

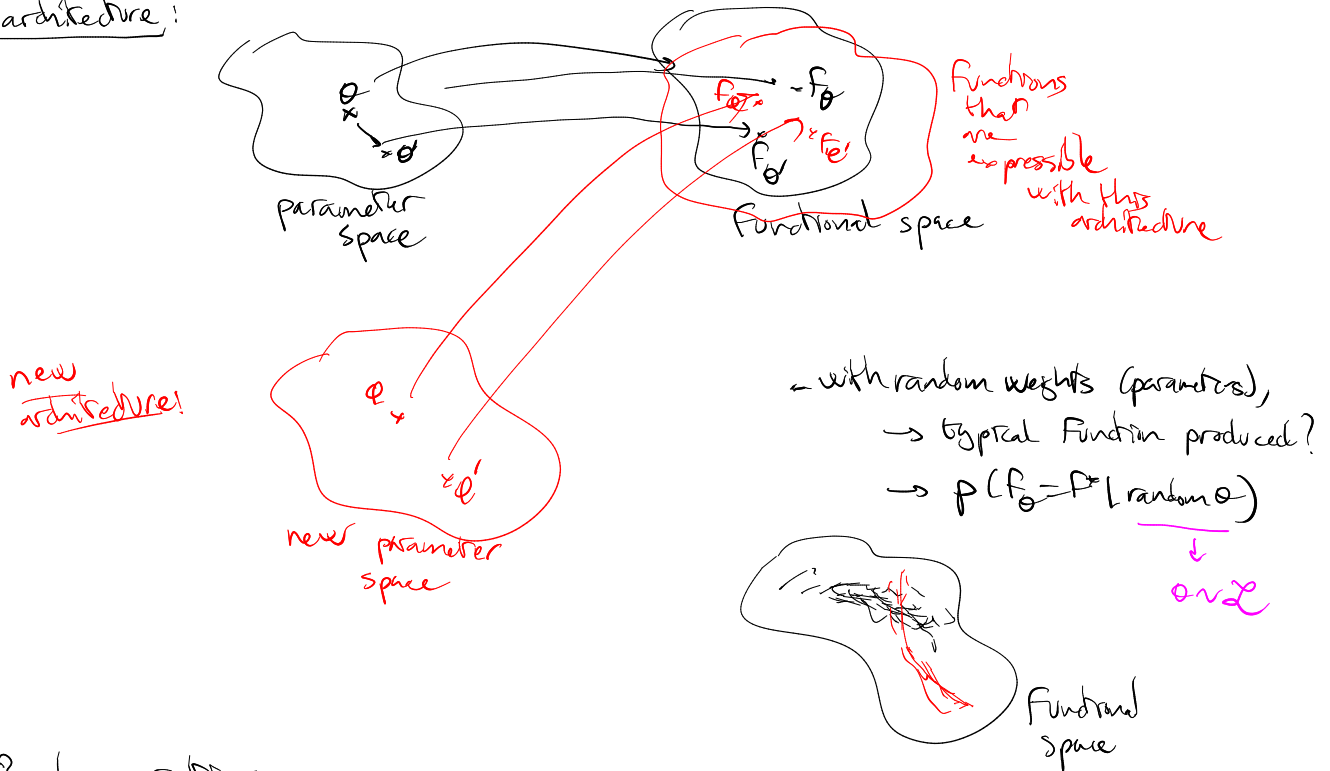
Architectures

I Architectures = prior on the function space

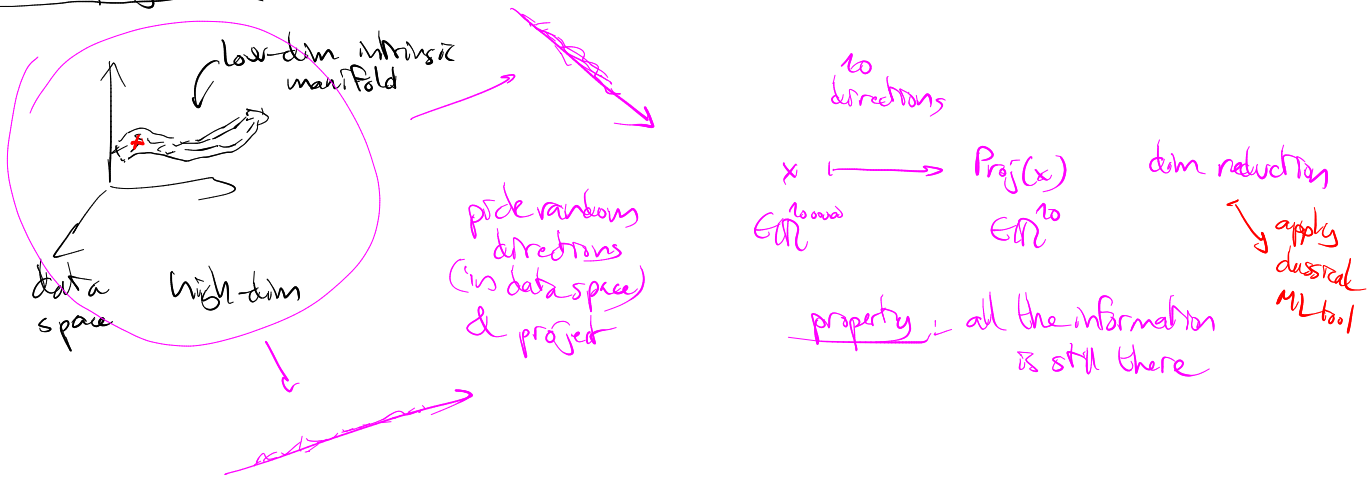
change of paradigm =

- classical ML: design features by hand
- deep DL: meta-design of features \rightarrow design architectures able to express
Learn features

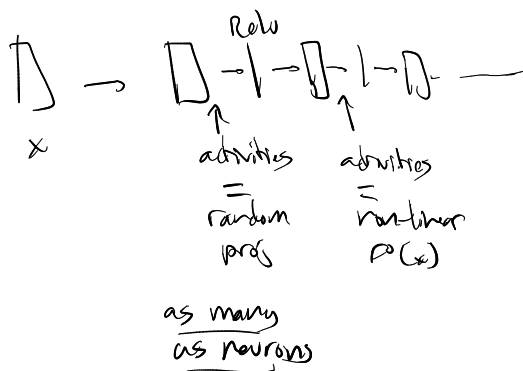
An architecture:



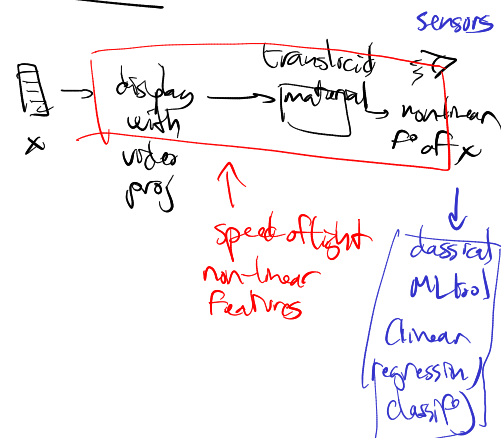
Random projections



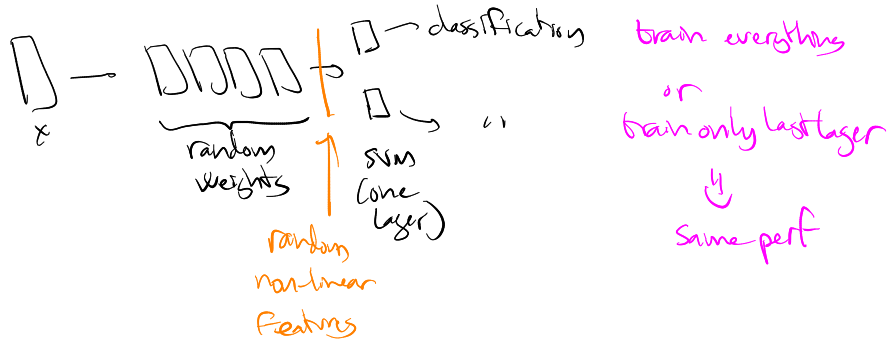
Random weights: \approx non-linear projection



Light on



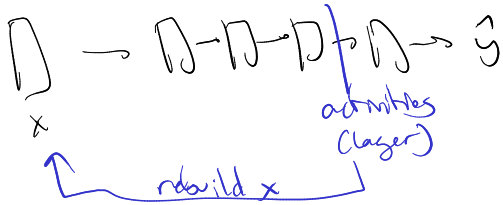
In some cases most of the performance is due to the architecture and not to the training



(Extreme ML)

~~works~~ work only for small architectures

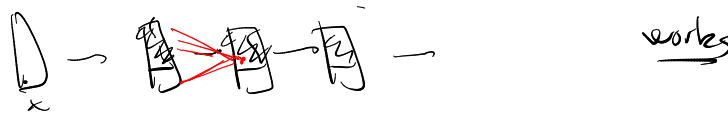
because the architecture was tuned beforehand



if random weights: good reconstruction of x → keep more information
 if trained: kind of work (not so good) → keep only information relevant for the task

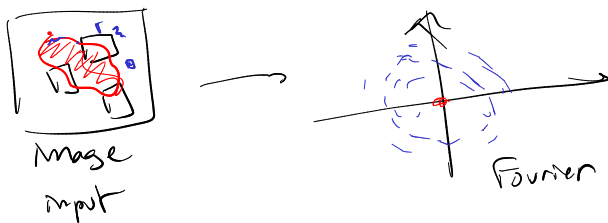
For deeper networks:

- in each layer: $\left\{ \begin{array}{l} 90\% \text{ of neurons: randomly initialized (and fixed)} \\ 10\% \text{ --- --- are trained} \end{array} \right.$



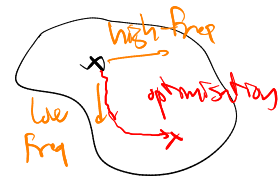
Architecture bias

CNN: move to Fourier space



shown to learn low-frequencies first

Looking at "the big picture" before at the details



add layers $M^t M^t M \dots$ no gain in expressive power

DD

But the distribution over functional space is different

$$M \sim \mathcal{L} \quad M_{ij} \sim \mathcal{N}(0,1) \rightarrow P(F)$$

$$M^t M$$

Random initialize

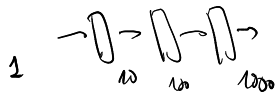
→ random: but which law?

Exploding / Vanishing gradients



activation amplitude \approx input amplitude
 $\|x'\| \approx \gamma \|x\|$

$\gamma = 10$



factor γ^L (where L is #layers)
 if 10 layers
 20
 100
 \Rightarrow issues
 $F \approx 100$

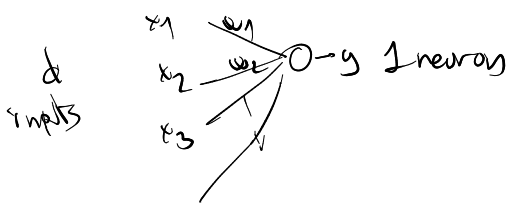
$\gamma = 0.1 \Rightarrow F \approx 0$

mitigation law
 st. $\gamma \approx 1$

same issue with gradients

$\nabla_{\theta} \approx \pm \infty \Rightarrow$ not possible to train
 $\nabla_{\theta} \approx 0$

Xavier Glorot's initialization



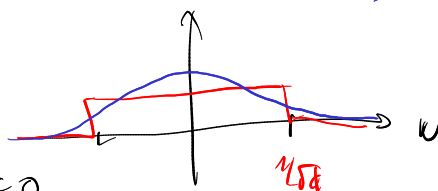
inputs $x_i \sim \mathcal{M}(0, 1)$ or just mean=0 variance=1
 $w_i \sim \mathcal{L}$

goal: find \mathcal{L} st. variance(output) = 1

solution: $w \sim \mathcal{U}\left[\frac{-1}{\sqrt{d}}, \frac{+1}{\sqrt{d}}\right]$ or $w \sim \mathcal{M}(0, \sigma^2 = \frac{1}{d})$

Proof: $y = \sum_i w_i x_i$

mean: $\mathbb{E}_x \mathbb{E}_w [y] = \mathbb{E}_w \left[\sum_i w_i \mathbb{E}[x_i] \right] = 0$



variance: $\mathbb{E}_x \mathbb{E}_w [y^2] = \mathbb{E}_w \left[\mathbb{E}_x \left[\sum_i w_i^2 x_i^2 \right] \right] = \mathbb{E}_w \left[\sum_i w_i^2 \underbrace{\mathbb{E}[x_i^2]}_1 \right]$
 $= d \mathbb{E}_w [w^2]$

With activation function:

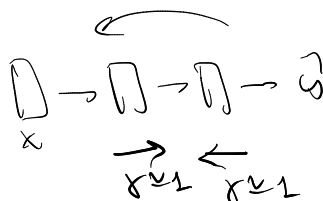


\Rightarrow activation sp-specific correction factor



Biases: initialized to 0

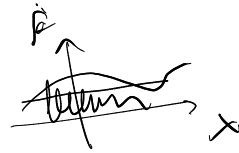
Other approaches: Kaiming He et al



for forward pass
 for backprop
 mitigate: $\frac{1}{\sqrt{\text{#inputs}}}$ and $\frac{1}{\sqrt{\text{#outputs}}}$

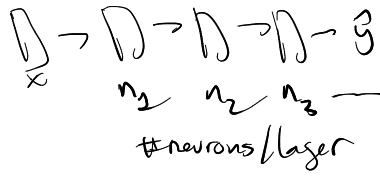
[Paris Ann 2018] Jacobian properties

$J \approx \frac{dF}{dx}$ well behaved at initialize



regularizer $\| \frac{dF}{dx} \|^2$

maximize dim outputs & dim inputs



upon initialize
 $var\left(\frac{dF}{dx}\right) = \frac{1}{n_0}$

$var\left(\left(\frac{dF}{dx}\right)^2\right) = E\left[\left(\frac{dF}{dx}\right)^4\right]$

$\alpha e^{\sum \frac{1}{n_l}} \leq \dots \leq \beta e^{\sum \frac{1}{n_l}}$

→ what matters! $\sum \frac{1}{n_l}$ as small as possible

→ better → as many neurons as possible in each layer

if on a budget: consider same-width layers (better than thin ones)



→ avoid thin layers (unless necessary, e.g. auto-encoder)

Design easy to train architectures

- deeper architectures: more difficult to train



parameter update?

$\frac{\partial L}{\partial G} \propto \frac{\partial L}{\partial x_1} \frac{\partial x_1}{\partial x_2} \frac{\partial x_2}{\partial x_3} \dots \frac{\partial x_n}{\partial w}$

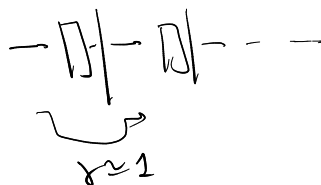
$g = f(x, b, w)$

Information $\frac{\partial L}{\partial G}$ is lost

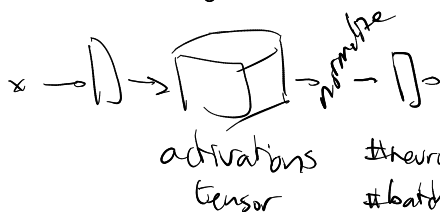
OK for few layers but not too many

- normalization

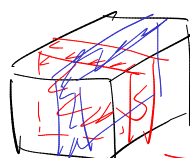
- SEW activation ρ : nice property: keeps the variance



- add scaling layers



activations tensor
 #neurons
 #batch size



normalize each slice of activities

compute (for each slice) mean μ s.d. σ

scaling parameter $\lambda = \frac{A_j - \mu}{\sigma} + b$

learned parameters

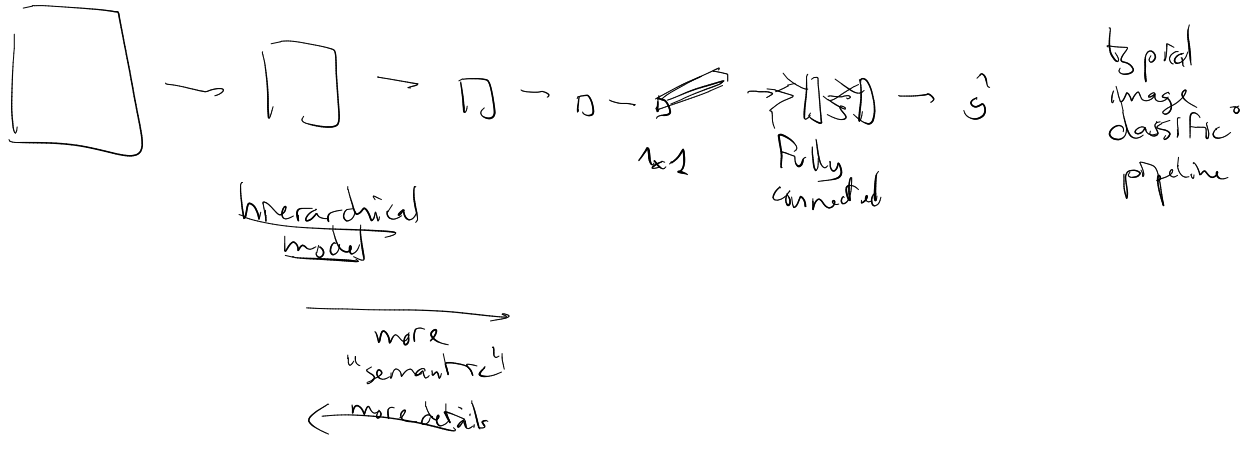
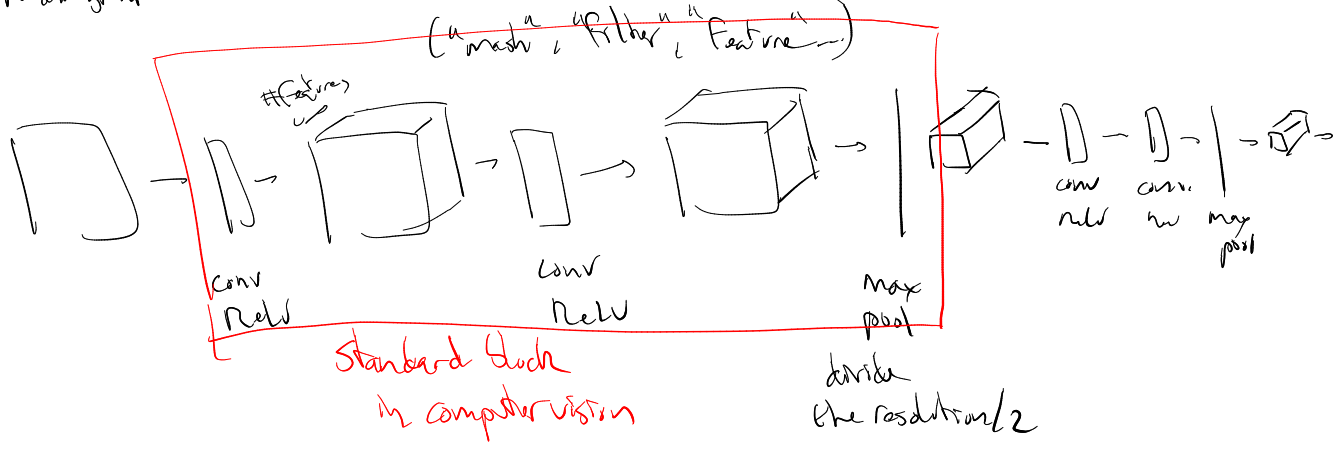
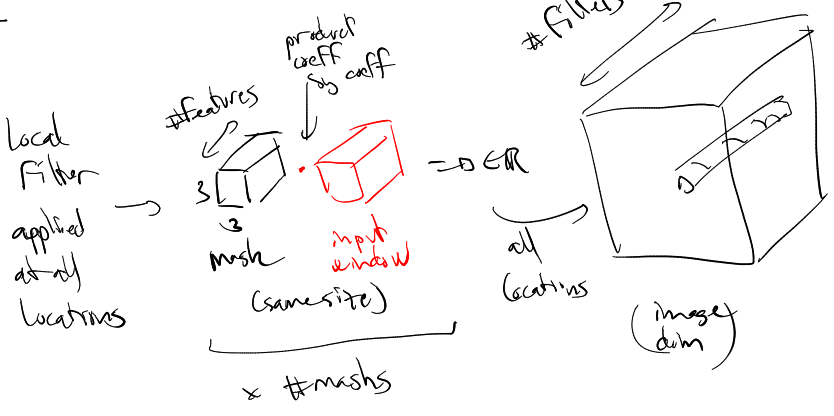
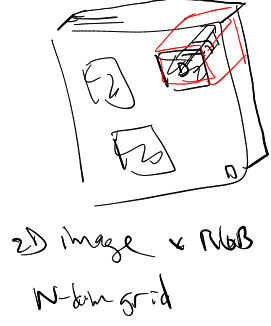
batch-norm
 instance-norm
 layer-norm

Fix up init

#input geometry (width/height/mage)

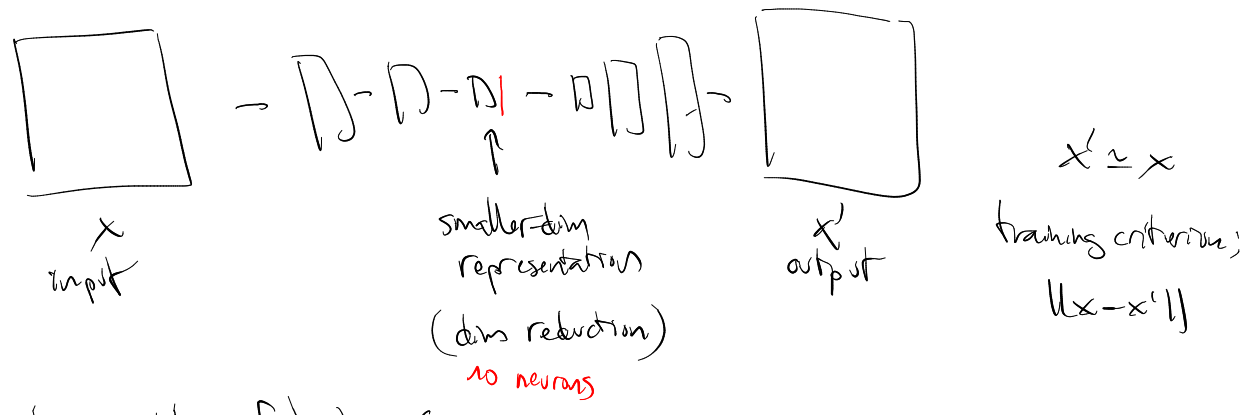
Architecture too

- CNN: convolutions

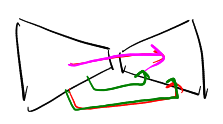


Auto-encoders: unsupervised setup

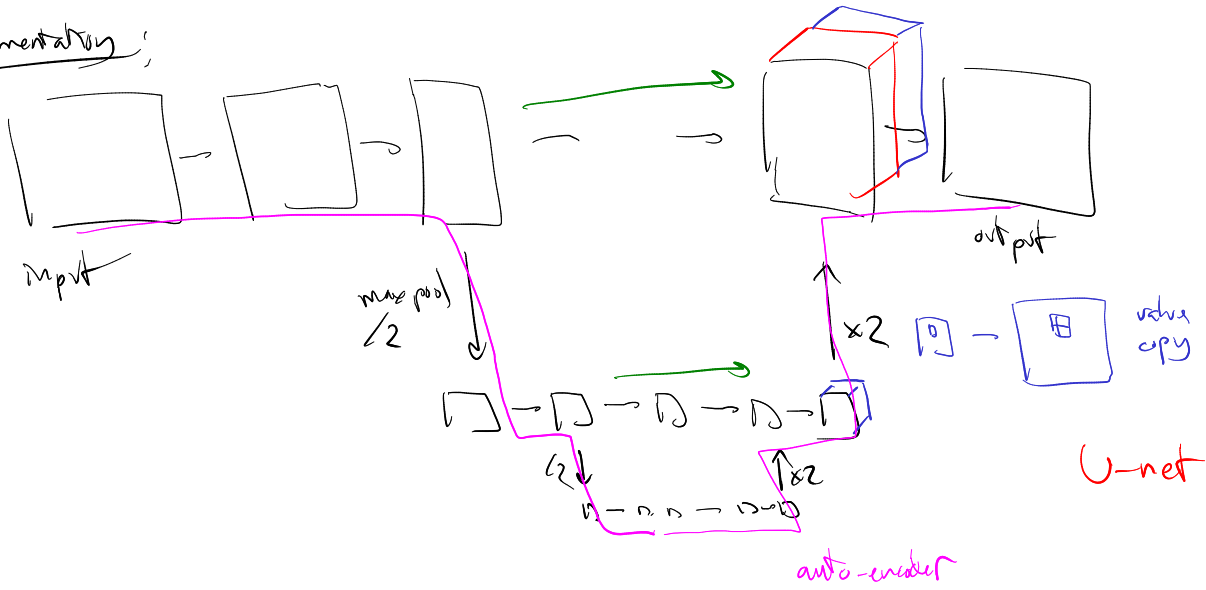
"self-supervised": build a supervised task from labels built from the input.



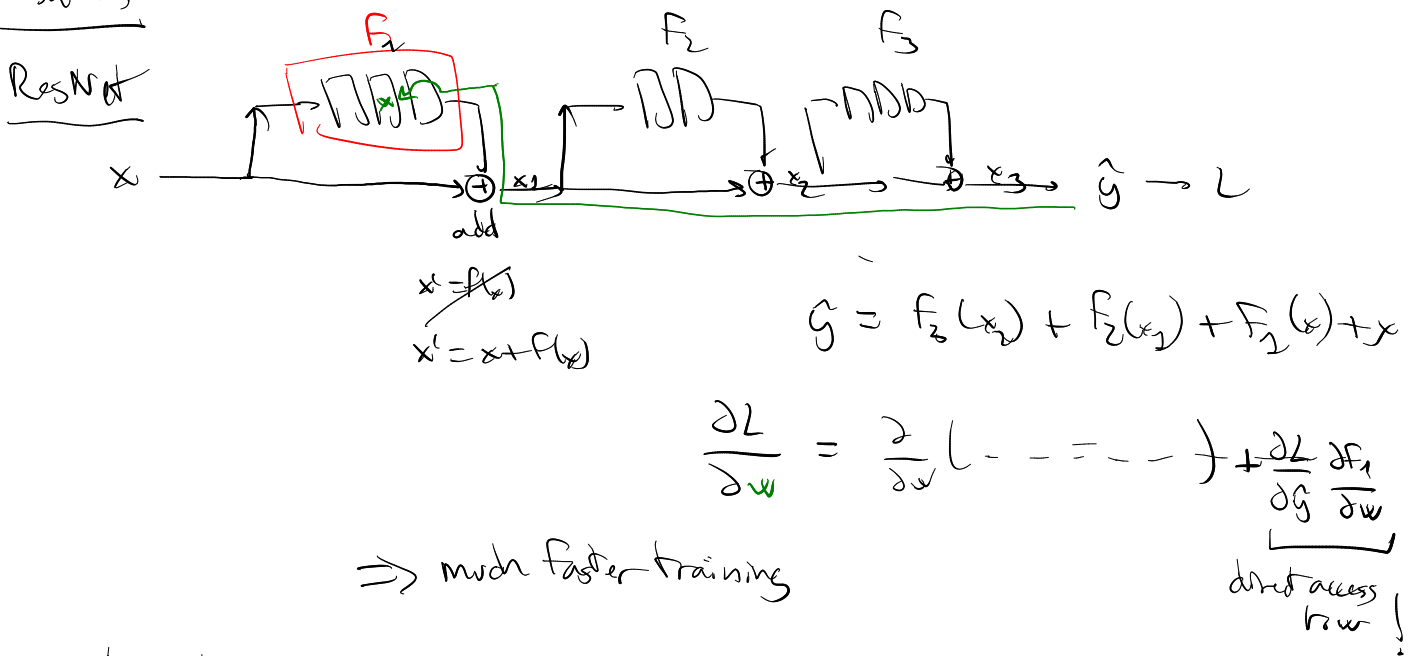
generative models = cf chapter 6



Segmentation:

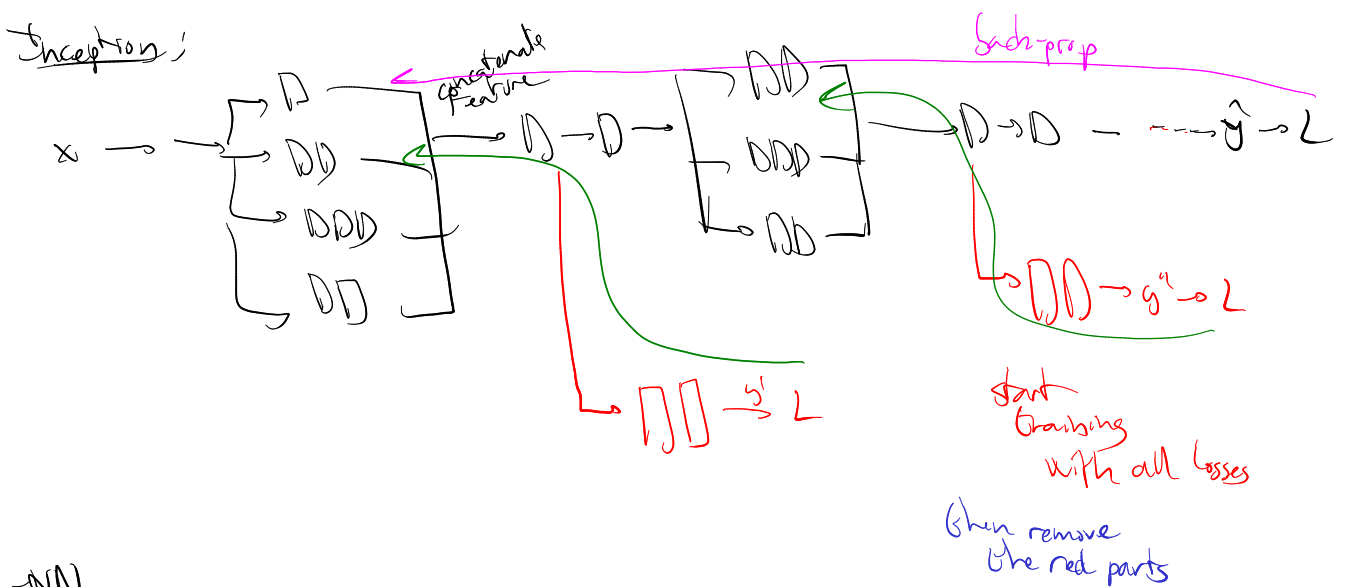


Optimisation



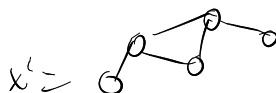
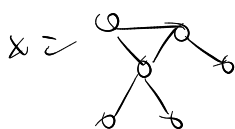
Auxiliary losses

Inception:



Graph-MN

task 1 input = a graph



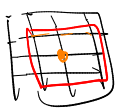
examples:

input = molecule

or: social graph

or: 3D simulation mesh

CNN

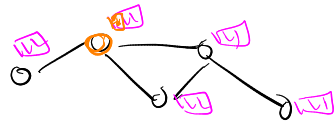
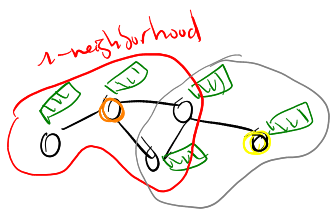
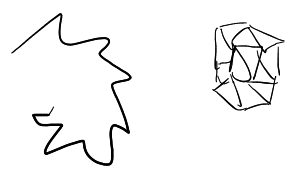


regular grid

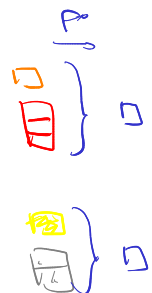


apply mask (Filter)

or: surface mesh



input x

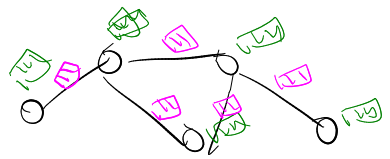


function flexible

output (of the layer)

$$y = w_{\text{center}} x_c + \sum_{n \in \text{Neighbors}} w_n x_n$$

#neighbors



$$y = w_c x_c + \sum_{n \in \text{Neighbors}} w(e_{n \rightarrow c}) x_n$$

edge-dependent

free design

very flexible

also: compute near edge information
 $e_{a \rightarrow b} = f(x_a, x_b, \{e_{n \rightarrow a}, e_{n \rightarrow b}\})$

Spectral methods

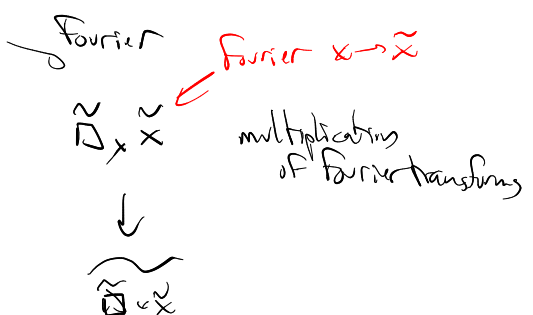
Fourier domain



$$(M * x)(p) = \int_{p'} M(p-p') x(p') dp'$$

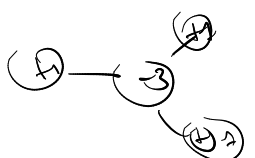
convolution

p'
 pixel location within the mesh



Fourier for graph?

↳ Laplacian Δ (graph)



Fourier

choice of mesh = multiplication weights for each frequency

↳ Split the filters: eigenvectors of Δ

↳ use $(\Delta \text{ matrix})^n$

↳ for surface meshing graph

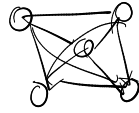
math sense only if edges are binary or distance-valued



Δ^n
 n-neighborhood

≡ fully connected graph

each word = a node

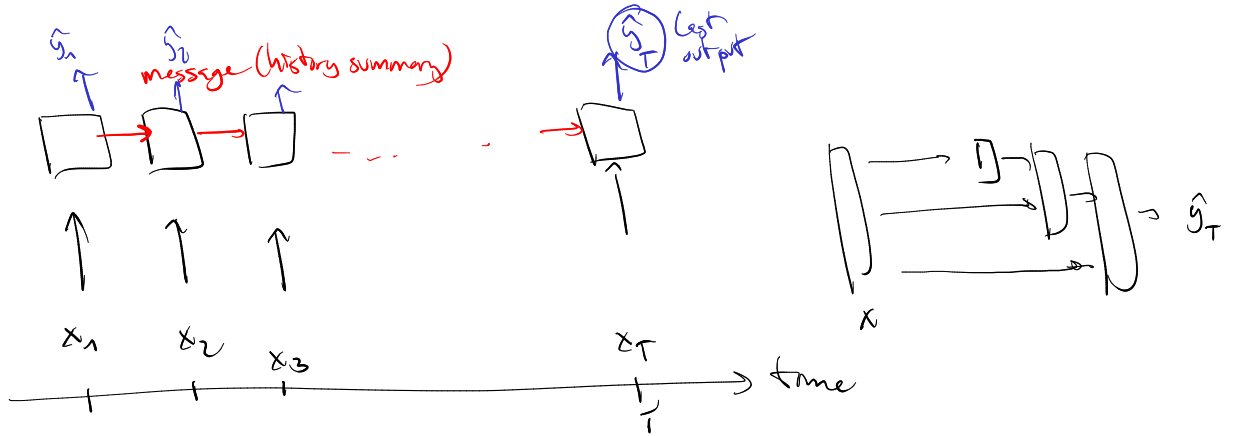


⇒ graph attention network (GAT)

⇒ listen to all nodes according to how similar they are to you
attention

Recurrent networks

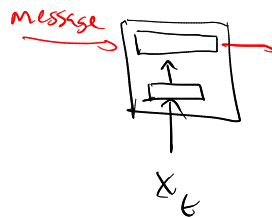
involve time : 1 input = time series $\vec{x} = (x_1, x_2, \dots, x_T)$



all blocks are identical & share the same parameters

Block = based on attention

GRU, LSTM



$$\alpha \text{ message} + (1 - \alpha) x_t$$

$$\alpha = g(\text{message}, x_t)$$

"updating" the message
 $\alpha \leq 1$

or $\alpha \leq 0$ if very important new event