

How to Scale Your EMA

Что такое EMA?

- EMA - Exponential Moving Average - взвешенное усреднение модели

$$\zeta_{t+1} = \rho \zeta_t + (1 - \rho) \theta_t$$

- У EMA более широкий район оптимума.
- EMA - простой способ получить модель, отличающуюся от исходной. Например, для дистилляции.
- Авторы статьи предлагают при увеличении батча в k раз, заменять ρ на ρ^k

EMA для SGD

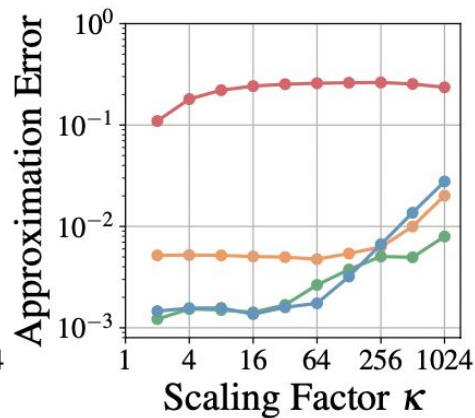
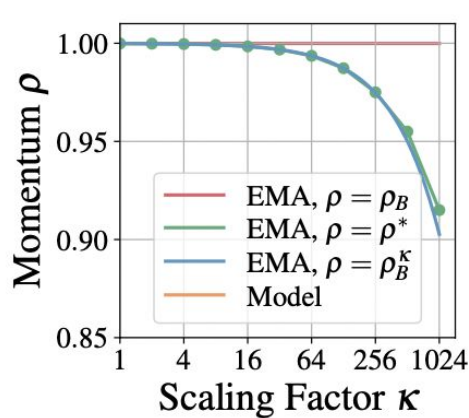
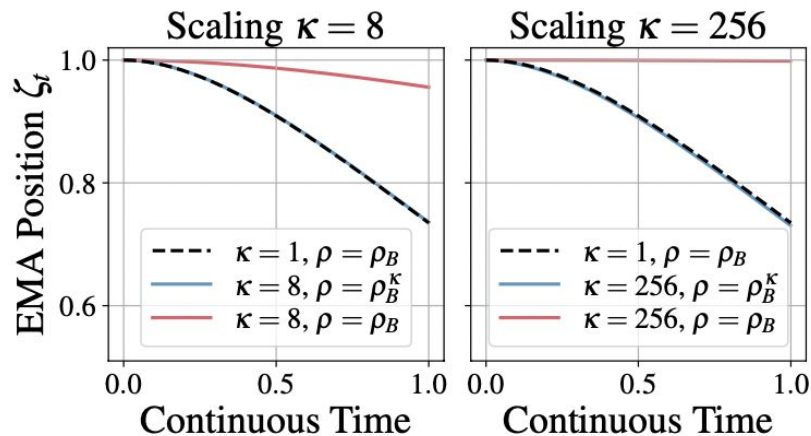
$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \times \frac{1}{B} \sum_{x \in \mathbb{B}} \nabla_{\boldsymbol{\theta}} \mathcal{L}(x; \boldsymbol{\theta}_t),$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(x; \boldsymbol{\theta}_{t+j}, \boldsymbol{\zeta}_{t+j}) \approx \nabla_{\boldsymbol{\theta}} \mathcal{L}(x; \boldsymbol{\theta}_t, \boldsymbol{\zeta}_t) \approx \mathbf{g},$$

$$\begin{bmatrix} \boldsymbol{\theta}_{t+\kappa} \\ \boldsymbol{\zeta}_{t+\kappa} \\ \mathbf{g} \end{bmatrix} = \begin{bmatrix} 1 & 0 & -\eta \\ (1 - \rho) & \rho & 0 \\ 0 & 0 & 1 \end{bmatrix}^{\kappa} \cdot \begin{bmatrix} \boldsymbol{\theta}_t \\ \boldsymbol{\zeta}_t \\ \mathbf{g} \end{bmatrix} = \begin{bmatrix} \boldsymbol{\theta}_t - \eta \kappa \mathbf{g} \\ \rho^{\kappa} \boldsymbol{\zeta}_t + (1 - \rho^{\kappa}) \boldsymbol{\theta}_t + \mathcal{O}(\eta \times \beta_{\rho}) \\ \mathbf{g} \end{bmatrix}.$$

Polyak Ruppert averaging

На inference используем модель с ЕМА весами



Semi-supervised learning

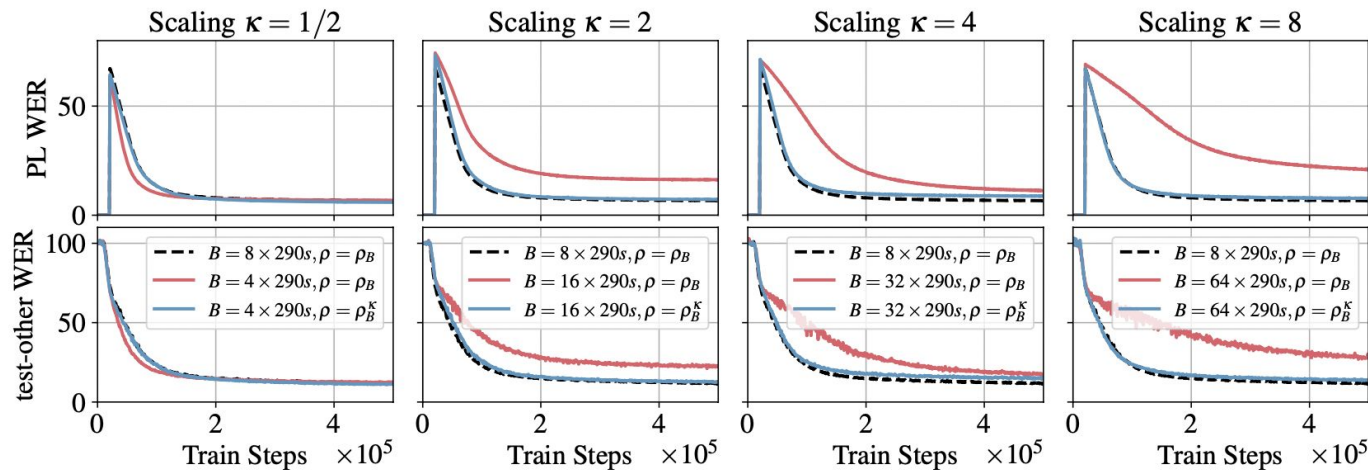


Figure 4: *Transformer pseudo-labeling on LibriSpeech* with different scalings κ . The baseline ($\kappa = 1$, black dashed) is trained with Adam at a *dynamic* batch size of 8×290 seconds, which corresponds to a single train step on the x -axis. The model EMA (*teacher*) is updated with momentum $\rho_B = 0.9999$. We investigate dynamic batch sizes down to $B = 4 \times 290$ s (left) and up to $B = 64 \times 290$ s (right), with (blue, $\rho = \rho_B^\kappa$) and without (red, $\rho = \rho_B$) the EMA Scaling Rule. The Adam Scaling Rule (Malladi et al. (2022), Definition C.3) is used throughout. For $\kappa \leq 2$, we start pseudo-labeling after $20\text{k}/\kappa$ training steps; while for $\kappa > 2$, we start when pre-training WER matches the baseline WER.

Self-supervised

