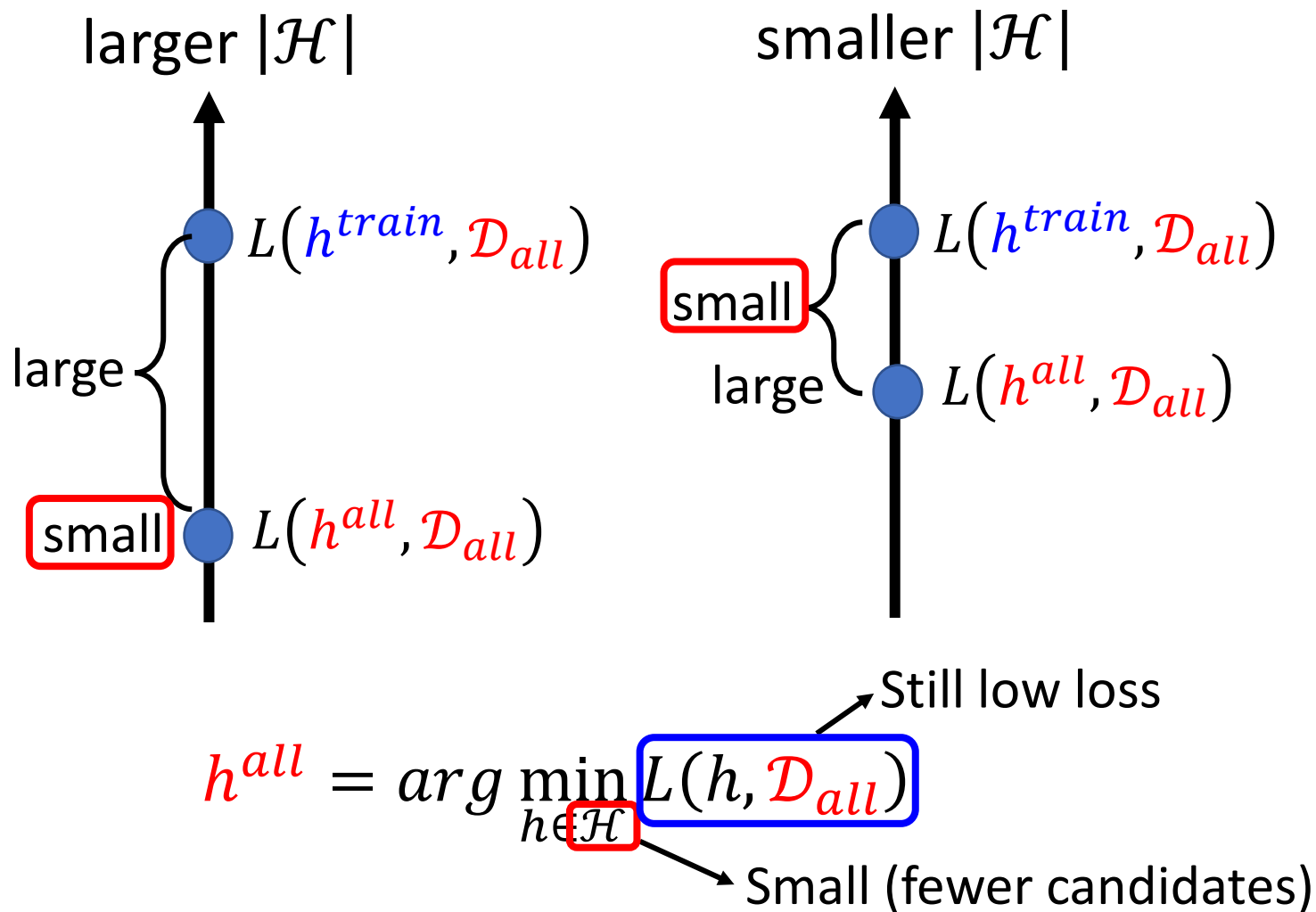


Why Deep Learning?

李宏毅

Hung-yi Lee

魚與熊掌可以兼得嗎？

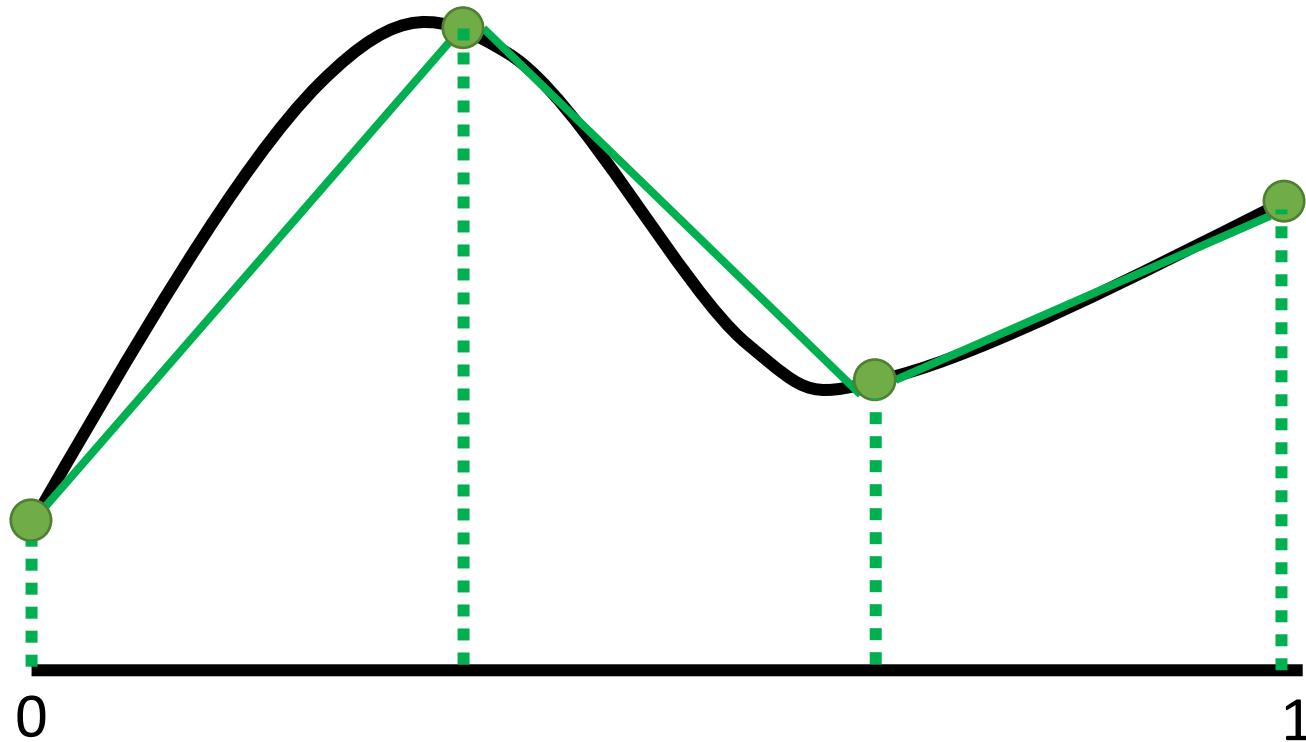




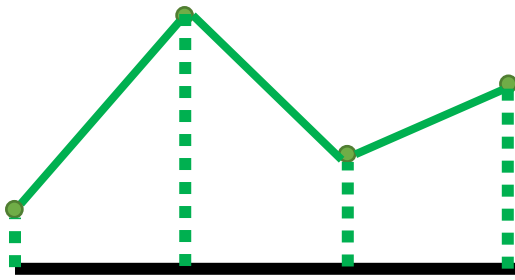
Review: Why Hidden Layer?

Piecewise Linear

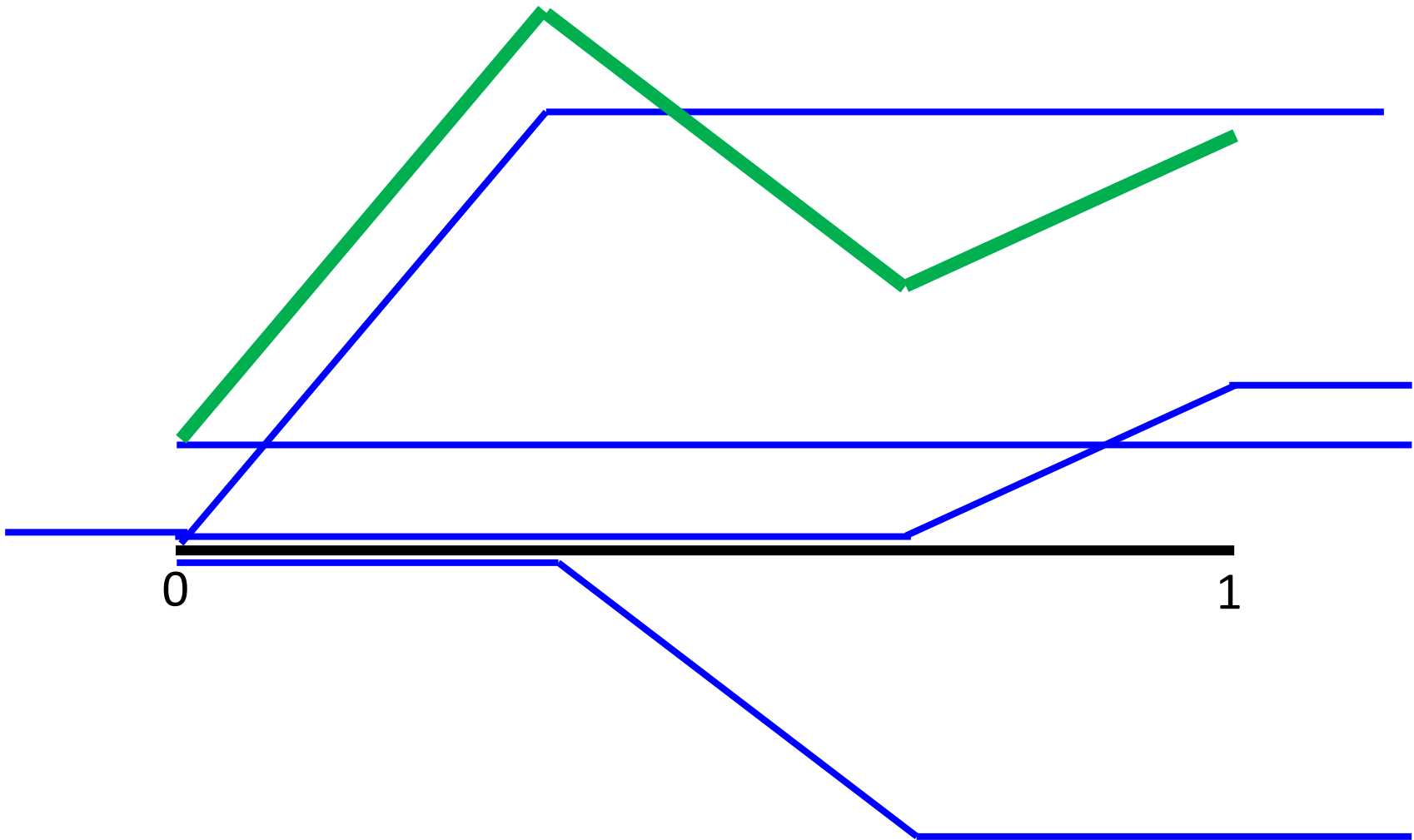
We can have good approximation with sufficient pieces.



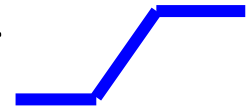
piecewise
linear



= constant +
sum of a set of

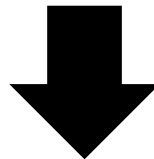


Piecewise linear = constant + sum of a set of



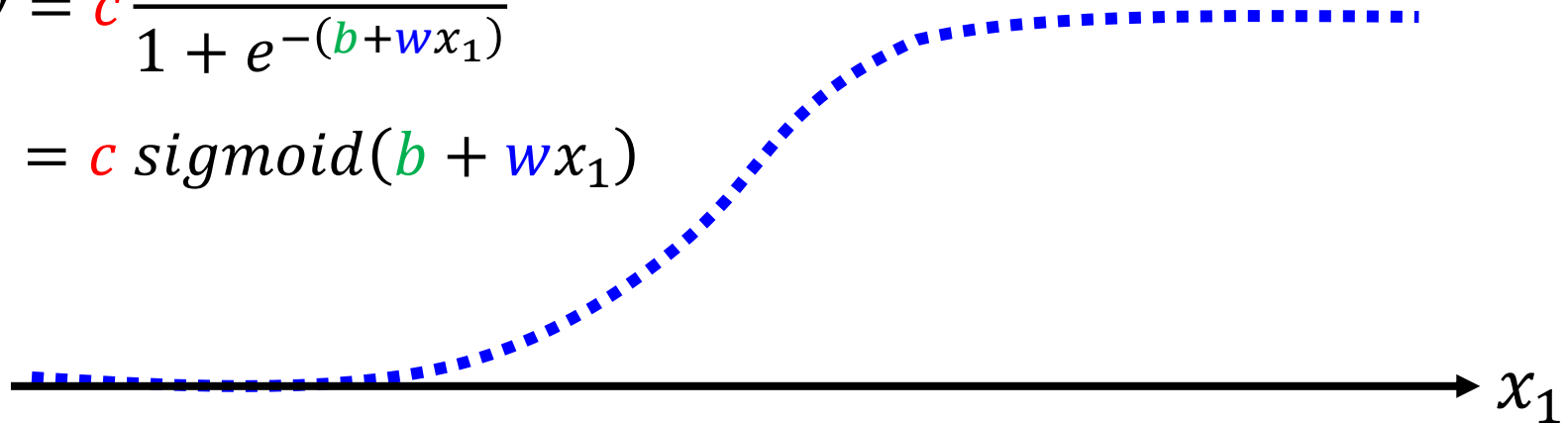
How to represent
this function?

Hard Sigmoid

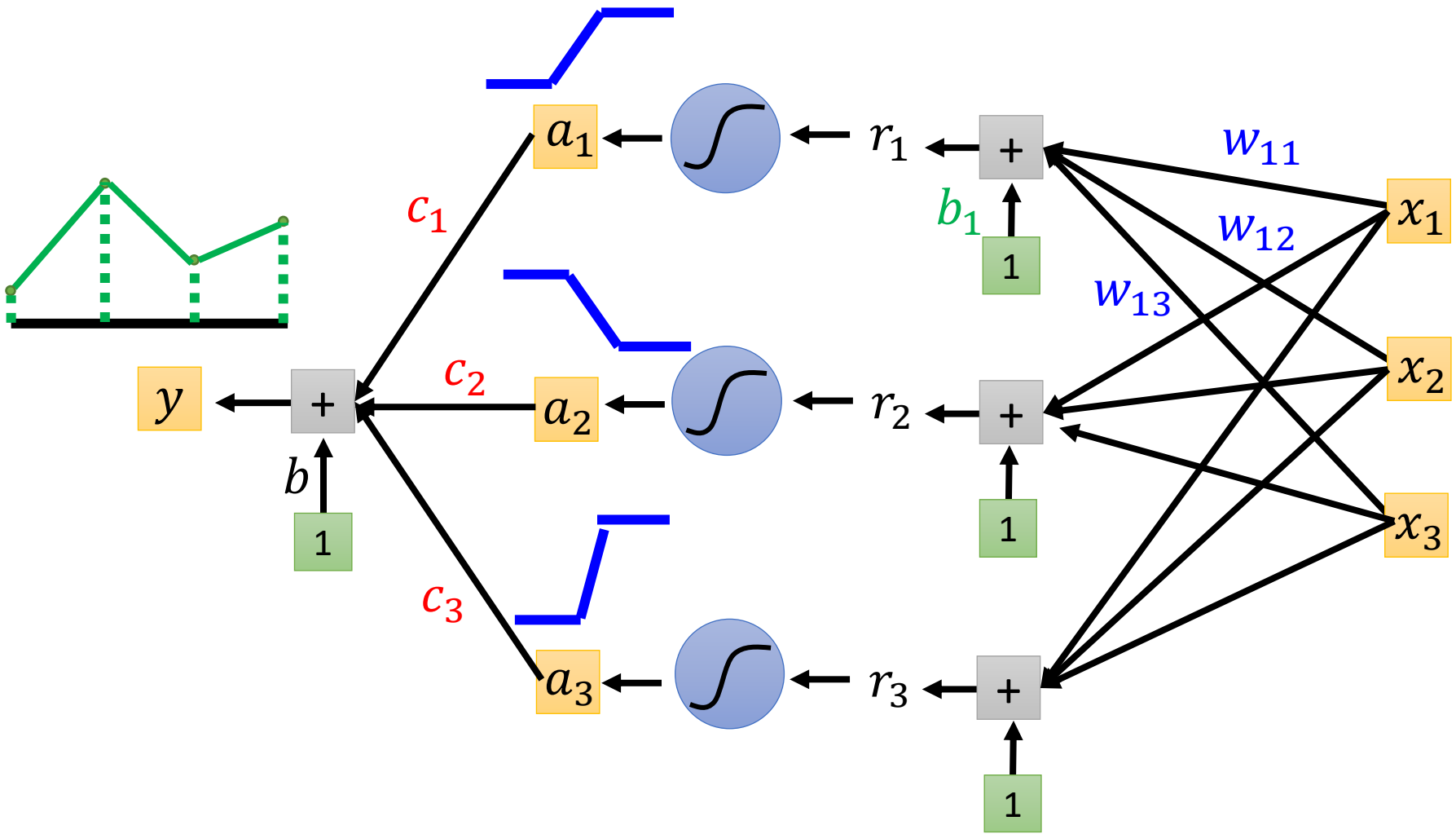
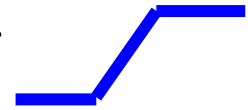


Sigmoid Function

$$y = c \frac{1}{1 + e^{-(b+wx_1)}}$$
$$= c \text{ sigmoid}(b + wx_1)$$

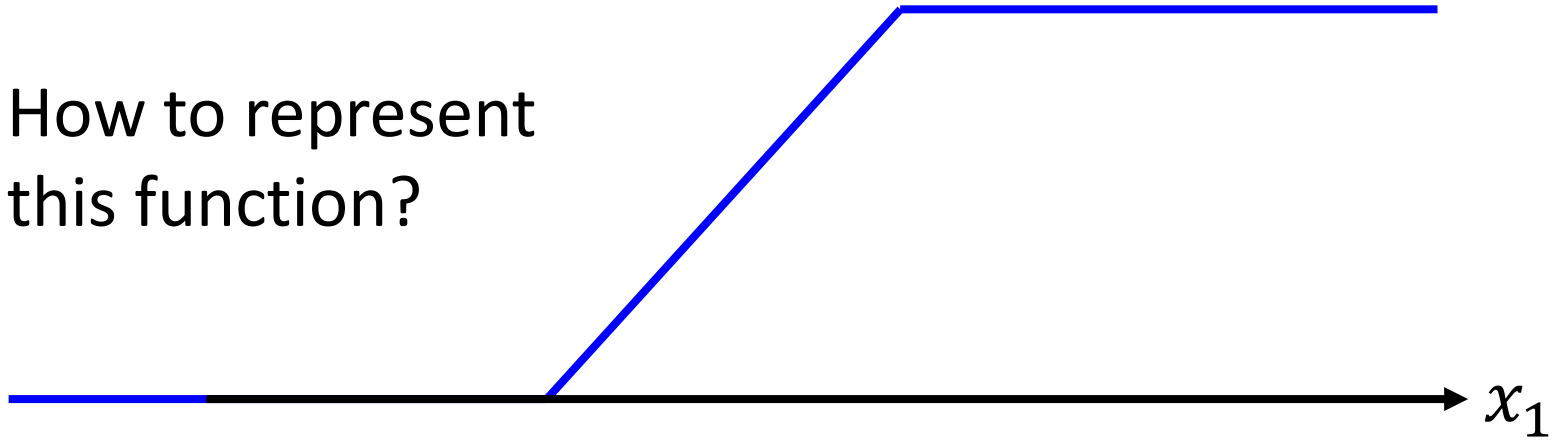


Piecewise linear = constant + sum of a set of

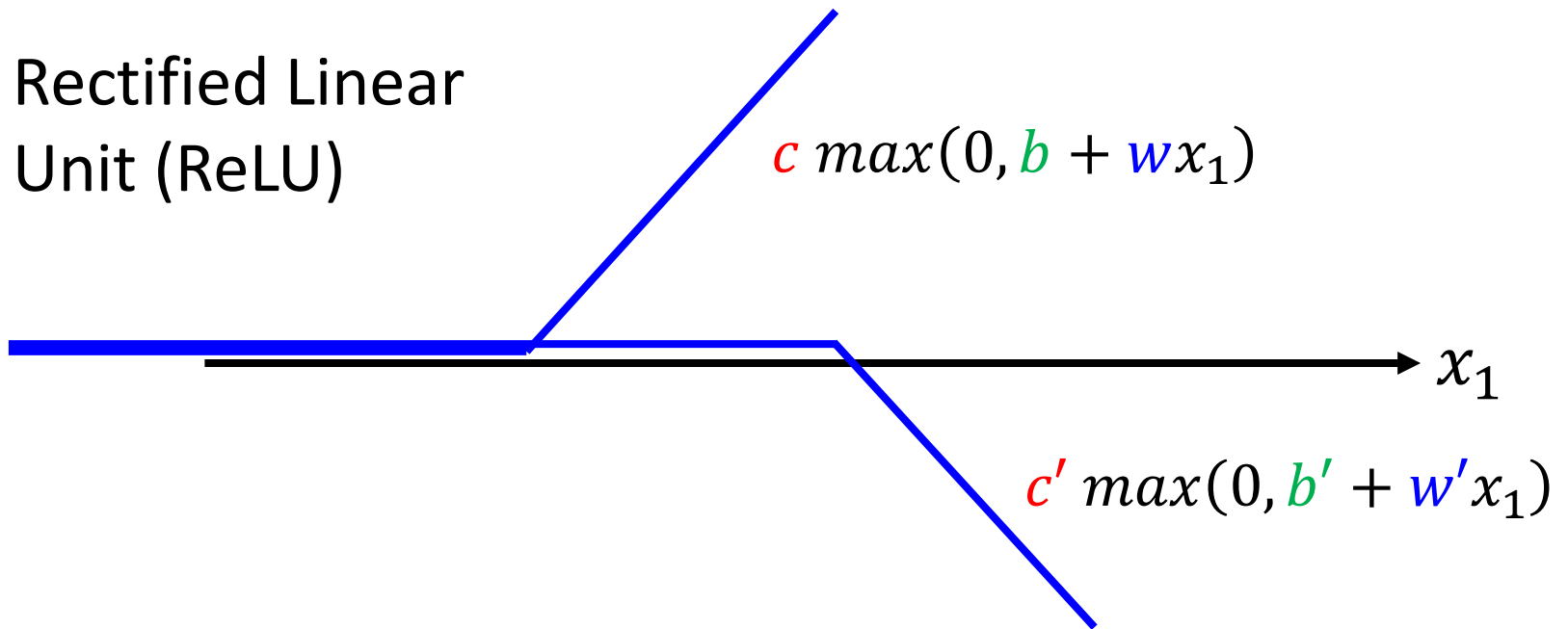


Hard Sigmoid \rightarrow ReLU

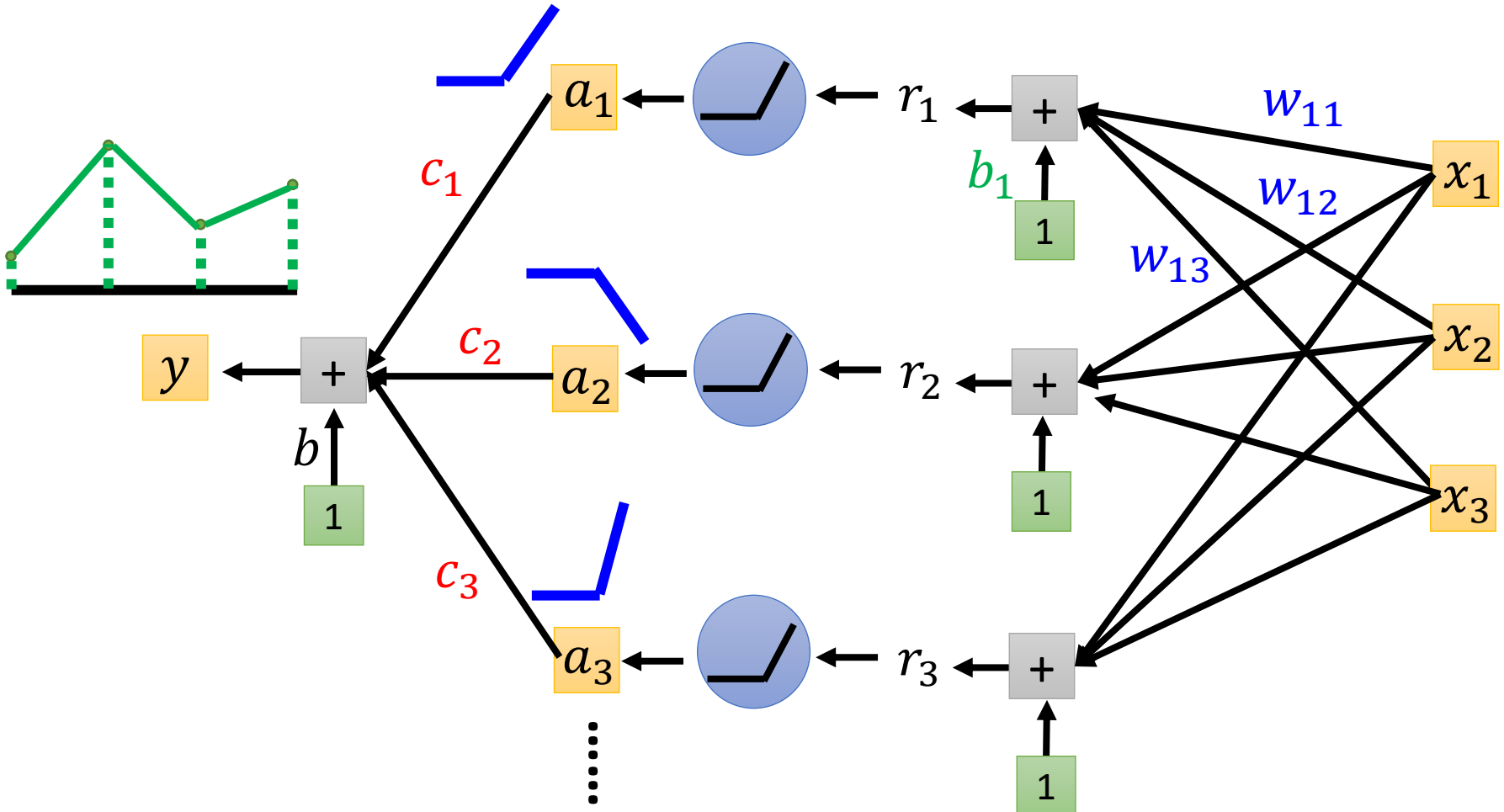
How to represent
this function?



Rectified Linear
Unit (ReLU)



Piecewise linear = constant + sum of a set of



Why we want “**Deep**” network, not “**Fat**” network?

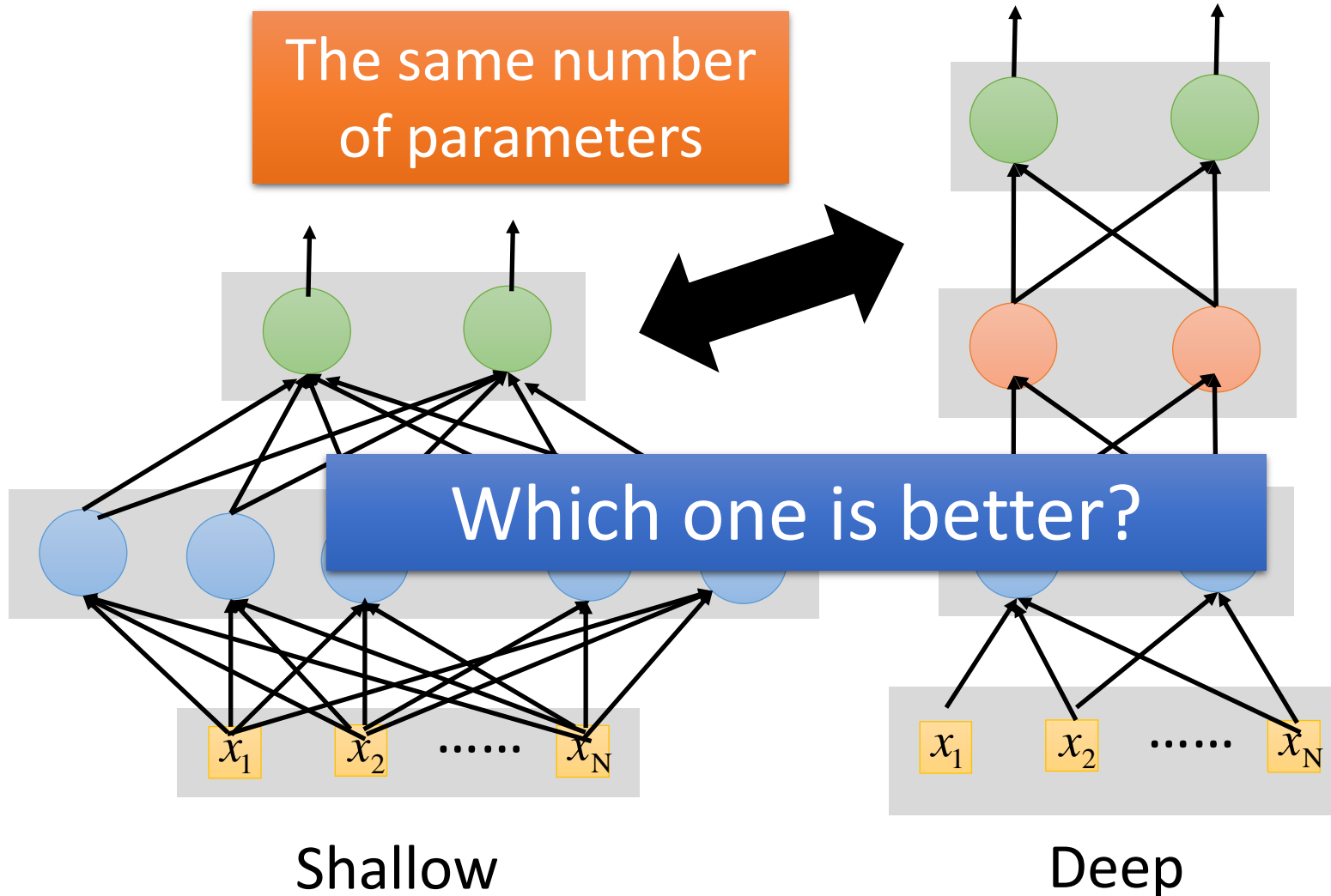
Deeper is Better?

Layer X Size	Word Error Rate (%)
1 X 2k	24.2
2 X 2k	20.4
3 X 2k	18.4
4 X 2k	17.8
5 X 2k	17.2
7 X 2k	17.1

Not surprised, more parameters, better performance

Seide Frank, Gang Li, and Dong Yu. "Conversational Speech Transcription Using Context-Dependent Deep Neural Networks." *Interspeech*. 2011.

Fat + Short v.s. Thin + Tall



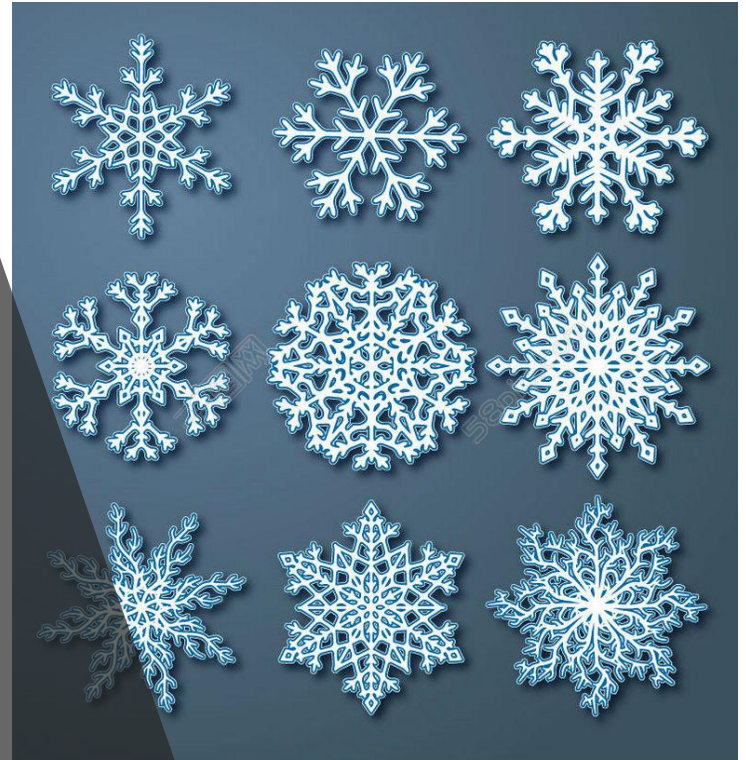
Fat + Short v.s. Thin + Tall

Layer X Size	Word Error Rate (%)	Layer X Size	Word Error Rate (%)
1 X 2k	24.2		
2 X 2k	20.4		
3 X 2k	18.4		
4 X 2k	17.8		
5 X 2k	17.2	1 X 3772	22.5
7 X 2k	17.1	1 X 4634	22.6
		1 X 16k	22.1

Why?

Seide Frank, Gang Li, and Dong Yu. "Conversational Speech Transcription Using Context-Dependent Deep Neural Networks." *Interspeech*. 2011.

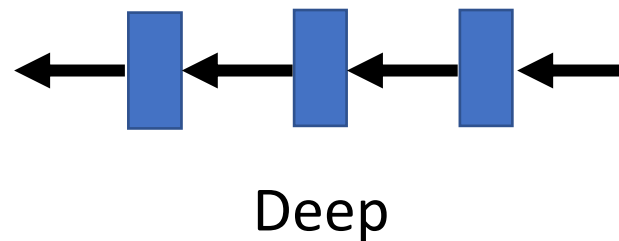
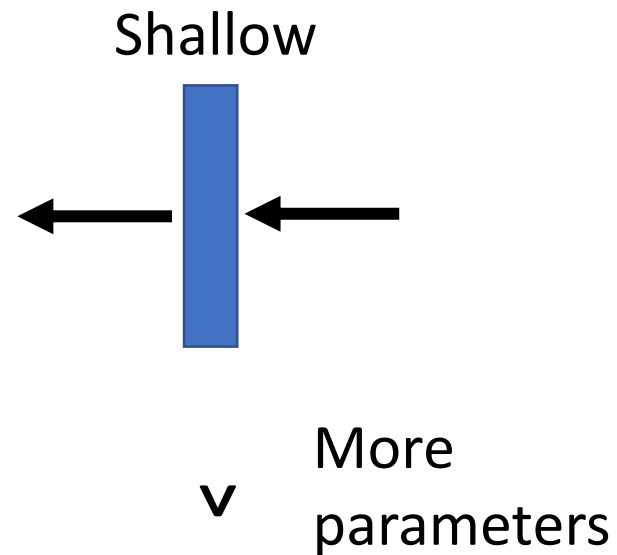
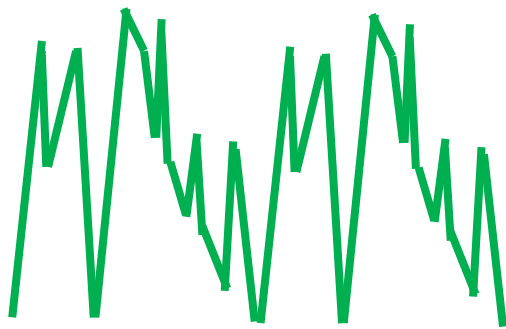
Why we need deep?



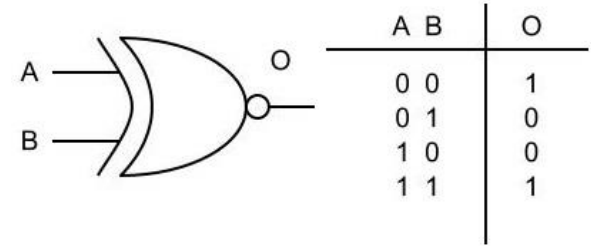
Yes, one hidden layer can represent any function.

However, using deep structure is more effective.

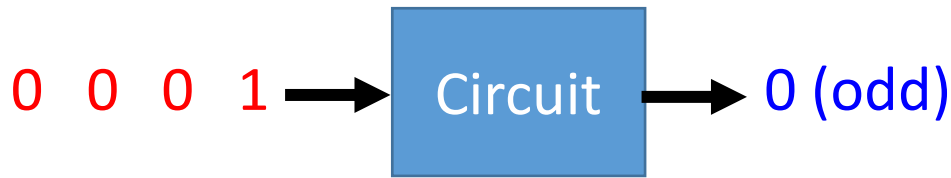
Why we need deep?



Analogy – Logic Circuits

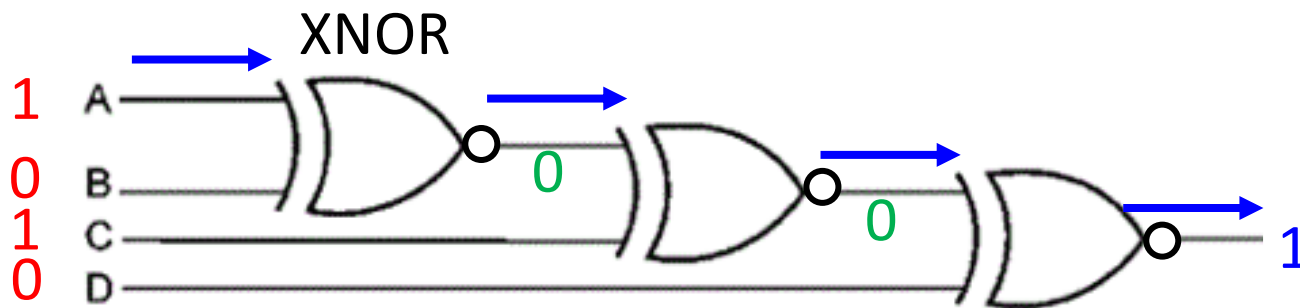


- E.g., parity check



For input sequence with d bits,

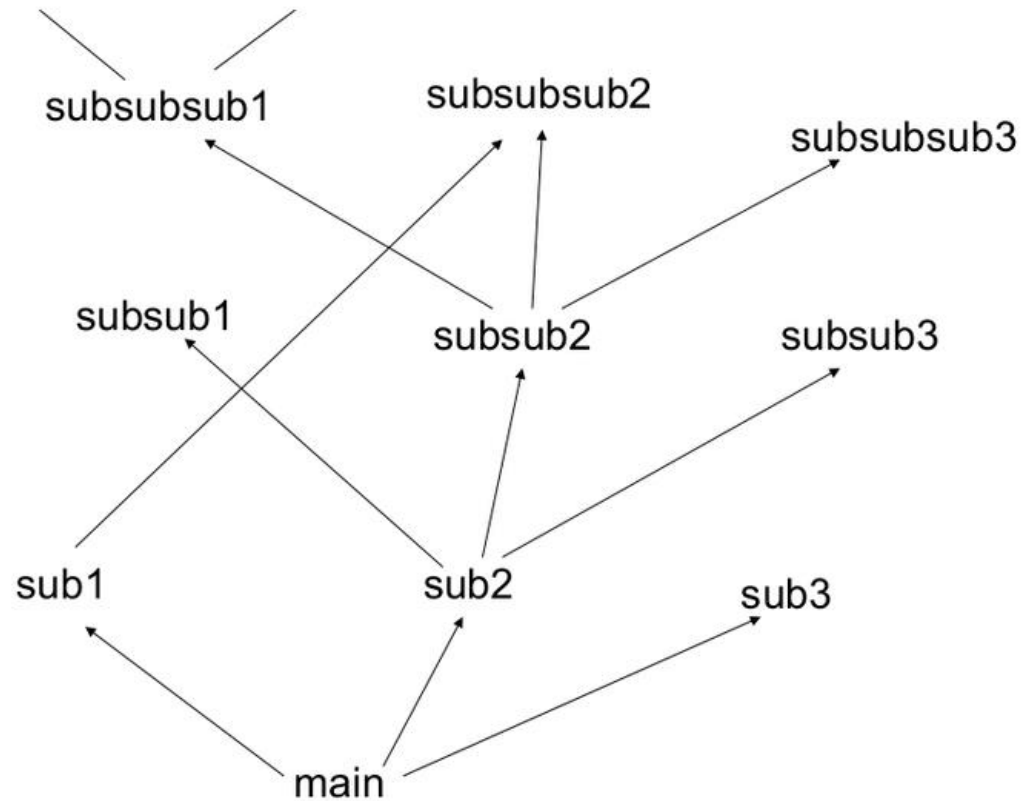
Two-layer circuit need $O(2^d)$ gates.



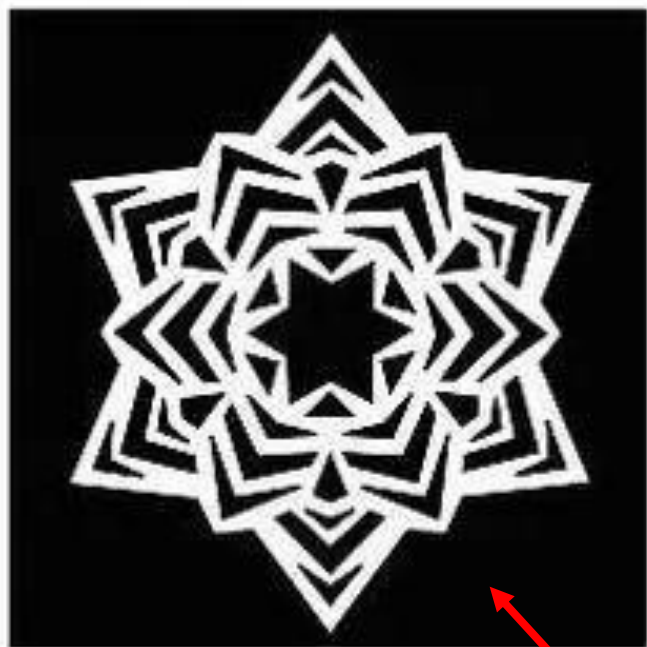
With multiple layers, we need only $O(d)$ gates.

Analogy – Programming

Don't put everything in your main function.

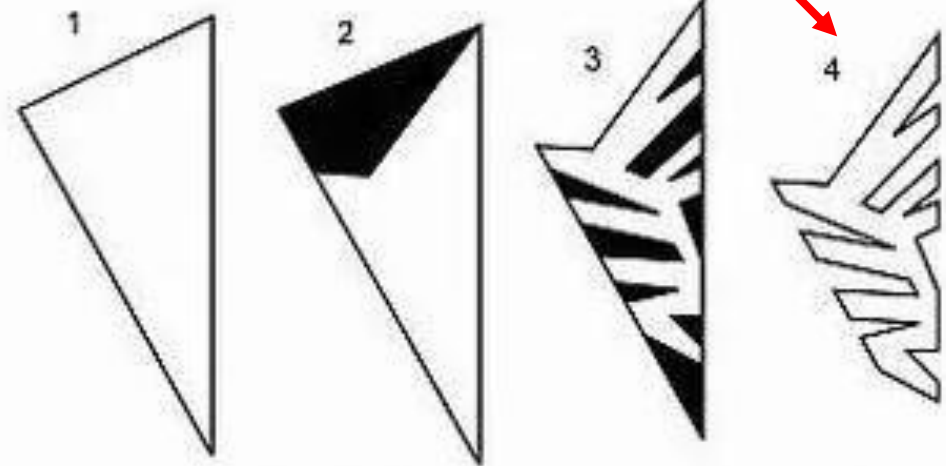


More Analogy

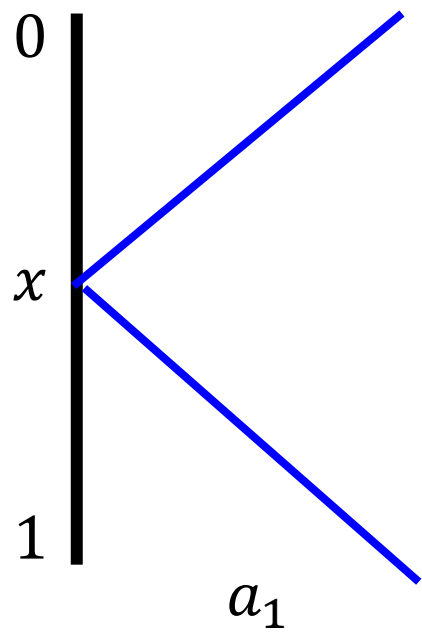
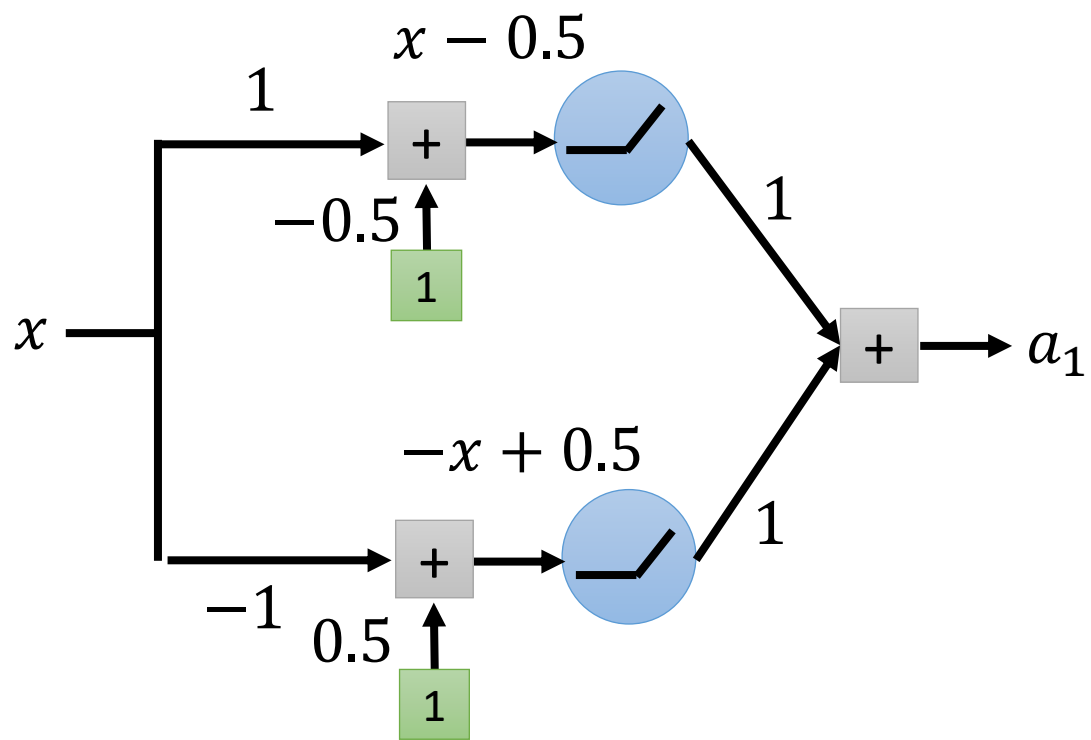


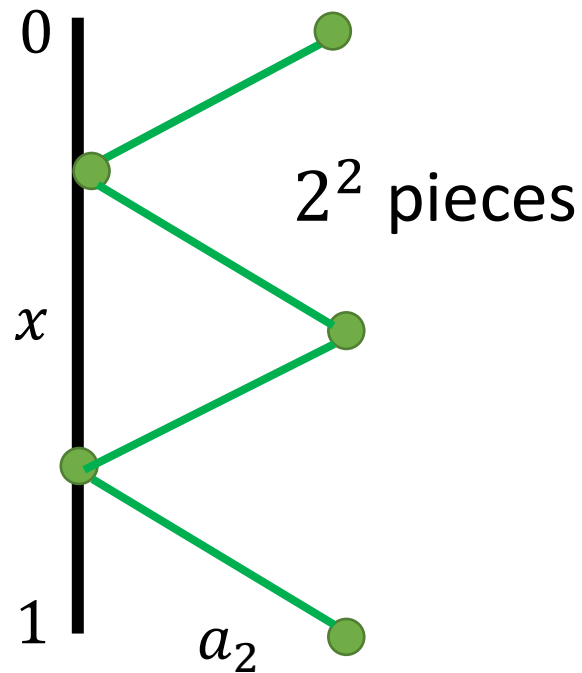
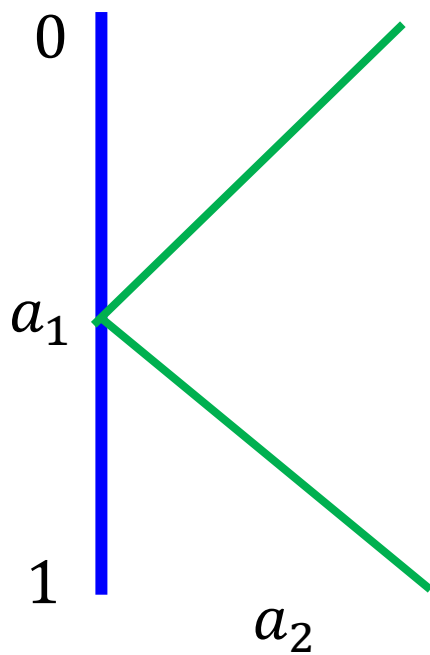
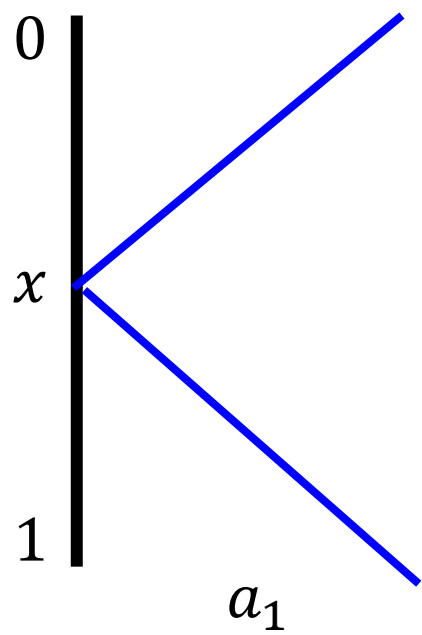
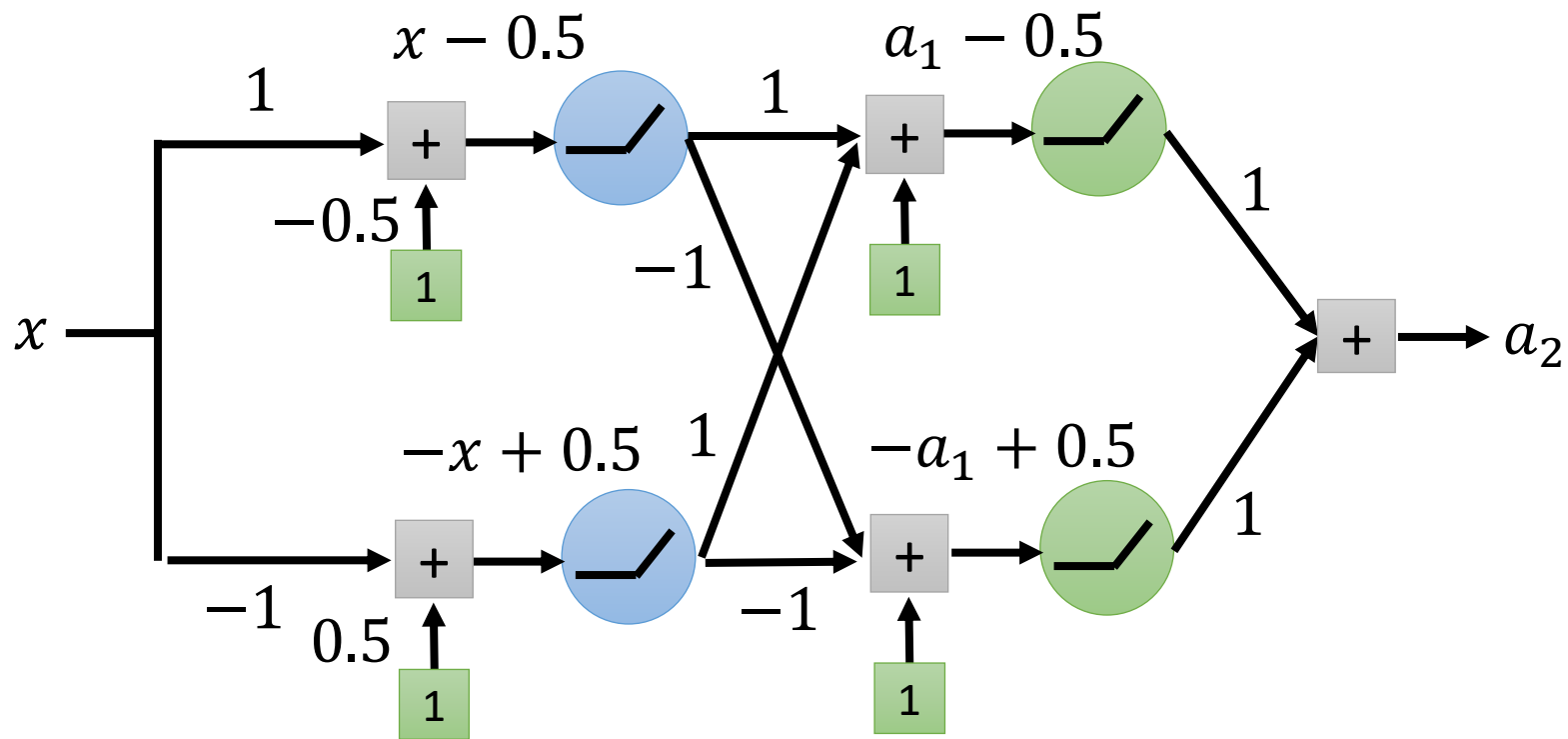
剪很多刀

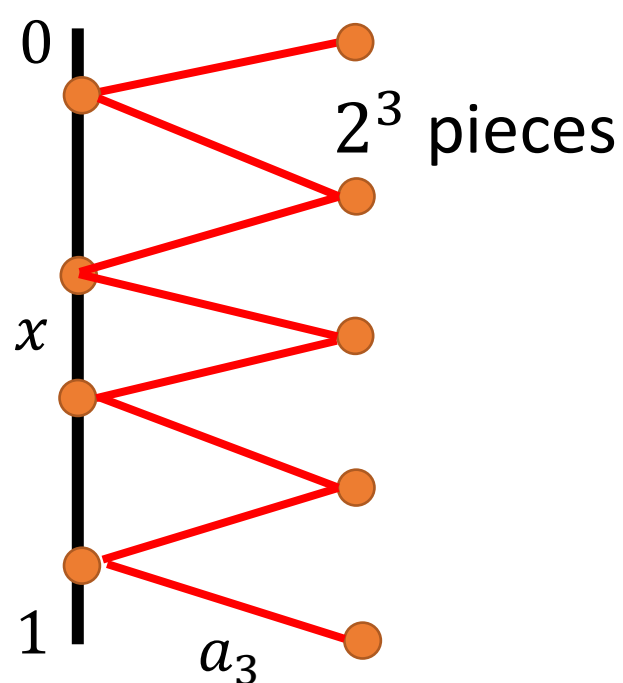
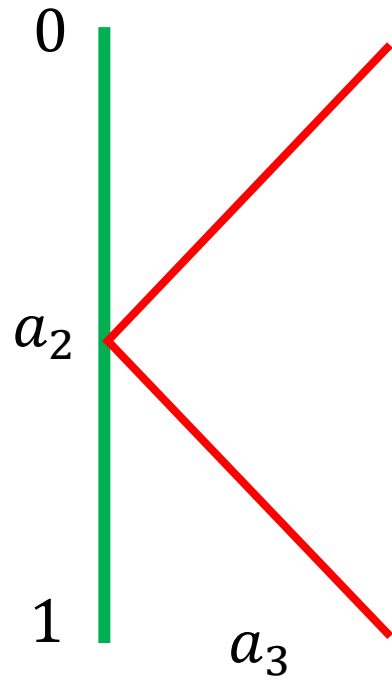
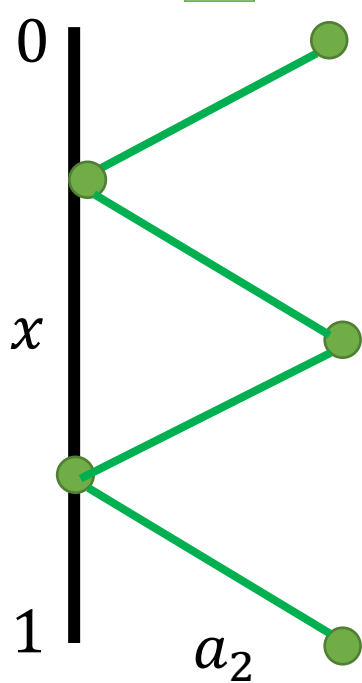
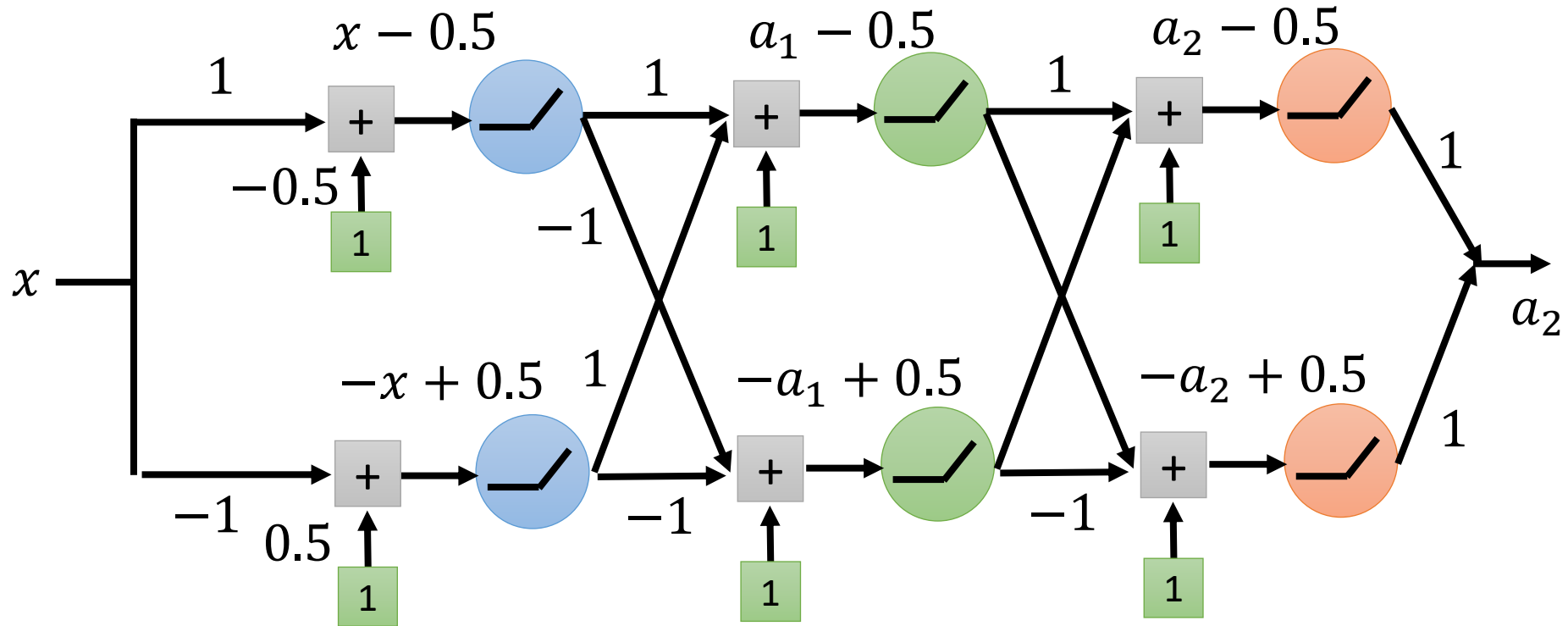
比較有效率



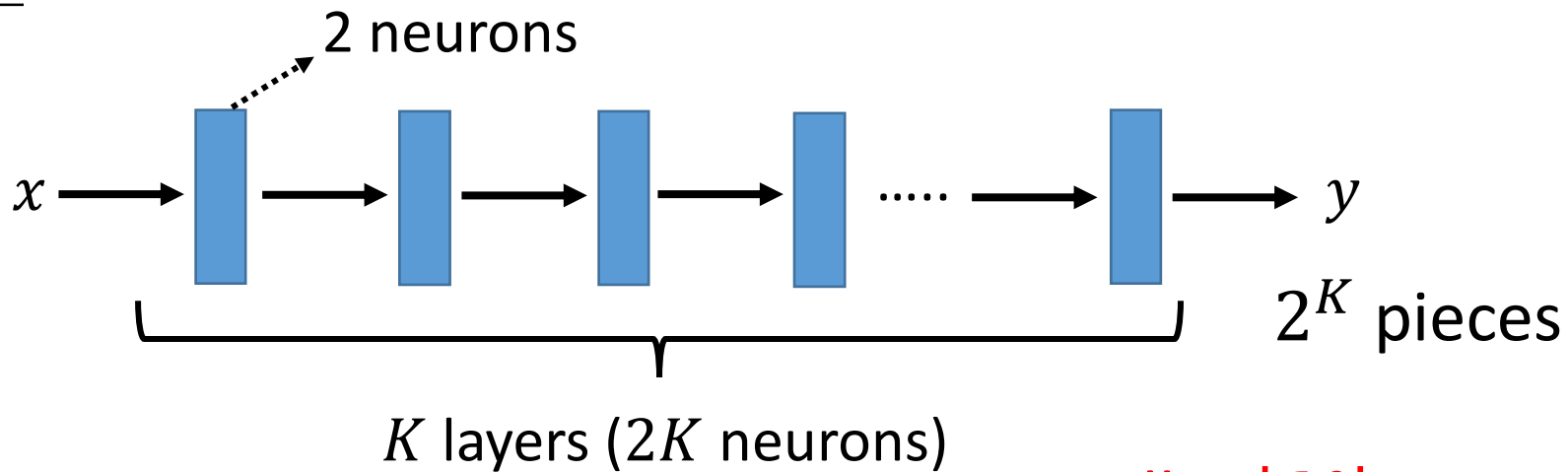
头条号 / 幼师宝典





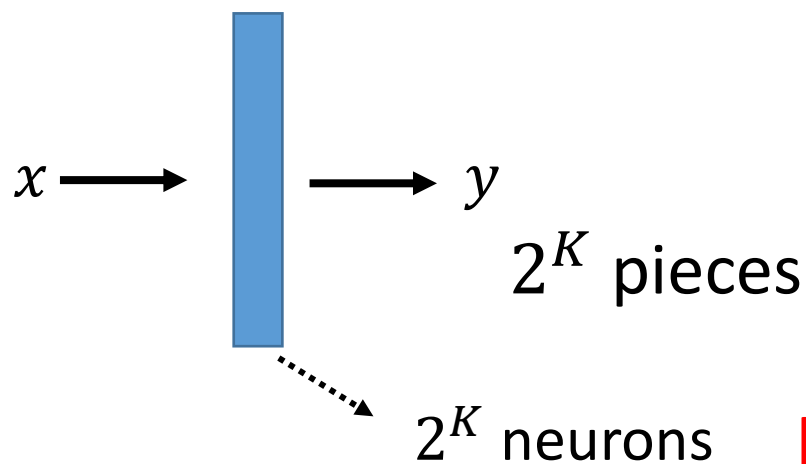


Deep



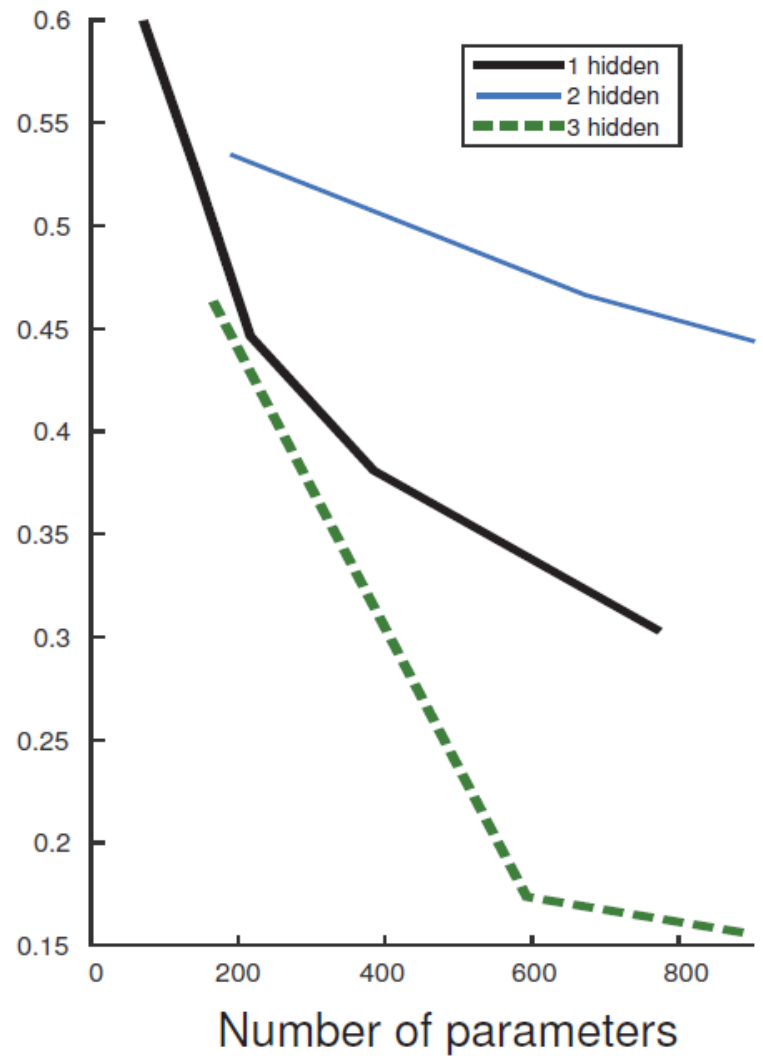
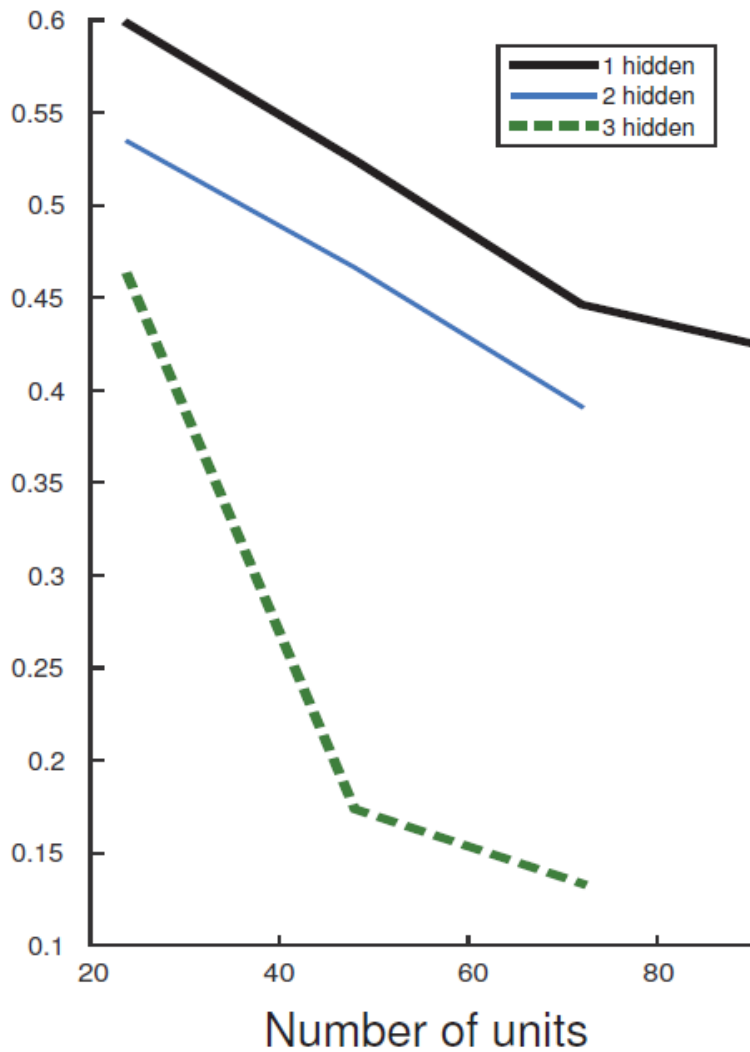
smaller $|\mathcal{H}|$

Shallow



larger $|\mathcal{H}|$

$$f(x) = 2(2\cos^2(x) - 1)^2 - 1$$



Source of image:

<https://www.aaai.org/ocs/index.php/AAAI/AAAI17/paper/viewPaper/14849>

Thinks more

- Deep networks outperforms shallow ones when the required functions are complex and regular.

Image, speech, etc. have this characteristics.

- Deep is exponentially better than shallow even when $y = x^2$.



<https://youtu.be/FN8jclCrqY0>



<https://youtu.be/qpuLxXrHQB4>