leftarrow \theta_{t-1} - \alpha_t r_t \widehat{m_t}/\widehat{v_t}$ (Update parameters with adaptive momentum)

13      **else**

14         $\theta_t \leftarrow \theta_{t-1} - \alpha_t \widehat{m_t}$ (Update parameters with un-adapted momentum)
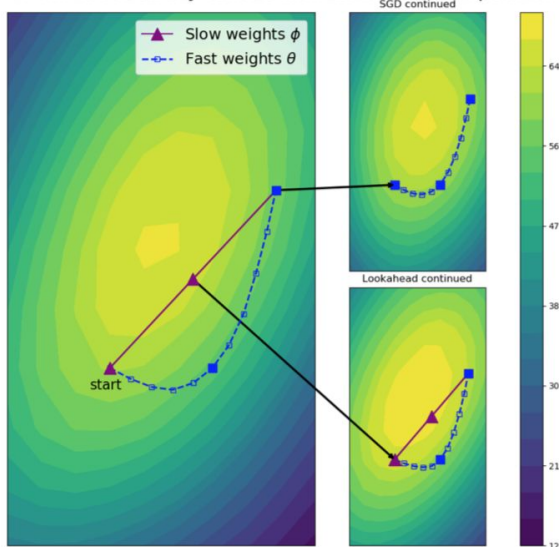
15   **return** $\theta_T$

---

# RAdam, Adam, and SGD

# LookAhead

- Iteratively updates two sets of weights:
  - "slow weights" get updated by looking ahead at the sequence of "fast weights" generated by another optimizer



CIFAR-100 accuracy surface with Lookahead interpolation

Slow weights $\phi$
Fast weights $\theta$

**Algorithm 1** Lookahead Optimizer:

**Require:** Initial parameters $\phi_0$, objective function $L$
**Require:** Synchronization period $k$, slow weights step size $\alpha$, optimizer $A$
**for** $t = 1, 2, \ldots$ **do**
    Synchronize parameters $\theta_{t,0} \leftarrow \phi_{t-1}$
    **for** $i = 1, 2, \ldots, k$ **do**
        sample minibatch of data $d \sim \mathcal{D}$
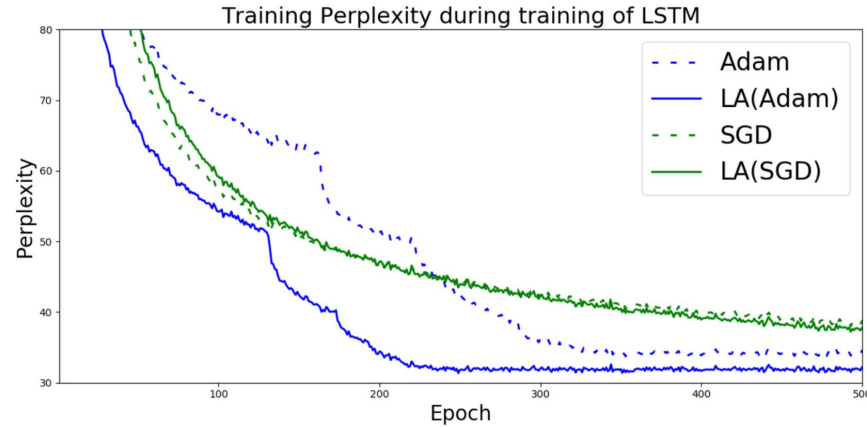        $\theta_{t,i} \leftarrow \theta_{t,i-1} + A(L, \theta_{t,i-1}, d)$
    **end for**
    Perform outer update $\phi_t \leftarrow \phi_{t-1} + \alpha(\theta_{t,k} - \phi_{t-1})$
**end for**
**return** parameters $\phi$

https://arxiv.org/pdf/1907.08610.pdf

Vienna
**Deep Learning**
Meetup

# LookAhead + ?



Training Perplexity during training of LSTM

A better combination?

## → Ranger = LookAhead + RAdam

https://arxiv.org/pdf/1907.08610.pdf &
https://medium.com/@lessw/new-deep-learning-optimizer-ranger-synergistic-combination-of-radam-lookahead-for-the-best-of-2dc83f79a48d

Vienna
**Deep Learning**
Meetup