# Discovering the Idea of Neural Networks

Or, teaching computers to read.

# A Reading Machine

An image of
a hand-drawn
letter



A letter of
the alphabet

# A Reading Machine



An image of a hand-drawn letter

A probability distribution on the alphabet

The alphabet has 26 letters, so we should expect to output a probability for each letter.



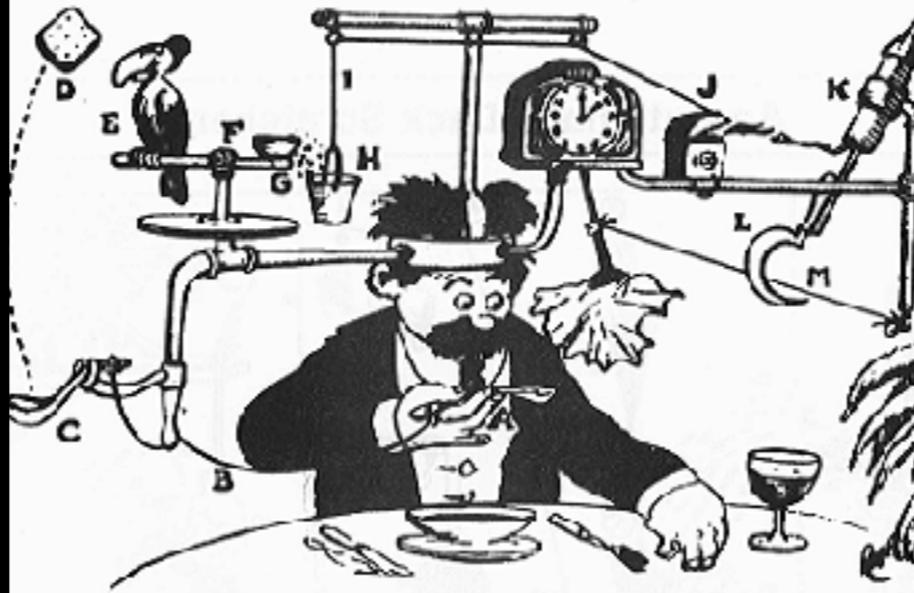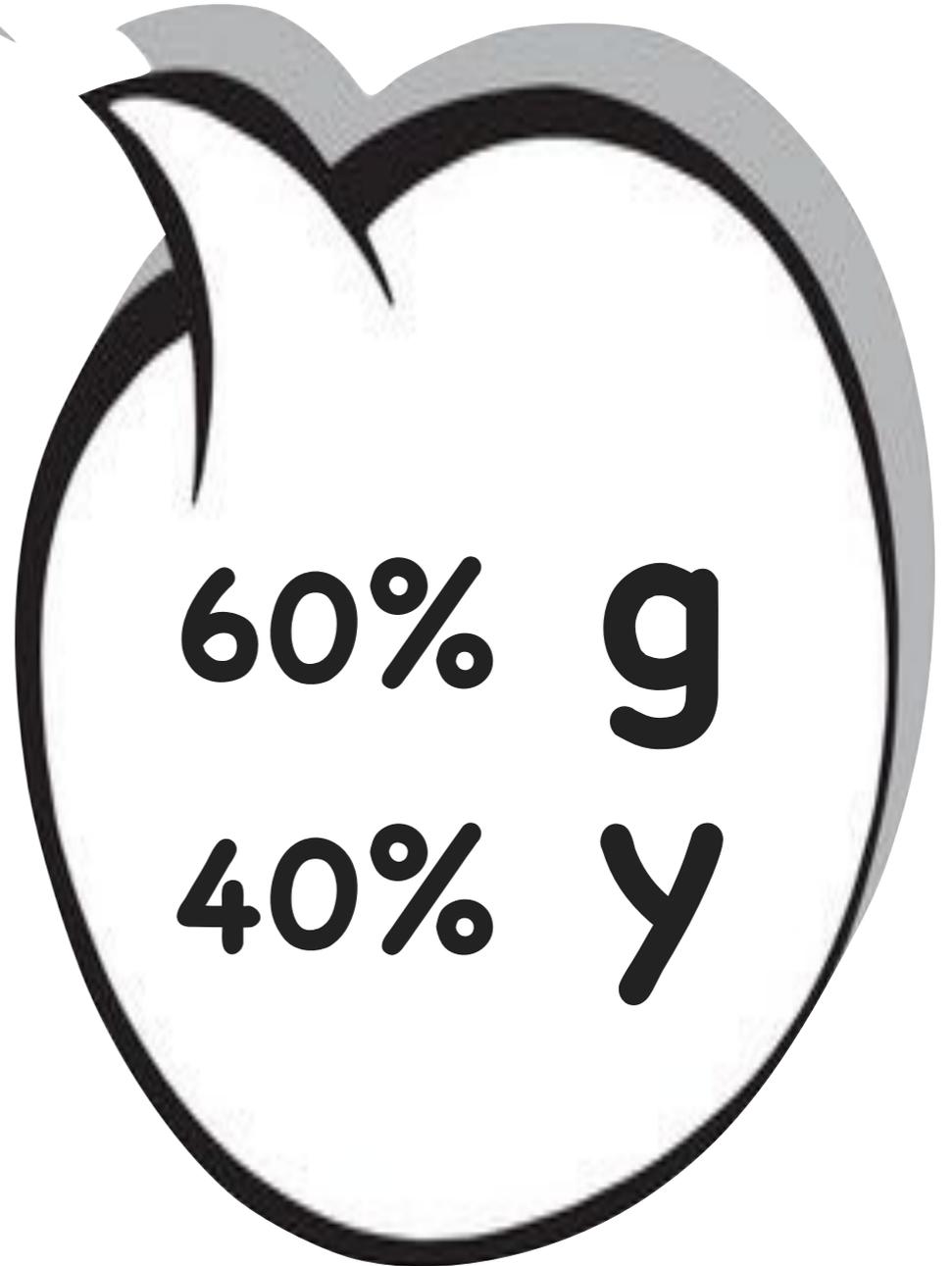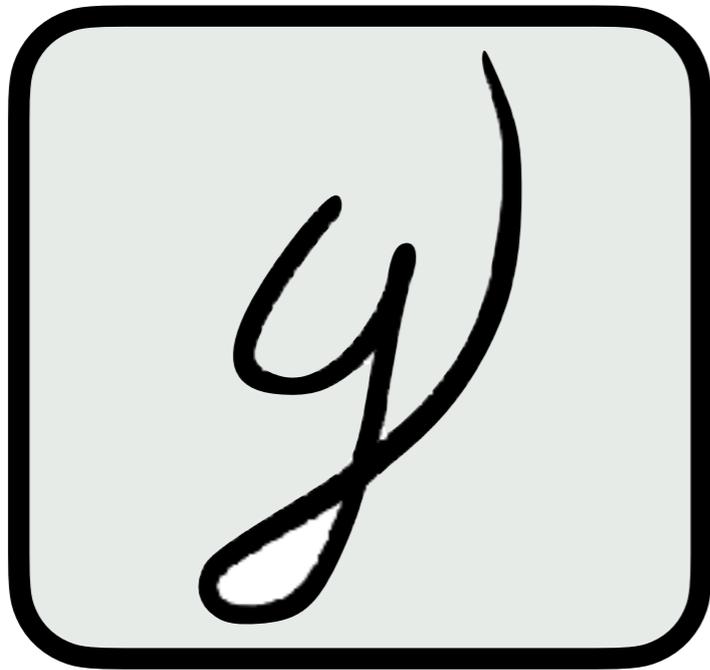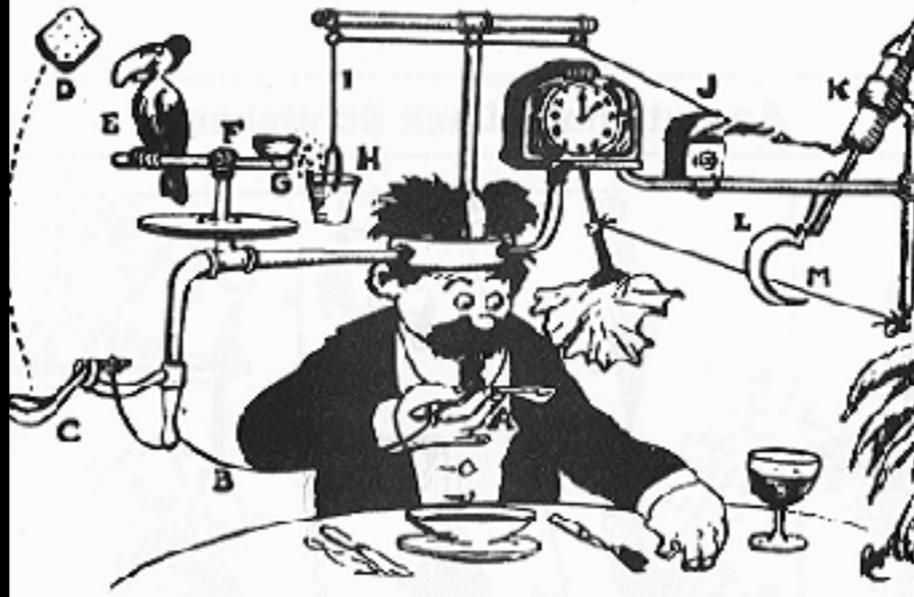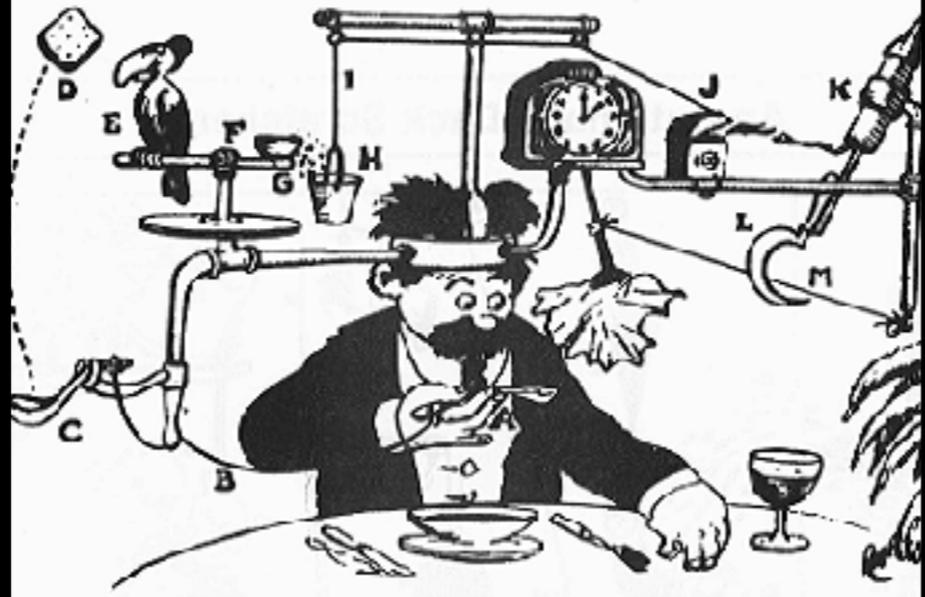This means our output should be a list of 26 nonnegative numbers which sum to 1.

# A Reading Machine
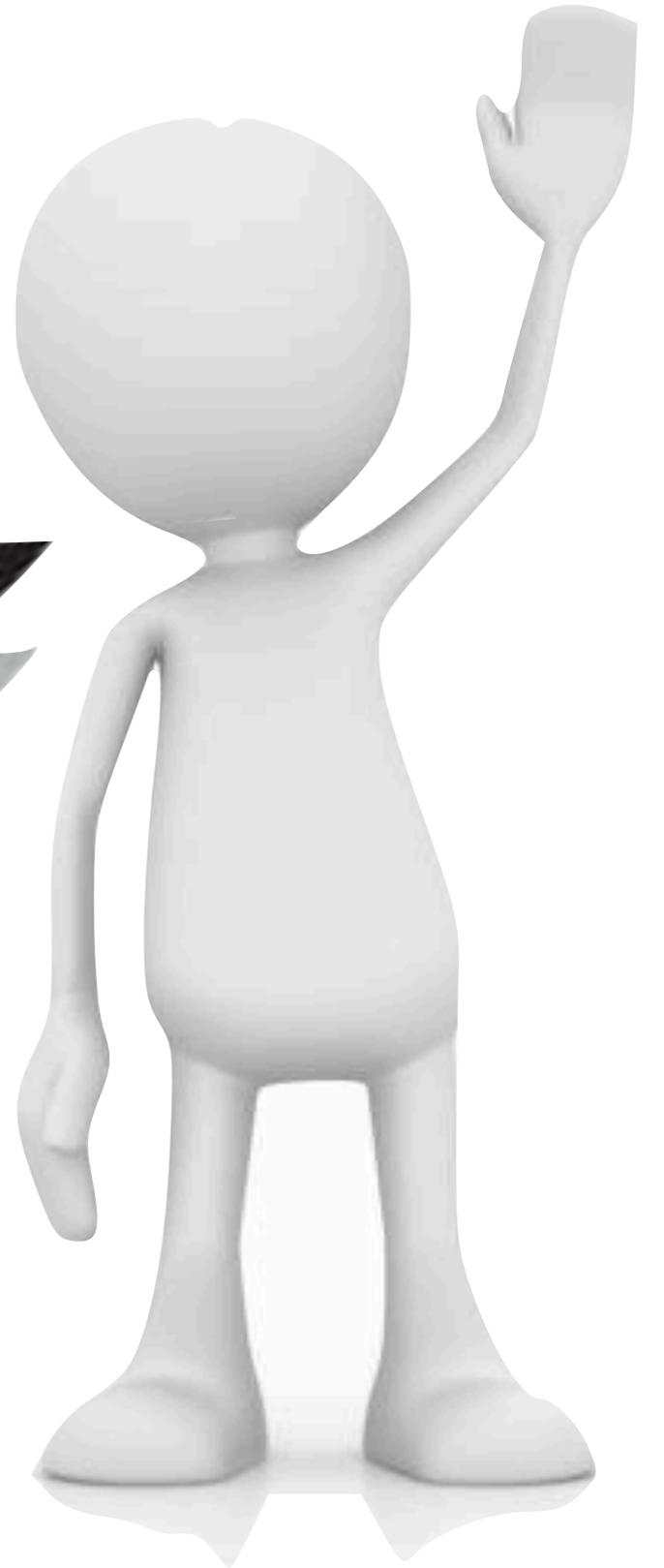
An image of
a hand-drawn
letter



A vector
$(p_1, p_2, p_3, \ldots, p_{26})$
with $p_i \geq 0$
and
$$\sum_{i=1}^{26} p_i = 1$$

If you
look
really
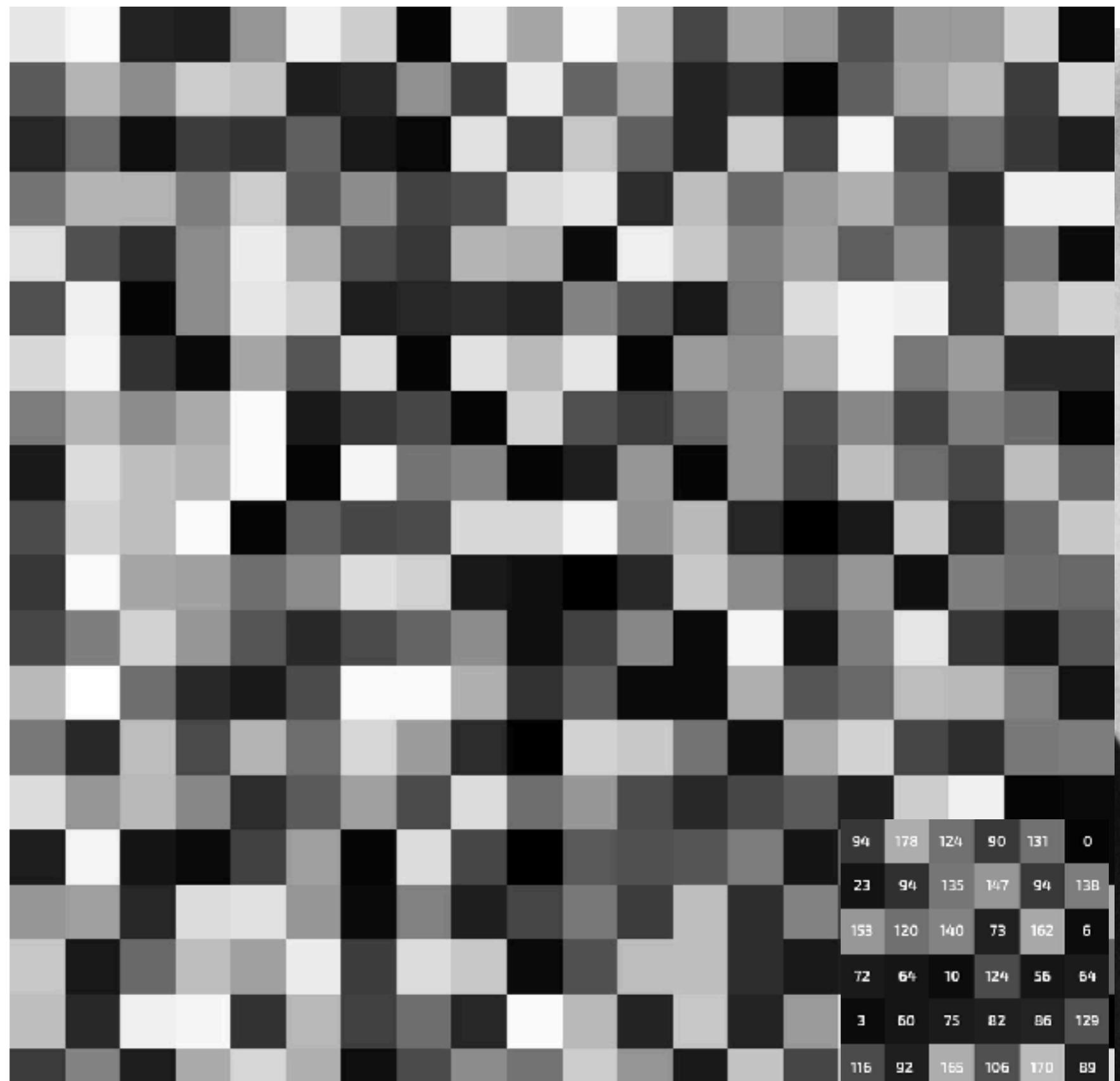close…

Images are just arrays of pixels

| 94 | 178 | 124 | 90 | 131 | 0 |
| 23 | 94 | 135 | 147 | 94 | 138 |
| 153 | 120 | 140 | 73 | 162 | 6 |
| 72 | 64 | 10 | 124 | 56 | 64 |
| 3 | 60 | 75 | 82 | 86 | 129 |
| 116 | 92 | 165 | 106 | 170 | 89 |

And,
Pixels
are
just
numbers

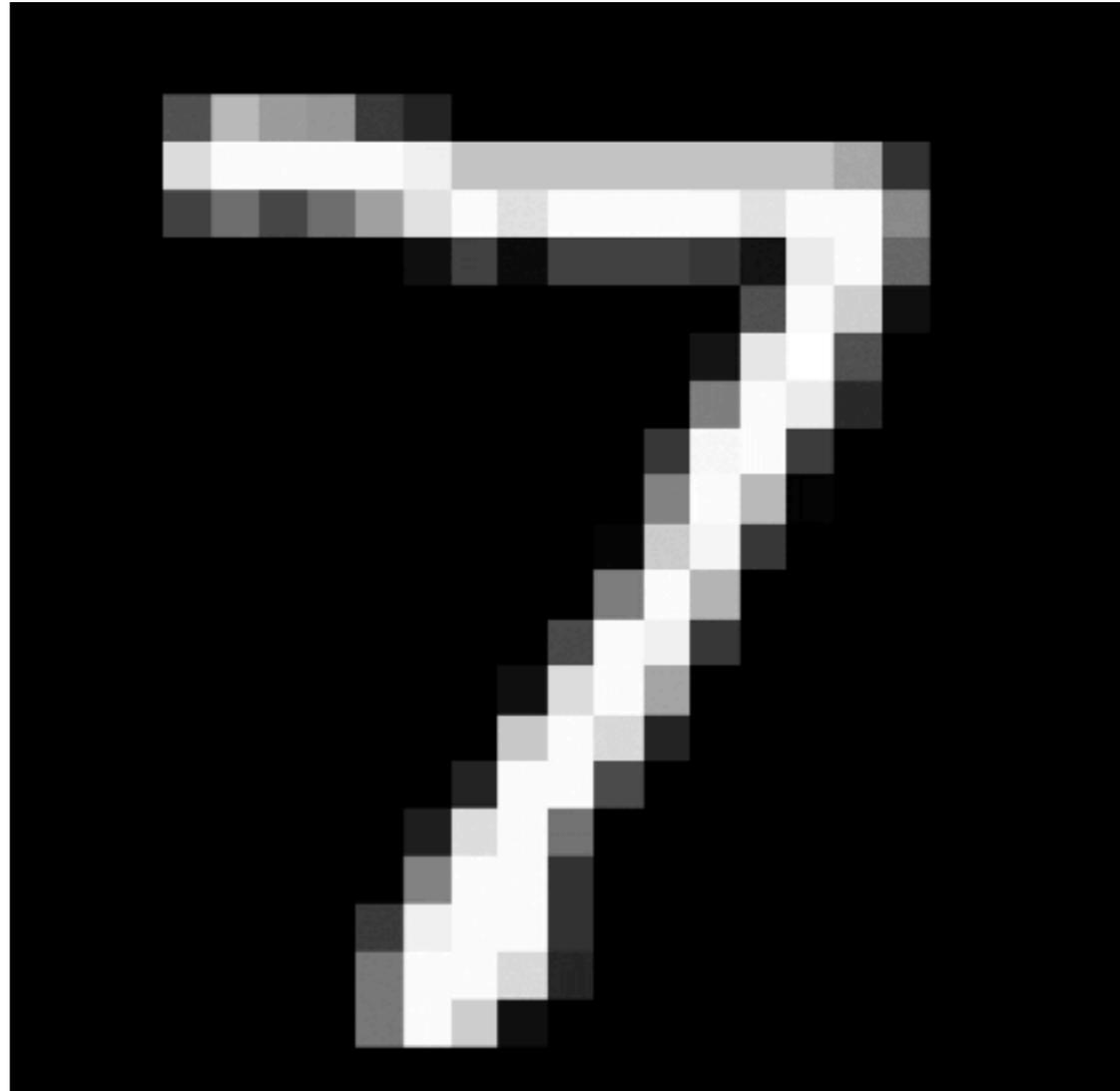| | | | | | |
|---|---|---|---|---|---|
| 94 | 178 | 124 | 90 | 131 | 0 |
| 23 | 94 | 135 | 147 | 94 | 138 |
| 153 | 120 | 140 | 73 | 162 | 6 |
| 72 | 64 | 10 | 124 | 56 | 64 |
| 3 | 60 | 75 | 82 | 86 | 129 |
| 116 | 92 | 165 | 106 | 170 | 89 |

# If an image is p pixels tall and q wide



## It's just an ordered list of pxq numbers!

For example...an HD image is an ordered list of 1920x1080, or 2,073,600 numbers



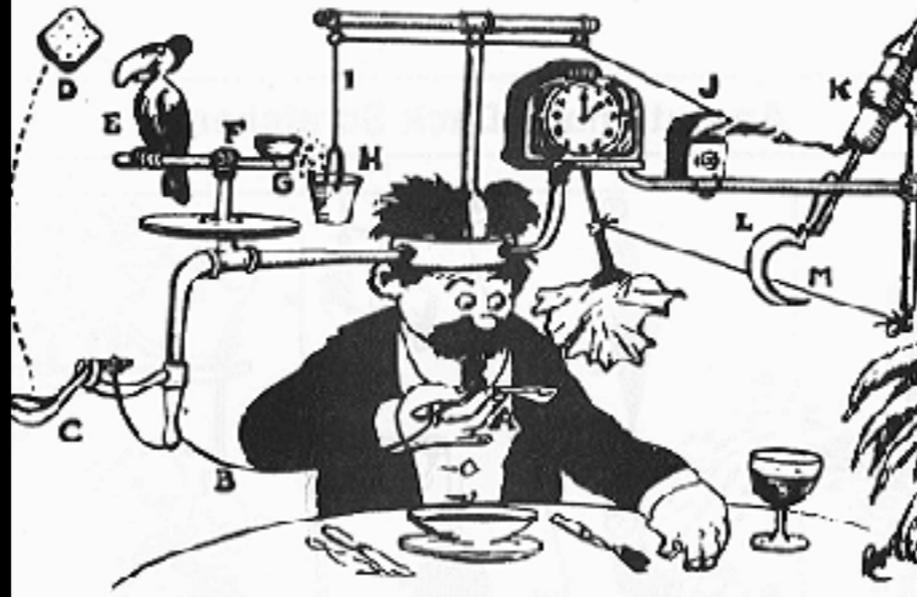This is a vector in 2 million dimensional space!

# An image in the handwriting dataset is 28x28 pixels



This is a vector in 784-dimensional space.

# A Reading Machine

A 784-dimensional vector



A vector
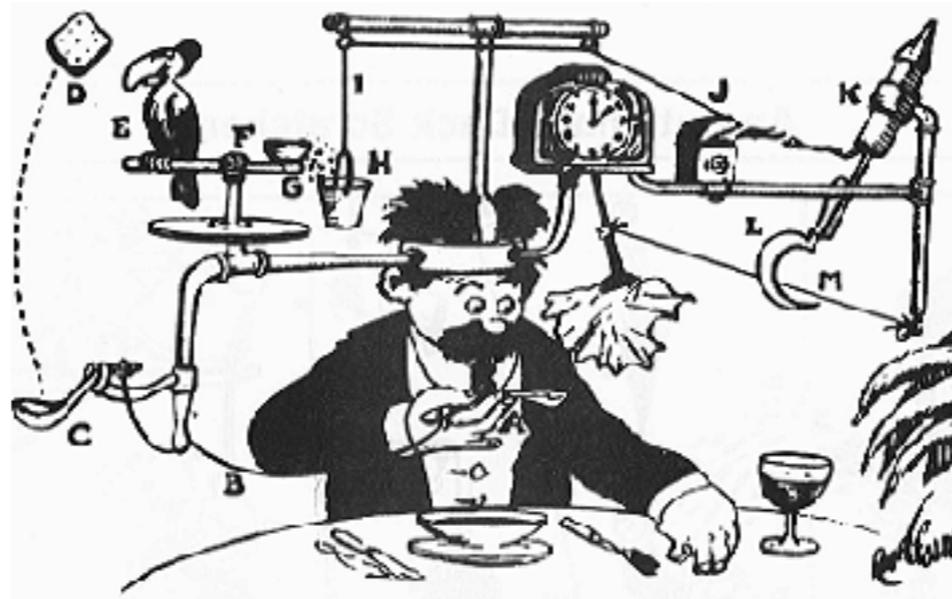$$(p_1, p_2, p_3, \ldots, p_{26})$$
with $p_i \geq 0$
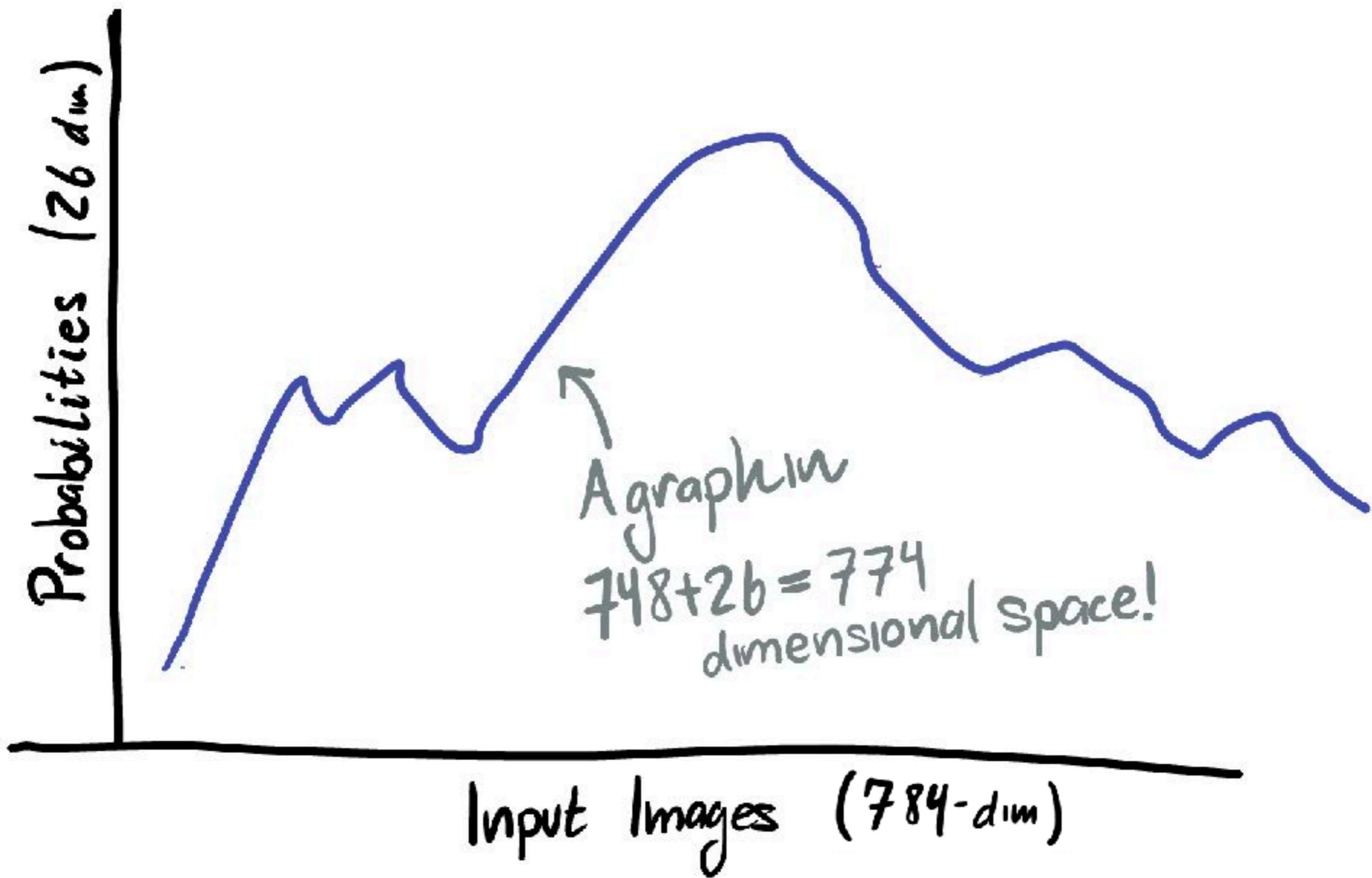and
$$\sum_{i=1}^{26} p_i = 1$$

This sounds really daunting: there are so many high dimensions floating around.
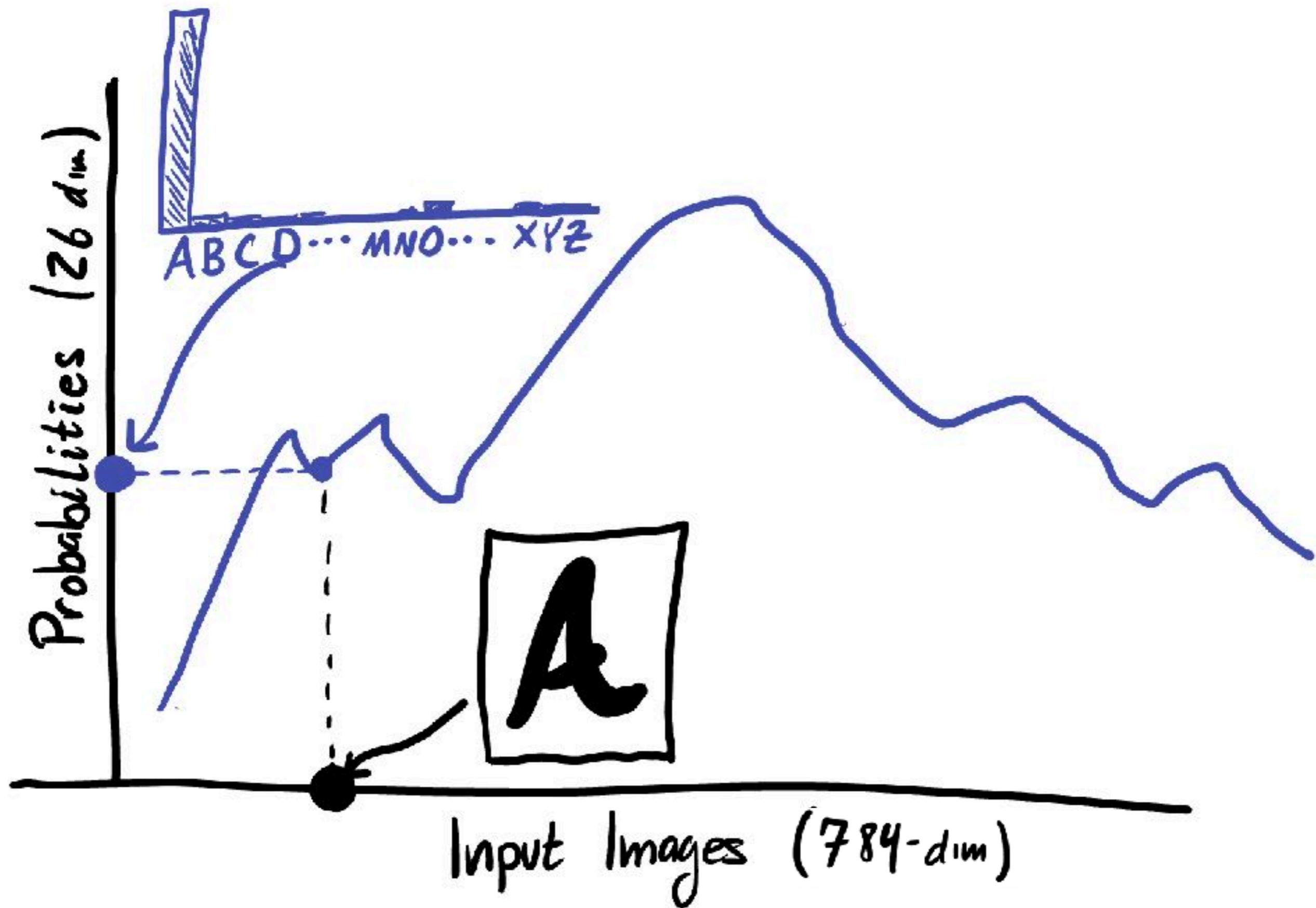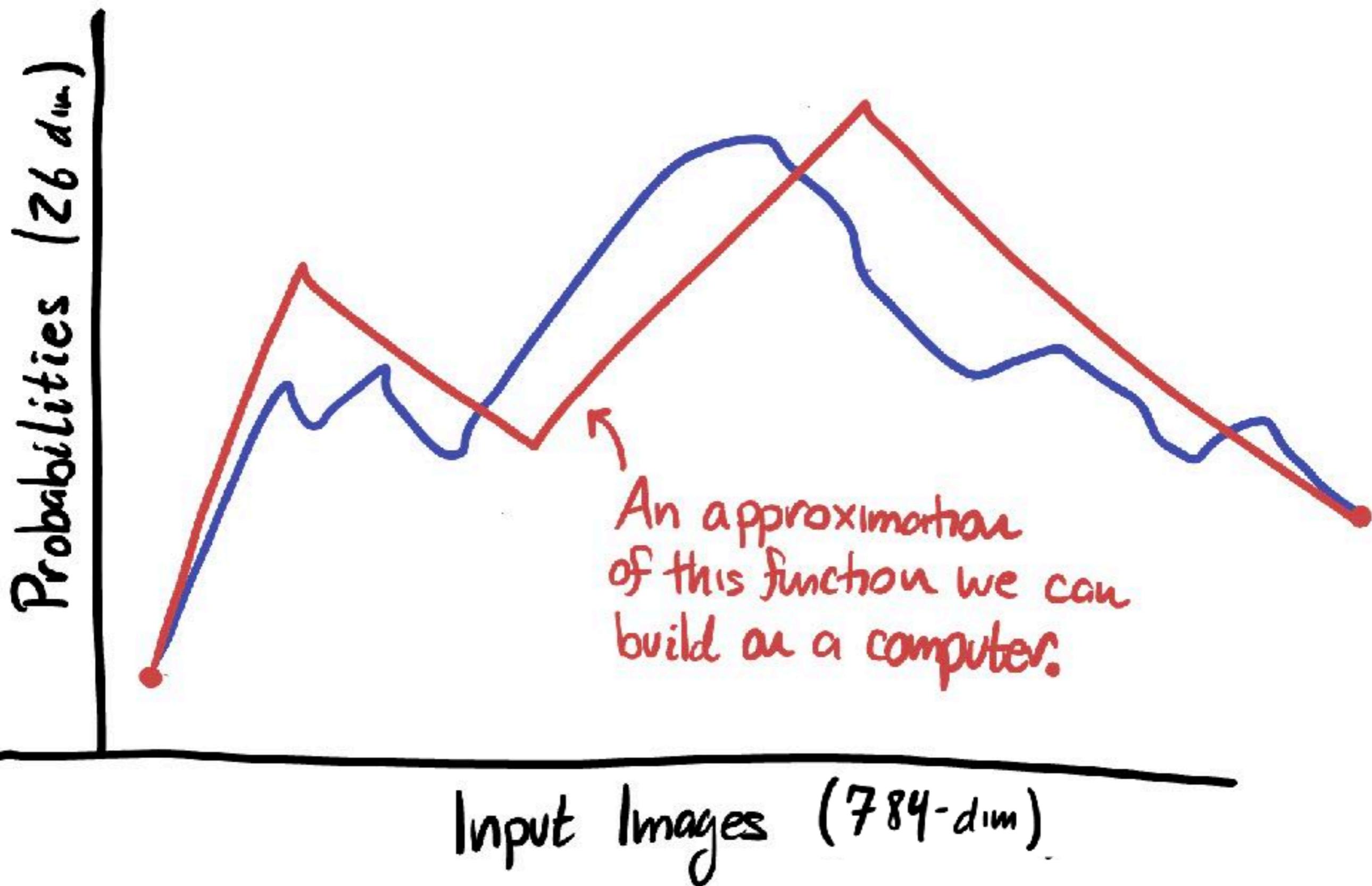


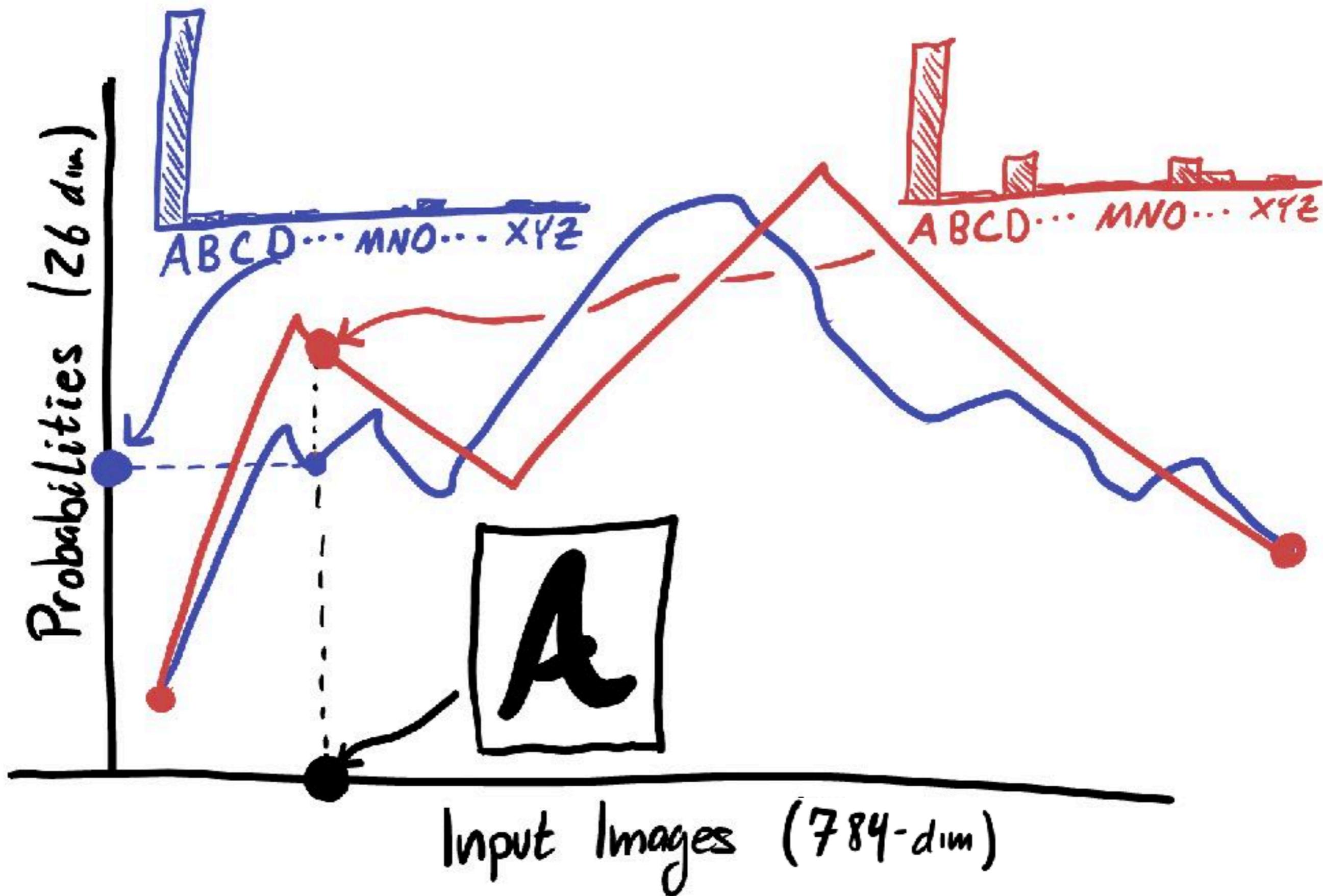But your brain performs this function without trouble!

How can we draw a schematic of
**The Reading Function**
to reason about?

Probabilities (26 dim)
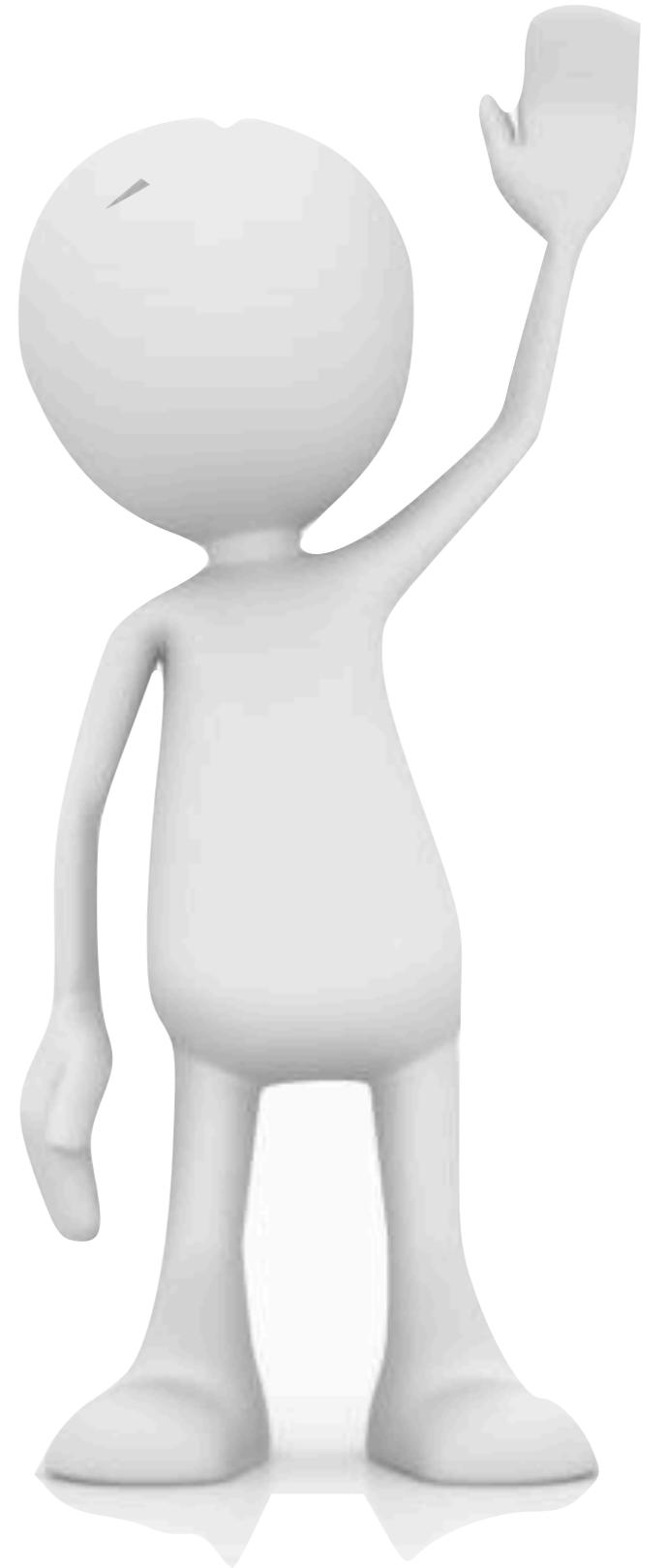
A B C D ··· M N O ··· X Y Z

Input Images (784-dim)

Probabilities (26 dim)

Input Images (784-dim)

An approximation of this function we can build on a computer.

Probabilities (26 dim)

Input Images (784-dim)

"Training Data"

Probabilities (26 dim)

ABCD···MNO····XYZ

Input Images (784-dim)

**Top chart:** Probabilities (26 dim) vs. Input Images (784-dim)

This approximation's pretty good!

**Bottom chart:** Probabilities (26 dim) vs. Input Images (784-dim)

This one is _NOT_

# The Loss Function

Probabilities (26 dim)

Input Images (784-dim)

Goal: Measure the distance between the collection of blue dots and red dots.

Each dot is a
**Probability Distribution!**



Probabilities (26 dim)

A B ··· G H I ··· M N ··· Y Z

A B ··· G H I ··· M N ··· Y Z

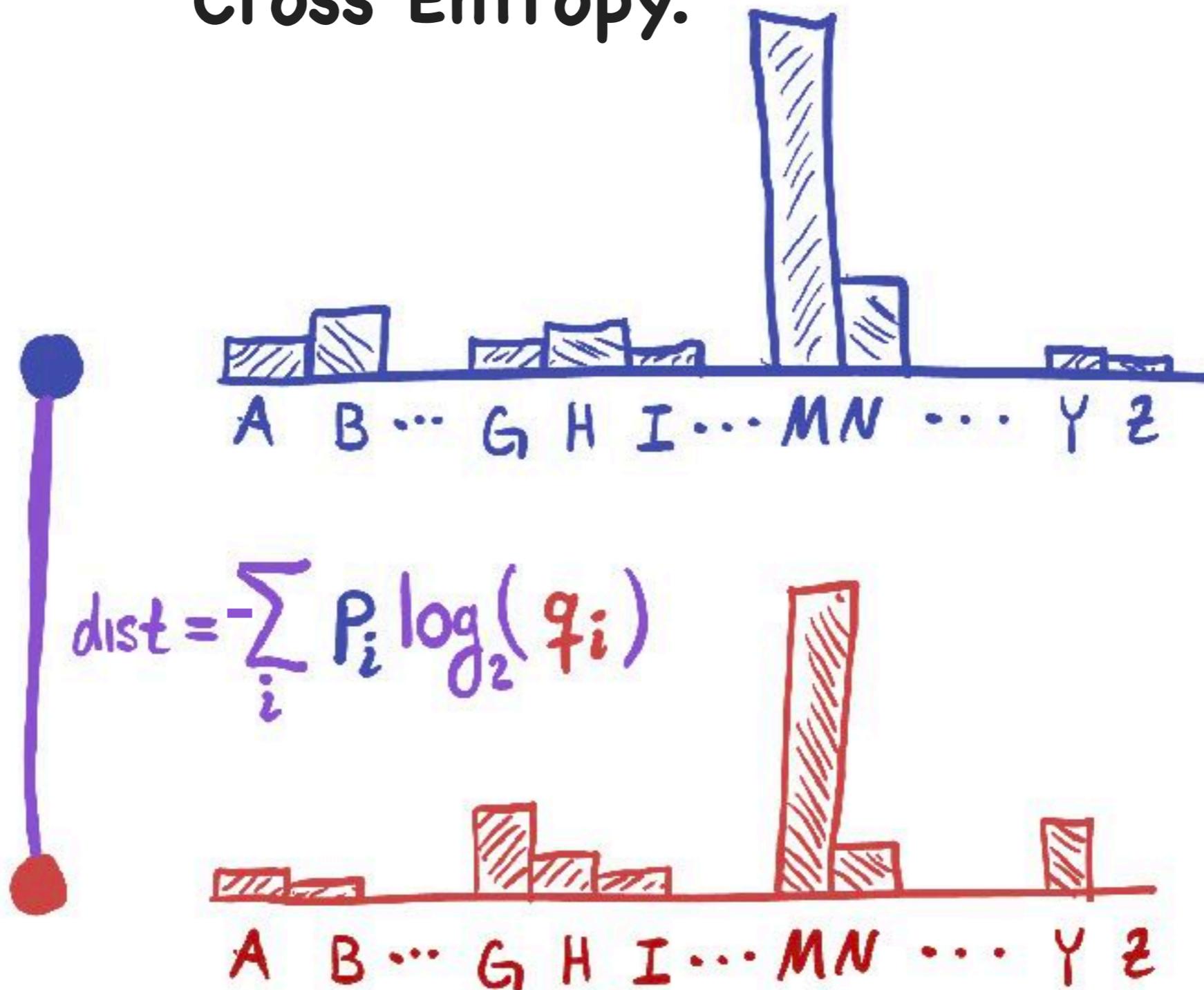Luckily, mathematicians know how to compare probability distributions.  One option is the **Cross Entropy.**



$$dist = -\sum_i P_i \log_2(q_i)$$

Probabilities (26 dim)

A B ⋯ G H I ⋯ M N ⋯ Y Z

A B ⋯ G H I ⋯ M N ⋯ Y Z

The difference between our "true function" and the approximation we are considering is called the "loss" – its how much info our approximation has lost.
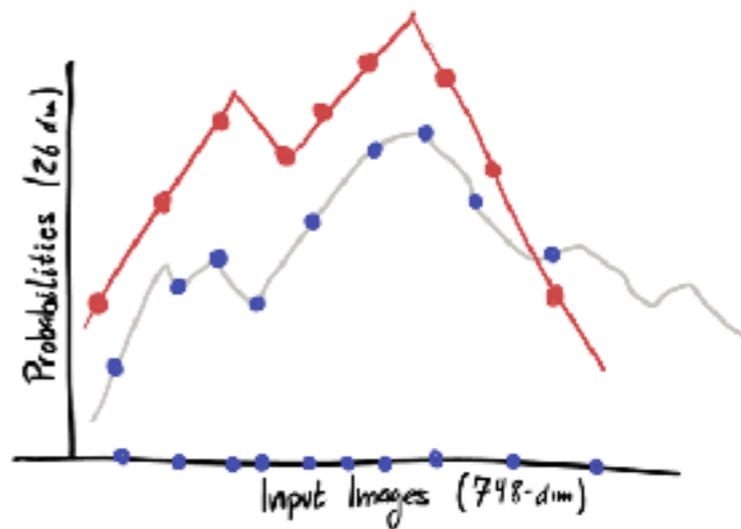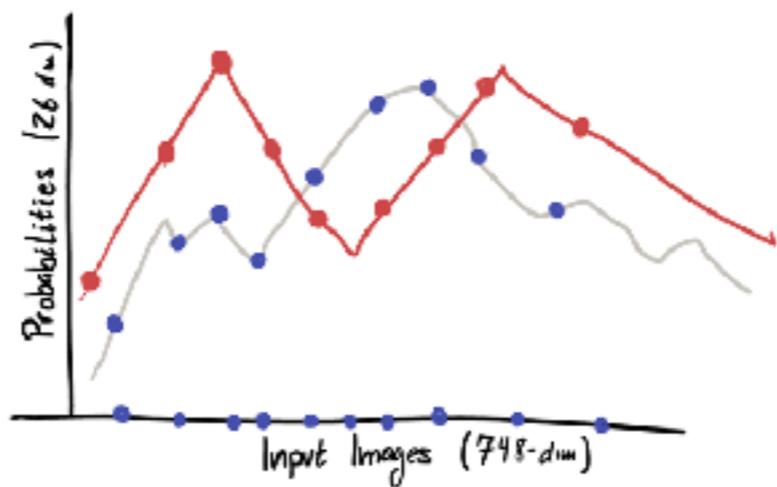


Loss of the approximation

The best approximation is the one that minimizes the loss function.

8

# Function composition is a way to build complicated functions from easy pieces:

$$\left. \begin{array}{l} e^x \\ \sin x \\ \cos x \\ \sqrt{x} \end{array} \right\} \longrightarrow \sqrt{e^{\sin\left(\sqrt{\cos e^x}\right)}}$$

# Function composition is a way to build complicated functions from easy pieces:

$x$ → $\boxed{\exp}$ → $\boxed{\cos}$ → $\boxed{\sqrt{\phantom{x}}}$ → $\boxed{\sin}$

$\sqrt{e^{\sin\sqrt{\cos x}}}$ ← $\boxed{\sqrt{\phantom{x}}}$ ← $\boxed{\exp}$ ←

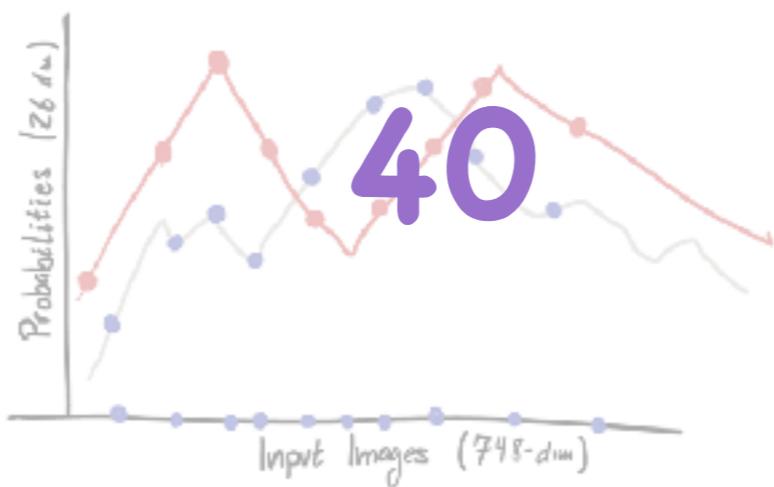Because composition "strings functions together" we call the result a **network.**

# A network of linear pieces:

# A network of linear pieces:

Probabilities (26-dim)

Input Images (798-dim)

Image

W

LINEAR MAP

Probability

ABC --- MNO ... WXYZ

# What happens if we surround some linear functions with the absolute value, before composing them?

$f(x) = |3x - 2|$

$g(x) = \left| \dfrac{x}{3} - 2 \right|$

$h(x) = |x - 3|$

# The result is much more complicated than any of the individual building blocks!

$f(x) = |3x - 2|$

$g(x) = \left| \dfrac{x}{3} - 2 \right|$

$h(x) = |x - 3|$

$g\big(h(f(x))\big)$

# A proposed "brick"



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

$$\sigma = |x| \qquad \sigma = e^x/(1 + e^x)$$

# Building our Castle out of Bricks:
The architecture of a Neural Network



Each linear/nonlinear combo is called a **layer.**

**First:** make sure the final linear map is of the right dimension! So for us, it needs to map into 26-dimensional space.

**Second:** vectors can have any real numbers as entries, but probabilities must be positive. We much choose a way to make all entries positive.

**Third:** Probabilities add to 1: we must rescale the vector so that its entires add to 1.

# The Softmax Function

$$\vec{v} = \langle v_1, v_2, \ldots, v_{26} \rangle$$

$$\langle e^{v_1}, e^{v_2}, \ldots, e^{v_{26}} \rangle$$

$$\text{softmax}(\vec{v}) = \frac{\langle e^{v_1}, e^{v_2}, \ldots, e^{v_{26}} \rangle}{e^{v_1} + e^{v_2} + \cdots + e^{v_{26}}}$$

# Building our Castle out of Bricks:
## The architecture of a Neural Network

# Universal Approximation

## Theorem:

Using enough* layers and large enough* matrices, you can approximate any* function to whatever accuracy** you need



*Warning: I'm ignoring a TON of technical details here...
**With accuracy measured by the Loss function!

Somewhere in the space of all such functions, there are ones that approximate as accurately as you wish.

# A geometric perspective



Neural Networks are Like Disco Balls

# A geometric perspective



To approximate a complicated shape, like the graph of a classifier function,

# A geometric perspective

You can use simple linear pieces joined together nonlinearly..

# A geometric perspective



And doing this enough times gives a really good approximation!

This approximation scheme is **universal** so it can approximate any shape you like!

It's possible for computers to learn to read!!!!!!

By universal approximation, there is some sequence of layers that approximates the actual "reading function" to arbitrary precision.

**Two different architectures.**

An **architecture** is a particular way of stringing together maps: you've fixed a number of linear maps, the size of each linear map, and the type of nonlinearity.

A linear map can be represented by a **matrix**. So each linear map (layer) of our function (network) can be specified by a list of numbers.

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

$$\begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{bmatrix}$$

$$\begin{bmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{bmatrix}$$

Doing this to all the linear maps, an entire network can be specified by a list of numbers. These are called the **weights of the network.**

$$\begin{bmatrix} a_{11} & a_{12} & a_{21} & a_{22} & b_{11} & b_{12} & b_{21} & b_{22} & b_{31} & b_{32} & c_{11} & c_{12} & c_{21} & c_{22} \end{bmatrix}$$

Amazing idea! The **space of all networks of a fixed architecture** is a vector space!

For example, this architecture is specified by 14 numbers, so is a point in 14-dimensional space.

$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

$$\begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \\ b_{31} & b_{32} \end{bmatrix}$$

$$\begin{bmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{bmatrix}$$

$$\begin{bmatrix} a_{11} & a_{12} & a_{21} & a_{22} & b_{11} & b_{12} & b_{21} & b_{22} & b_{31} & b_{32} & c_{11} & c_{12} & c_{21} & c_{22} \end{bmatrix}$$

An example
architecture we
will actually
build!



```python
class SimpleNeuralNetwork(torch.nn.Module):

  def __init__(self):

    super().__init__()
    self.dense1 = torch.nn.Linear(784,100)
    self.dense2 = torch.nn.Linear(100,300)
    self.dense3 = torch.nn.Linear(300,10)
```

This has 784x100+ 100x300 + 300x10
= 111,400 weights!

Slightly tweaking a neural network is as easy as slightly changing the numbers that determine it!

**Weights B1**

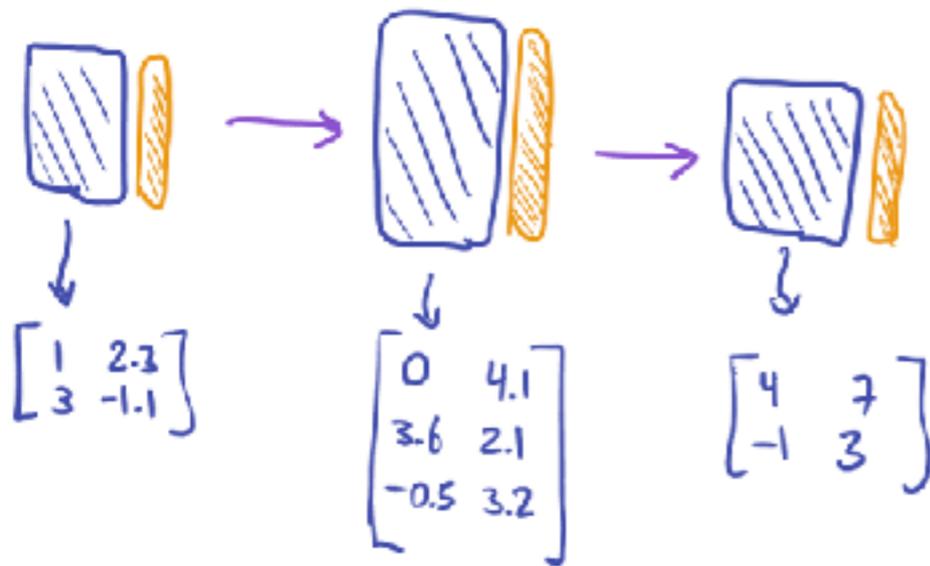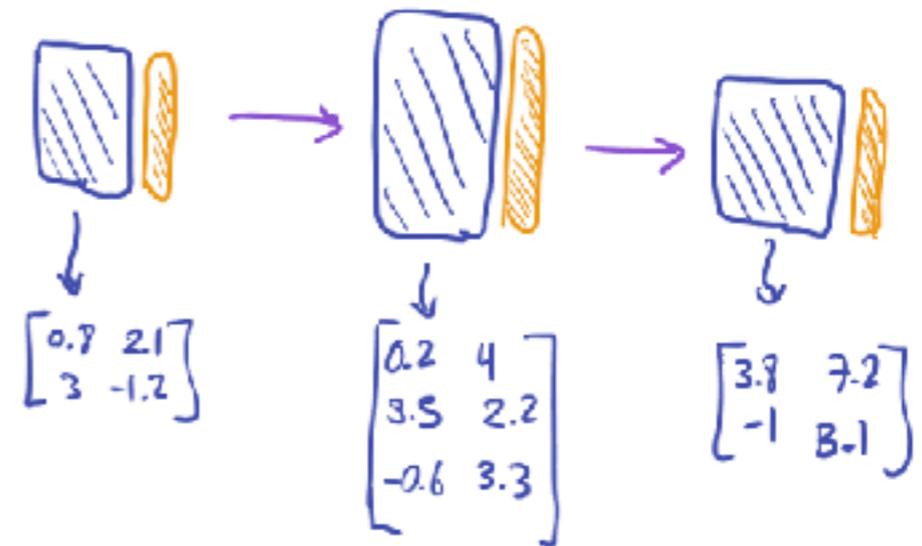$$\begin{bmatrix} 1 & 2.3 \\ 3 & -1.1 \end{bmatrix} \quad \begin{bmatrix} 0 & 4.1 \\ 3.6 & 2.1 \\ -0.5 & 3.2 \end{bmatrix} \quad \begin{bmatrix} 4 & 7 \\ -1 & 3 \end{bmatrix}$$

**Weights B2**

$$\begin{bmatrix} 0.8 & 2.1 \\ 3 & -1.2 \end{bmatrix} \quad \begin{bmatrix} 0.2 & 4 \\ 3.5 & 2.2 \\ -0.6 & 3.3 \end{bmatrix} \quad \begin{bmatrix} 3.8 & 7.2 \\ -1 & 8.1 \end{bmatrix}$$
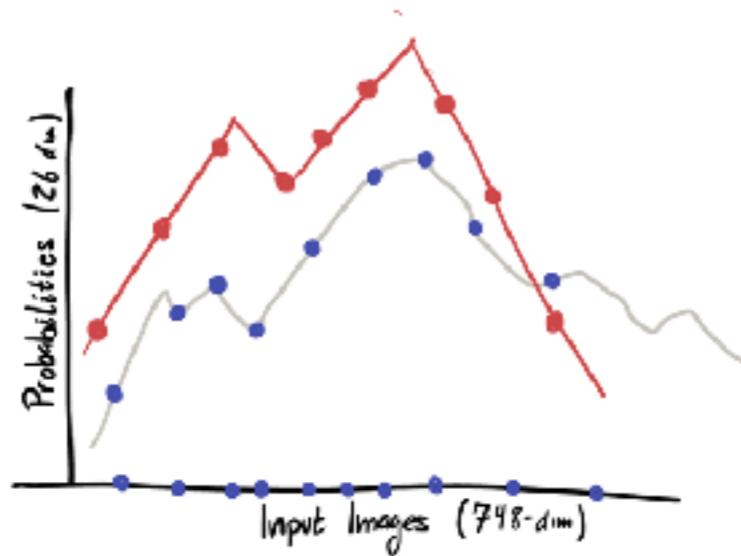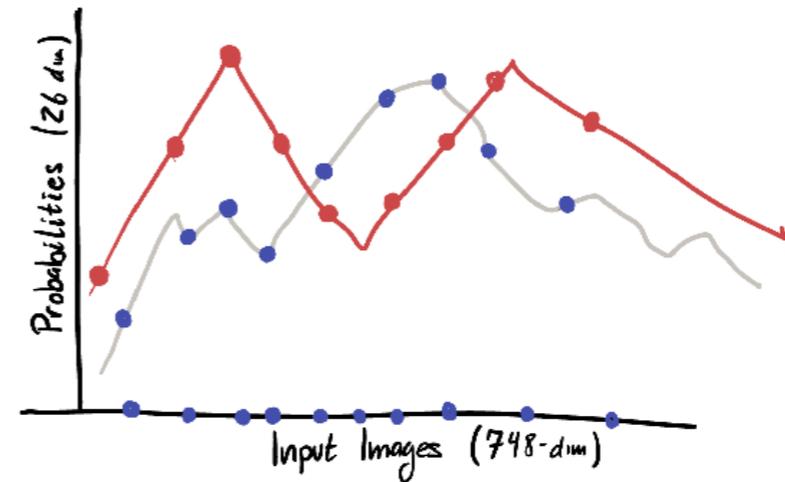
So, we can try to change the numbers a bit, and see if the new network does better or worse with respect to our loss function.

Slightly tweaking a neural network is as easy as slightly changing the numbers that determine it!
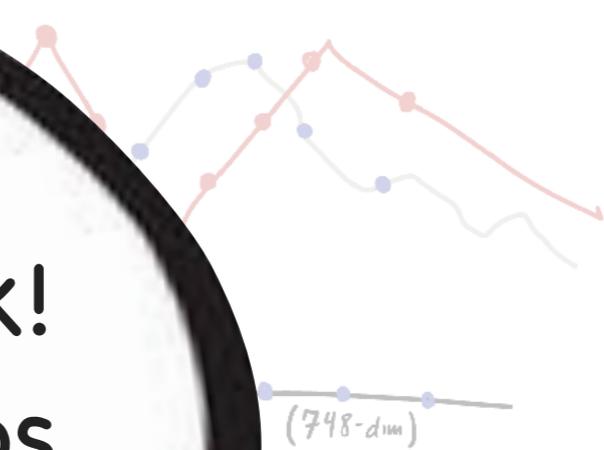


**loss(B1) = 6**



**loss(B2) = 3**

So, we can try to change the numbers a bit, and see if the new network does better or worse with respect to our loss function.

Original
Network

B

Slightly adjusted
Network

B+h

| Original Network | Slightly adjusted Network |
|:---:|:---:|
| $B$ | $B+h$ |

| Original Loss | Slightly adjusted Loss |
|:---:|:---:|
| $loss(B)$ | $loss(B+h)$ |

Net improvement of the change
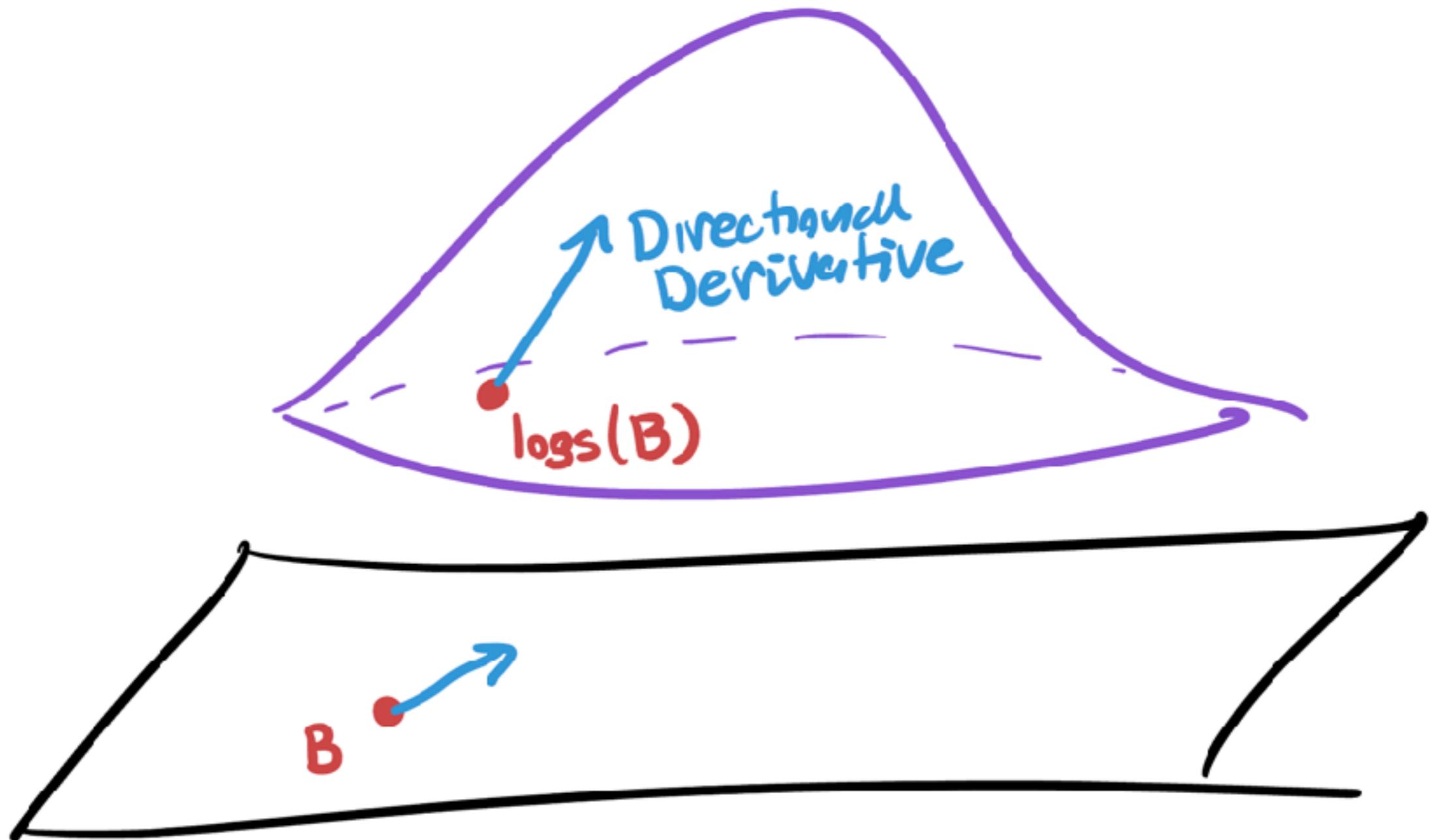(normalized by size of change).

$$\frac{\text{loss(B+h)}-\text{loss(B)}}{h}$$

Net improvement of the change (normalized by size of change).

$$\frac{\text{loss}(B+h) - \text{loss}(B)}{h}$$

This is the **directional derivative of loss in the direction h,** on the space of networks!

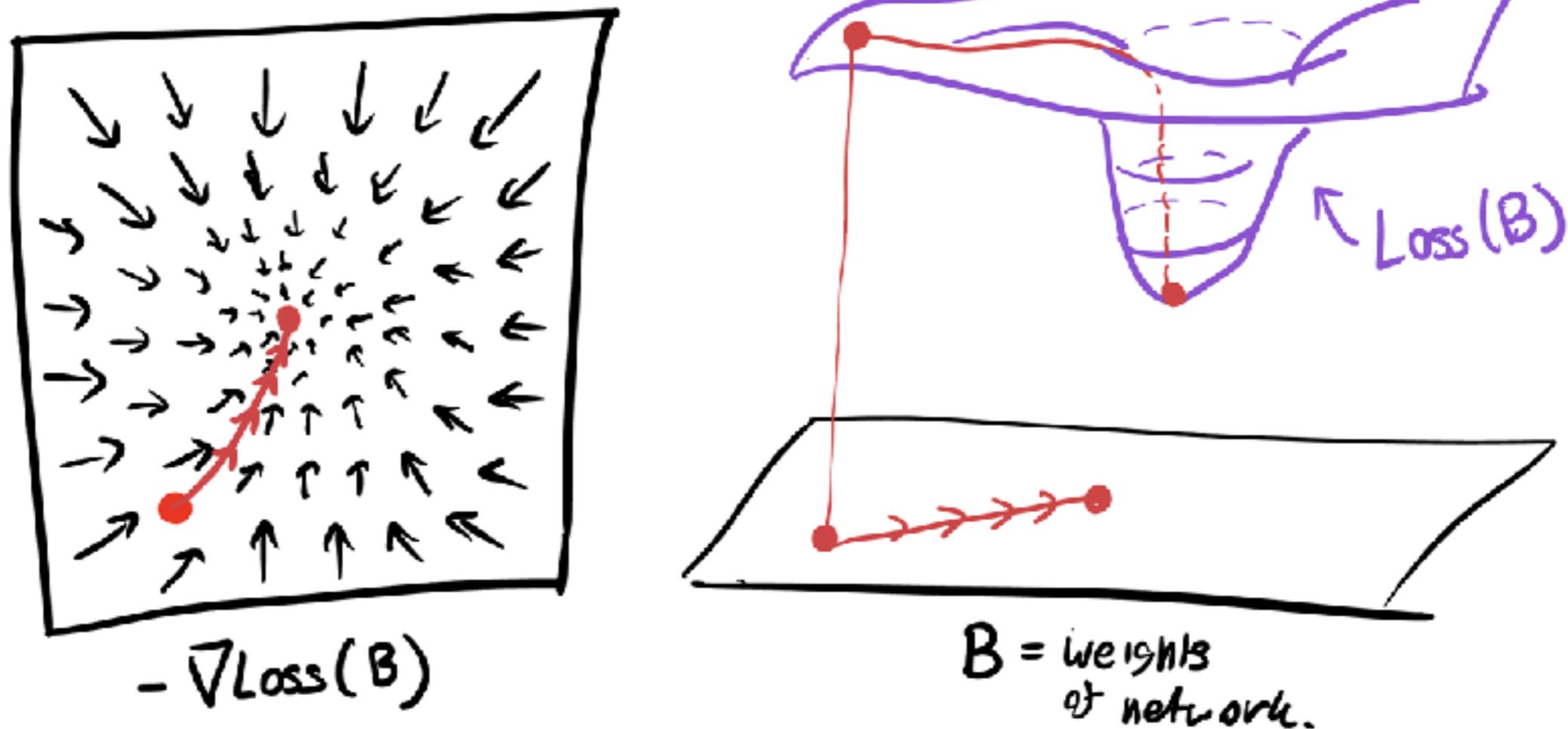This is the **directional derivative of loss in the direction h,** on the space of networks!

# **Gradient Descent:** Following the negative gradient leads to a minimum - a set of weights which minimize the loss!



$$-\nabla Loss(B)$$

Loss(B)
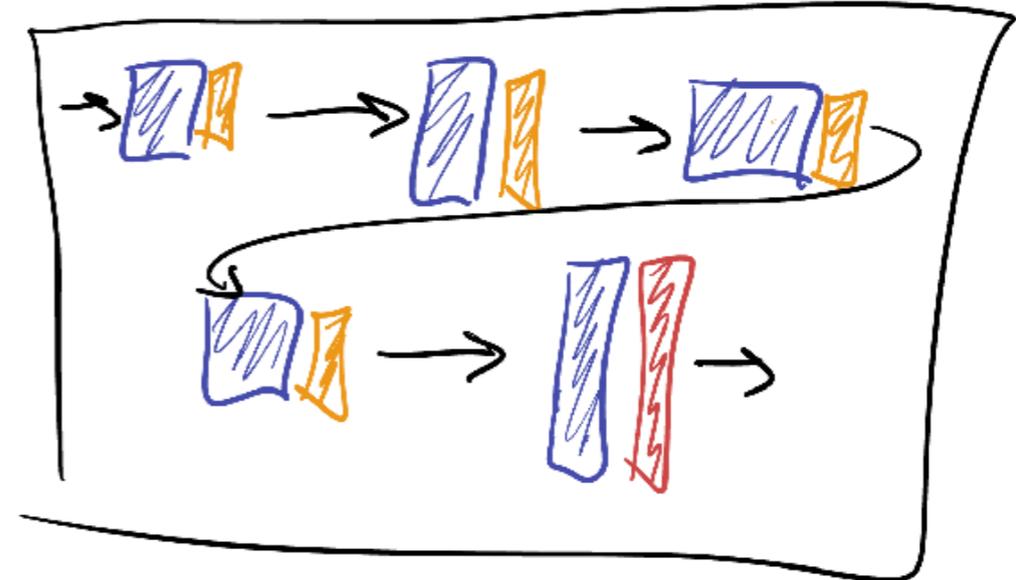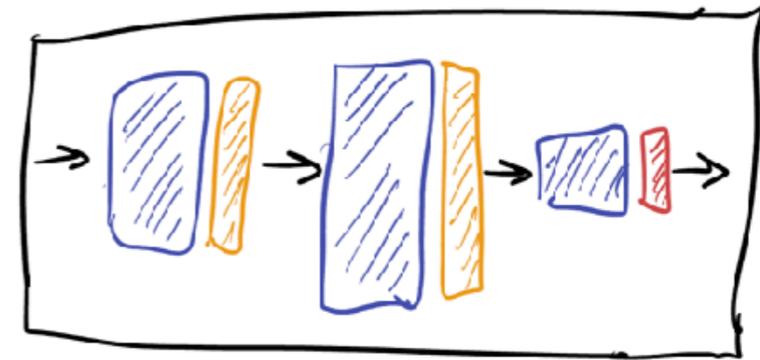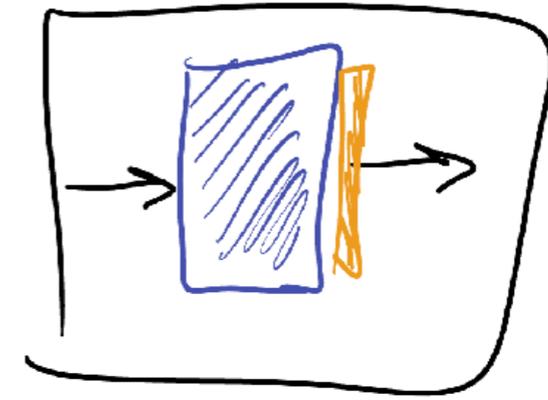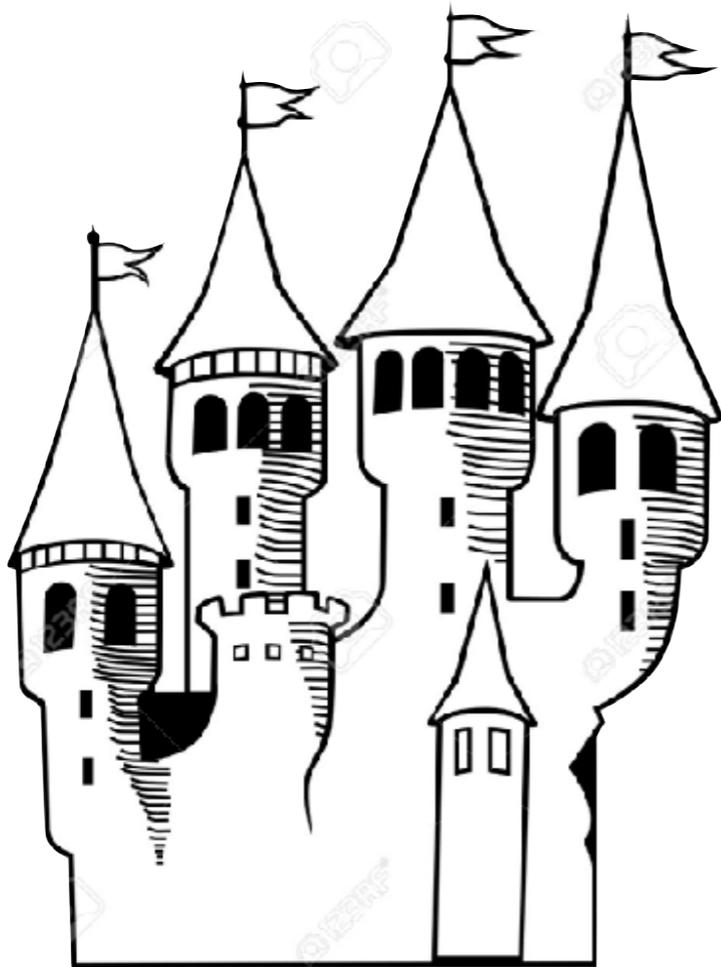
$$B = \text{weights of network.}$$

Universal Approximation doesn't tell us how many bricks (layers) or what size (number of weights) are needed...just that there is some combination that works.