

OVERVIEW

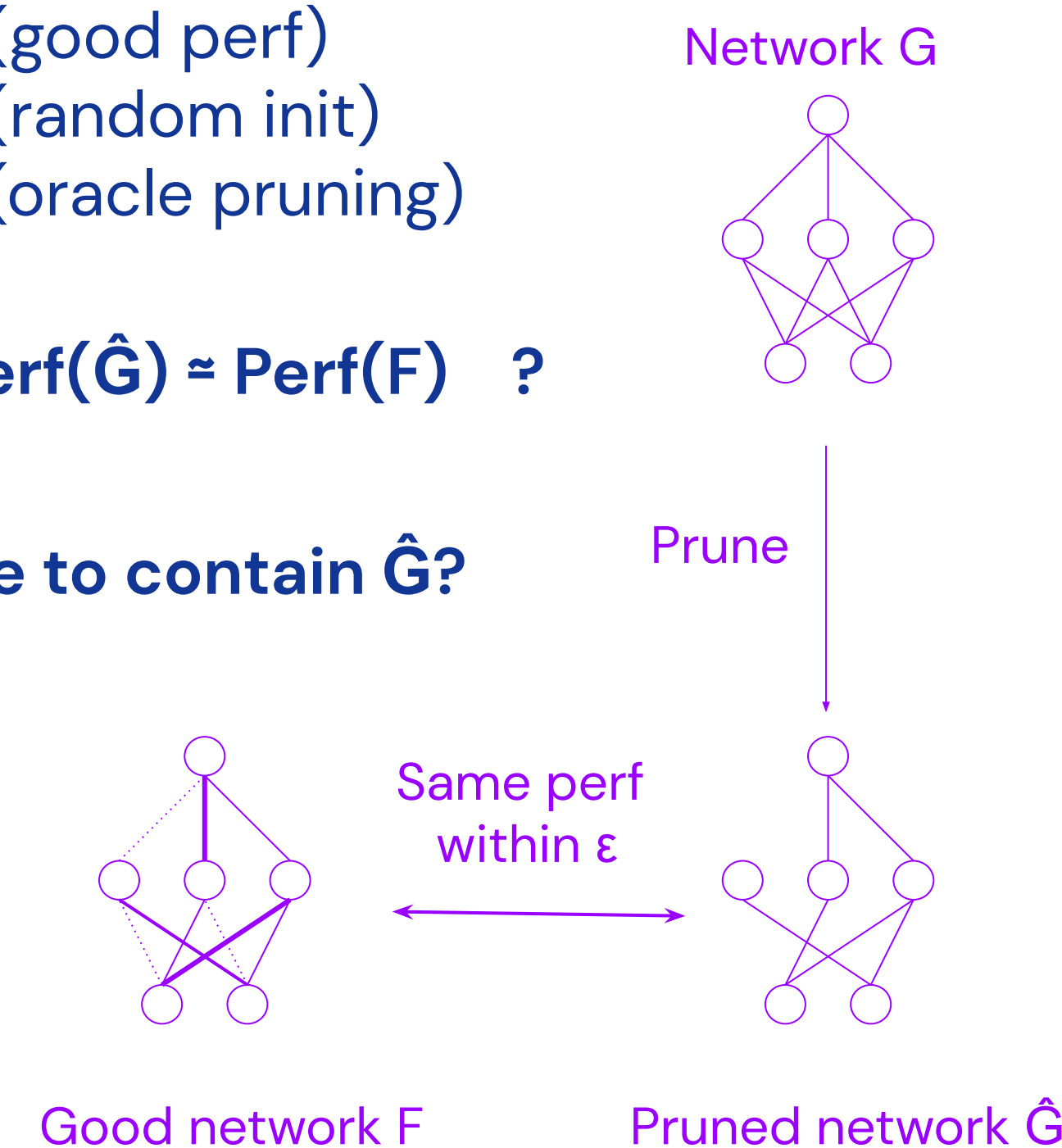
Problem description

- Conjecture by Ramanujan et al. 2019

- F**: Small ReLU NN (good perf)
- G**: Big ReLU NN (random init)
- \hat{G}** = Prune(G) (oracle pruning)

$\exists G, \forall F, \exists \hat{G}$ s.t. $\text{Perf}(\hat{G}) = \text{Perf}(F)$?

- How big must G be to contain \hat{G} ?



Previous result (Malach et al., ICML 2020)

- Idea: Add an intermediate ReLU layer**
 $\# \text{layers}(G) = 2\ell$
 $\# \text{neurons in } G \text{ per target weight of } F$:

$$O\left(\frac{n^3 \ell^2}{\varepsilon^2} \log \frac{n\ell}{\delta}\right)$$

- Strong assumptions**

- $\|W_i\|_2 \leq 1$ (weights at layer i)
- $\|W_i\|_\infty \leq 1/\sqrt{n}$
- $\|W_i\|_0 \leq n$
- $\|x\|_2 \leq 1$ (inputs)

Our result

- Assumptions
 - Same as Malach et al.
 - Hyperbolic distribution of the initial weights

#neurons in G per target weight of F:

$$\tilde{O}\left(\log \frac{n\ell}{\varepsilon}\right)$$

n : #neurons(F) per layer
 ℓ : #layers(F)
 ε : approximation error

Our result (more general)

- Assumptions
 - Hyperbolic distribution of the initial weights

#neurons in G per target weight of F:

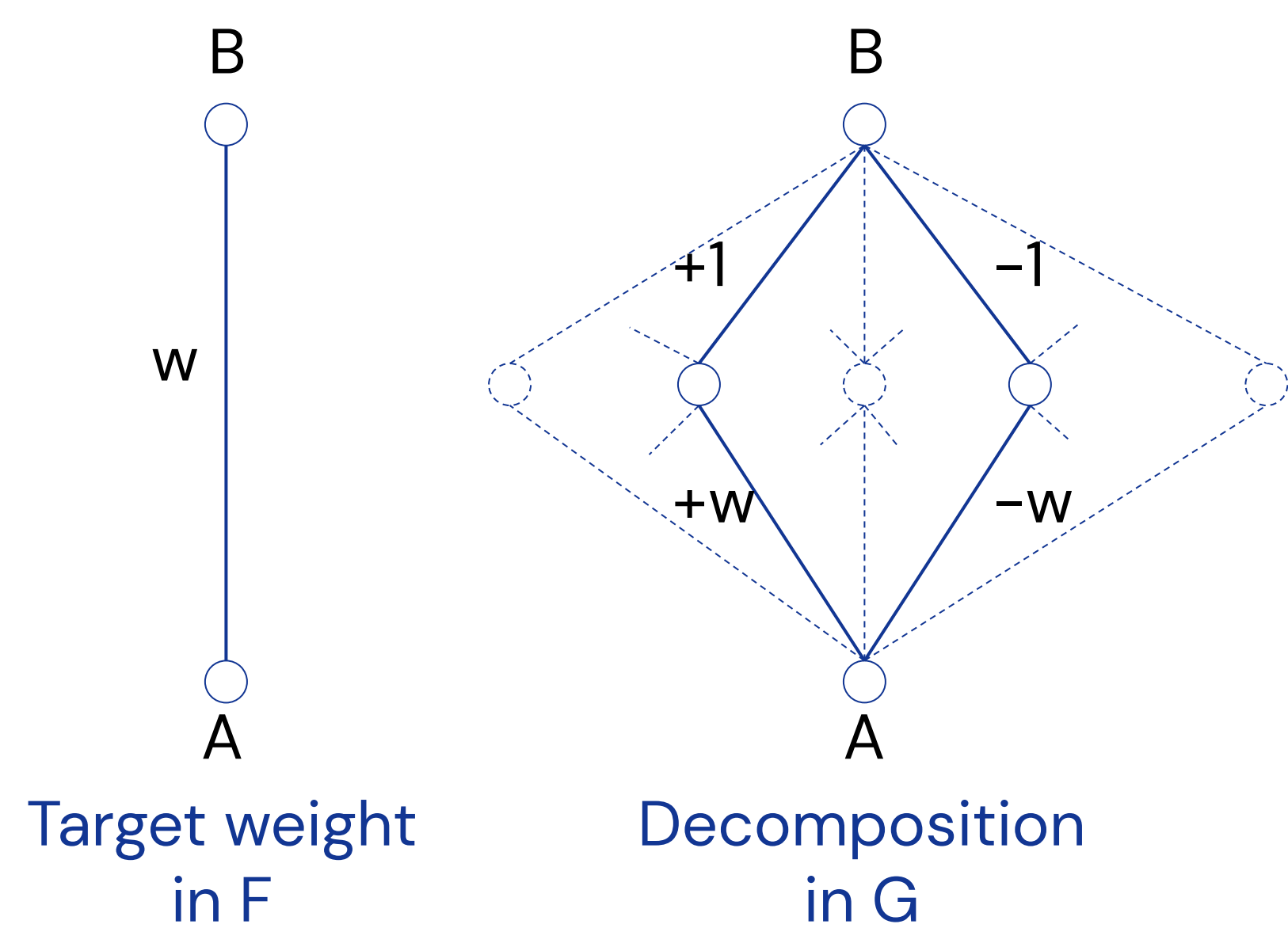
$$\tilde{O}\left(\log\left(\frac{n\ell}{\varepsilon} w_{\max} F_{\max}\right) + \sum_{i=1}^{\ell} \log \max\{1, \|W_i\|_2\}\right)$$

F_{\max} : max activation of any neuron
 w_{\max} : max weight
 W_i : matrix weight at layer i

TECHNICAL IDEAS

Weight Decomposition (Malach et al. 2020)

- Simulate one weight with 2 ReLU neurons



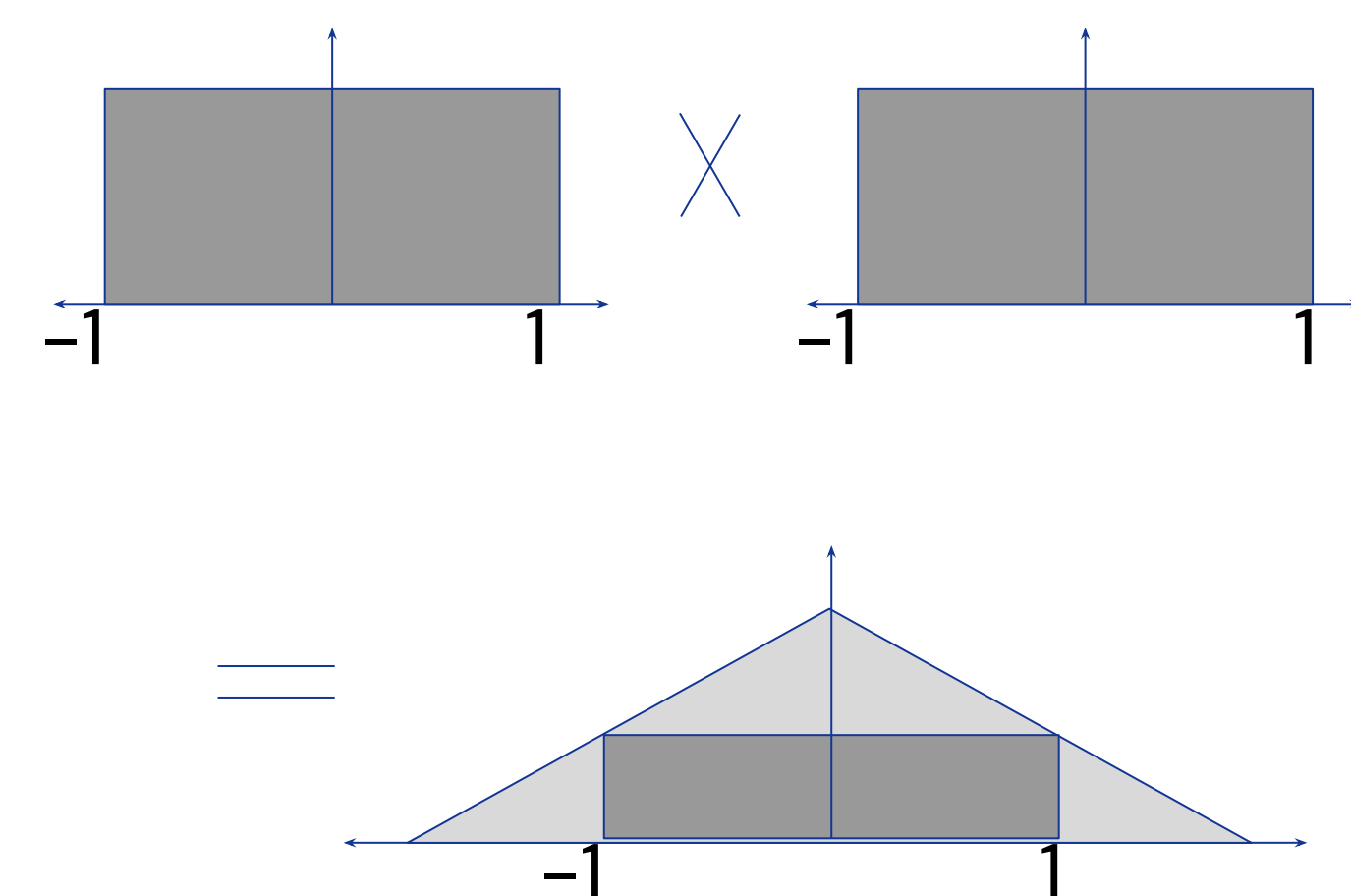
$$w = 1 \times \sigma(w) + (-1) \times \sigma(-w)$$

$$\sigma(x) = \max\{0, x\}$$

Sample intermediate neurons until
 Input and outputs are all within ε

Takes $O(1/\varepsilon^2 \log 1/\delta)$ samples

Product weights



$$w \approx wa \sigma(wb) + wc \times \sigma(wd)$$

Sample intermediate neurons until

$$wa \times wb = +w \text{ \& \; } \text{sgn}(wa) = +1$$

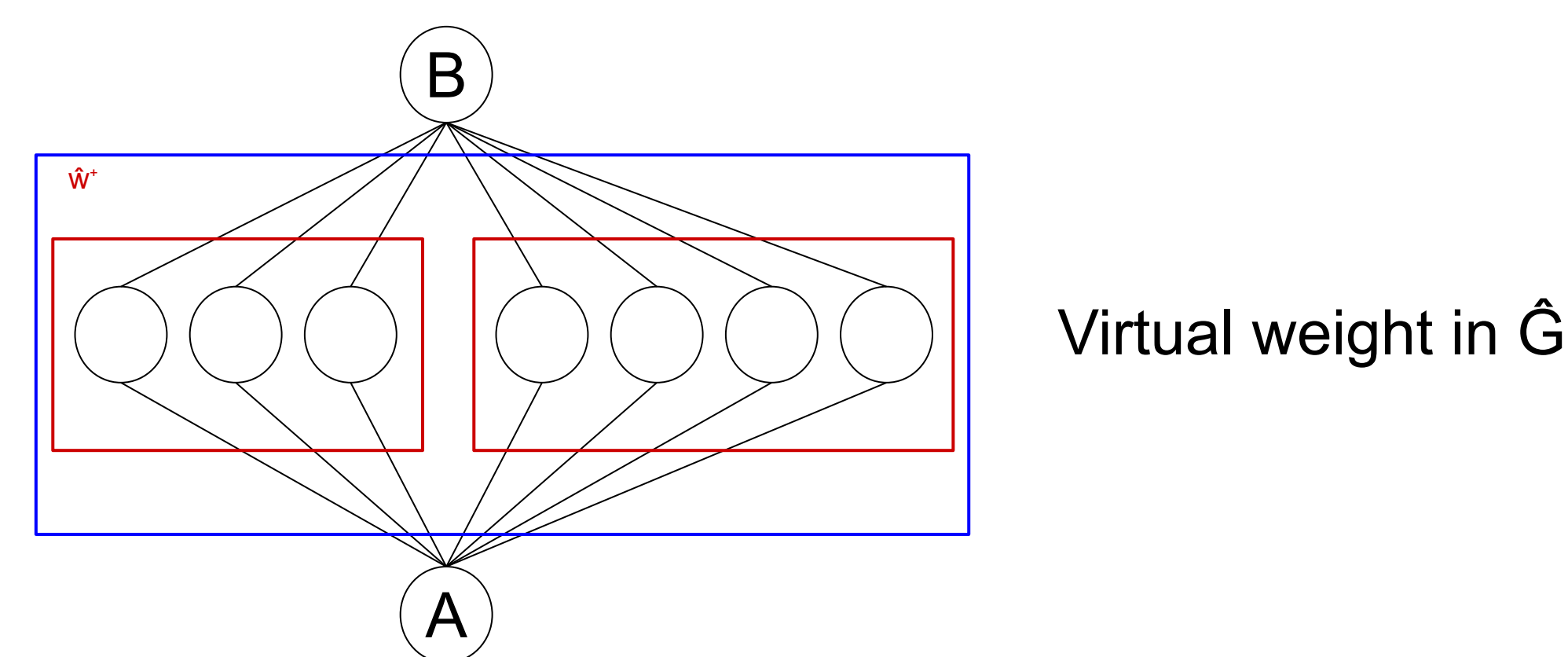
$$wc \times wd = -w \text{ \& \; } \text{sgn}(wc) = -1$$

Requires $O(1/\varepsilon \log 1/\delta)$ samples

(similar for hyperbolic distribution
 by a change of variable argument)

Golden Ratio Decomposition

- Take advantage of the sum in the neuron function
- Binary decomposition: requires $\log 1/\varepsilon$ intermediate neurons
- Weights are sampled from hyperbolic $P(w) = 1/w$
- Base 2 not possible, use base 3/2 instead (or $\varphi = (1+\sqrt{5})/2$)



Binary decomposition: $\sum_{i=1}^k b_i 2^{-i}$

Golden-ratio decomposition: $\sum_{i=1}^k b_i x_i^{-i}$, $x_i \in [\varphi^{-i-1}, \varphi^{-i}]$

$$\varphi = \frac{1+\sqrt{5}}{2} \text{ or } \varphi = \frac{3}{2}$$

→ Hyperbolic sampling: $P([\varphi^{-i-1}, \varphi^{-i}]) \geq c \quad \forall i$; $P_w(w) \propto 1/w$

Need only $\tilde{O}(\log 1/\varepsilon \log 1/\delta)$ samples

Batch sampling

- Don't throw away samples that can be reused elsewhere
- Fill k disjoint categories each with n samples w.p. $1-\delta$
 - $P(\text{any cat.}) \geq c$
- Needs #samples M : (Chernoff-Hoeffding)

$$M = \left\lceil \frac{2}{c} \left(m + \ln \frac{k}{\delta} \right) \right\rceil$$

m : #weights(F) per layer
 k : #neurons to decompose a weight = $O(\log 1/\varepsilon)$
 c : probability of one of the k segments

CONCLUSION

- Is hyperbolic sampling worth trying in practice?
- What about uniform sampling? A lower bound?
 - Pensia et al., Neurips 2020
 "Optimal lottery tickets via subset-sum: Logarithmic over-parameterization is sufficient."