

Chinchilla scaling law:

$L$  = LM test loss (avg. cross-entropy loss)

$D$  = dataset size (# training tokens)

$N$  = # model parameters (FF layers, Self-Attn, etc)

$C$  = compute budget, deterministic  
FLOPS ( $N, D$ ), fn. of data size  
and model size

Given a fixed FLOPS budget  $C$ ,

find

$$\operatorname{argmin}_{N, D} L(N, D)$$

$$N, D \text{ s.t. } \text{FLOPS}(N, D) = C$$

$$L(N, D) = \frac{A}{N^\alpha} + \frac{B}{D^\beta} + E$$

↳ contribution of model size

↳ contribution of data

Loss w/ perfect LM

after  
training a lot of diff models  
and fitting this eqn, we get

$$\alpha = 0.34, \beta = 0.28, E = 1.69$$

$A \sim B \sim 400$

Two models; same compute C

- Gopher (280B params, 300B tokens)
- Chinchilla (70B params, 1.4T tokens)

$$L(\text{Gopher}) = 1.993$$

$$L(\text{Chinchilla}) = 1.936$$

to put in perspective, according to those  
scaling laws, with this compute budget C,  
no model trained w/ 300B tokens could ever  
be better than Chinchilla!