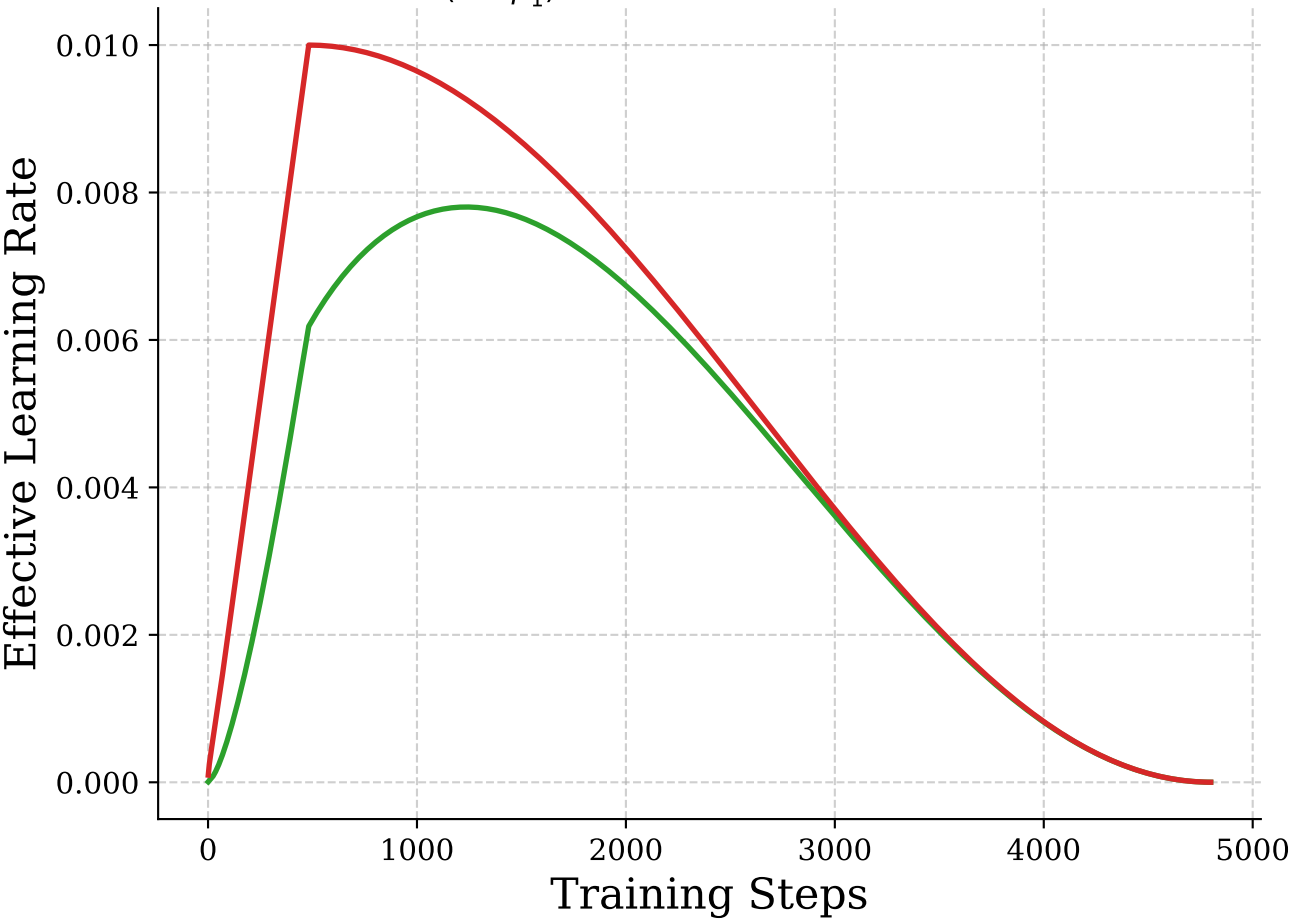


With Learning Rate Scheduler

$$\frac{\sqrt{1 - \beta_2^t}}{(1 - \beta_1^t)} \cdot \text{warmup\_cosine}(t)$$



With Constant Learning Rate

$$\frac{\sqrt{1 - \beta_2^t}}{(1 - \beta_1^t)} \cdot \text{lr}$$

