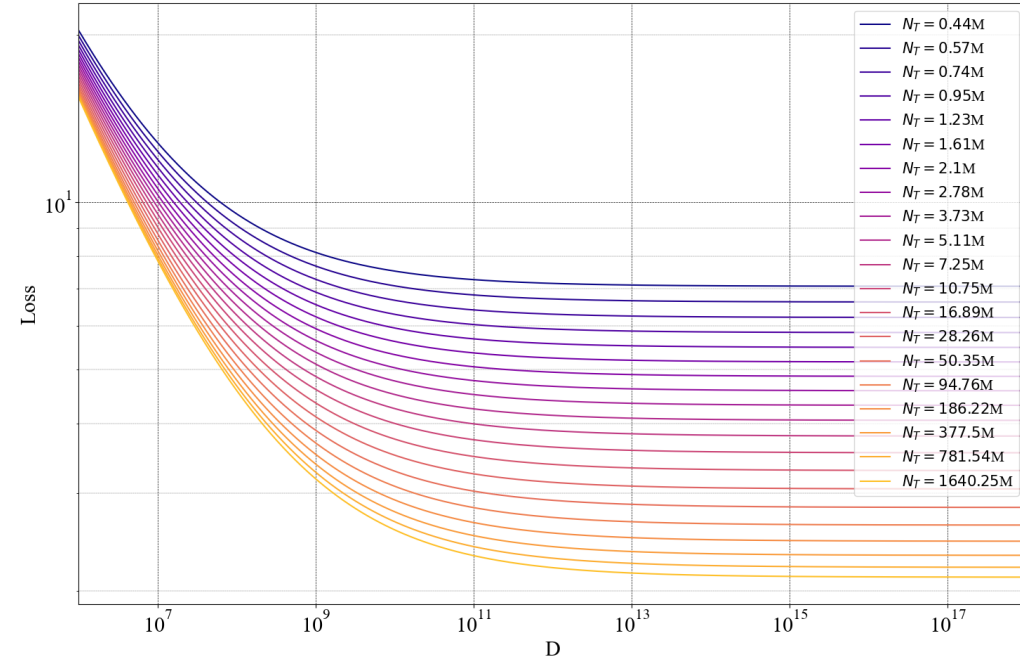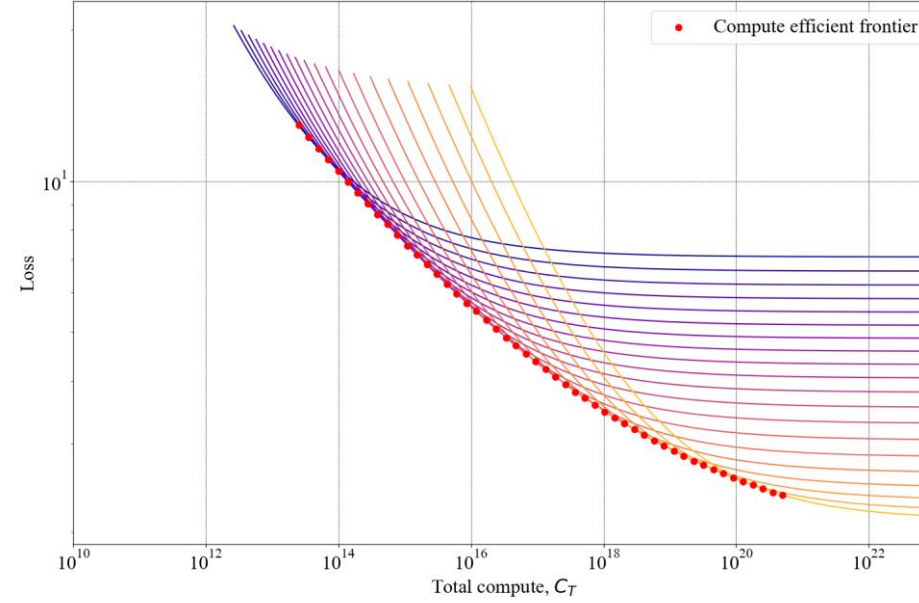Start from fitted model of **Chinchilla's** training curves

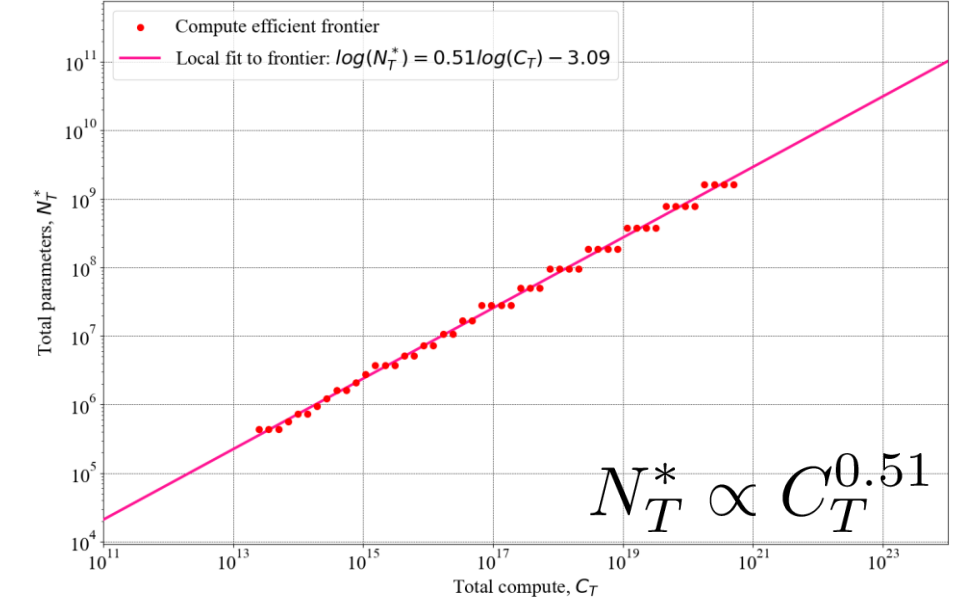$$\text{Loss}(N_T, D) = \frac{482}{N_T^{0.35}} + \frac{2085}{D^{0.37}} + 1.82$$

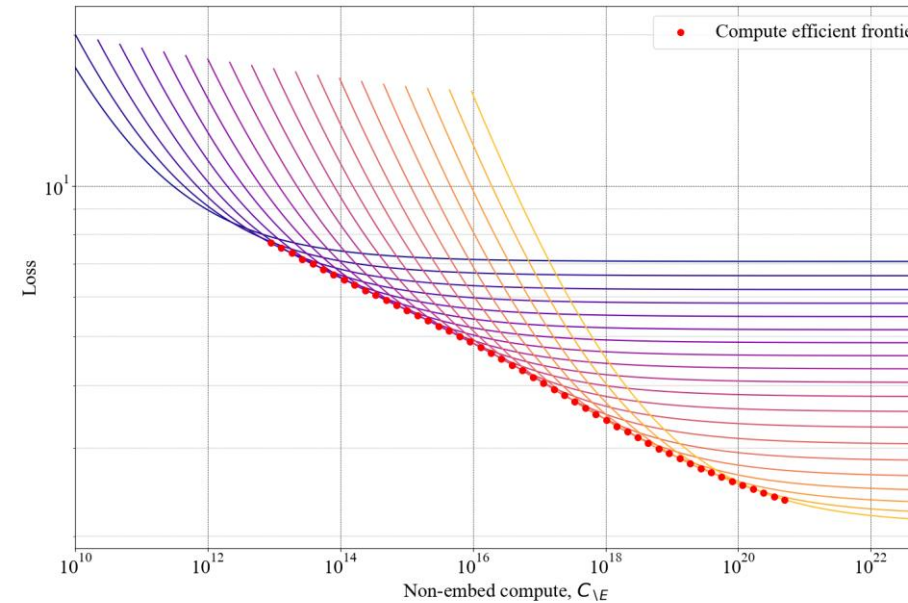Generate training curves for model sizes used in **Kaplan's** study (1k to 1.5B params)



Find compute optimal frontier in terms of **total** parameters $N_T$ and compute as in **Chinchilla**



Find power law scaling coefficient of 0.51, close to Chinchilla's 0.50



$$N_T^* \propto C_T^{0.51}$$

Find compute optimal frontier in **non-embedding** parameters $N_{\backslash E}$ and compute as in **Kaplan**



Find **local** power law scaling coefficient of 0.78, close to Kaplan's 0.73



$$N_{\backslash E}^* \propto C_{\backslash E}^{0.78}$$