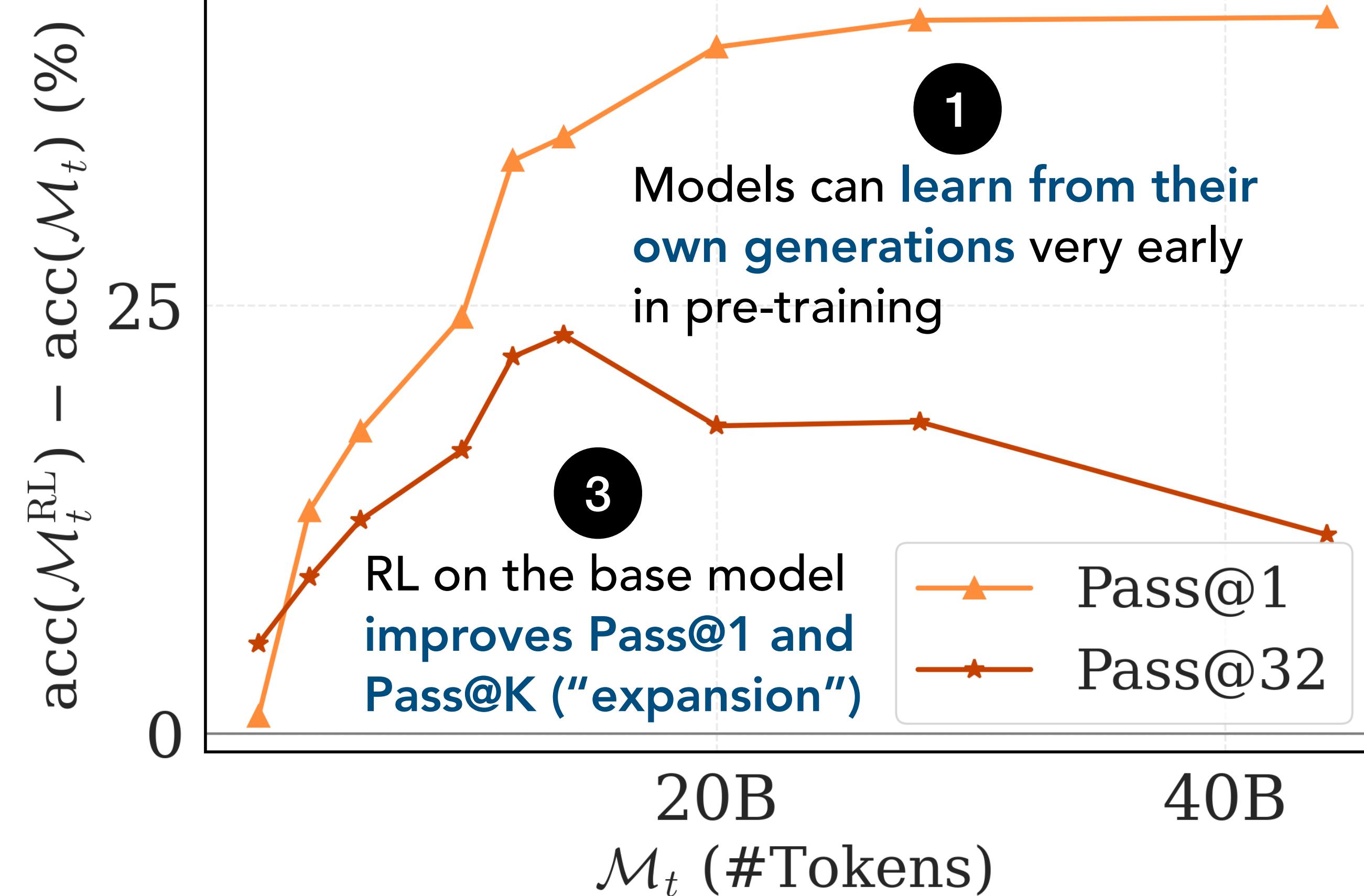


## RL on Pre-training Checkpoints: $\mathcal{M}_t^{\text{RL}}$



## RL in Standard Pipeline: $\mathcal{M}_t^{\text{SFT} \rightarrow \text{RL}}$

