

CPH 100A: Modeling Images and Volumes: Convolutional Neural Networks

Instructor: Adam Yala, PhD (yala@berkeley.edu)

Agenda

Recap

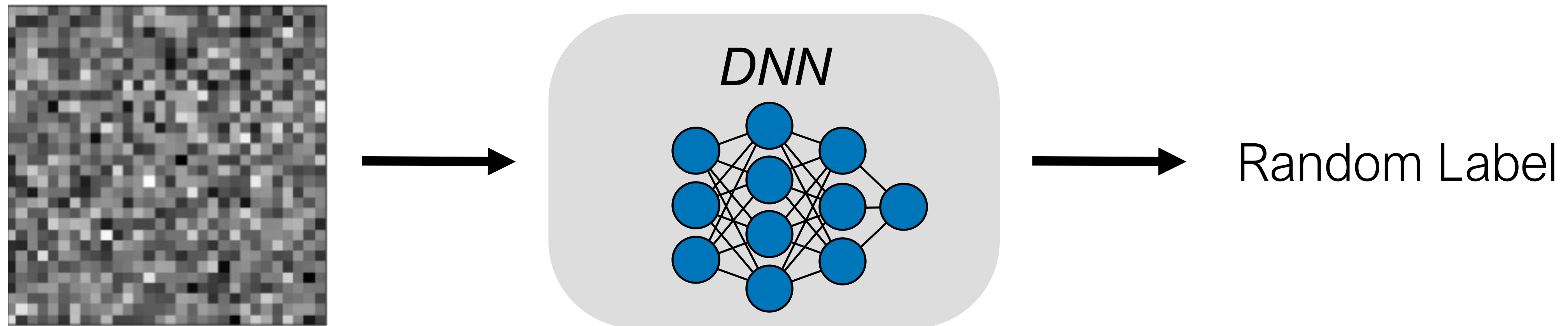
Failure modes of fully-connected neural networks

Convolutions

Pooling

CNNs across modalities

Well-trained Networks can learn anything

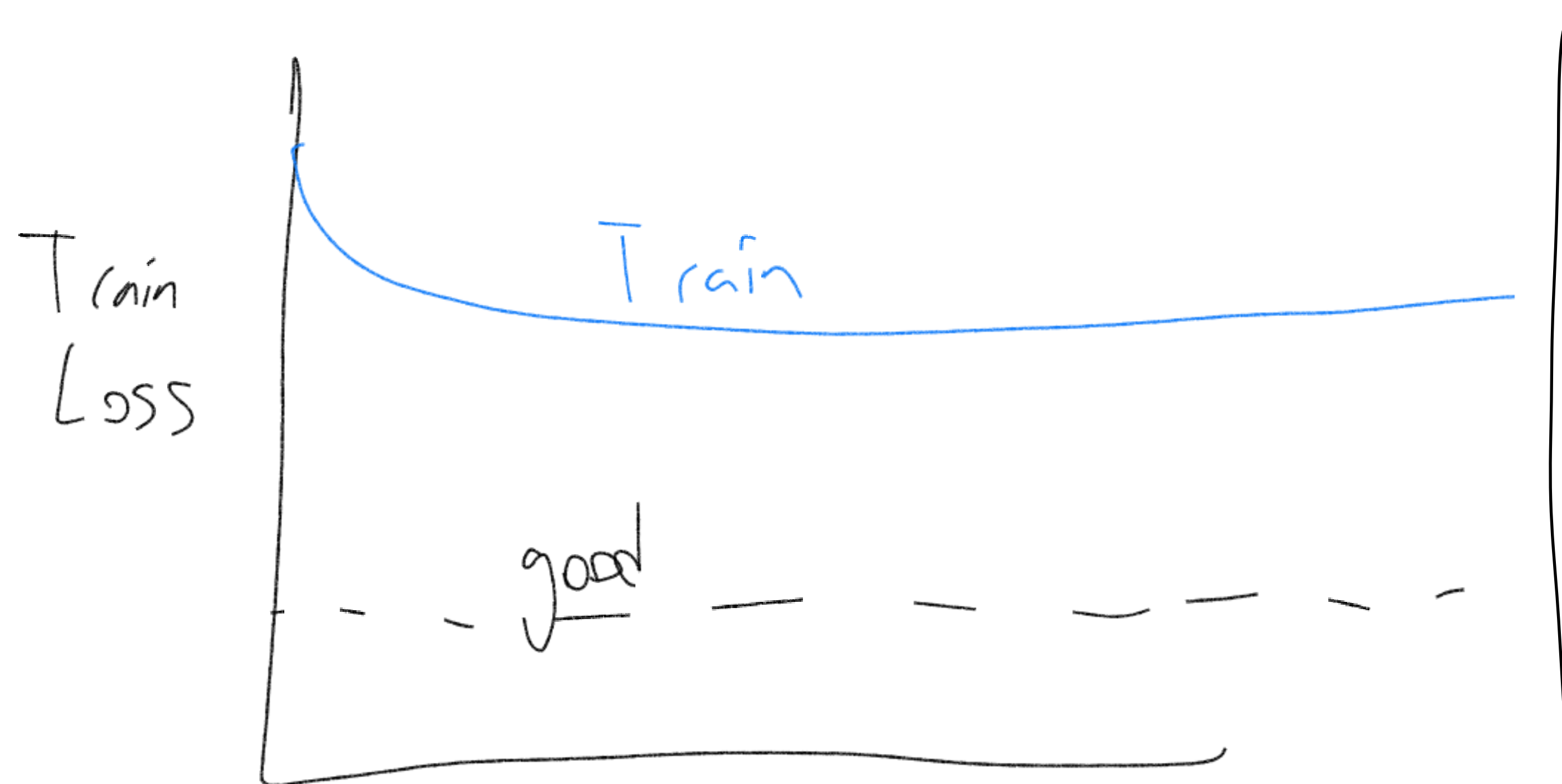


Understanding deep learning (still) requires rethinking generalization

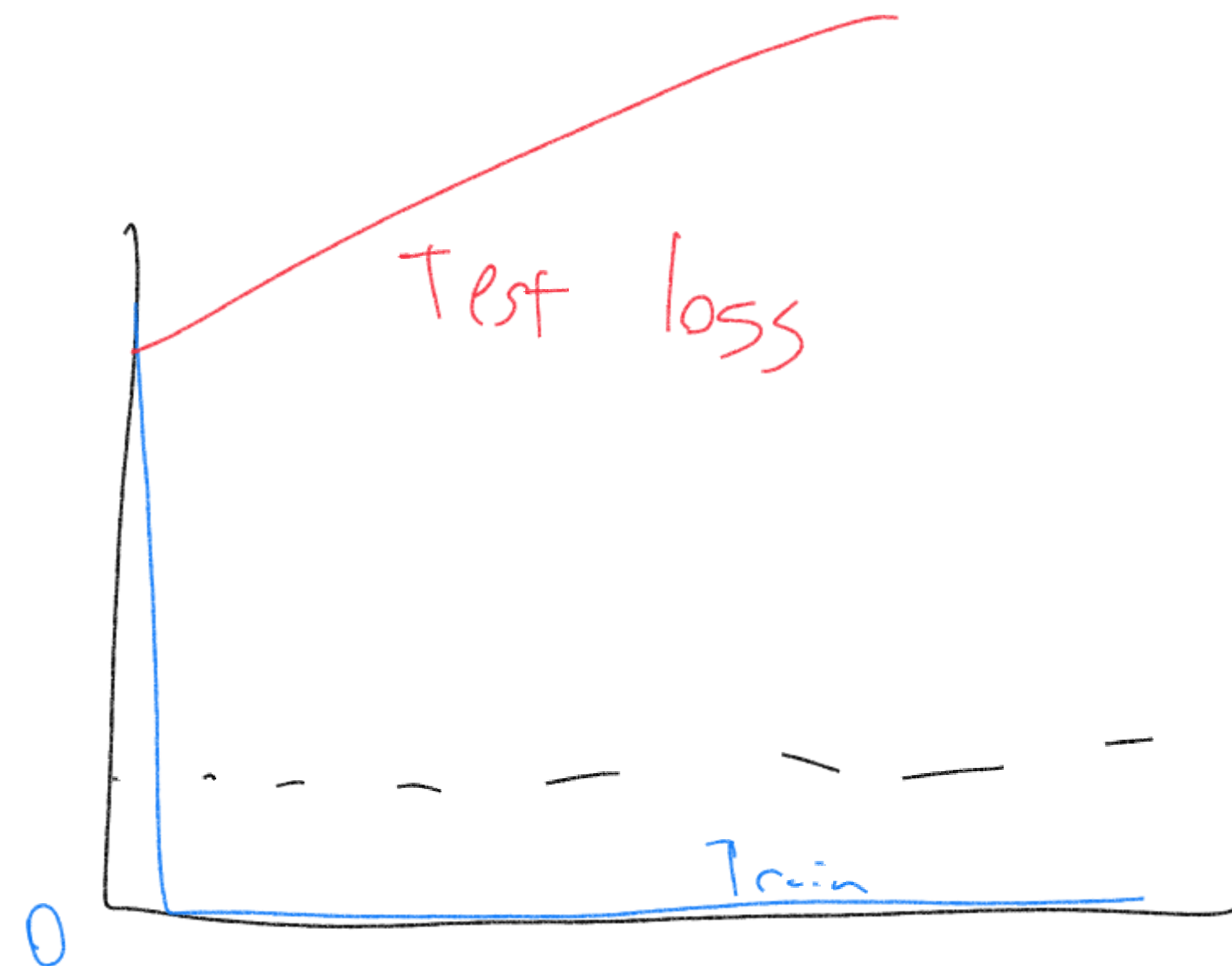
Authors:  [Chiyuan Zhang](#),  [Samy Bengio](#),  [Moritz Hardt](#),  [Benjamin Recht](#),  [Oriol Vinyals](#) [Authors Info & Claims](#)

Communications of the ACM, Volume 64, Issue 3 • March 2021 • pp 107–115 • <https://doi.org/10.1145/3446776>

Failure modes of Neural Network Optimization



Underfitting



Overfitting

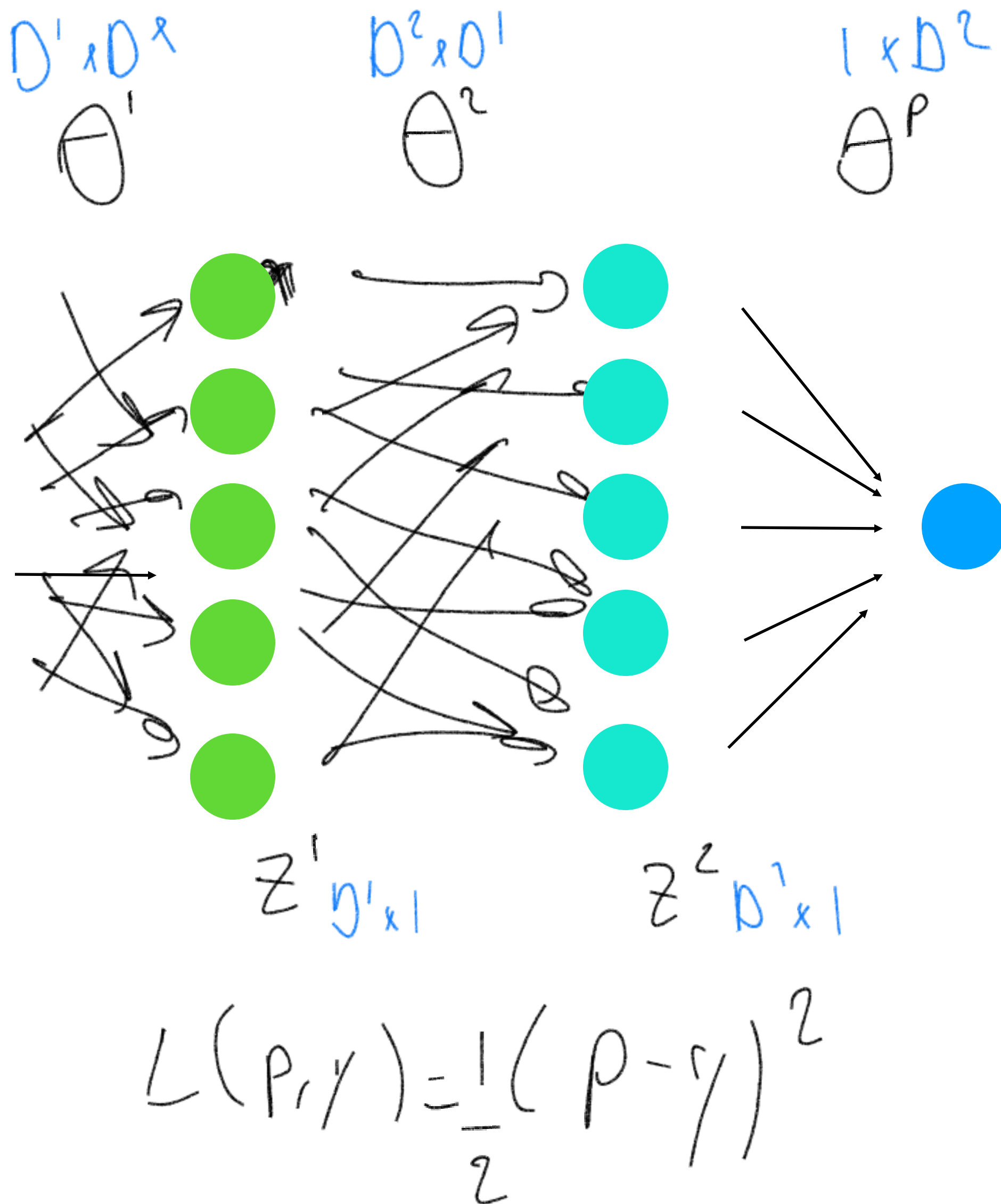
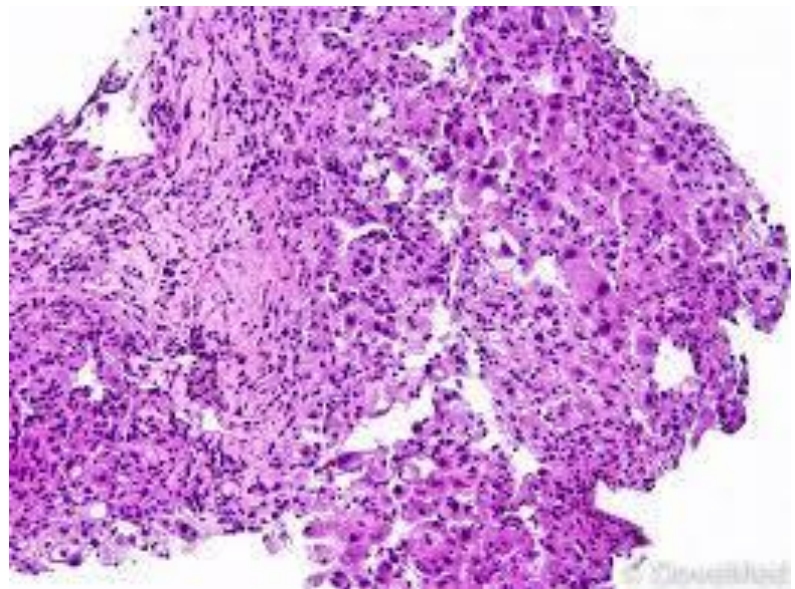
Why did the neural networks fail to train?

Complex Interaction between:

- Initialization
- Hypothesis Class
- Optimizer
- Learning Objective
- Data



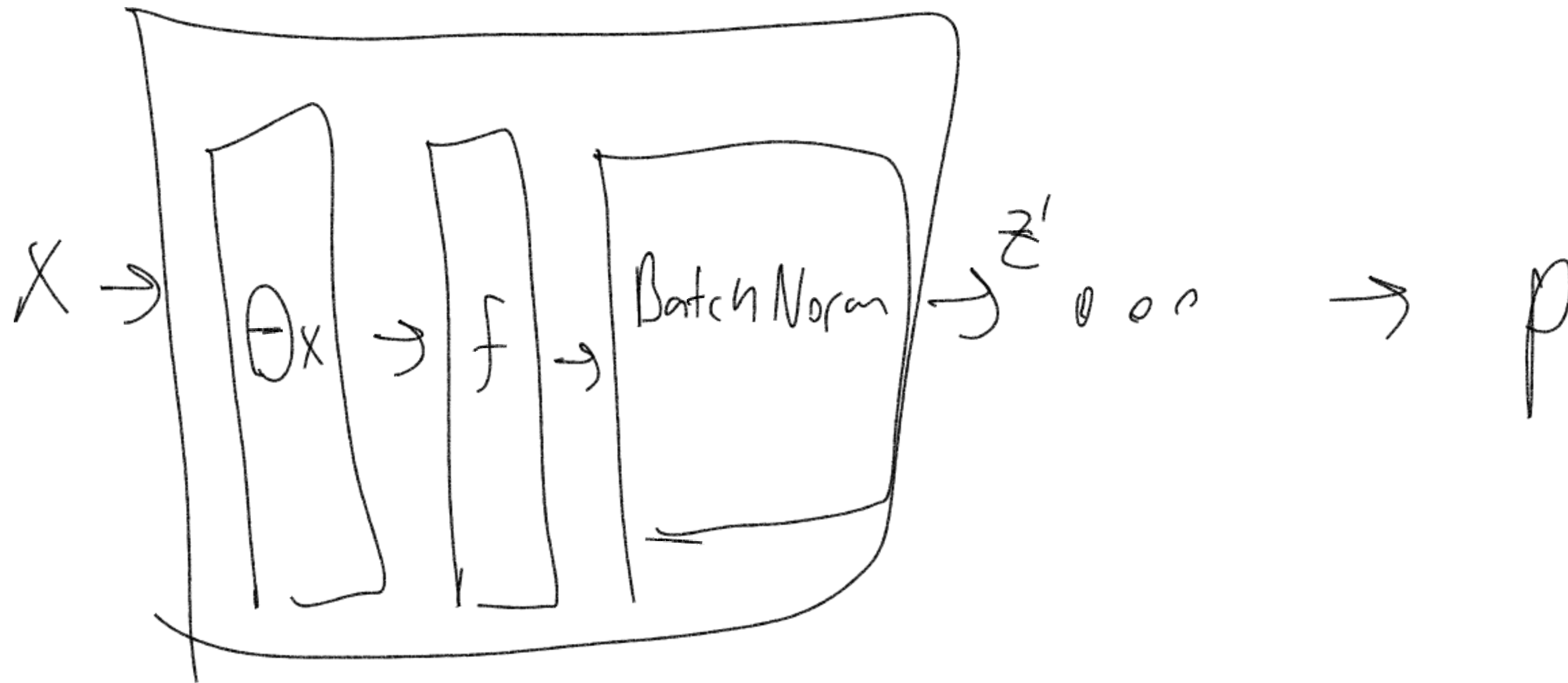
Choosing scales of random init



Control $\|z^e\|$
 Control $\left\| \frac{\partial L}{\partial \theta} \right\|$
 — depend on f
 He init: $\in \text{ReLU}$
 $\theta \sim N(0, \sqrt{2/D})$

BatchNorm: Preventing activation collapse

$$\text{Batchnorm}(X) = \frac{X - \mu_X}{\sqrt{\sigma_X^2 + \epsilon}}$$



Use $\forall z^l$!

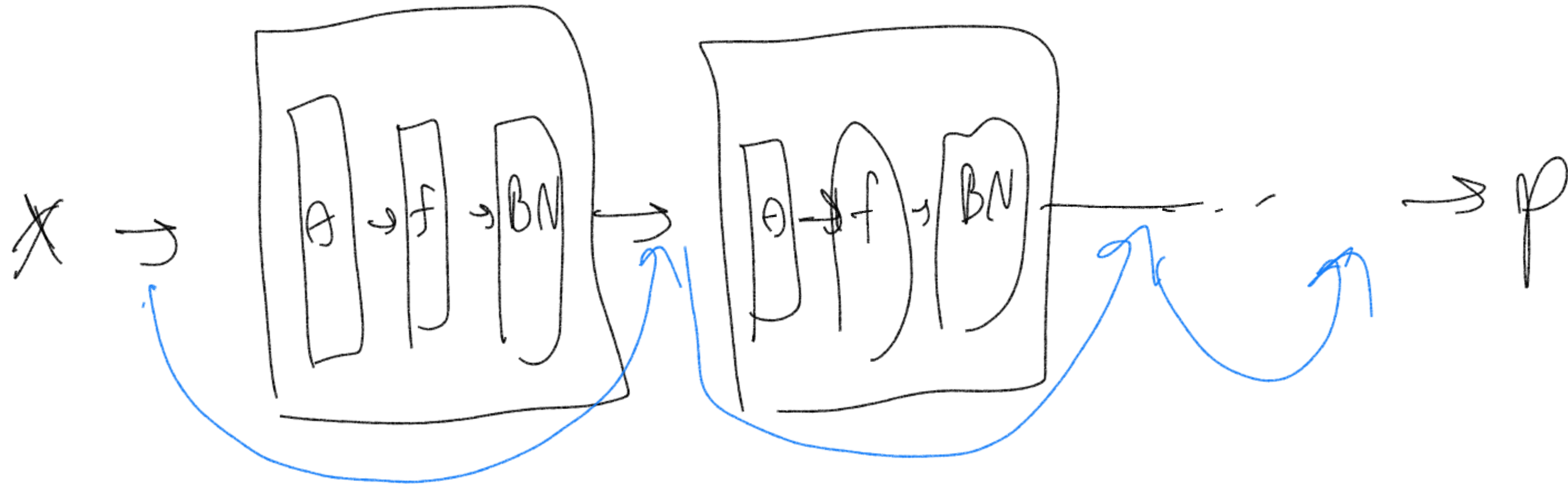
Now z^l can't
explode!

Residual connections: Skipping training bottlenecks

Add Skip / Residual

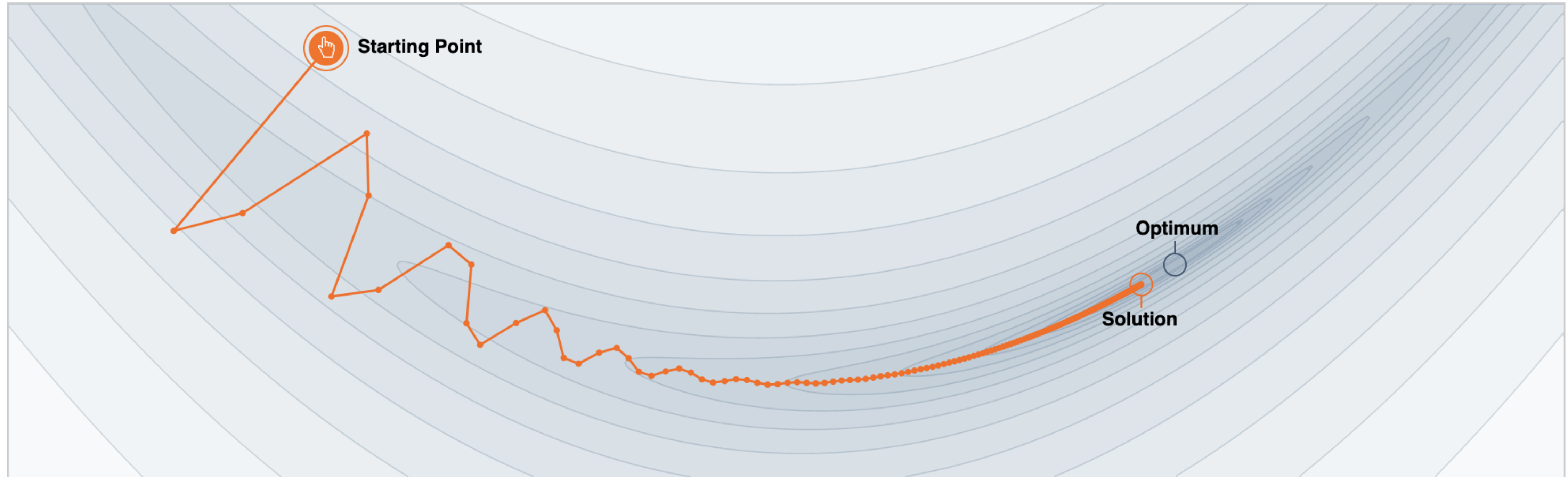
$$z^l = f(\underbrace{\Theta z^{l-1} + b}_{a^l}) + z^{l-1}$$

$$\frac{\partial z^l}{\partial z^{l-1}} = f'(a^l) \Theta' \boxed{+ I}$$

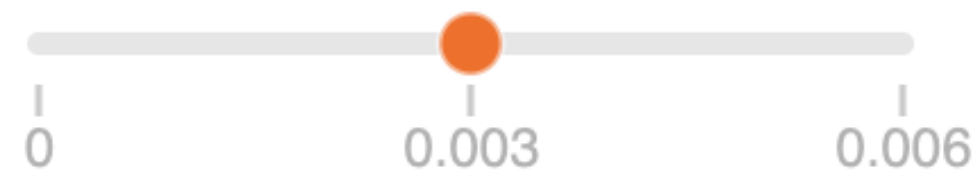


Allows training
deep NNs!

Momentum: Accelerating SGD



Step-size $\alpha = 0.02$



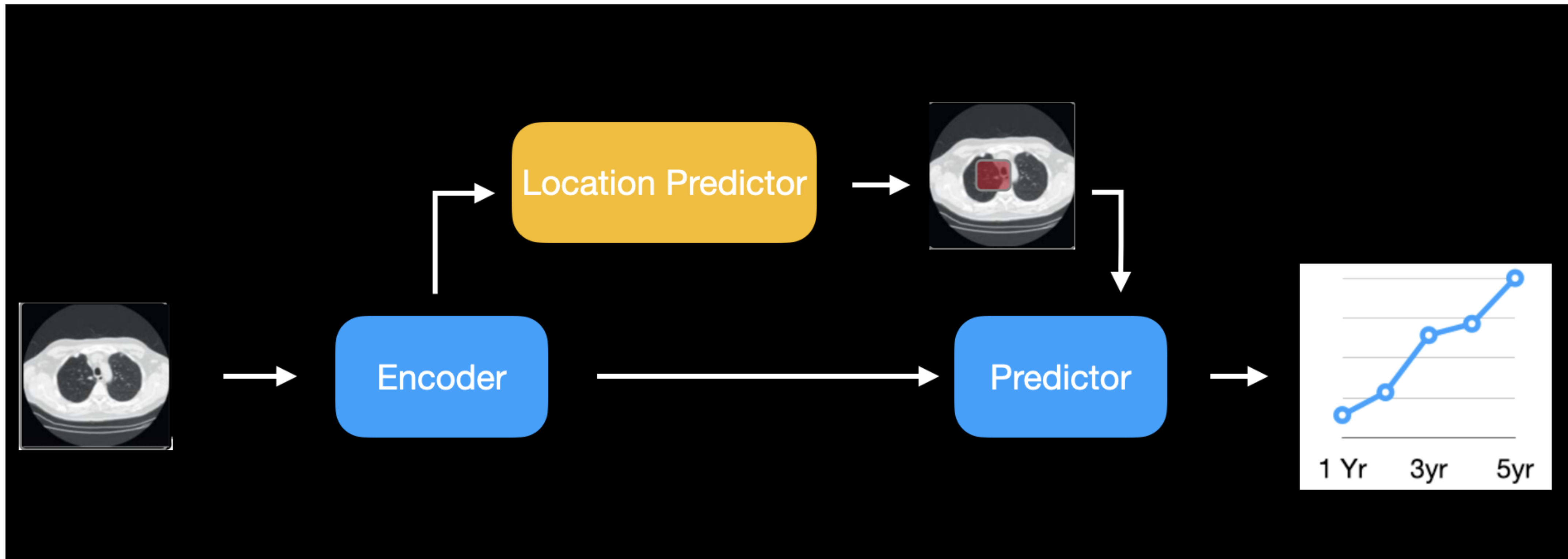
Momentum $\beta = 0.99$



We often think of Momentum as a means of dampening oscillations and speeding up the iterations, leading to faster convergence. But it has other interesting behavior. It allows a larger range of step-sizes to be used, and creates its own oscillations. What is going on?

<https://distill.pub/2017/momentum/>

Managing overfitting: Adding additional losses



$$\mathcal{L} = \mathcal{L}_{\text{Risk}} + \mathcal{L}_{\text{Location}} + \mathcal{L}_{\text{Weight Decay}} \leftarrow \|\Theta\|^2$$

Managing overfitting: Data Augmentation



→ Aug



Agenda

Recap

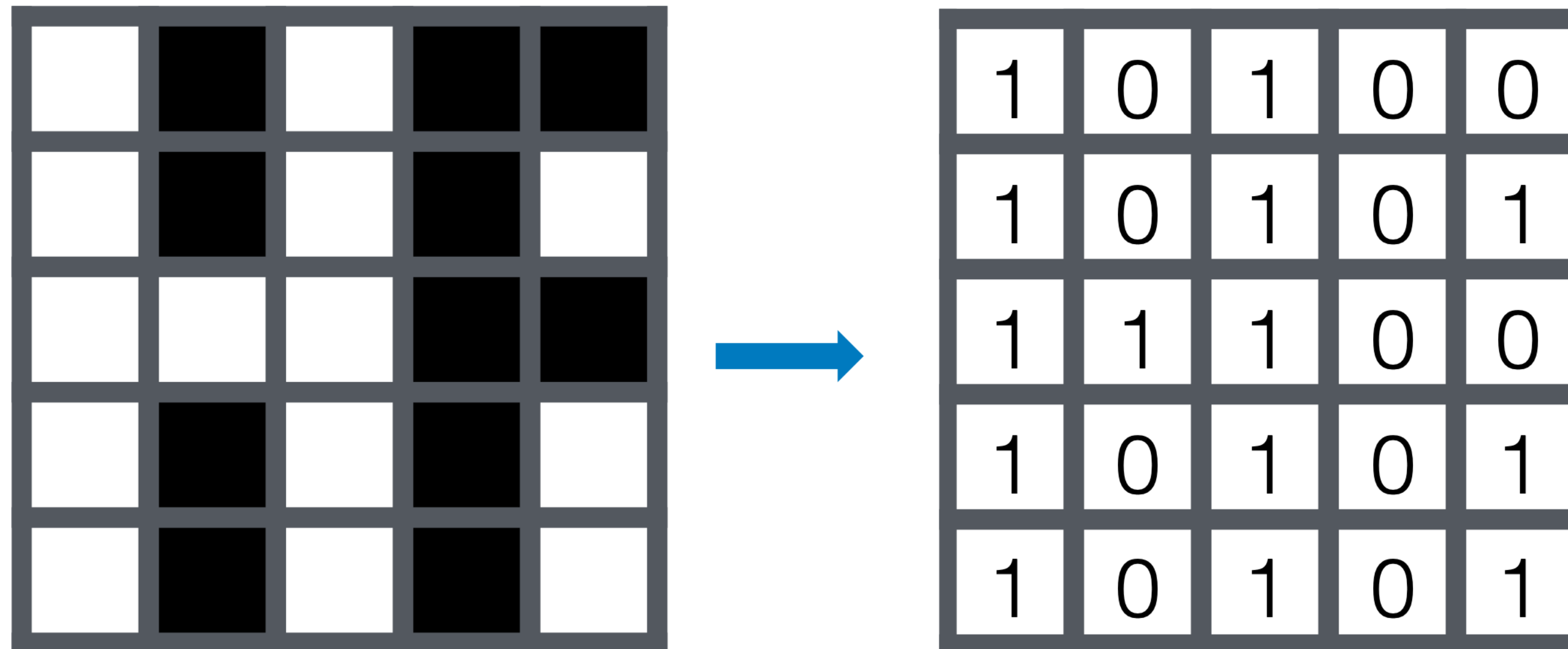
Failure modes of fully-connected neural networks

Convolutions

Pooling

CNNs across modalities

Images are tensors

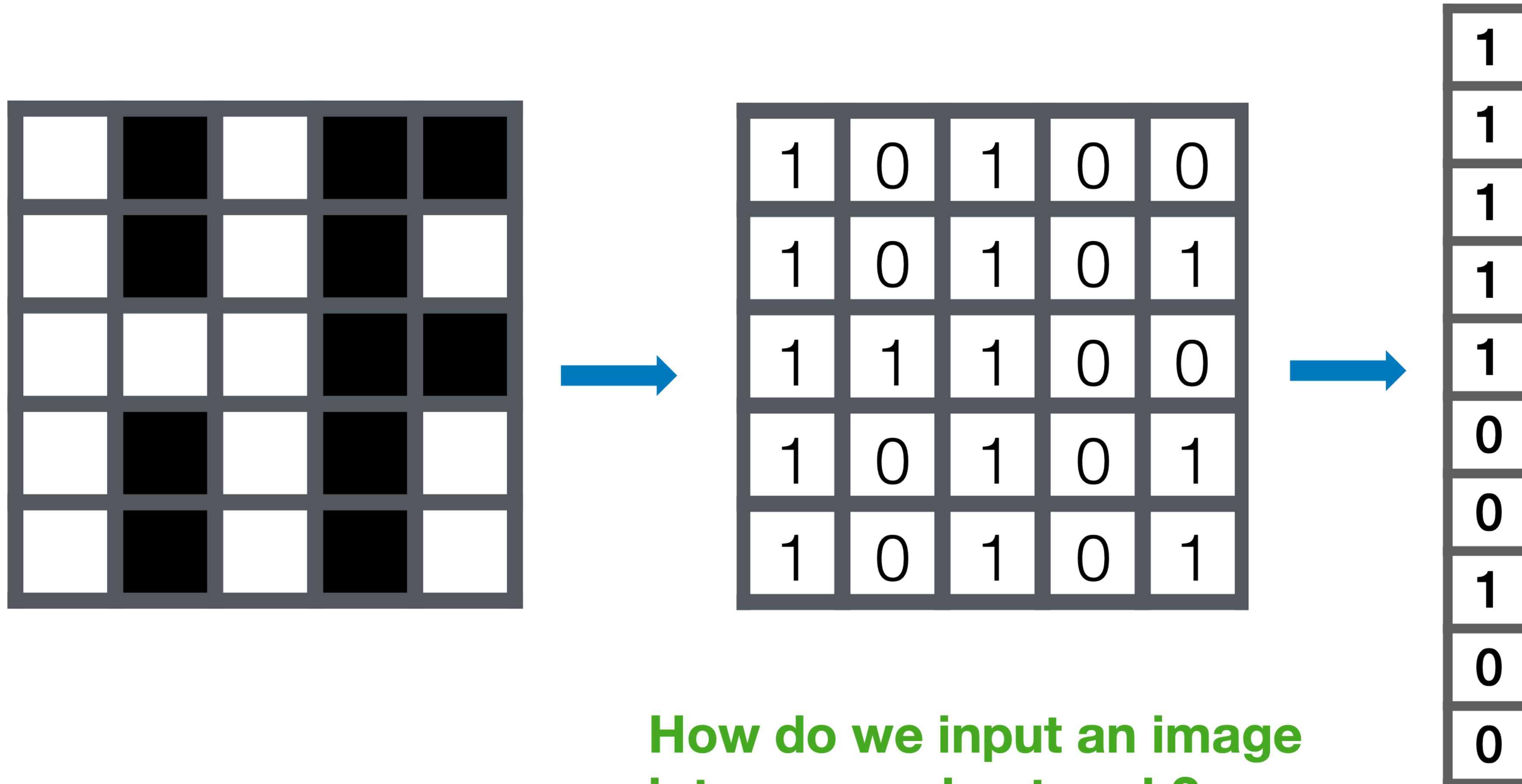


**How do we input an image
into a neural network?**

...

Image credit: Tamara Broderick

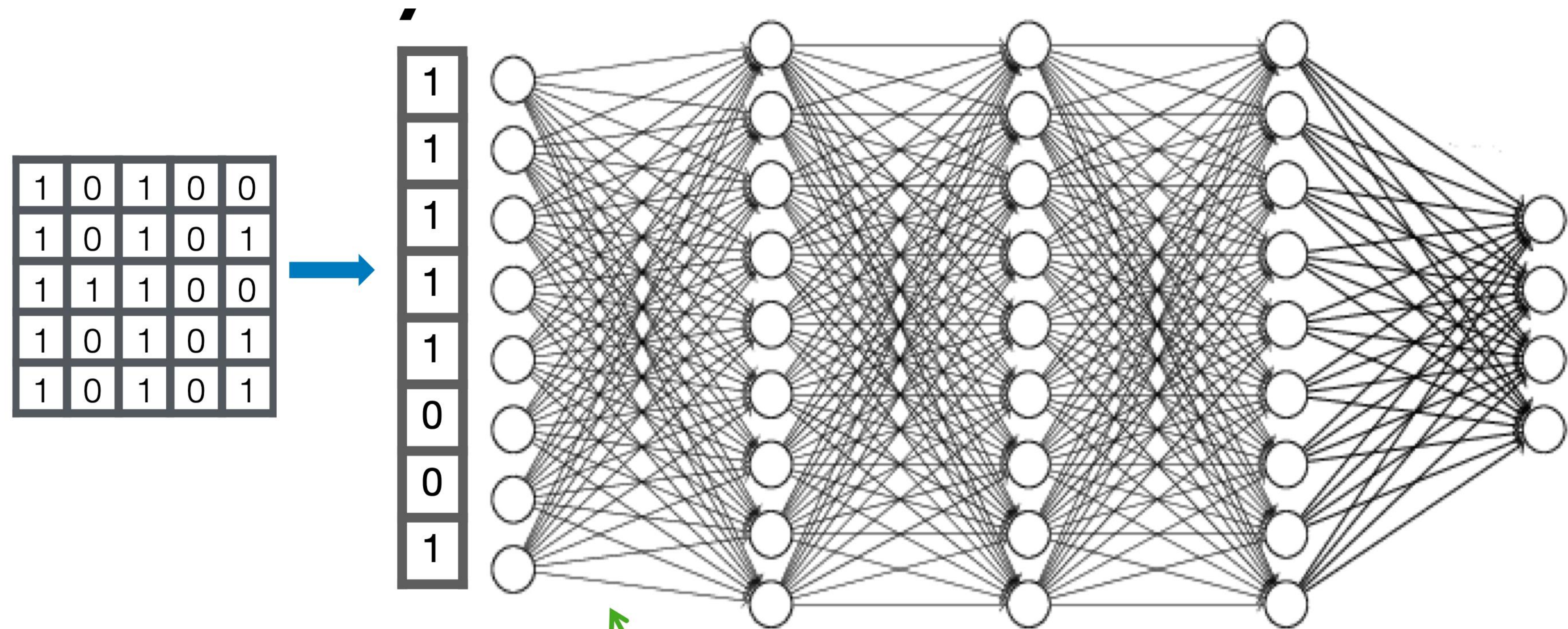
Images in our class so far..



**How do we input an image
into a neural network?**

Idea: make it a vector and use FNN!...

Using FNNs on tiny images



Example:

512*512 image

~250K hidden units

62.5B parameters per layer <- Bigger than GPT3...

What's wrong with FNNs?



$$x^0 = [0, 1, 1, \dots, 0]$$



$$x^1 = [0, 0, 1, \dots, 0]$$

Small padding, very different feature vectors!

Can it learn?

Yes

Will it be easy?

No.

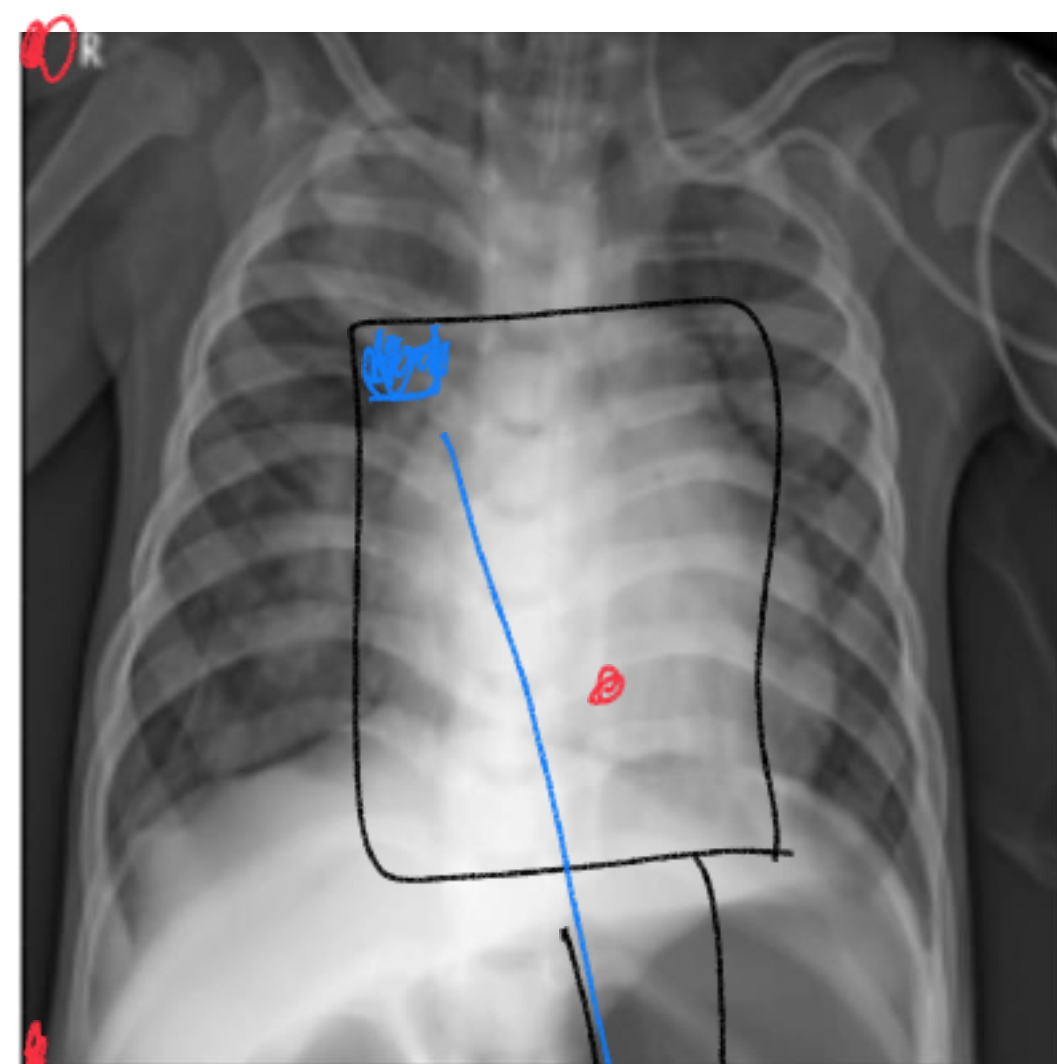
→ Learn invariance is hard.

→ Memory prohibitive

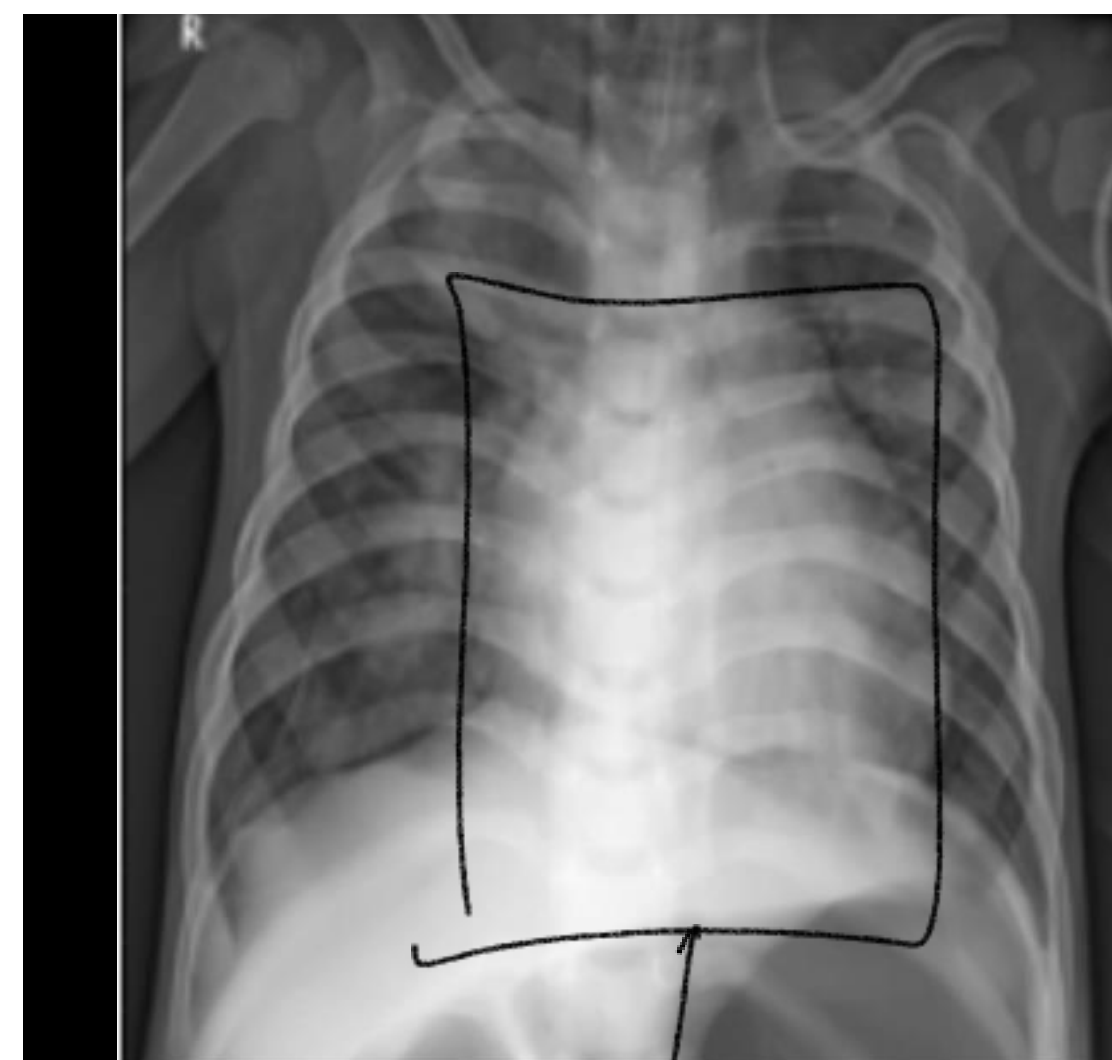
Conceptual role of Hypothesis class:
Choose your Mountain range



What do we know?



Locality



Translations
don't matter!

Desired properties for a good Hypothesis Class

Capture spatial dependencies: Pixel positions and locality matter!

Handle Translations / Nuisance variations: Objects of interest can appear anywhere

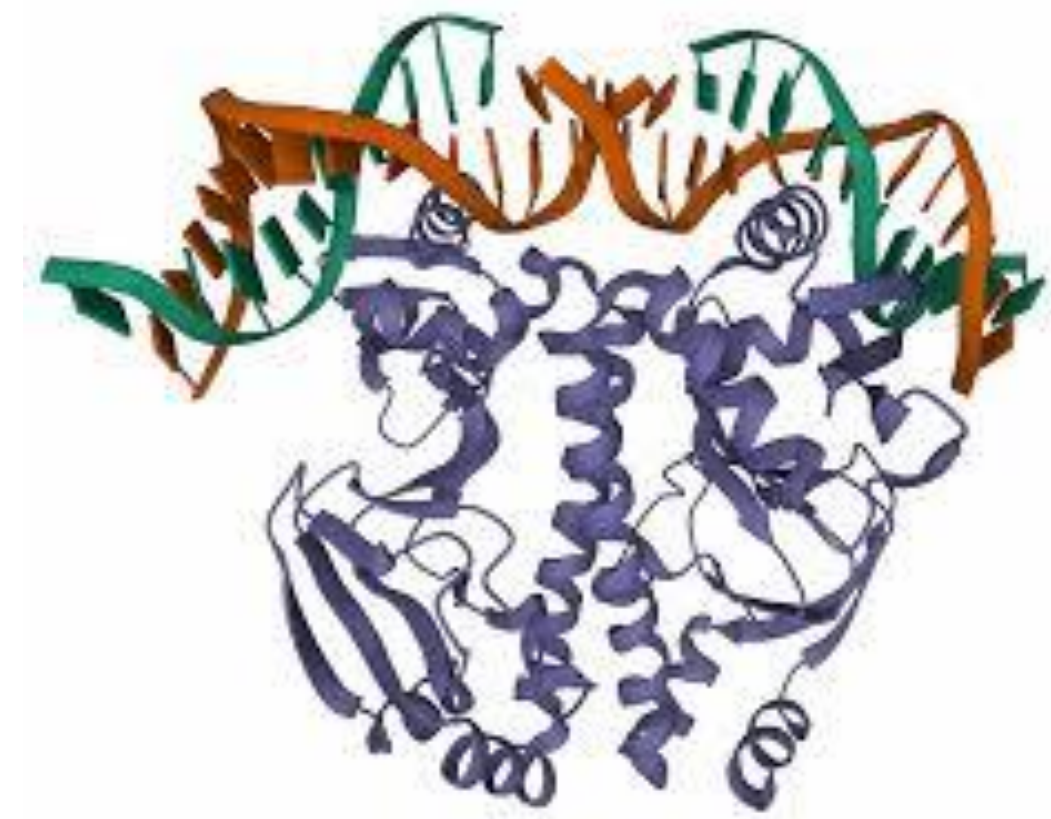
Scale: Allow efficient computation for large inputs



Images/ Volumes



Text



Graphs

Agenda

Recap

Failure modes of fully-connected neural networks

Convolutions: Capturing locality and positional-equivariance

Pooling

CNNs across modalities

Convolution: Definition

$$\text{Conv}(x, \theta) = (x * \theta)[i] = \sum_{j=0}^k x[i-j] \theta[j]$$

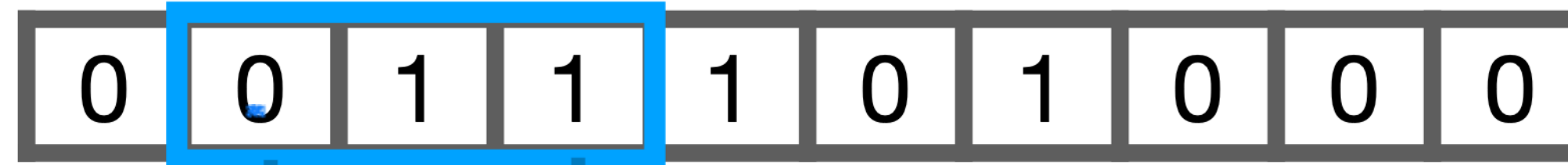
$$x \in \mathbb{R}^N$$

$$\theta \in \mathbb{R}^k$$

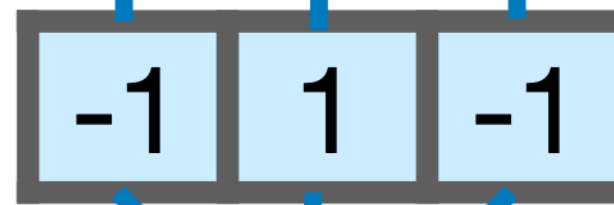
Sliding window dot product!

Convolutional Layer: 1D example

- 1D image



- Filter

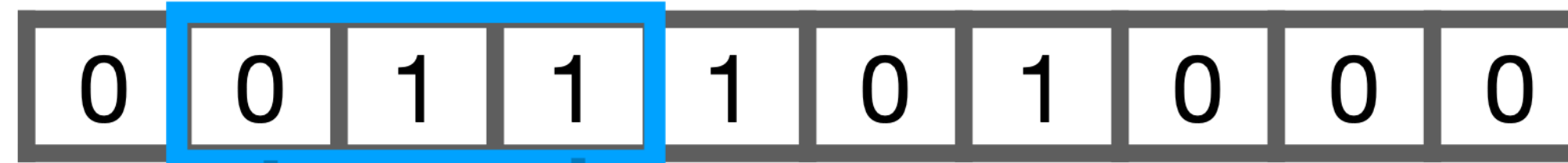


- After convolution

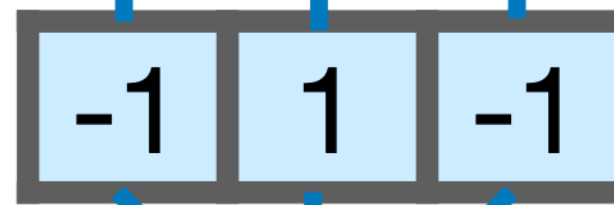


Convolutional Layer: 1D example

- 1D image



- Filter



Stride: by how much
do we shift the filter

- After
convolution



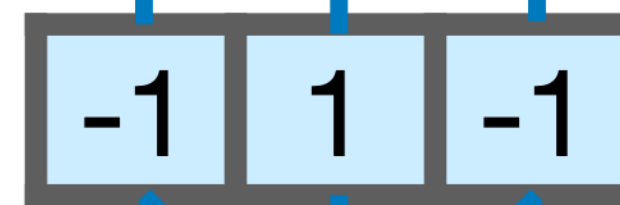
Weight sharing:
same filter applied
to both (all) patches

Convolutional Layer: 1D example

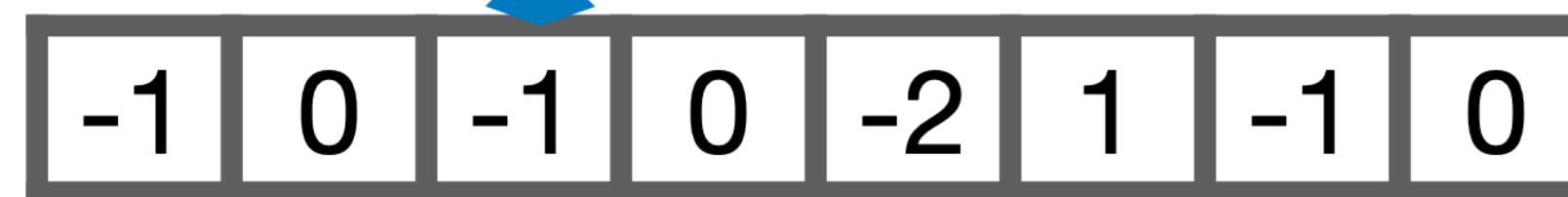
- 1D image



- Filter



- After convolution



Big advantage: due to weight sharing, needs much fewer weights than a fully connected network

Translation equivariance!

Convolutional Layer: 1D example

- 1D image

0	0	1	1	1	0	1	0	0	0
---	---	---	---	---	---	---	---	---	---

- Filter

-1	1	-1
----	---	----

- After
convolution

-1	0	-1	0	-2	1	-1	0
----	---	----	---	----	---	----	---

- After ReLu

--	--	--	--	--	--	--	--

Padding

- 1D image

0	0	1	1	1	0	1	0	0	0
---	---	---	---	---	---	---	---	---	---

- Filter

-1	1	-1
----	---	----

- After
convolution

-1	0	-1	0	-2	1	-1	0
----	---	----	---	----	---	----	---

- After ReLu

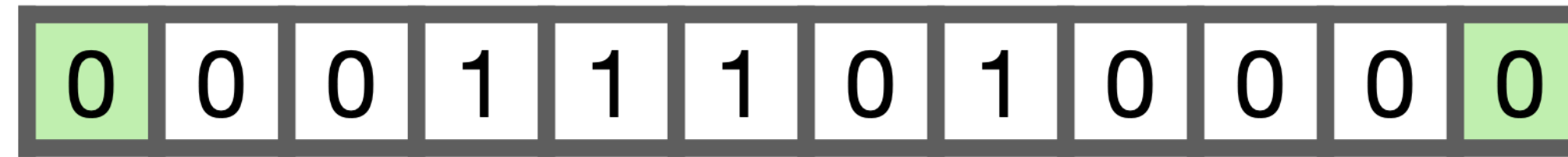
0	0	0	0	0	1	0	0
---	---	---	---	---	---	---	---

Output is smaller! (why?)

Remedy: pad input with zeros

Padding

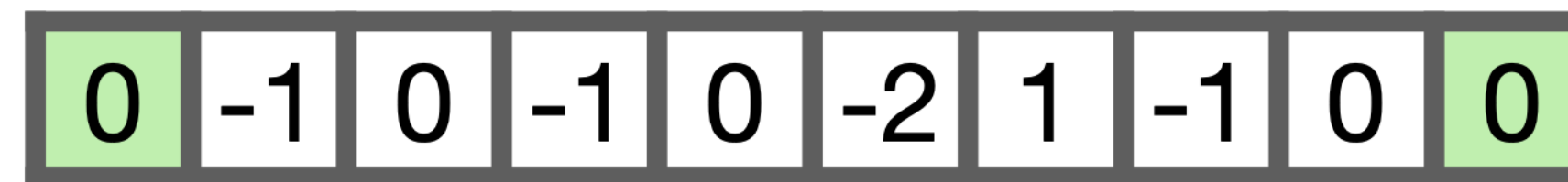
- 1D image



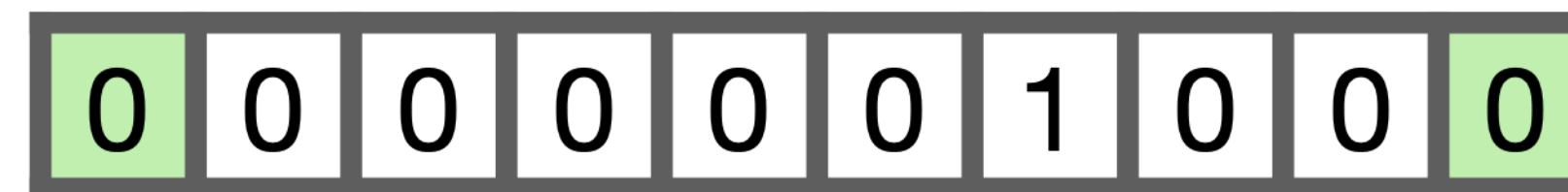
- Filter



- After convolution



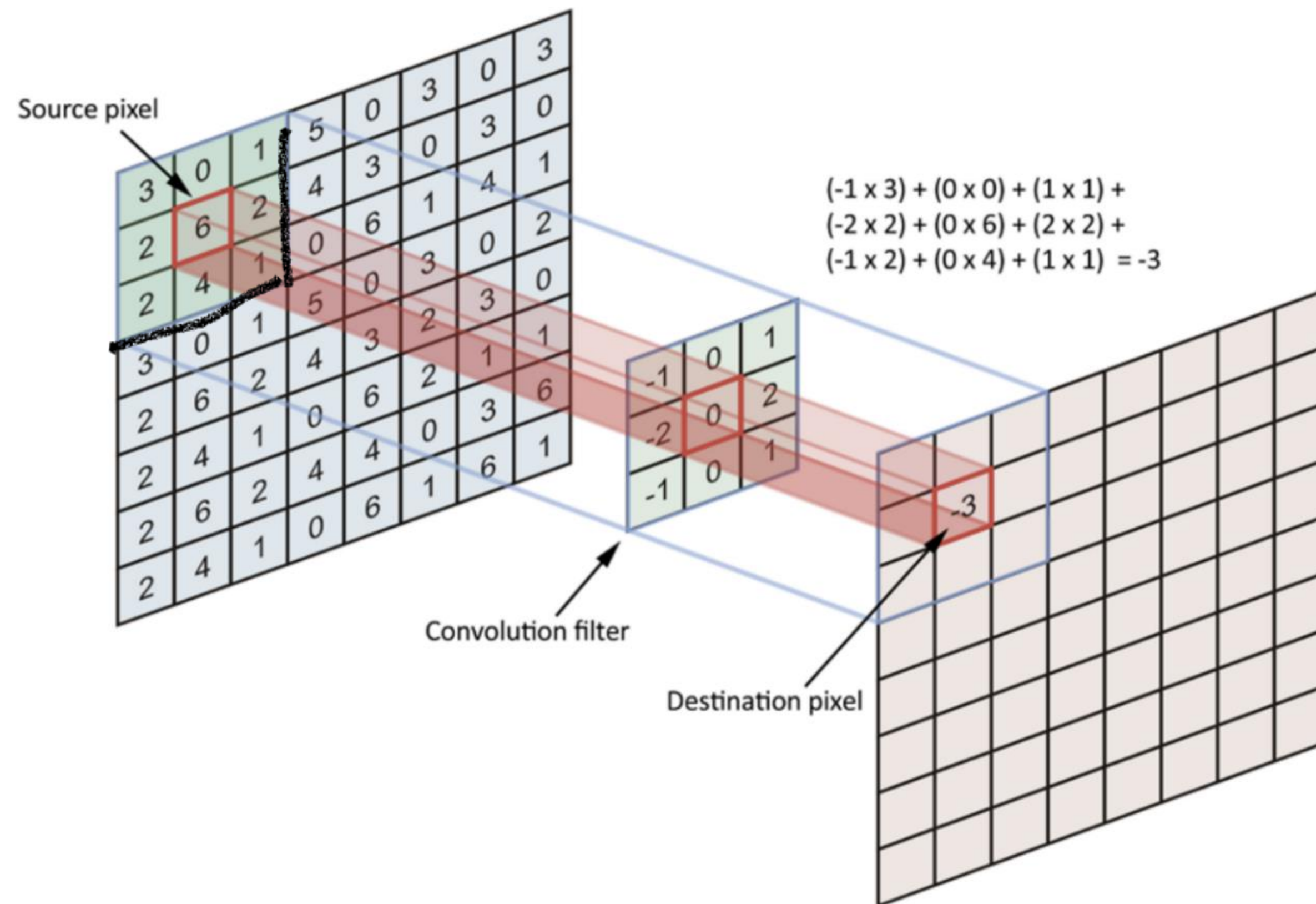
- After ReLu



Output is smaller! (why?)

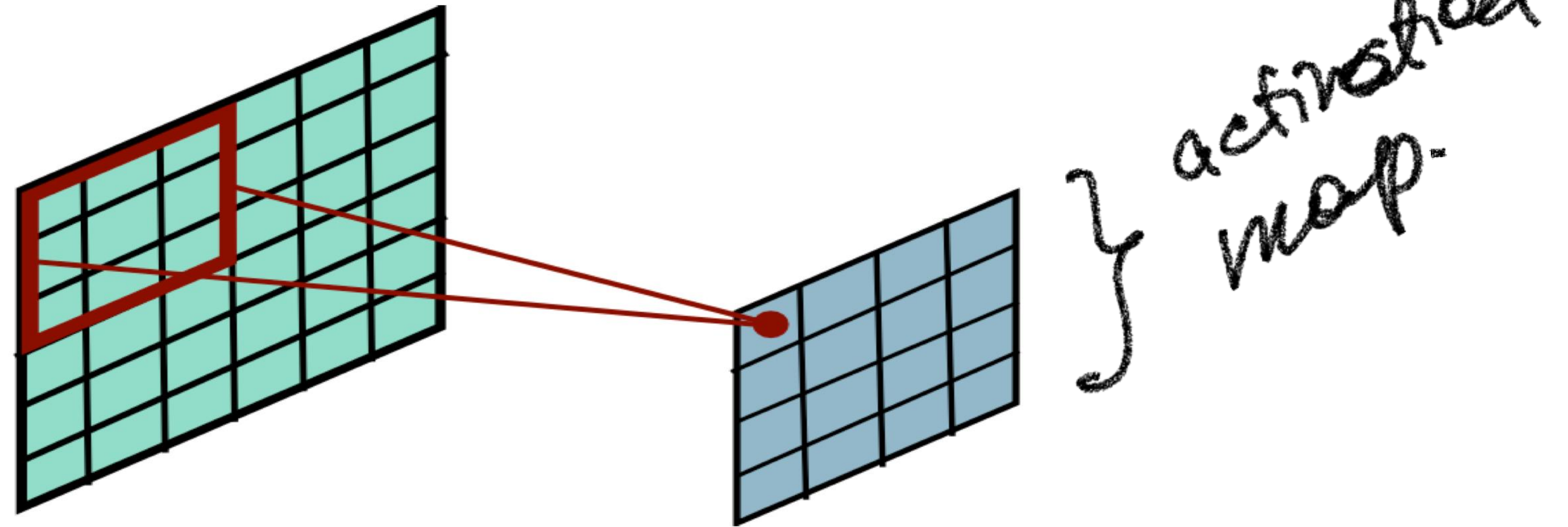
Remedy: pad input with zeros

2D Convolutions

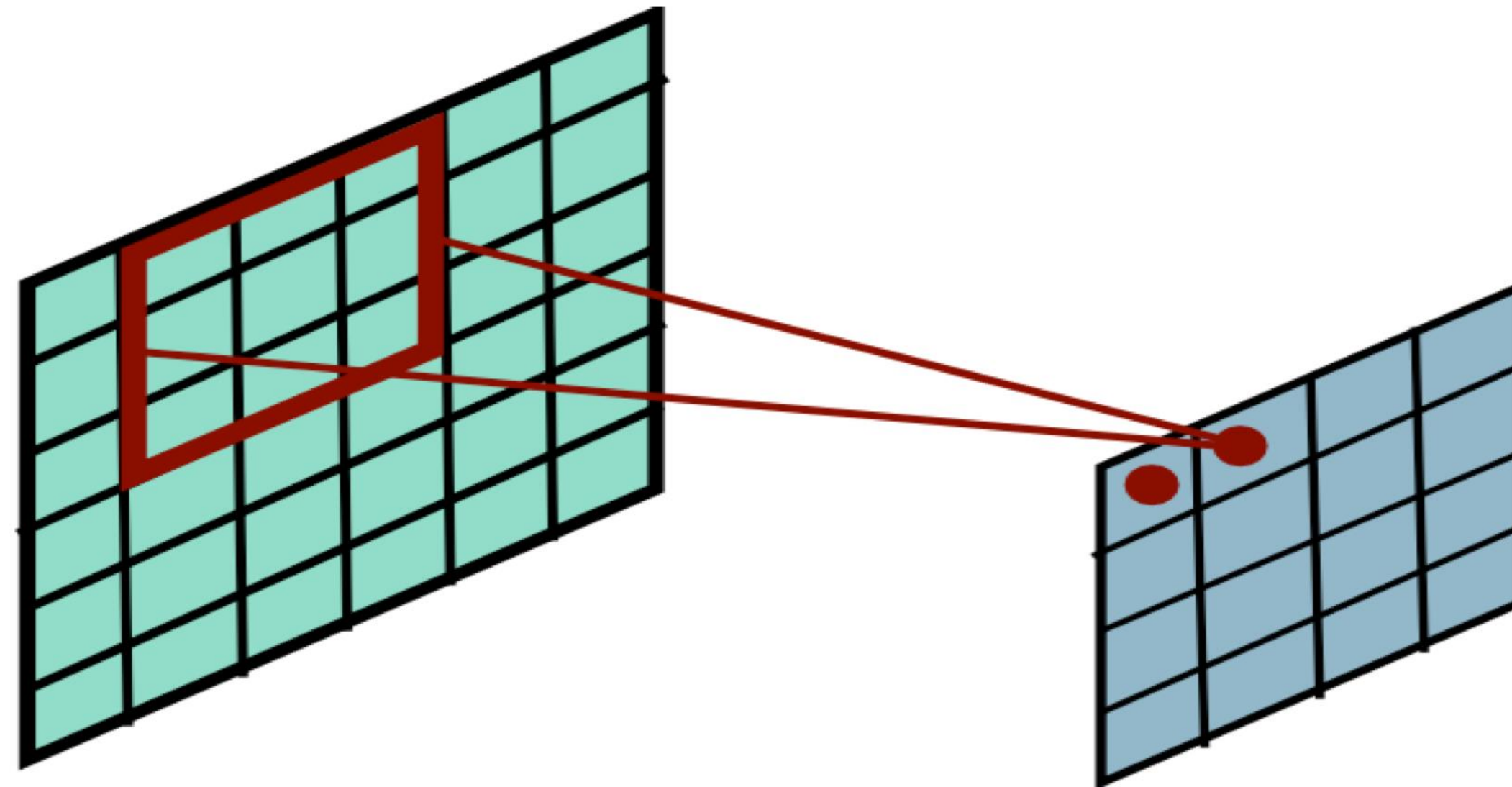


The convolution operation.

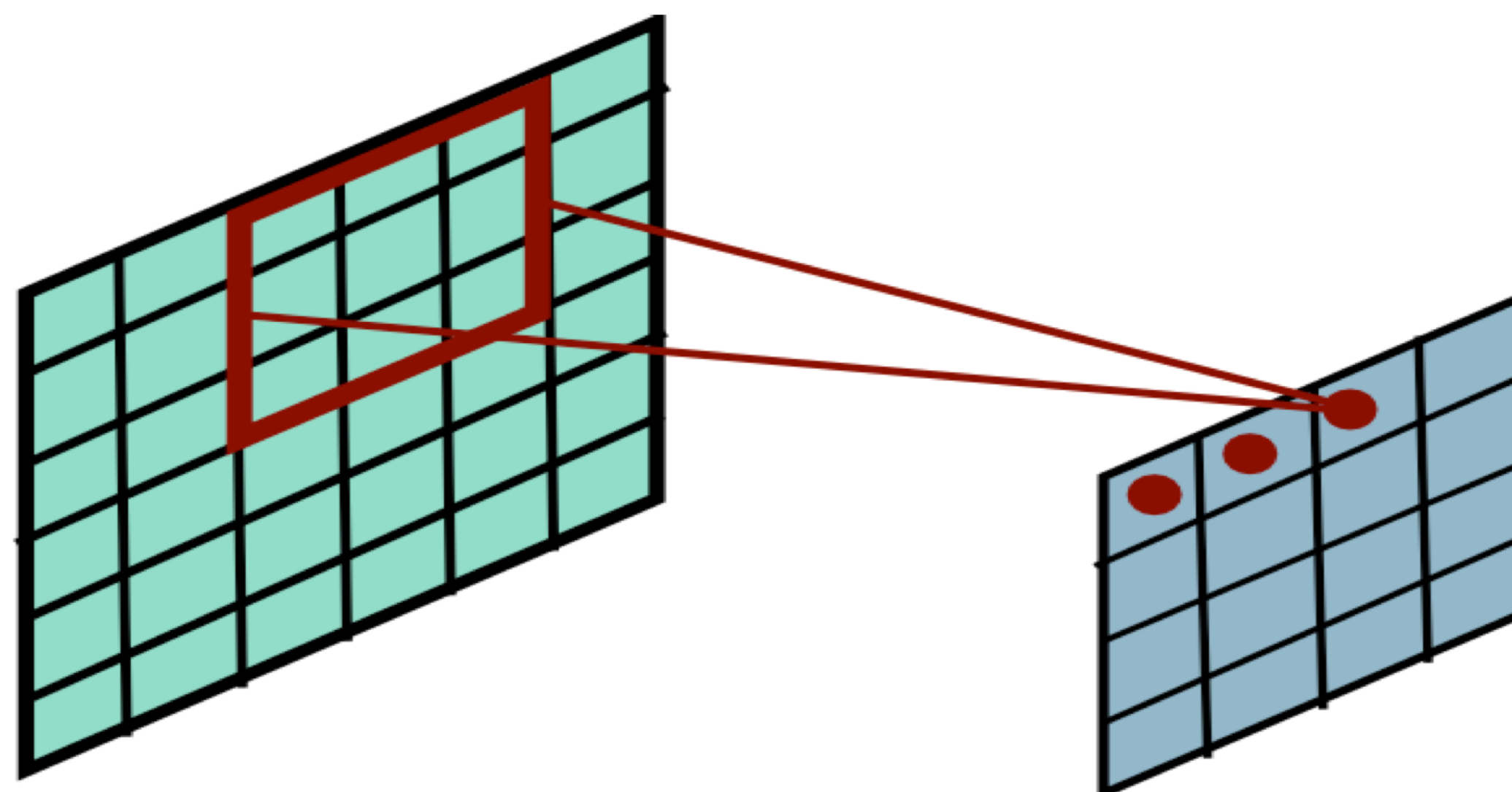
Convolutional Layer



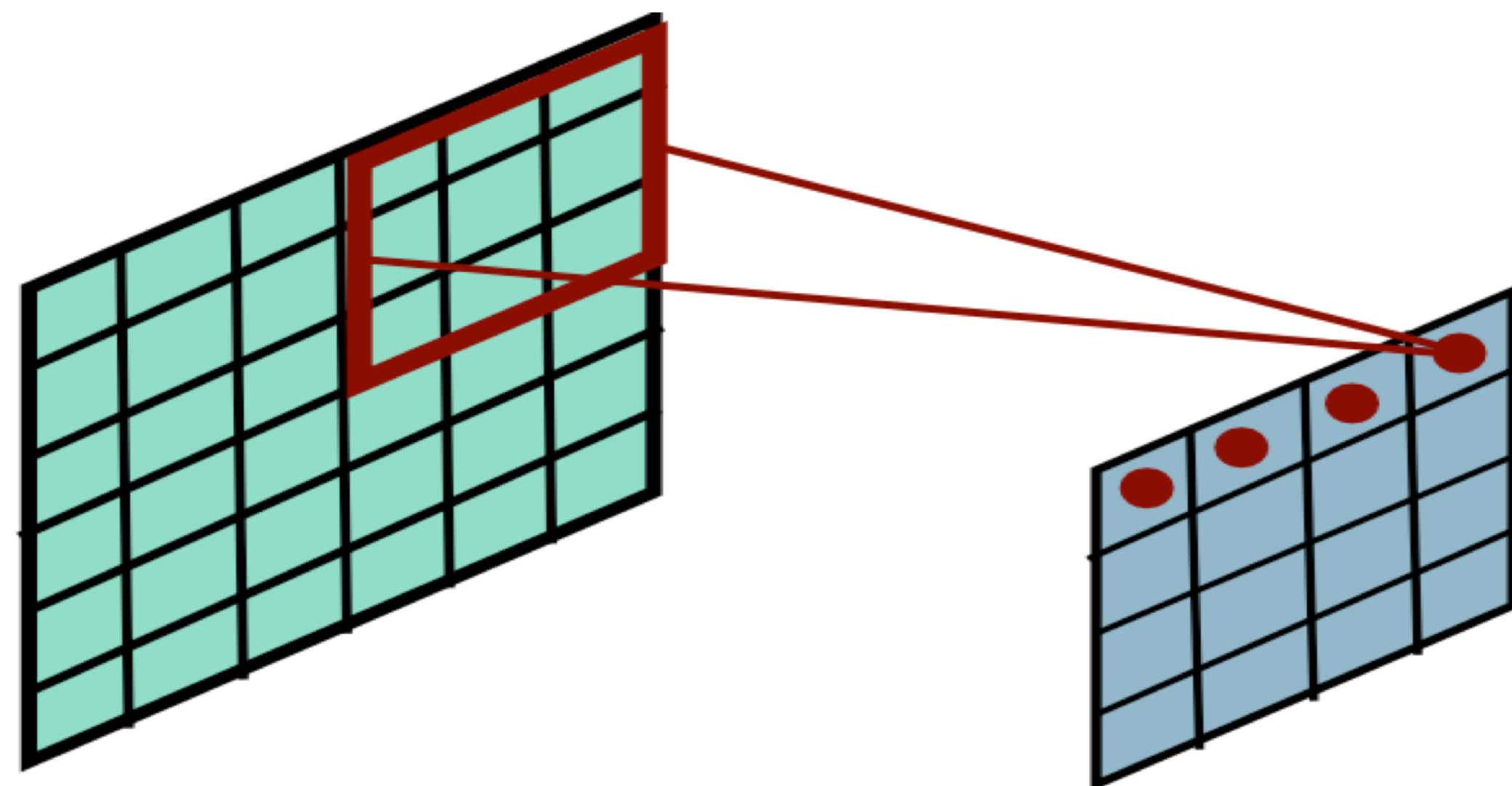
Convolutional Layer



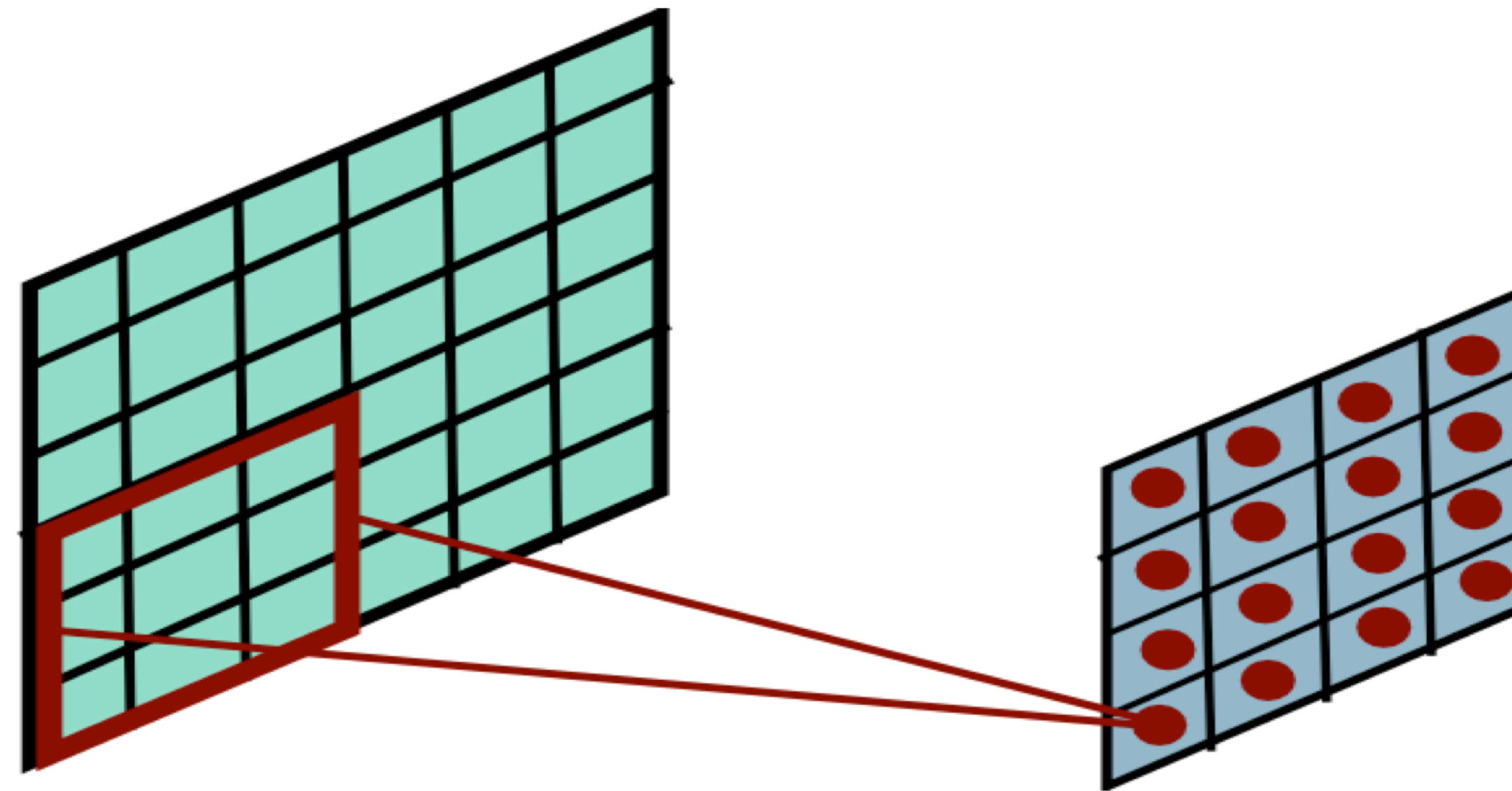
Convolutional Layer



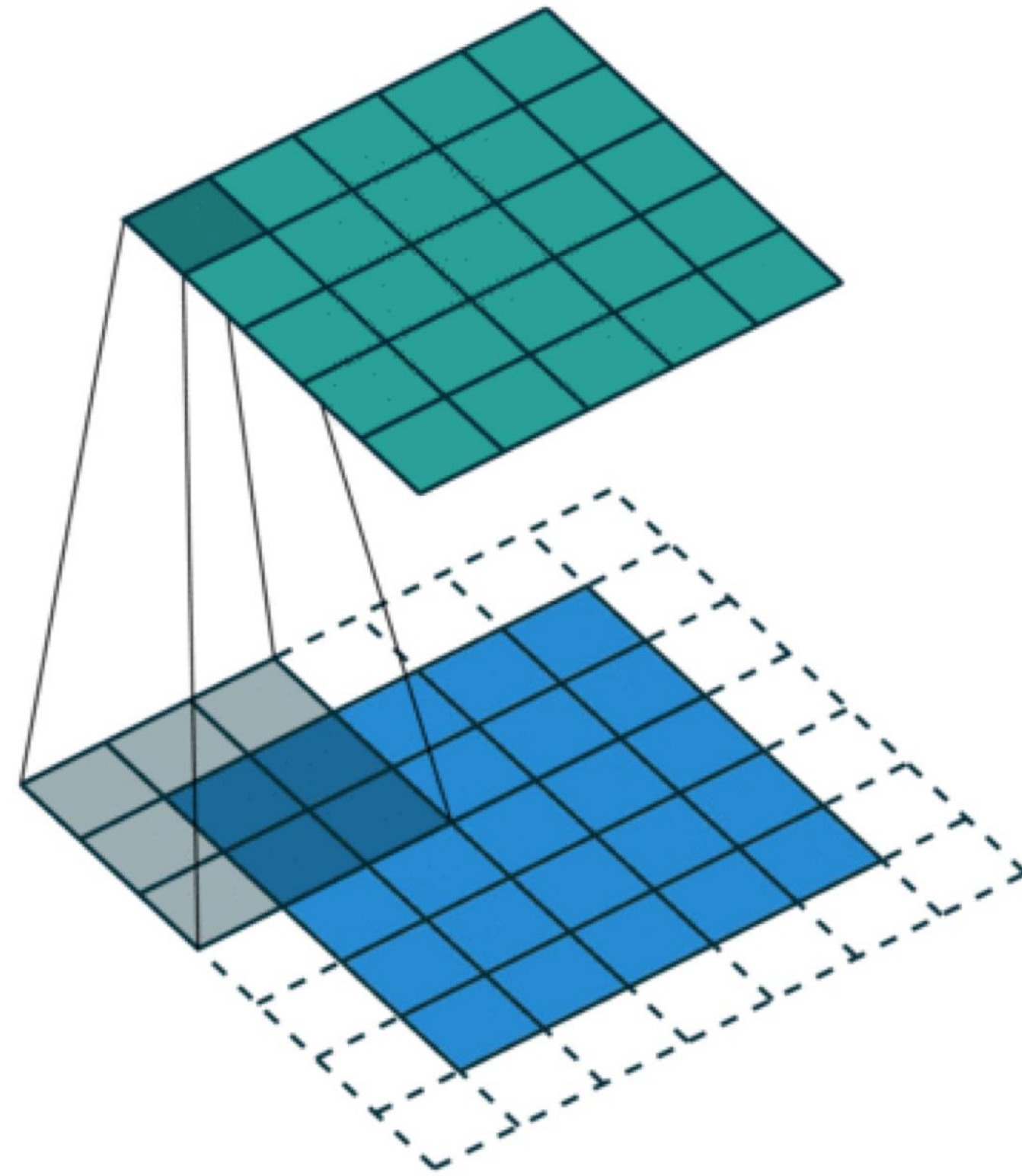
Convolutional Layer



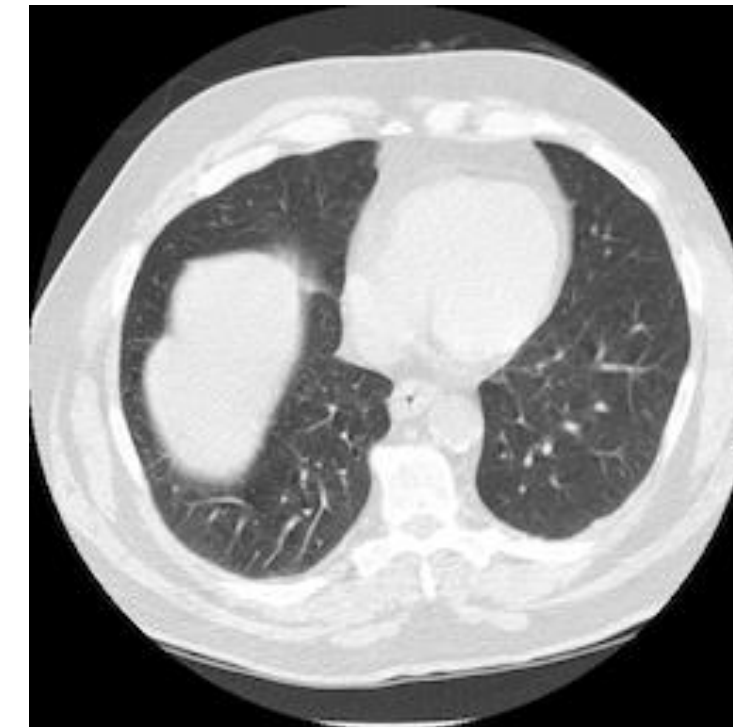
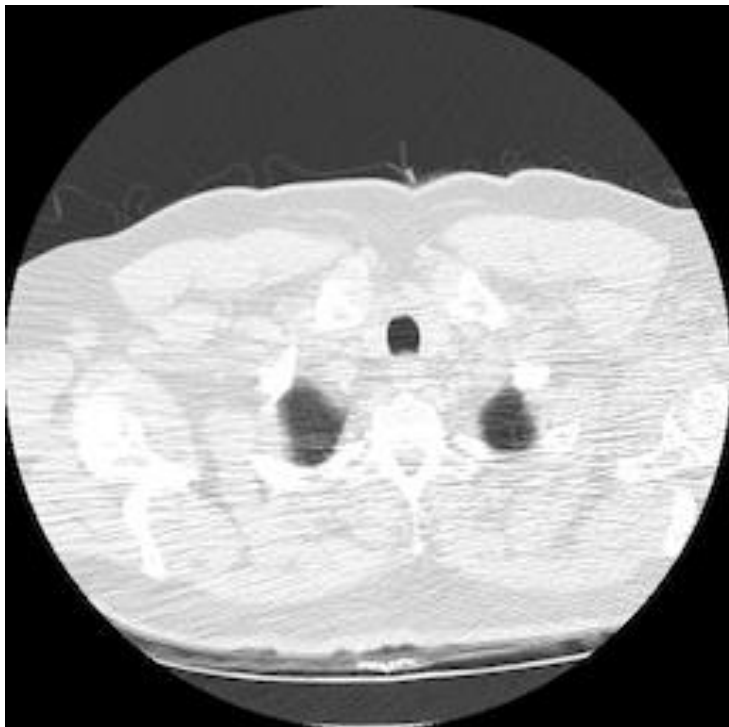
Convolutional Layer



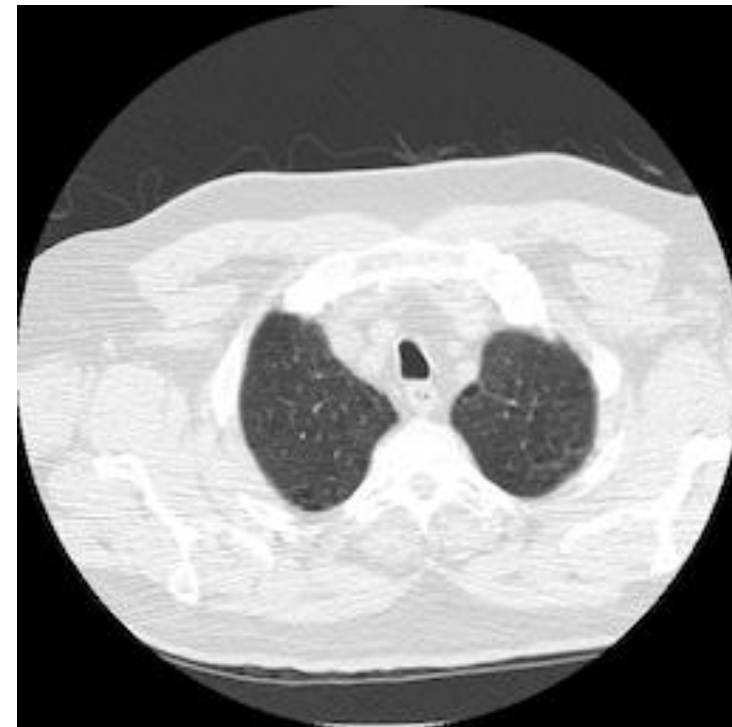
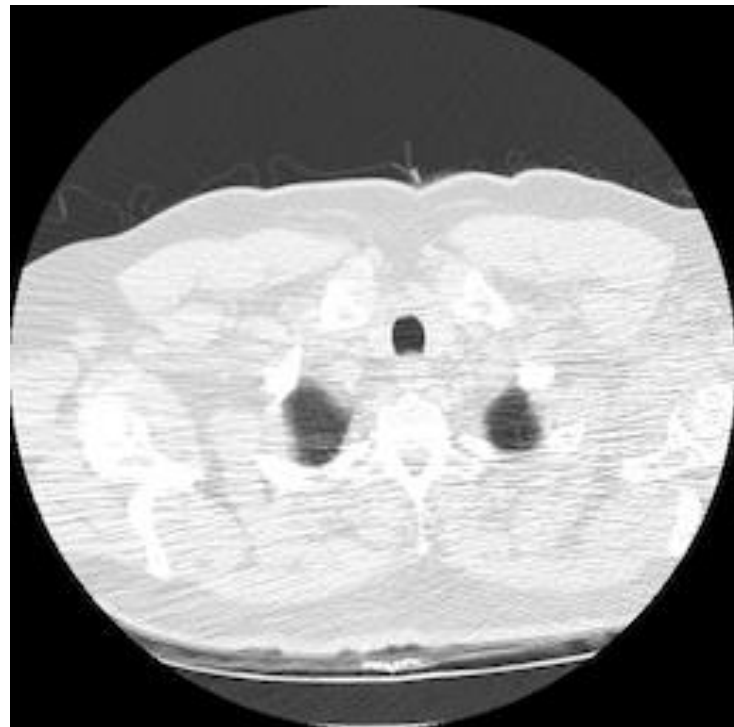
Convolution with Padding



Question: How would you apply this idea to a CT-scan?

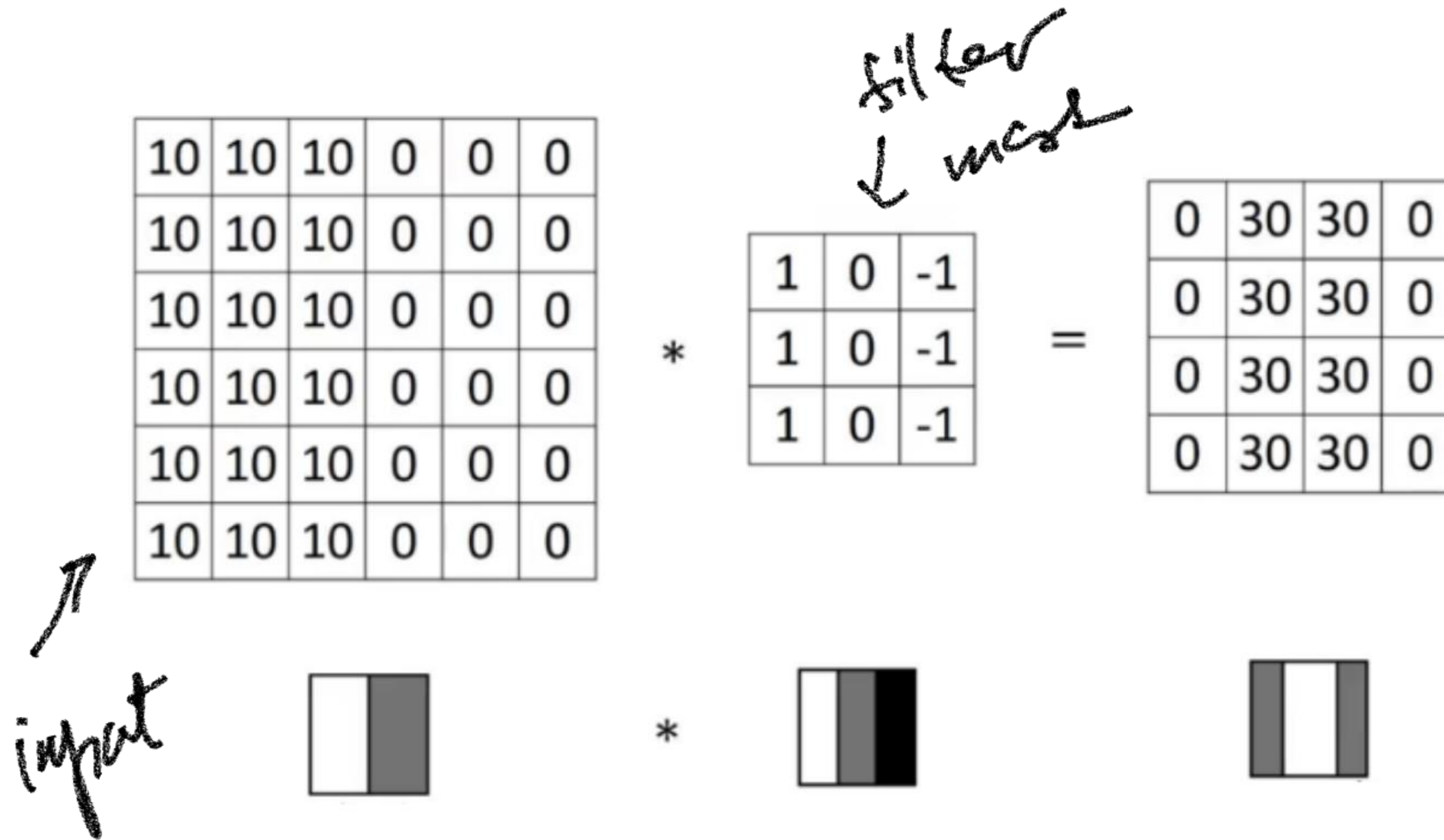


Question: How would you apply this idea to a CT-scan?

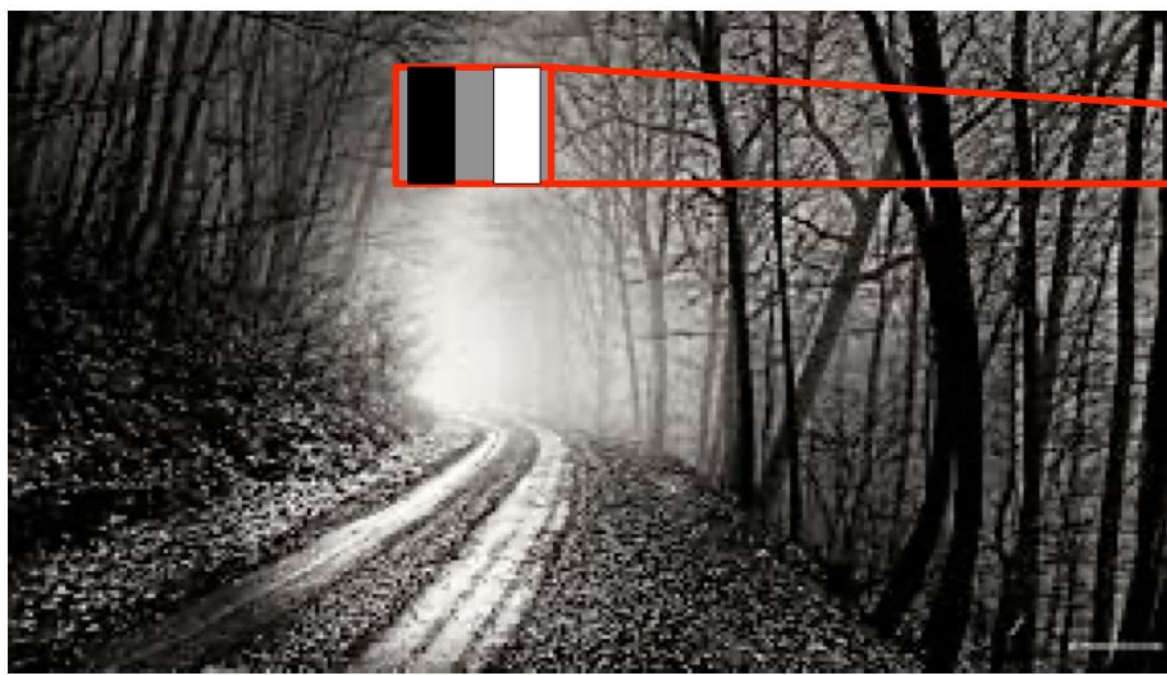


3D Convolutions

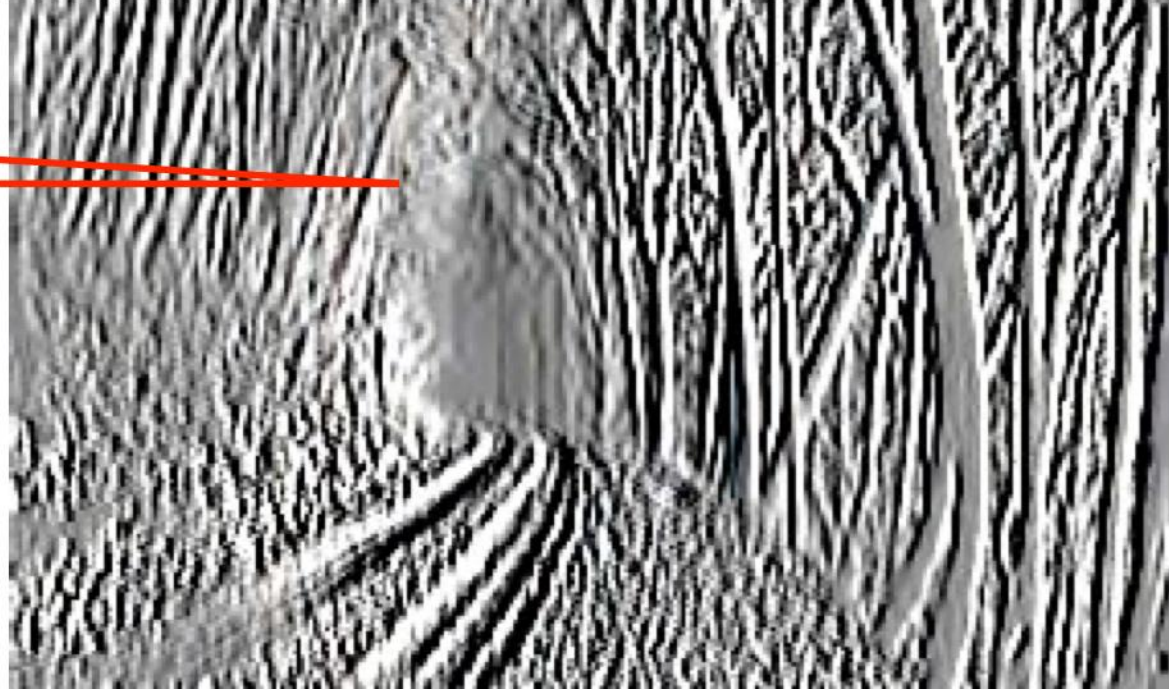
Examples of Convolutions



Examples of Convolutions


$$\ast \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix} =$$



mask



(vector, NOT a matrix!)



Examples of Convolutions

Edge detection

 * $\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$ = 

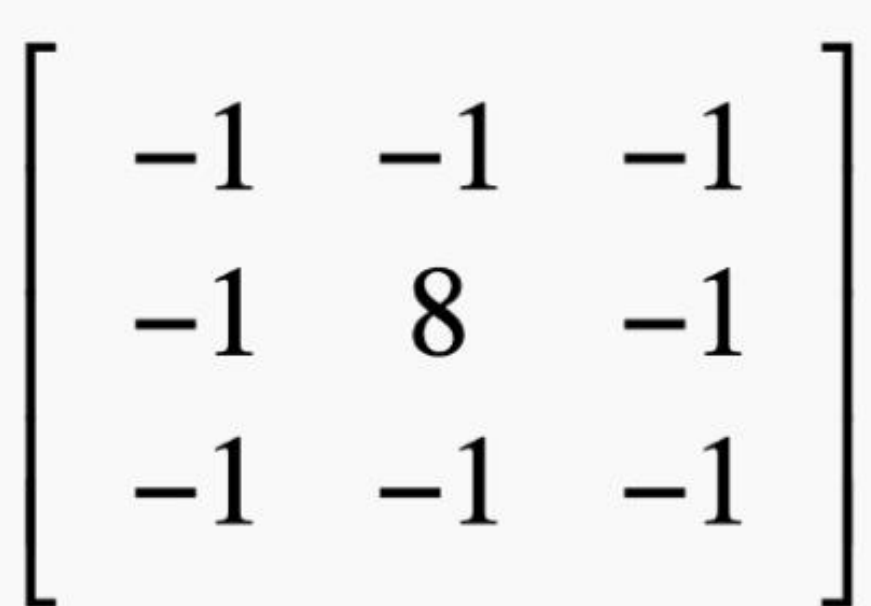
Kernel

Sharpen



 * $\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$ = 

Question: How can we make convolutions more expressive?

Edge detection

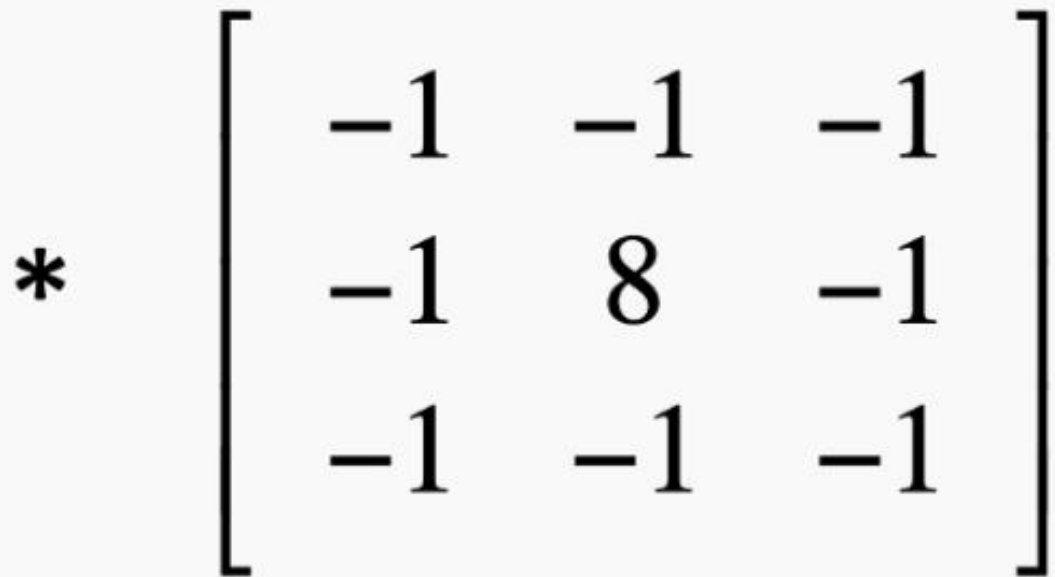

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

Kernel



$$*$$
$$=$$


Question: How can we make convolutions more expressive?

Edge detection


$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

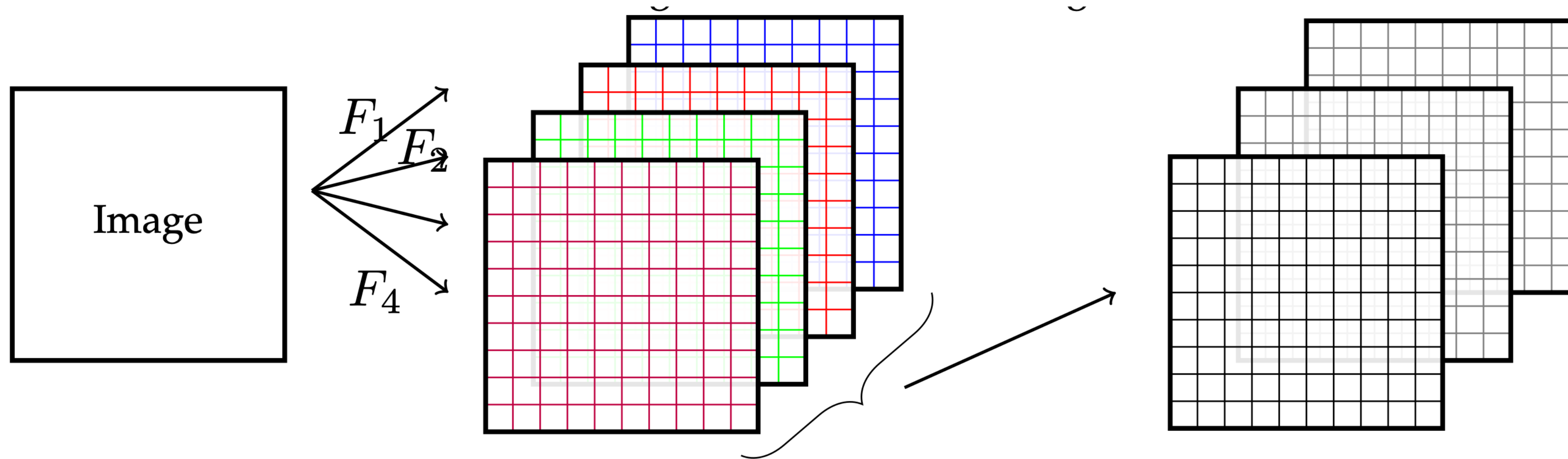
=



Kernel

Width: Many kernels in parallel
Depth: Composing kernels

Multiple channels/filters

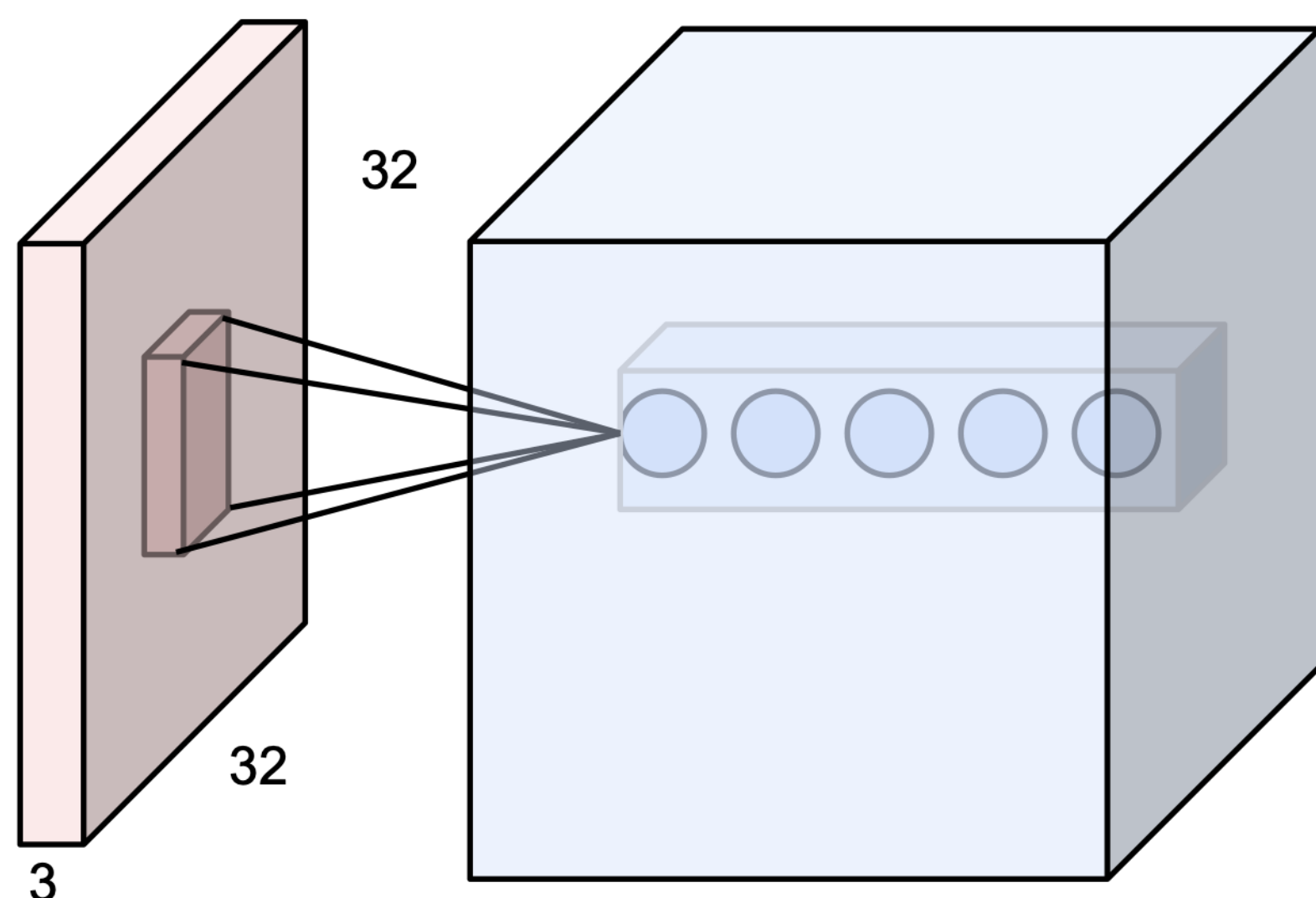


Filter bank:
Collection of
filters in a layer

Channels:
outputs of
convolution

$$F_i \rightarrow \Theta_i \in \mathbb{R}^{(c, w, h)} \text{ \& kernel size}$$

Naming Conventions



Hidden layer of “depth” 5:
five neurons all looking at
the same patch; five
different masks.

Apply the same 5 masks to
each patch. Five neurons
per patch.

Conv w param Θ

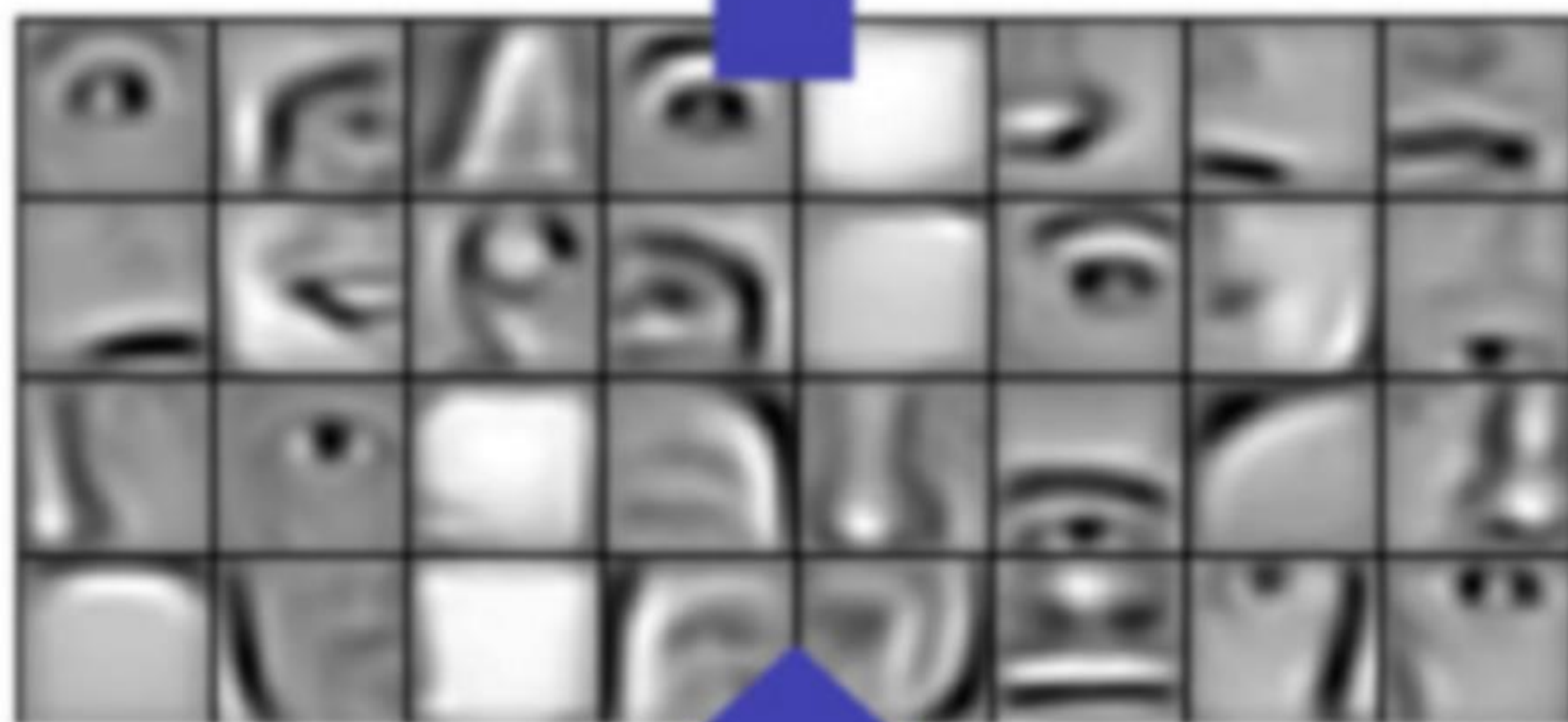
$$\dim(\Theta) \rightarrow [C_{in}, w, h, C_{out}]$$

\uparrow	\uparrow	\uparrow	\uparrow
1	3	3	128

Examples
of
Composition
in the
wild!



Layer 3

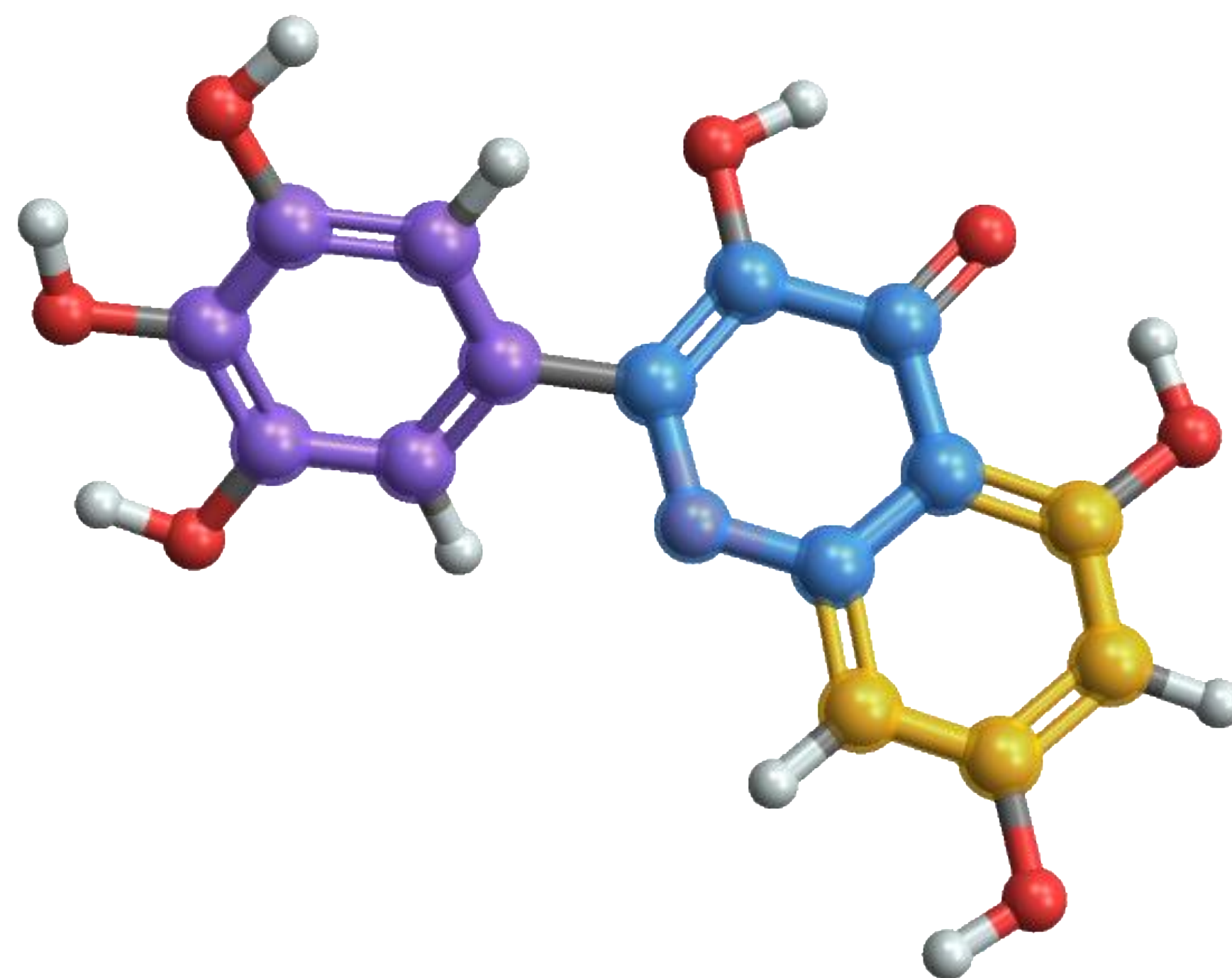


Layer 2

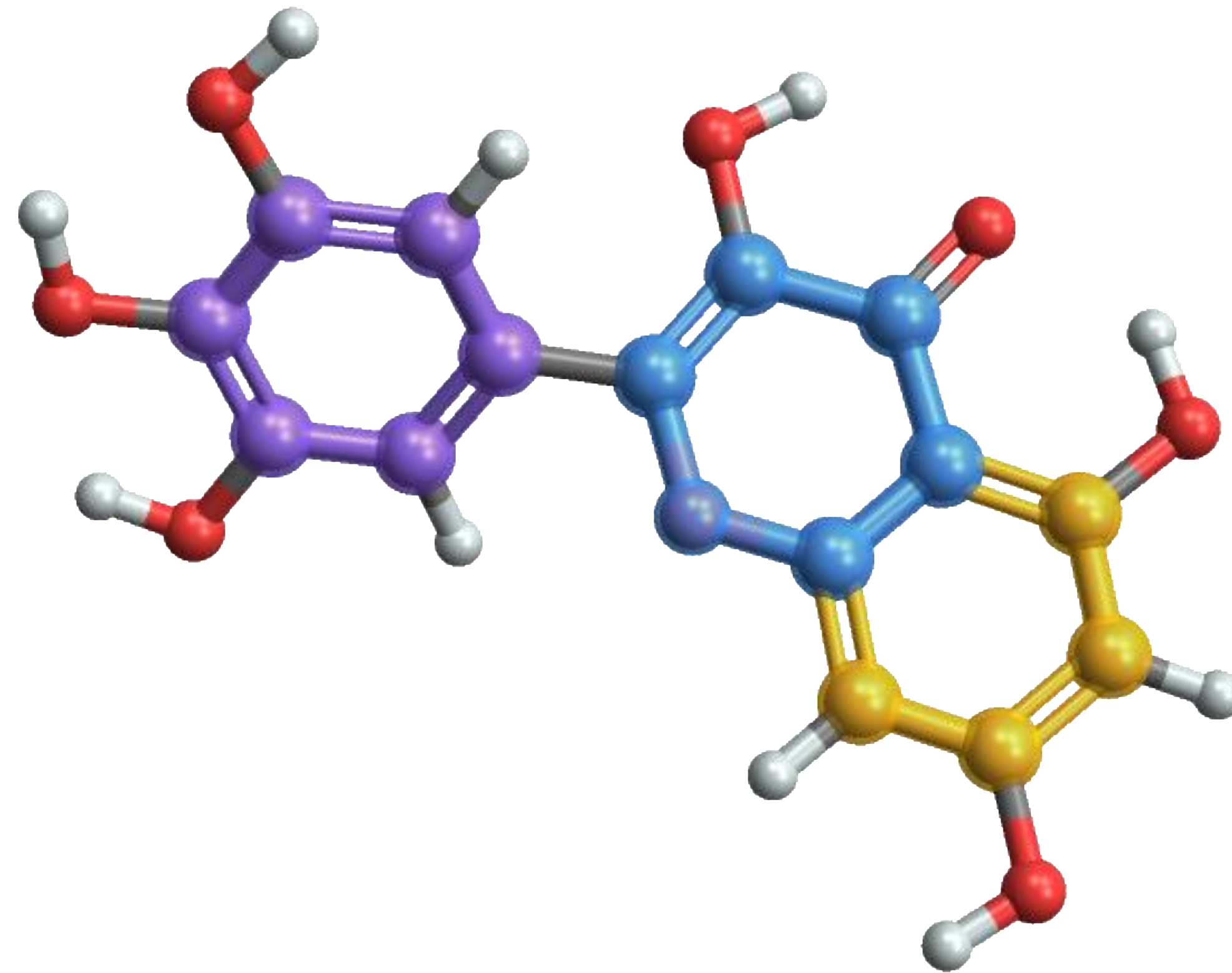


Layer 1

Question: How would you apply this idea to a graph?

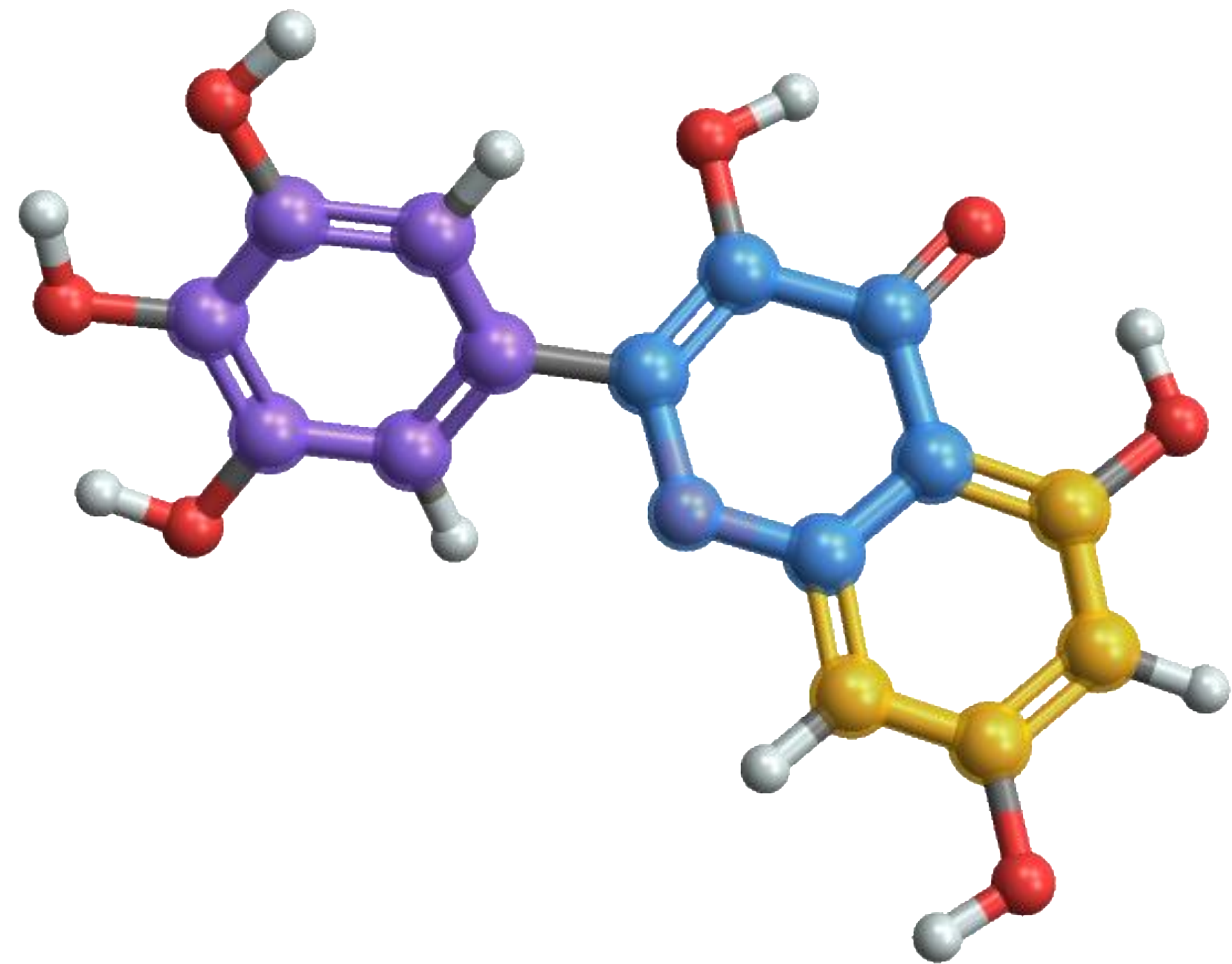


Question: How would you apply this idea to a graph?



Graph Convolutions: Leverage neighboring nodes

Question: How would you apply this idea to a graph?



connected neighbors

$$\text{Graph Conv}(X)[i] = \theta_b \otimes i + \sum_j \theta_{ij} \otimes X[j]$$

Graph Convolutions: Leverage neighboring nodes

Agenda

Recap

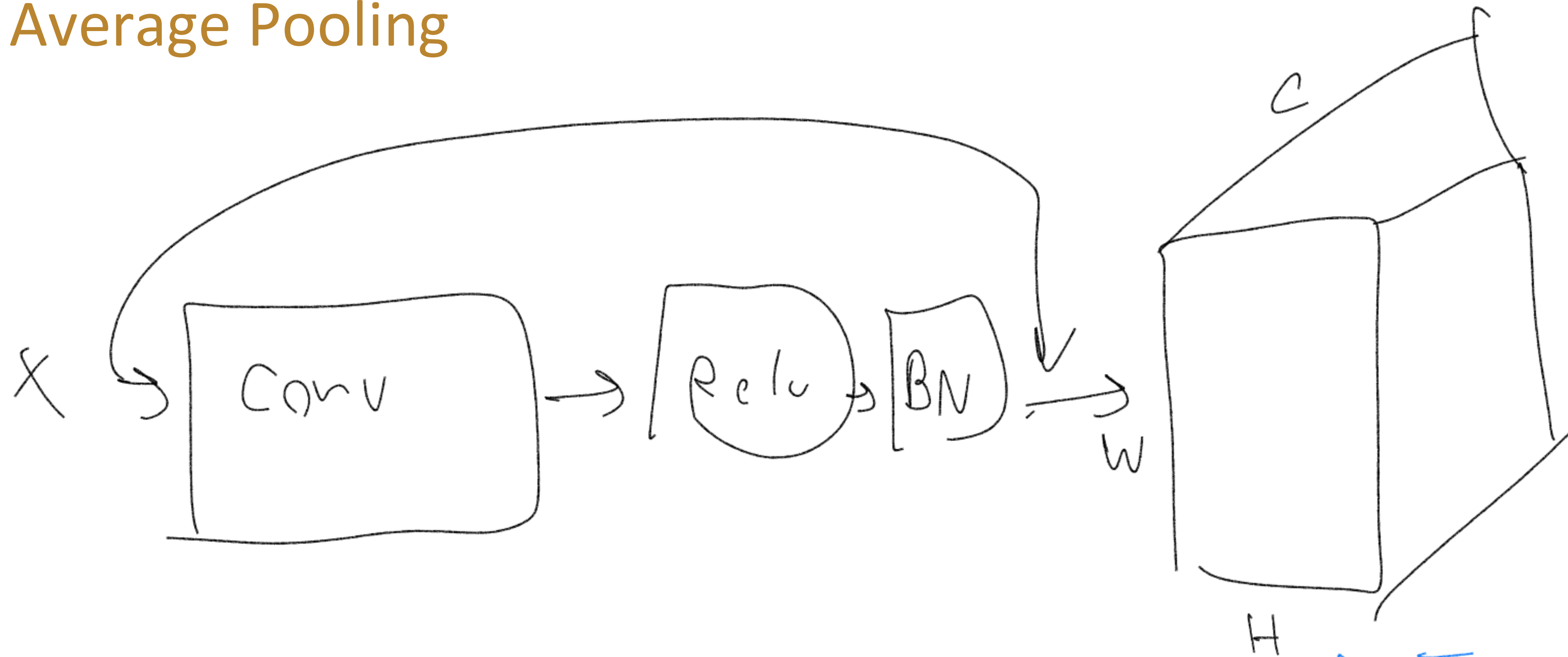
Failure modes of fully-connected neural networks

Convolutions

Pooling: Aggregating features and location invariance

CNNs across modalities

Average Pooling

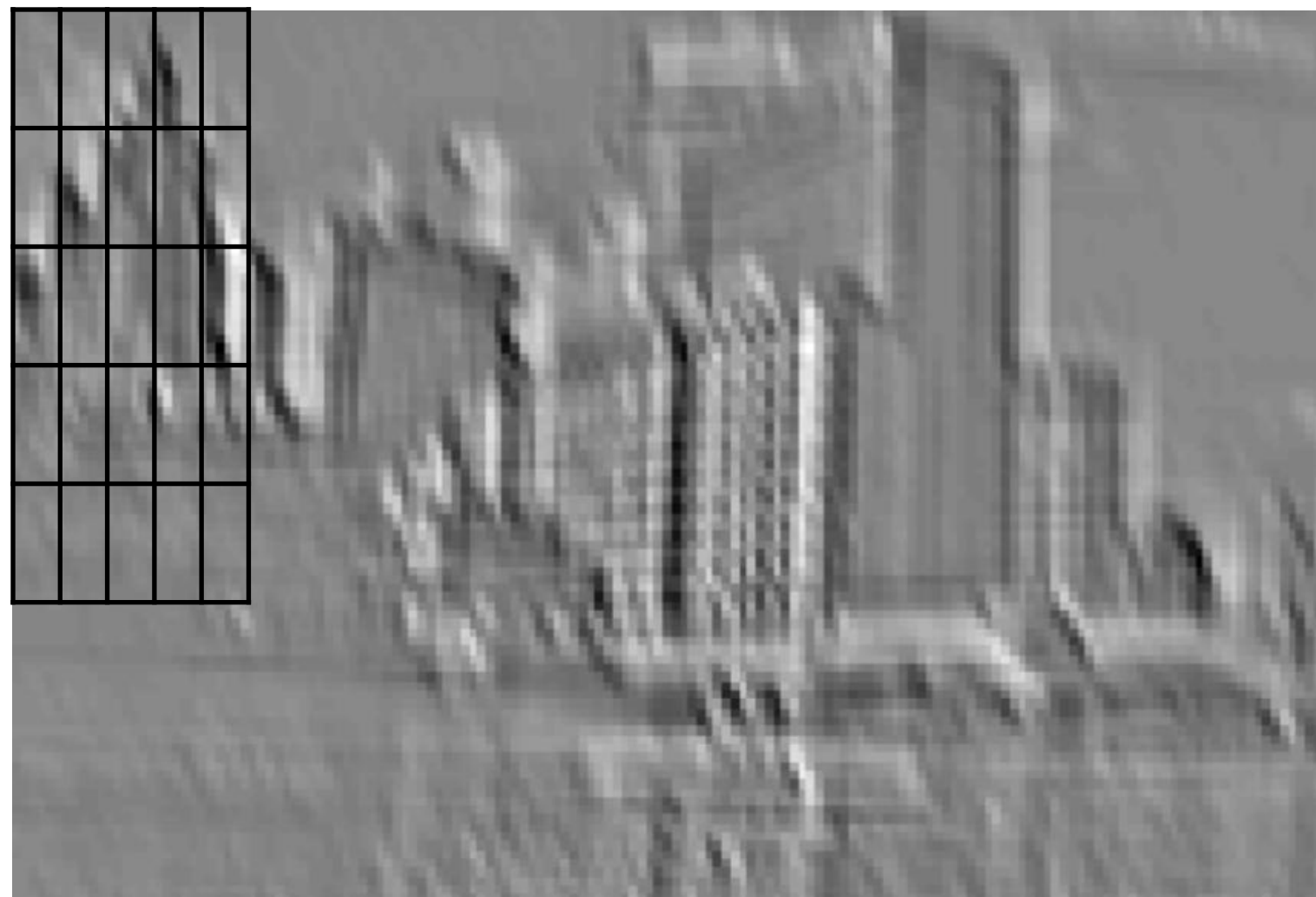


$$\text{Avg Pool}(z) = \frac{1}{H \cdot W} \sum_i^H \sum_j^W z[i, j]$$

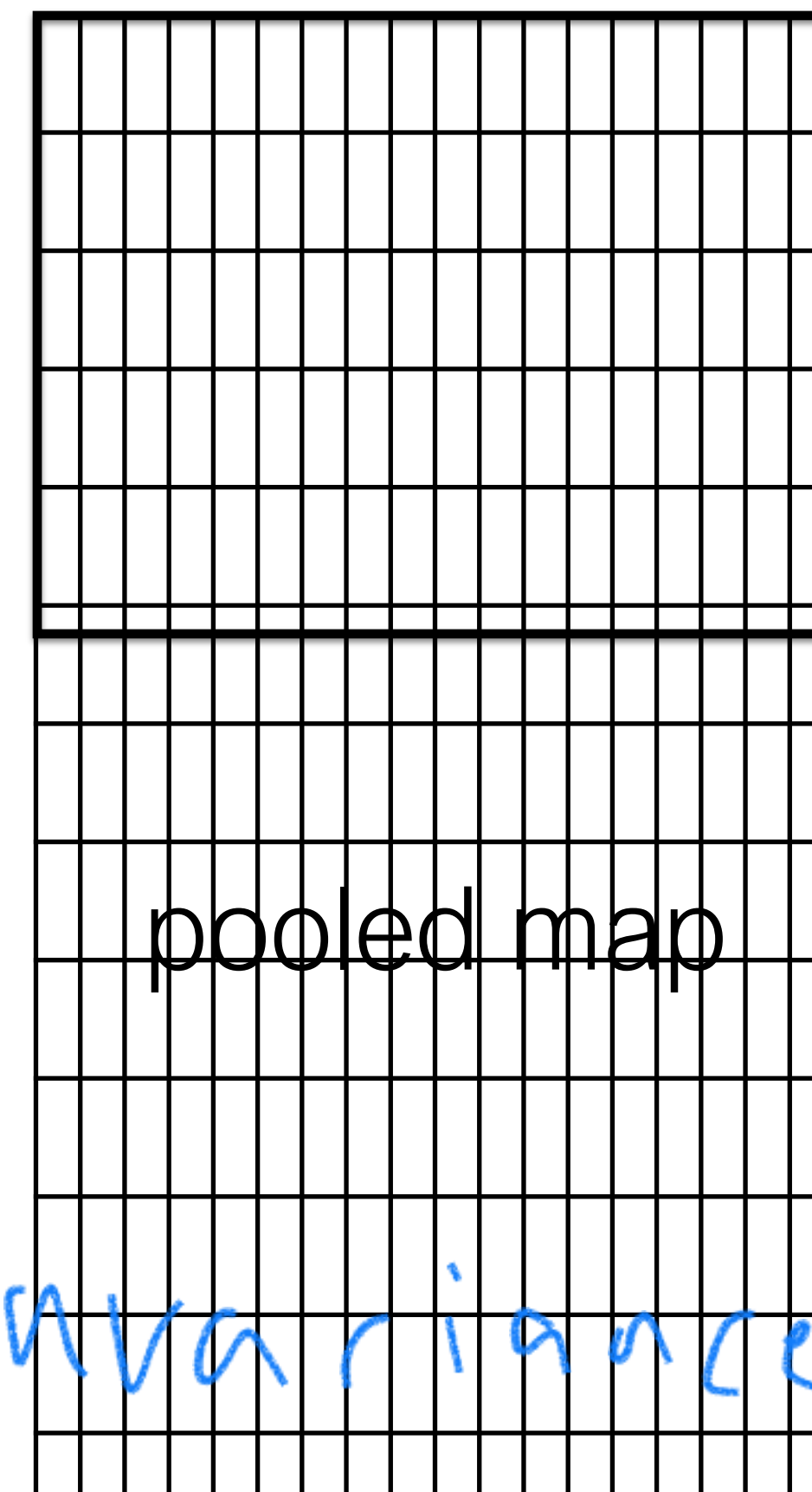
Translation
invariance!

Pooling

- We wish to know whether a feature was there but not exactly where it was



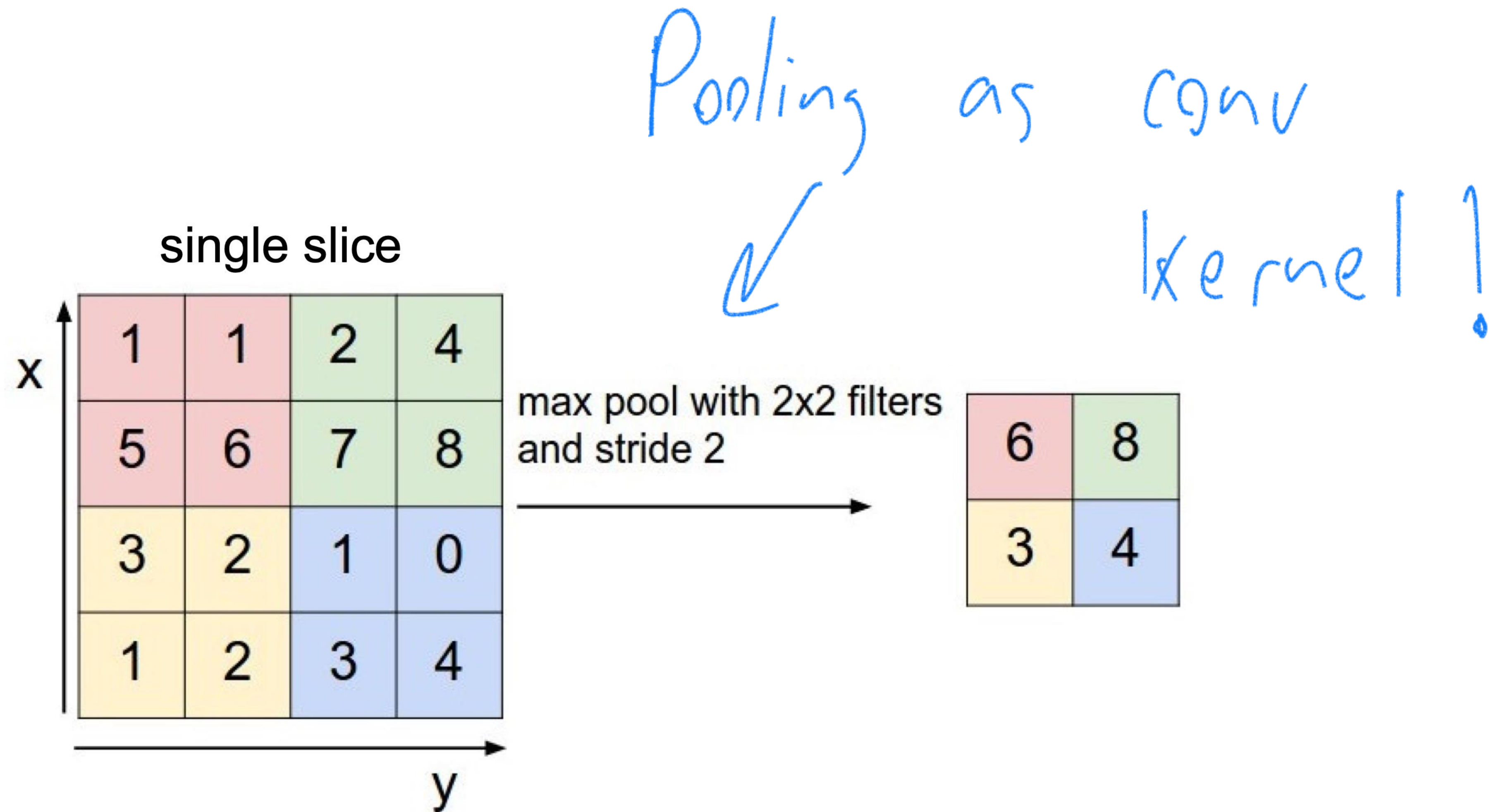
feature map



pooled map

Translation invariance!

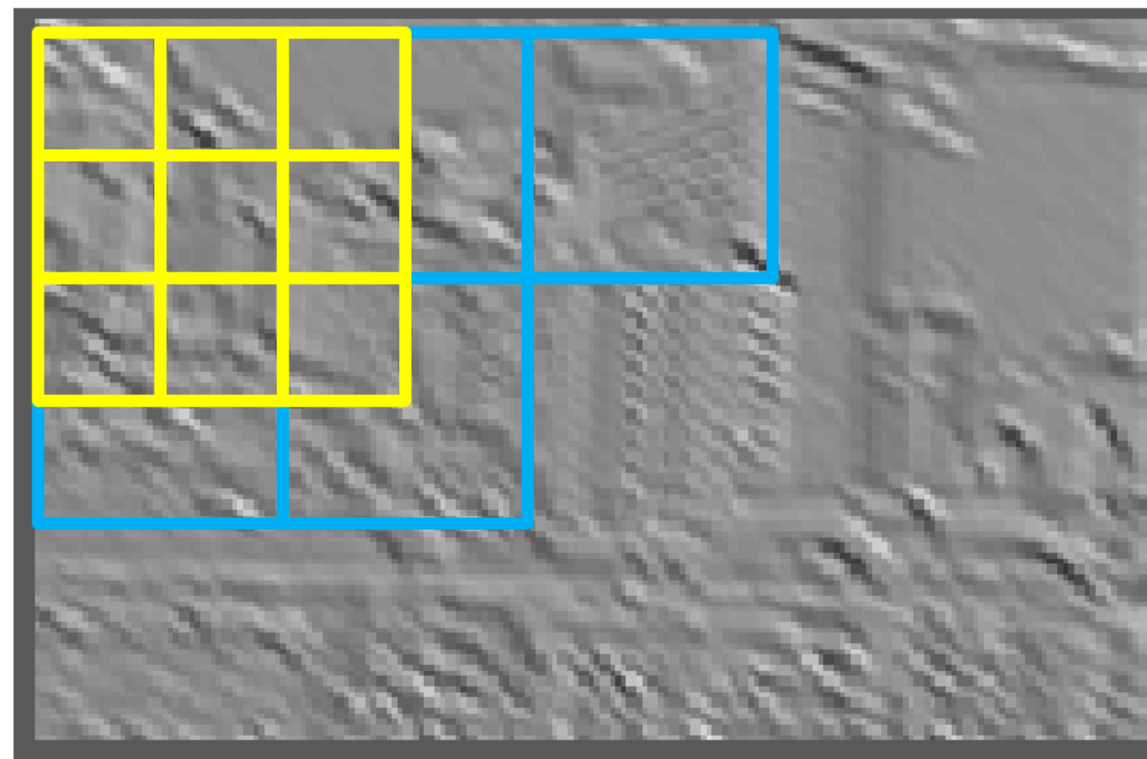
Max Pooling



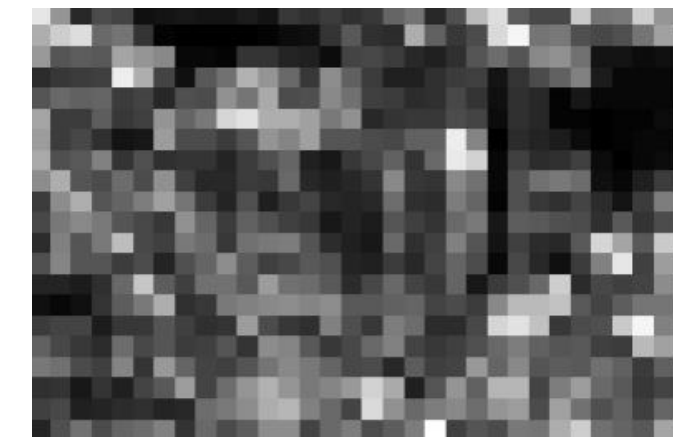
- Similar to filtering, but output the maximum entry instead of a weighted sum

Pooling (max)

- Pooling region and “stride” may vary
 - pooling induces translation invariance at the cost of spatial resolution
 - stride reduces the size of the resulting feature map

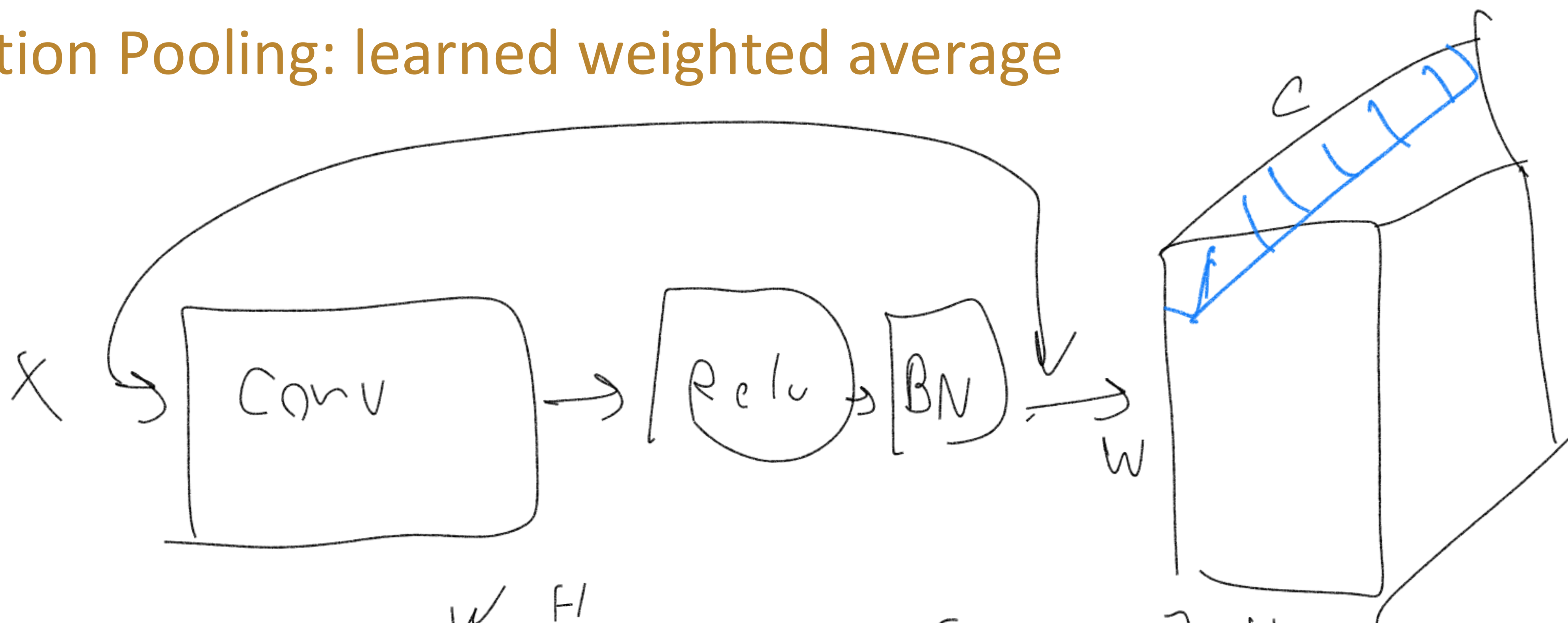


feature map



feature map
after max pooling

Attention Pooling: learned weighted average

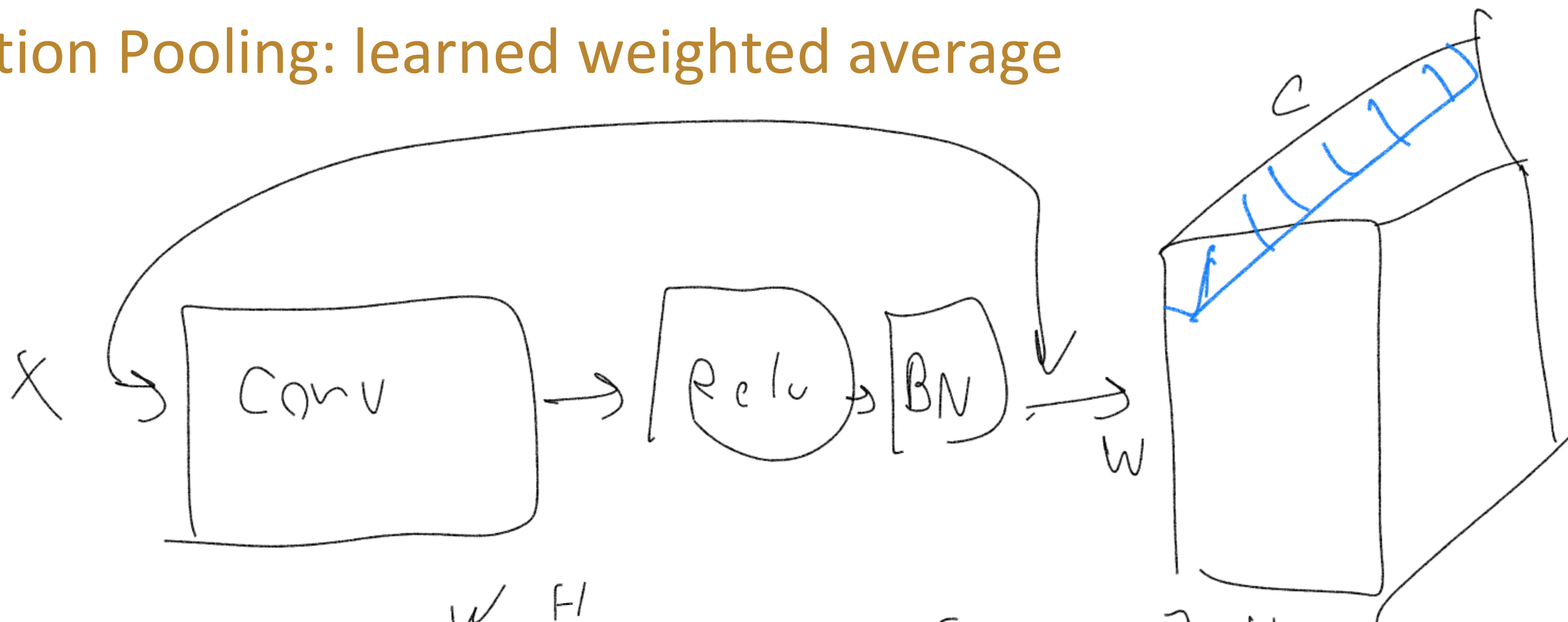


$$\text{Attn Pool}(z) = \sum_i^W \sum_j^H \alpha_{ij} z[i, j]^H$$

$$\alpha = \text{Softmax}(f(z))$$

$$|f(z)| \rightarrow (W, H, 1)$$

Attention Pooling: learned weighted average



$$\text{Attn Pool}(z) = \sum_i^W \sum_j^H \alpha_{ij} z[i, j]$$

$$\alpha = \text{softmax}(f(z))$$

$$\text{softmax}(a) = \frac{e^a}{\sum_i e^{a_i}}$$

Multi-Head Attention Pooling: more shots on goal

$$\text{Attn Pool}(z) = \sum_i^W \sum_j^F \alpha_{ij} z[i, j]$$

$$\text{MultiHead Attn}(z) = \left[\underbrace{\text{Attn Pool}(z) ; \text{Attn Pool}(z) \dots}_{\text{unique } \alpha_s!} \right]$$

Focus on different things

Agenda

Recap

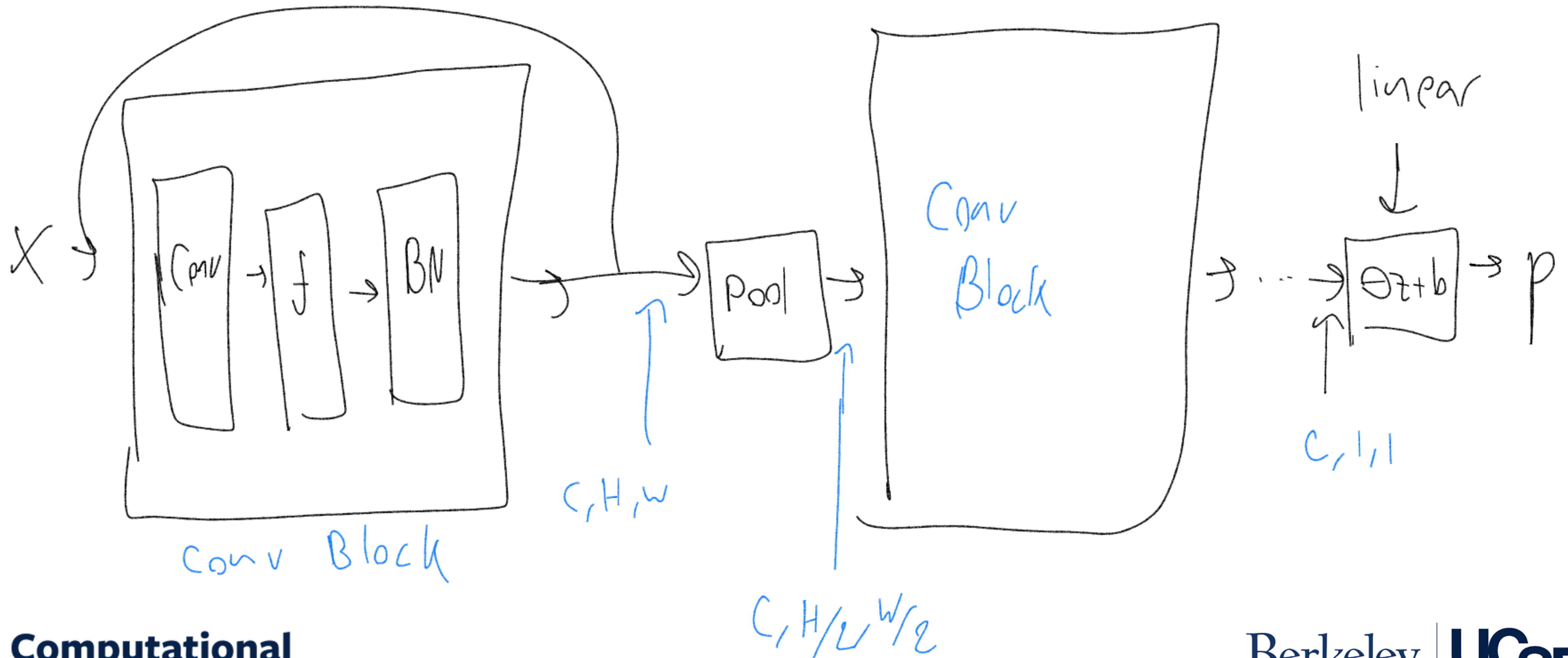
Failure modes of fully-connected neural networks

Convolutions

Pooling

CNNs across modalities

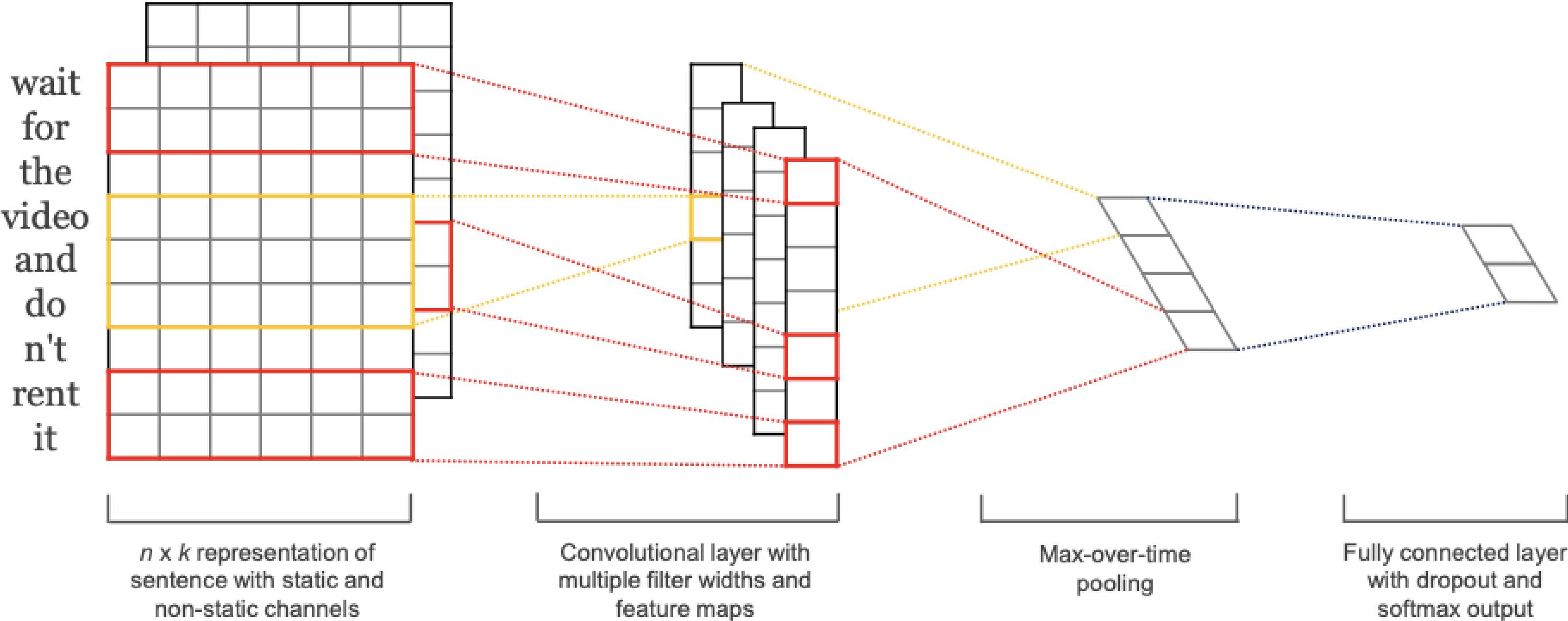
Putting it together: CNNs



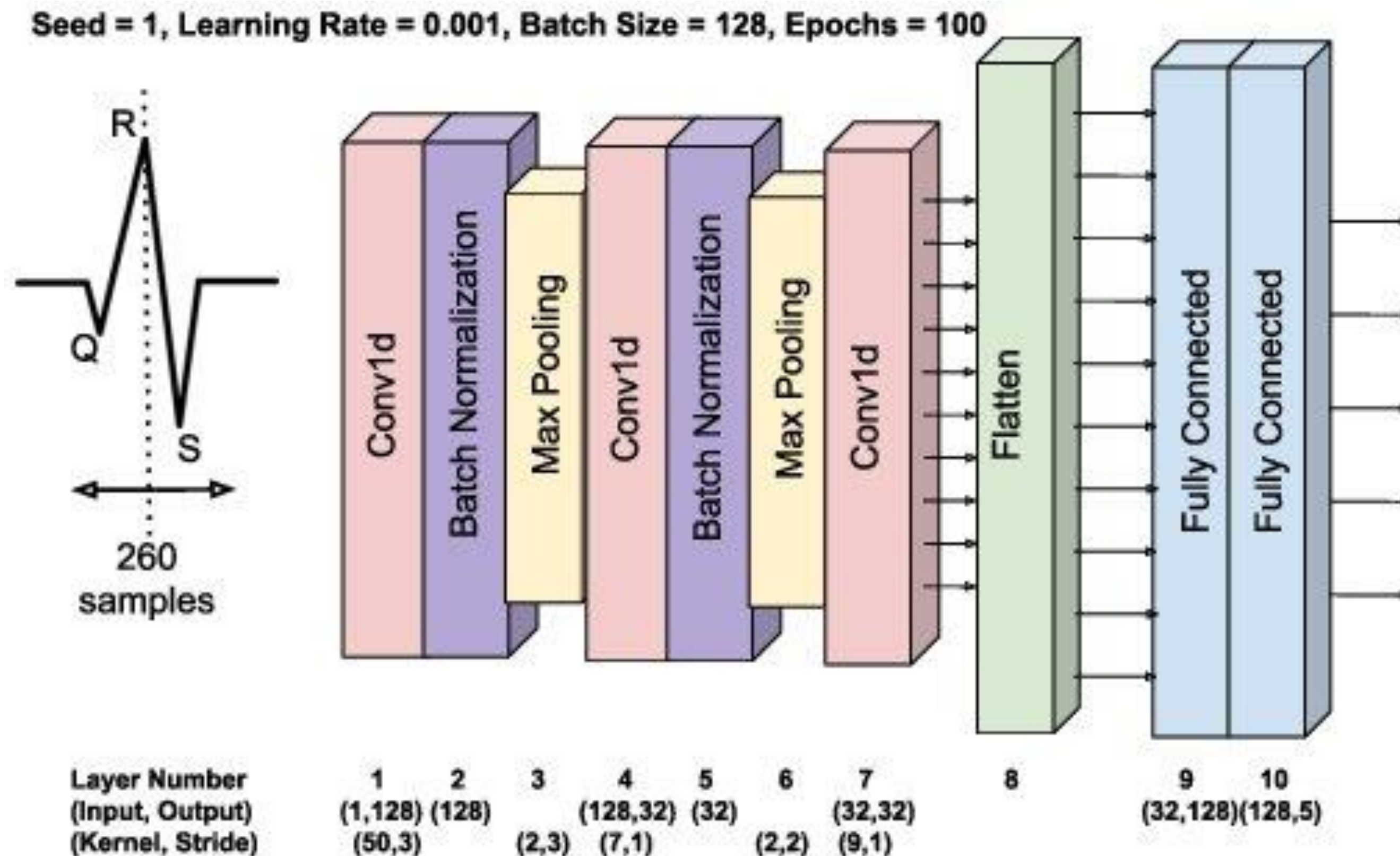
Popular 1D CNNs: Text

Convolutional Neural Networks for Sentence Classification

Yoon Kim
New York University
yhk255@nyu.edu

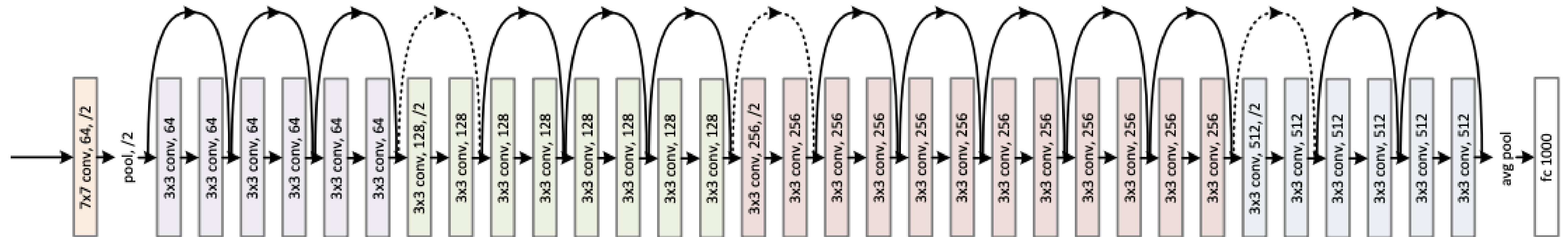


Popular 1D CNNs: Wave Forms



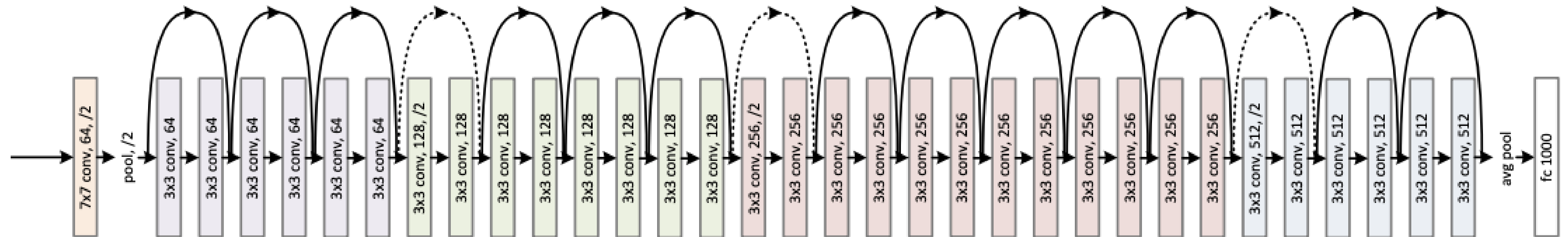
Xiaolin, Li, Barry Cardiff, and Deepu John. "A 1d convolutional neural network for heartbeat classification from single lead ecg." *2020 27th IEEE International Conference on Electronics, Circuits and Systems (ICECS)*. IEEE, 2020.

Popular 2D CNNs: ResNets



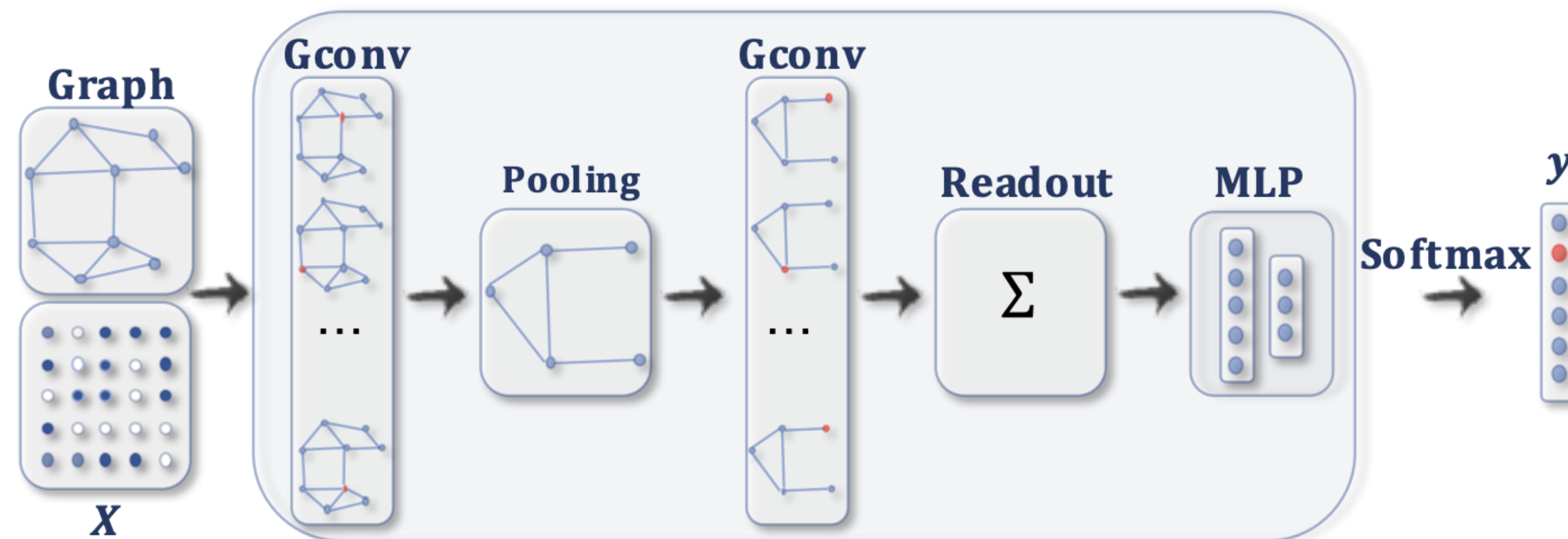
He, Kaiming, et al. "Deep residual learning for image recognition." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016.

Popular 3D CNNs: ResNet3D



Just make the kernels all 3D. Used in Sybil and many other 3D models.

Popular GNNs: Convolutions on graphs



(b) A ConvGNN with pooling and readout layers for graph classification [21]. A graph convolutional layer is followed by a pooling layer to coarsen a graph into sub-graphs so that node representations on coarsened graphs represent higher graph-level representations. A readout layer summarizes the final graph representation by taking the sum/mean of hidden representations of sub-graphs.

Wu, Zonghan, et al. "A comprehensive survey on graph neural networks." *IEEE transactions on neural networks and learning systems* 32.1 (2020): 4-24.

Summary

FFNs are **wildly** inefficient

Convolutions: Capture **local patterns** data

Pooling: **Spatial invariant** method to summarize features

Attention Pooling: Parameterized Weighted Averages

CNNs: NNs with Conv and Pooling building blocks

Applications to text, images, graphs, and more

Questions?