

CSE8803/CX4803

Machine Learning in Computational Biology

Lecture 6: Deep learning for sequence data

Yunan Luo

CSE8803/CX4803 in-class survey

Please complete this 1-min survey

<https://tinyurl.com/cse8803DLsurvey>

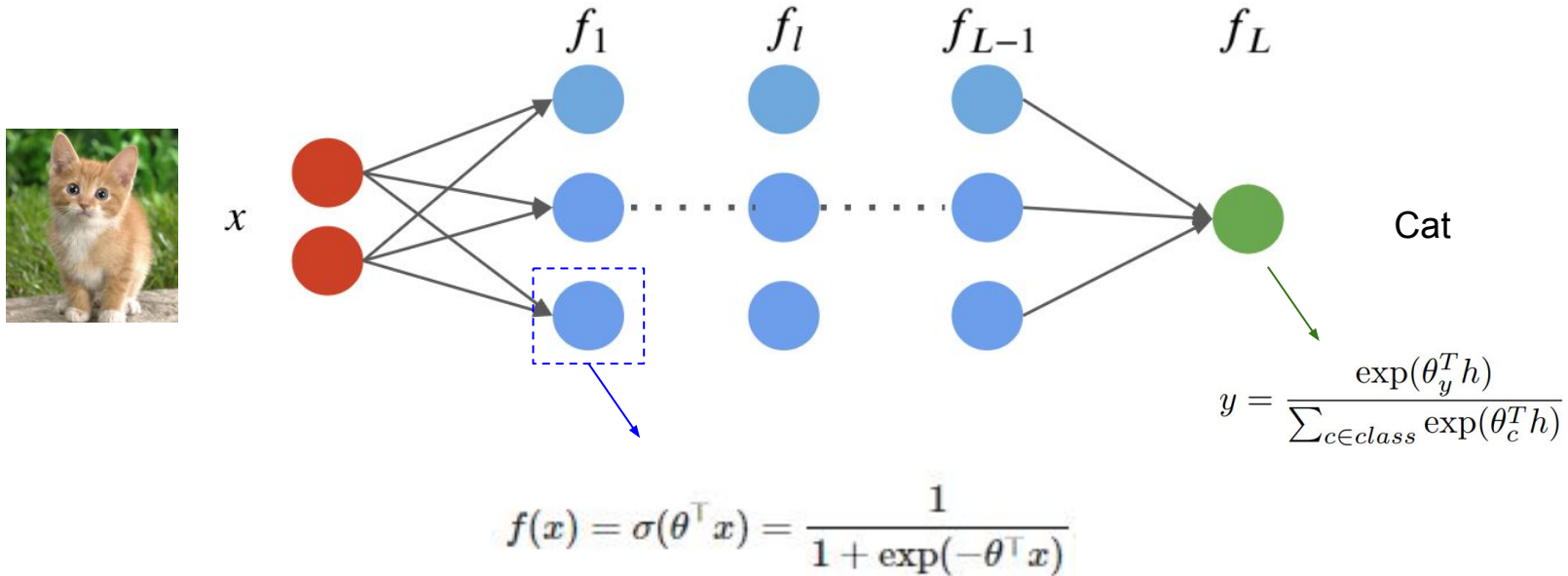


Today's plan

- A short primer on deep learning
- Biological sequences
- Deep learning for sequence data in biomedicine
 - Supervised learning
 - CNN, RNN, LSTM, Transformer
 - Unsupervised learning
 - Protein language modeling

Recap: neural network

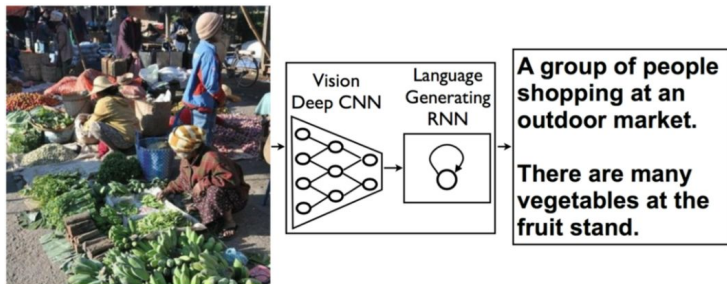
A (fully connected) neural network is a computational model that consists of a composition of multiple neural network layers



Deep learning

A class of machine learning methods that emphasizes:

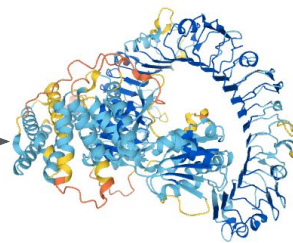
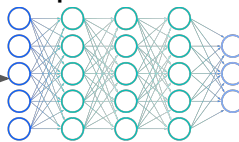
- *Deep* neural networks
- Large unstructured datasets, especially images, text, and audio
- Modern computational resources, like GPUs



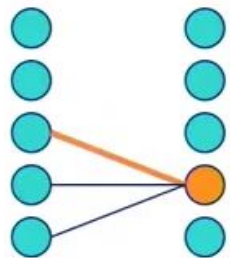
Ref: Learning CNN-LSTM Architectures for Image Caption Generation: <https://cs224d.stanford.edu/reports/msoh.pdf>

```
MAGELVSFAVNKLWDLISHEYTLFGQVEDQVAE  
CVEEIKDIVYDAEDVLETFVQKEKLGTTSGIRK  
RVIRDMQSFVGQMIYDDYMHPLNREREIRRT  
NYQVVSITGMGGLGKTTLARQVFNHDMVTKKFD  
EETKEEEKKILEMTEYTLQRELYQLLEMSKSL  
LLTSRNESIVAPTNTKYFNFKEPCLKTDDSWKL  
IEHCGGLPLAIKVLGGMLAEKYTSHDWRLSEN  
FEELPSYLKHCFLYLAHFPEDEYIKVENLSYYN  
VRRNMVISERDVKTSRFETCHLHDMIREVCLLK  
LVYQYPTTLHVEKDINNPKLRSLVVTLGSWNM  
SCIGKLIHLRYSLEYAEVTHIPYSLGNLKLLI  
LALPSLIERKTKLESLNLVKLETLENFSTKNSS
```

AlphaFold 2

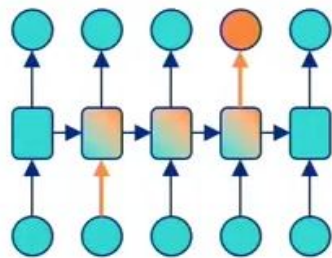


Deep learning models and implementations



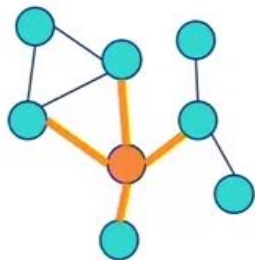
Convolutional Networks
(e.g. computer vision)

- data in regular grid
- information flow to local neighbours
- AlphaFold 1



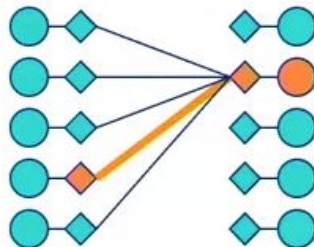
Recurrent Networks
(e.g. language)

- data in ordered sequence
- information flow sequentially



Graph Networks (e.g. recommender systems or molecules)

- data in fixed graph structure
- information flow along fixed edges



Attention Module (e.g. language)

- data in unordered set
- information flow dynamically controlled by the network (via keys and queries)

 PyTorch

 TensorFlow

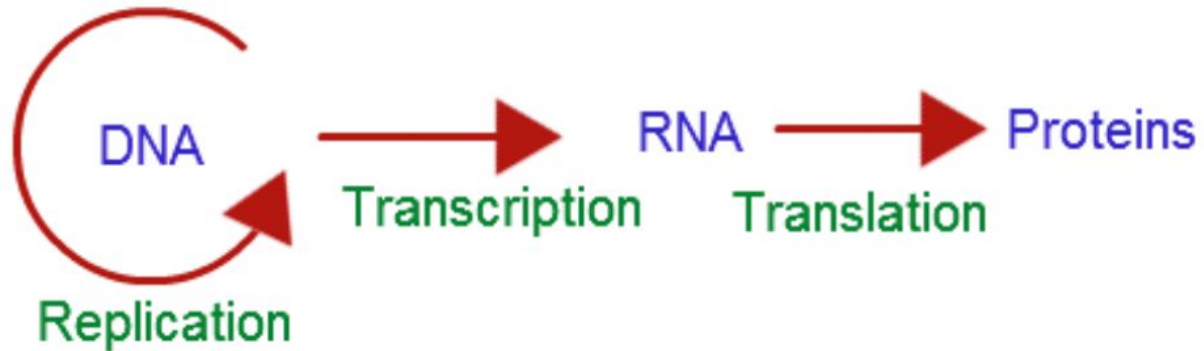


 Keras



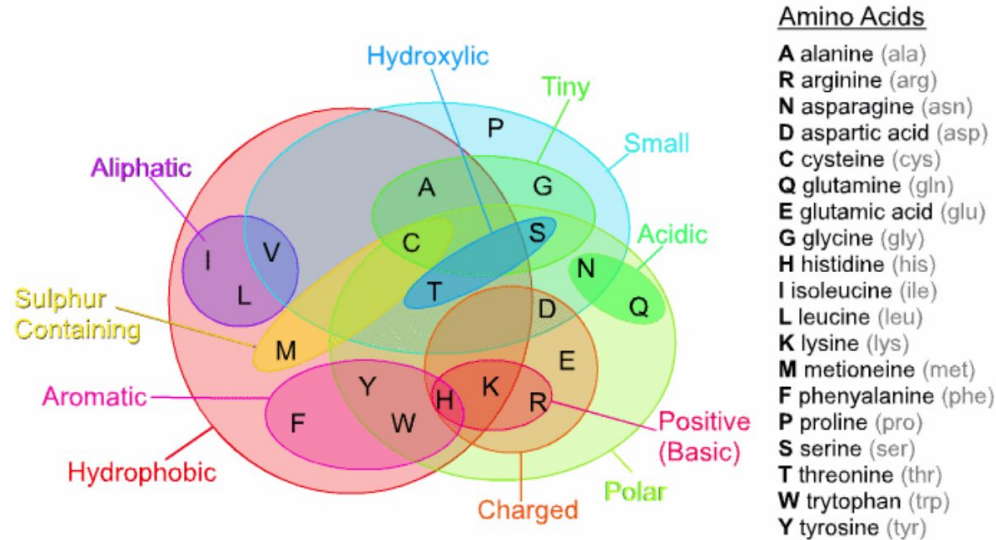
Biological sequences

Central dogma



Protein and amino acid

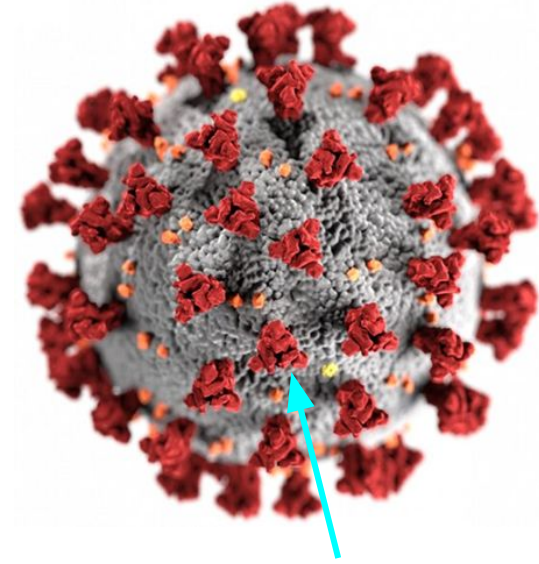
In CS language: protein sequence is a string s of length L over the alphabet Σ of 20 characters



Alphabet Σ : amino acids

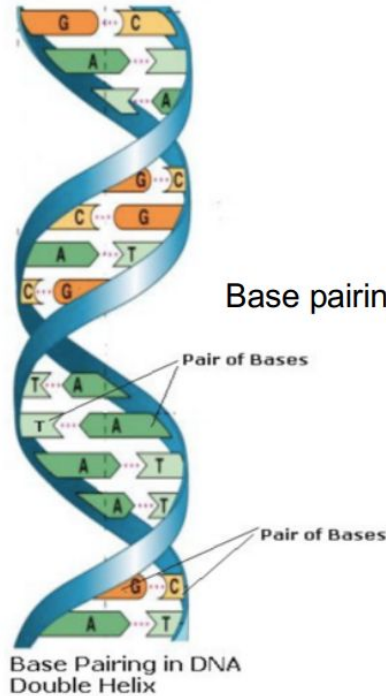
Protein sequence

MFVFLVLLPLVSSQCVNLTTRTQLPPAYTNSFTRGVYYPDKVFRSSVLHSTQDLFLPFFS
NVTWFHAIHVSGTNGTKRFDNVLPLFNDGVYFASTEKSNIIRGWIFGTTLDSKTQSLLIV
NNATNVVIKVCEFQFCNDPFLGVYYHKNNKSWMESEFRVYSSANNCTFEYVSQPFLMDLE
GKQGNFKNLREFVFKNIDGYFKIYSKHTPINLVRDLPGGFSALEPLVDLPIGINITRFQT
LLALHRSYLTPGDSSSGWTAGAAAYYVGYLQPRTFLLKYNENGTITDAVDCALDPLSETK
CTLKSFTVEKGIYQTSNFRVQPTESIVRFPNITNLCPFGEVFNATRFASVYAWNRRKISN
CVADYSVLVNSASFSTFKCYGVSPTKLNLDLCFTNVYADSFVIRGDEVQRQIAPGQTGKIAD
YNYKLPDDFTGCVIAWNSNNLDSKVGNNYLYRLFRKSNLKPFERDISTEIQAGSTPC
NGVEGFNCYFPLQSYGFQPTNGVGYPYRVVLSFELLHAPATVCGPKKSTNLVKNKCVN
FNFNGLTGTGVLTESNKKFLPFQQFGRDIADTTDAVRDPQTLIELDITPCSFGGVSVITP
GTNTSNQVAVLYQDVNCTEVPVAIHADQLTPTWRVYSTGSNVFQTRAGCLIGAEHVNNSY
ECDIPIGAGICASYQTQTSNPRRARSVASQSI IAYTMSLGAENSVAYSNNIAIPTNFTI
SVTTEILPVSMTKTSVDCTMYICGDNTECSNLLLQYGSFCTQLNRALTGIAVEQDKNTQE
VFAQVKQIYKTPPIKDFGGFNFSQILPDPSKPSKRSFIEDLLFNKVTLADAGFIKQYGDC
LGDIAARDLICAQKFNGLTVLPLLTDEMIAQYTSALLAGTITSGWTFGAGAAALQIPFAM
QMAYRFNGIGVTQNVLYENQKLIANQFNSAIGKIQDSLSTASALGKLQDVVNQNAQALN
TLVKQLSSNFGAISSVLNDILSRDLKVEAEVQIDRLITGRLQSLQTYVTQQLIRAAEIRA
SANLAATKMSECVLGQSKRVDFCGKGHYLMSFPQSAPHGVVFLHVTYVPAQEKNTTAPA
ICHDGKAHFPPREGVVFVSNGTHWFVTQRNFYEPQIITDNTFVSGNCDVVIGIVNNTVYDP
LQPELDSFKEELDKYFKNHTSPDVLGDISGINASVVNIQKEIDRLNEVAKNLNESLIDL
QELGKYEQYIKWPWYIWLGFIAGLIAIVMVTIMLCCMTSCCSCCLKGCCSCGSCCKFDEDD



SARS-CoV-2 spike protein

DNA



Base pairing property

=

The DNA Molecule

5'

G -- C
A -- T
T -- A
G -- C
C -- G
G -- C
T -- A
G -- C
T -- A
T -- A
A -- T
A -- T
C -- G
T -- A

3'

DNA sequence

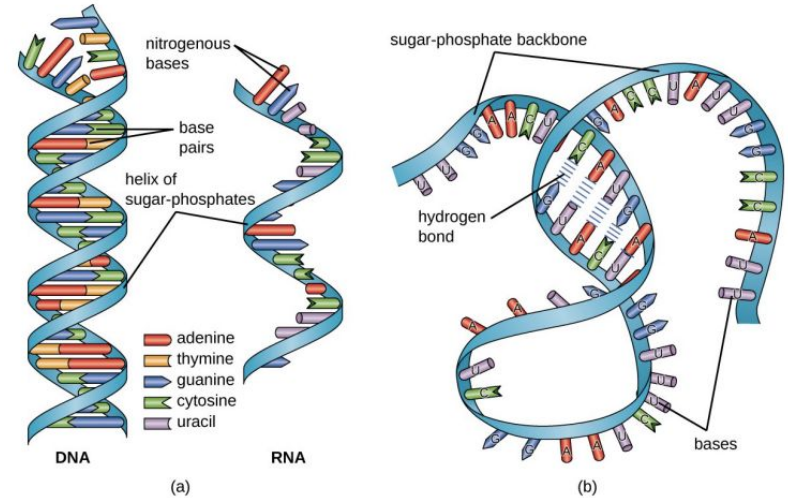
A string \mathbf{s} of length L over the alphabet $\Sigma=\{\text{A,C,G,T}\}$

[illegible]

RNA

RNA = ribonucleic acid

- “U” instead of “T”
- Usually single stranded
- Has base-pairing capability
 - Can form simple non-linear structures
- Life may have started with RNA

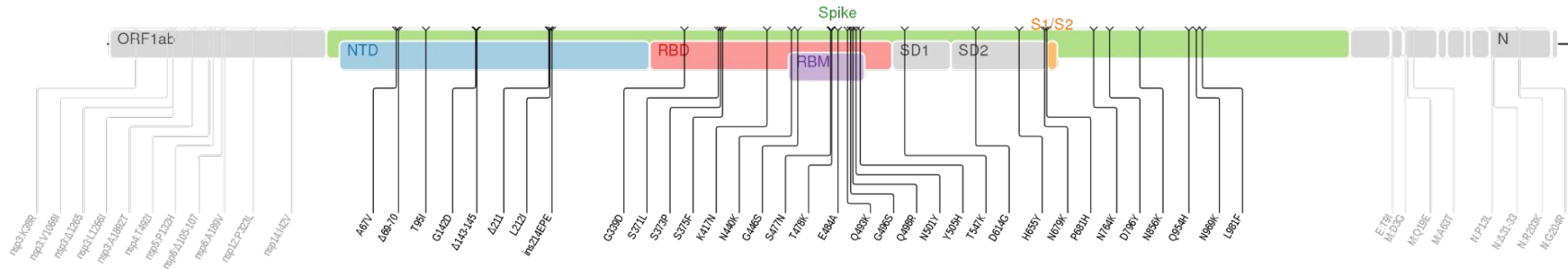


Summary: biological sequences

- DNA = nucleotide sequence
 - Alphabet size = 4 (A,C,G,T)
- RNA (single stranded)
 - Alphabet size = 4 (A,C,G,U)
- Protein sequence
 - Alphabet size = 20

Why study biological sequences

- **DNA** stores the genetic information
- **RNA** is the intermediate molecule for protein synthesis
- **Proteins** carry out the biochemical functions in cell



Example: genomic sequence SARS-CoV-2 Omicron

Protein sequence-structure-function

New Variant Shows More Mutations in Spike Protein

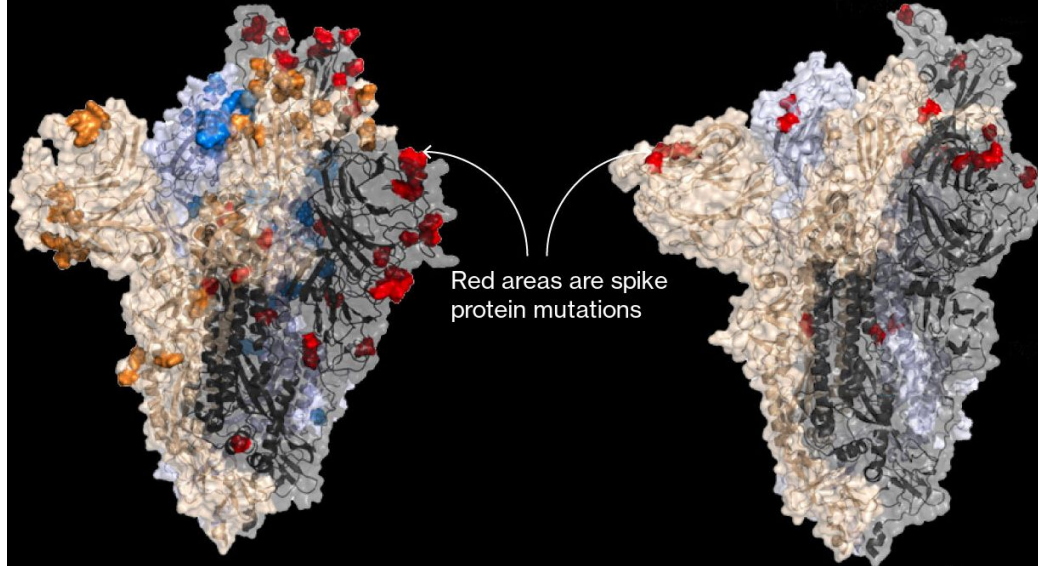
The omicron variant has more than three times as many mutations as delta, but it's too early to know if omicron is more dangerous

Omicron (B.1.1.529)

32 mutations in spike protein

Delta (B.1.617.2)

9 mutations in spike protein

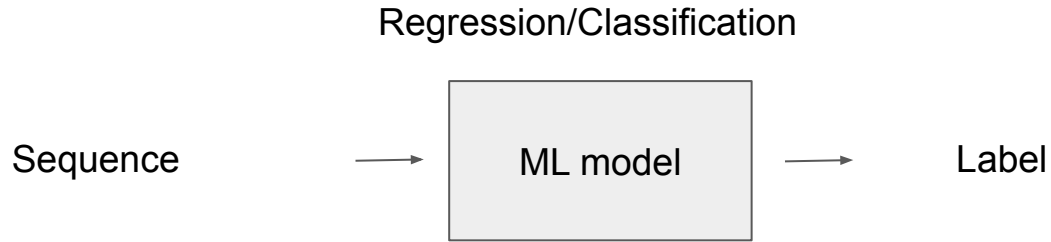


Source: The COVID-19 Genomics UK Consortium

Bloomberg

Deep learning for sequence data

Problem formulation



Example: protein stability prediction

Input	Output
...DNGVDGEWTYDDATKTFTVTE	1.0
...DNGCDGEWTYDDATKTFTVTE	-0.2
...DNGVWGEWTYDDATKTFTVTE	3.9
...DNGVWGEWTYDDATKTFTFTE	5.4
...DNGVMGEWTYDDATKTFTDTE	-0.1

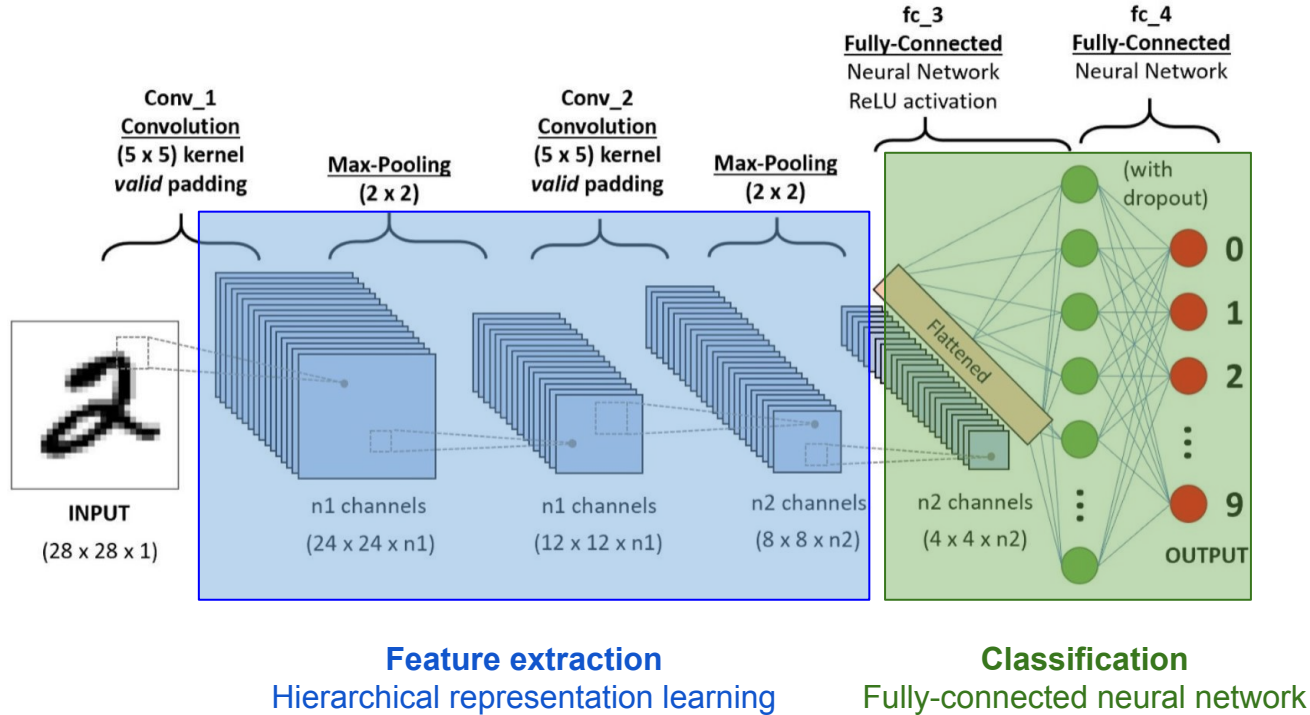
Sequence encoding

- one-hot encoding

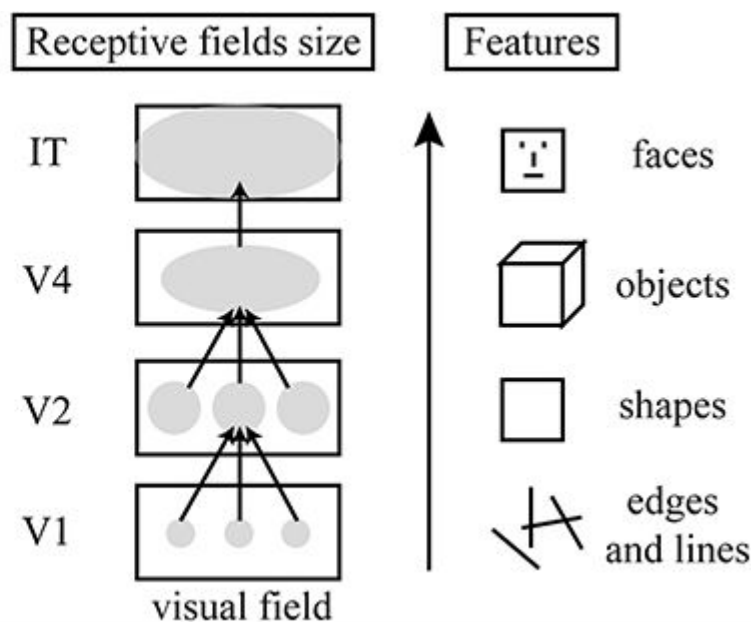
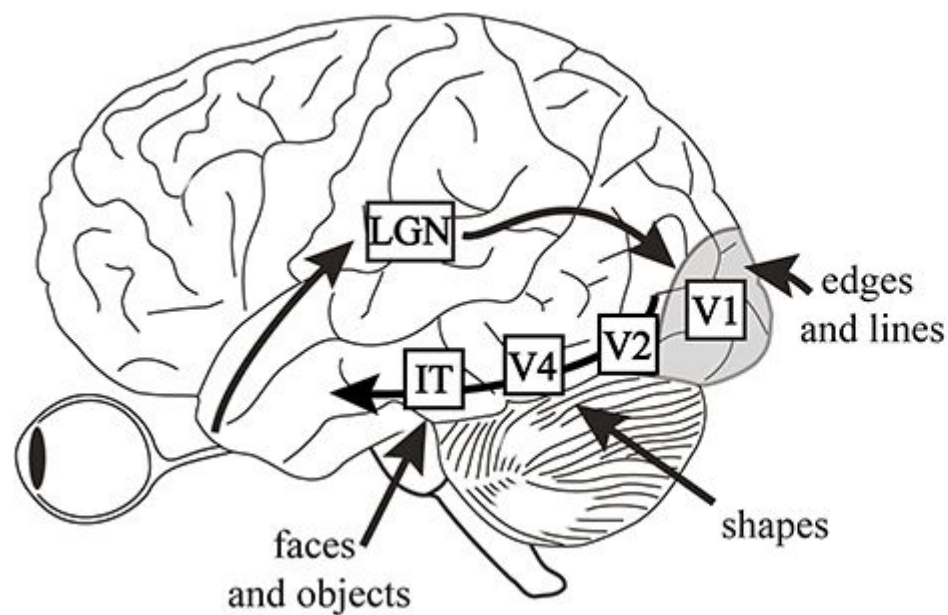
	C	G	A	T	A	A	C	C	G	A	T	A	T
A	0	0	1	0	1	1	0	0	0	1	0	1	0
C	1	0	0	0	0	0	1	1	0	0	0	0	0
G	0	1	0	0	0	0	0	0	1	0	0	0	0
T	0	0	0	1	0	0	0	0	0	0	1	0	1

- contextual embedding (language models)
 - Rives, Alexander, et al. "Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences." *Proceedings of the National Academy of Sciences* 118.15 (2021).

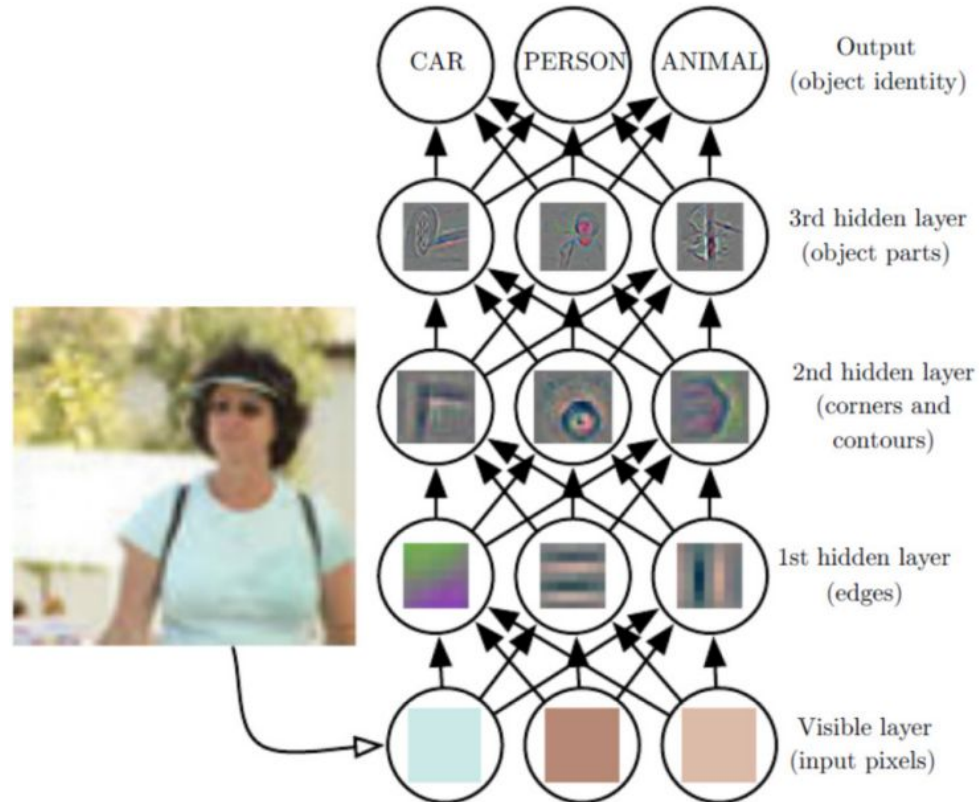
Model #1: Convolutional Neural Network (CNN)



Inspired by human vision



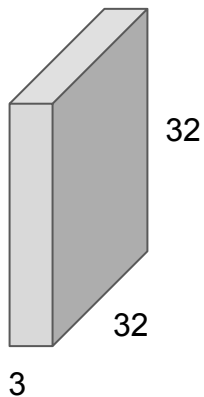
Hierarchical feature extraction in CNN



Convolution layer

Input image

$32 \times 32 \times 3$



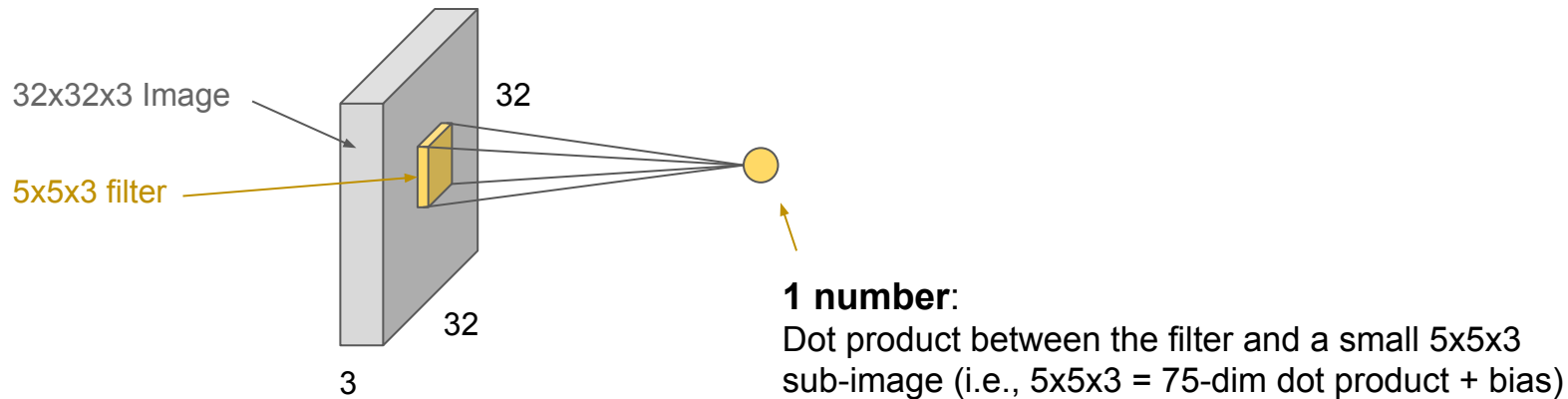
Filter/kernel

$5 \times 5 \times 3$



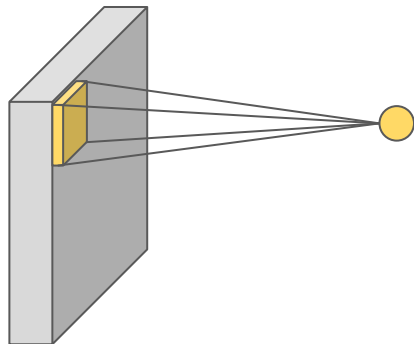
Convolution: slide the filter over the image spatially and compute the dot products

Convolution layer

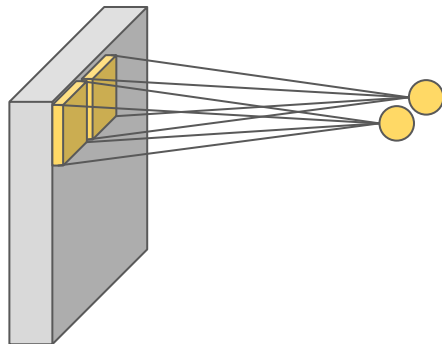


Convolution: slide the filter over the image spatially and compute the dot products

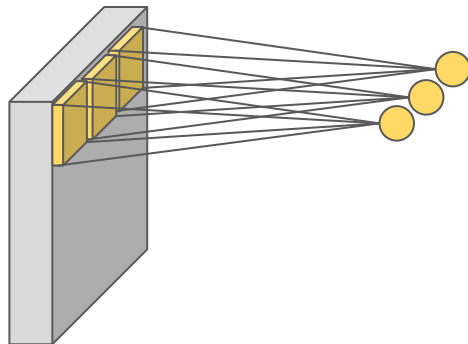
Convolution layer



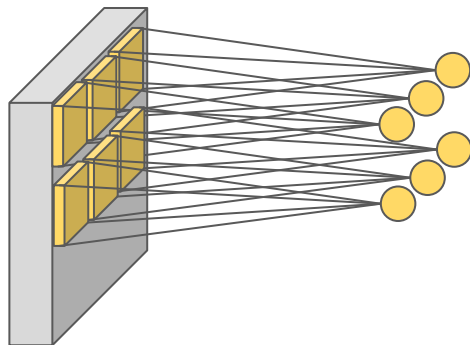
Convolution layer



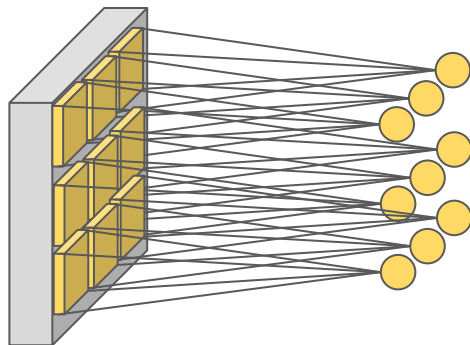
Convolution layer



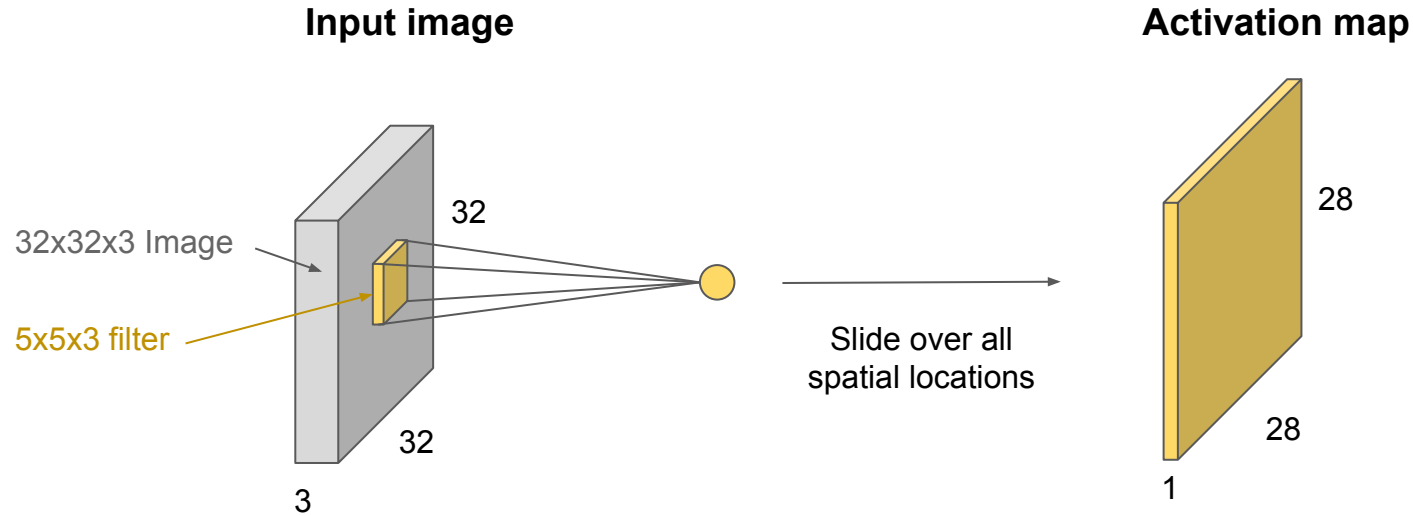
Convolution layer



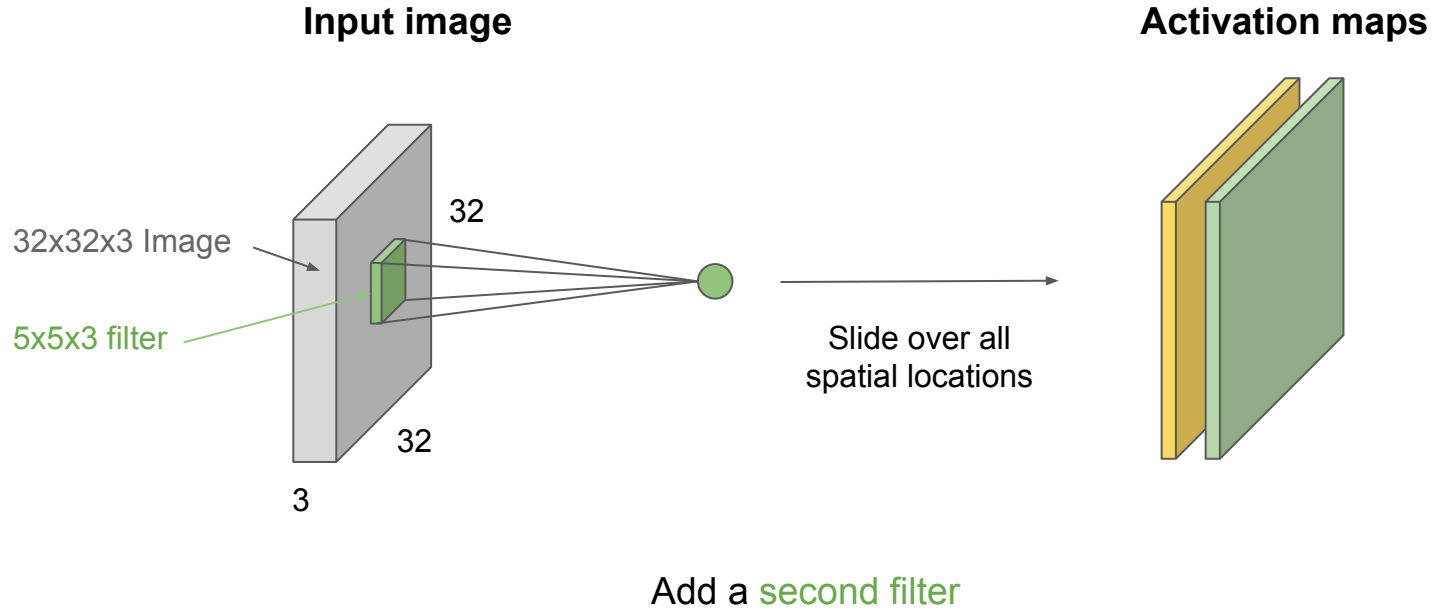
Convolution layer



Convolution layer

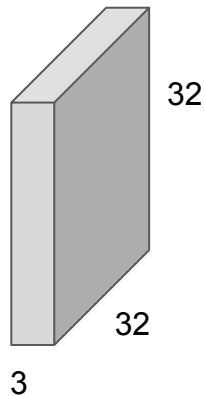


Convolution layer



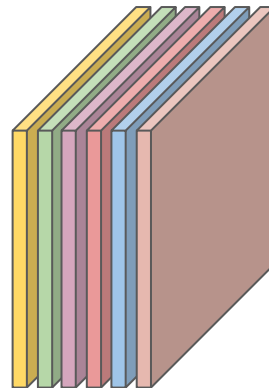
Convolution layer

Input image

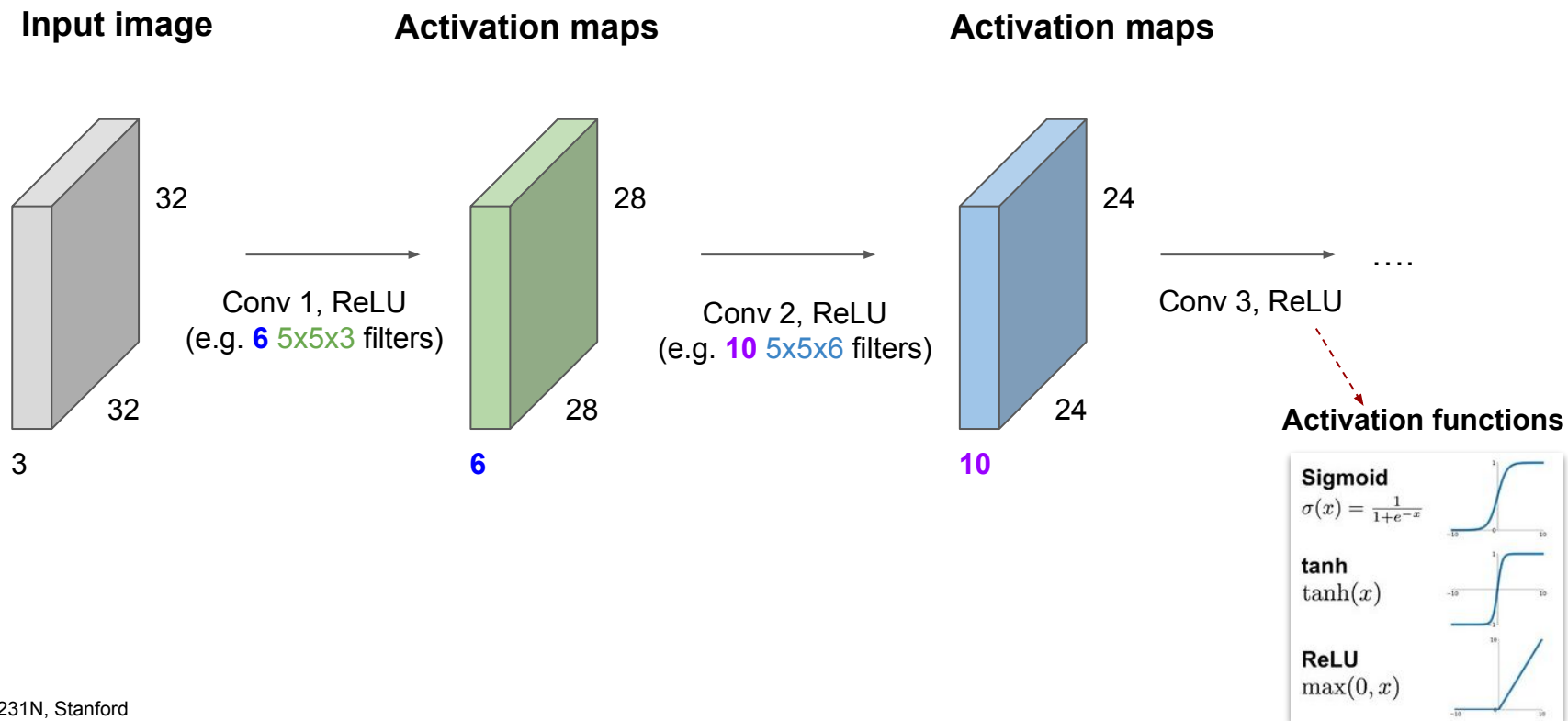


A convolution layer with 6 filters

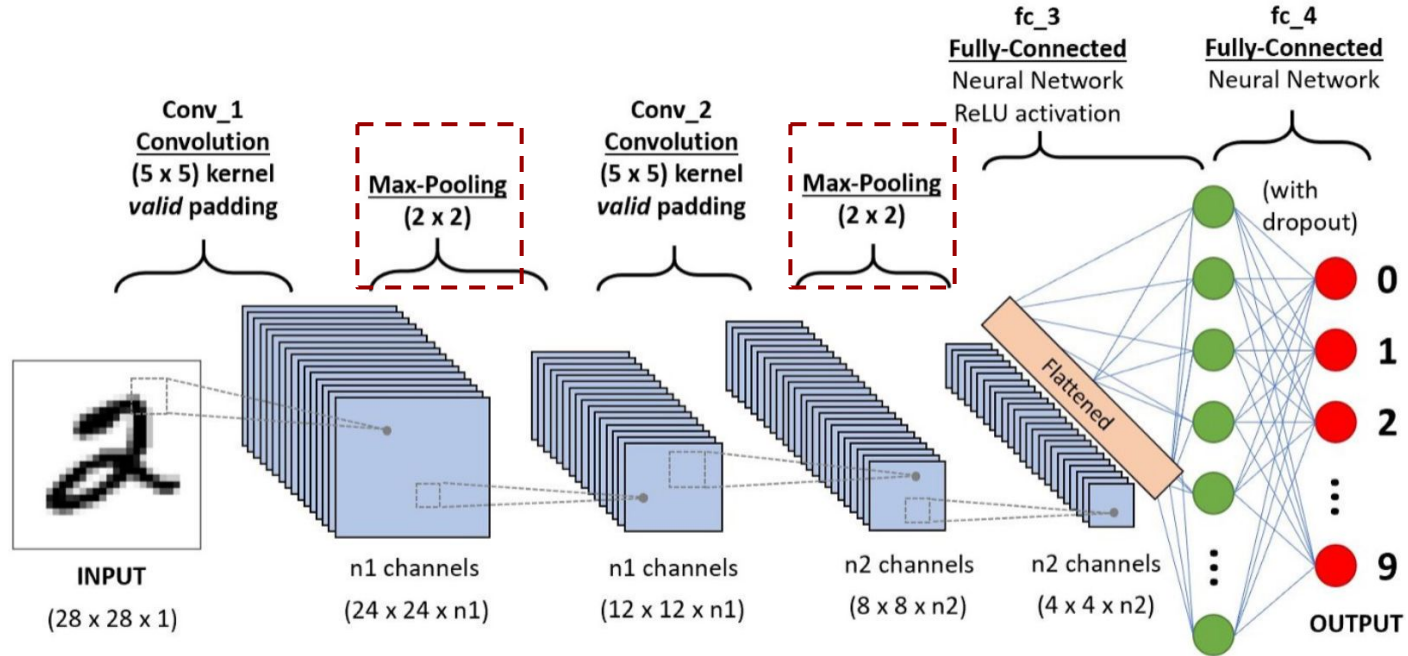
Activation maps



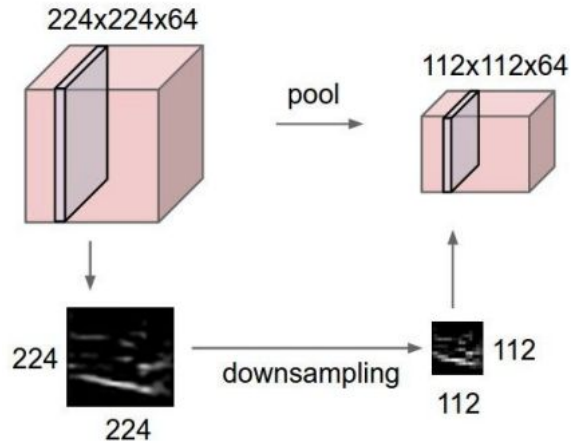
CNN is a sequence of convolution layers, interspersed with activation functions



Pooling layer

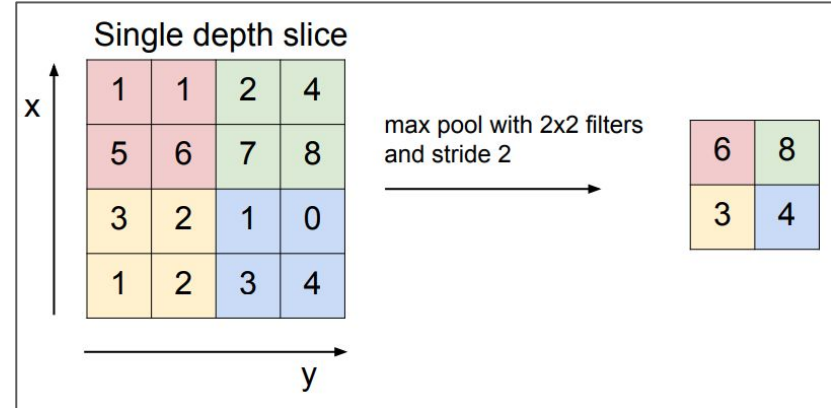


Pooling layer



Pooling layer makes the activation map smaller and more manageable

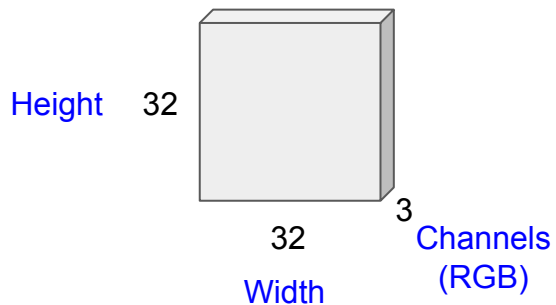
Max pooling



CNN for biological sequences

Image (2D CNN)

Input dimension: e.g., 32x32x3



Sequence (1D CNN)

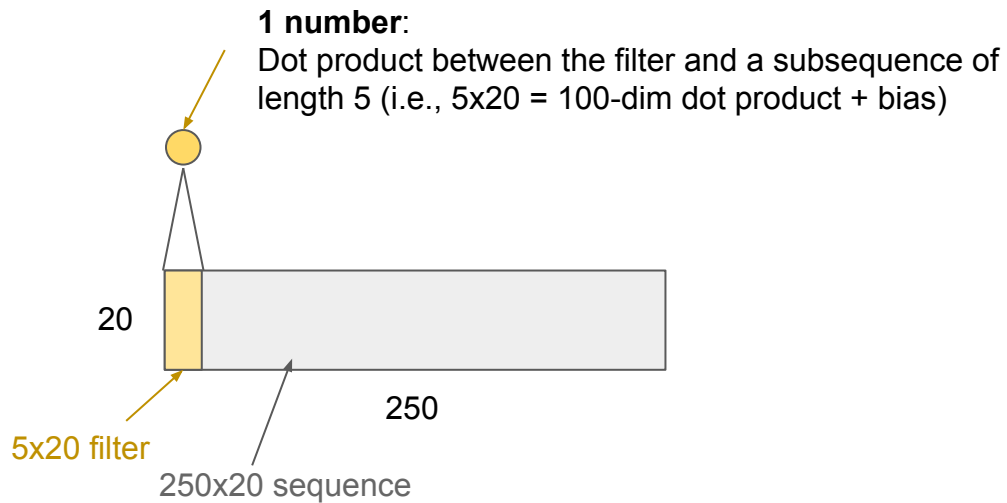
Input dimension: e.g., 250x20

Channels
(20 amino acids) 20

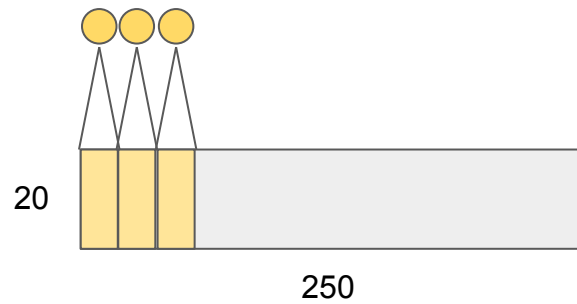


Alphabet size:
4 for DNA/RNA
20 for protein

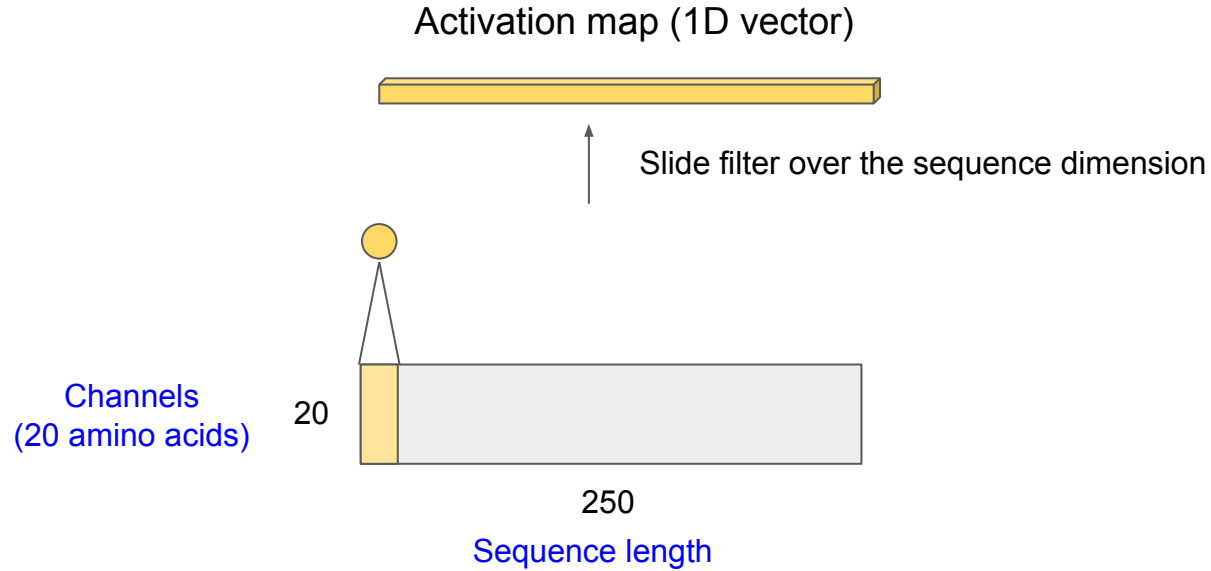
1D convolution



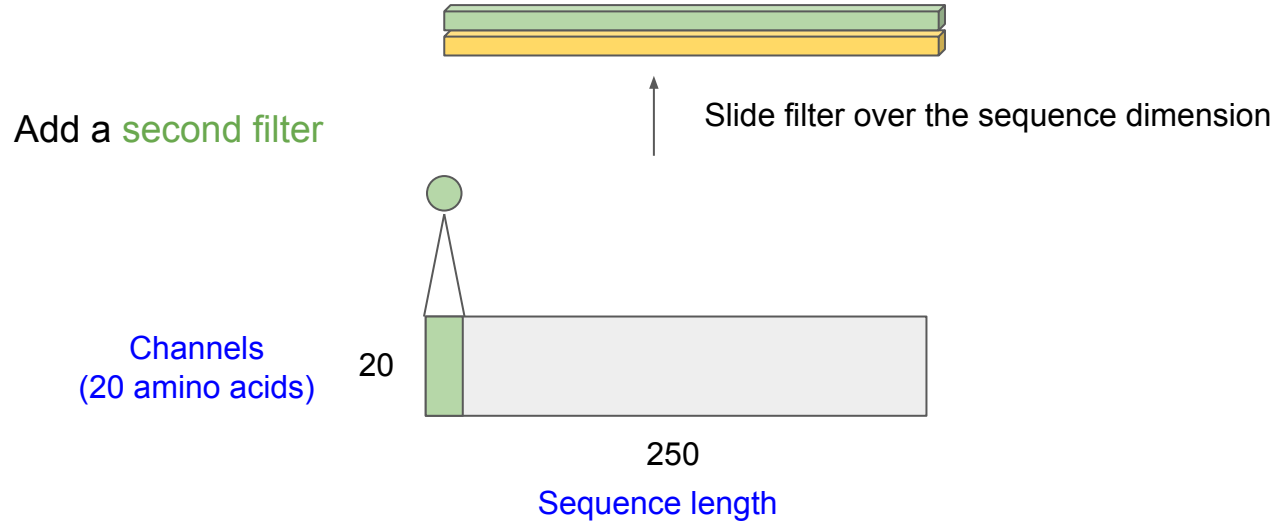
1D convolution



1D convolution



1D convolution



Convolution layers in PyTorch



Convolution Layers

`nn.Conv1d`

Applies a 1D convolution over an input signal composed of several input planes.

`nn.Conv2d`

Applies a 2D convolution over an input signal composed of several input planes.

`nn.Conv3d`

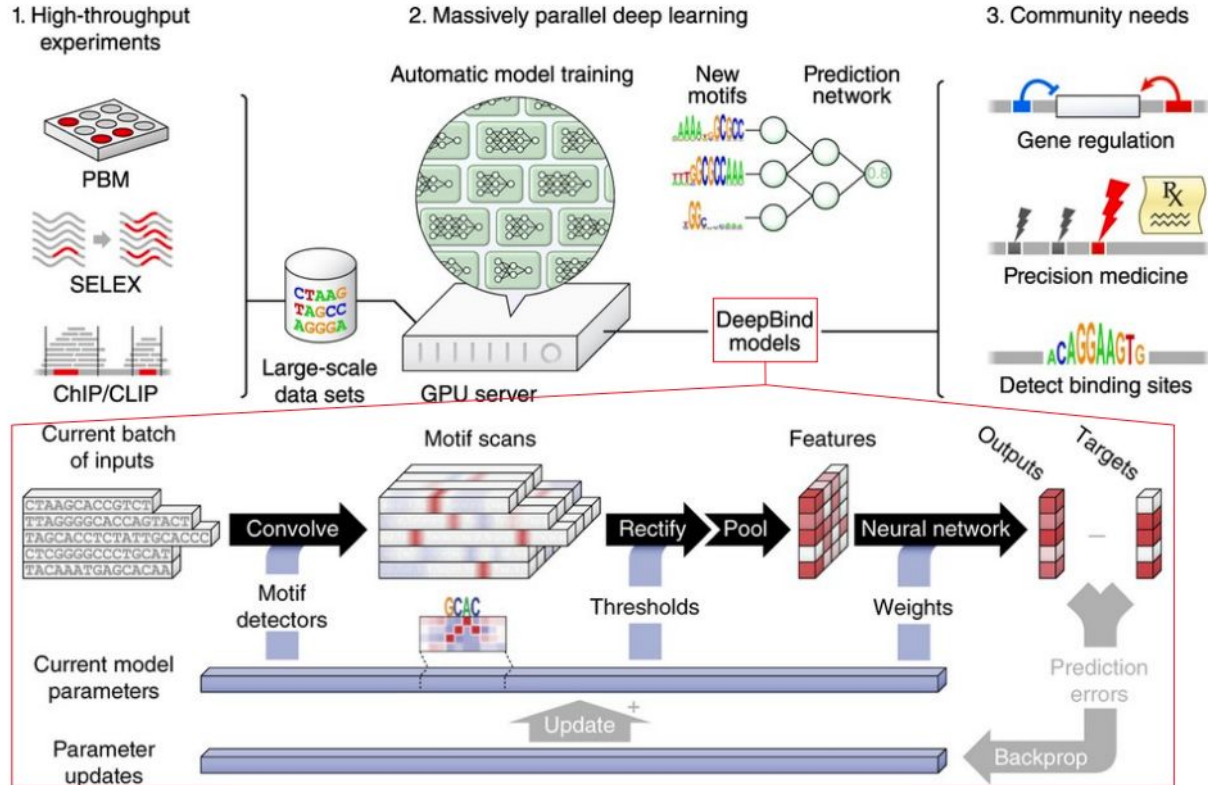
Applies a 3D convolution over an input signal composed of several input planes.

An example of 2-layer 2D CNN for RGB image

```
class Net(nn.Module):
    def __init__(self):
        super().__init__()
        self.conv1 = nn.Conv2d(3, 6, 5)
        self.pool = nn.MaxPool2d(2, 2)
        self.conv2 = nn.Conv2d(6, 16, 5)
        self.fc1 = nn.Linear(16 * 5 * 5, 120)
        self.fc2 = nn.Linear(120, 84)
        self.fc3 = nn.Linear(84, 10)

    def forward(self, x):
        x = self.pool(F.relu(self.conv1(x)))
        x = self.pool(F.relu(self.conv2(x)))
        x = torch.flatten(x, 1) # flatten all dimensions except batch
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = self.fc3(x)
        return x
```

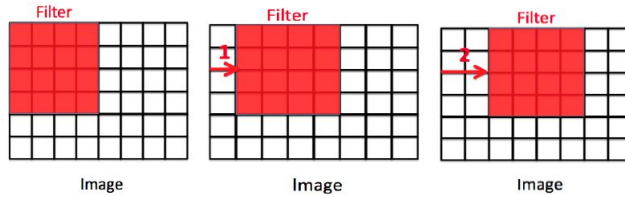
Example: predicting DNA/RNA-protein binding



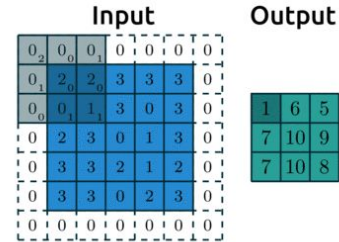
Further readings of CNN

Other important concepts of CNN:

- *Deep Learning*, [Chapter 5](#)



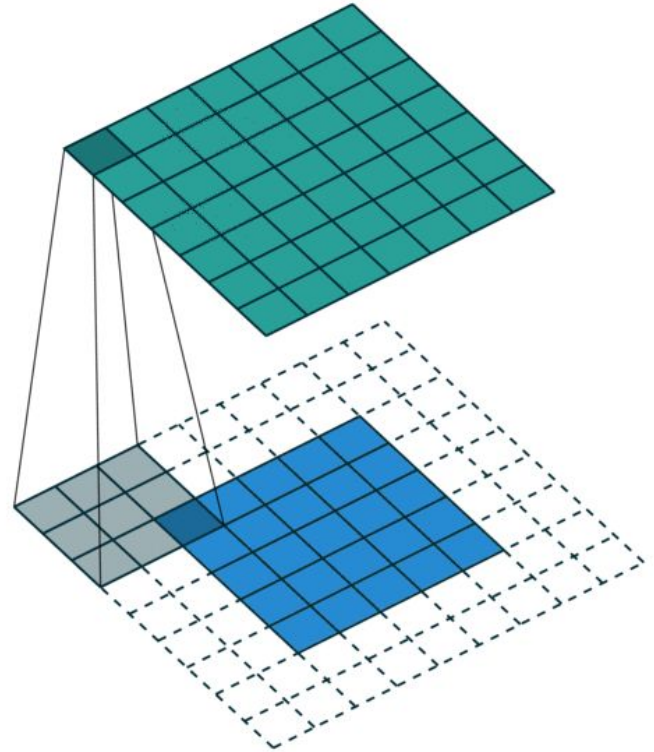
Stride



Padding

Summary: key ideas of CNNs

- Use a set of filters to extract **local** features
- Use **multiple** filters to extract different features
- Spatially **share** parameters of each filter



ResNet

Deep Residual Learning for Image Recognition

Kaiming He Xiangyu Zhang Shaoqing Ren Jian Sun

Microsoft Research

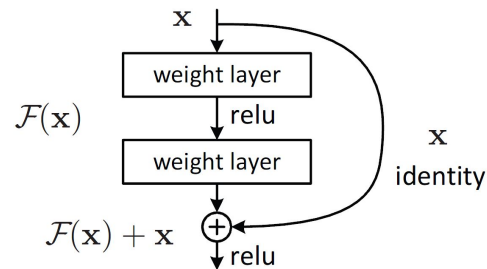
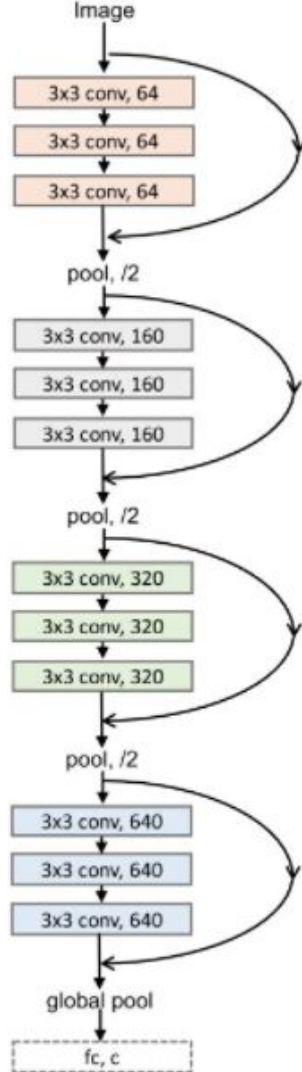
{kahe, v-xiangz, v-shren, jiansun}@microsoft.com

Deep residual learning for image recognition

[K He](#), [X Zhang](#), [S Ren](#), [J Sun](#) - ... and pattern recognition, 2016 - openaccess.thecvf.com

Deeper neural networks are more difficult to train. We present a residual learning framework to ease the training of networks that are substantially deeper than those used previously. We explicitly reformulate the layers as learning residual functions with reference to the layer ...

☆ Save 📄 Cite Cited by 103017 Related articles All 63 versions 🔗



Model #2: Recurrent Neural Network (RNN)

Example: sequence labeling problems

- Part of speech

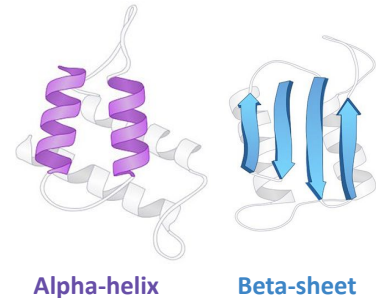
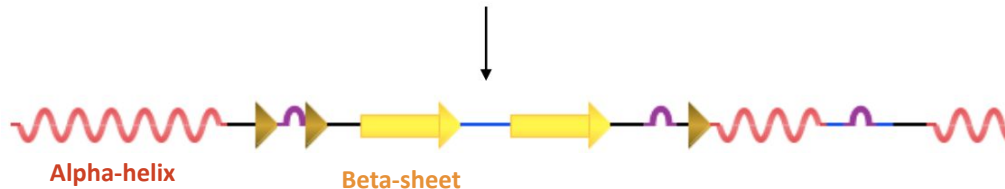
Article Noun Verb Preposition Article Noun
| | | | |
The cat sat on the mat.

- Handwriting recognition

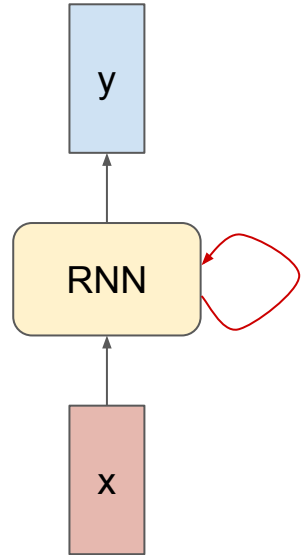
Foreign Minister. → FOREIGN MINISTER.

- Protein secondary structure prediction

NKEILDEAYVMASVDNPHVCRLLGICLTSTVQLITQLMPFGCLLDYVREHKDNIGSQYLL

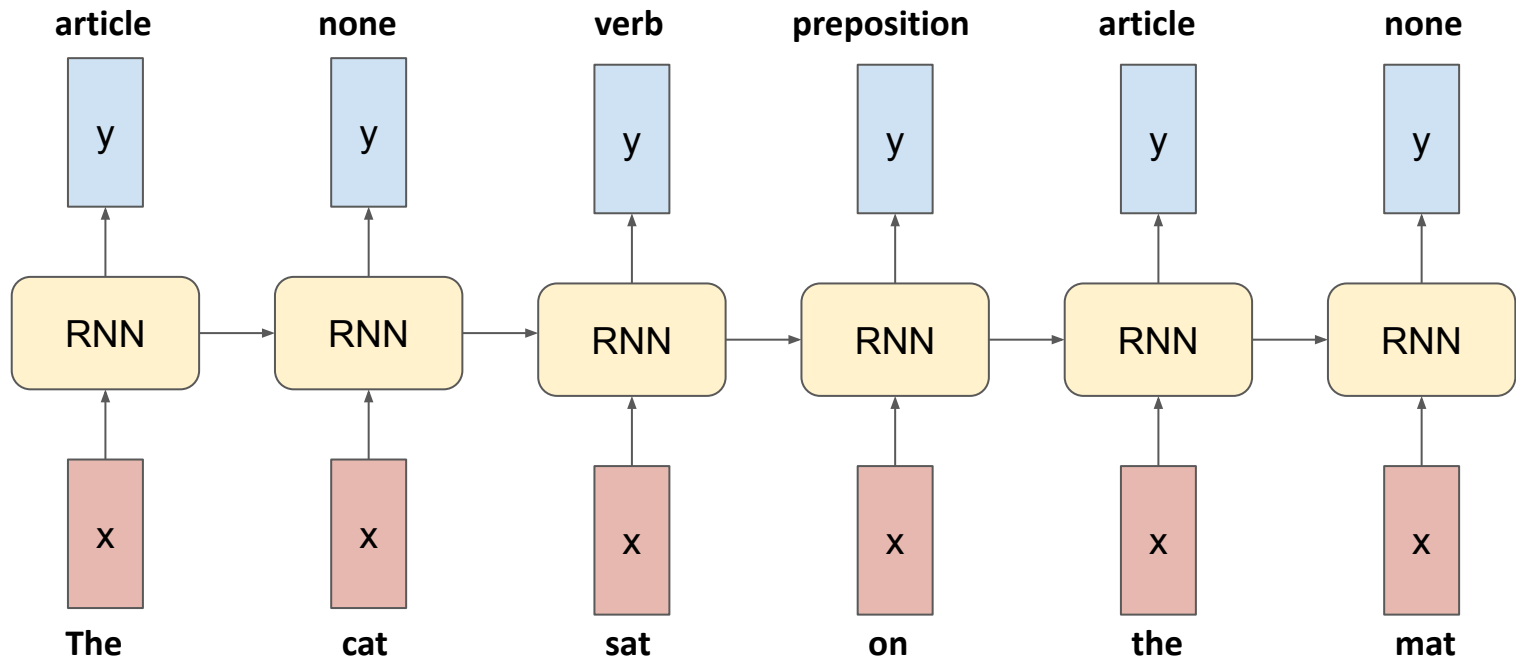


Recurrent Neural Network (RNN)

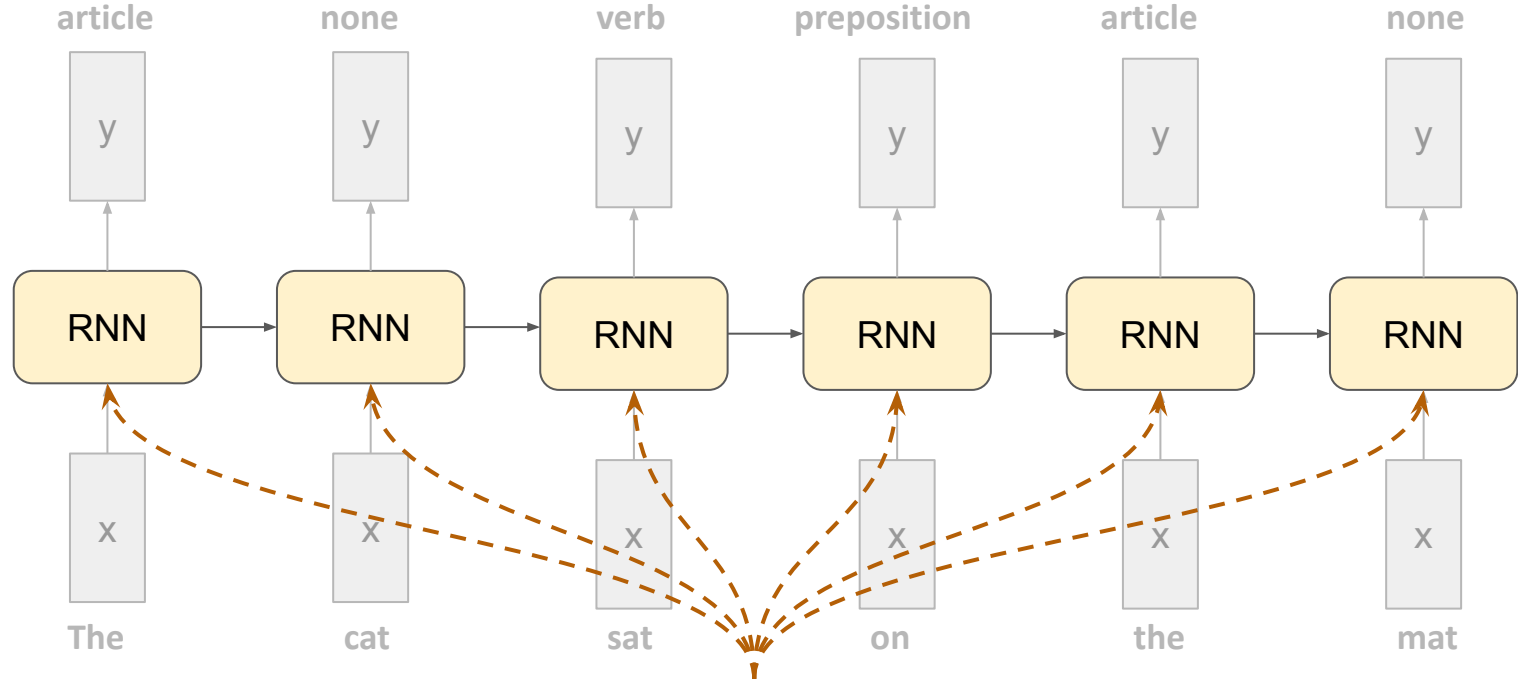


Key idea: RNNs have an “internal state” that is updated as the input sequence is processed

Unrolled RNN



Unrolled RNN

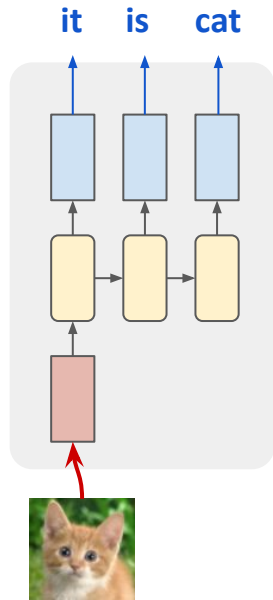


The same set of function and the same set of parameters are used at every time step

RNNs for sequence processing

one to many

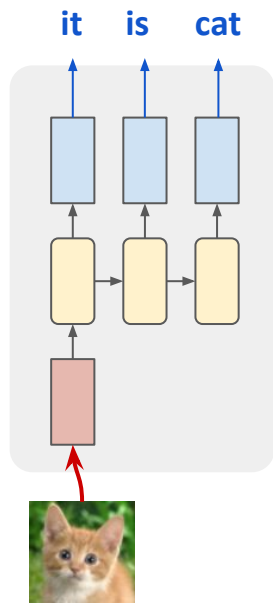
(e.g., image
captioning)



RNNs for sequence processing

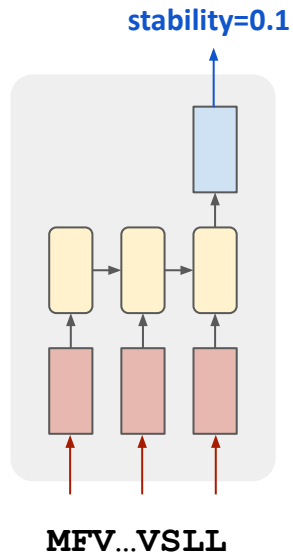
one to many

(e.g., image captioning)



many to one

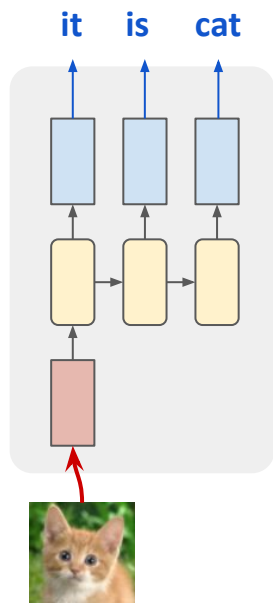
(e.g., protein function prediction)



RNNs for sequence processing

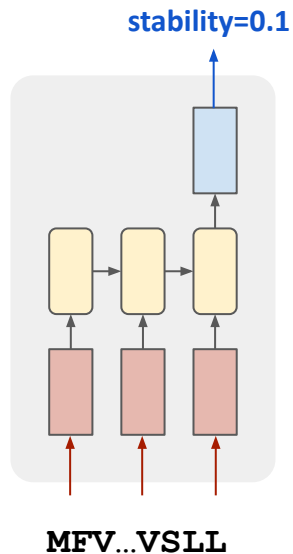
one to many

(e.g., image captioning)



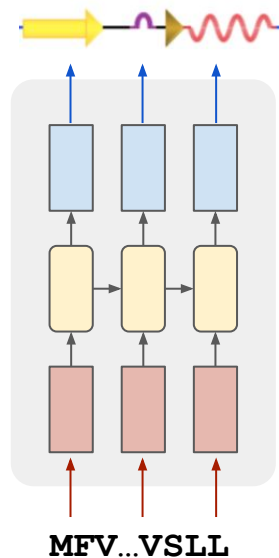
many to one

(e.g., protein function prediction)



many to many

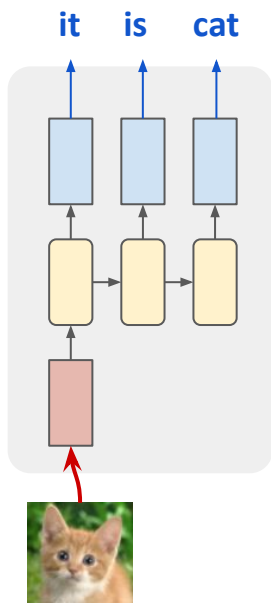
(e.g., protein structure prediction)



RNNs for sequence processing

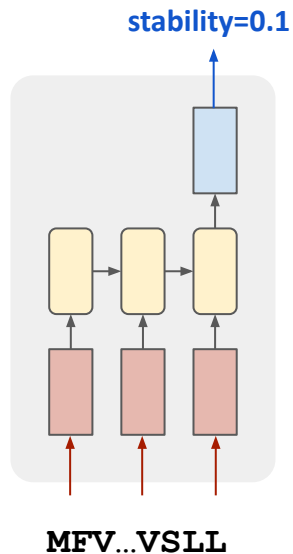
one to many

(e.g., image captioning)



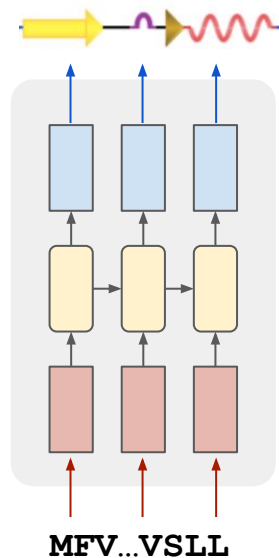
many to one

(e.g., protein function prediction)



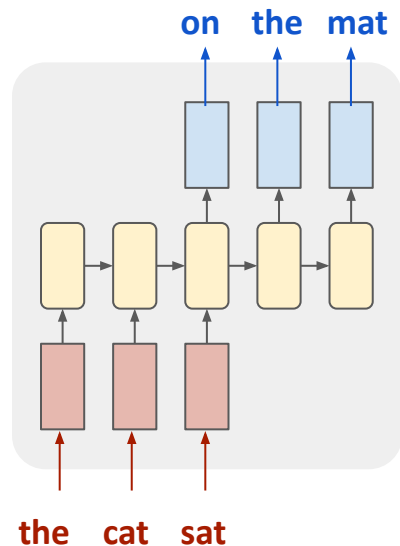
many to many

(e.g., protein structure prediction)



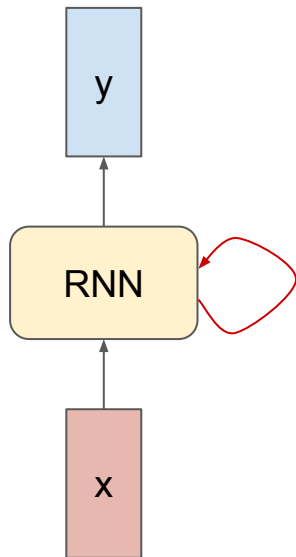
many to many

(e.g., auto completion)



RNN hidden state update

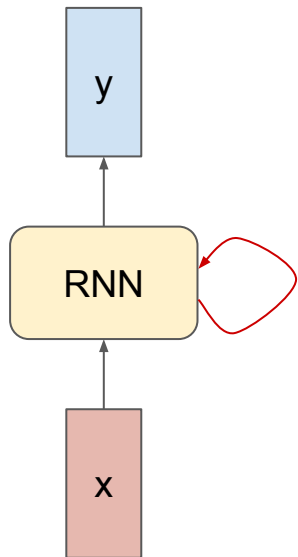
At every step t , the hidden state is updated based on the previous state and the current input



$$\boxed{h_t} = \boxed{f_W}(\boxed{h_{t-1}}, \boxed{x_t})$$

new state function with parameter W previous state Input vector at time step t

RNN output



At every step t , the output is generated based on the current state

Another function with
parameter θ

$$\boxed{y_t} = \boxed{f_{\theta}}(\boxed{h_t})$$

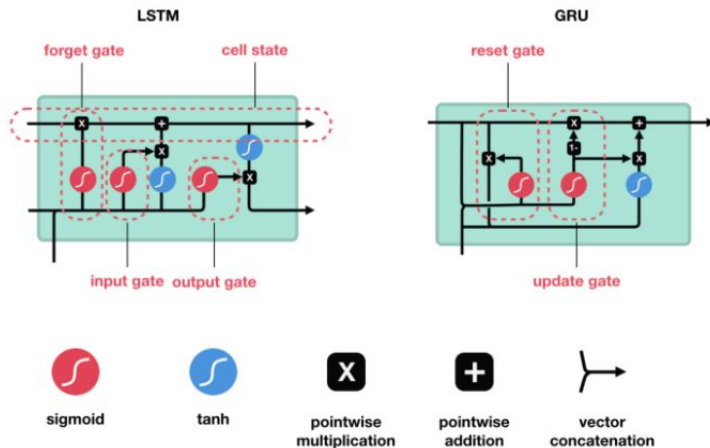
new state new state

Further readings of RNN

- *Deep Learning*, [Chapter 6](#)
- What is function $f()$?
 - $f()$ is usually called “unit” in RNN
 - It defines a “computational graph” the produces h_t based on h_{t-1} and x_t
- Popular RNN variants
 - Long short-term memory (LSTM)
 - Gated recurrent unit (GRU)

function with
parameter W

$$h_t = f_W(h_{t-1}, x_t)$$



[\(image source\)](#)

Recurrent layers in PyTorch



Recurrent Layers

`nn.RNNBase`

`nn.RNN`

Applies a multi-layer Elman RNN with `tanh` or `ReLU` non-linearity to an input sequence.

`nn.LSTM`

Applies a multi-layer long short-term memory (LSTM) RNN to an input sequence.

`nn.GRU`

Applies a multi-layer gated recurrent unit (GRU) RNN to an input sequence.

`nn.RNNCell`

An Elman RNN cell with `tanh` or `ReLU` non-linearity.

`nn.LSTMCell`

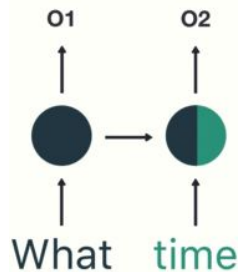
A long short-term memory (LSTM) cell.

`nn.GRUCell`

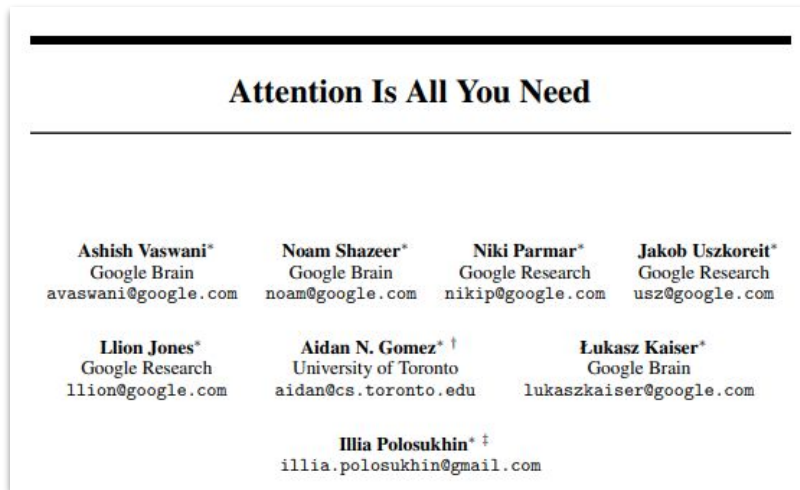
A gated recurrent unit (GRU) cell

Summary: key ideas of RNNs

- Process sequence data with **variable lengths**
 - DNA/RNA/protein sequences, text, audio, time series data
- Capture **sequential** (temporal) information/dependencies in the data
- Parameters **shared** over time steps
- Common to use LSTM or GRU

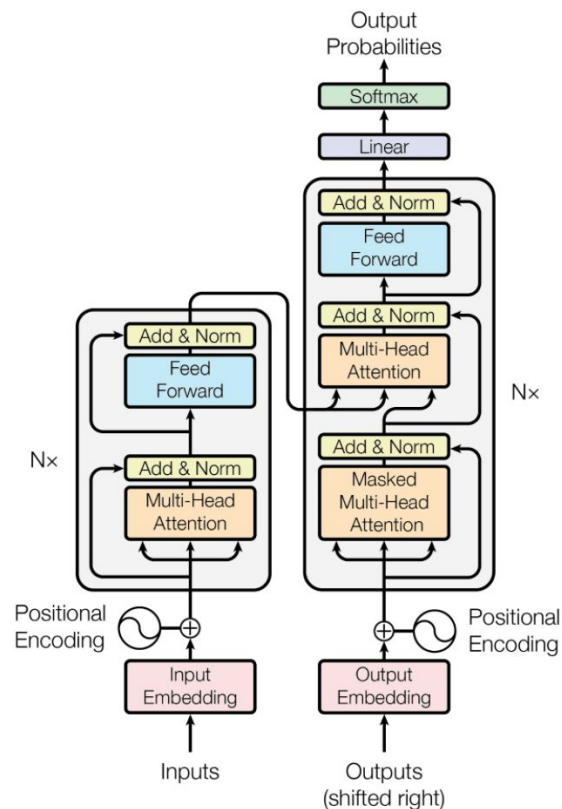


Model #3: Transformers



(NeurIPS 2017)

- Encoder-Decoder
- Sequence-to-sequence
- Transforms one sequence into another sequence, using full context of each



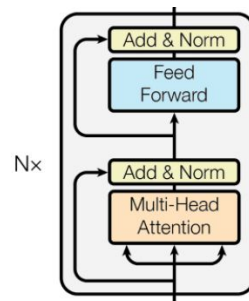
Building blocks of Transformer

N blocks, each has

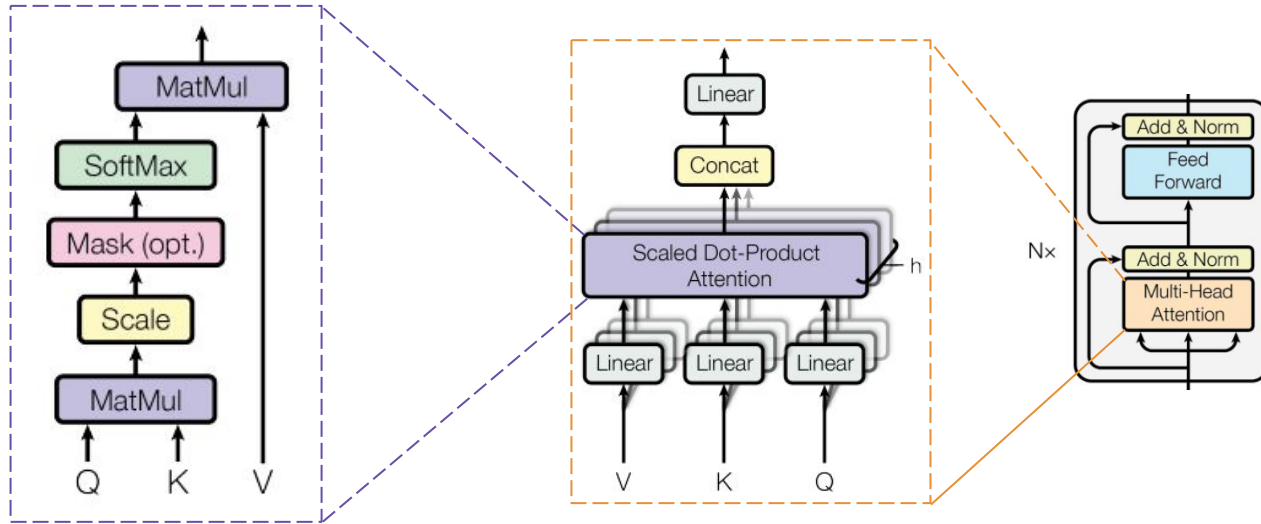
- Multi-head self-attention layer
- Two-layer feed-forward neural nets

Residual connection and layer normalization are used

- Reading: LayerNorm (<https://arxiv.org/pdf/1607.06450.pdf>)

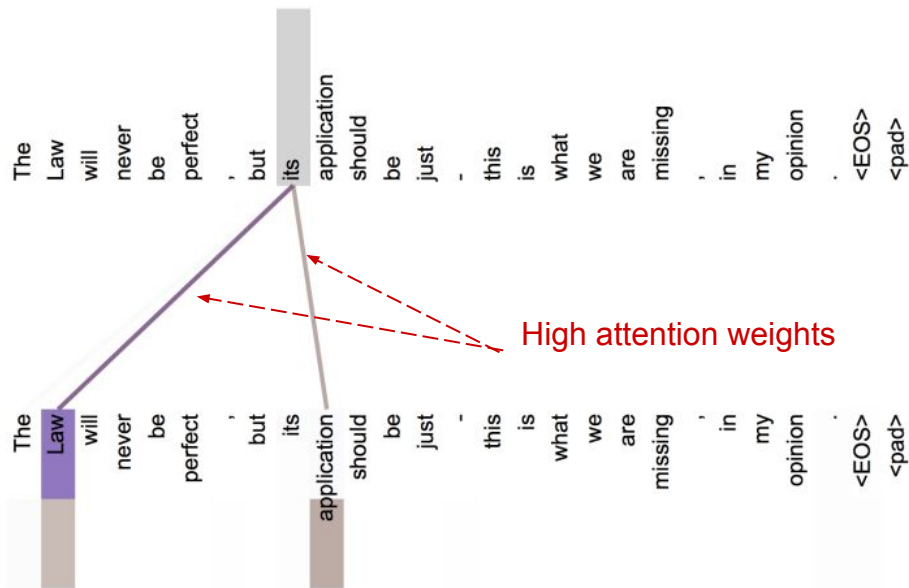


Building blocks of Transformer



Key ideas: self-attention layer

- **Attention layer:** a layer to learn the dependency between words in the input. The dependency is quantified using “attention weights”
- For each word, a new representation is computed by weight-averaging the old representations of all words, where the weight is the learned attention weight



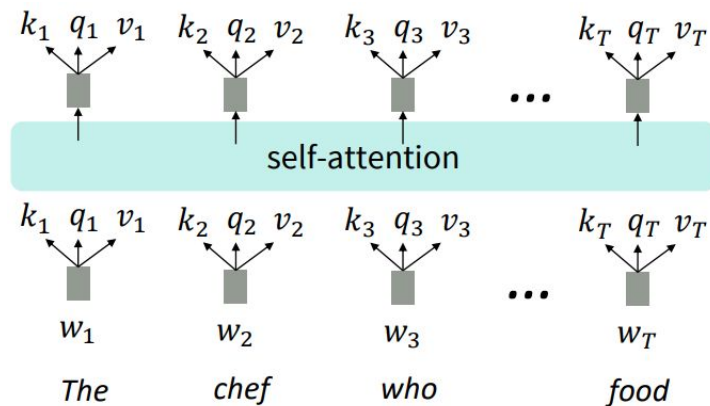
Key ideas: self-attention layer

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Q: “query” matrix, a vector representation for each word

K: “key” matrix, a vector representation for each word

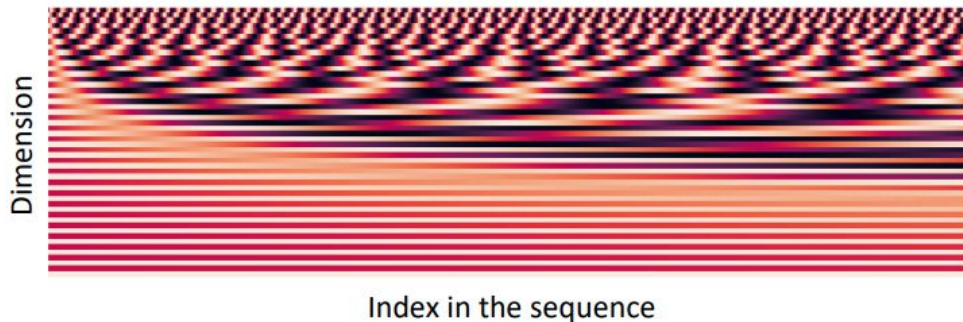
V: “value” matrix, a vector representation for each word



Key ideas: positional encoding

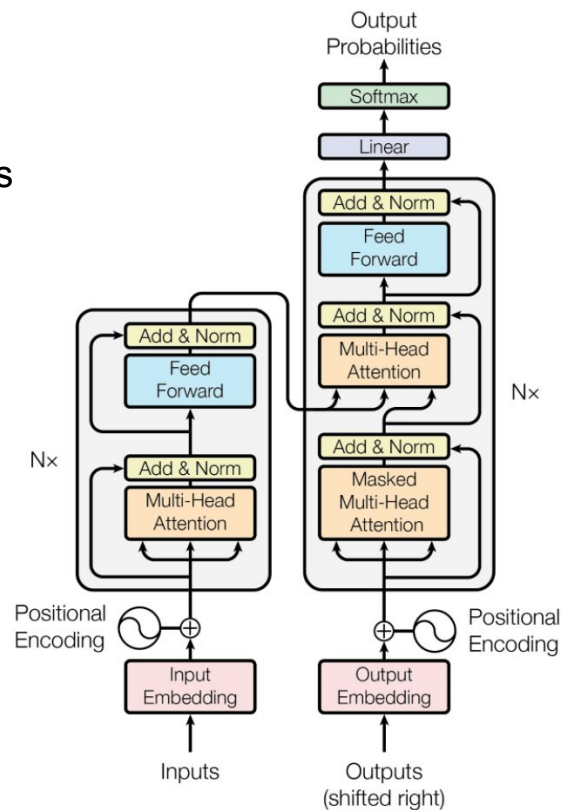
- Self-attention does not know the order of input words
- Positional encodings are added to the word representations, so same words at different locations have different overall representations

$$PE_{(pos, 2i)} = \sin(pos/10000^{2i/d_{\text{model}}})$$
$$PE_{(pos, 2i+1)} = \cos(pos/10000^{2i/d_{\text{model}}})$$



Summary of Transformer

- Learning temporal relationships without unrolling and without RNNs
- Encoder/Decoder framework, multi-head self-attention modules
- Widely used in state-of-the-art NLP models
- Readings:
 - “Attention is all you need” (<https://arxiv.org/abs/1706.03762>)
 - PyTorch [implementation](#) and [tutorial](#) of Transformer

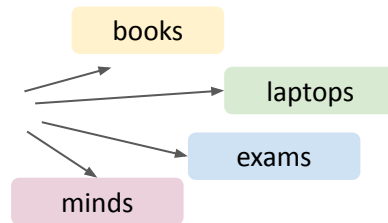


Language modeling of protein sequence

Language modeling in natural language

- **Language modeling** is the task of predicting what word comes next

“The students opened their _____”



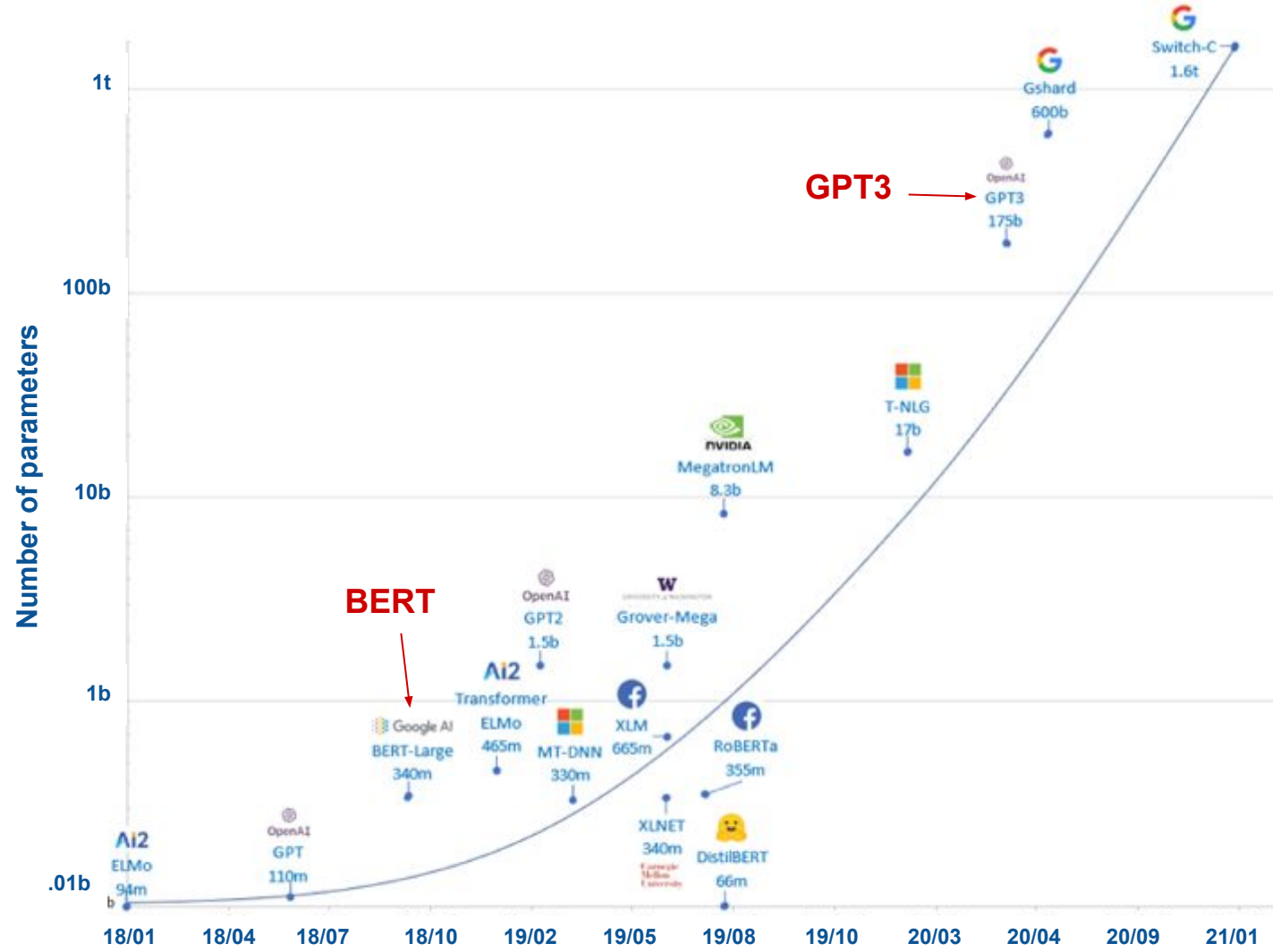
- More formally: given a sequence of words $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)}$, compute the probability distribution of the next word $\mathbf{x}^{(t+1)}$:

$$P(\mathbf{x}^{(t+1)} \mid \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})$$

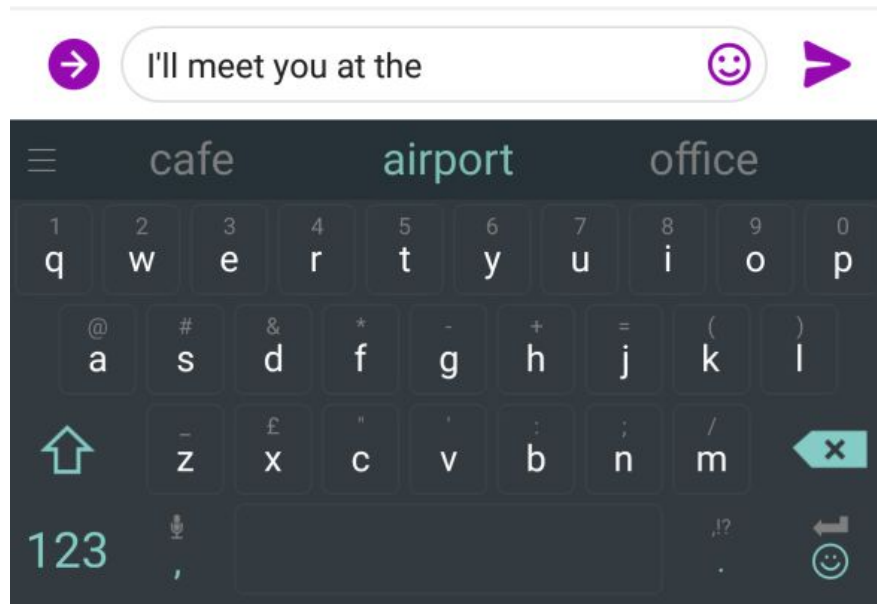
where $\mathbf{x}^{(t+1)}$ can be any word in the vocabulary $V = \{\mathbf{w}_1, \dots, \mathbf{w}_{|V|}\}$.

- A system that does this is called a **language model (LM)**

Large LMs




You use LM every day!



You use LM every day!



what is the | 

what is the **weather**
what is the **meaning of life**
what is the **dark web**
what is the **xfl**
what is the **doomsday clock**
what is the **weather today**
what is the **keto diet**
what is the **american dream**
what is the **speed of light**
what is the **bill of rights**

You use LM every day!

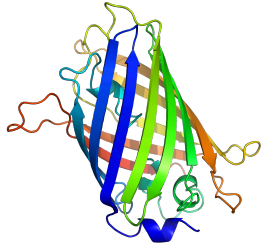
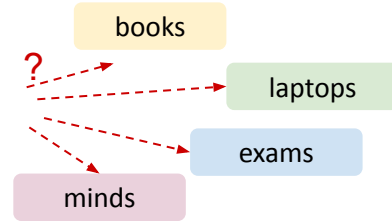


GitHub Copilot

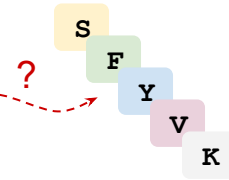
<https://copilot.github.com/>

Learning the language of proteins

The students opened their _____

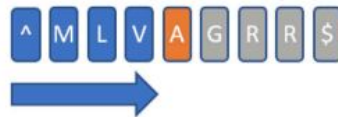
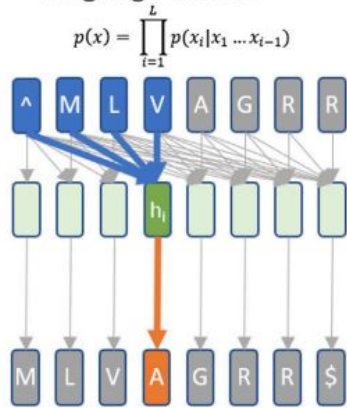


MSKGEELFTGVVPILVELDGDVNGHKFSV_



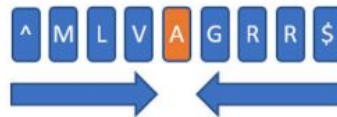
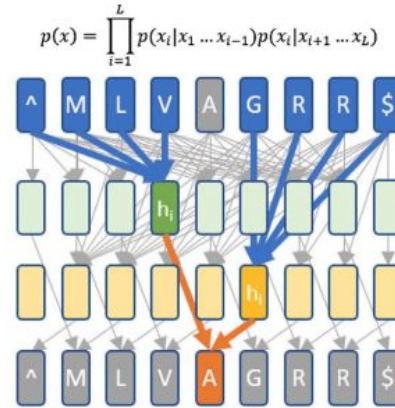
Different language modeling approaches

A Autoregressive language model



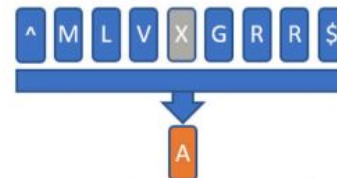
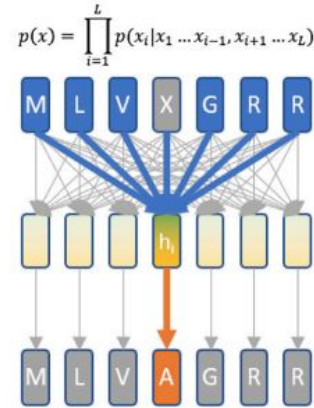
$$p(x_i = A | x_1 \dots x_{i-1})$$

B Bidirectional language model



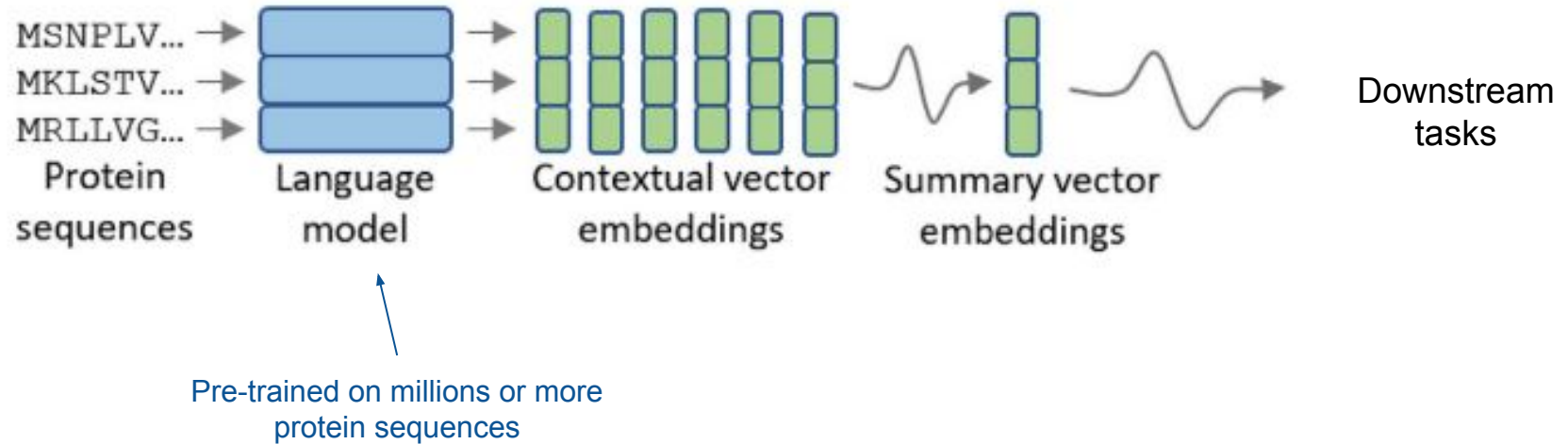
$$p(x_i = A | x_1 \dots x_{i-1}) p(x_i = A | x_{i+1} \dots x_L)$$

C Masked language model

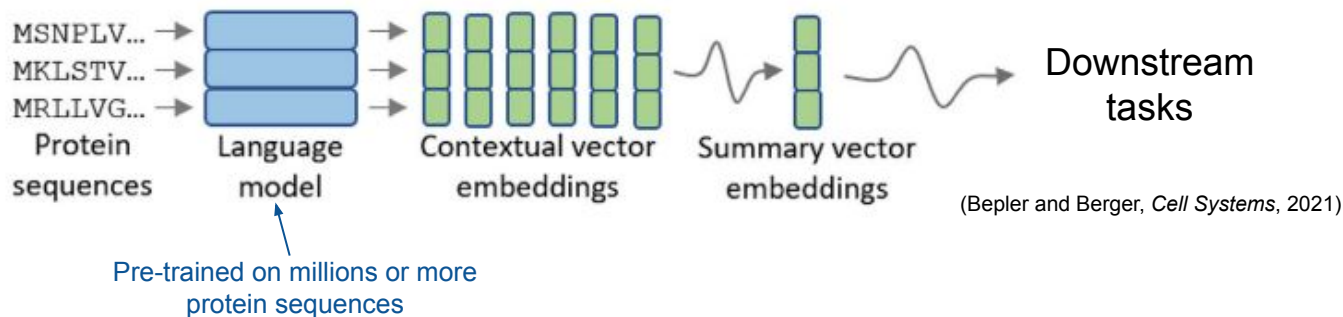


$$p(x_i = A | x_1 \dots x_{i-1}, x_{i+1} \dots x_L)$$

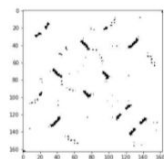
Protein language models



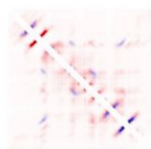
Application: improving downstream prediction tasks



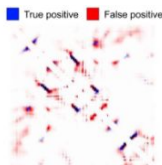
Protein contact prediction
(Rao et al., NeurIPS, 2019)



Ground Truth

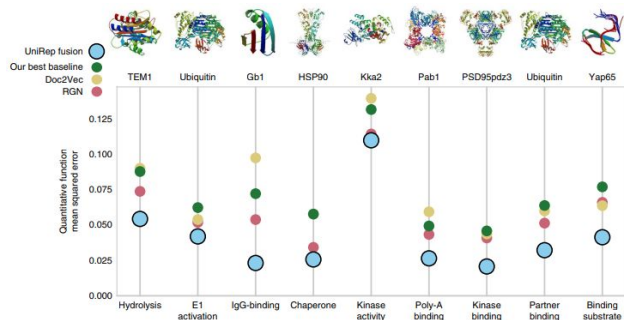


No Pretraining



Pretrained

Protein function prediction
(Alley et al., Nature Methods, 2019)

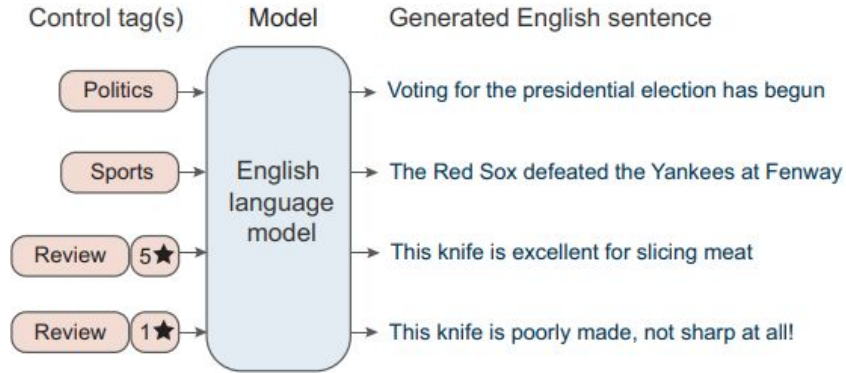


Remote homology detection
(Rives et al., PNAS, 2021)

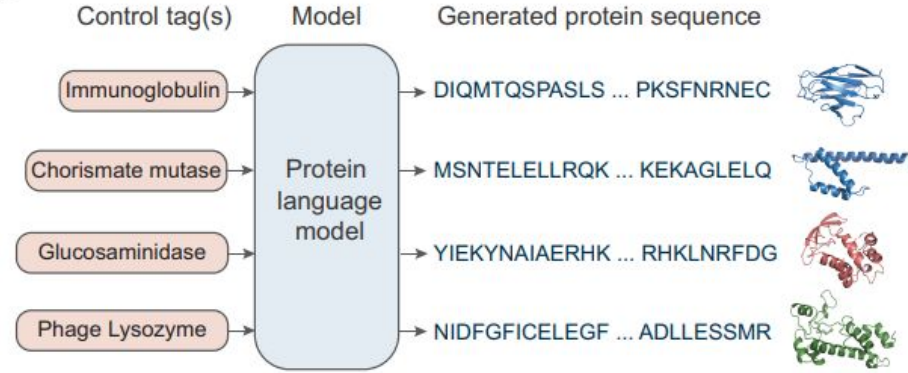
	Pretraining	Hit-10		AUC	
		Fold	SF	Fold	SF
HHblits*		0.584	0.965	0.831	0.951
LSTM (S)	UR50/S	0.558	0.760	0.801	0.863
LSTM (L)	UR50/S	0.574	0.813	0.805	0.880
Transformer-6	UR50/S	0.653	0.878	0.768	0.901
Transformer-12	UR50/S	0.639	0.915	0.778	0.942
Transformer-34	(None)	0.481	0.527	0.755	0.807
Transformer-34	UR100	0.599	0.841	0.753	0.876
Transformer-34	UR50/D	0.617	0.932	0.822	0.932
Transformer-34	UR50/S	0.639	0.931	0.825	0.933
ESM-1b	UR50/S	0.532	0.913	0.770	0.880

Application: generating novel functional protein sequences

a

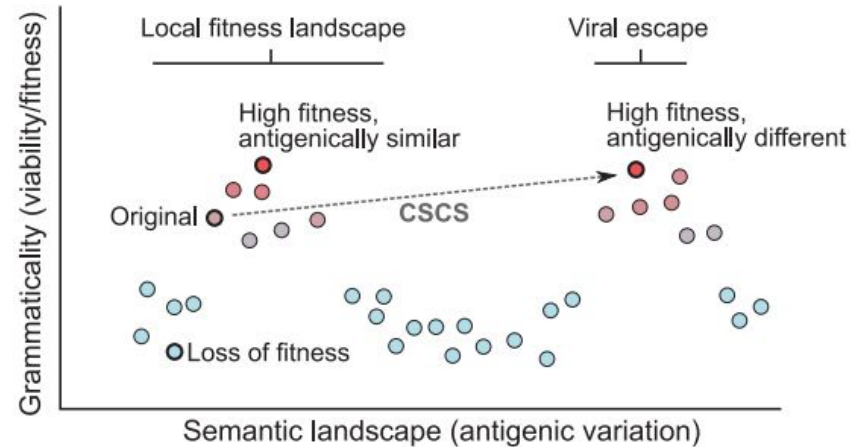
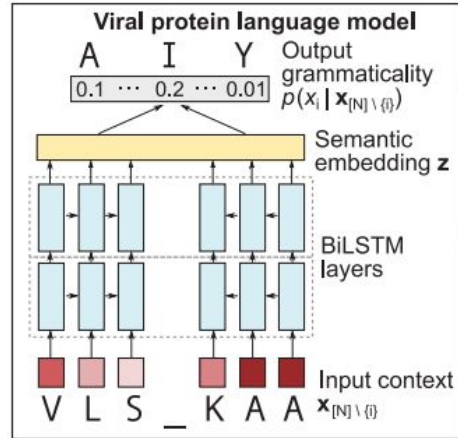


b



Controllable generation of protein sequences
(Madani et al., bioRxiv, 2021)

Application: unsupervised prediction of viral escape



LM predicts viral escape
(Hie et al., Science, 2021)

Summary of today

- Biological sequences
 - DNA, RNA, protein
- CNNs
 - Extract spatial features using convolution filters
- RNNs
 - Capture temporal dependency
 - LSTM, GRU
- Transformers
 - Self-attention layer
 - Widely used in recent state-of-the-art NLP models
- Example: language model of protein sequences
 - LM: predicting a word given context
 - Protein LM used for supervised predictions, sequence generation, unsupervised prediction (e.g., viral escape)