



3050571 Practical Clin Data Sci

Session 14: Artificial neural network designs

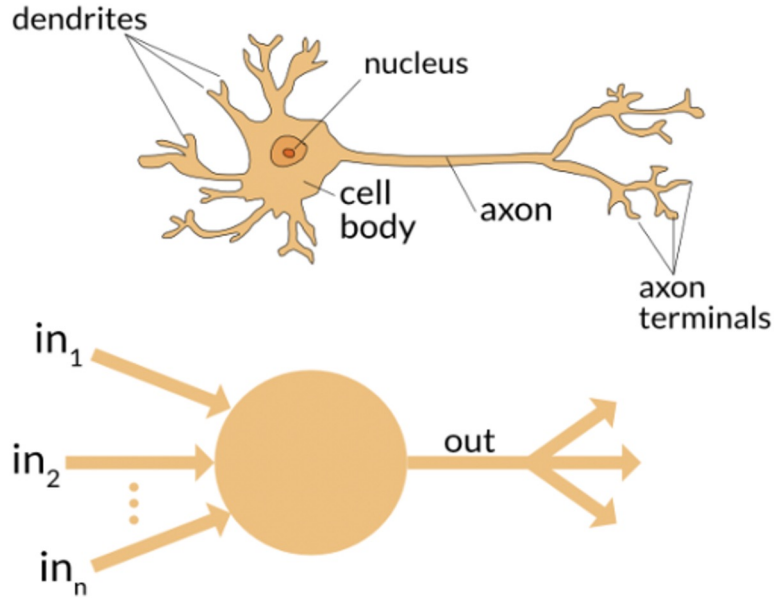
February 29, 2024



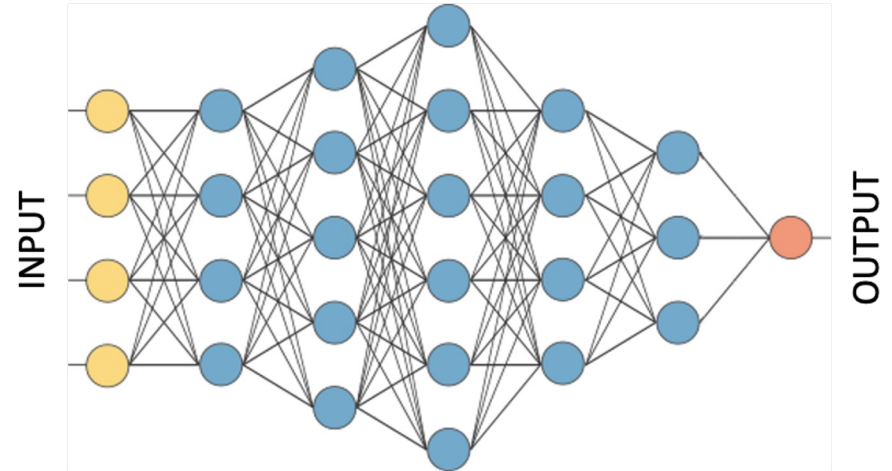
Sira Sriswasdi, PhD

- Research Affairs
- Center of Excellence in Computational Molecular Biology (CMB)
- Center for Artificial Intelligence in Medicine (CU-AIM)

Artificial neural network

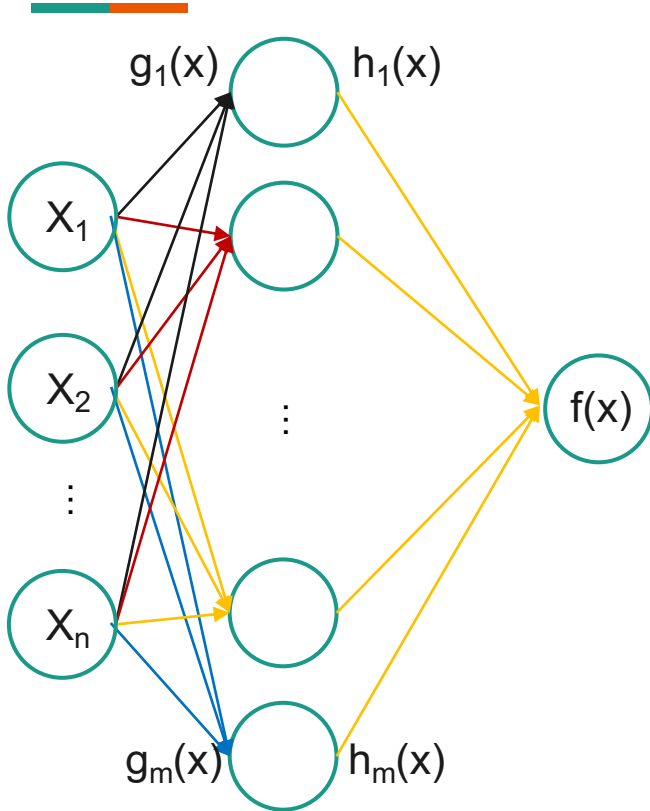


Artificial Neural Network



- Network of simple computation nodes: $out = f(w_1in_1 + w_2in_2 + ... + w_nin_n)$

Calculations inside neural network



Linear neuron input

- $g_1(x) = w_{1,1}x_1 + \dots + w_{1,n}x_n$
- $g_m(x) = w_{m,1}x_1 + \dots + w_{m,n}x_n$

Sigmoid activation

- $h_1(x) = \frac{1}{1+e^{-g_1(x)}}$
- $h_m(x) = \frac{1}{1+e^{-g_m(x)}}$

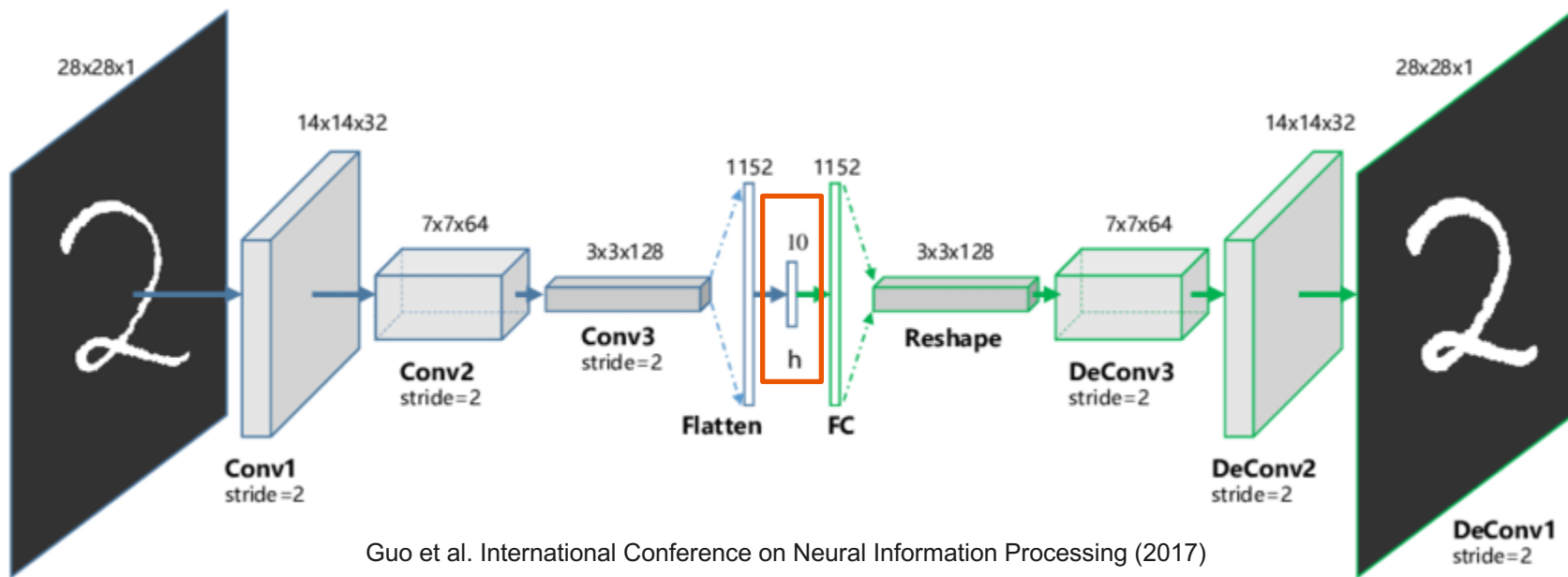
Linear aggregated output

- $f(x) = u_1h_1(x) + \dots + u_mh_m(x)$



Autoencoder

Representation learning via self-reconstruction



- Similar to dimensionality reduction

Denoising autoencoder

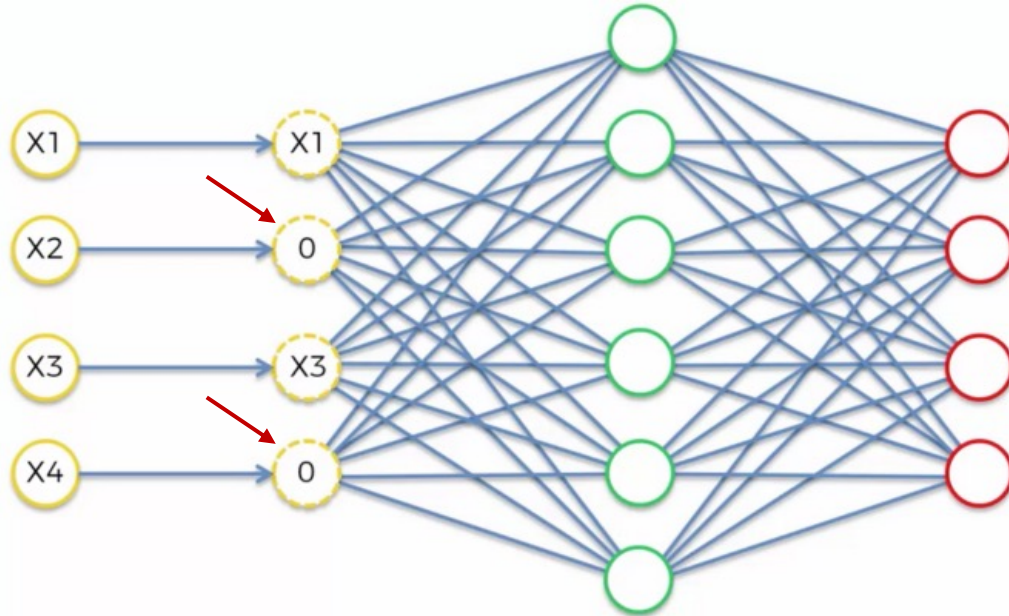


Image from towardsdatascience.com/denoising-autoencoders-explained-dbb82467fc2

- Randomly set some inputs to zero → robust representation

Variational autoencoder (VAE)

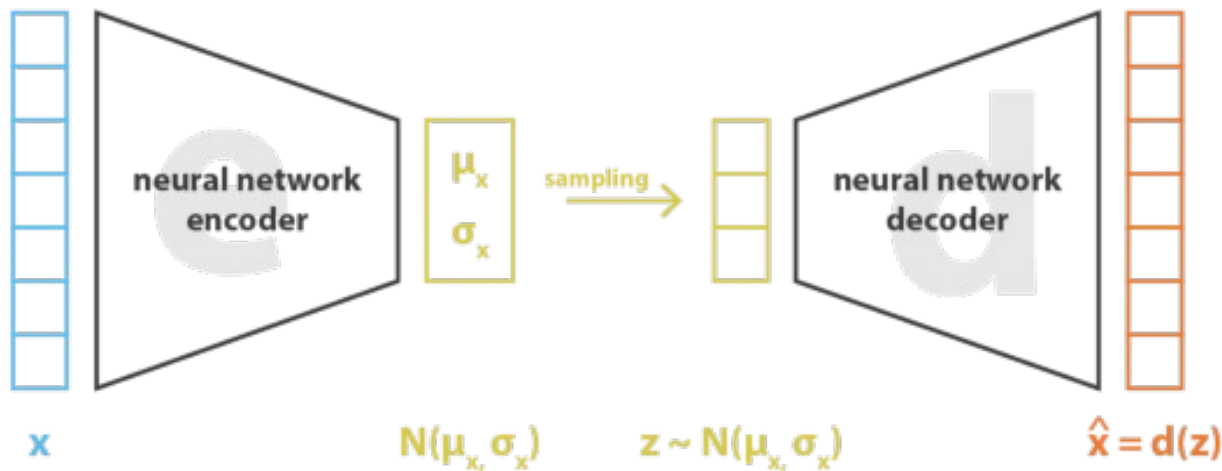
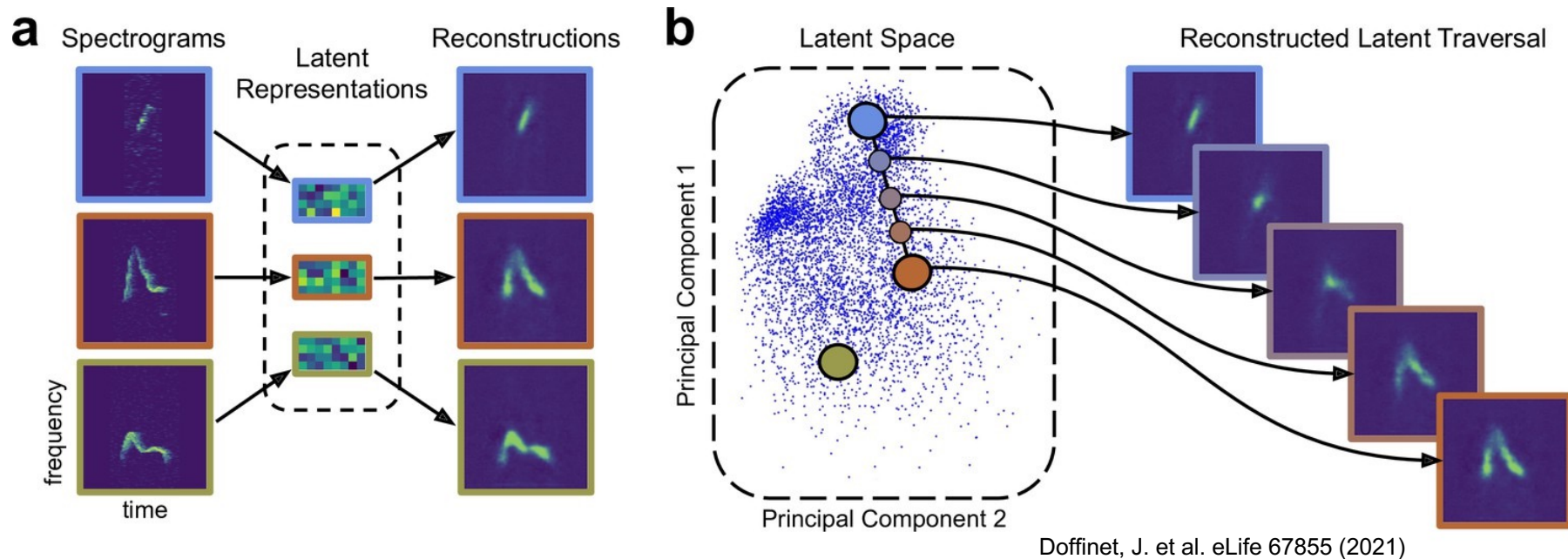


Image from www.jeremyjordan.me/variational-autoencoders/

- Learned representation = parameters for distribution
- Decoder is robust to small changes in the representation
 - Smooth representation space

VAE generates smoother representation space



- VAE learn representation distribution, not just individual vectors



Convolutional neural network

Extracting contextual pattern with filter



input image

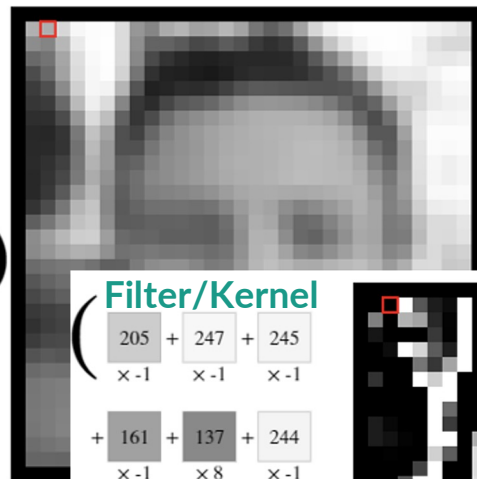
Filter/Kernel

$$\begin{pmatrix} 205 & 247 & 245 \\ \times 0.0625 & \times 0.125 & \times 0.0625 \\ + 161 & 137 & 244 \\ \times 0.125 & \times 0.25 & \times 0.125 \\ + 154 & 75 & 200 \\ \times 0.0625 & \times 0.125 & \times 0.0625 \end{pmatrix}$$

= 175

kernel: blur

<https://seco3a.io/en/image-kernels/>



Filter/Kernel

$$\begin{pmatrix} 205 & 247 & 245 \\ \times -1 & \times -1 & \times -1 \\ + 161 & 137 & 244 \\ \times -1 & \times 8 & \times -1 \\ + 154 & 75 & 200 \\ \times -1 & \times -1 & \times -1 \end{pmatrix}$$

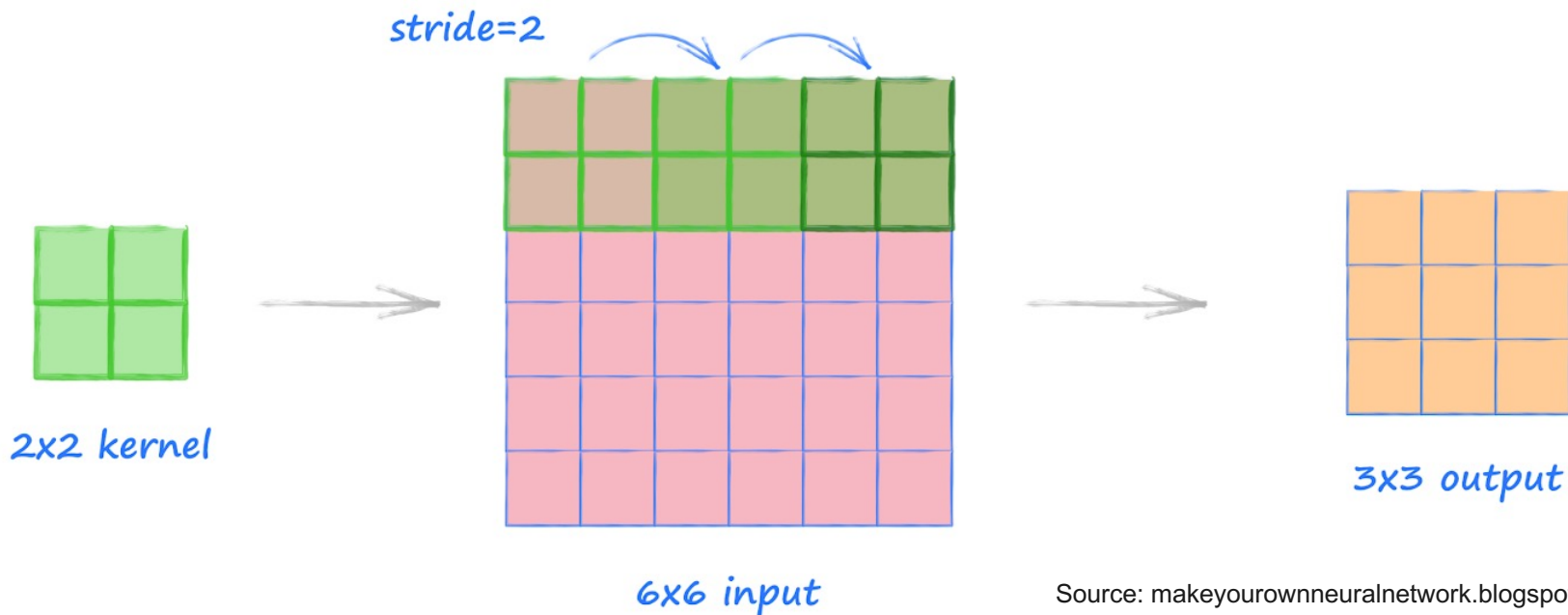
= -435

kernel: outline



output image

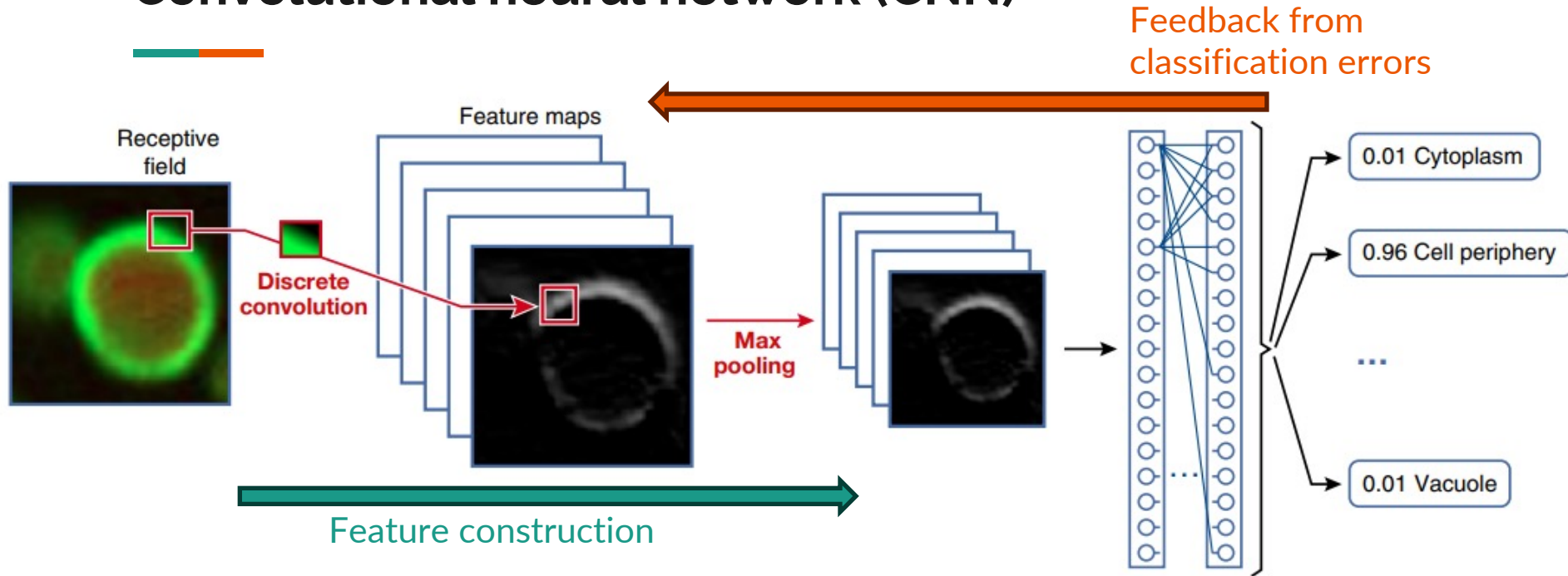
Convolutional operation



Source: makeyourownneuralnetwork.blogspot.com

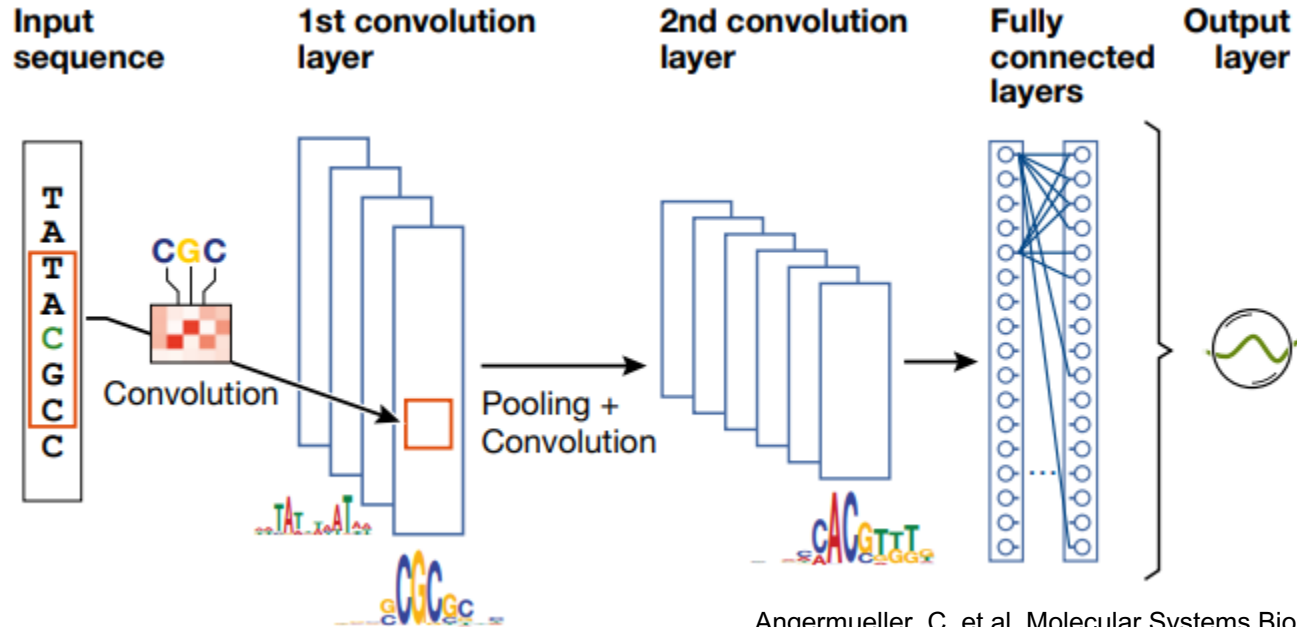
- Linear combination of values in nearby pixels – applied throughout

Convolutional neural network (CNN)



- Instead of using human-define filters to extract contextual pattern, CNN learns the best filters from the data

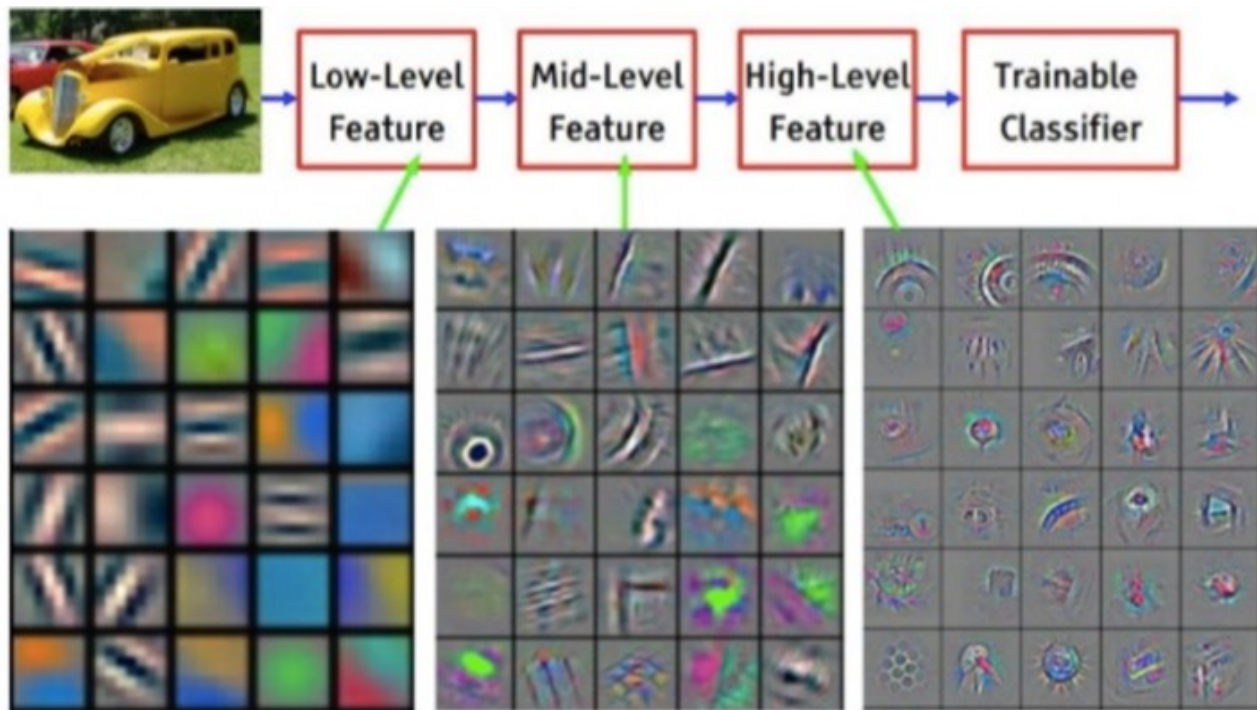
Convolution for DNA sequences



Angermueller, C. et al. Molecular Systems Biology 12:878 (2016)

- Motif = contextual pattern on DNA sequence

Hierarchical feature assembly inside CNN





Some CNN designs

Vanishing and exploding gradient problems



$$\text{Gradient: } \frac{\delta L}{\delta w_{i,j}} = \frac{\delta L}{\delta f} \frac{\delta f}{\delta h_i} \frac{\delta h_i}{\delta g_i} \frac{\delta g_i}{\delta w_{i,j}} = (f(x) - y) \cdot u_i \cdot g_i(x)(1 - g_i(x)) x_j$$

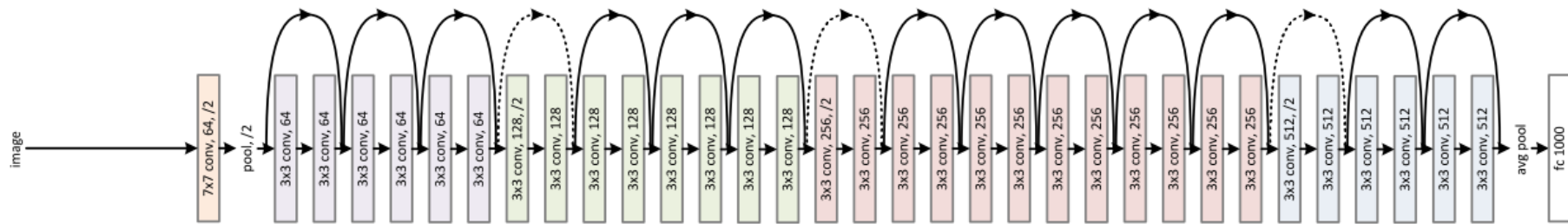
The number of multiplicative terms scales with the number of layers

What would happen if all values are $\ll 1$ or $\gg 1$?

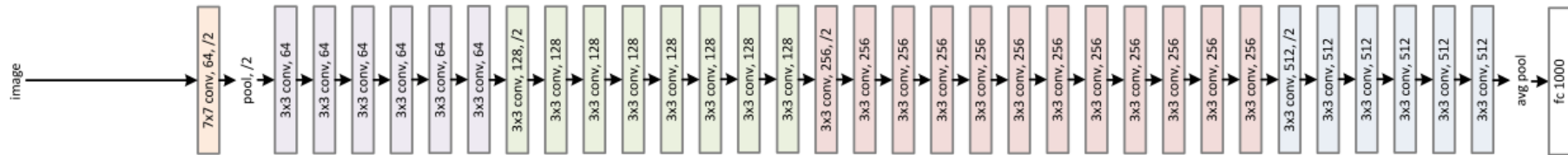
- Gradient became **very small** \rightarrow No weight update
- Gradient became **very large** \rightarrow Unstable

Residual network (ResNet)

34-layer residual



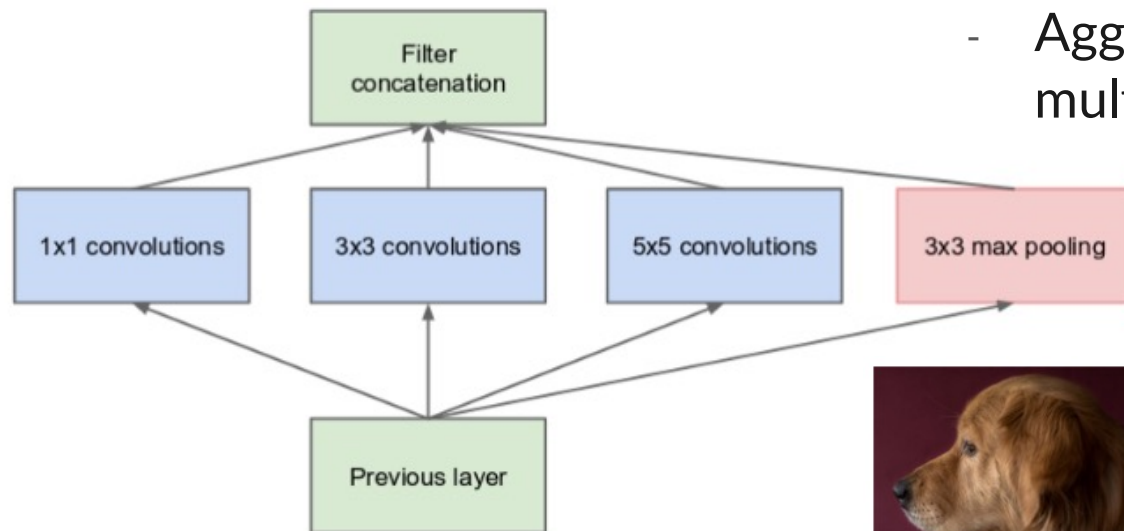
34-layer plain



Source: medium.com

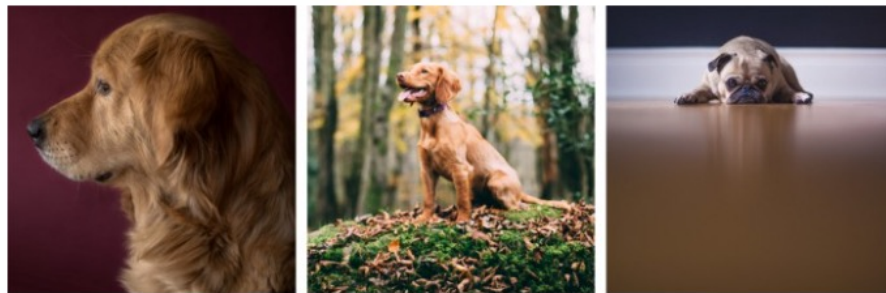
- Adding skip connections jumping over blocks of convolutional layers
- Reduce the number of terms in gradient of early weights

Inception = multi-resolution layer



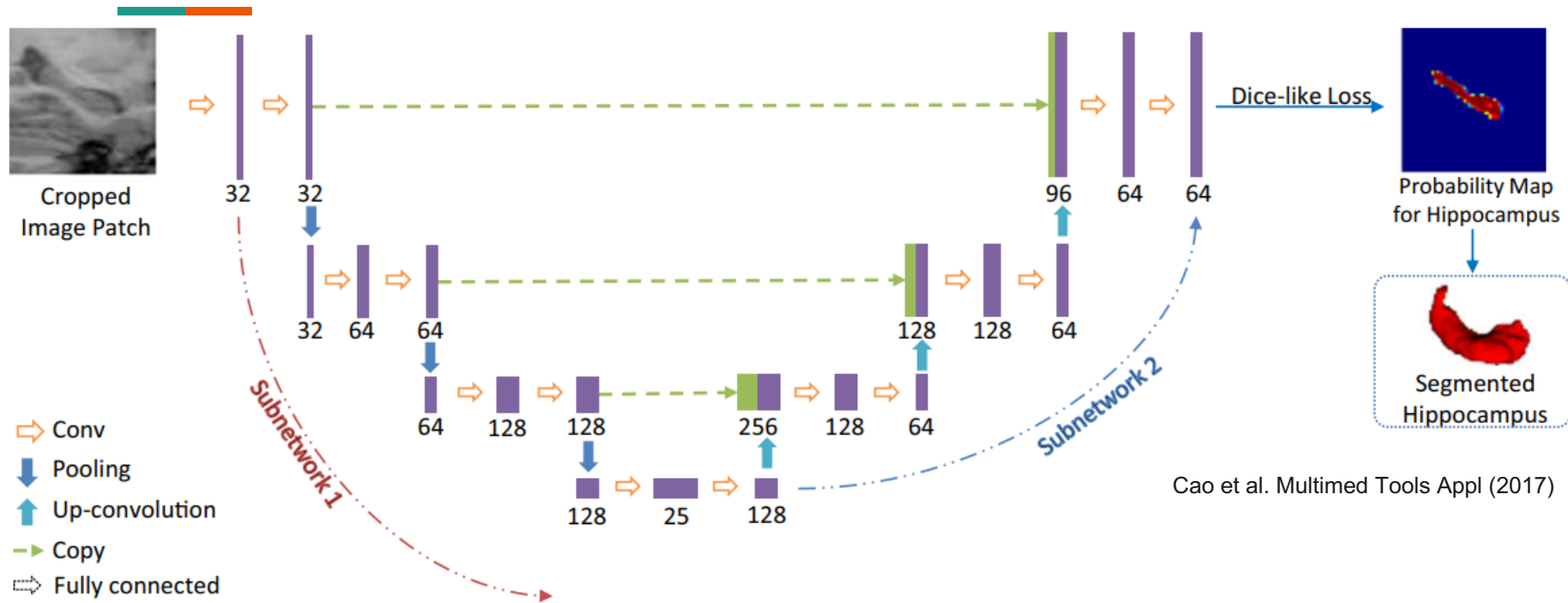
- Aggregate information from multiple resolutions together

- Convolution with multiple filter sizes per layer



From left: A dog occupying most of the image, a dog occupying a part of it, and a dog occupying very little space (Images obtained from [Unsplash](#)).

U Net = producing image from image



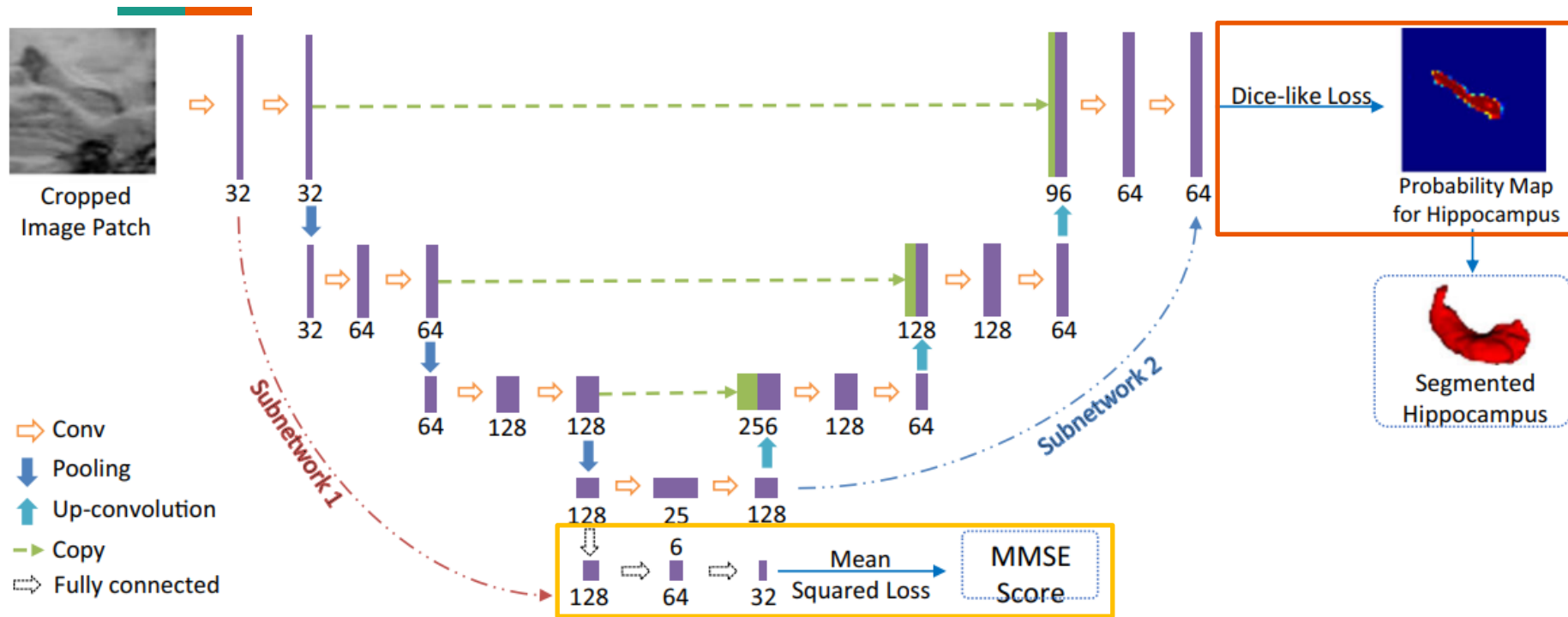
Cao et al. Multimed Tools Appl (2017)

- Make prediction for every pixel \rightarrow output size = input size



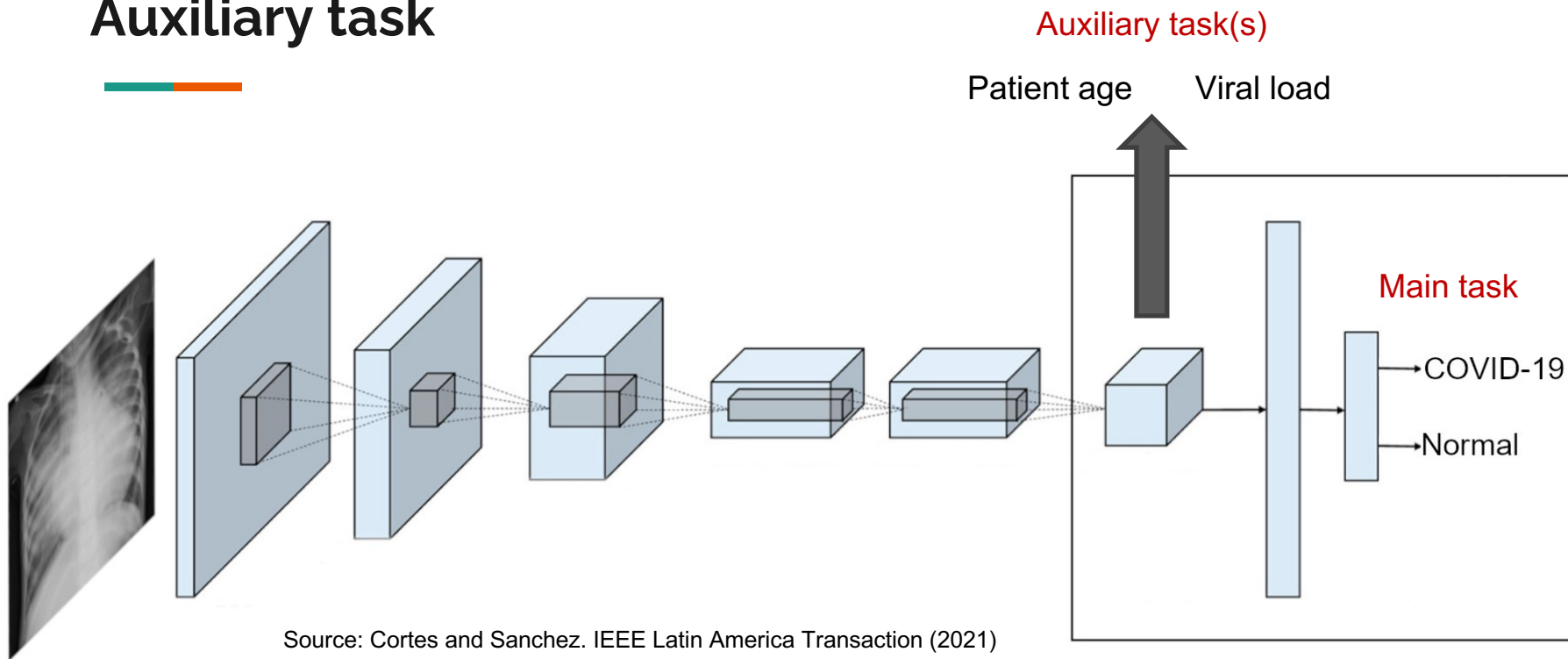
Multitasking

Simultaneous segmentation & classification



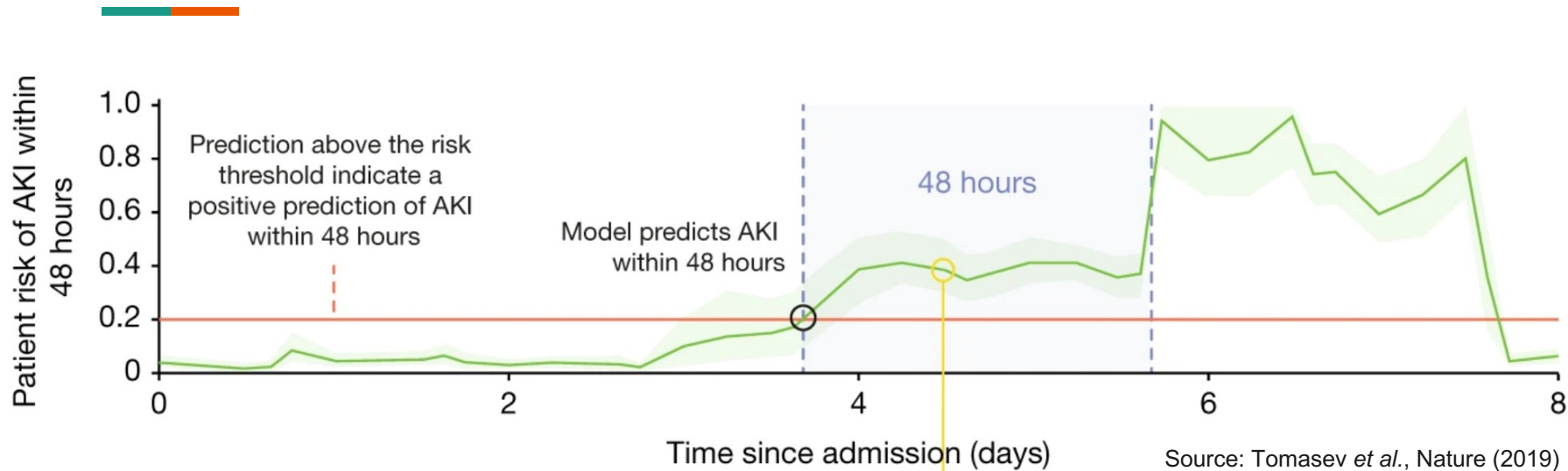
- Combine gradients from both tasks

Auxiliary task



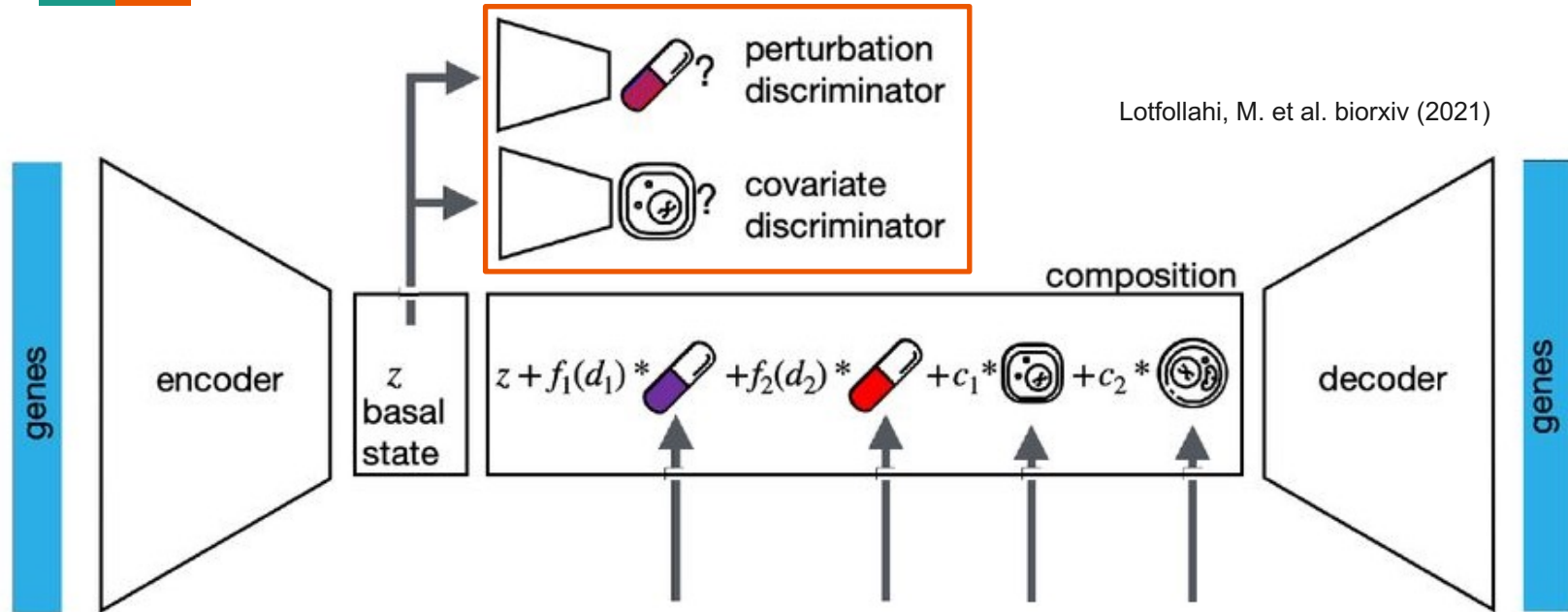
- Encourage the learned representation to include more information

Acute kidney injury prediction



- **Main task:** Occurrence of acute kidney injury within 48 hours
- **Auxiliary tasks:** Maximal values of 7 key lab tests within 48 hours
 - Provide more feedback on what the model gets wrong

Decoupling / debiasing

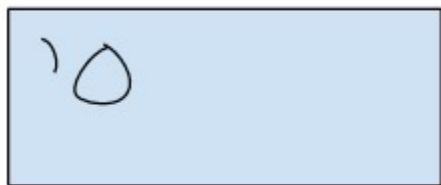


- Deconvolute cell basal state from perturbation and covariate
- Update weights in the opposite direction of gradient



Generative model

Why generative model?

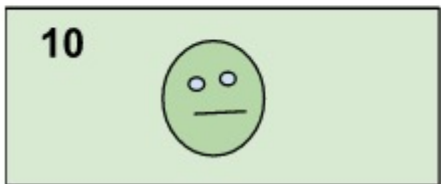


FAKE

REAL



https://developers.google.com/machine-learning/gan/gan_structure



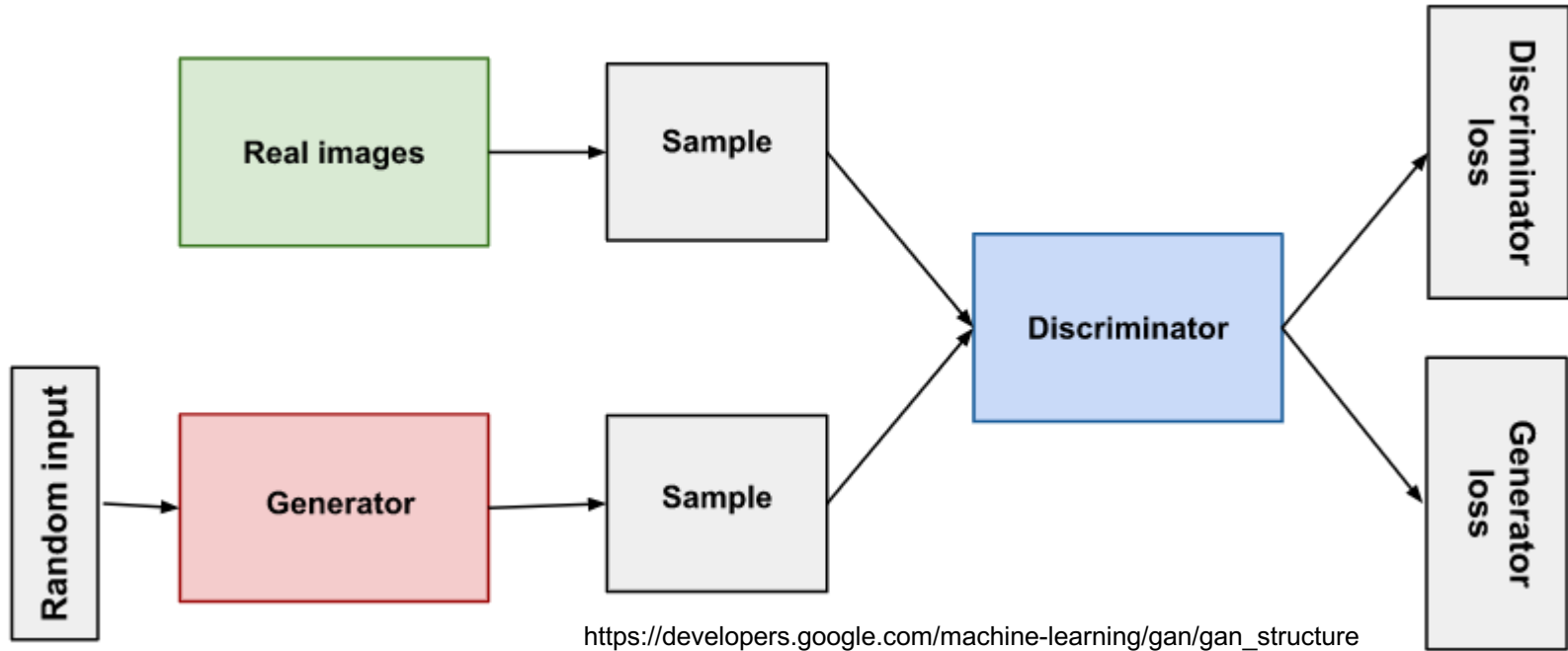
FAKE

REAL



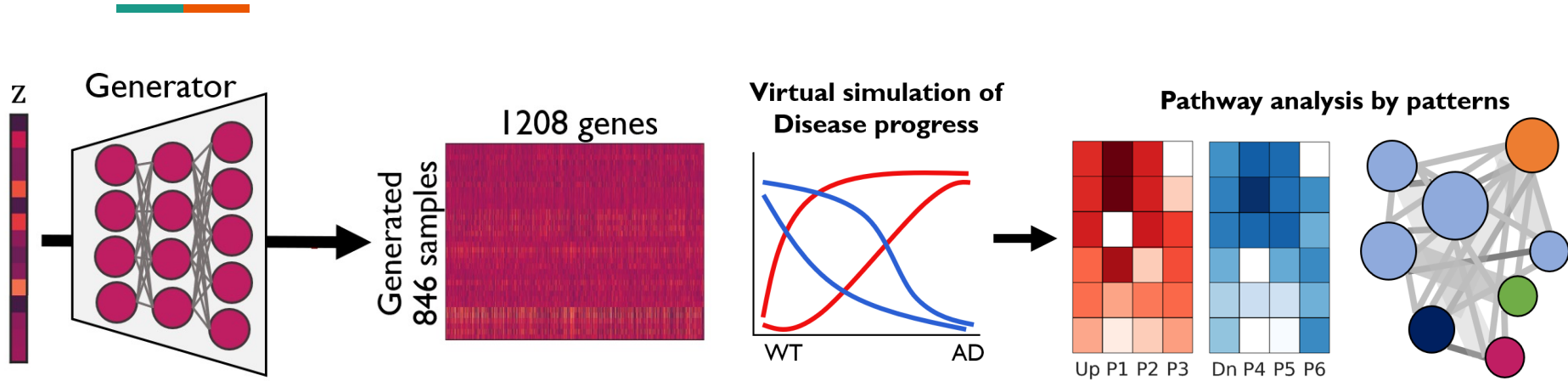
- Models that **generate realistic data** can tell us about the underlying mechanisms of the system

Generative adversarial network (GAN)



- Simultaneous training of **generator** and **discriminator**

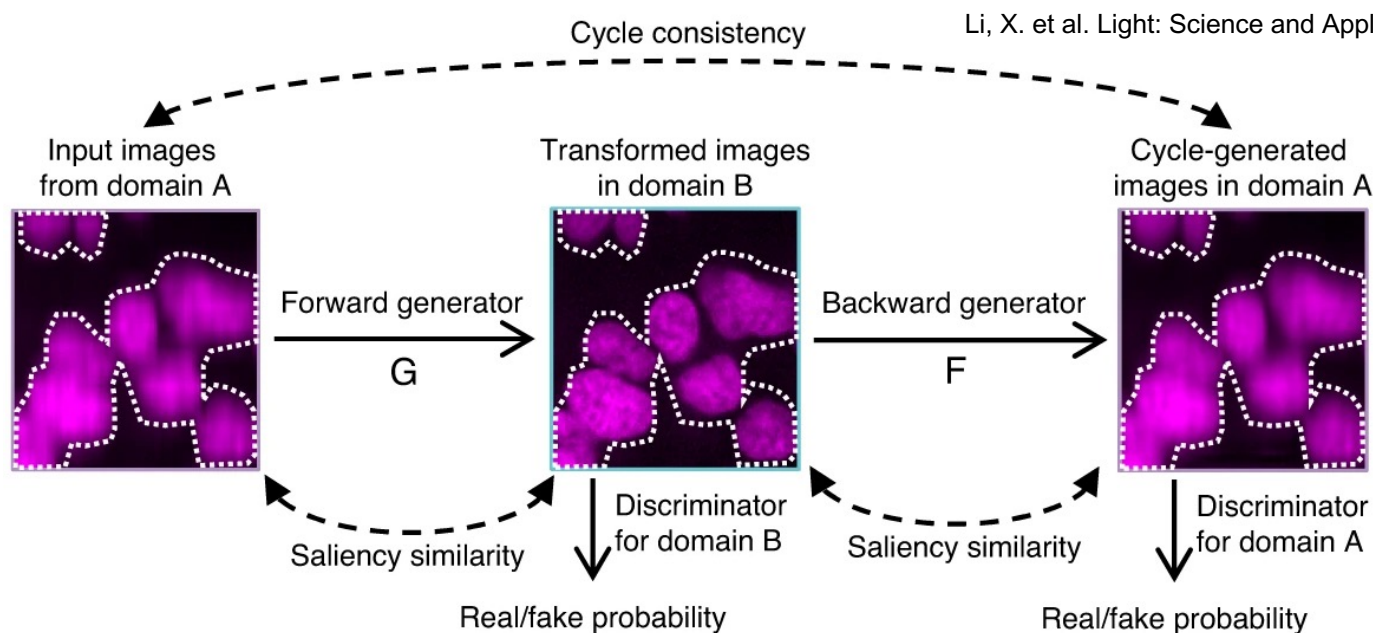
Knowledge from simulated data



Park, J. et al. PLoS Computational Biology 16:e1008099 (2020)

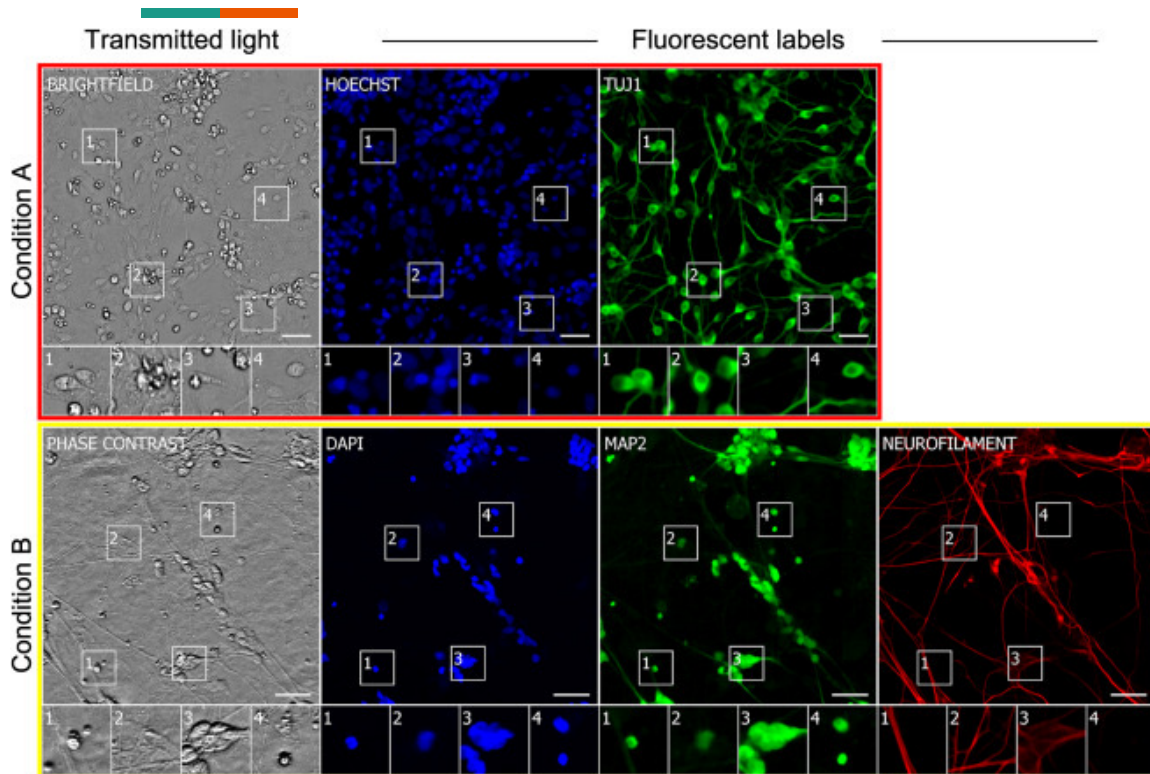
- Train a generator with data from small-scale experiment
- Simulate time-course gene expression profiles
- Perform usual bioinformatics analyses to infer biological knowledge

Cycle GAN for transforming image

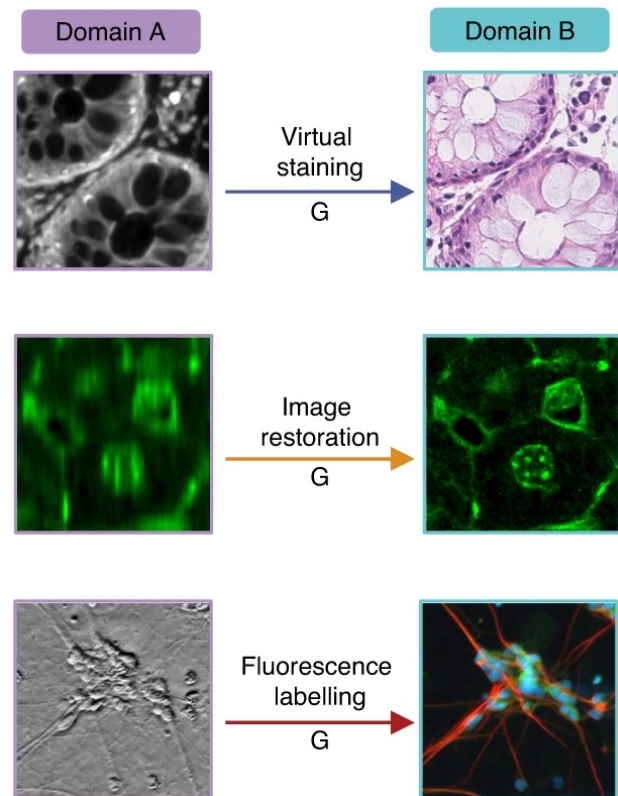


- Generate **sharpened image** from **blurry image** and back

Virtual staining



Christiansen, E.M. et al. Cell 173:792-803.e19 (2018)

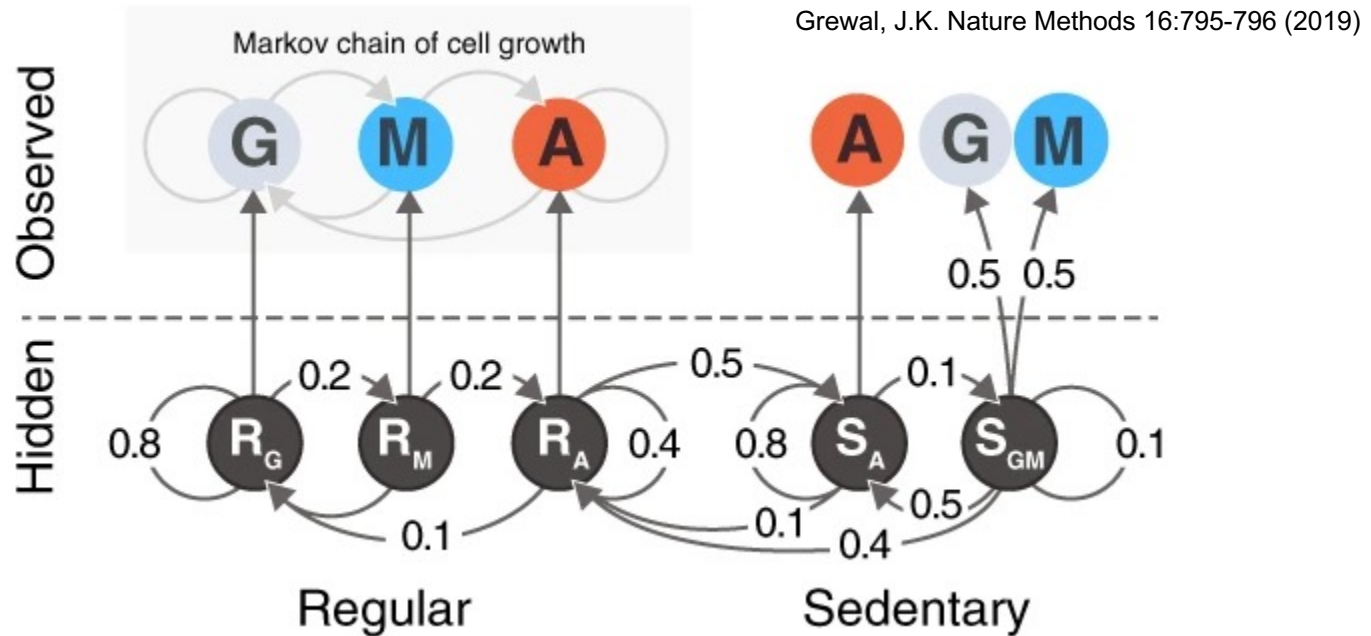


Li, X. et al. Light: Science and Applications 10:44 (2021)



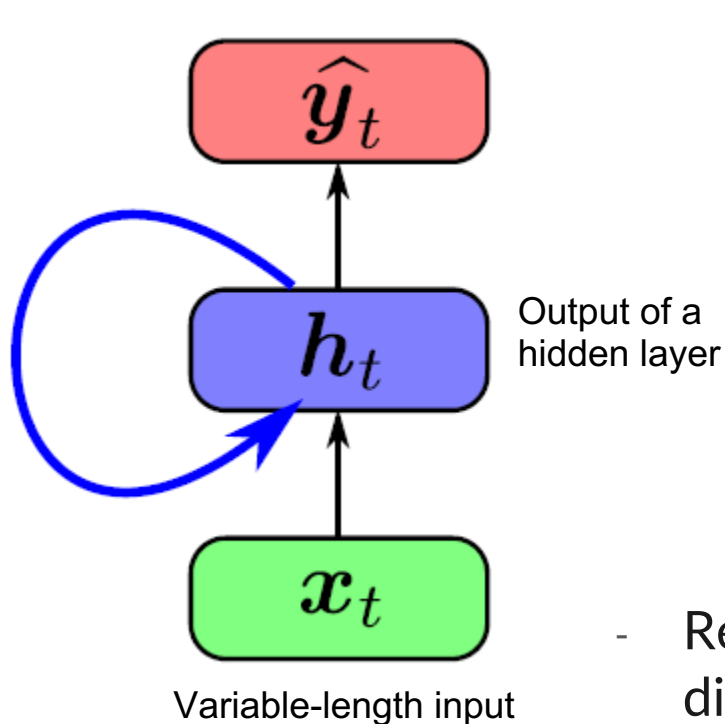
Recurrent neural network

Hidden Markov Model



- Sequence of observations, each generated from a model

Recurrent neural network



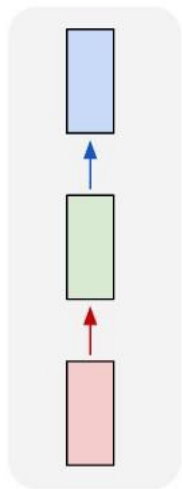
Shared weights!

$$\begin{aligned}h_1 &= f(\mathbf{u} \cdot x_1 + \mathbf{v} \cdot h_0 + c) \\h_2 &= f(\mathbf{u} \cdot x_2 + \mathbf{v} \cdot h_1 + c) \\&\dots \\h_t &= f(\mathbf{u} \cdot x_t + \mathbf{v} \cdot h_{t-1} + c) \\ \hat{y}_t &= \mathbf{w} \cdot h_t + b\end{aligned}$$

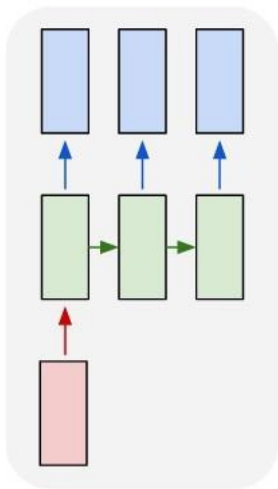
- Reuse a single layer (weights) over time with different input

Sequence-to-sequence capability

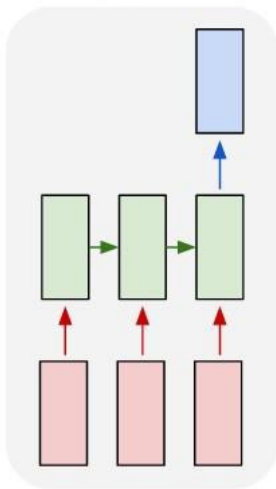
one to one



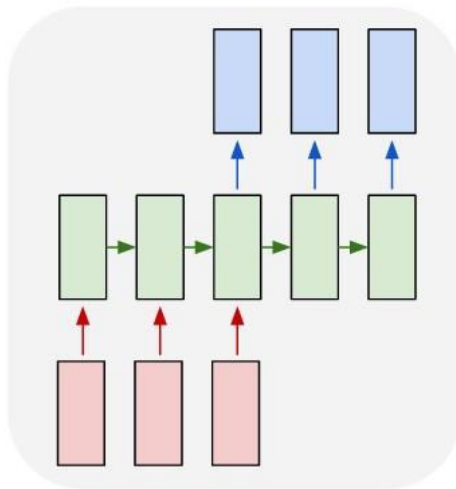
one to many



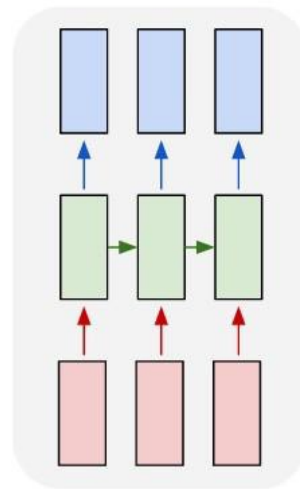
many to one



many to many



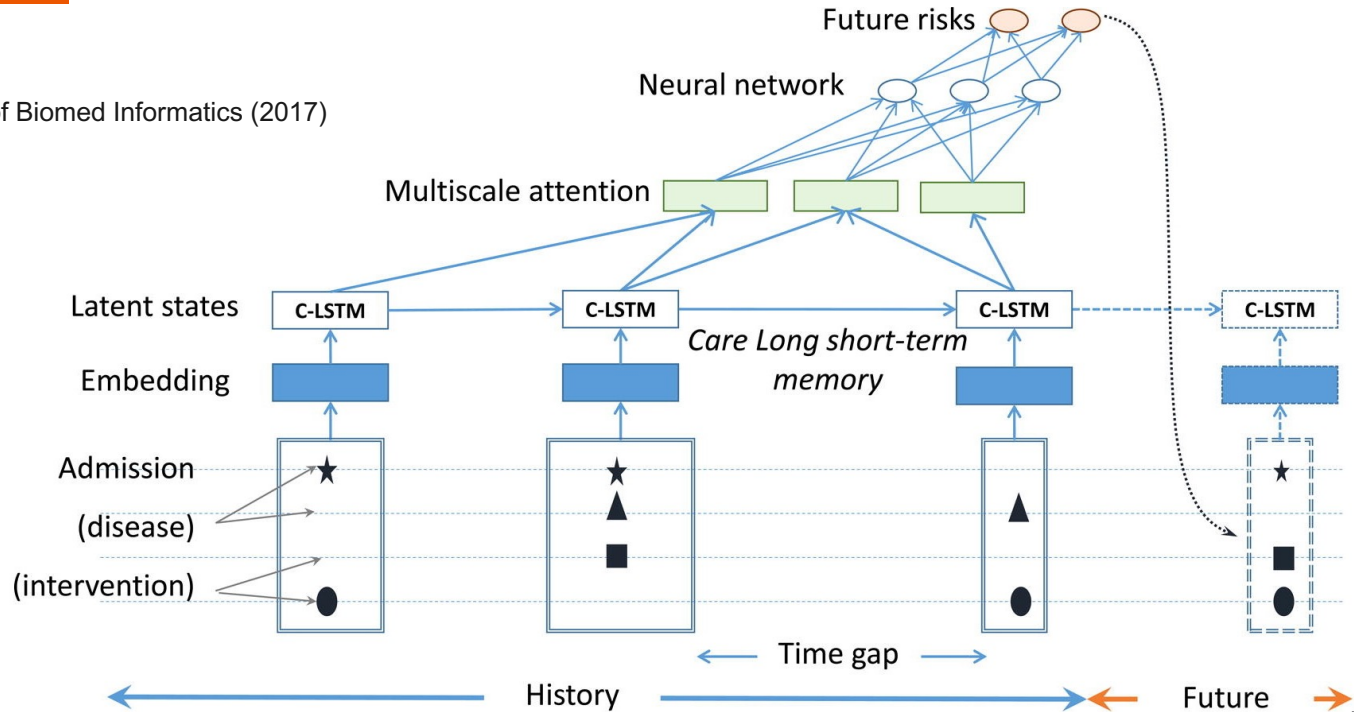
many to many



- Handle variable input and output

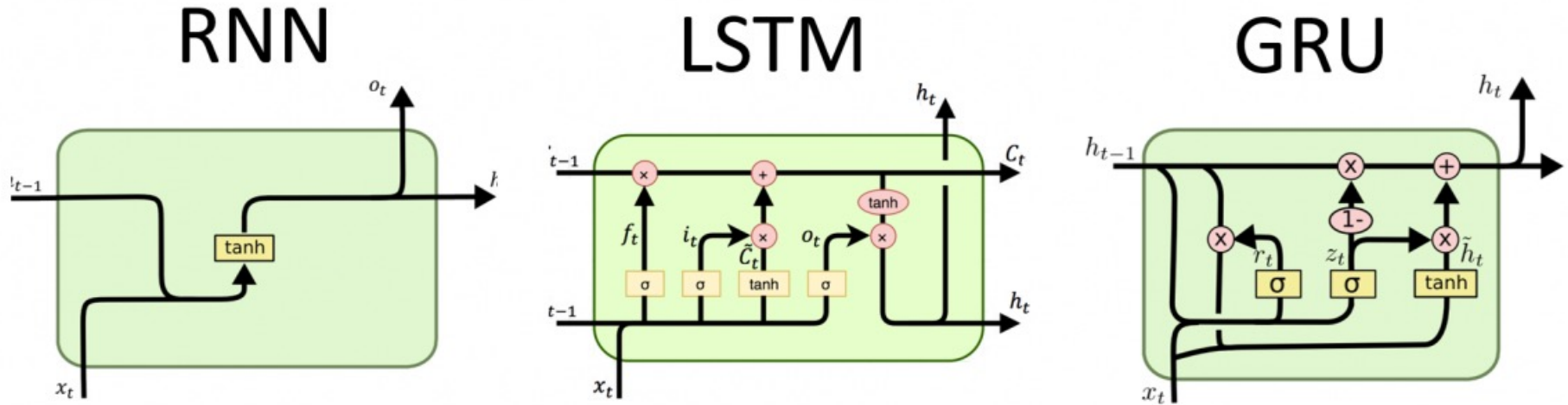
RNN on medical history

Pham et al. J of Biomed Informatics (2017)



- Aggregate information across time to make prediction

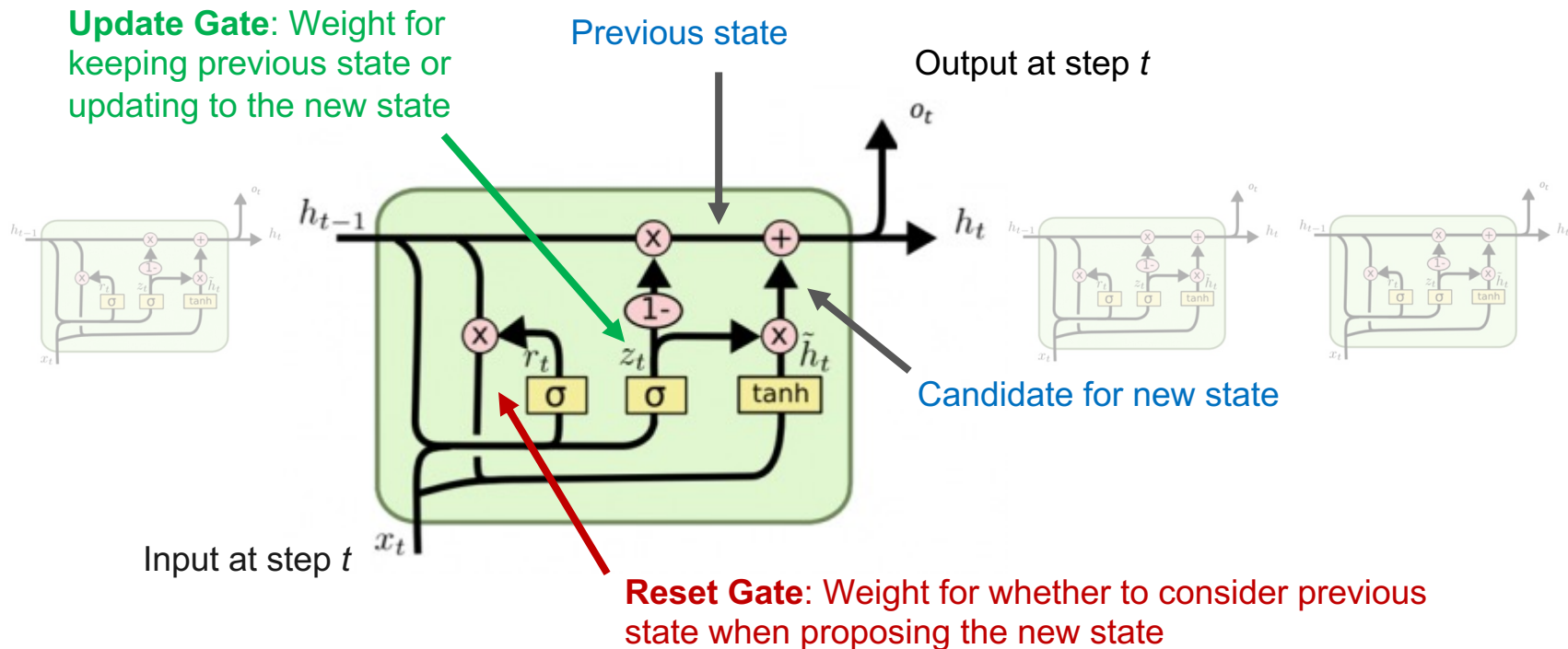
RNN architecture



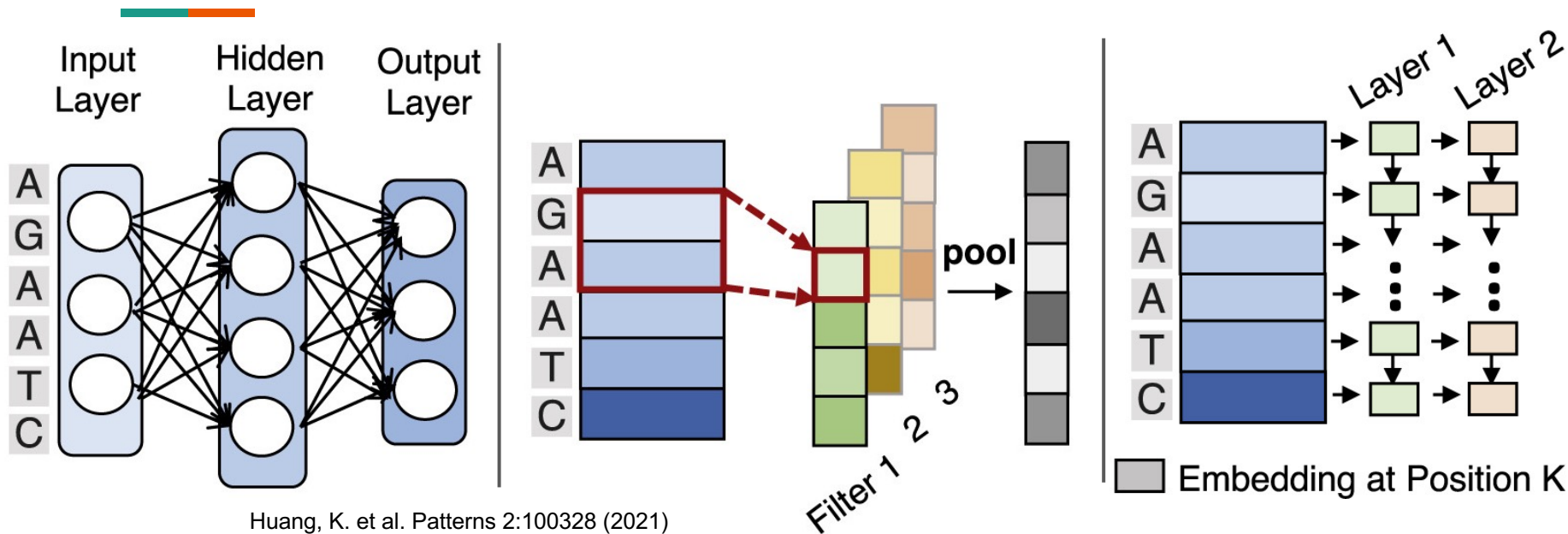
Source: www.linkedin.com/pulse/recurrent-neural-networks-rnn-gated-units-gru-long-short-robin-kalia

- Allow the model to **retain** / **forget** information from earlier time points
- Include **shortcuts for gradient calculation** – similar to ResNet

Gated recurrent unit (GRU)



Picking the right model requires domain knowledge

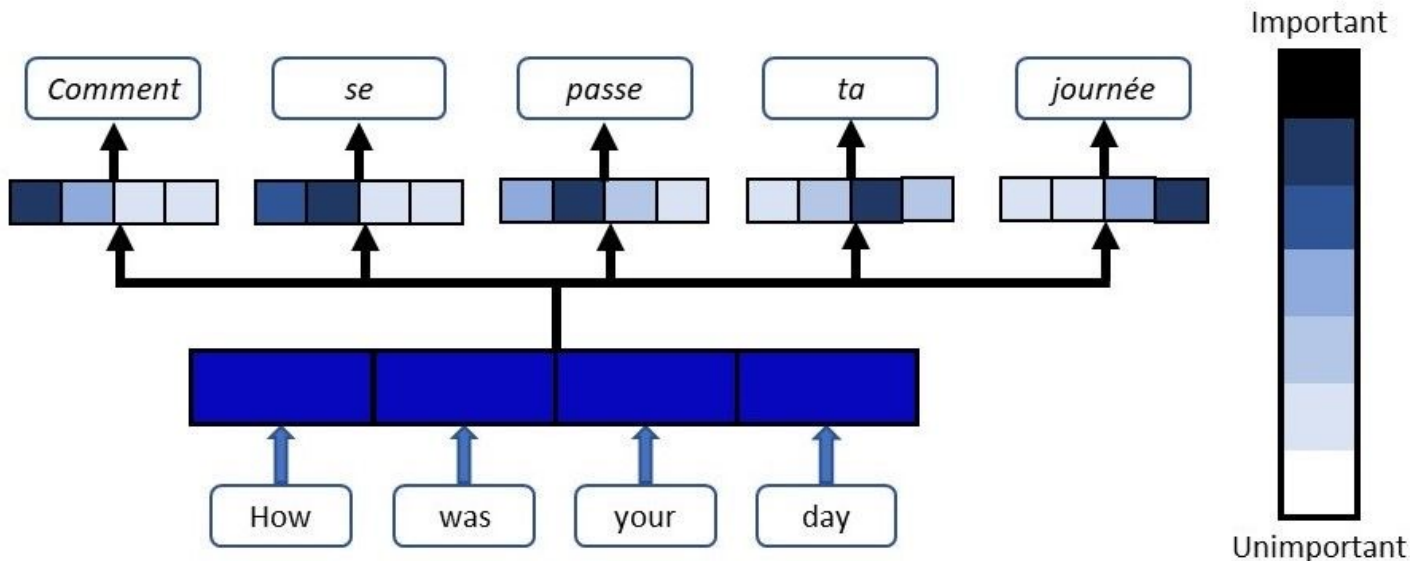


- Choosing the “right” model depends on the interpretation of the task and the underlying mechanisms



Transformer and attention

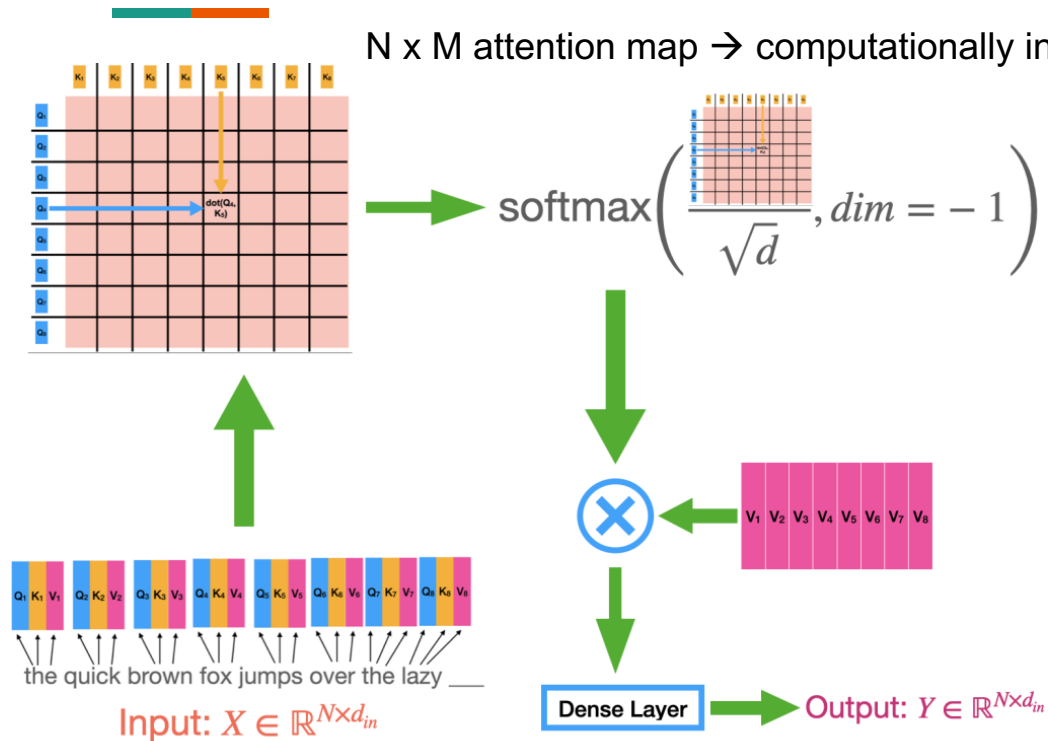
Attention explicitly model the contribution of each input



<https://blog.floydhub.com/attention-mechanism/>

- The model learn to estimate contribution directly during training, not as a post-hoc explainability

Attention vs regular layers



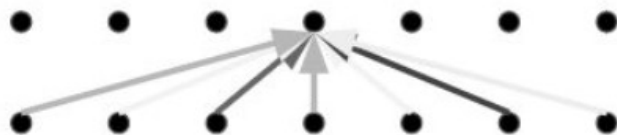
Convolution

Fixed weights regardless of input

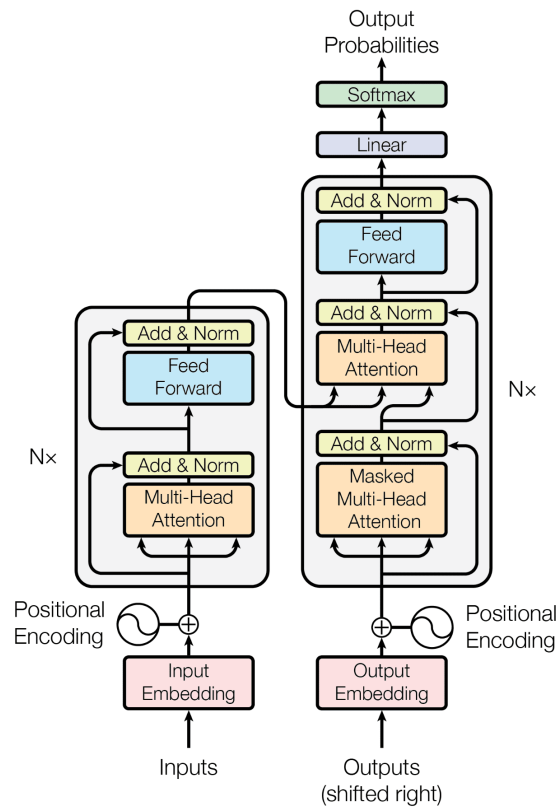
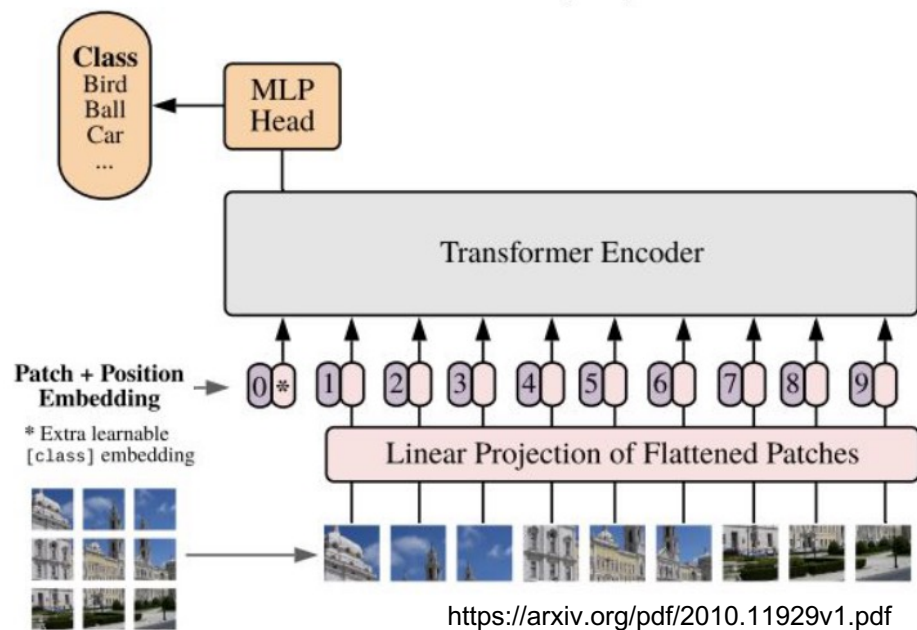


Self-Attention

Adaptive weights depending on input

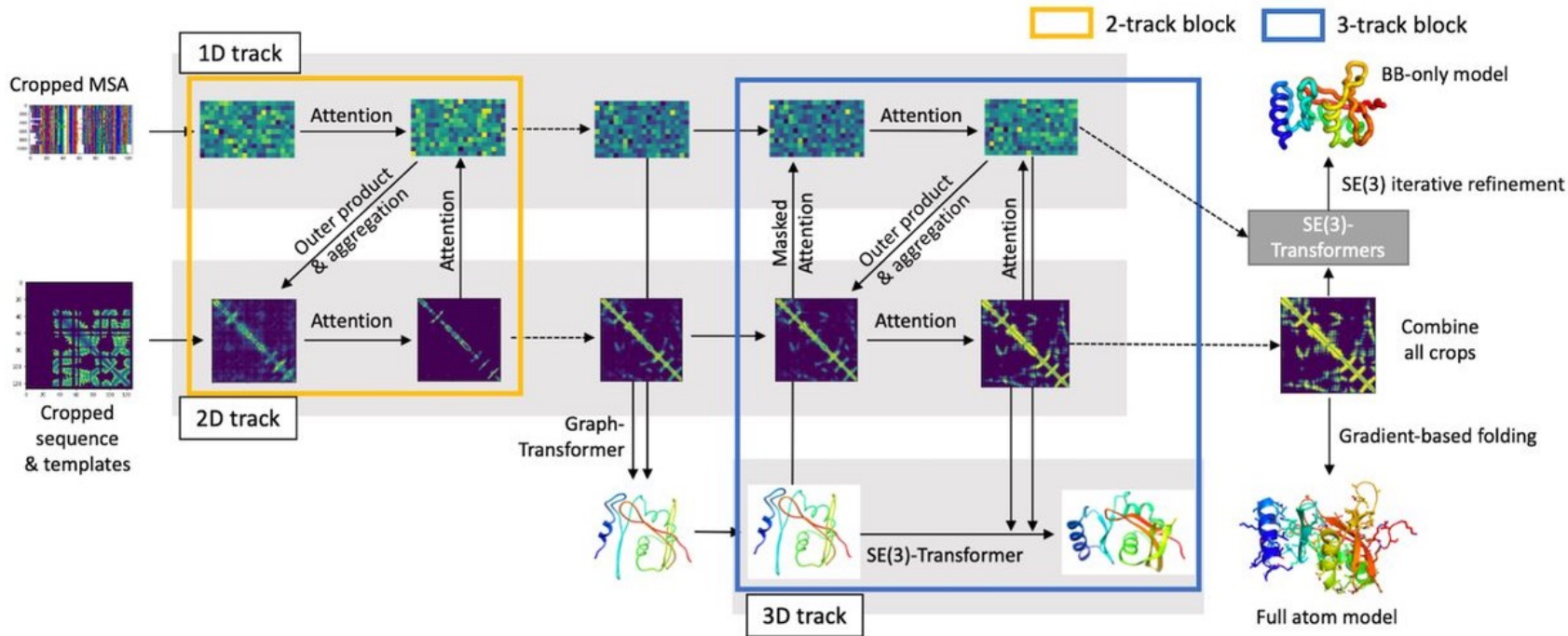


Transformer architecture



- Multiple attention layers replace conventional networks

Transformer model for protein folding prediction



- Attention is all you need

Any questions?



See you on March 1st