

Deep Learning

Lecture 6:

Intro to DL

Young & Yandex



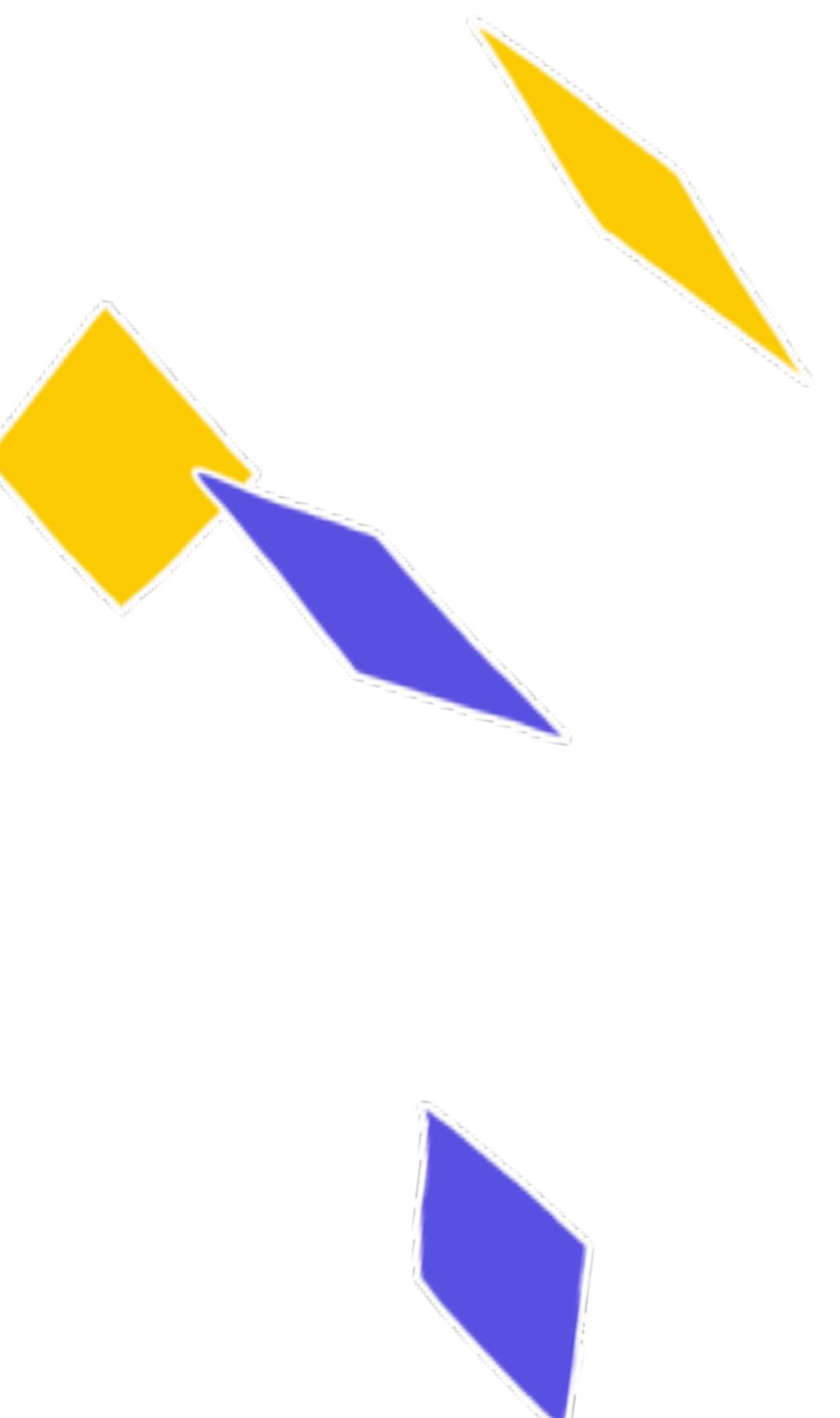
girafe
ai

Radoslav Neychev

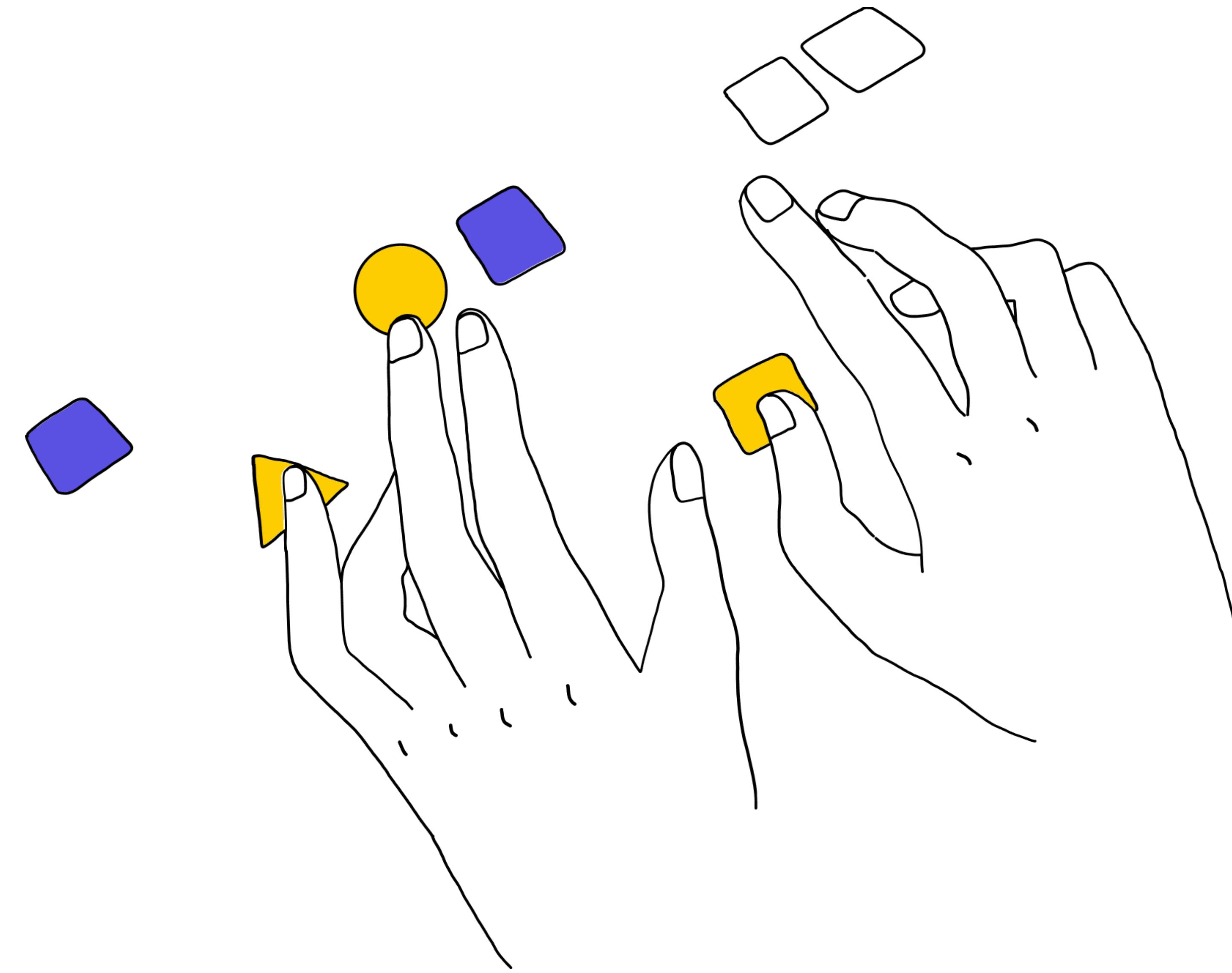


Outline

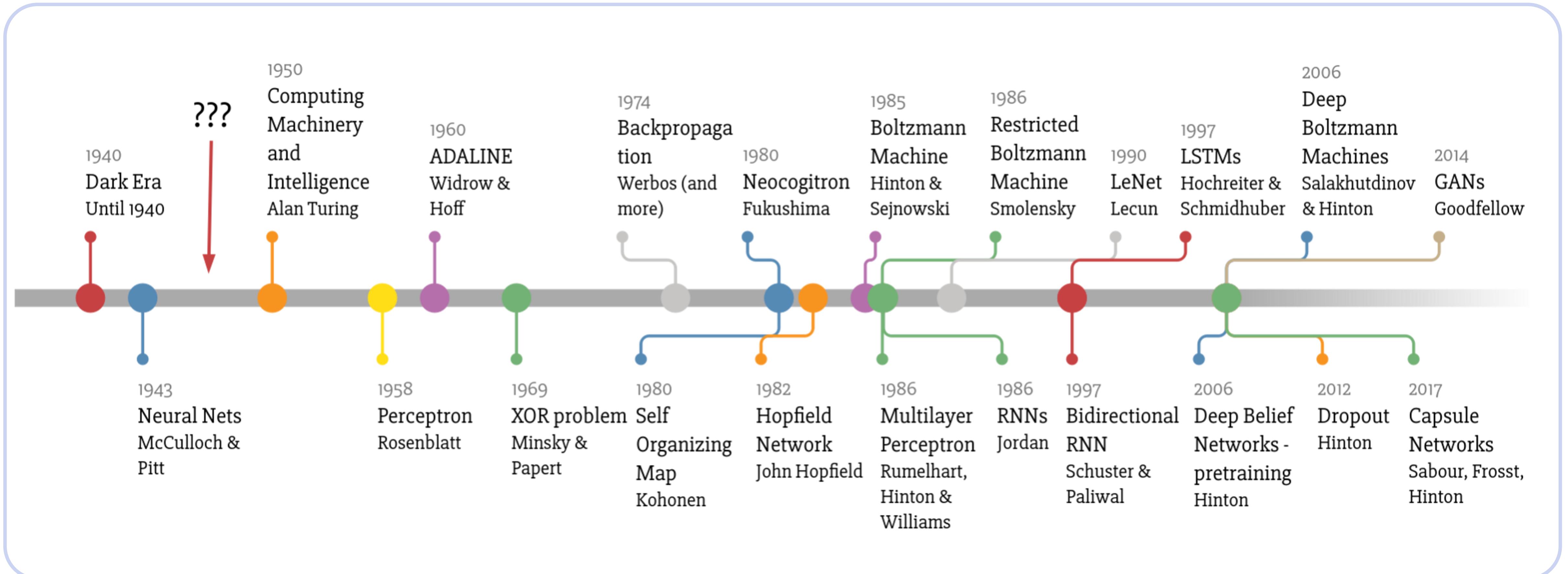
- 01** Neural Networks in different areas.
Historical overview
- 02** Backpropagation
- 03** More on backpropagation
- 04** Activation functions
- 05** Playground



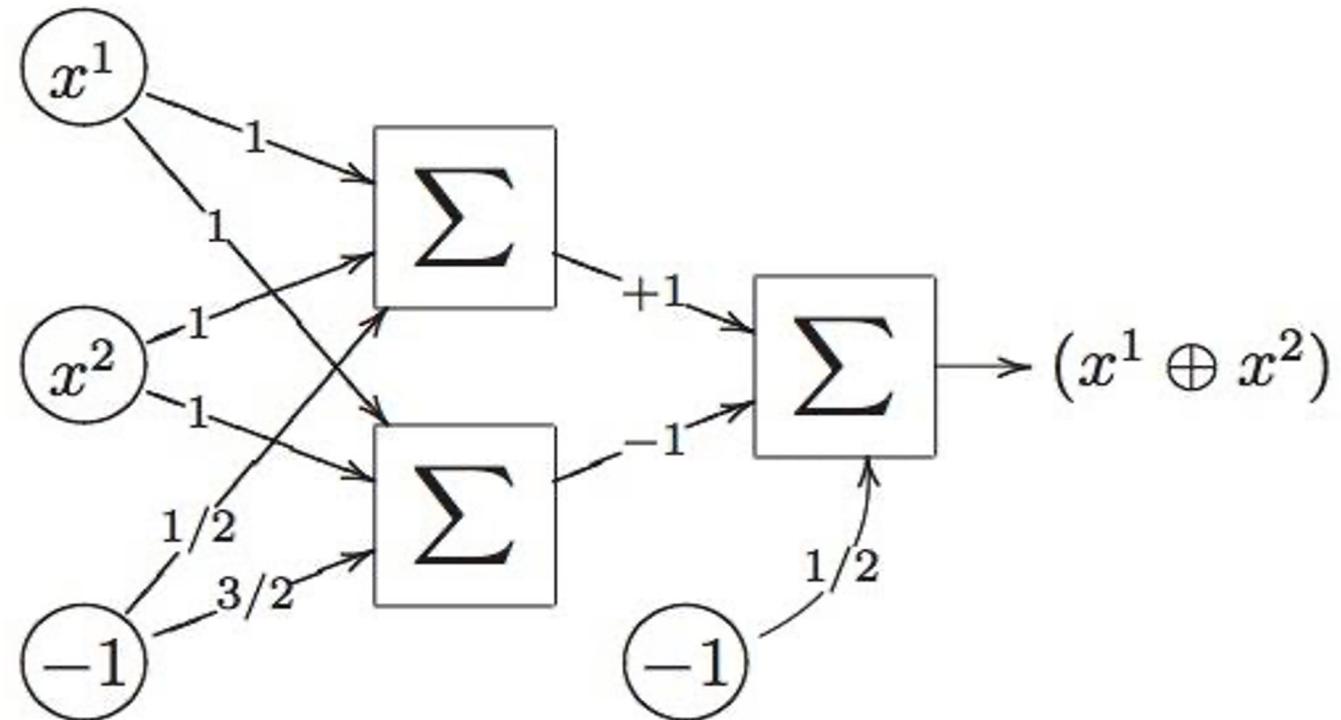
History of Deep Learning



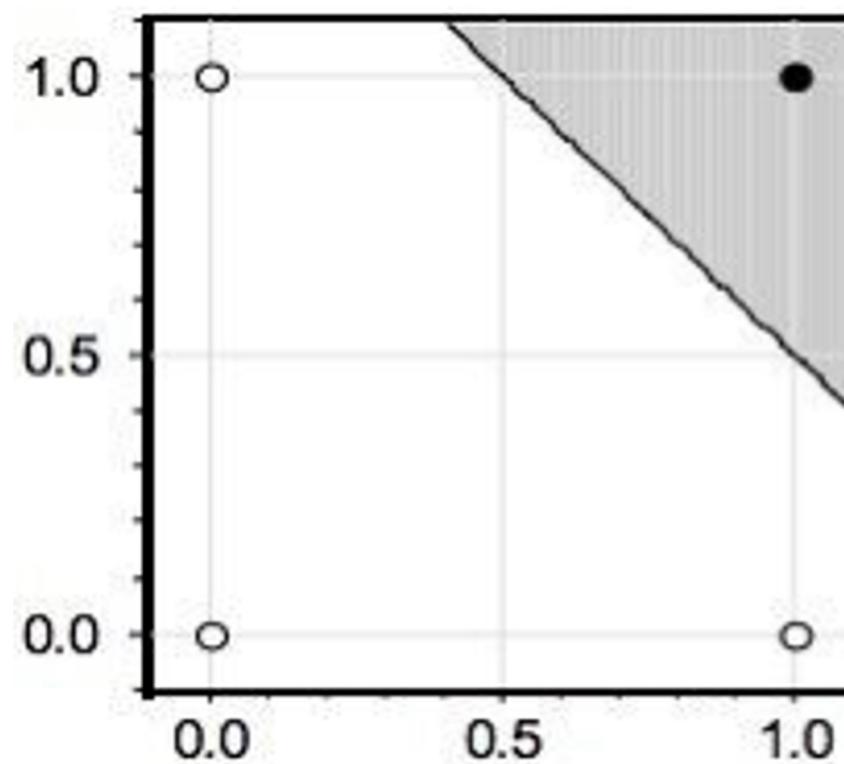
Deep Learning Timeline



XOR problem

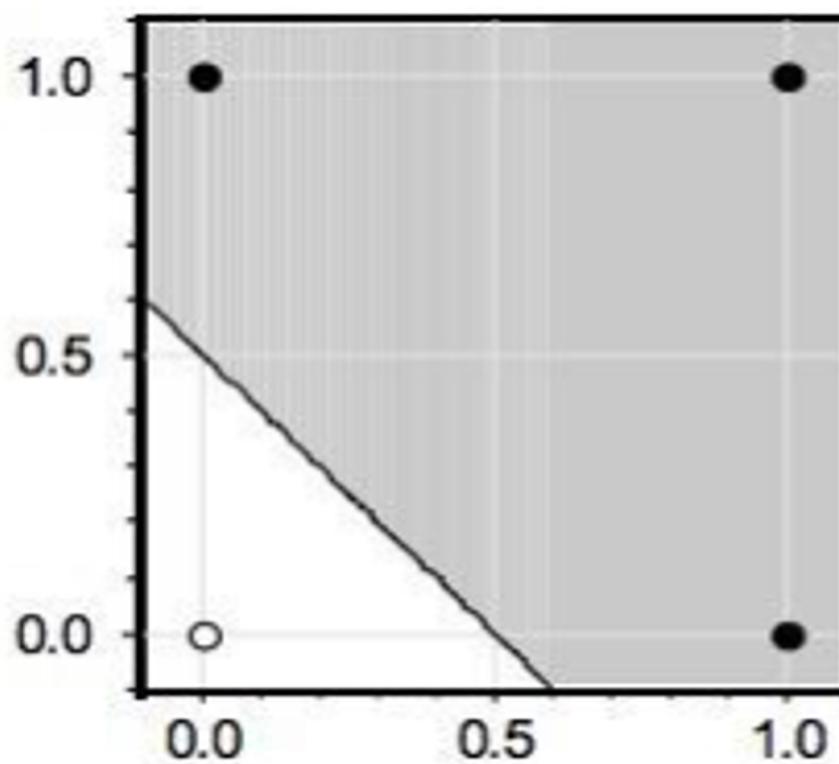


- This 2-layer NN (on the left) implements XOR with only x_1 and x_2 features.
- 1-layer NN also can succeed, but only with extra feature $x_1 \cdot x_2$

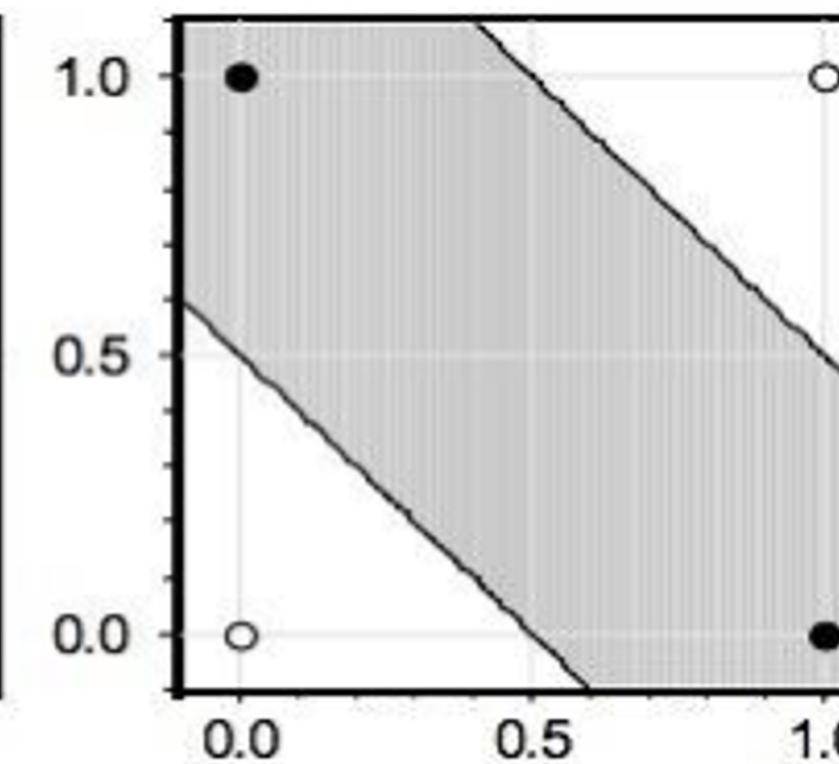
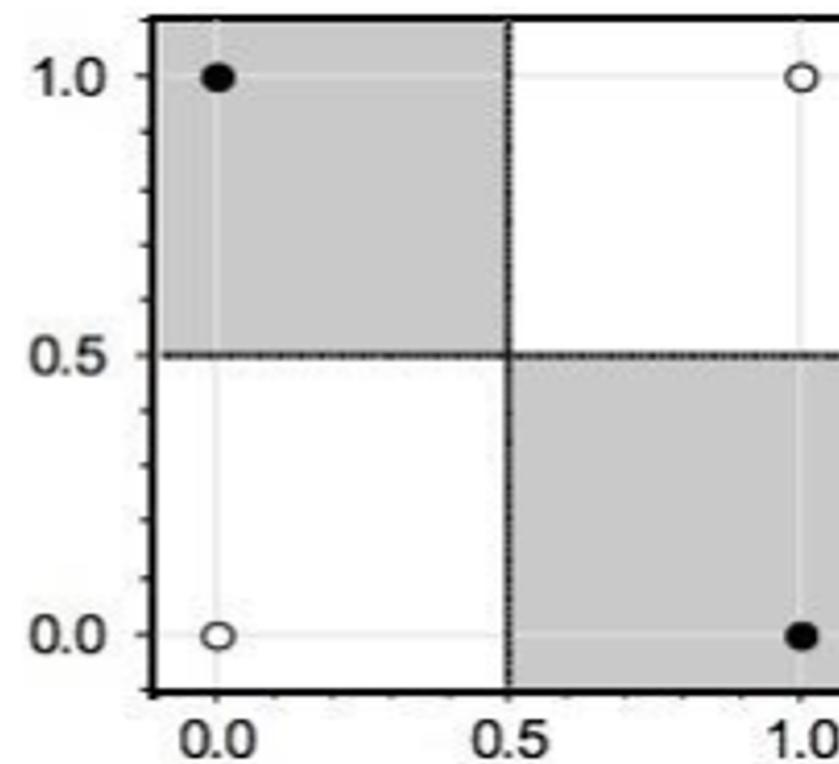


AND

XOR(with $x_1 \cdot x_2$)

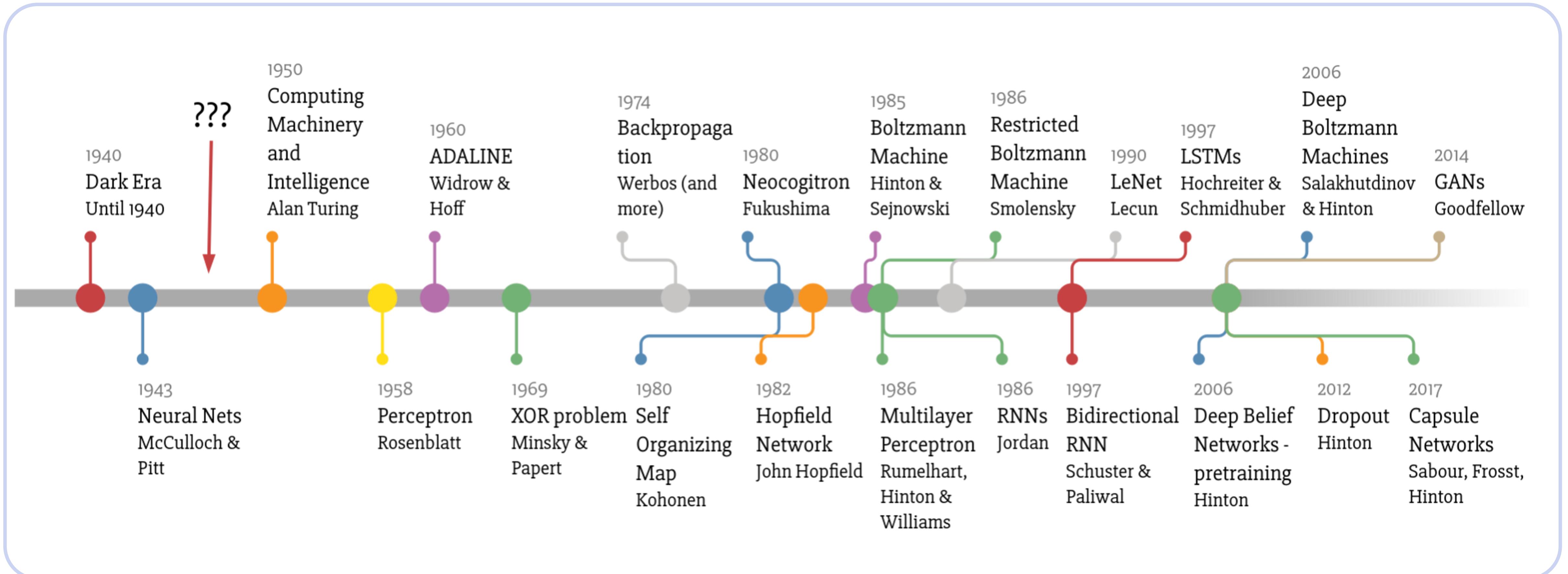


XOR

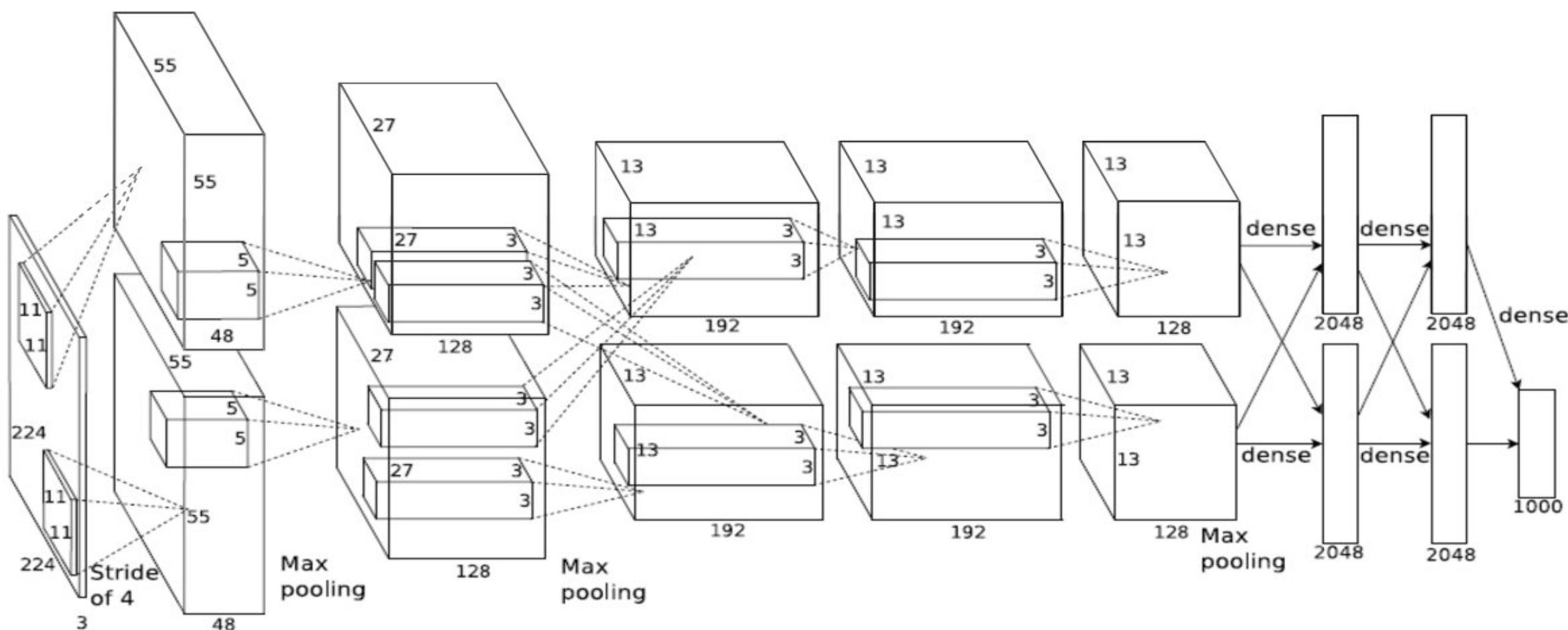


OR

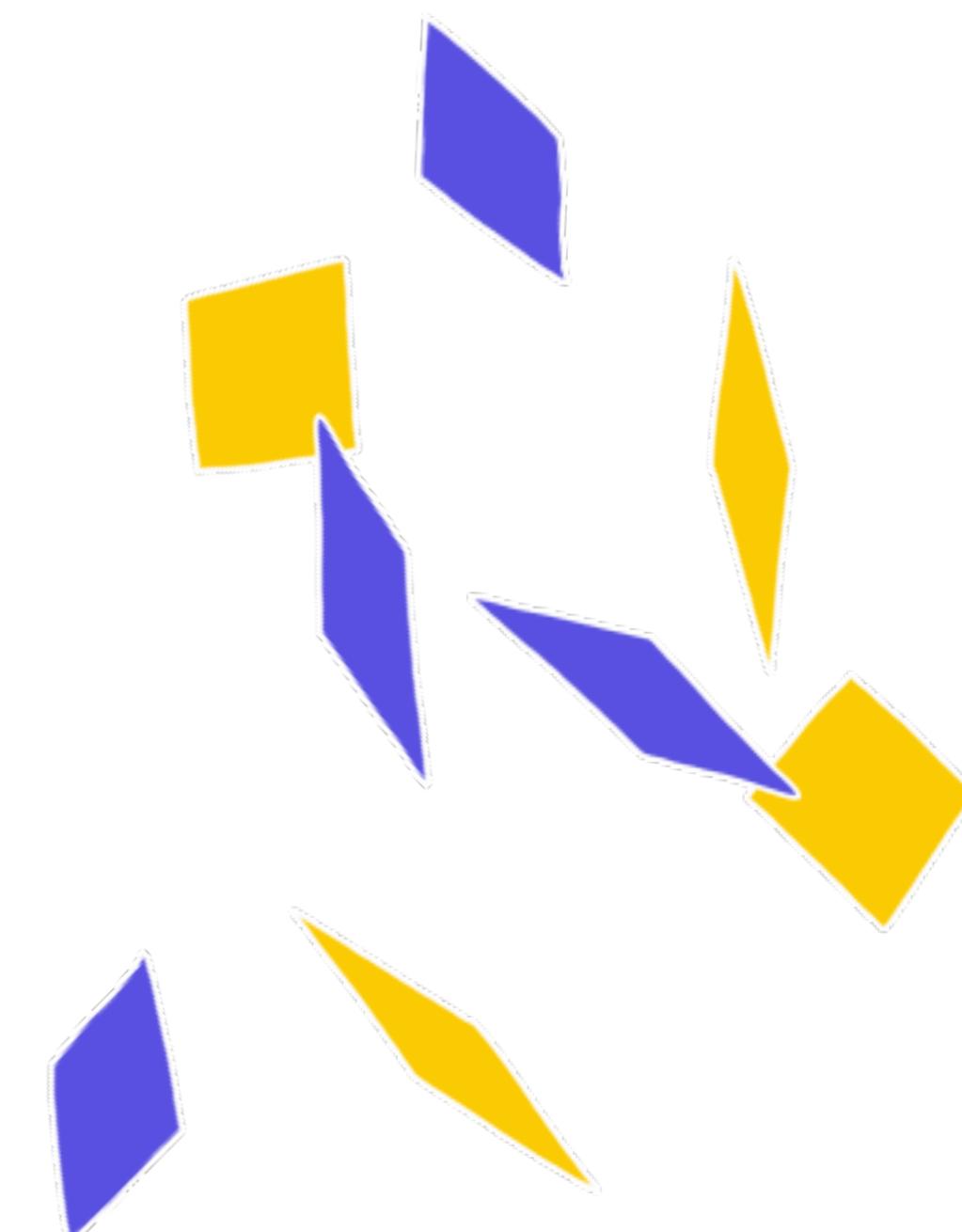
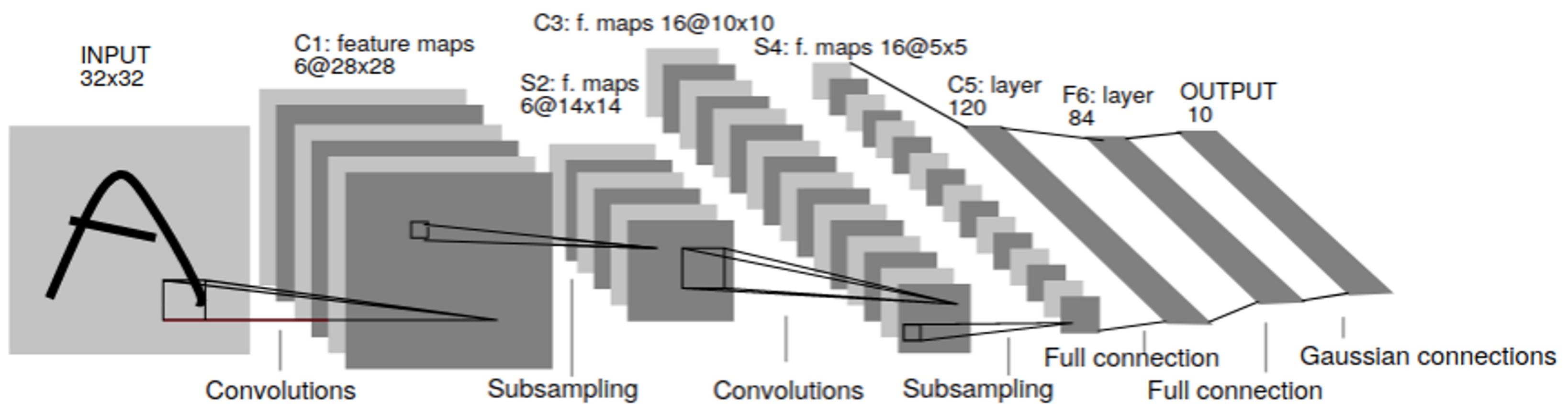
Deep Learning Timeline



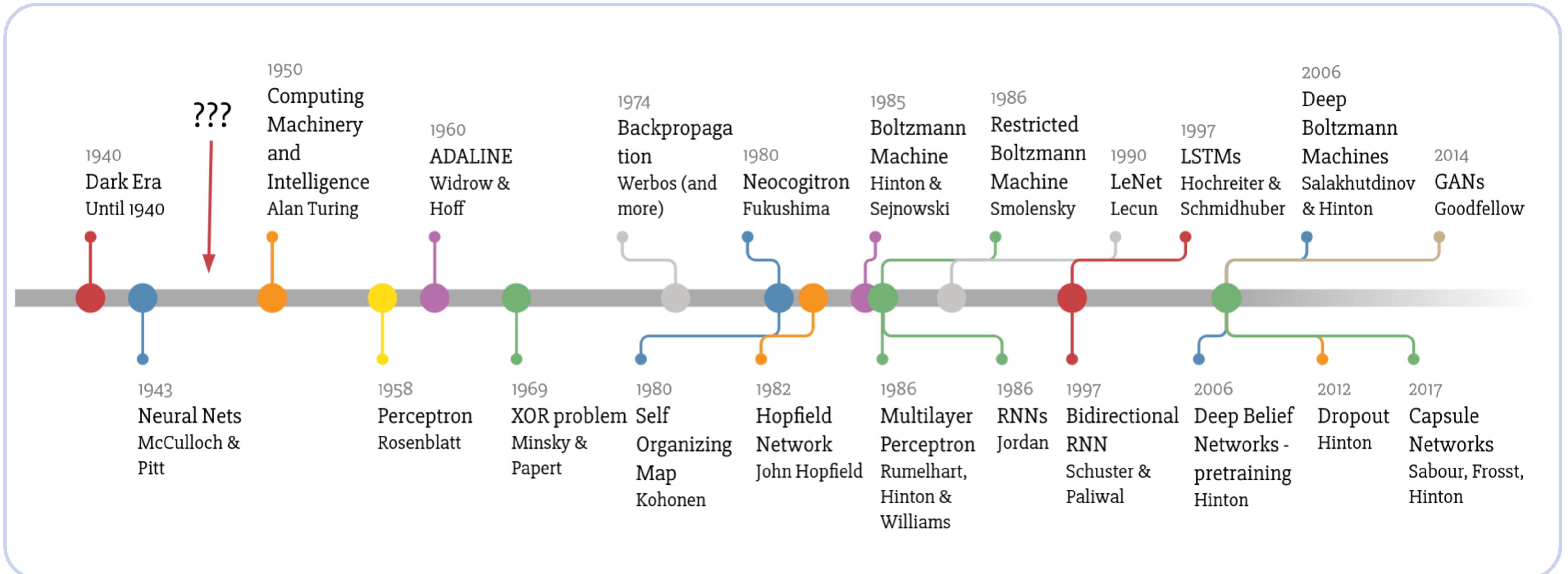
AlexNet, 2012



LeNet-5, 1998

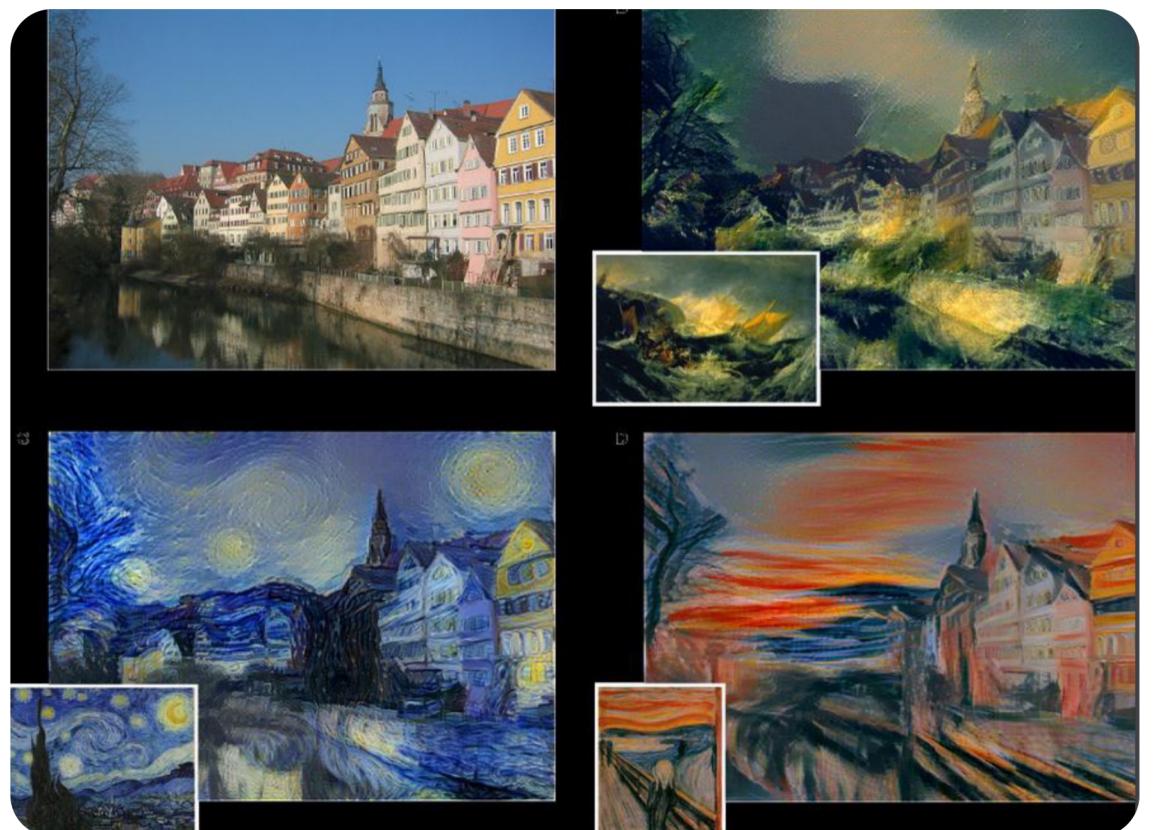
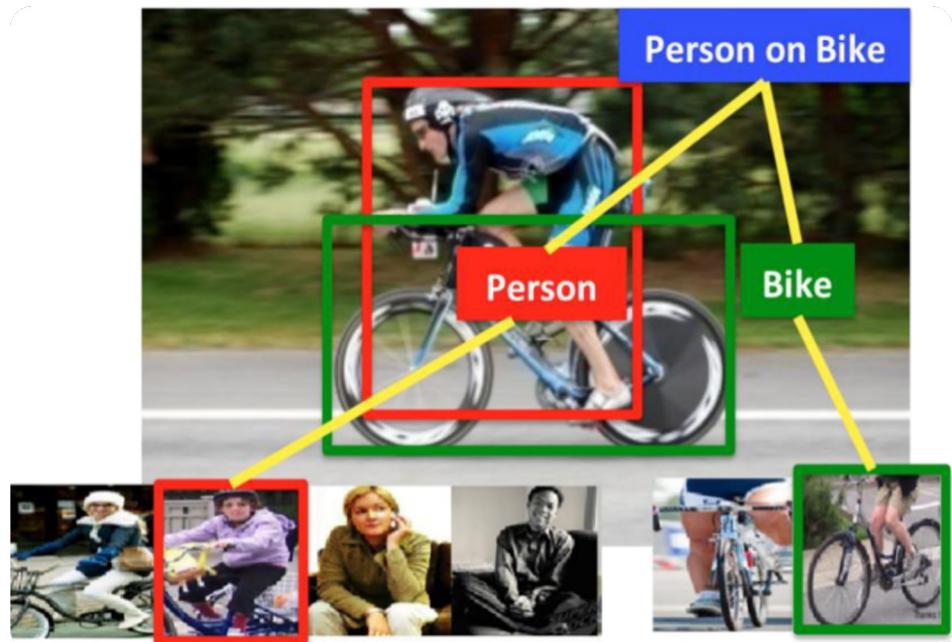
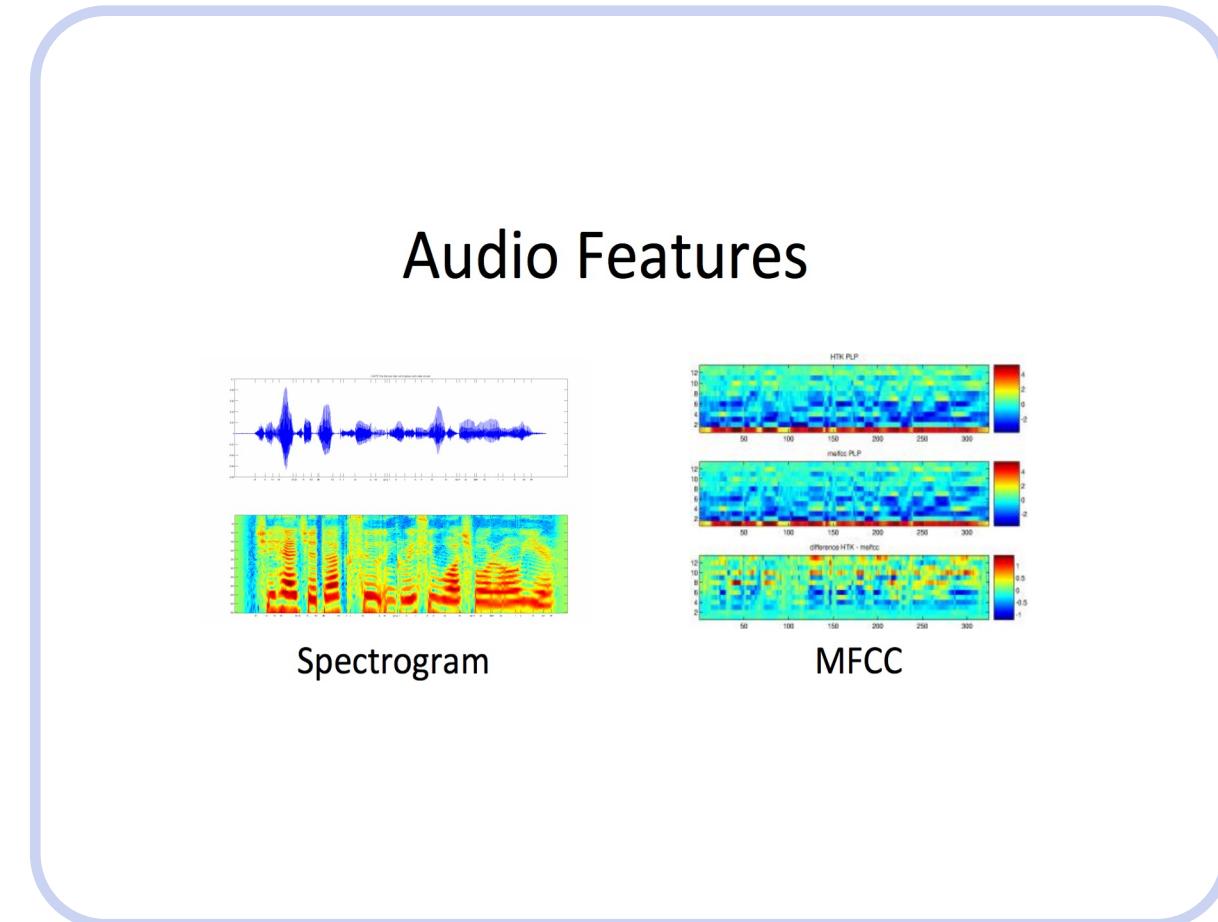


Deep Learning Timeline



Real world applications

- 01 Object detection
- 02 Action classification
- 03 Image captioning
- 04 ...



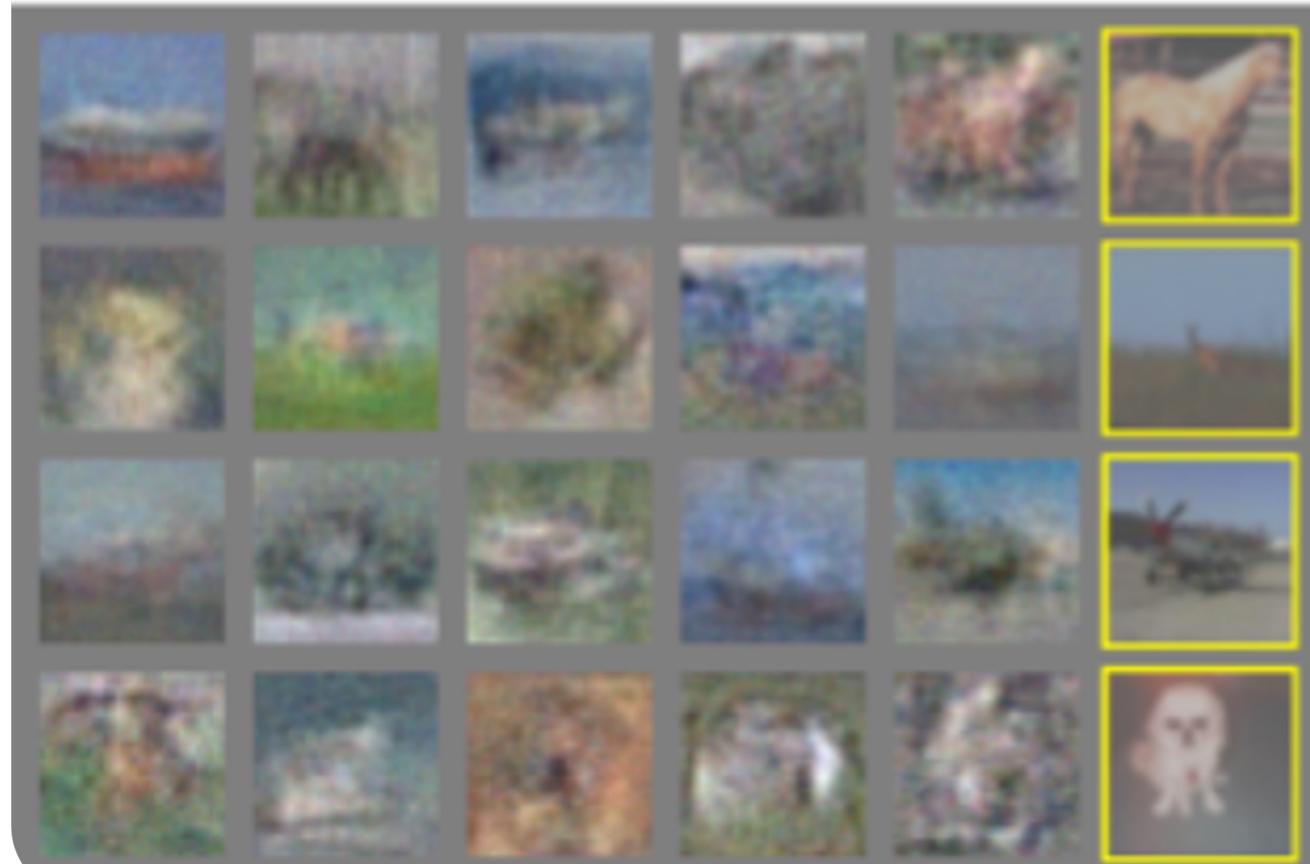
GANs, 2014+



a)



b)



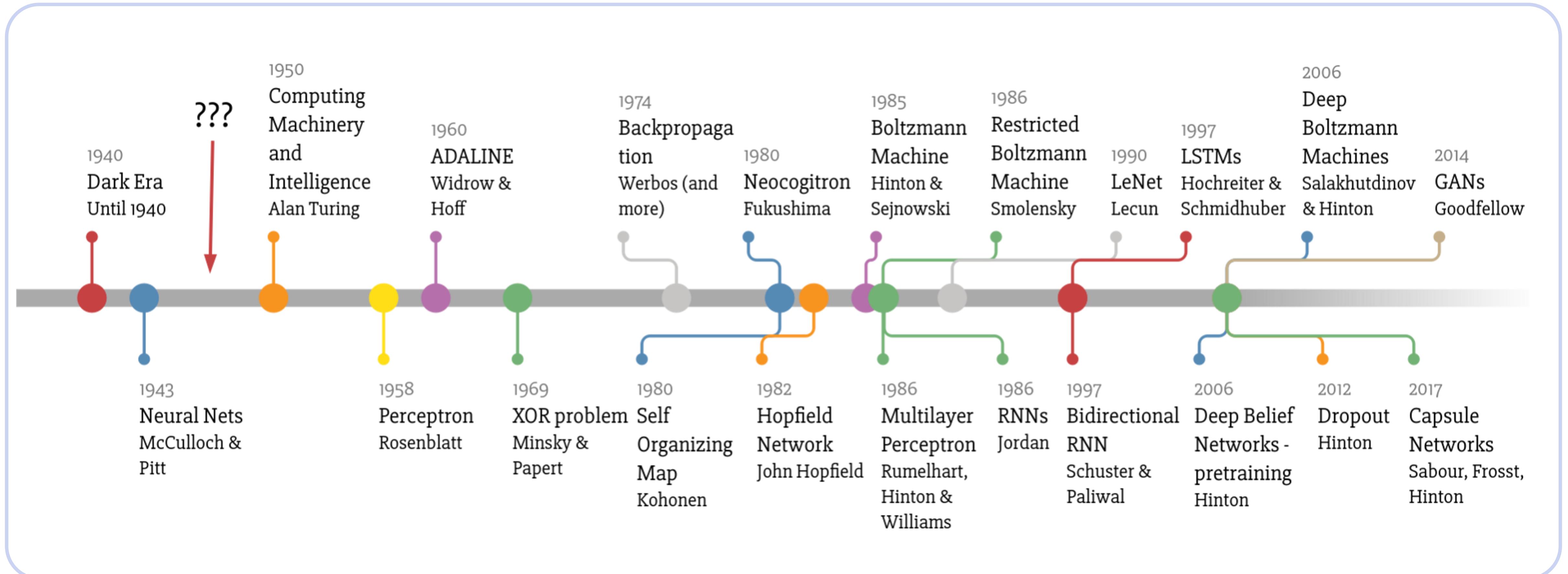
c)



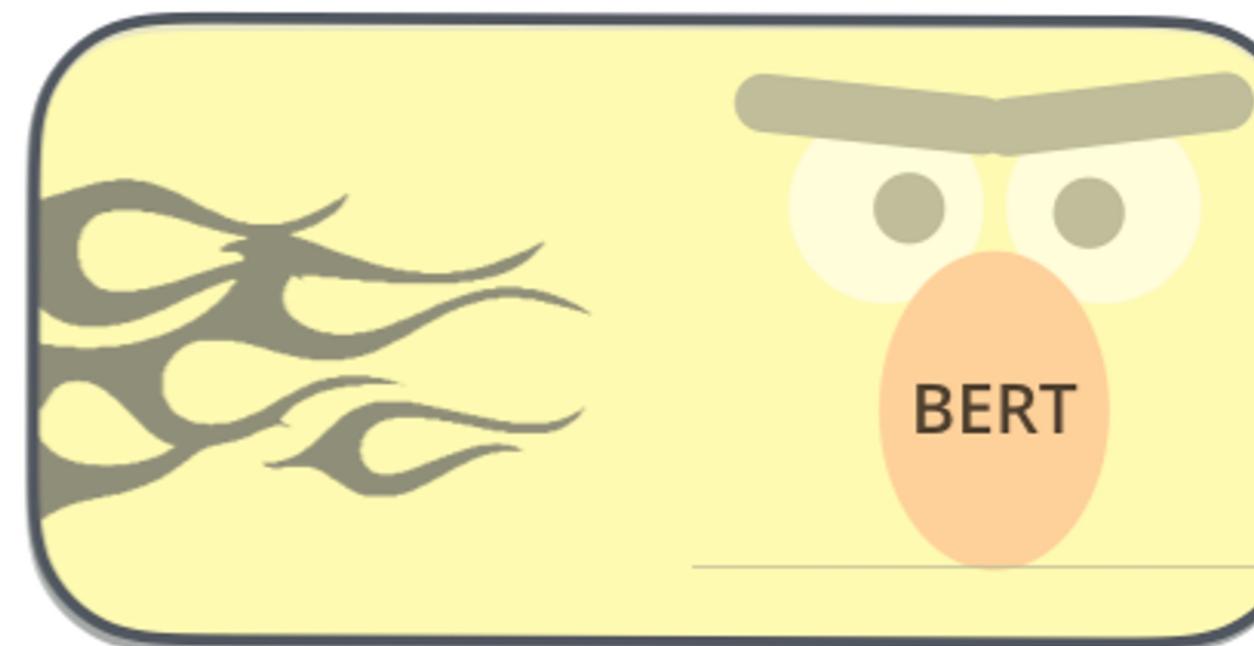
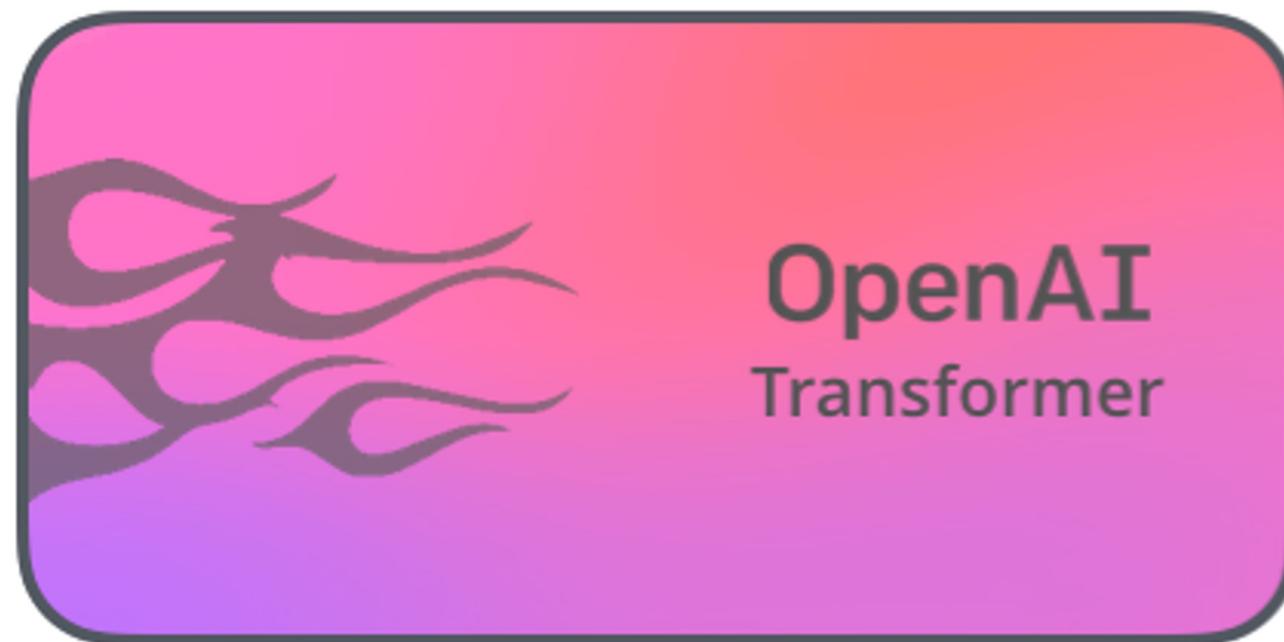
d)



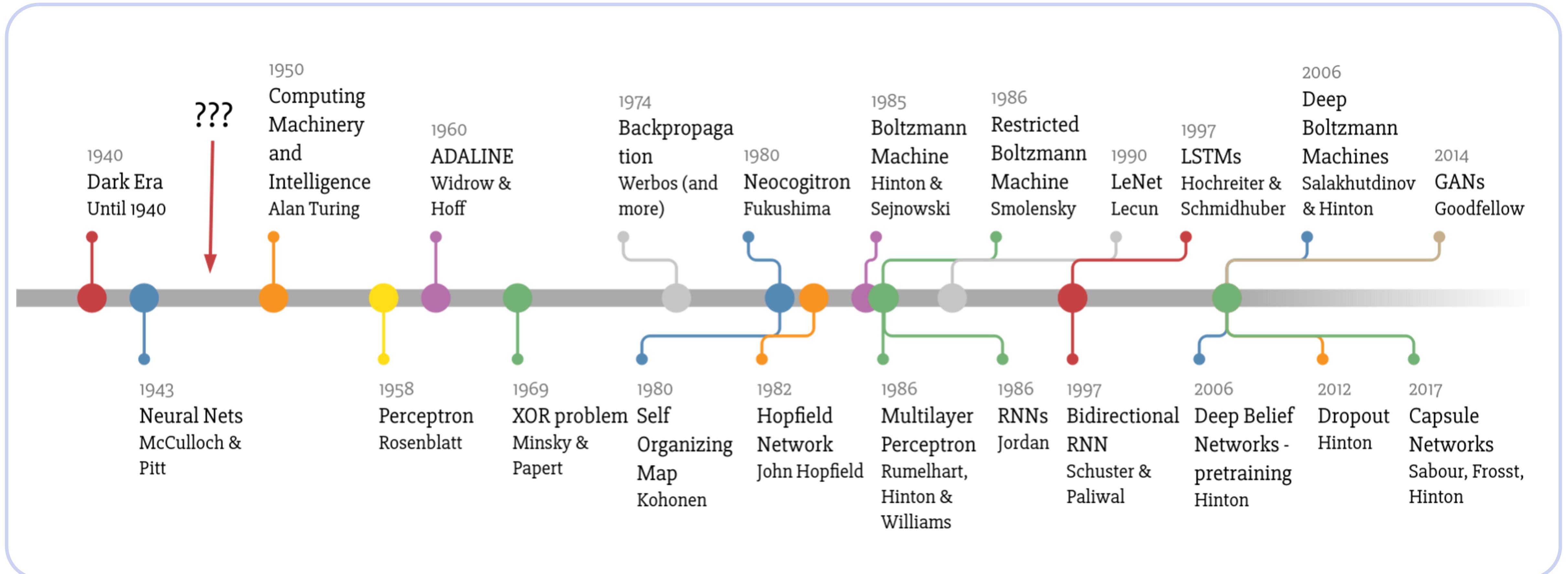
Deep Learning Timeline



Transformer, BERT, GPT-2 and more, 2017+



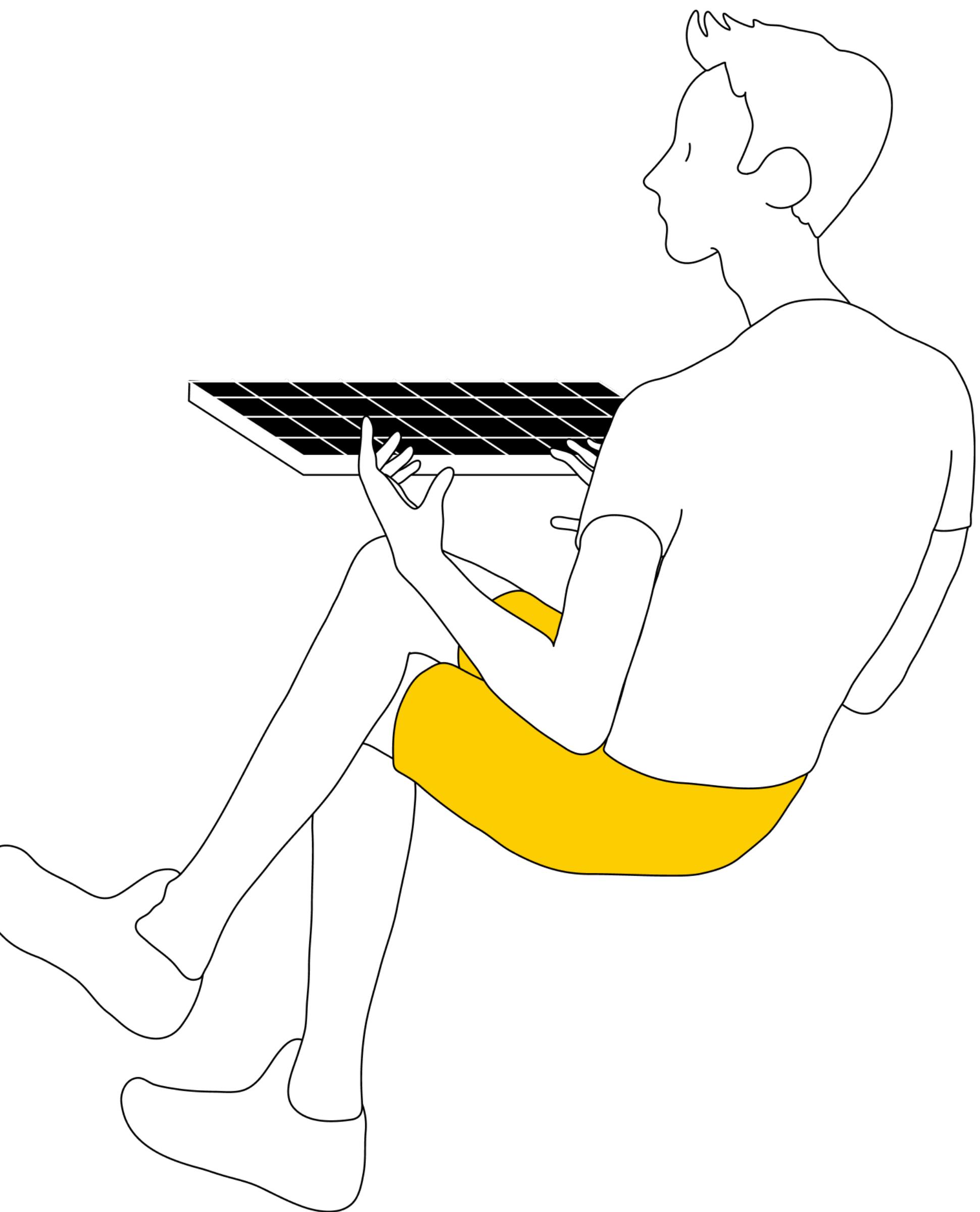
Deep Learning Timeline



2022: (Instruct)ChatGPT, StableDiffusion
2023: GPT boom: GPT-4, YaGPT, GigaChat,
LLaMa

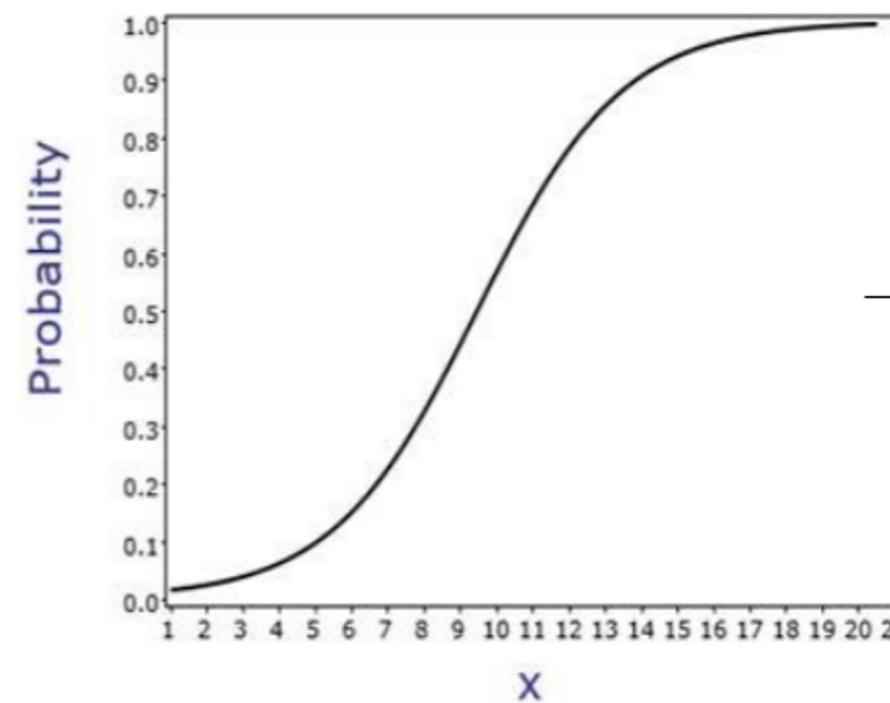


Deep Learning: intuition



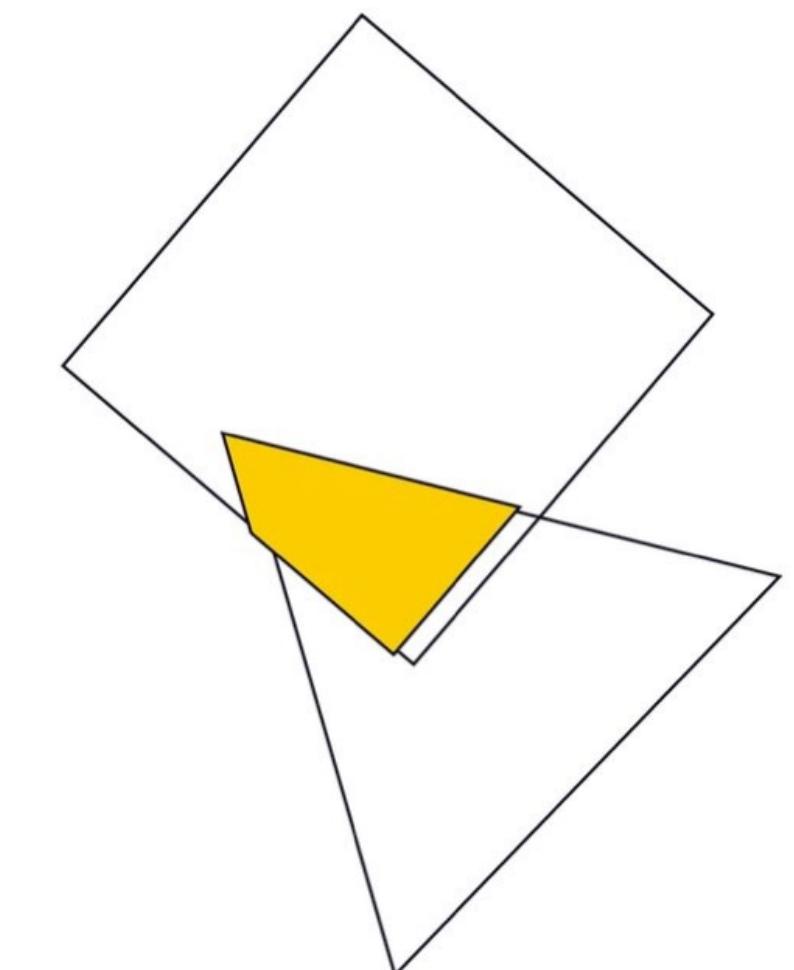
Logistic regression

$$x \rightarrow Wx + b \rightarrow P(y)$$



$$P(y = 1|x) = \sigma(w^\top x + b)$$

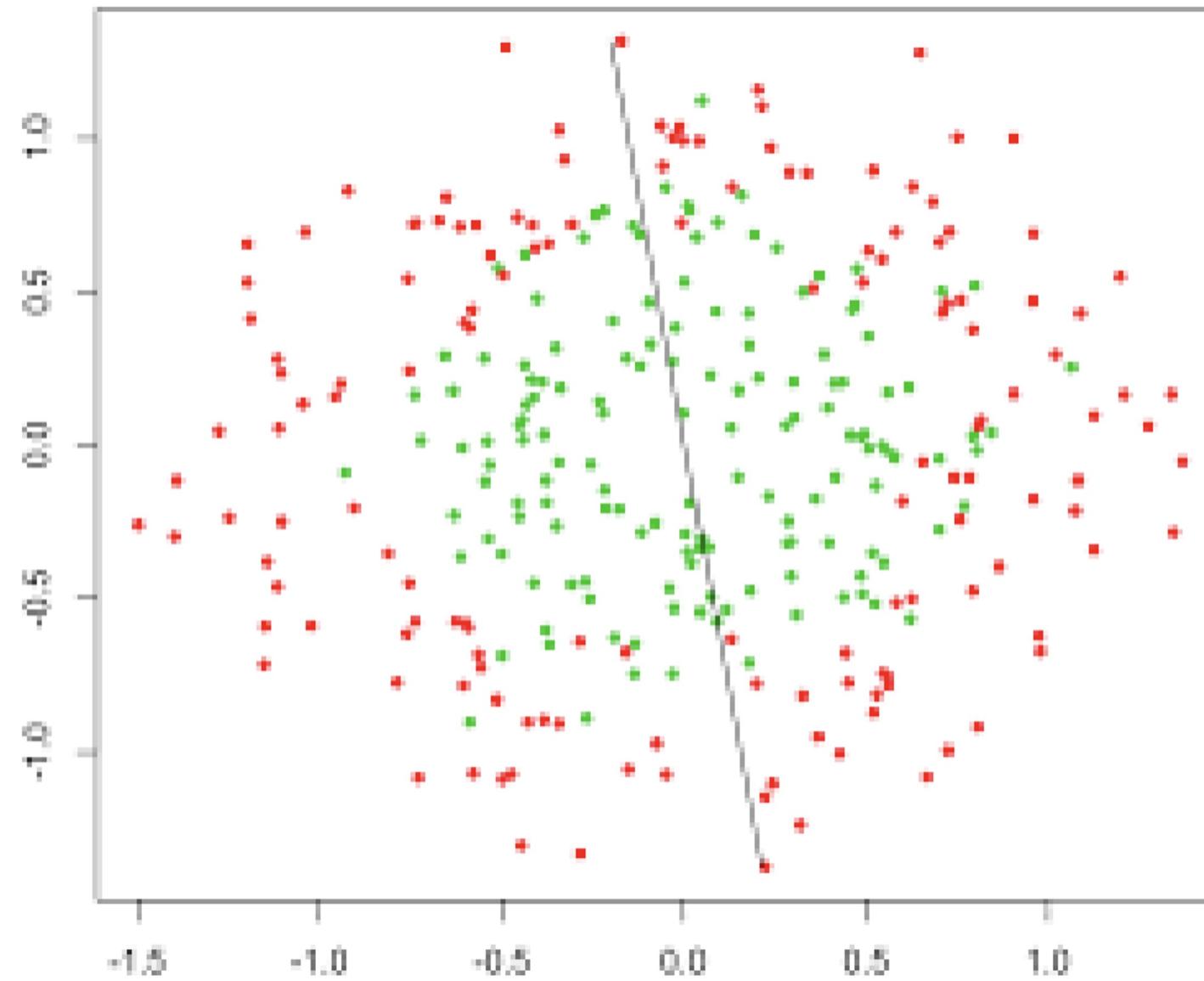
$$L = - \sum_i \left(y^{(i)} \log P(y = 1|x^{(i)}) + (1 - y^{(i)}) \log (1 - P(y = 1|x^{(i)})) \right)$$



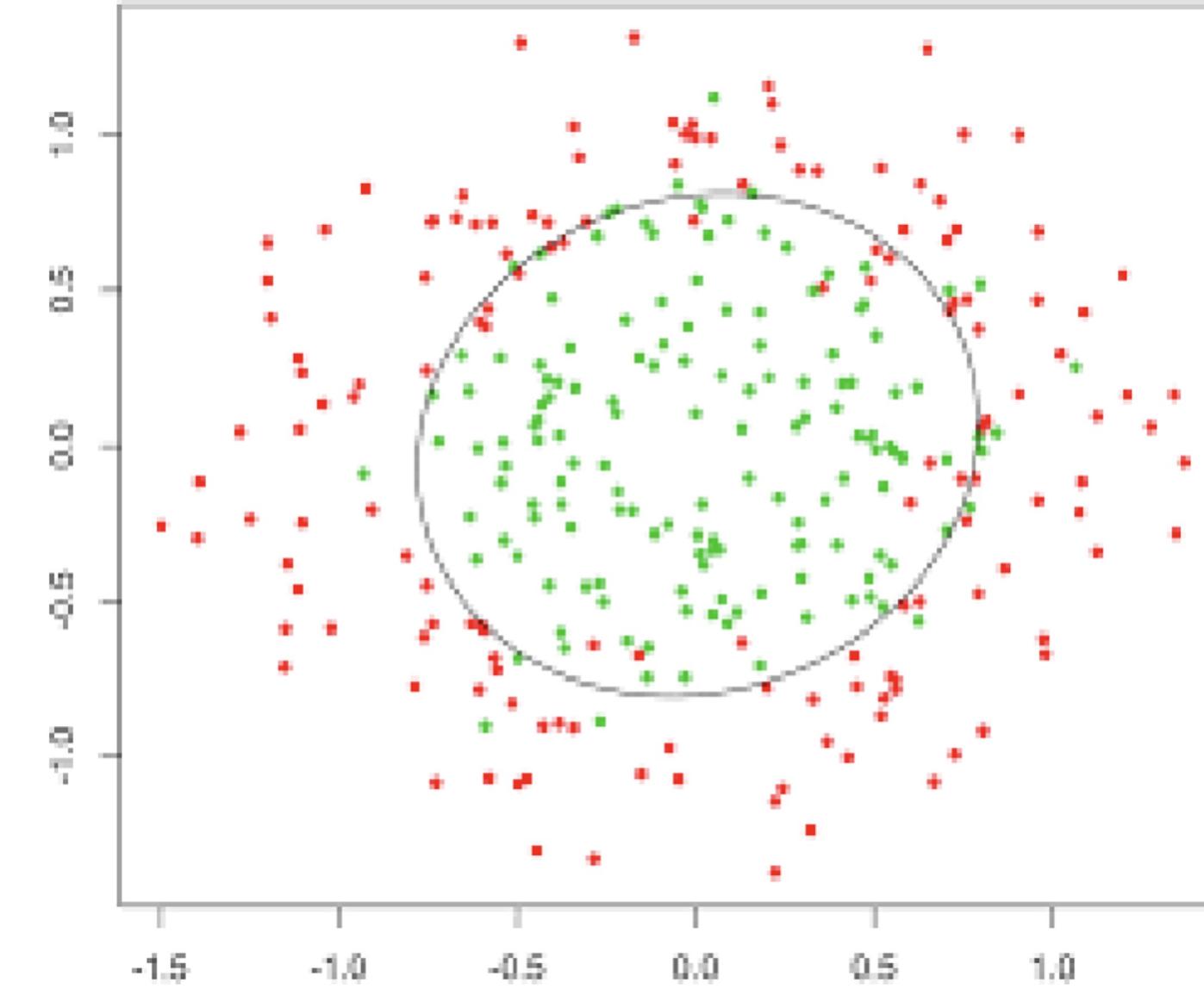
Problem: nonlinear dependencies

Logistic regression (generally, linear model) need feature engineering to show good results

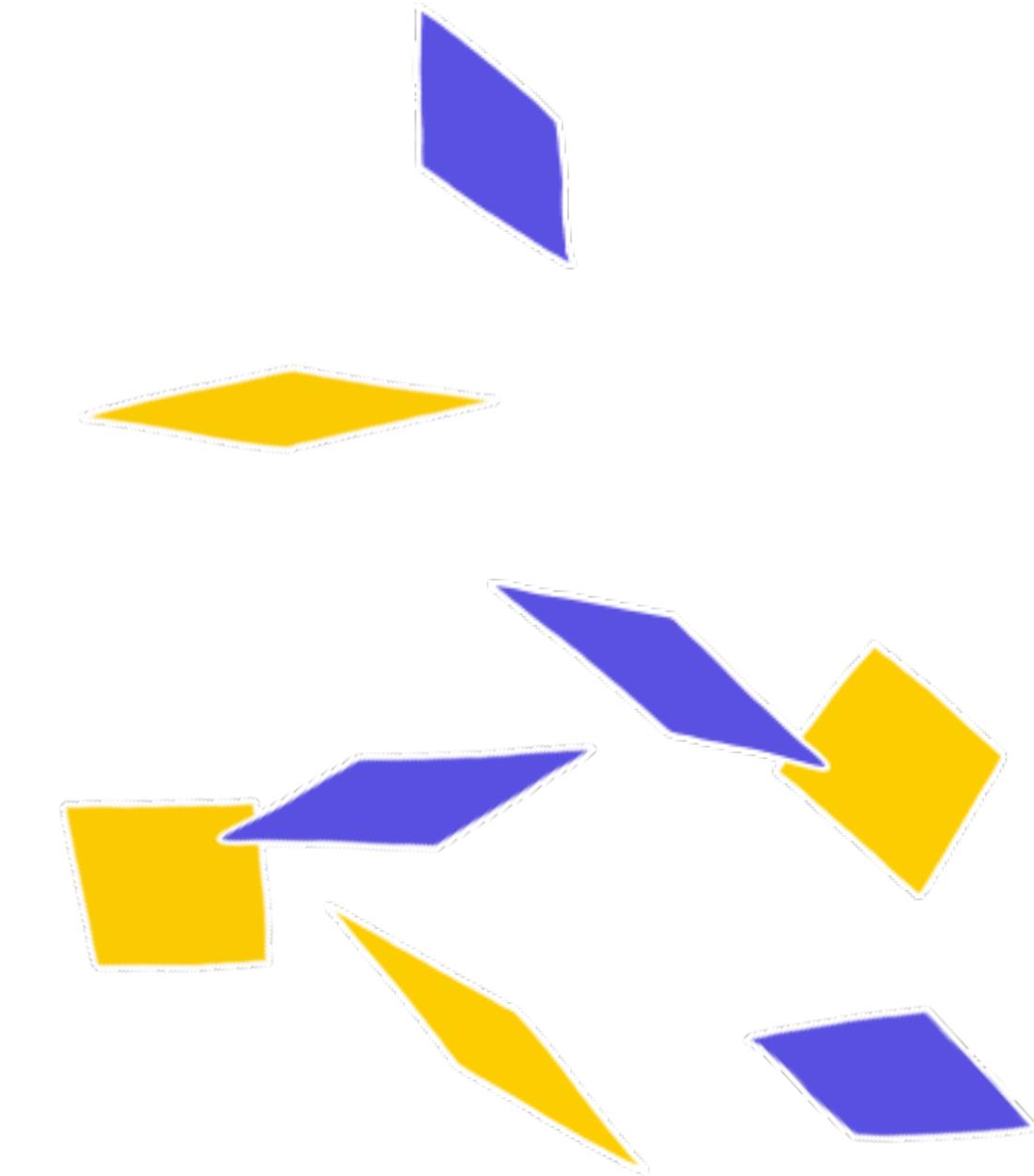
What we have



What we want



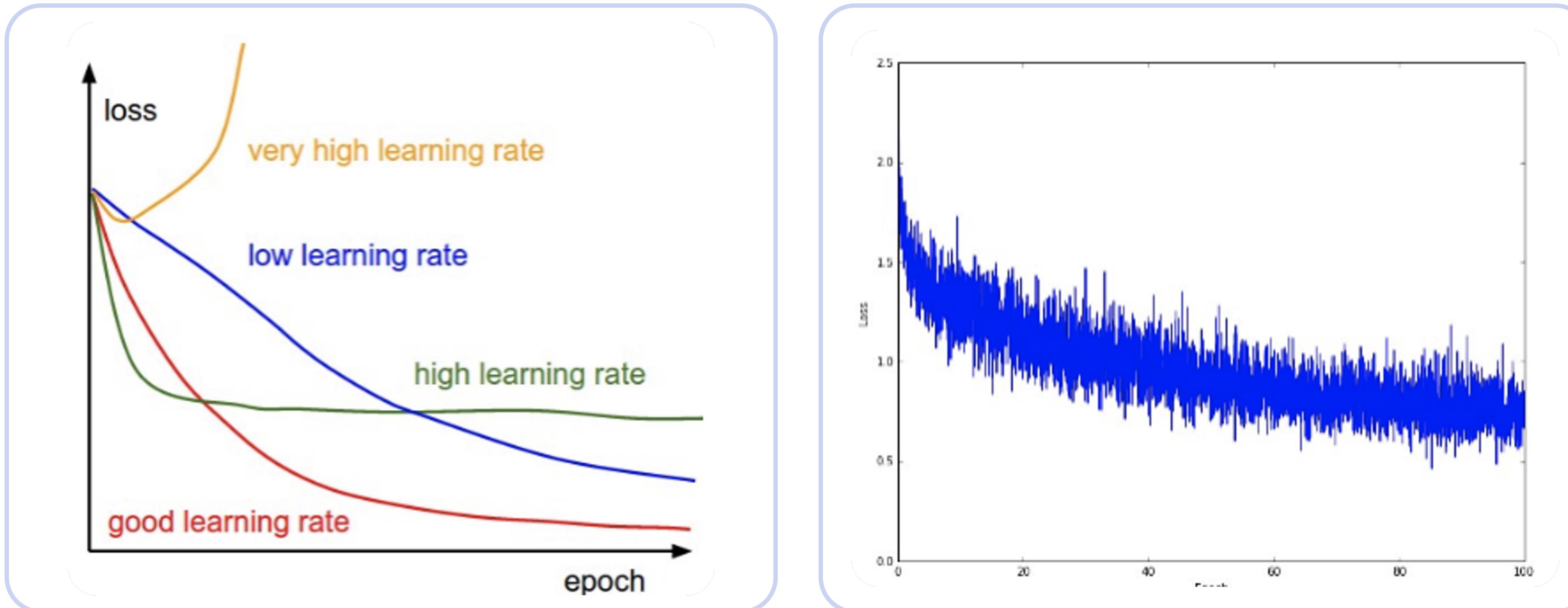
And feature engineering is an art



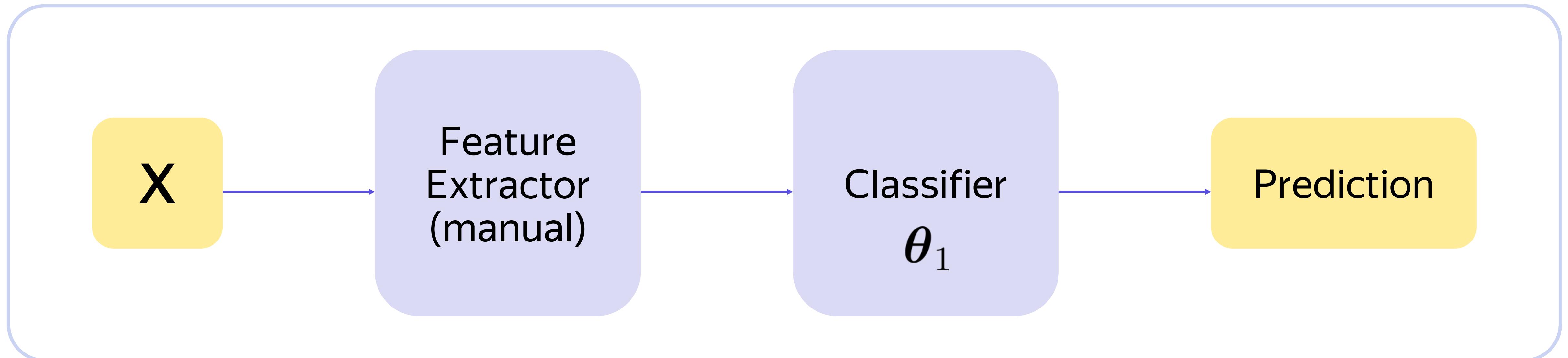
Gradient optimization

Stochastic gradient descent is used to optimize NN parameters

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \frac{dL}{d\mathbf{w}}$$

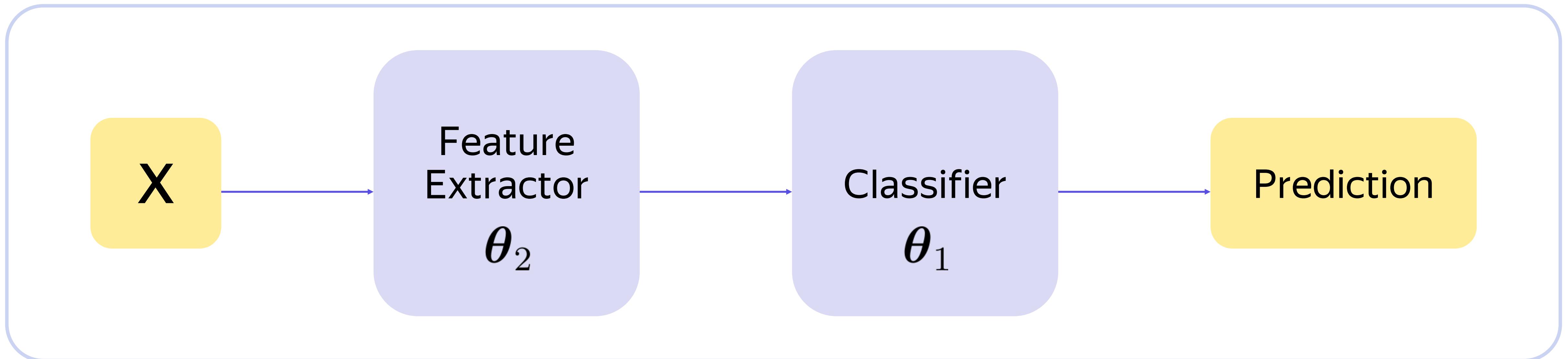


Classic pipeline



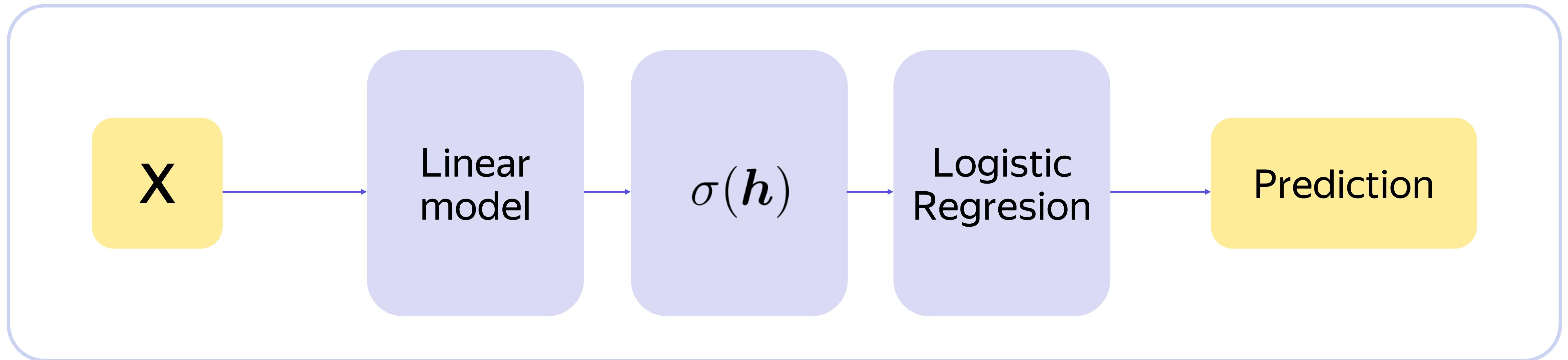
Handcrafted features, generated by experts

NN pipeline



Automatically extracted features

NN pipeline: example



E.g. two logistic regressions one after another.

Actually, it's a neural network

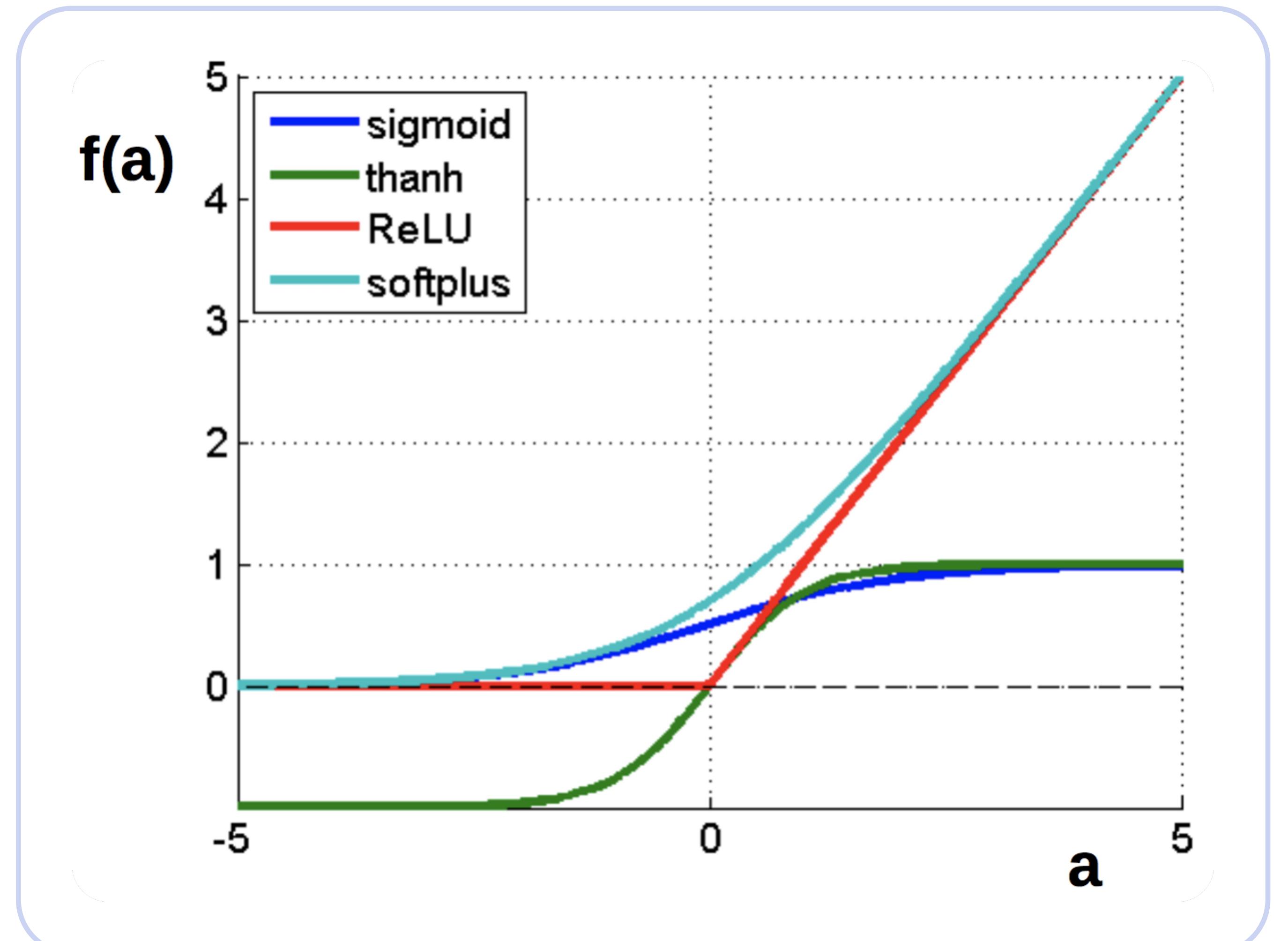
Activation functions: nonlinearities

$$f(a) = \frac{1}{1 + e^{-a}}$$

$$f(a) = \tanh(a)$$

$$f(a) = \max(0, a)$$

$$f(a) = \log(1 + e^a)$$



Some generally accepted terms

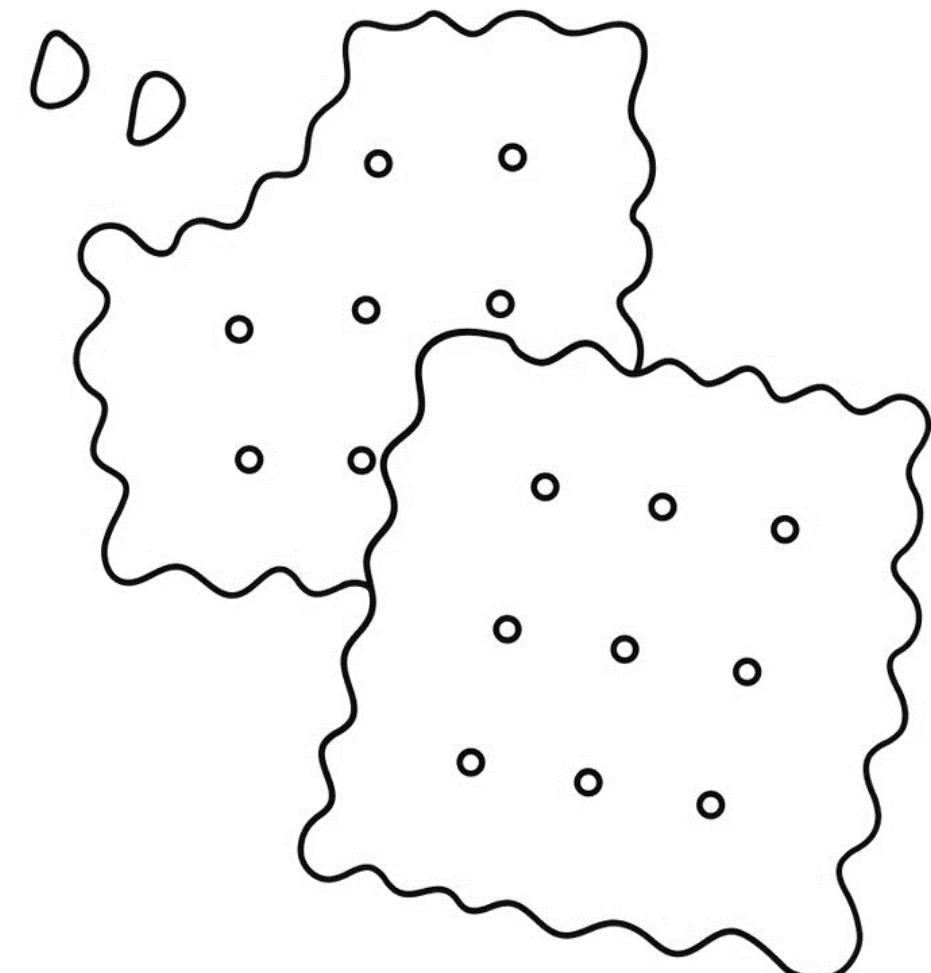
01 Layer — a building block for NNs :

- Dense/Linear/FC layer: $\mathbf{W}\mathbf{x} + \mathbf{b}$
- Nonlinearity layer, e.g. $\sigma(\mathbf{h})$
- Input layer, output layer
- etc.

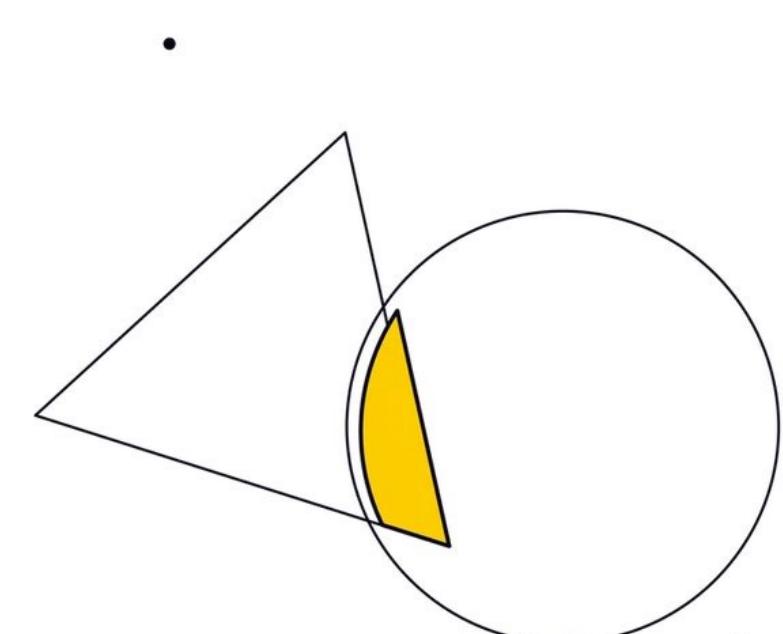
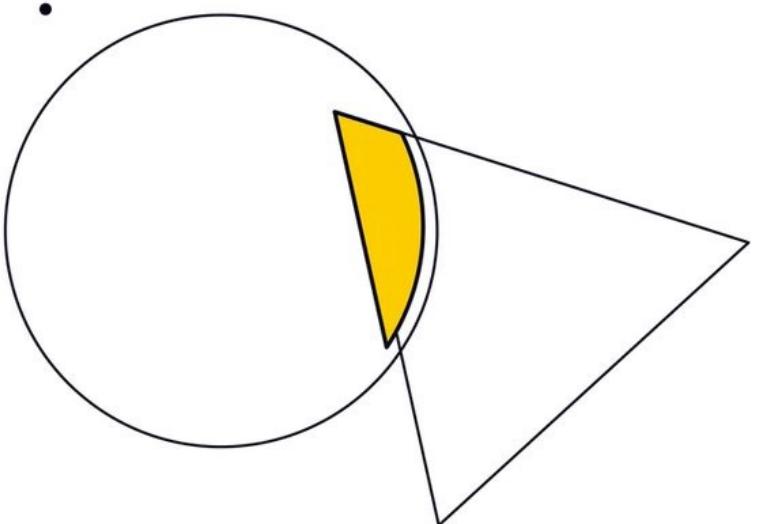
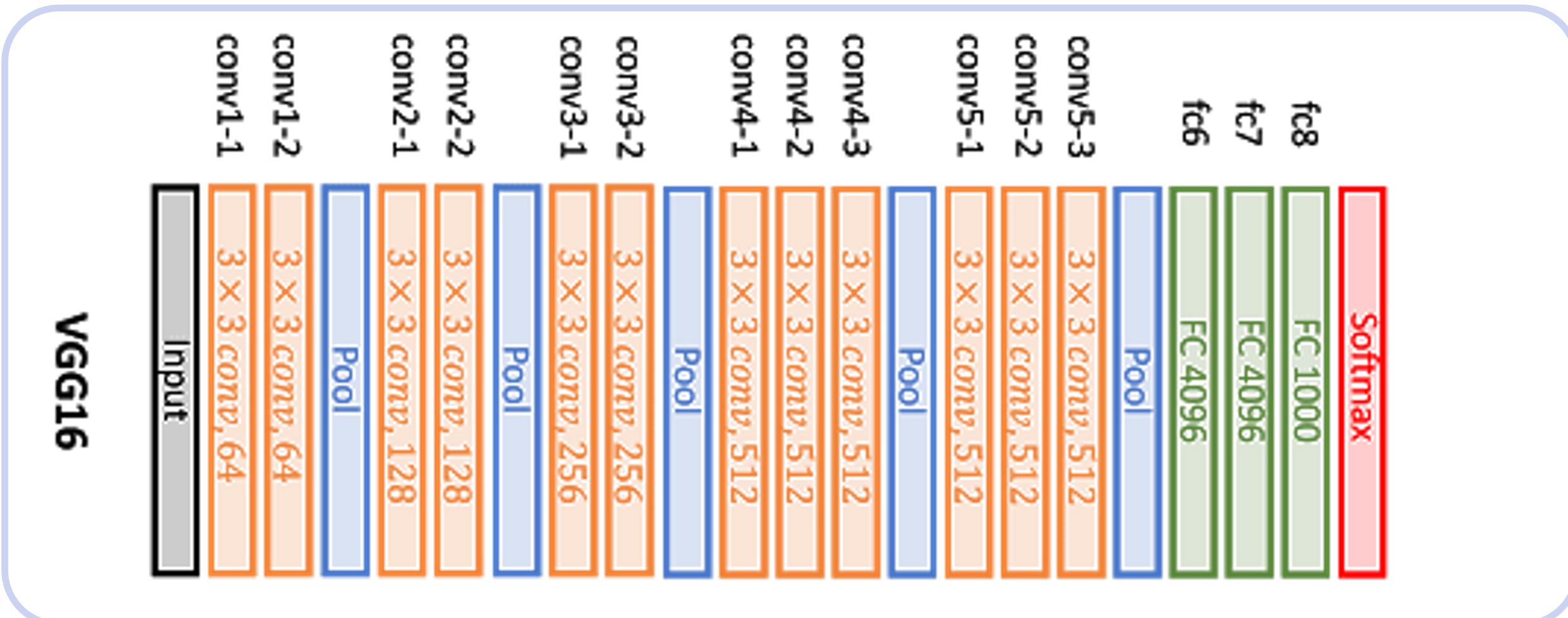
02 Activation function — function applied to layer output

- Sigmoid
- \tanh
- ReLU
- Any other function to get nonlinear intermediate signal in NN

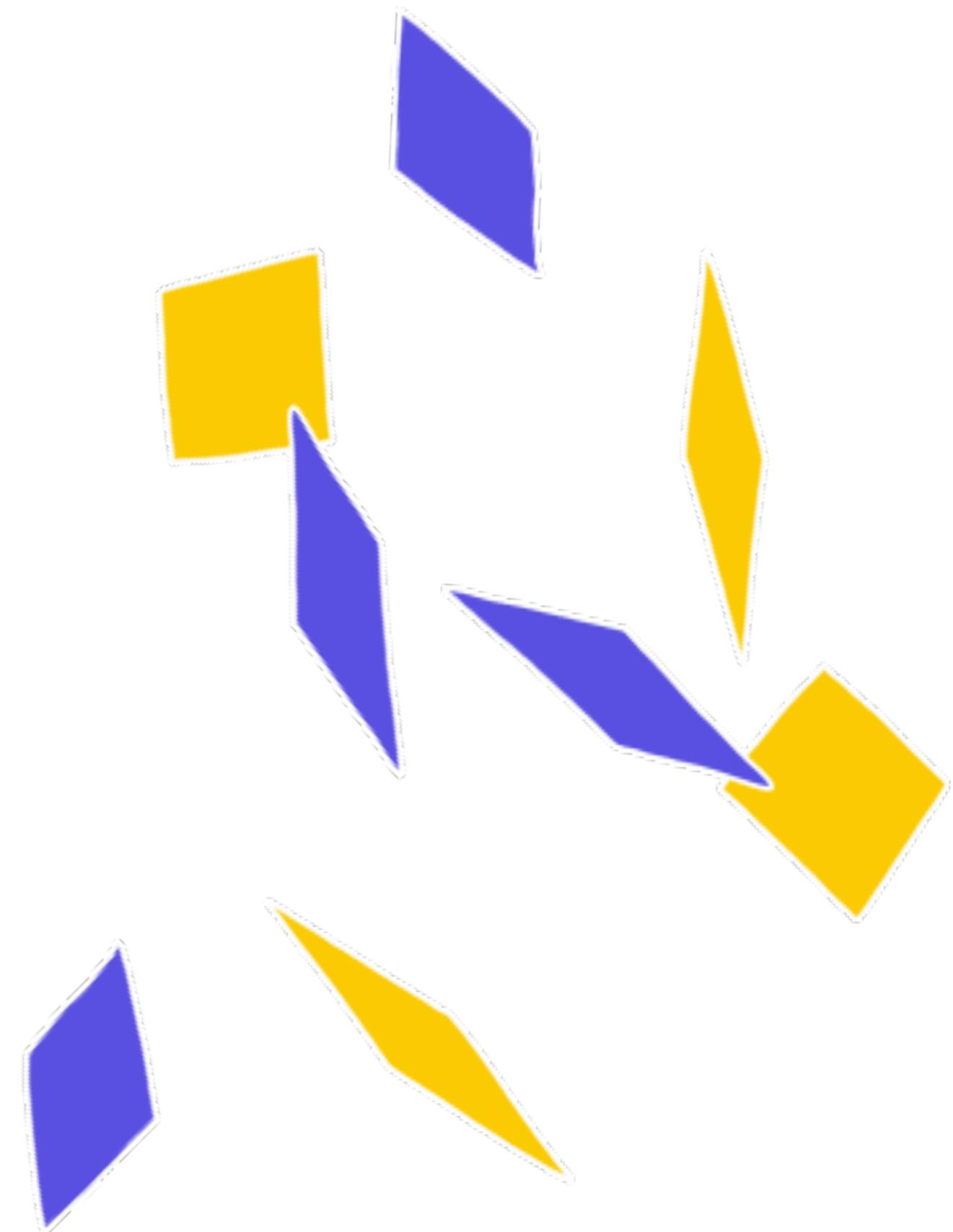
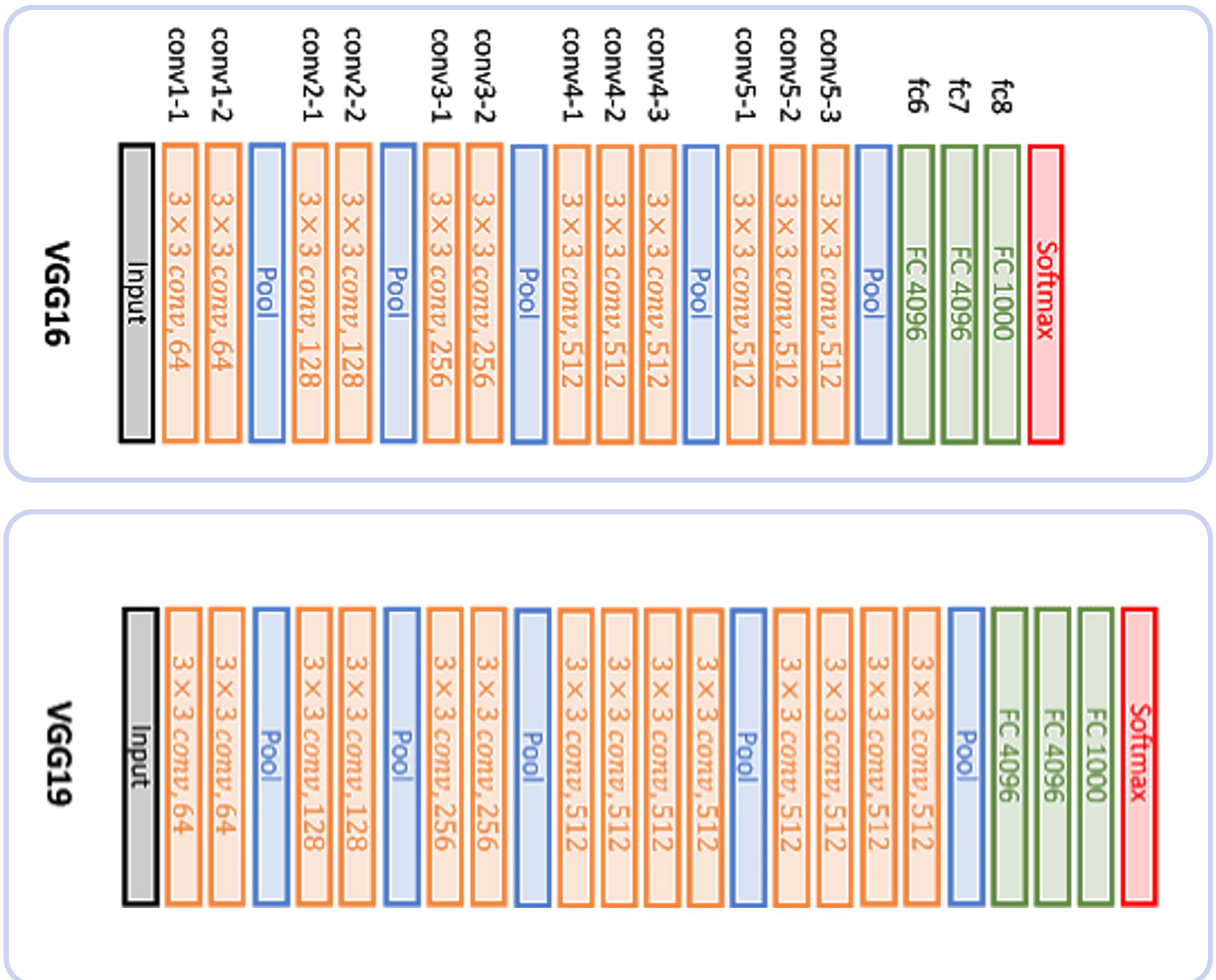
03 Backpropagation — a fancy word for “chain rule”



Actually, networks can be deep

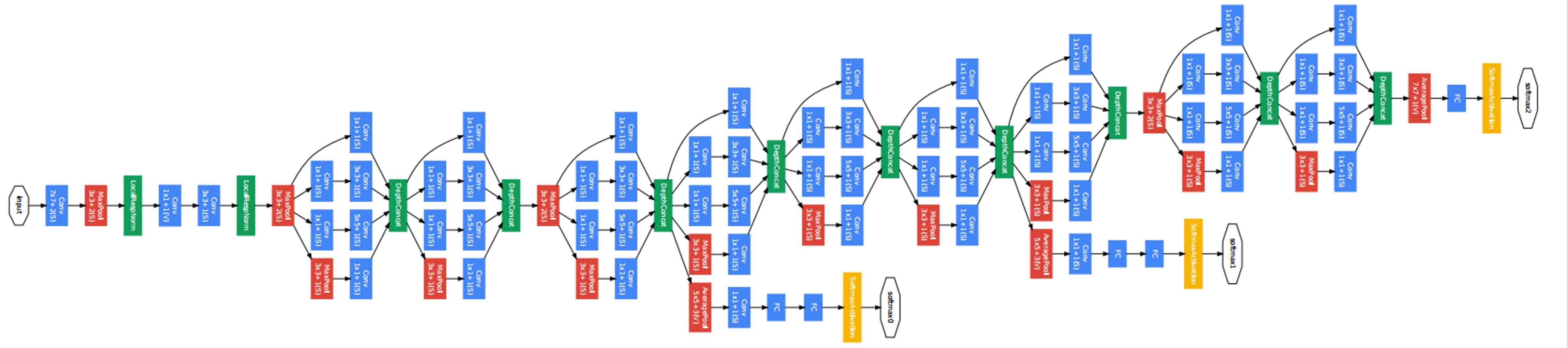


And deeper...



Much deeper...

How to train it?

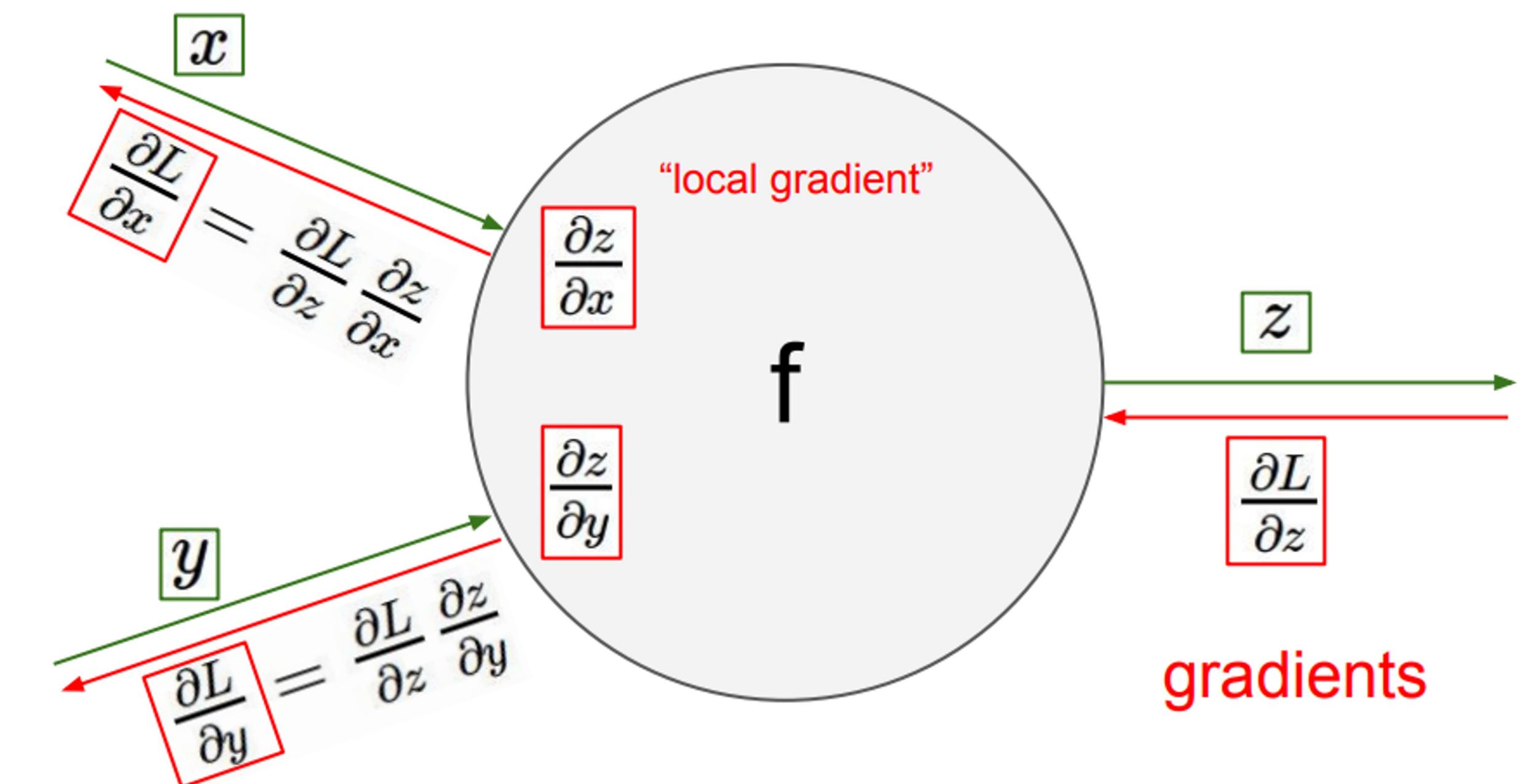


Backpropagation and chain rule

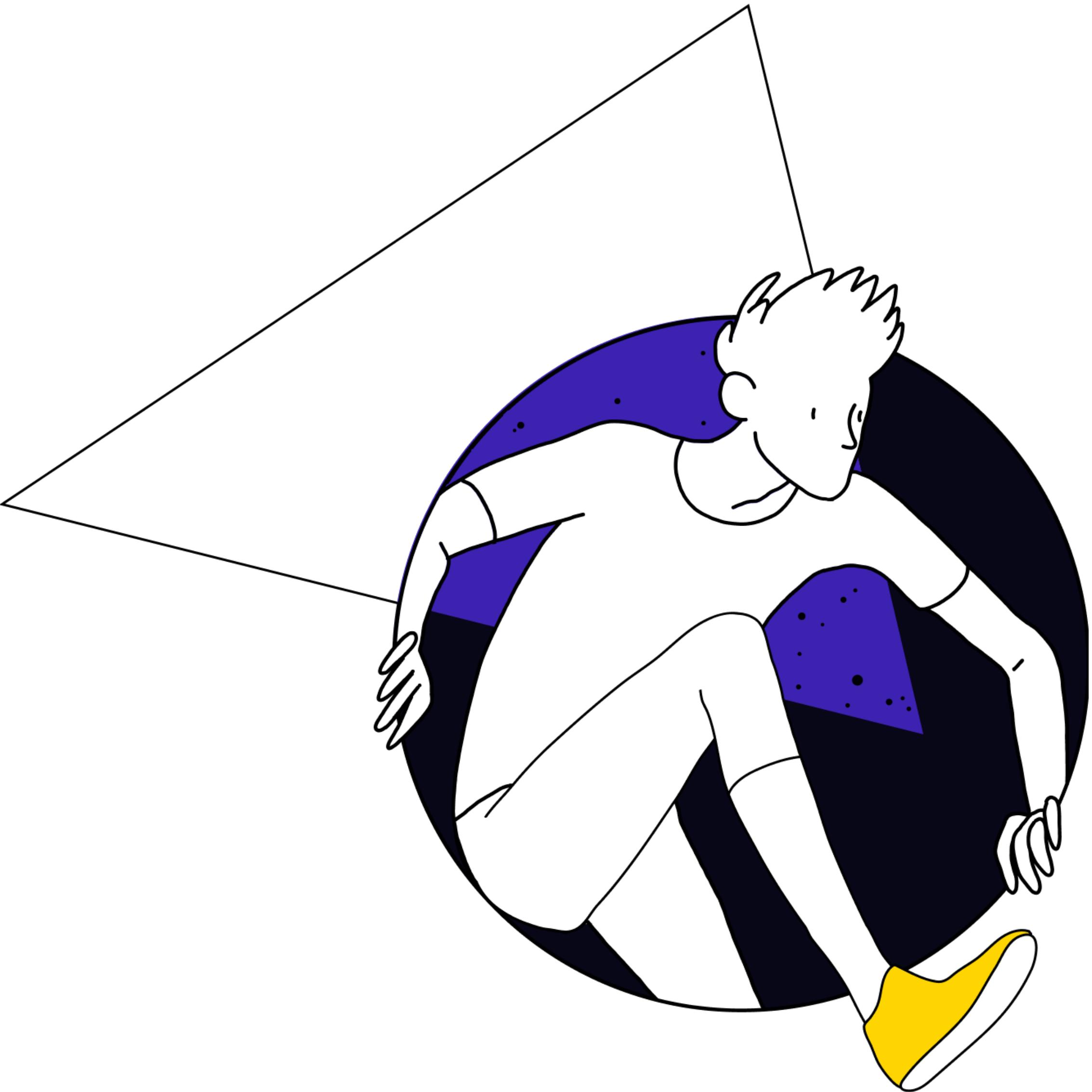
Chain rule is just simple math:

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x}$$

Backprop is just way
to use it in NN training

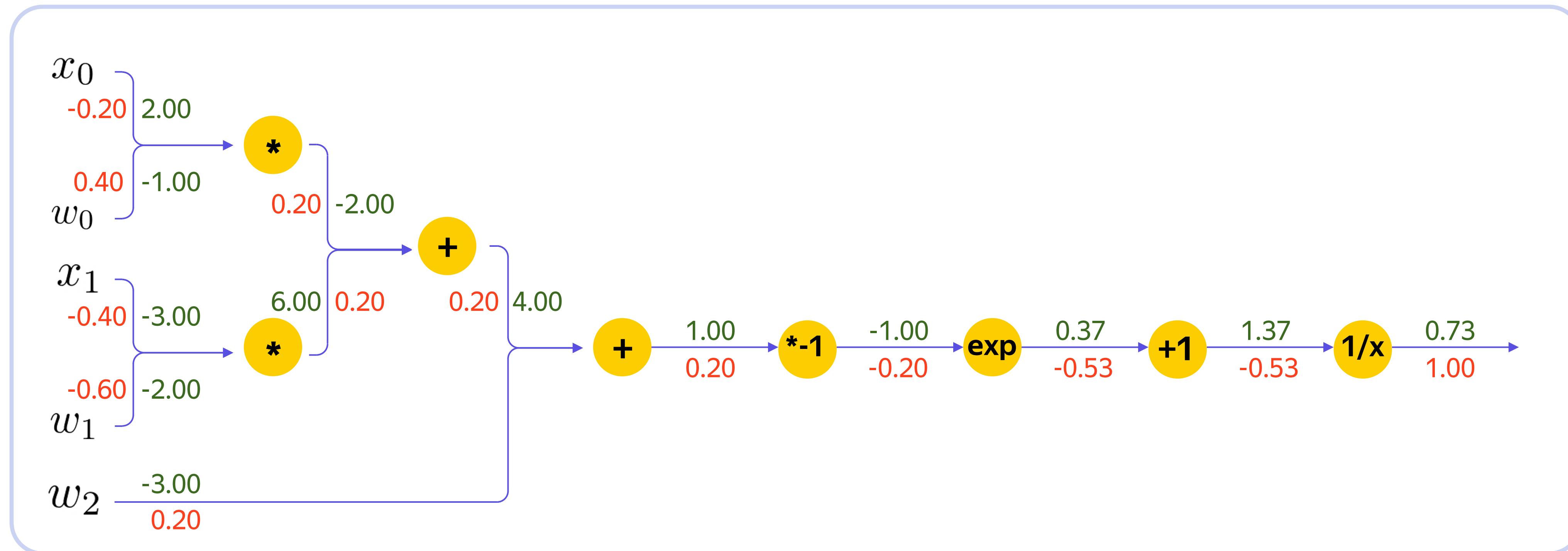


Backpropagation



Backpropagation example

$$L(\mathbf{w}, \mathbf{x}) = \frac{1}{1 + \exp(-(\mathbf{x} \cdot \mathbf{w} + b))}$$



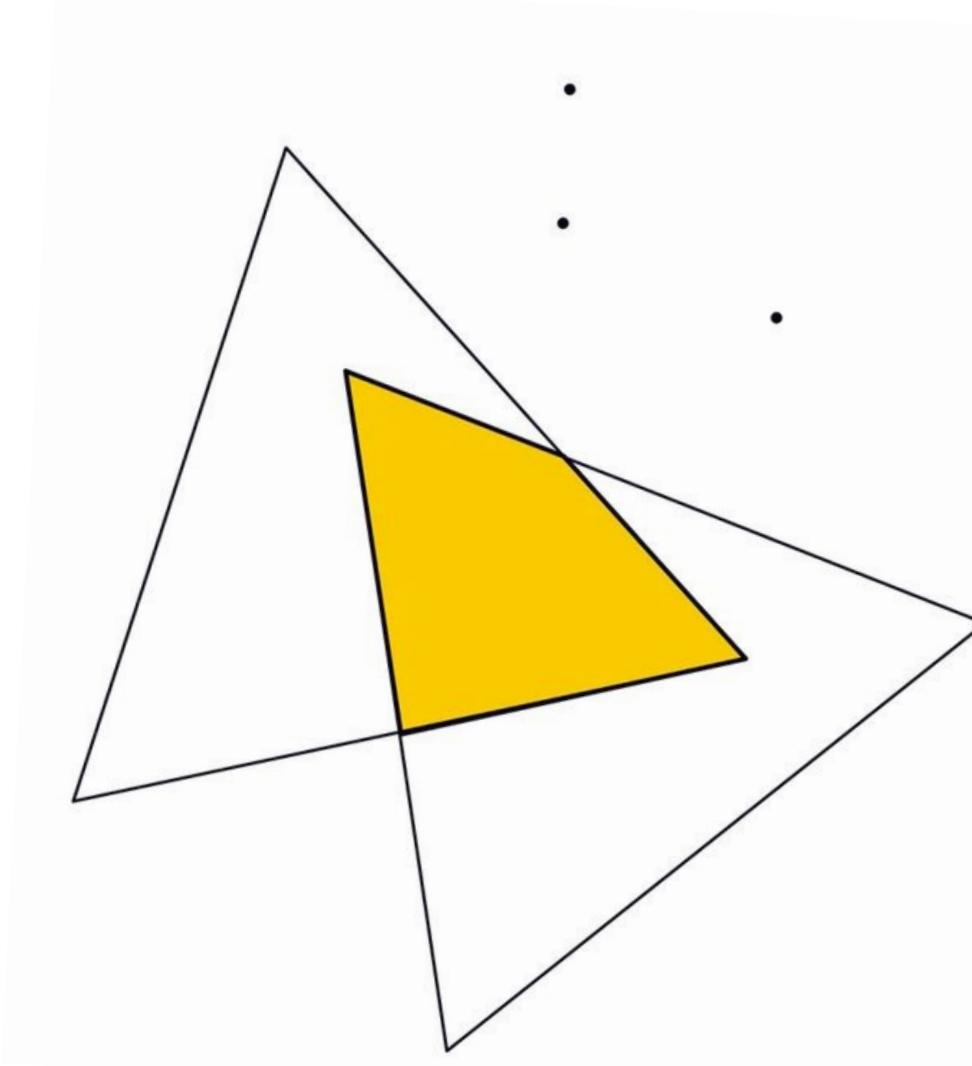
Backpropagation: matrix form

$$y_1 = f_1(\mathbf{x}) = x_1$$

$$y_2 = f_2(\mathbf{x}) = x_2$$

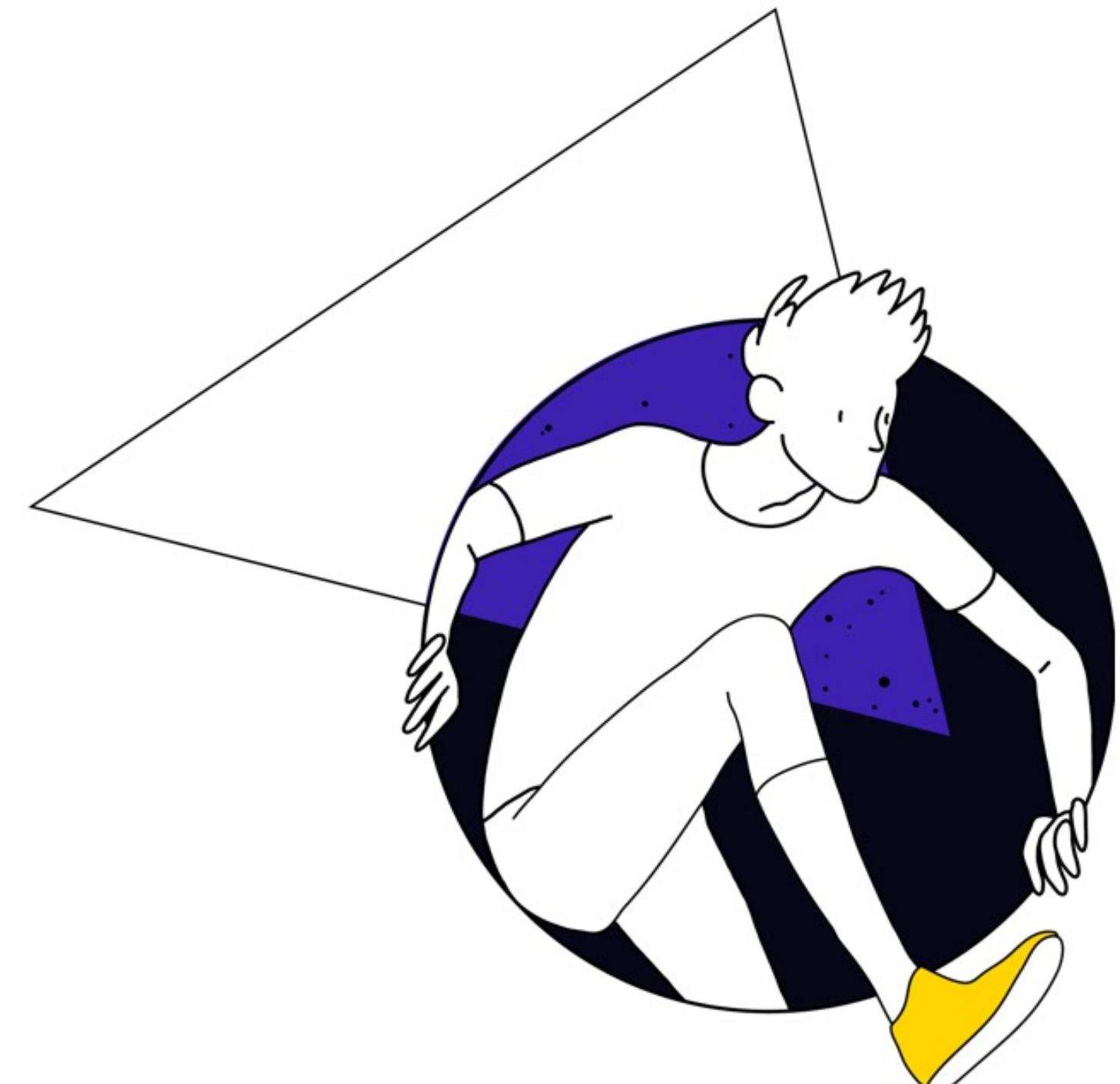
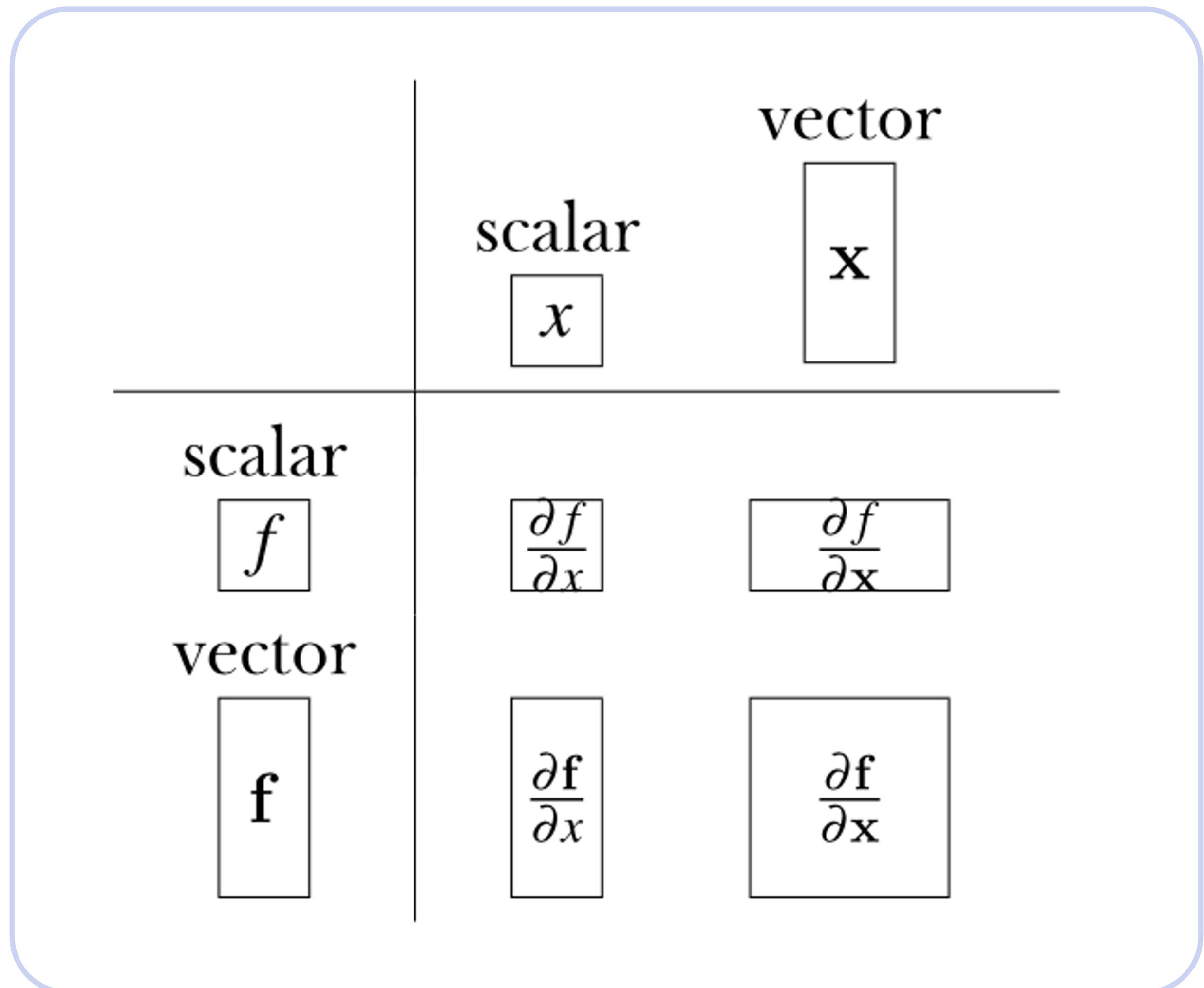
⋮

$$y_n = f_n(\mathbf{x}) = x_n$$



$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \nabla f_1(\mathbf{x}) \\ \nabla f_2(\mathbf{x}) \\ \dots \\ \nabla f_m(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial \mathbf{x}} f_1(\mathbf{x}) \\ \frac{\partial}{\partial \mathbf{x}} f_2(\mathbf{x}) \\ \dots \\ \frac{\partial}{\partial \mathbf{x}} f_m(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_1} f_1(\mathbf{x}) & \frac{\partial}{\partial x_2} f_1(\mathbf{x}) & \dots & \frac{\partial}{\partial x_n} f_1(\mathbf{x}) \\ \frac{\partial}{\partial x_1} f_2(\mathbf{x}) & \frac{\partial}{\partial x_2} f_2(\mathbf{x}) & \dots & \frac{\partial}{\partial x_n} f_2(\mathbf{x}) \\ \dots \\ \frac{\partial}{\partial x_1} f_m(\mathbf{x}) & \frac{\partial}{\partial x_2} f_m(\mathbf{x}) & \dots & \frac{\partial}{\partial x_n} f_m(\mathbf{x}) \end{bmatrix}$$

Backpropagation: matrix form



Backpropagation: matrix form

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial}{\partial \mathbf{x}} f_1(\mathbf{x}) \\ \frac{\partial}{\partial \mathbf{x}} f_2(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial \mathbf{x}} f_m(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_1} f_1(\mathbf{x}) & \frac{\partial}{\partial x_2} f_1(\mathbf{x}) & \dots & \frac{\partial}{\partial x_n} f_1(\mathbf{x}) \\ \frac{\partial}{\partial x_1} f_2(\mathbf{x}) & \frac{\partial}{\partial x_2} f_2(\mathbf{x}) & \dots & \frac{\partial}{\partial x_n} f_2(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial x_1} f_m(\mathbf{x}) & \frac{\partial}{\partial x_2} f_m(\mathbf{x}) & \dots & \frac{\partial}{\partial x_n} f_m(\mathbf{x}) \end{bmatrix}$$

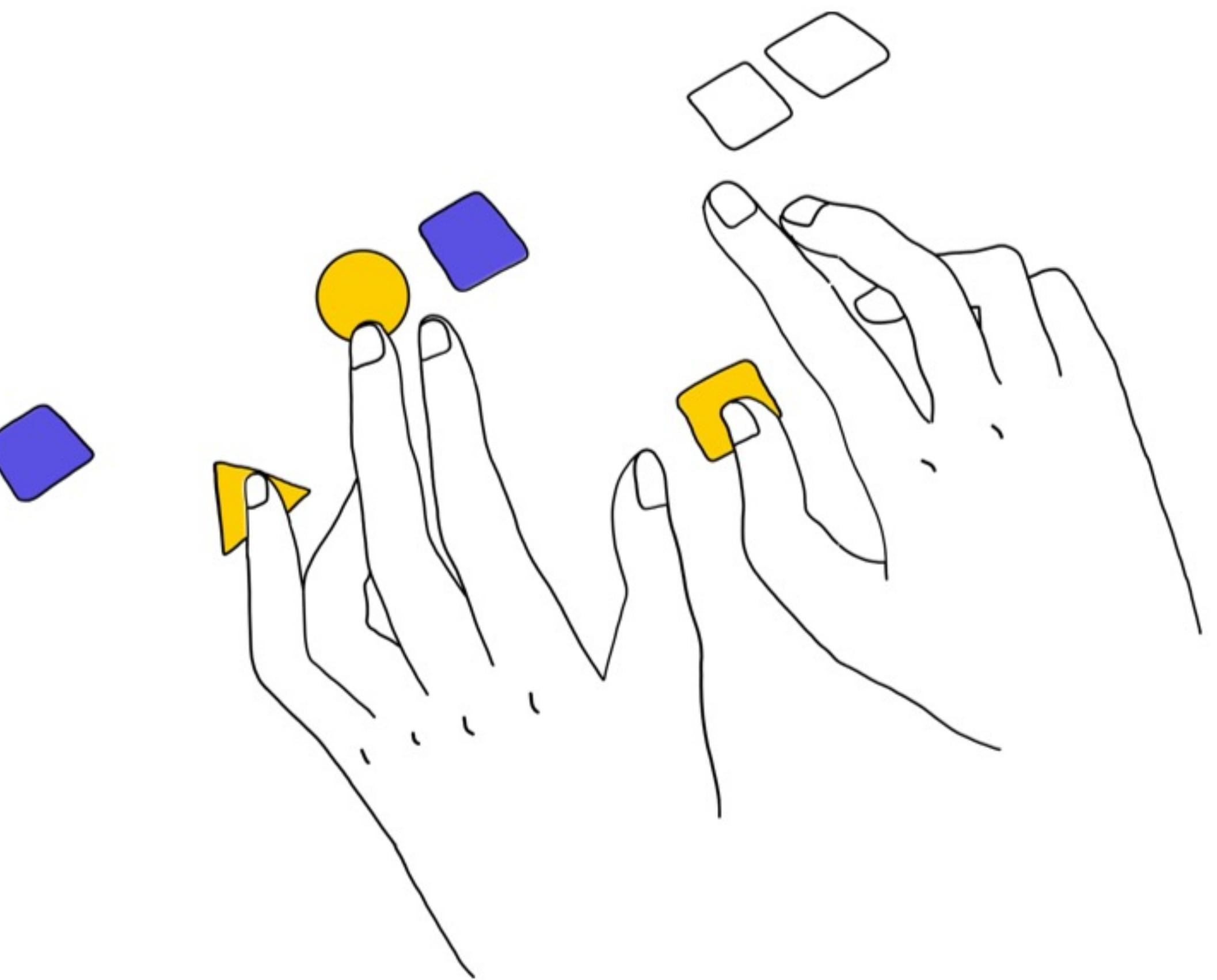
$$= \begin{bmatrix} \frac{\partial}{\partial x_1} x_1 & \frac{\partial}{\partial x_2} x_1 & \dots & \frac{\partial}{\partial x_n} x_1 \\ \frac{\partial}{\partial x_1} x_2 & \frac{\partial}{\partial x_2} x_2 & \dots & \frac{\partial}{\partial x_n} x_2 \\ \vdots \\ \frac{\partial}{\partial x_1} x_n & \frac{\partial}{\partial x_2} x_n & \dots & \frac{\partial}{\partial x_n} x_n \end{bmatrix}$$

(and since $\frac{\partial}{\partial x_j} x_i = 0$ for $j \neq i$)

$$= \begin{bmatrix} \frac{\partial}{\partial x_1} x_1 & 0 & \dots & 0 \\ 0 & \frac{\partial}{\partial x_2} x_2 & \dots & 0 \\ \ddots & & \ddots & \\ 0 & 0 & \dots & \frac{\partial}{\partial x_n} x_n \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \ddots & & \ddots & \\ 0 & 0 & \dots & 1 \end{bmatrix}$$

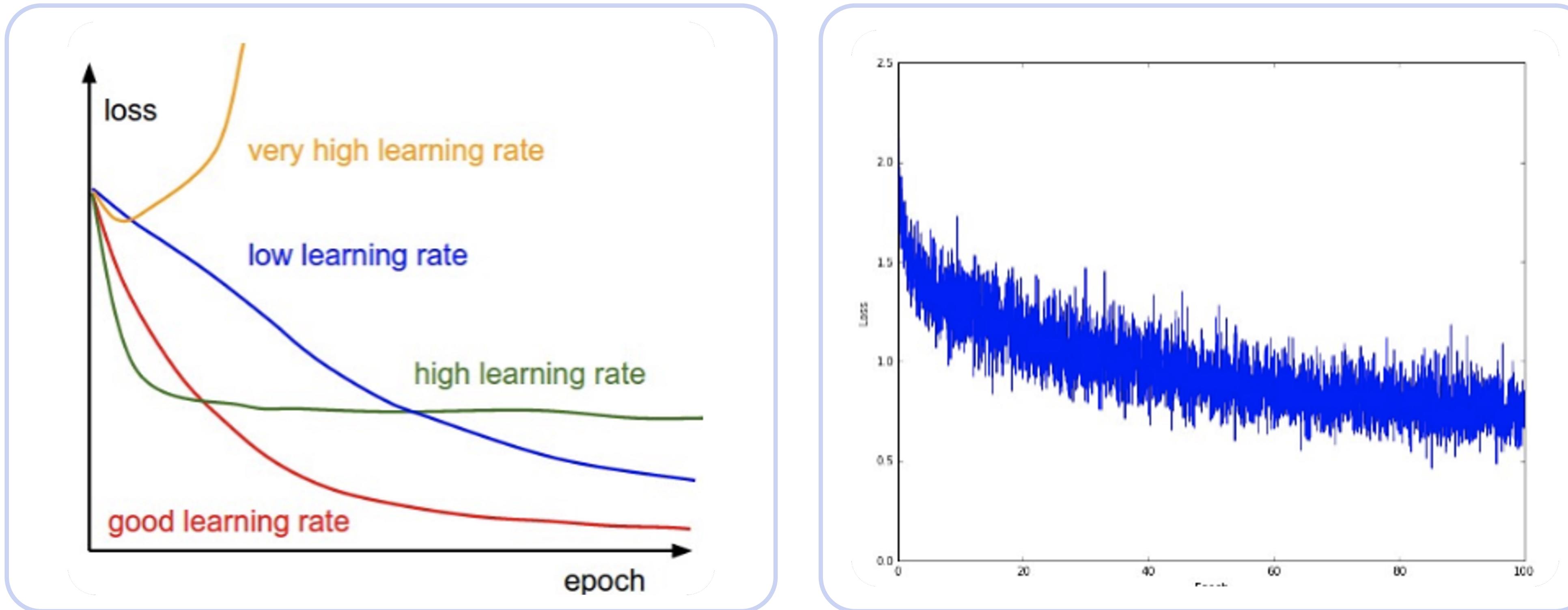
= I (I is the identity matrix with ones down the diagonal)



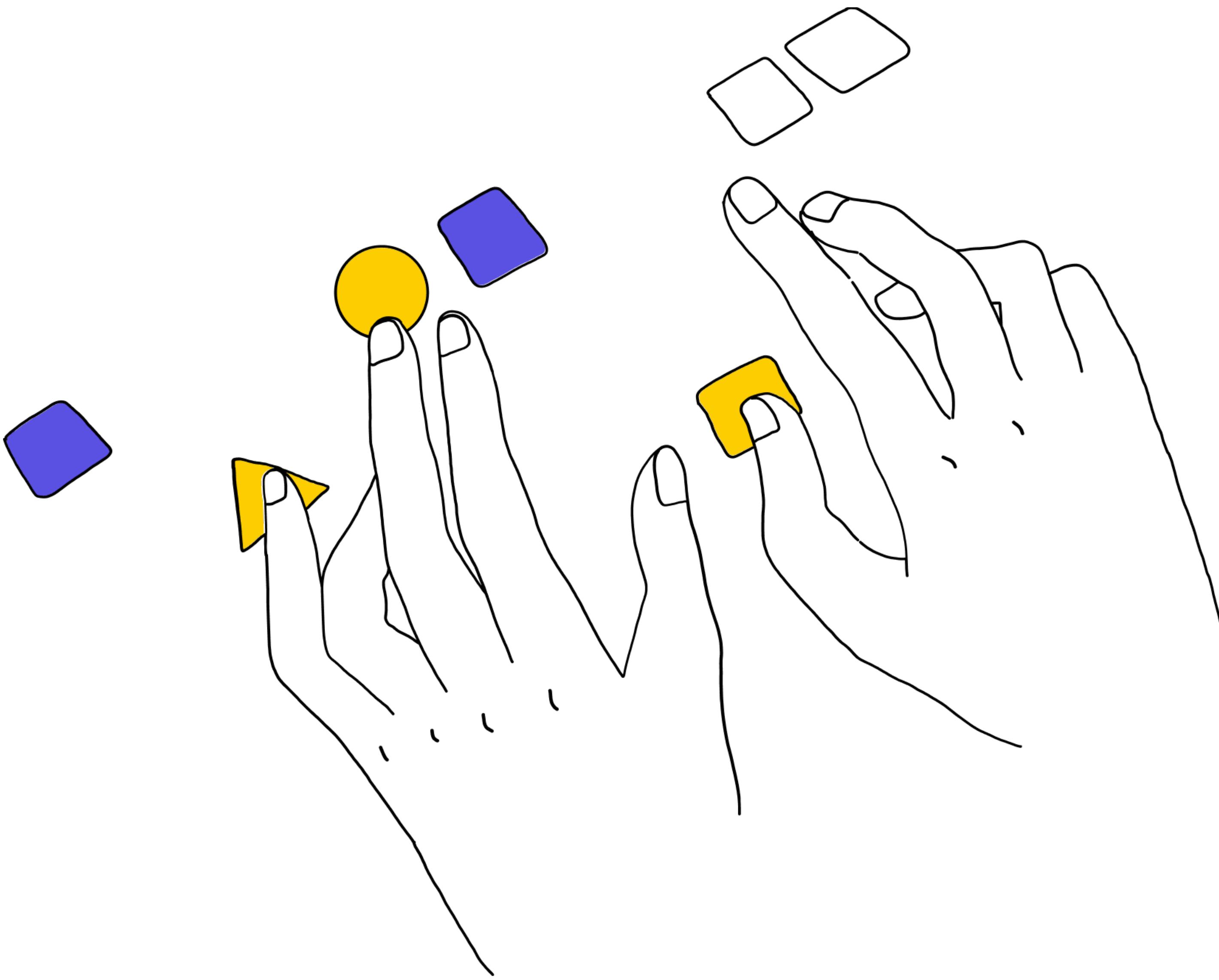
Gradient optimization

Stochastic gradient descent is used to optimize NN parameters

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \frac{dL}{d\mathbf{w}}$$



Activation functions



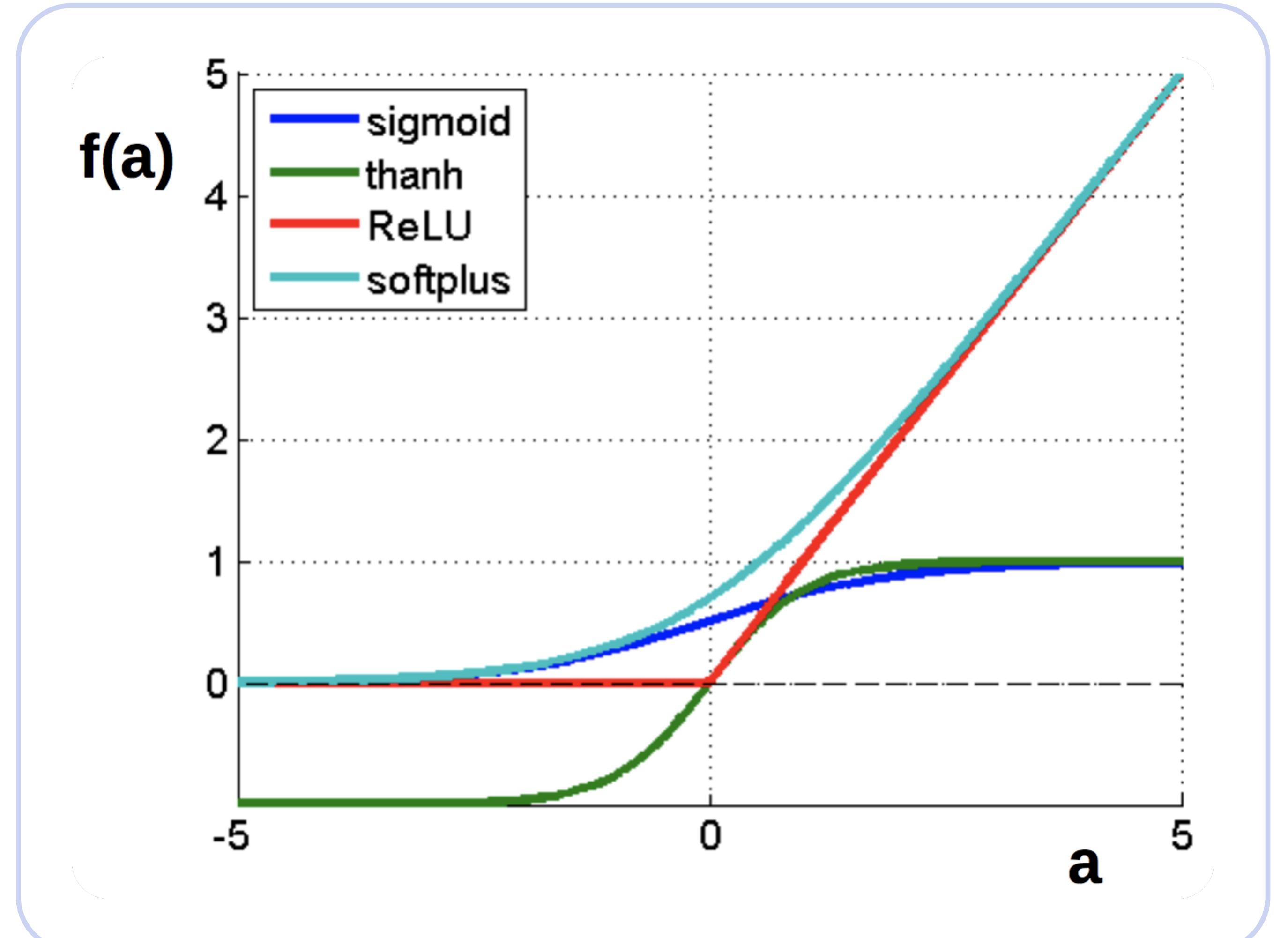
Once again: nonlinearities

$$f(a) = \frac{1}{1 + e^{-a}}$$

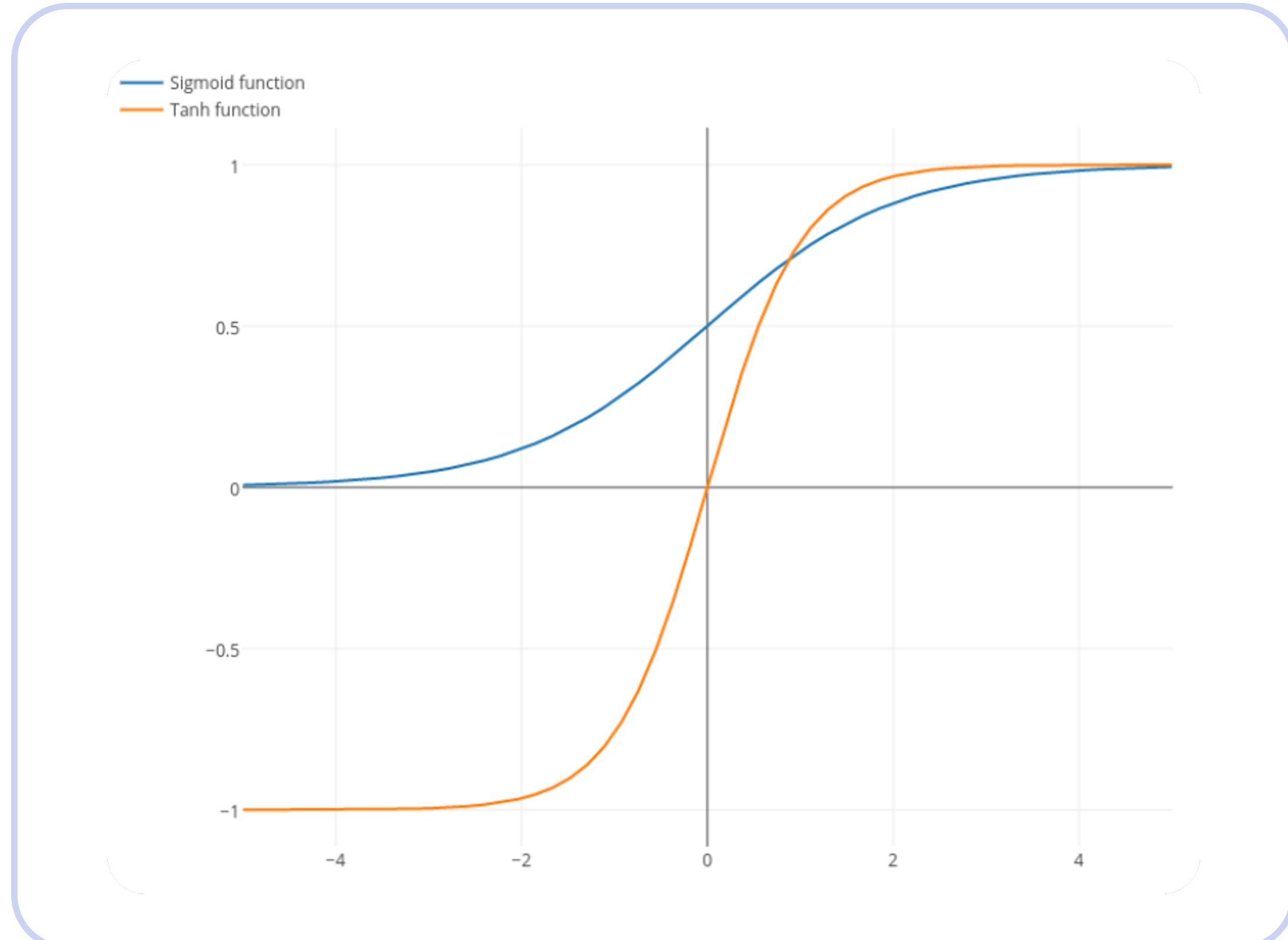
$$f(a) = \tanh(a)$$

$$f(a) = \max(0, a)$$

$$f(a) = \log(1 + e^a)$$



Activation functions: Sigmoid



Maps R to (0,1)

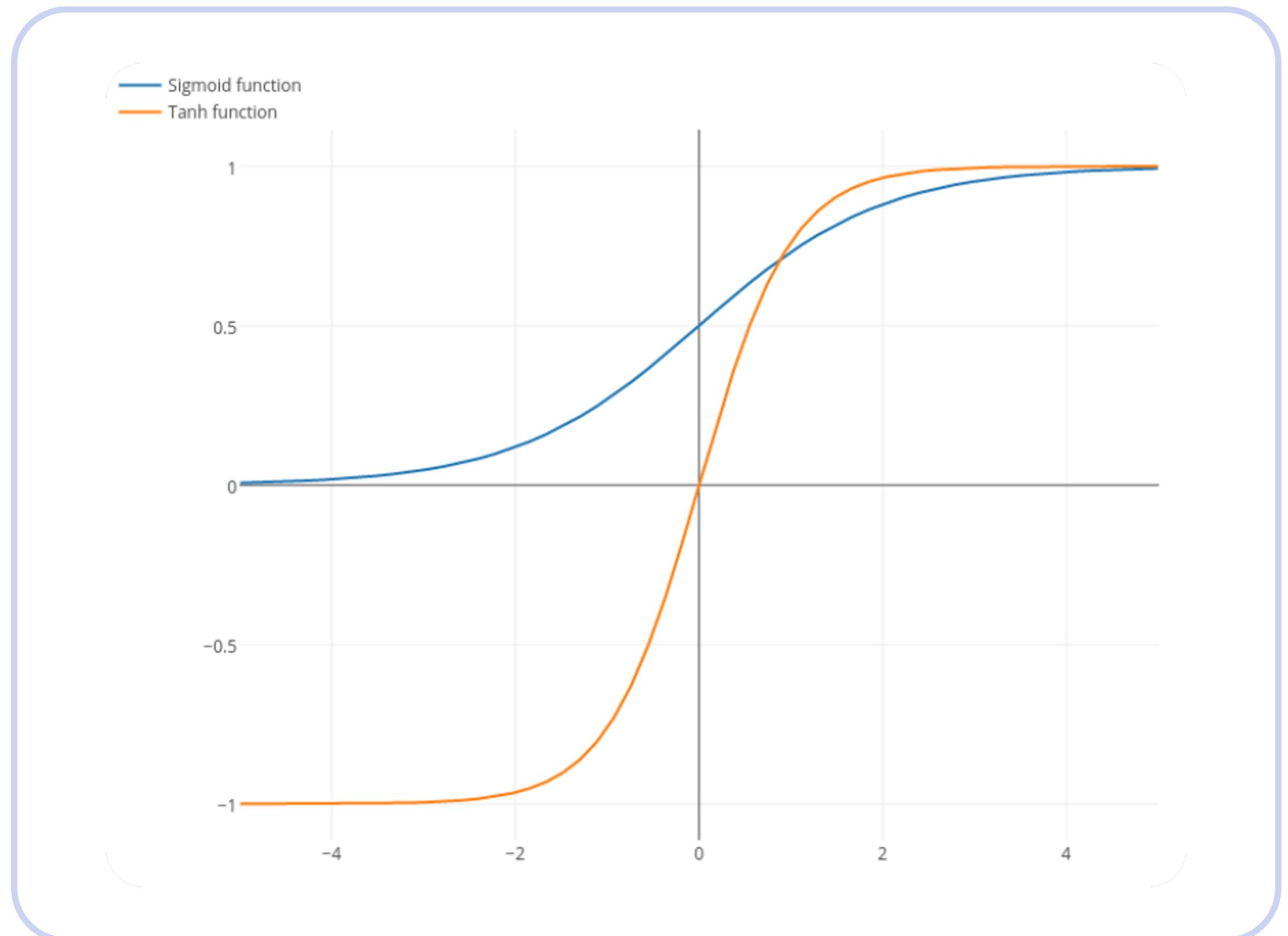
Historically popular, one of the first approximations of neuron activation

Problems:

- Almost zero gradients on the both sides (saturation)
- Shifted (not zero-centered) output
- Expensive computation of the exponent

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

Activation functions: tanh



Maps R to (0,1)

Similar to the Sigmoid in other ways

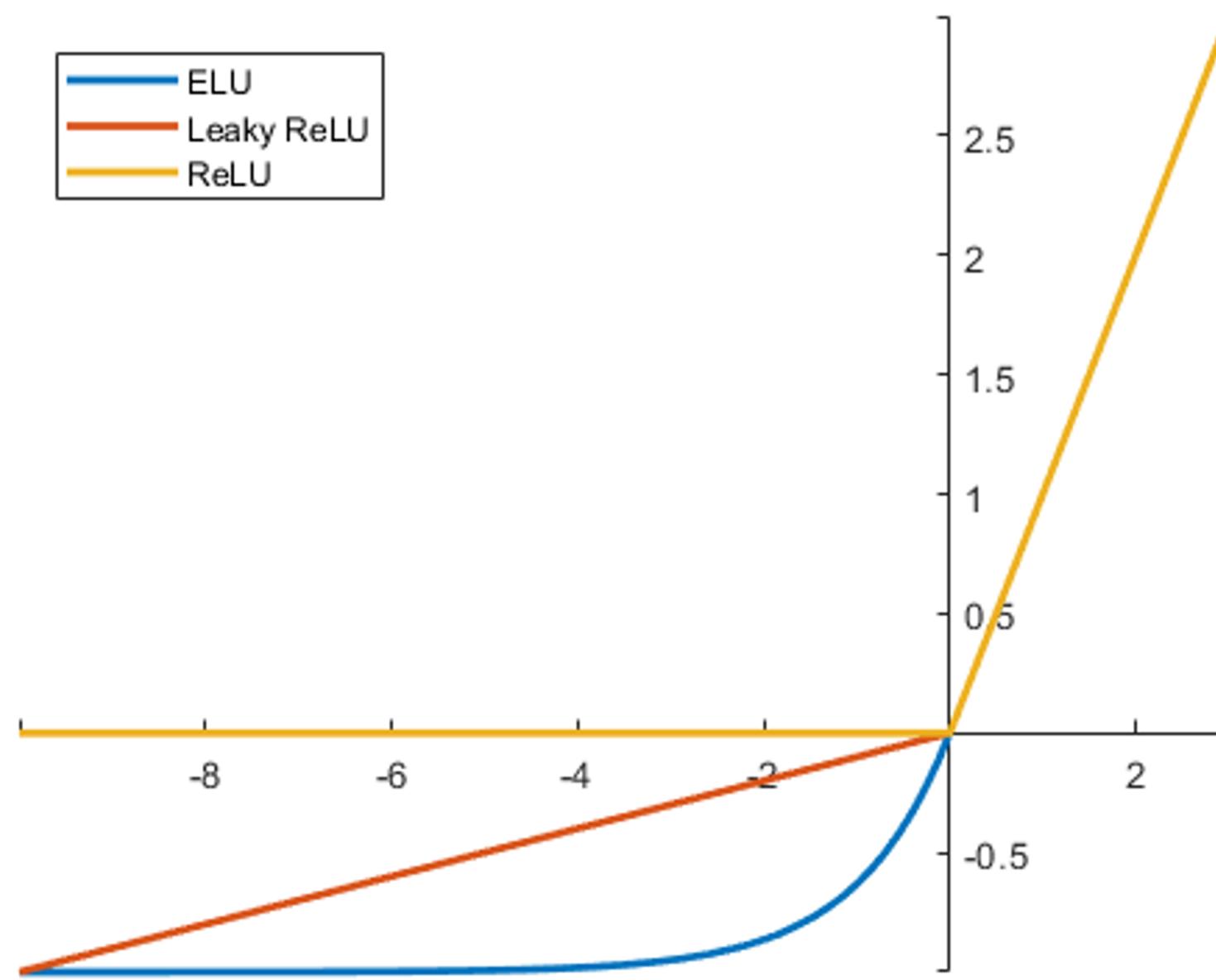
Problems:

- Almost zero gradients on the both sides (saturation)
- ~~Shifted (not zero-centered) output~~
- Expensive computation of the exponent

$$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$$

Activation functions: ReLU

ELU
Leaky ReLU
ReLU



Very simple to compute
(both forward and backward)

Up to 6 times faster than Sigmoid

Does not saturate when $x > 0$

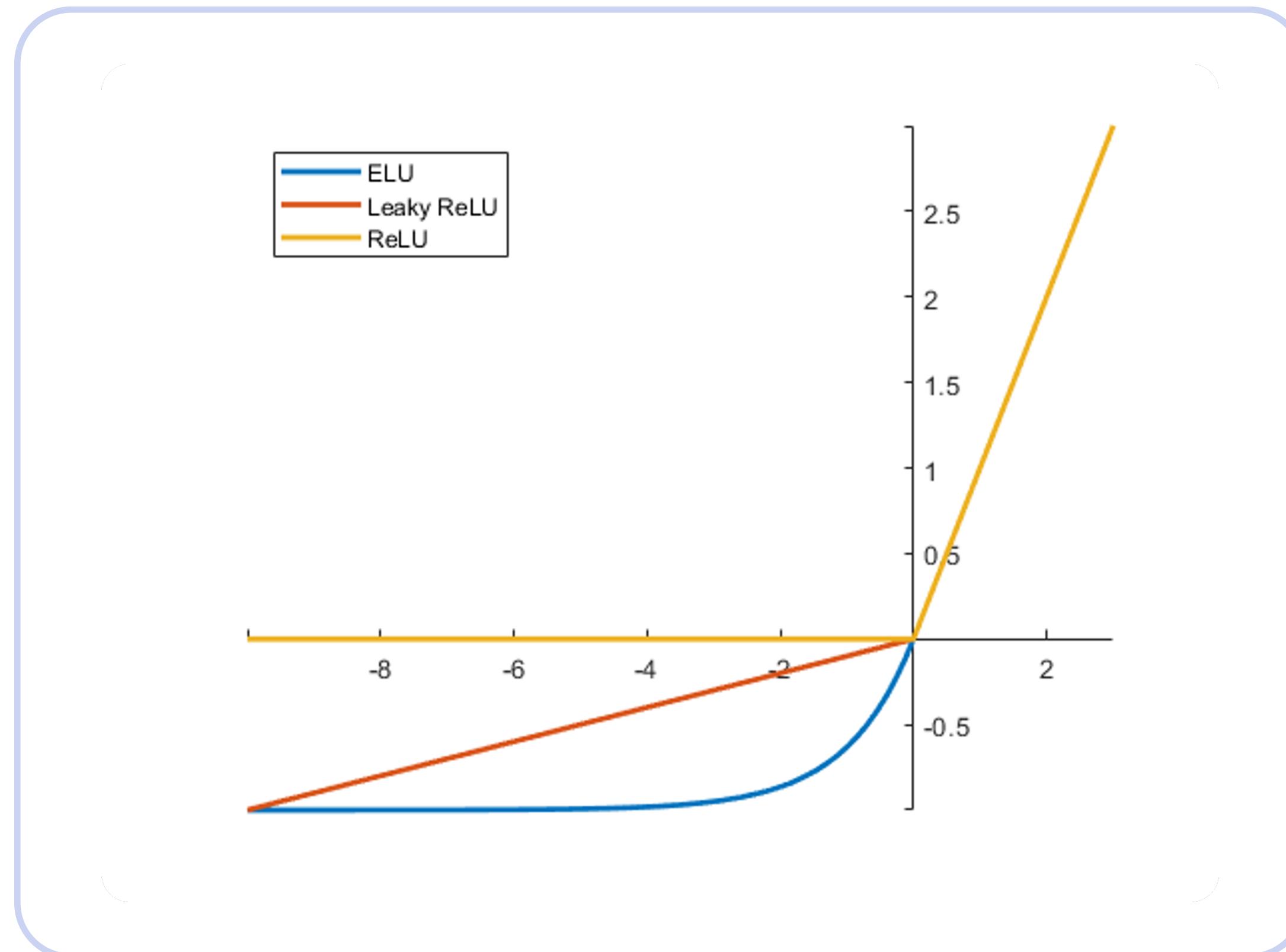
So the gradients are not 0

Problems:

- Zero gradients when $x < 0$
- Shifted (not zero-centered) output

$$\text{ReLU}(a) = \max(\{0, a\})$$

Activation functions: Leaky ReLU



Very simple to compute
(both forward and backward)

Up to 6 times faster than Sigmoid

Does not saturate when $x > 0$

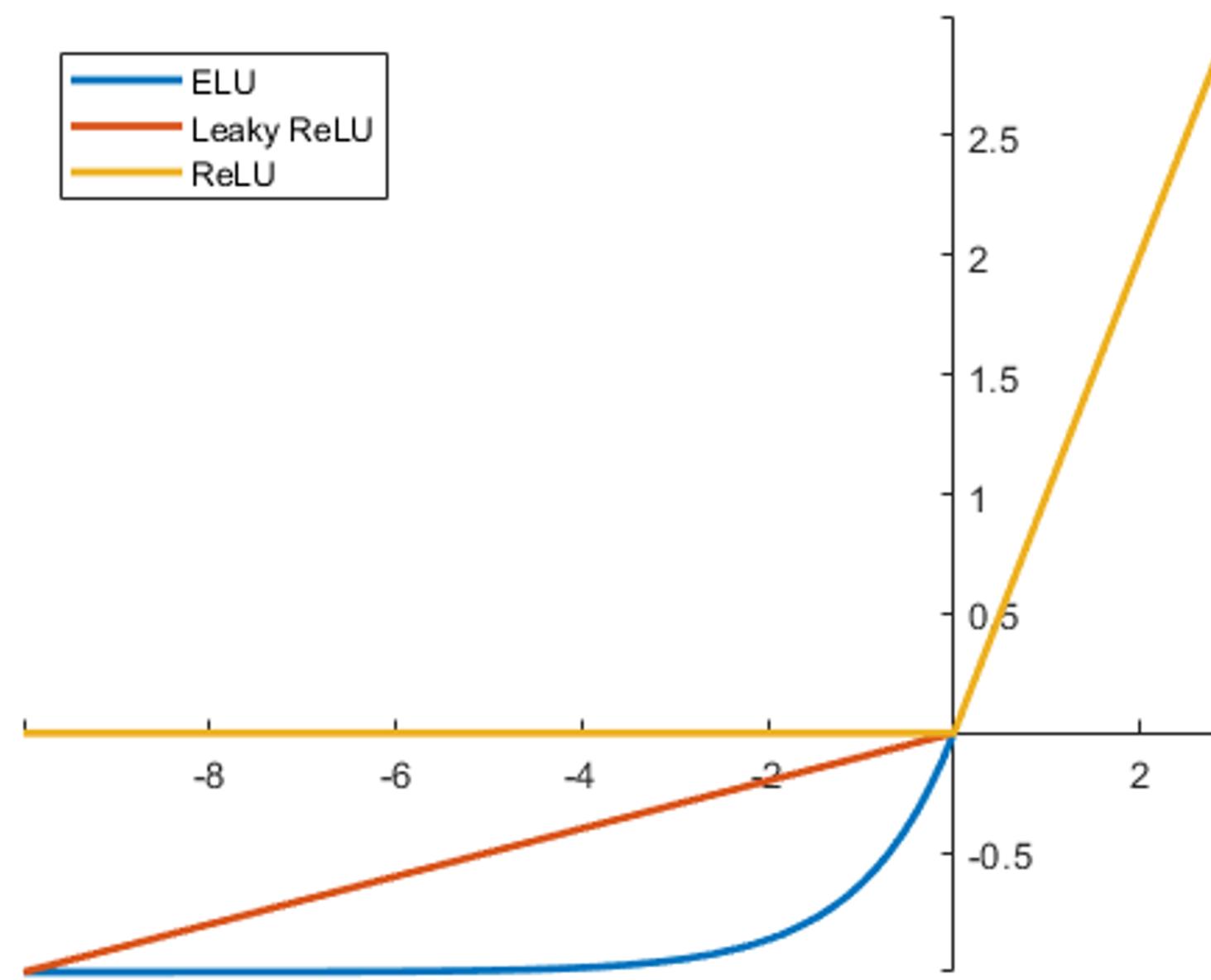
Problems:

- Shifted (not zero-centered) output

$$\text{Leaky ReLU}(a) = \max(\{0.01a, a\})$$

Activation functions: ELU

ELU
Leaky ReLU
ReLU



Similar to ReLU

Does not saturate

Close to zero mean outputs

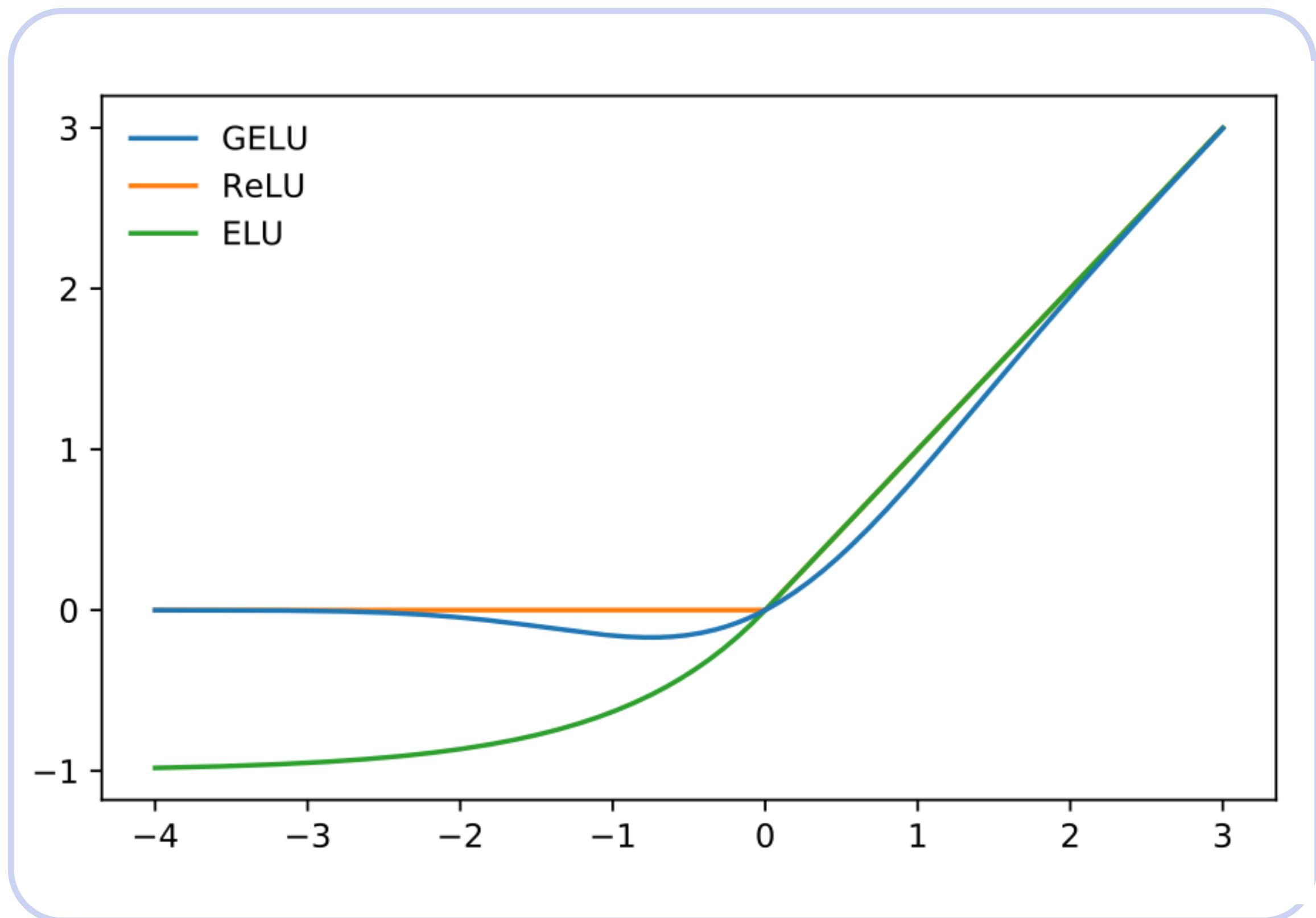
Problems:

Requires exponent computation

$$\text{ELU}(a) = \begin{cases} a & \text{if } a > 0 \\ \beta(\exp(a) - 1) & \text{if } a \leq 0 \end{cases}$$

Activation functions: GELU

Gaussian error linear unit



Outperforms ReLU and ELU according to the experiments

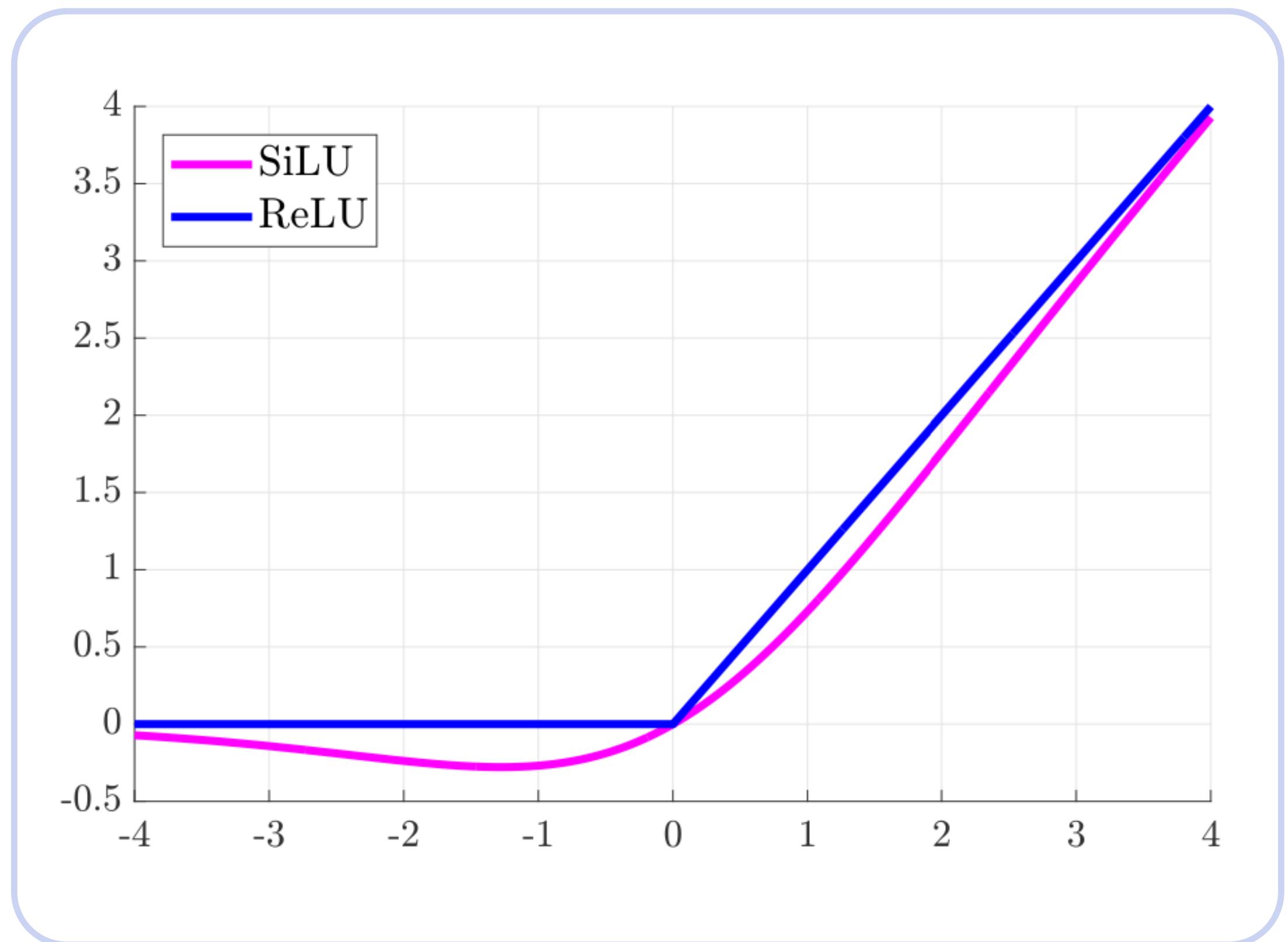
PyTorch's exact implementation is sufficiently fast such that these approximations may be unnecessary.

Problems:

Requires exponent computation

$$\text{GELU}(a) = aP(X \leq a) = x\Phi(a) = \frac{1}{2}a \left[1 + \operatorname{erf}\left(\frac{a}{\sqrt{2}}\right) \right], \quad X \in \mathcal{N}(0, 1)$$

Activation functions: SiLU Sigmoid-weighted linear unit



Special case of Swish activation function

Originally proposed for RL
algorithms

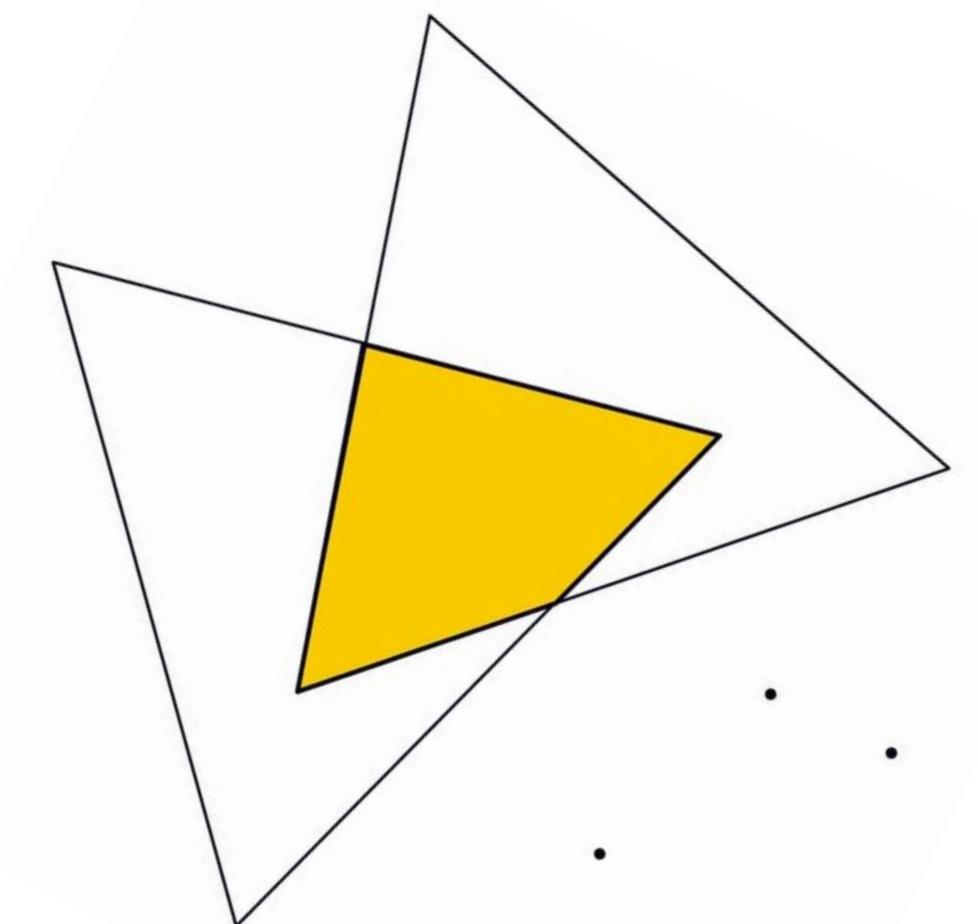
Problems:

Requires exponent computation

$$\text{SiLU}(a) = a\sigma(a)$$

Activation functions: sum up

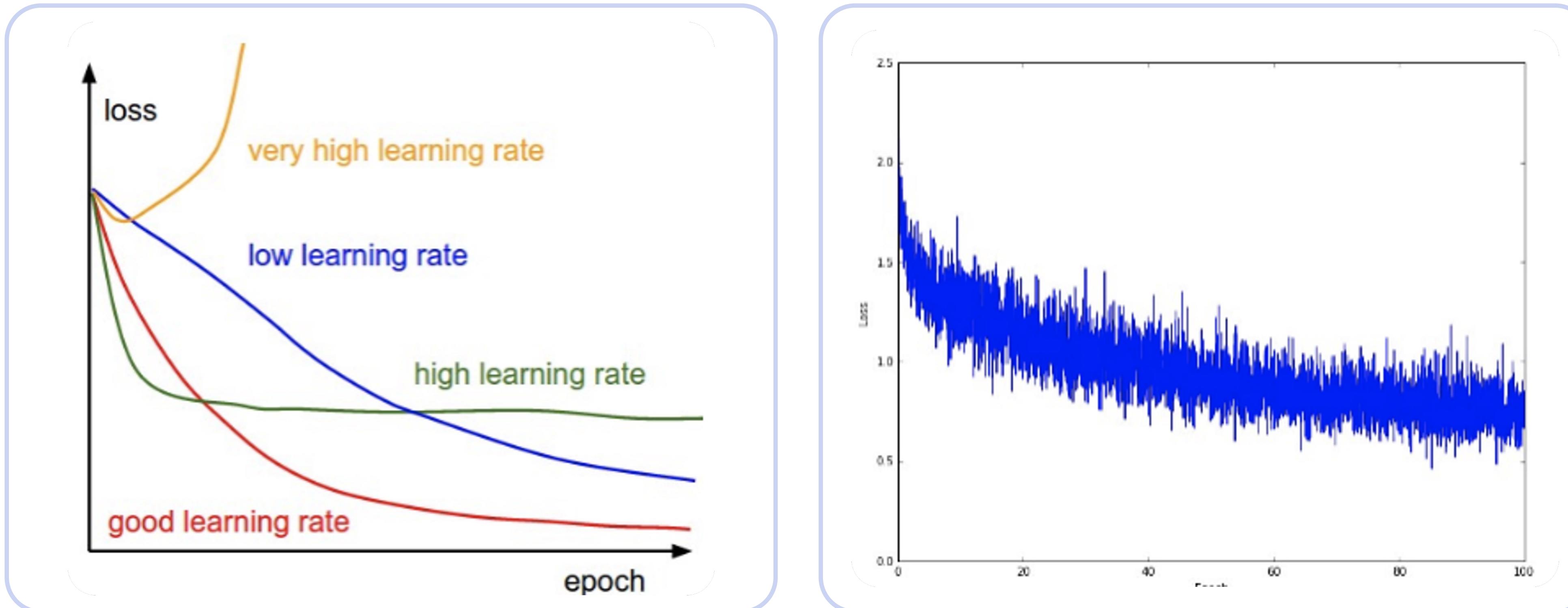
- Use **ReLU** as baseline approach
- Be careful with the learning rates
- Try out **Leaky ReLU** or **ELU**
- Try out **tanh** but do not expect much from it
- **Sigmoid** usually means binary classification somewhere close



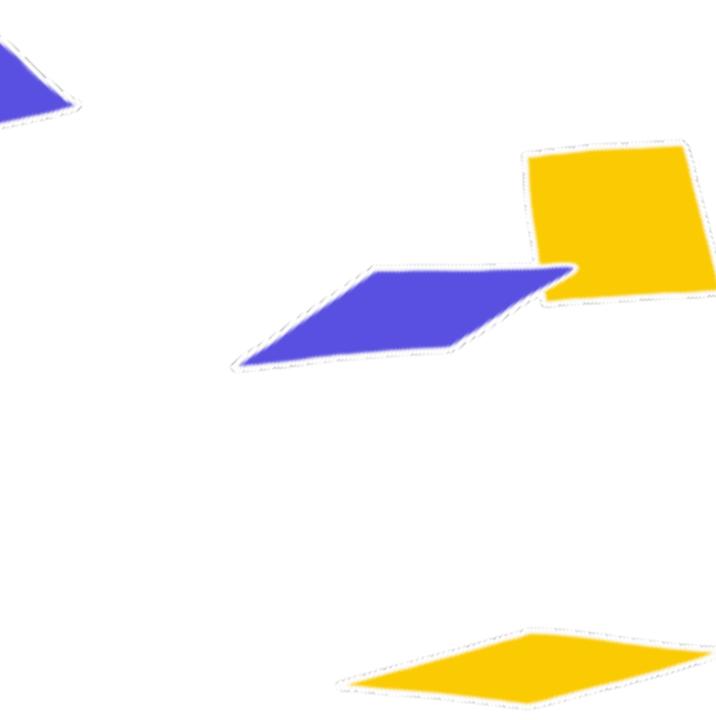
Once again: gradient optimization

Stochastic gradient descent is used to optimize NN parameters

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \frac{dL}{d\mathbf{w}}$$

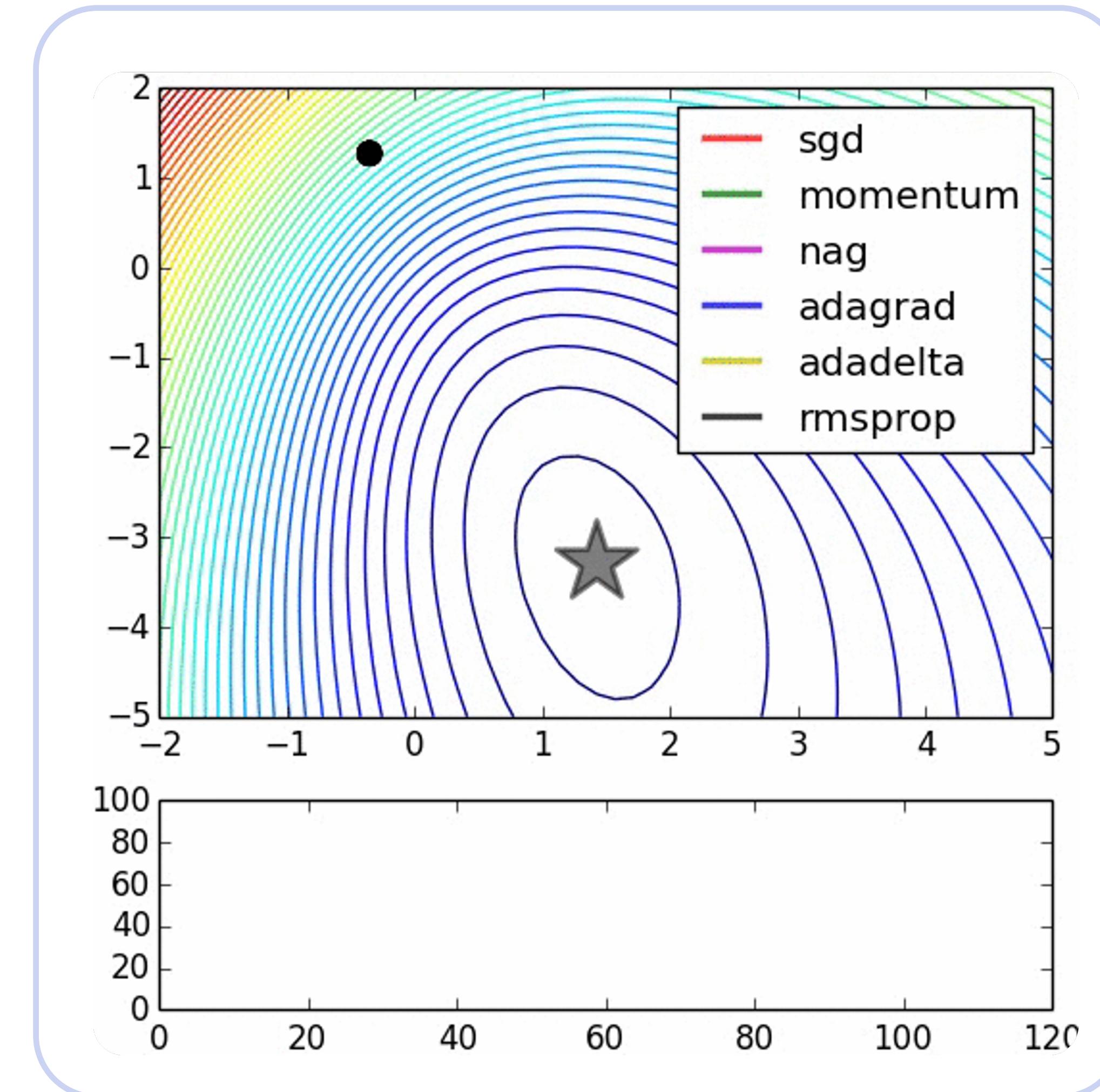


Optimizers



There are much more optimizers:

- Momentum
- Adagrad
- Adadelta
- RMSprop
- Adam
- ...
- even other NNs

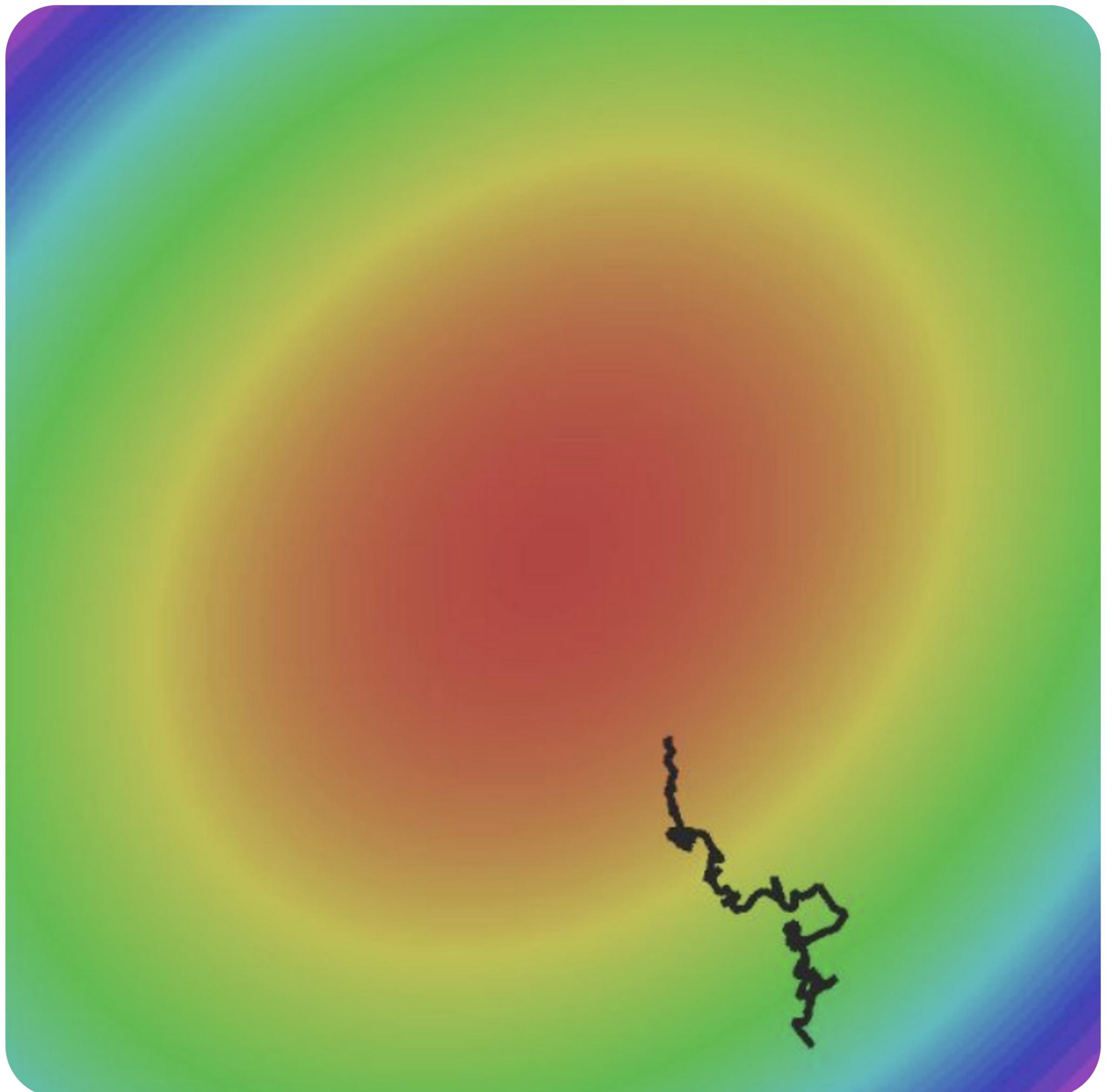


Optimization: SGD

$$L(\mathbf{w}) = \frac{1}{K} \sum_{i=1}^K L(\mathbf{x}^{(i)}, y^{(i)}, \mathbf{w})$$

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = \frac{1}{K} \sum_{i=1}^K \nabla_{\mathbf{w}} L(\mathbf{x}^{(i)}, y^{(i)}, \mathbf{w})$$

Averaging over minibatches leads to noisy gradient estimation



First idea: momentum

Simple SGD

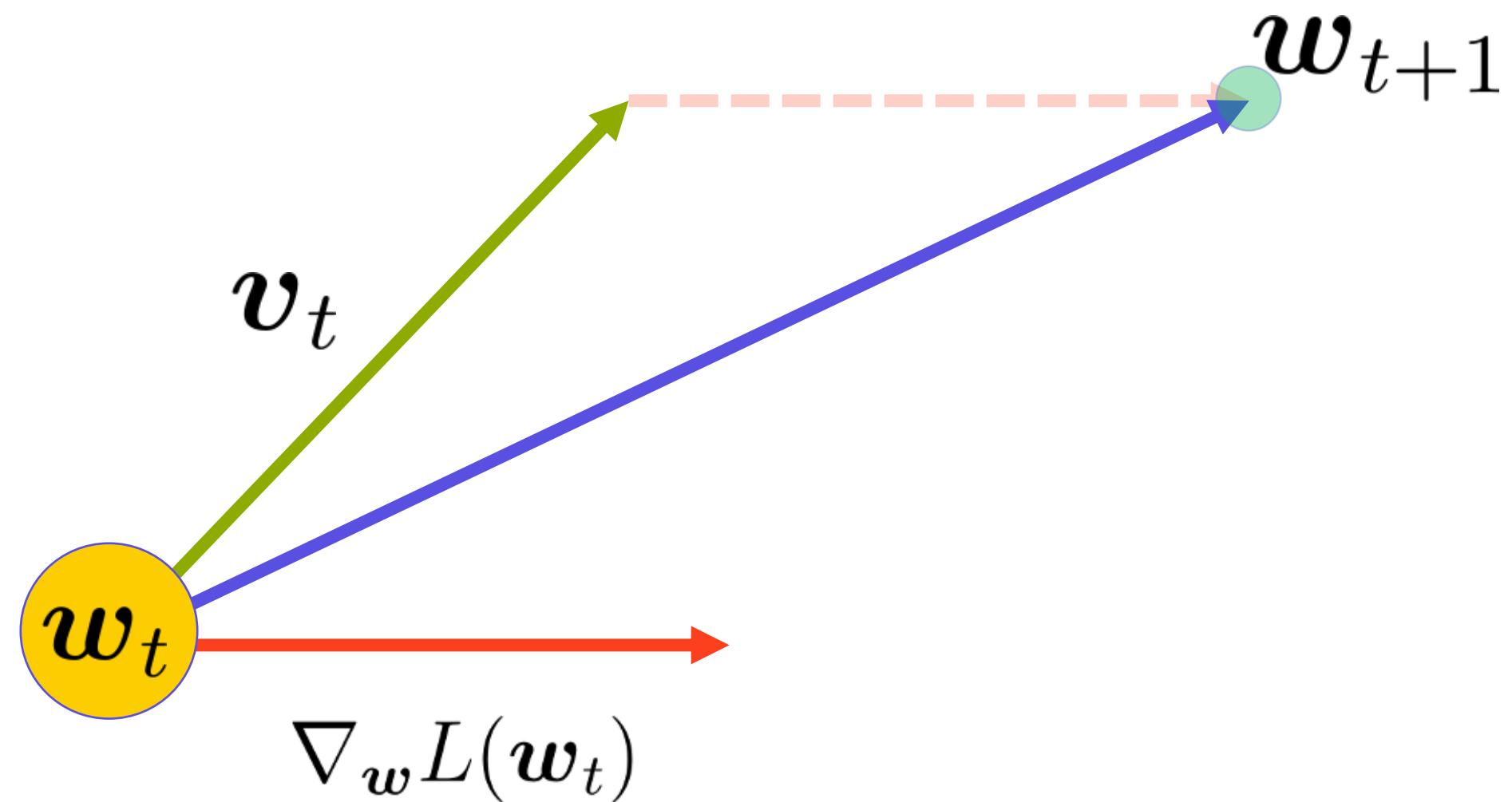
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla_{\mathbf{w}} L(\mathbf{w}_t)$$

SGD with momentum

$$\mathbf{v}_{t+1} = \rho \mathbf{v}_t + \nabla_{\mathbf{w}} L(\mathbf{w}_t)$$

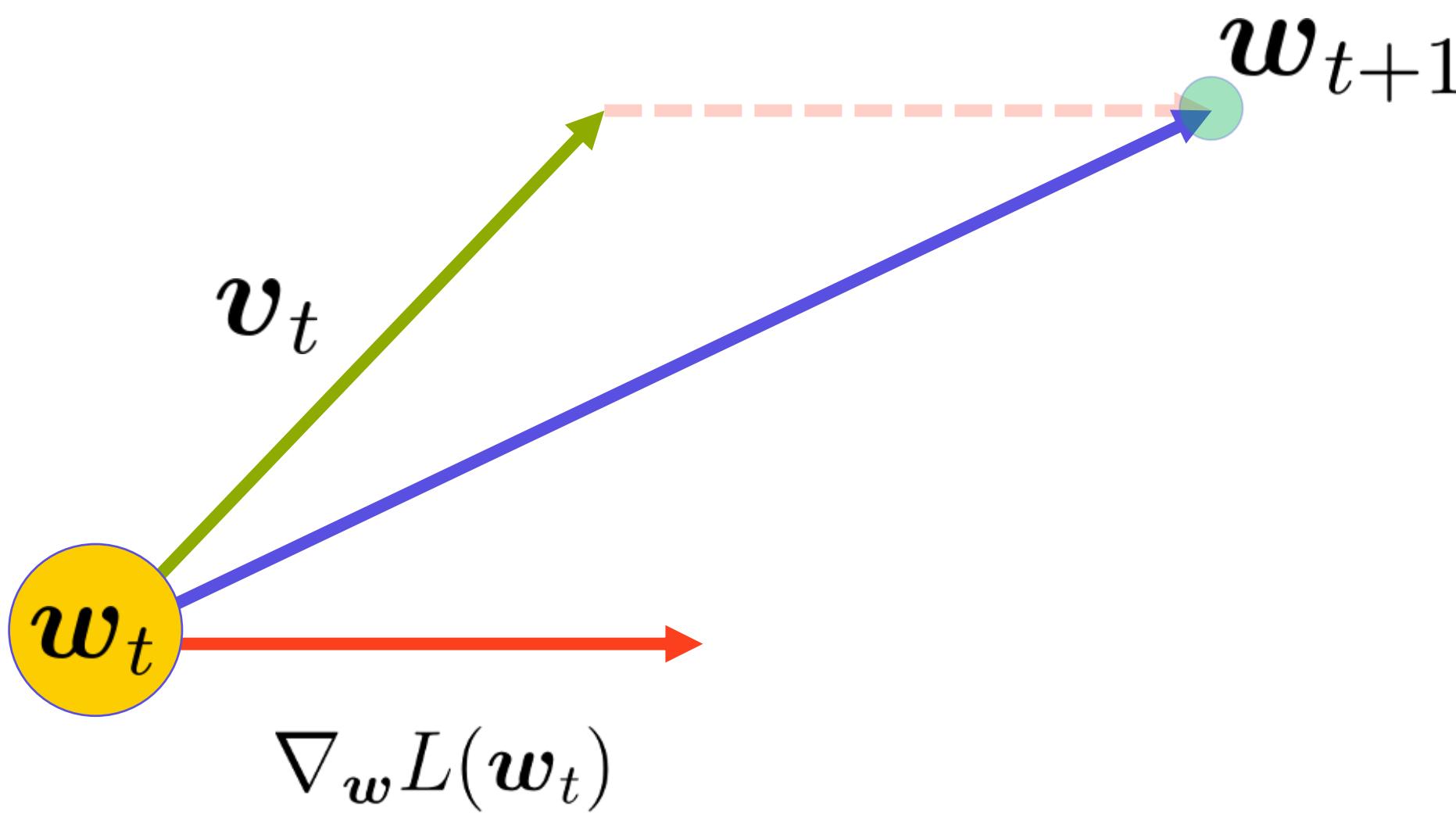
$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{v}_{t+1}$$

Momentum update:



Nesterov momentum

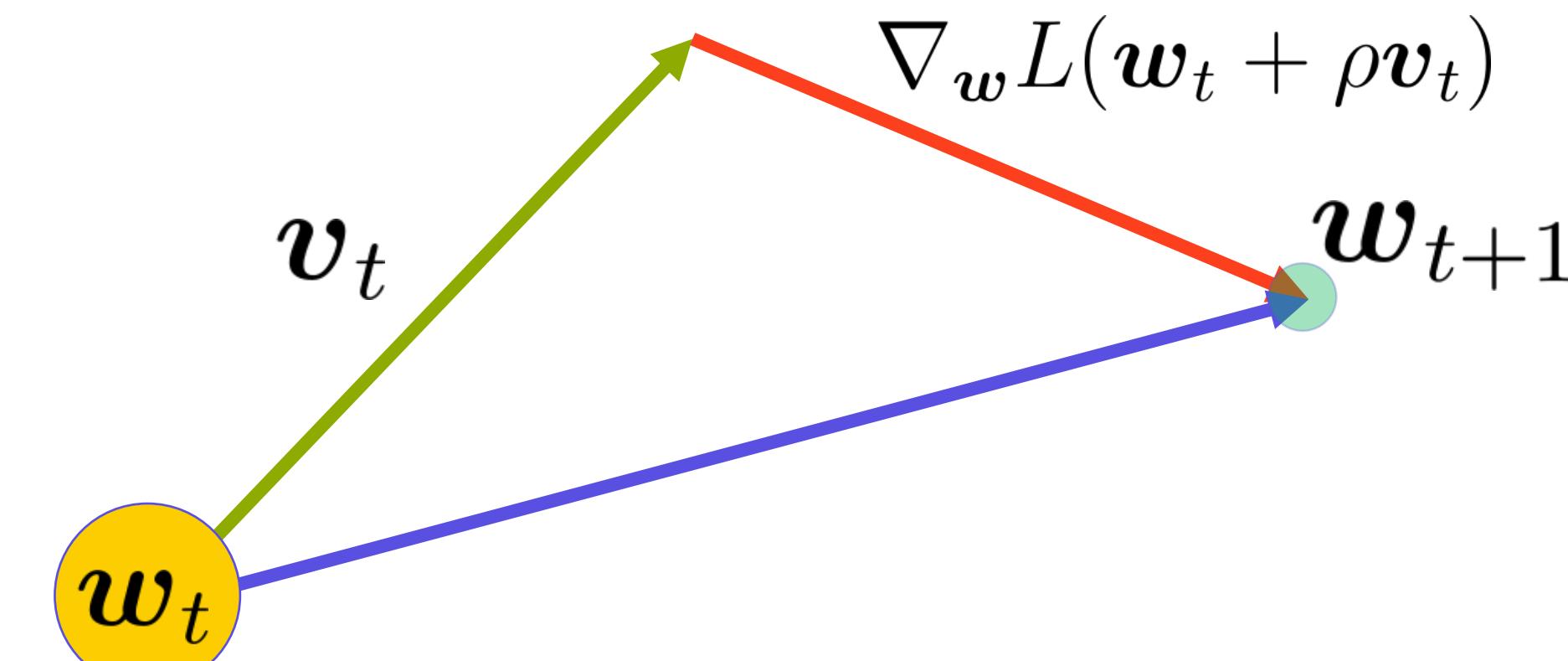
Momentum update:



$$v_{t+1} = \rho v_t + \nabla_w L(w_t)$$

$$w_{t+1} = w_t - \eta v_{t+1}$$

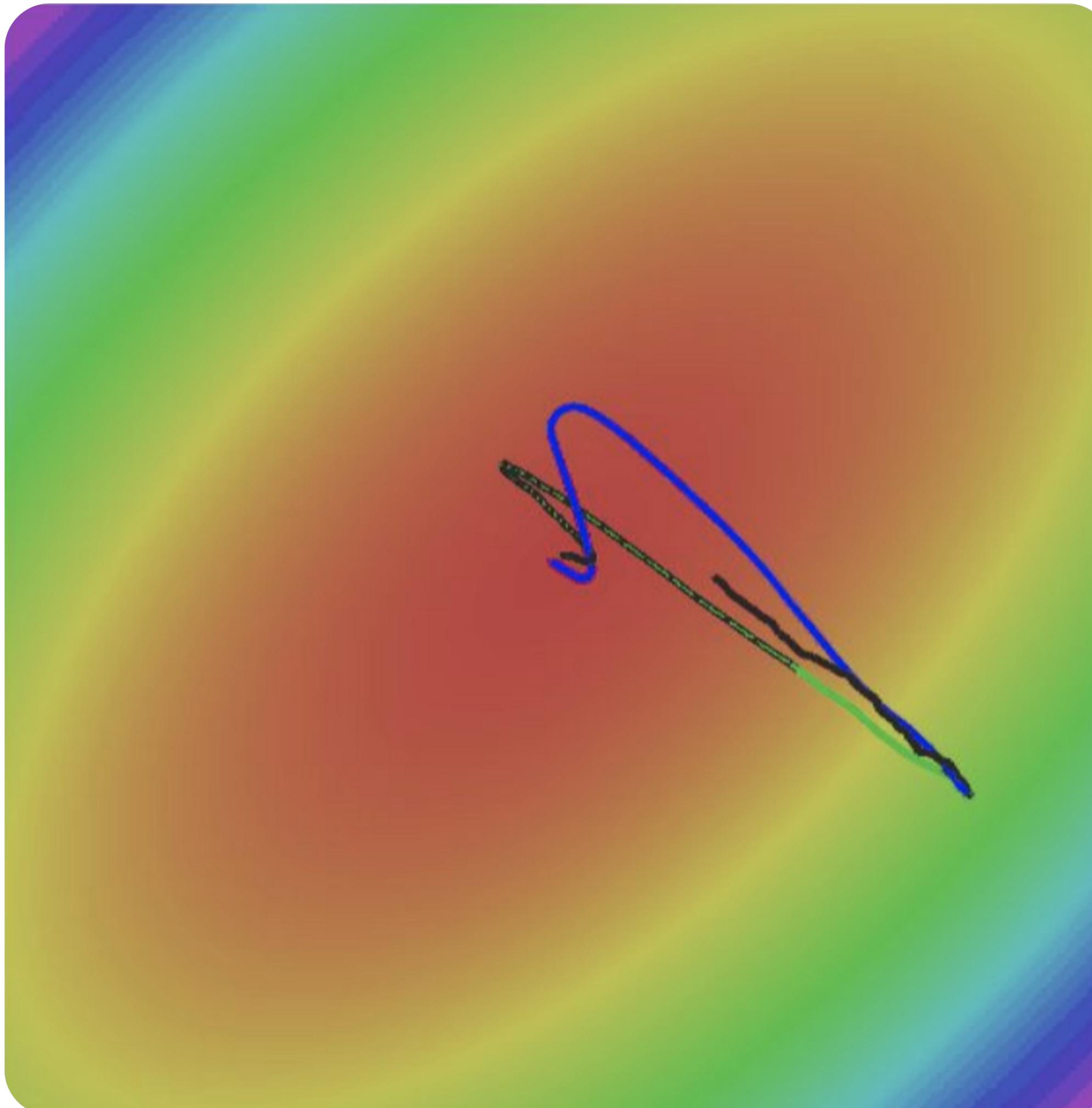
Nesterov Momentum



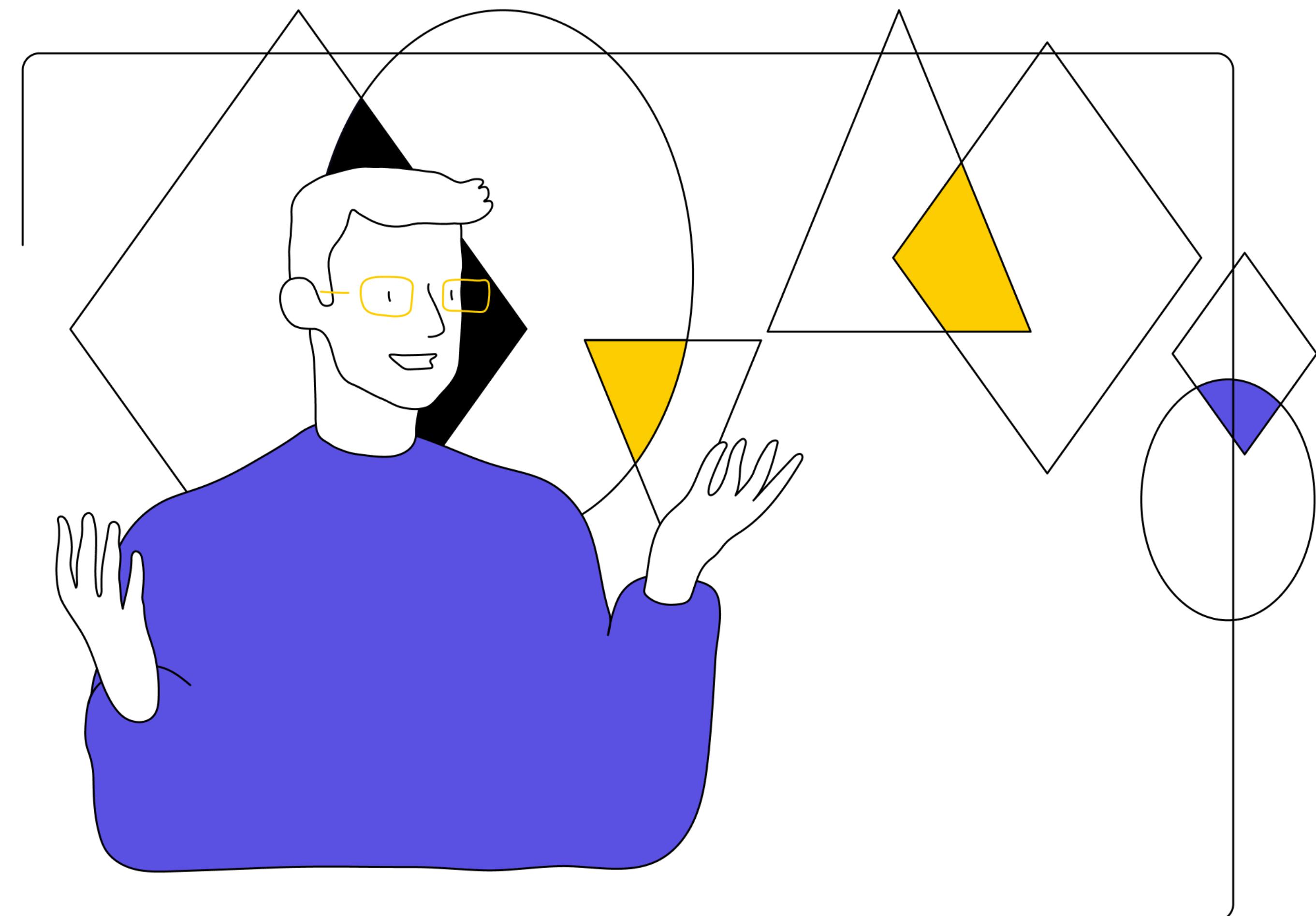
$$v_{t+1} = \rho v_t + \nabla_w L(w_t + \rho v_t)$$

$$w_{t+1} = w_t - \eta v_{t+1}$$

Comparing momentums



Fancy (tiny) neural networks



Shakespeare

PANDARUS:

Alas, I think he shall be come approached and the day
 When little strain would be attain'd into being never fed,
 And who is but a chain and subjects of his death,
 I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,
 Breaking and strongly should be buried, when I perish
 The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

Second Lord:

They would be ruled after this chamber, and
 my fair nues begun out of the fact, to be conveyed,
 Whose noble souls I'll have the heart of the wars.

Clown:

Come, sir, I will make did behold your worship.

VIOLA:

I'll drink it.

Algebraic Geometry (Latex)

Proof. Omitted. \square

Lemma 0.1. Let \mathcal{C} be a set of the construction.

Let \mathcal{C} be a gerber covering. Let \mathcal{F} be a quasi-coherent sheaves of \mathcal{O} -modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

Proof. This is an algebraic space with the composition of sheaves \mathcal{F} on $X_{\text{étale}}$ we have

$$\mathcal{O}_X(\mathcal{F}) = \{\text{morph}_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where \mathcal{G} defines an isomorphism $\mathcal{F} \rightarrow \mathcal{F}$ of \mathcal{O} -modules. \square

Lemma 0.2. This is an integer \mathcal{Z} is injective.

Proof. See Spaces, Lemma ??.

Lemma 0.3. Let S be a scheme. Let X be a scheme and X is an affine open covering. Let $\mathcal{U} \subset \mathcal{X}$ be a canonical and locally of finite type. Let X be a scheme. Let X be a scheme which is equal to the formal complex.

The following to the construction of the lemma follows.

Let X be a scheme. Let X be a scheme covering. Let

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

be a morphism of algebraic spaces over S and Y .

Proof. Let X be a nonzero scheme of X . Let X be an algebraic space. Let \mathcal{F} be a quasi-coherent sheaf of \mathcal{O}_X -modules. The following are equivalent

- (1) \mathcal{F} is an algebraic space over S .
- (2) If X is an affine open covering.

Consider a common structure on X and X the functor $\mathcal{O}_X(U)$ which is locally of finite type. \square

Linux kernel (source code)

```
/*
 * If this error is set, we will need anything right after that BSD.
 */
static void action_new_function(struct s_stat_info *wb)
{
    unsigned long flags;
    int lel_idx_bit = e->edd, *sys = ~((unsigned long) *FIRST_COMPAT);
    buf[0] = 0xFFFFFFFF & (bit << 4);
    min(inc, slist->bytes);
    printk(KERN_WARNING "Memory allocated %02x/%02x, "
        "original MLL instead\n"),
    min(min(multi_run - s->lem, max) * num_data_in),
    frame_pos, sz + first_seg);
    div_u64_w(val, insb_p);
    spin_unlock(&disk->queue_lock);
    mutex_unlock(&s->sock->mutex);
    mutex_unlock(&func->mutex);
    return disassemble(info->pending_bh);
}
```

Proof. Omitted.

□

Lemma 0.1. *Let \mathcal{C} be a set of the construction.*

Let \mathcal{C} be a gerber covering. Let \mathcal{F} be a quasi-coherent sheaves of \mathcal{O} -modules. We have to show that

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

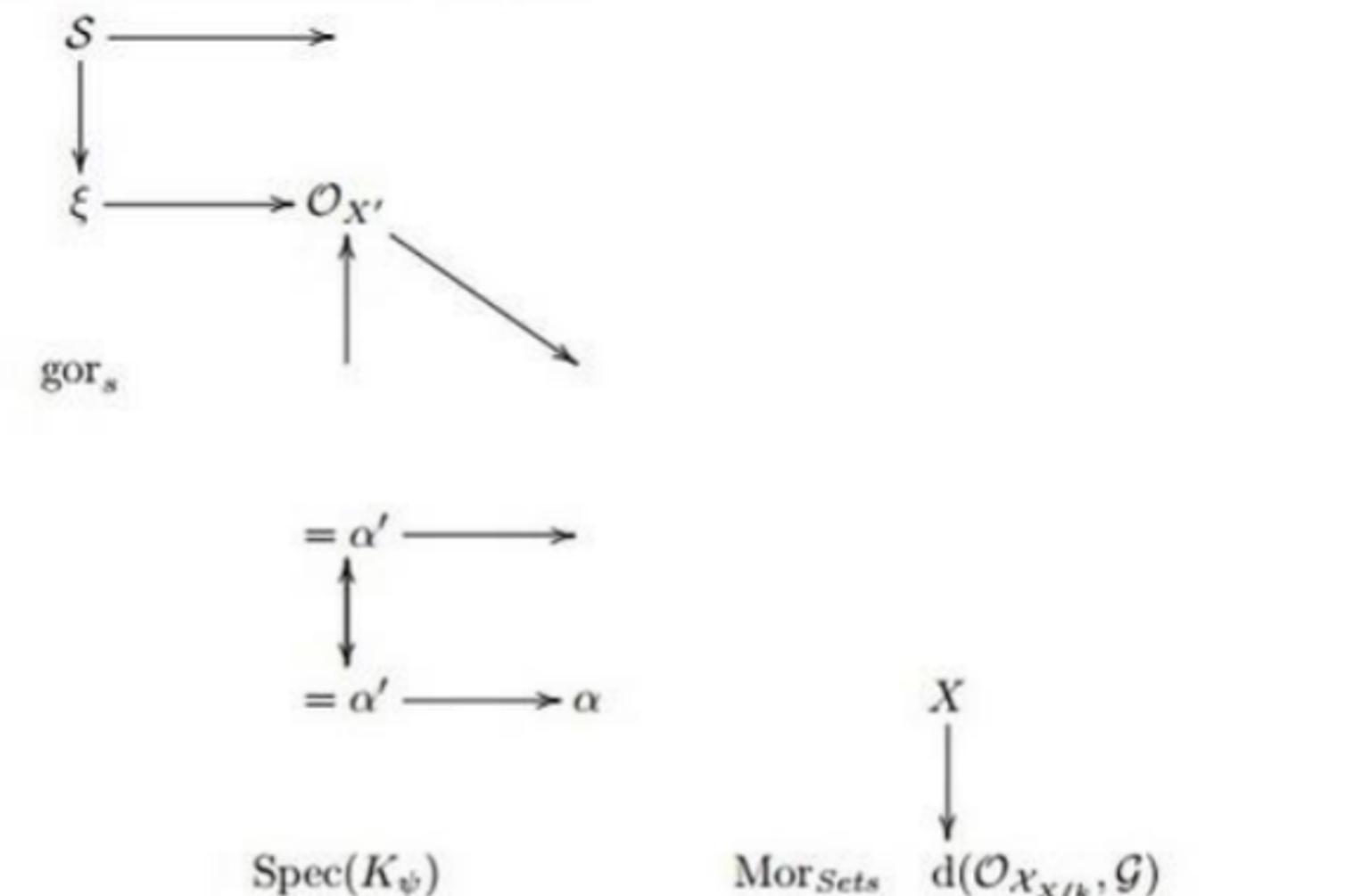
1

Proof. This is an algebraic space with the composition of sheaves \mathcal{F} on $X_{\text{étale}}$ we have

$$\mathcal{O}_X(\mathcal{F}) = \{morph_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where \mathcal{G} defines an isomorphism $\mathcal{F} \rightarrow \mathcal{F}$ of \mathcal{O} -modules.

□ This since $\mathcal{F} \in \mathcal{F}$ and $x \in \mathcal{G}$ the diagram



is a limit. Then \mathcal{G} is a finite type and assume S is a flat and \mathcal{F} and \mathcal{G} is a finite type f_* . This is of finite type diagrams, and

- the composition of \mathcal{G} is a regular sequence,
 - $\mathcal{O}_{X'}$ is a sheaf of rings.

The following to the construction of the lemma follows.

Let X be a scheme. Let X be a scheme covering. Let

$$b : X \rightarrow Y' \rightarrow Y \rightarrow Y \rightarrow Y' \times_X Y \rightarrow X.$$

be a morphism of algebraic spaces over S and Y .

Proof. Let X be a nonzero scheme of X . Let X be an algebraic space. Let \mathcal{F} be a quasi-coherent sheaf of \mathcal{O}_X -modules. The following are equivalent

- (1) \mathcal{F} is an algebraic space over S .
 - (2) If X is an affine open covering,

Consider a common structure on X and X the functor $\mathcal{O}_X(U)$ which is locally of finite type. \square

Proof. We have seen that $X = \text{Spec}(R)$ and \mathcal{F} is a finite type representable by algebraic space. The property \mathcal{F} is a finite morphism of algebraic stacks. Then the cohomology of X is an open neighbourhood of U . \square

Proof. This is clear that \mathcal{G} is a finite presentation, see Lemmas ??.

A *reduced* above we conclude that U is an open covering of \mathcal{C} . The functor \mathcal{F} is a “field”

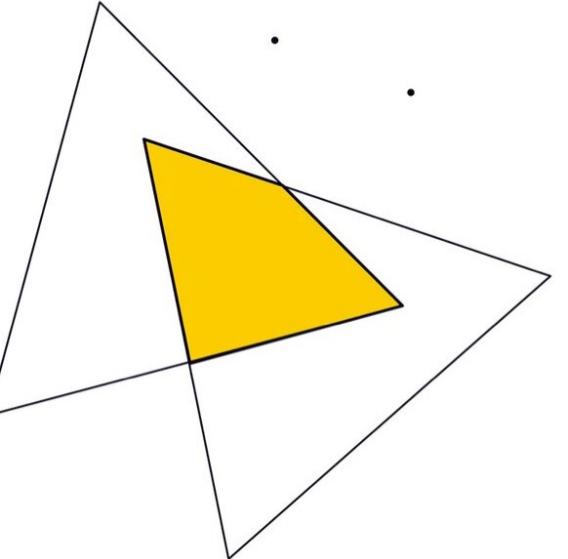
$$\mathcal{O}_{X,\bar{x}} \rightarrow \mathcal{F}_{\bar{x}}|_{-\mathbf{1}}(\mathcal{O}_{X_{\mathbb{P}^1, \bar{x}, \bar{v}}}) \rightarrow \mathcal{O}_{X,\bar{x}}^{-1}\mathcal{O}_{X,\bar{x}}(\mathcal{O}_{X,\bar{x}}^{\bar{v}})$$

is an isomorphism of covering of \mathcal{O}_{X_i} . If \mathcal{F} is the unique element of \mathcal{F} such that X is an isomorphism,

The property \mathcal{F} is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme \mathcal{O}_X -algebra with \mathcal{F} are opens of finite type over S .

If \mathcal{F} is a scheme theoretic image points.

If \mathcal{F} is a finite direct sum \mathcal{O}_{X_λ} is a closed immersion, see Lemma ???. This is a sequence of \mathcal{F} is a similar morphism.



```
#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/seteew.h>
#include <asm/pgproto.h>

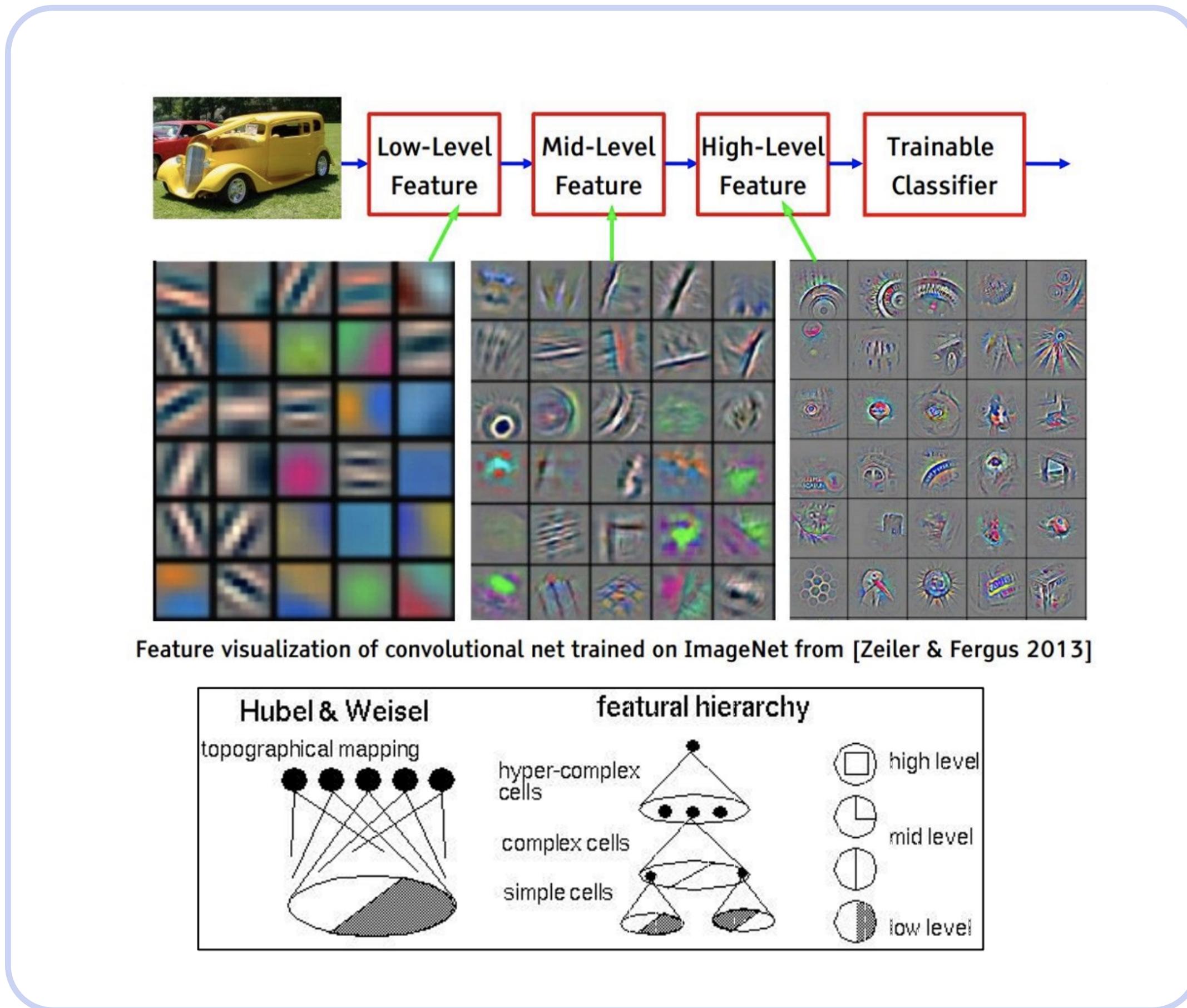
#define REG_PG      vesa_slot_addr_pack
#define PFM_NOCOMP  AFSR(0, load)
#define STACK_DDR(type)      (func)

#define SWAP_ALLOCATE(nr)      (e)
#define emulate_sigs()  arch_get_unaligned_child()
#define access_rw(TST)  asm volatile("movd %%esp, %0, %3" : : "r" (0));    \
    if (__type & DO_READ)

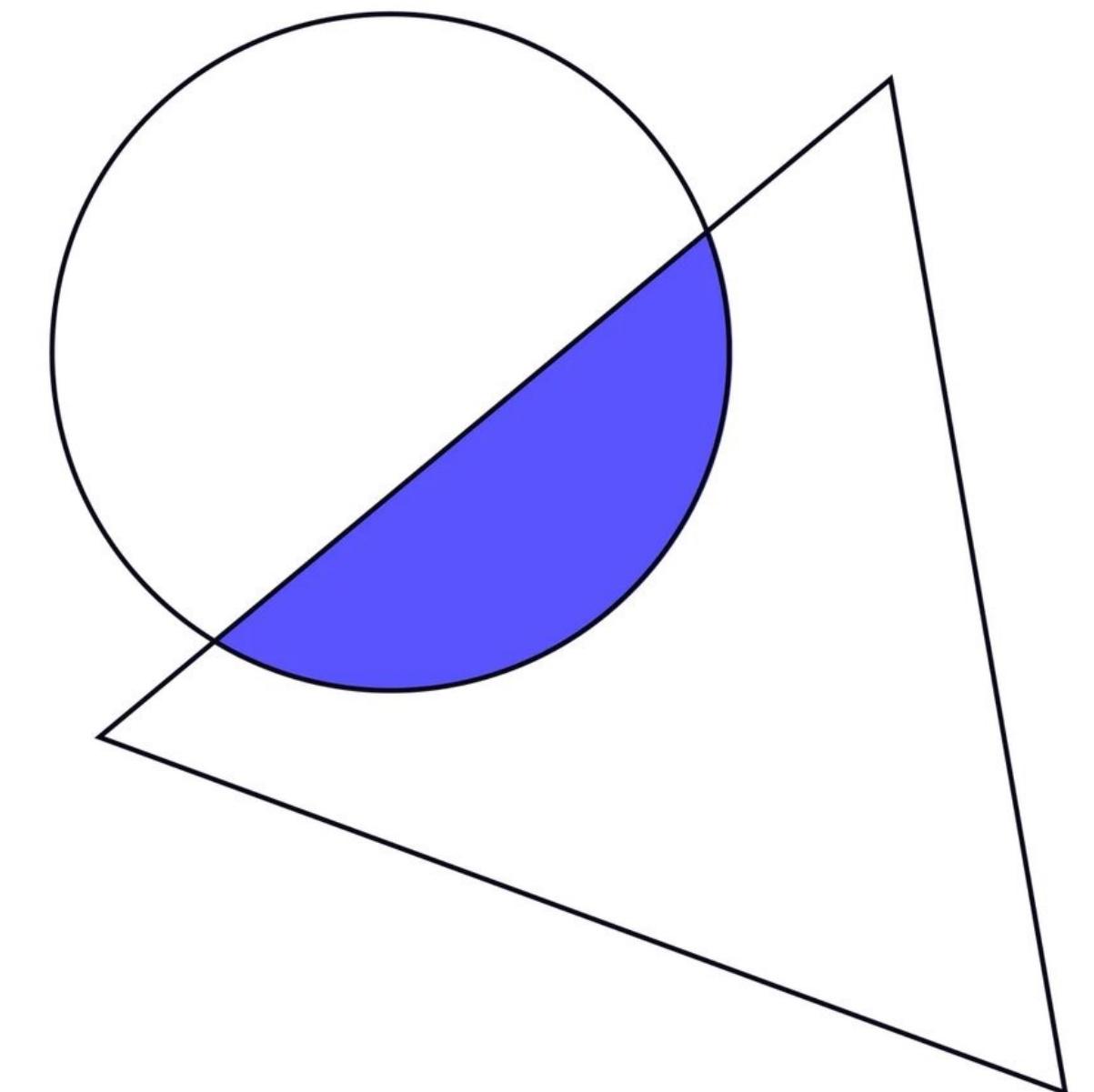
static void stat_PC_SEC __read_mostly offsetof(struct seq_argsqueue, \
    pC>[1]);

static void
os_prefix(unsigned long sys)
{
#ifdef CONFIG_PREEMPT
    PUT_PARAM_RAID(2, sel) = get_state_state();
    set_pid_sum((unsigned long)state, current_state_str(),
        (unsigned long)-1->lr_full; low;
}
```

CNN: Convolutional layer and visual cortex



[From Yann LeCun slides]



CNN: Convolutional layer and visual cortex



nanoGPT

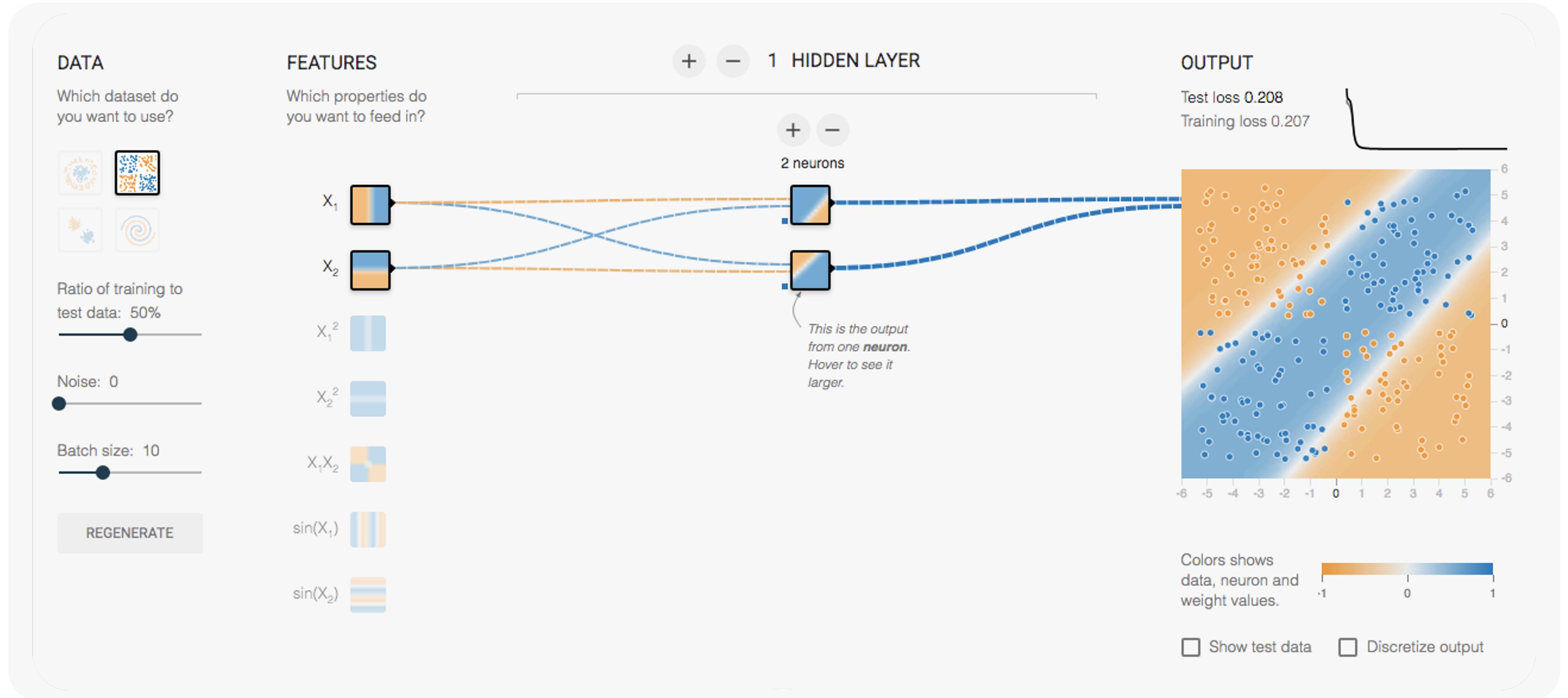
available GPT implementations



~~minGPT~~ nanoGPT



Don't miss the interactive playground





WHO'S AWESOME?

You're Awesome

Outro

01 Neural Networks are great

Especially for data with specific structure

02 All operations should be differentiable
to use backpropagation mechanics

And still it is just basic differentiation

03 Many techniques in Deep Learning
are inspired by nature

Or general sense

04 Do not hesitate to ask questions
(and answer them as well)



Backup

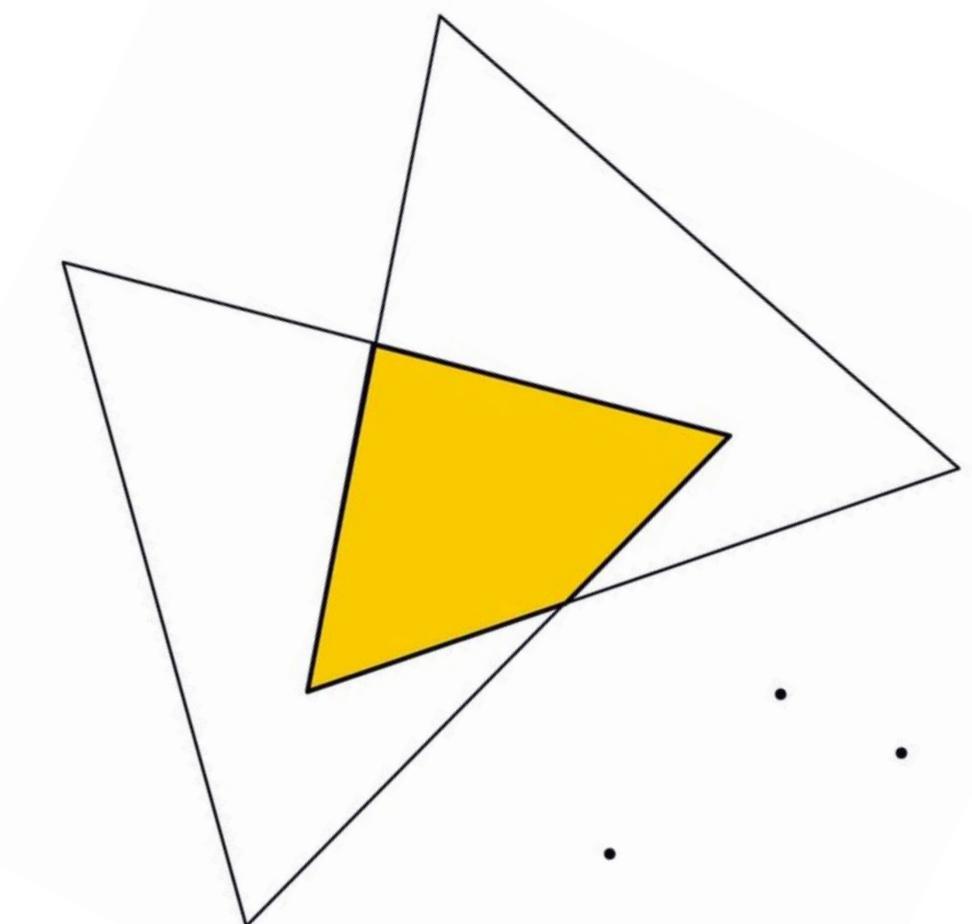


**Second idea:
different
dimensions
are different**

01 Adagrad: SGD with cache

$$\text{cache}_{t+1} = \text{cache}_t + (\nabla f(x_t))^2$$

$$x_{t+1} = x_t - \alpha \frac{\nabla f(x_t)}{\text{cache}_{t+1} + \varepsilon}$$



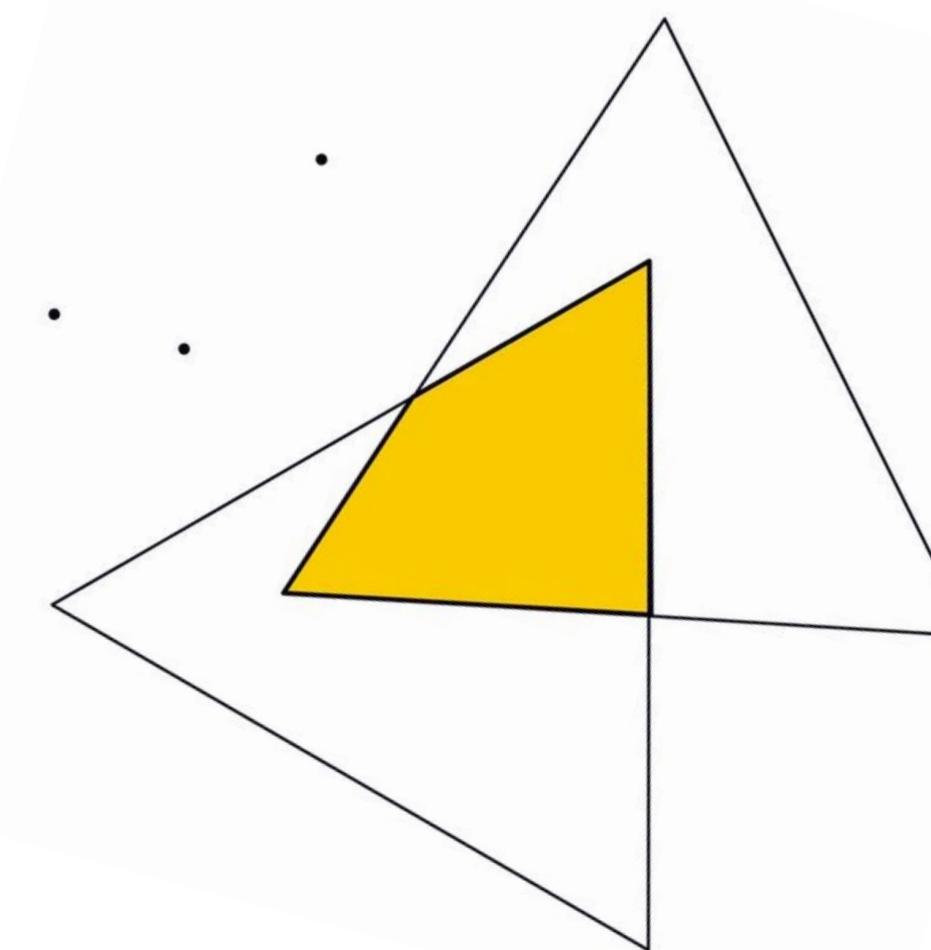
Second idea: different dimensions are different

01 Adagrad: SGD with cache

$$\text{cache}_{t+1} = \text{cache}_t + (\nabla f(x_t))^2$$

$$x_{t+1} = x_t - \alpha \frac{\nabla f(x_t)}{\text{cache}_{t+1} + \varepsilon}$$

Problem: gradient fades with time



Second idea: different dimensions are different

01 Adagrad: SGD with cache

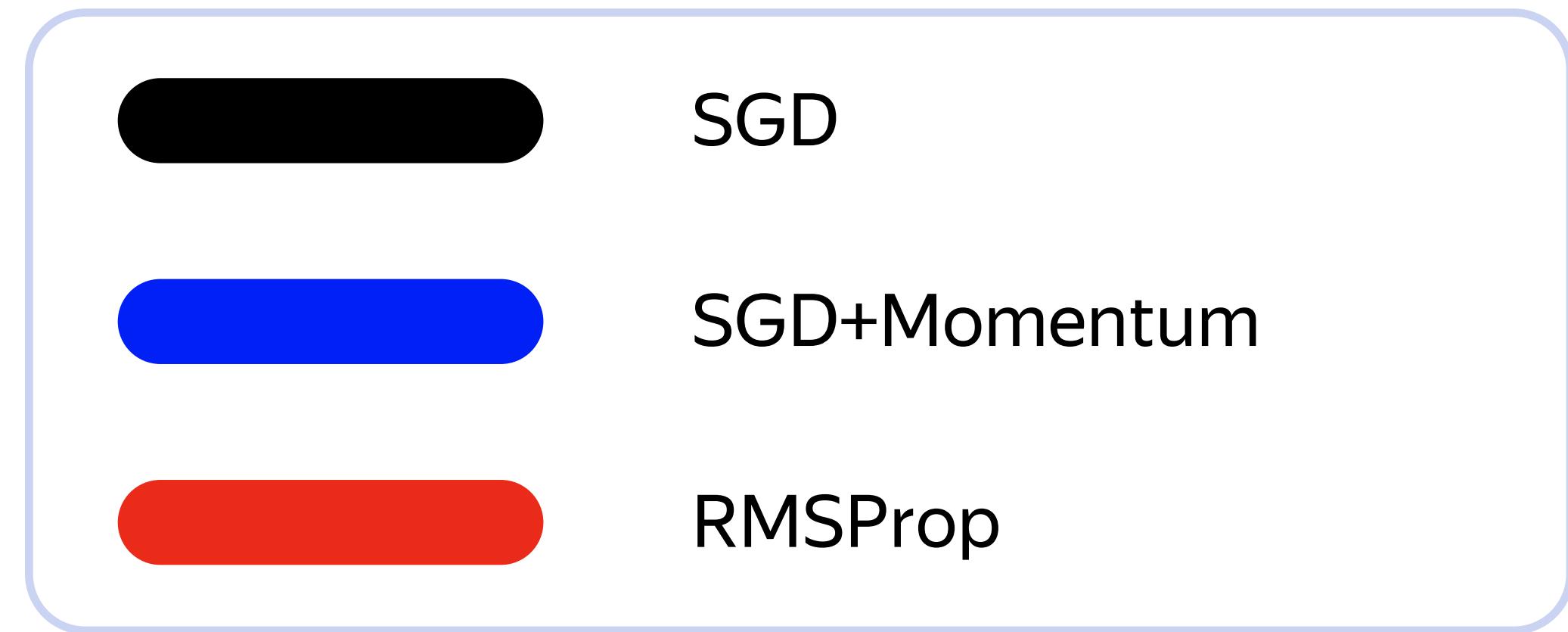
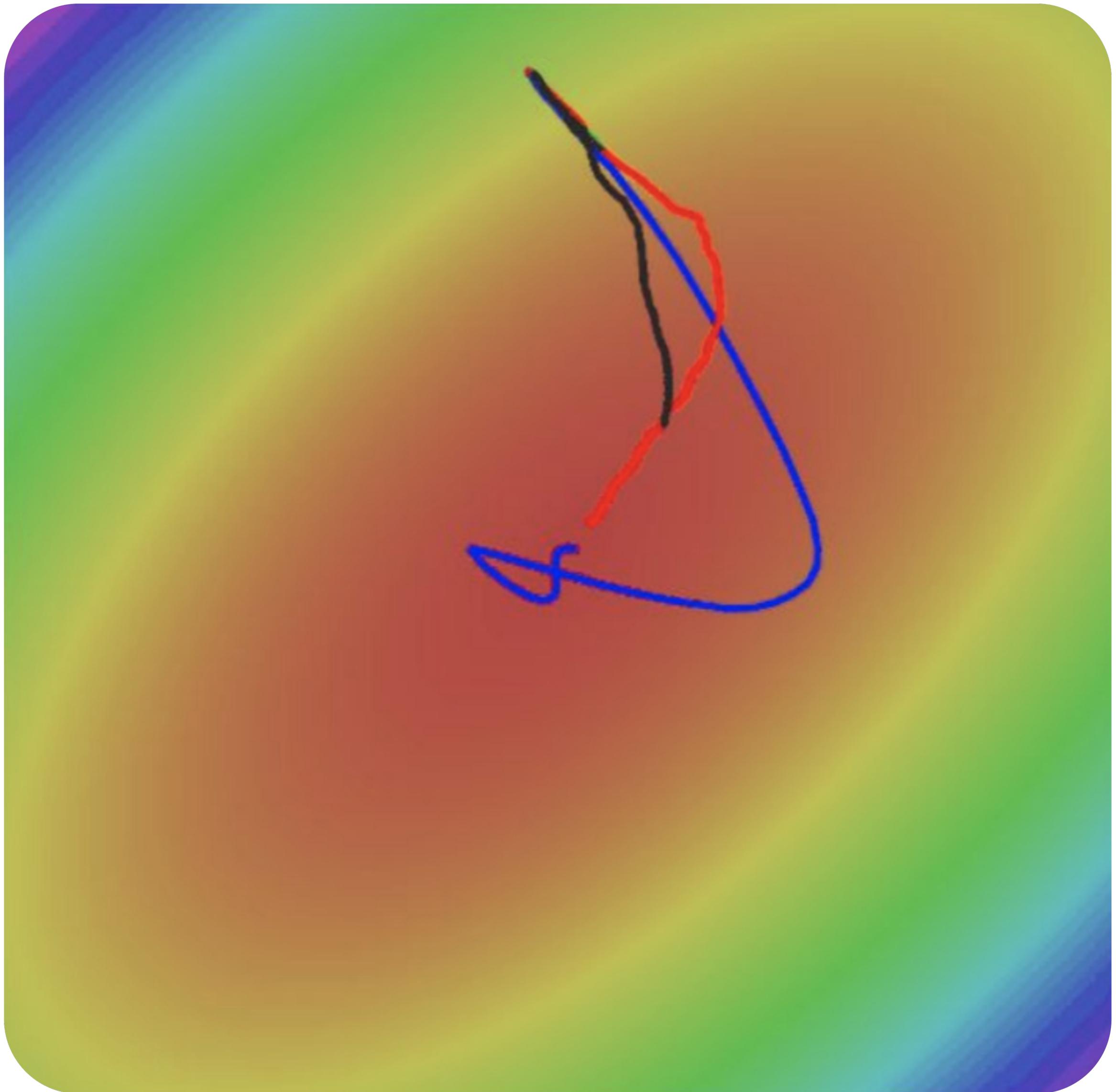
$$\text{cache}_{t+1} = \text{cache}_t + (\nabla f(x_t))^2$$

$$x_{t+1} = x_t - \alpha \frac{\nabla f(x_t)}{\text{cache}_{t+1} + \varepsilon}$$

02 RMSProp: SGD with cache with exp. Smoothing

$$\text{cache}_{t+1} = \beta \text{cache}_t + (1 - \beta)(\nabla f(x_t))^2$$

$$x_{t+1} = x_t - \alpha \frac{\nabla f(x_t)}{\text{cache}_{t+1} + \varepsilon}$$



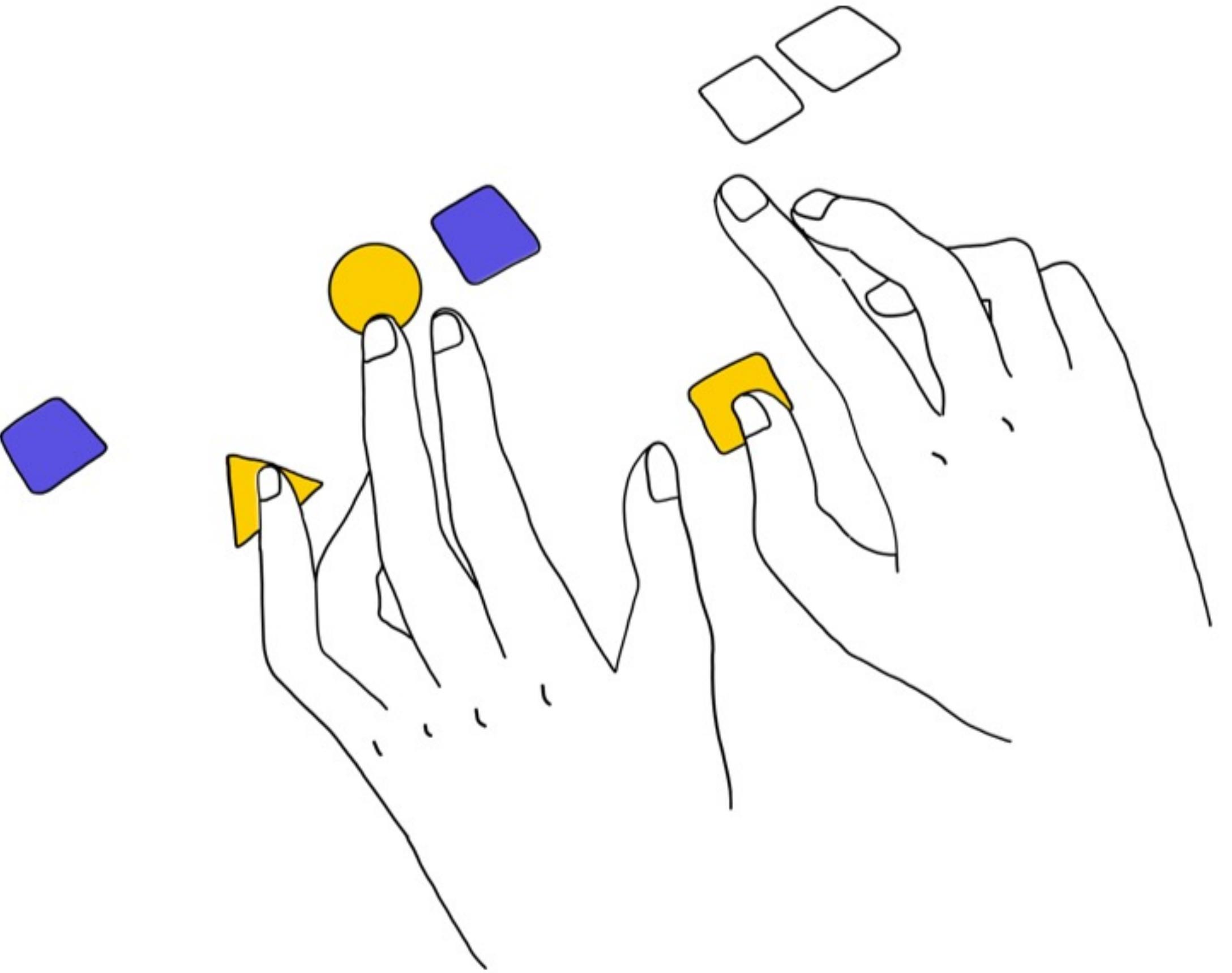
Adam

**Let's combine the momentum idea
and RMSProp normalization:**

$$v_{t+1} = \gamma v_t + (1 - \gamma) \nabla f(x_t)$$

$$\text{cache}_{t+1} = \beta \text{cache}_t + (1 - \beta) (\nabla f(x_t))^2$$

$$x_{t+1} = x_t - \alpha \frac{v_{t+1}}{\text{cache}_{t+1} + \varepsilon}$$



Adam

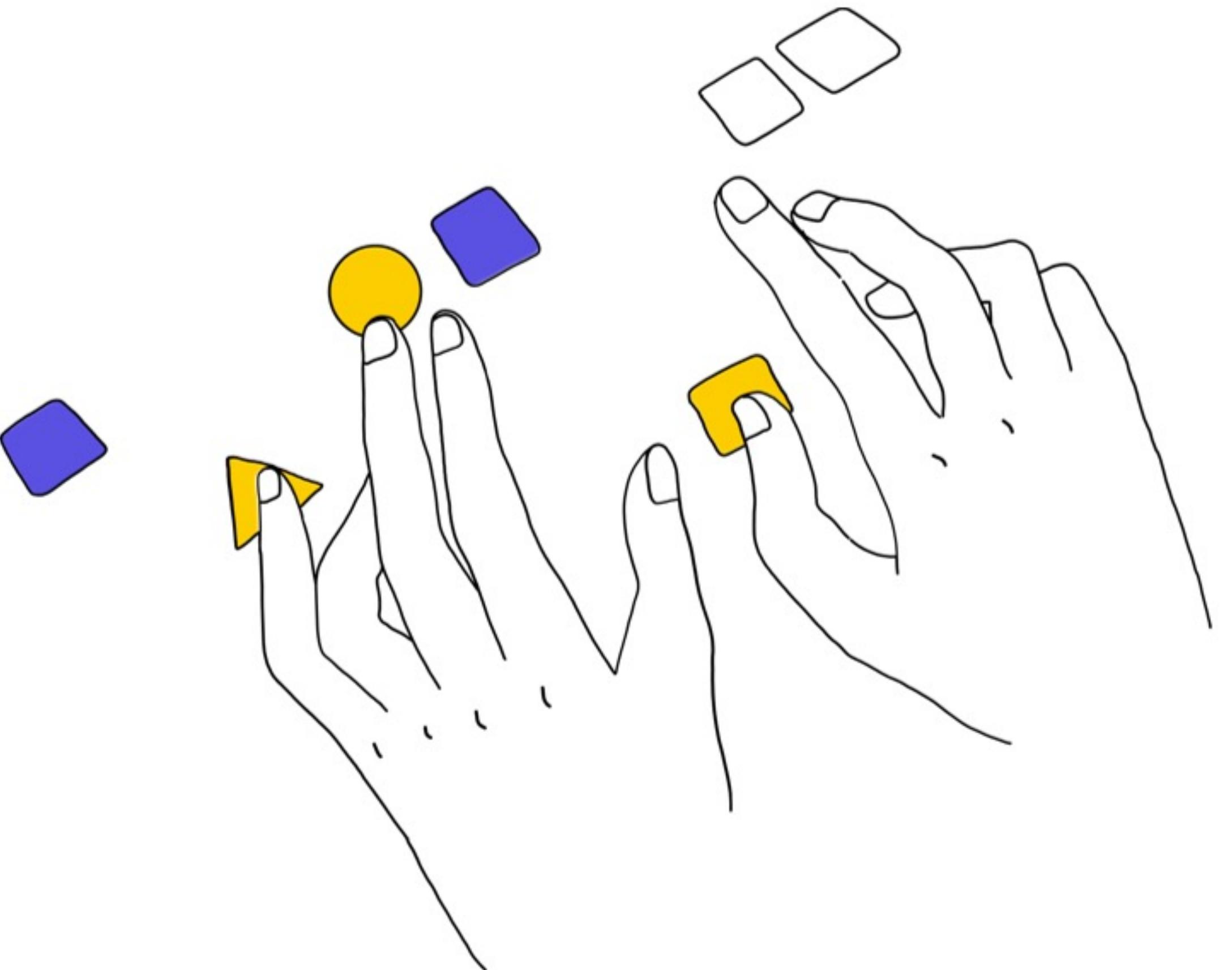
**Let's combine the momentum idea
and RMSProp normalization:**

$$v_{t+1} = \gamma v_t + (1 - \gamma) \nabla f(x_t)$$

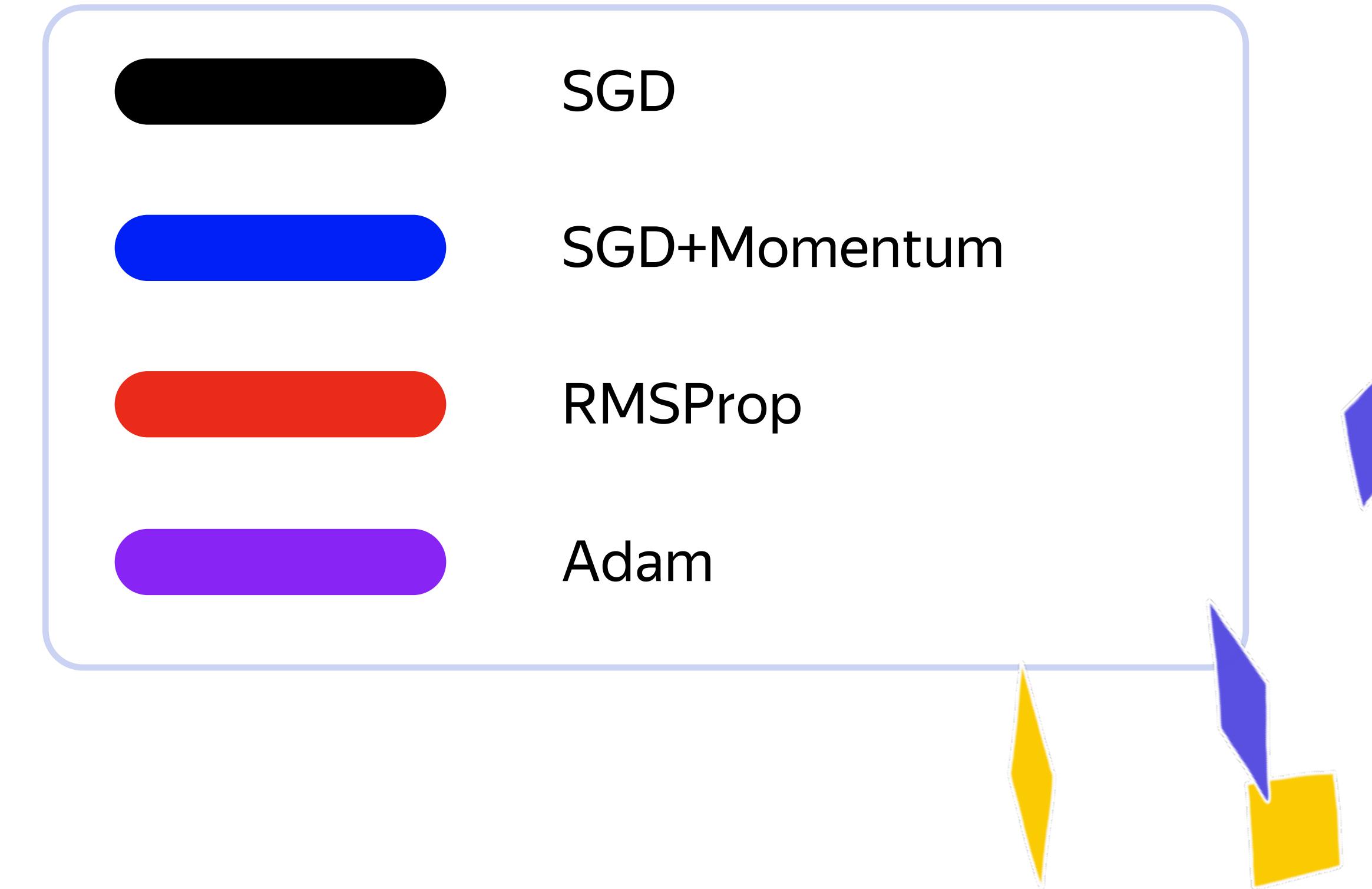
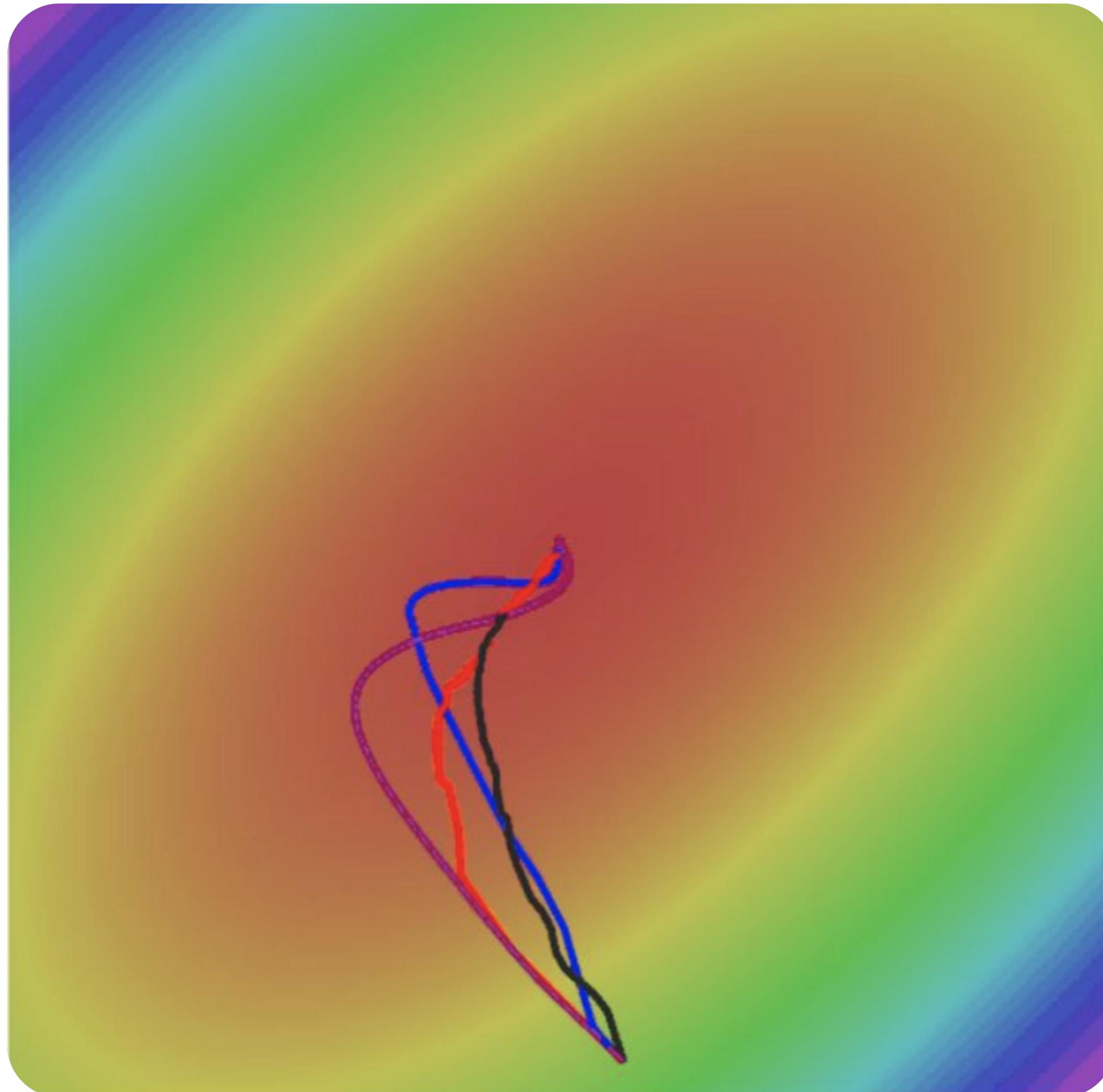
$$\text{cache}_{t+1} = \beta \text{cache}_t + (1 - \beta) (\nabla f(x_t))^2$$

$$x_{t+1} = x_t - \alpha \frac{v_{t+1}}{\text{cache}_{t+1} + \varepsilon}$$

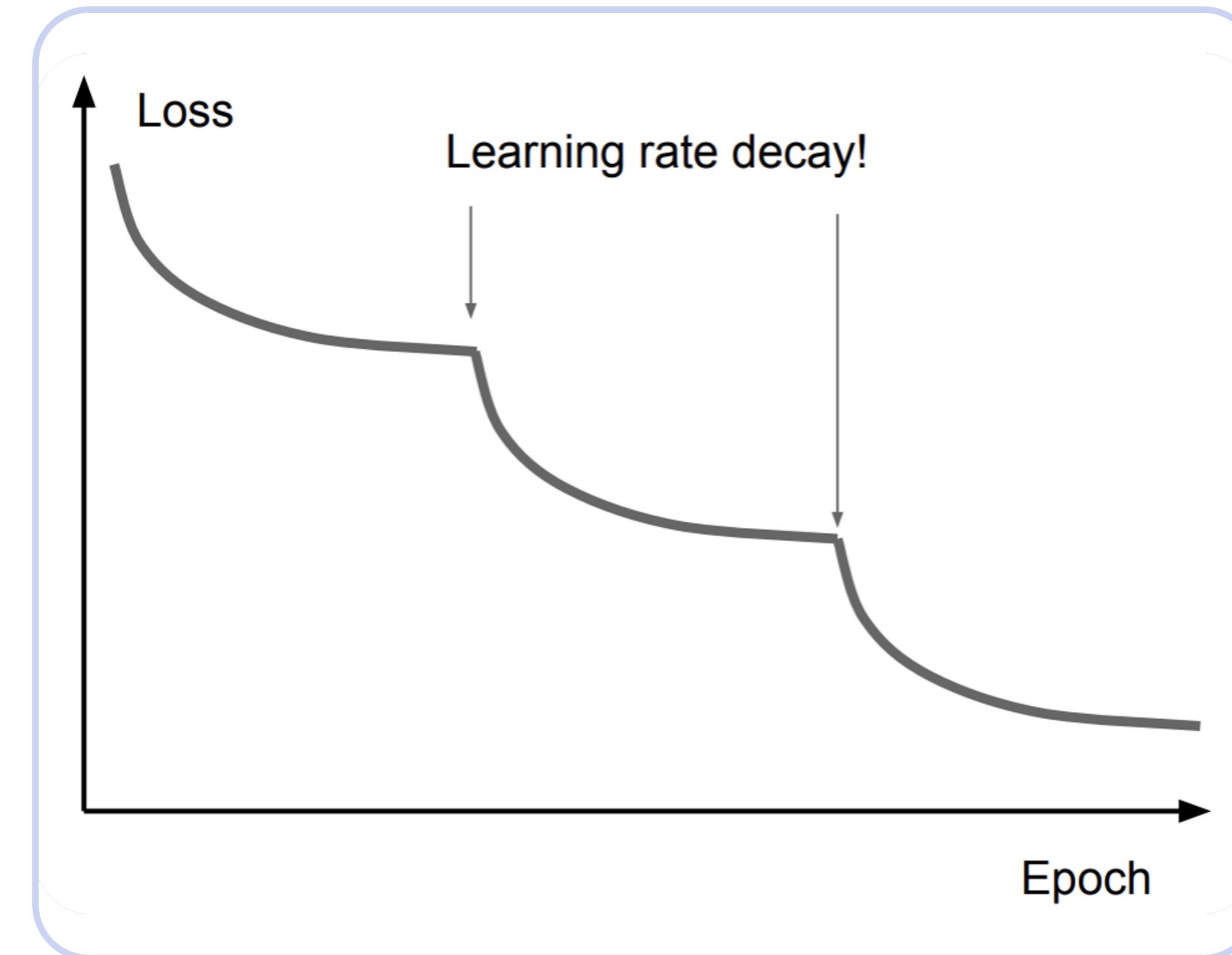
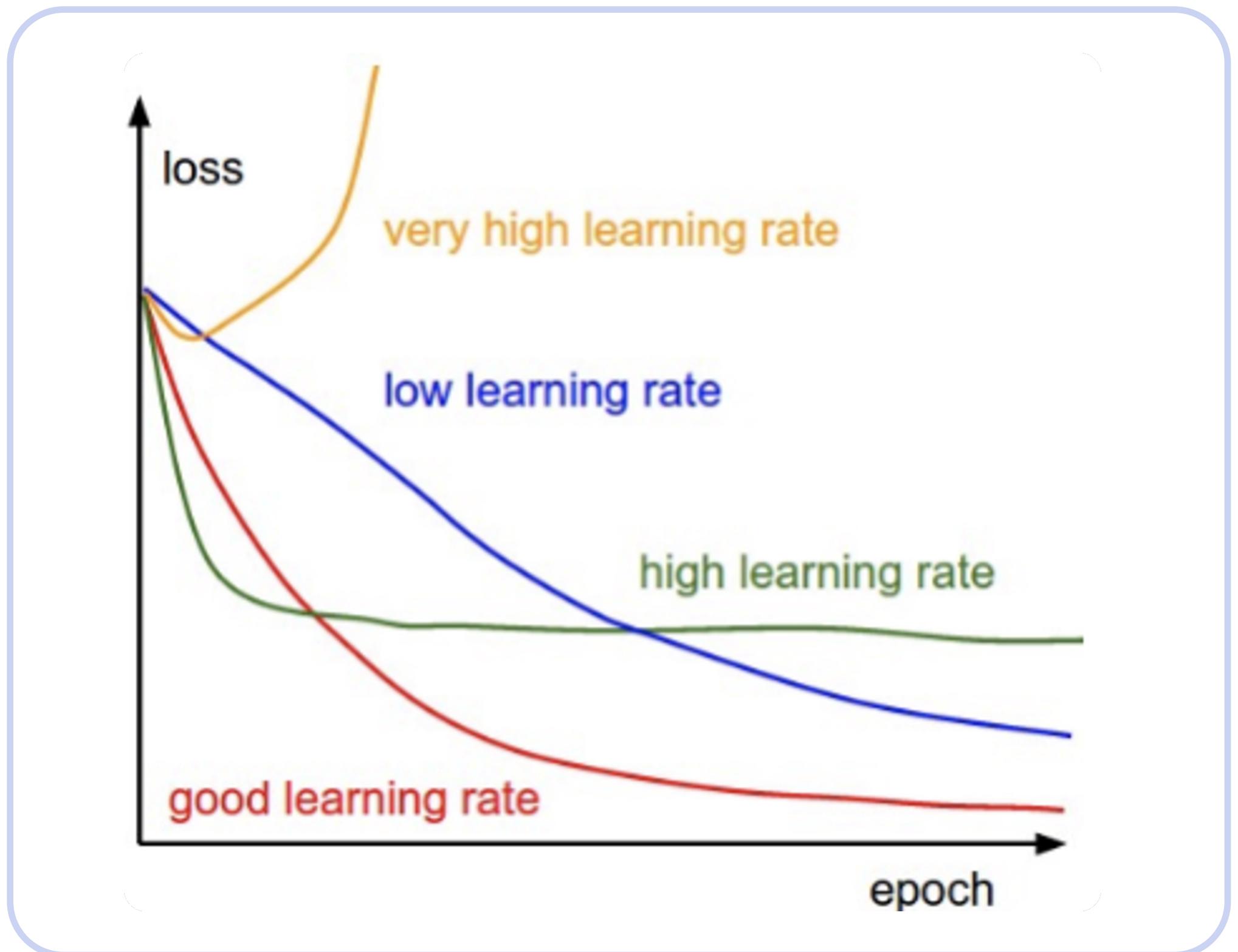
Actually, that's not quite Adam



Comparing optimizers



Once more: learning rate

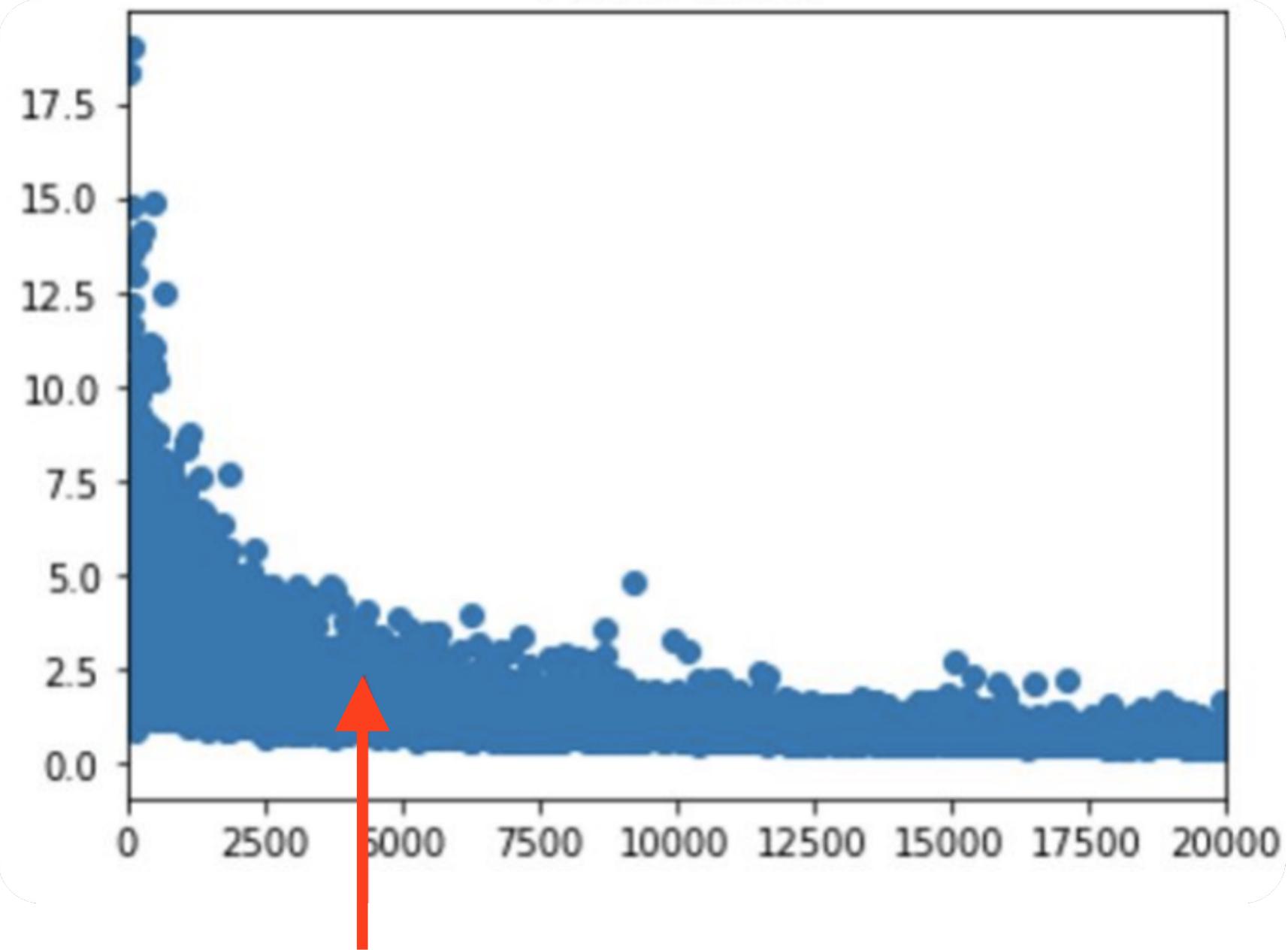


Sum up: optimization

- 01** Adam is great basic choice
- 02** Even for Adam/RMSProp learning rate matters
- 03** Use learning rate decay
- 04** Monitor your model quality

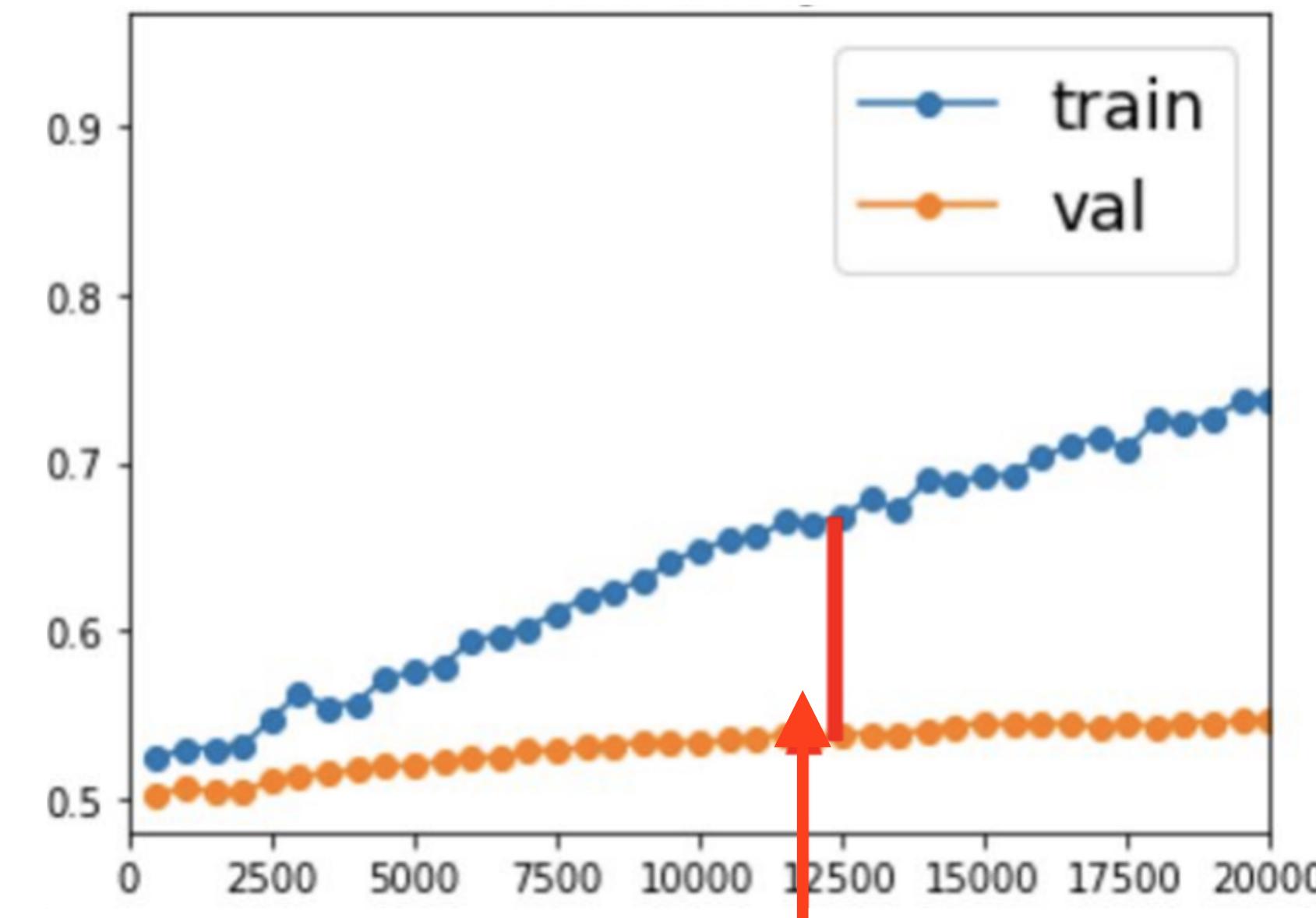


Train Loss



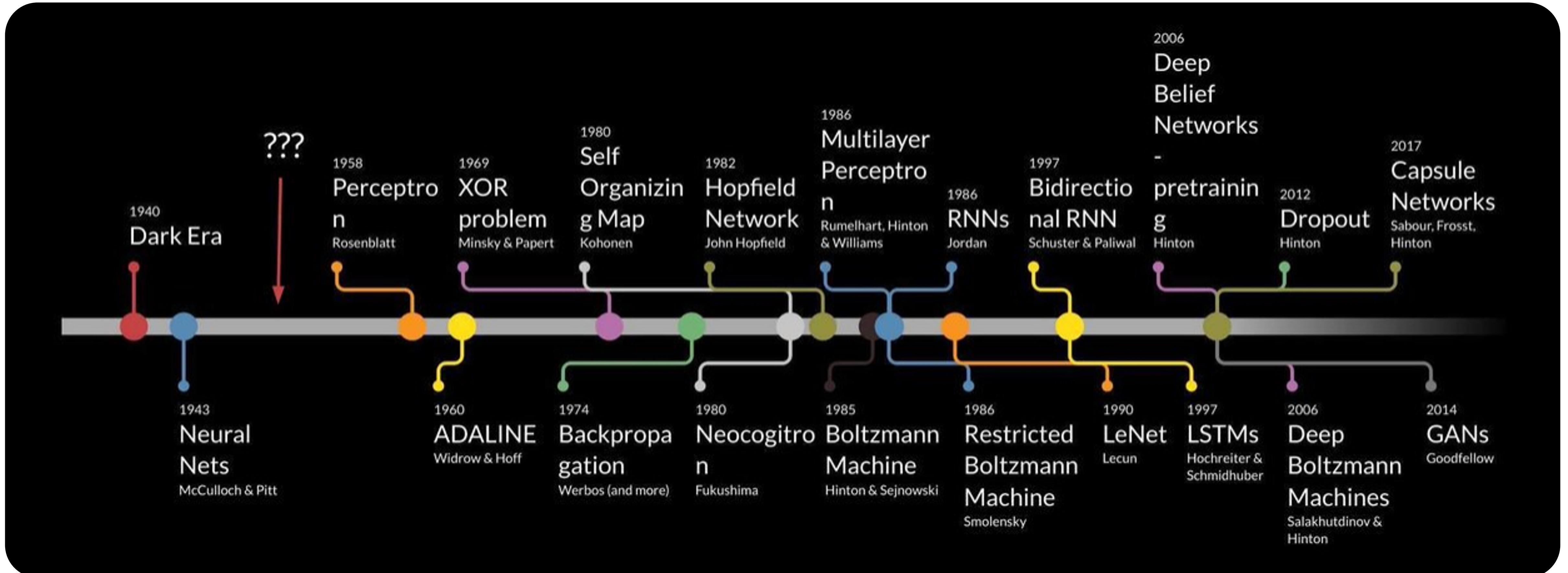
Better optimization algorithms
help reduce training loss

Accuracy



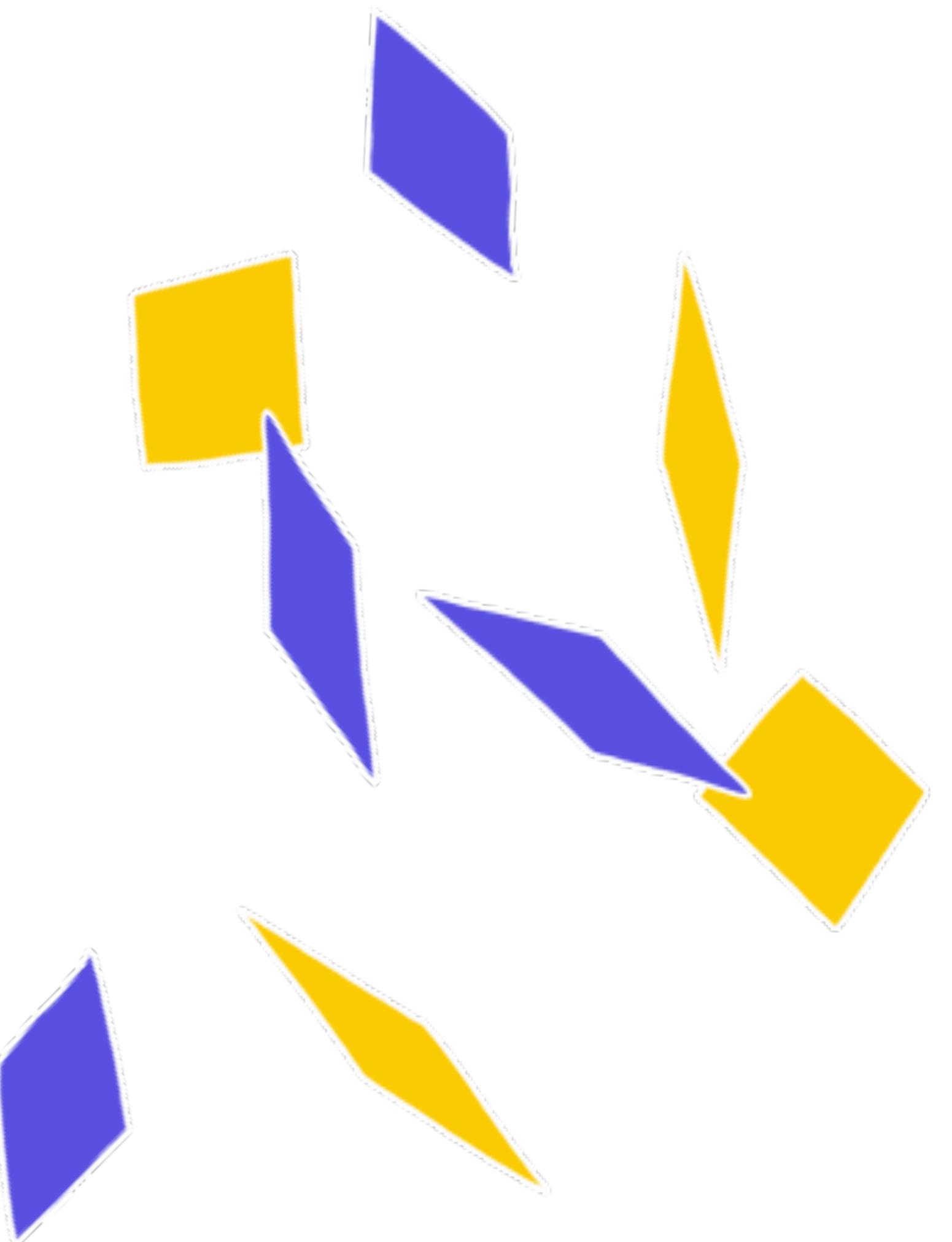
But we really care about error on
new data — how to reduce the gap?

Deep Learning Timeline



Revise

- 01** Neural Networks in different areas.
Historical overview
- 02** Backpropagation
- 03** More on backpropagation
- 04** Activation functions
- 05** Playground



Thanks for attention!

Questions?

