# Lecture 9
# CNNs and Spatial Processing

# 9. CNNs and Spatial Processing

- How to use deep nets for images

- New layer types: convolutional, pooling

- Feature maps and multichannel representations

- Popular architectures: Alexnet, VGG, Resnets

- Getting to know learned filters

- Unit visualization
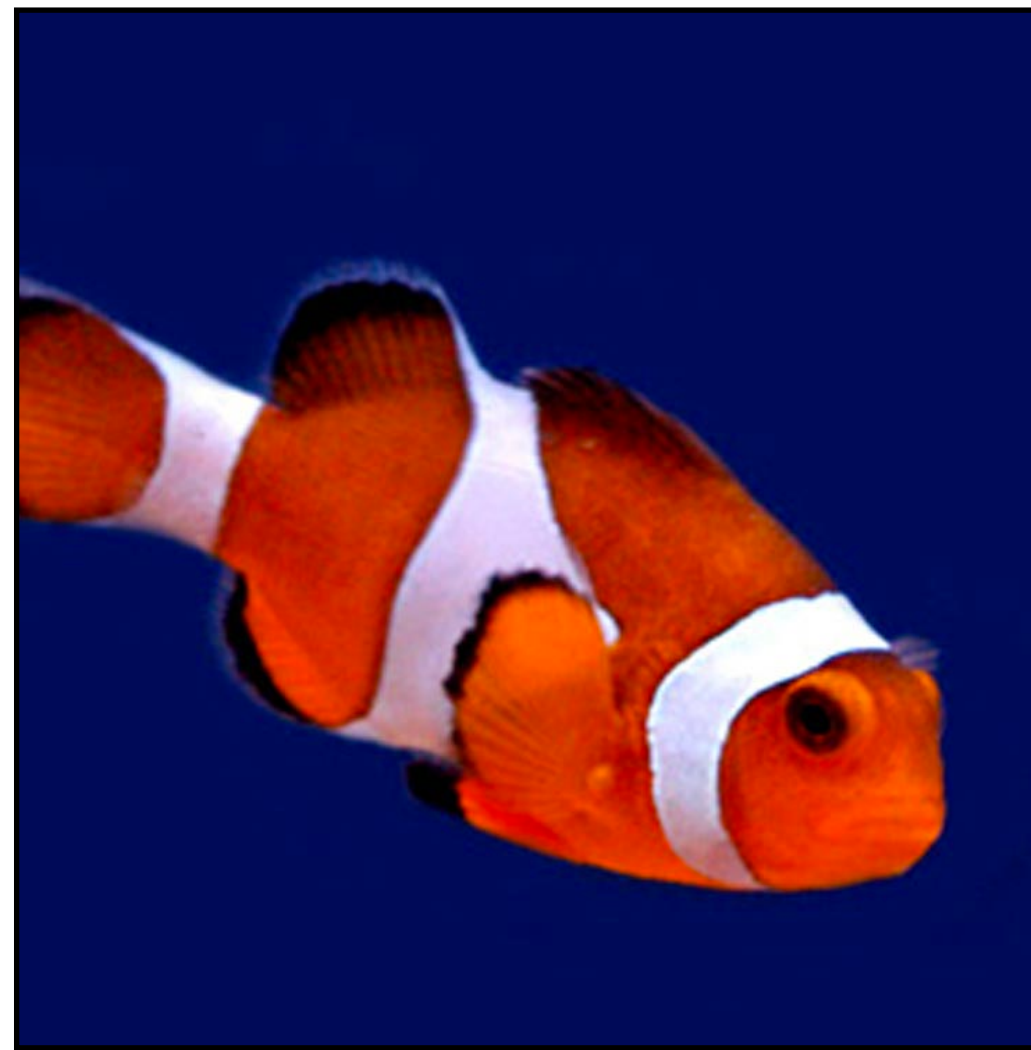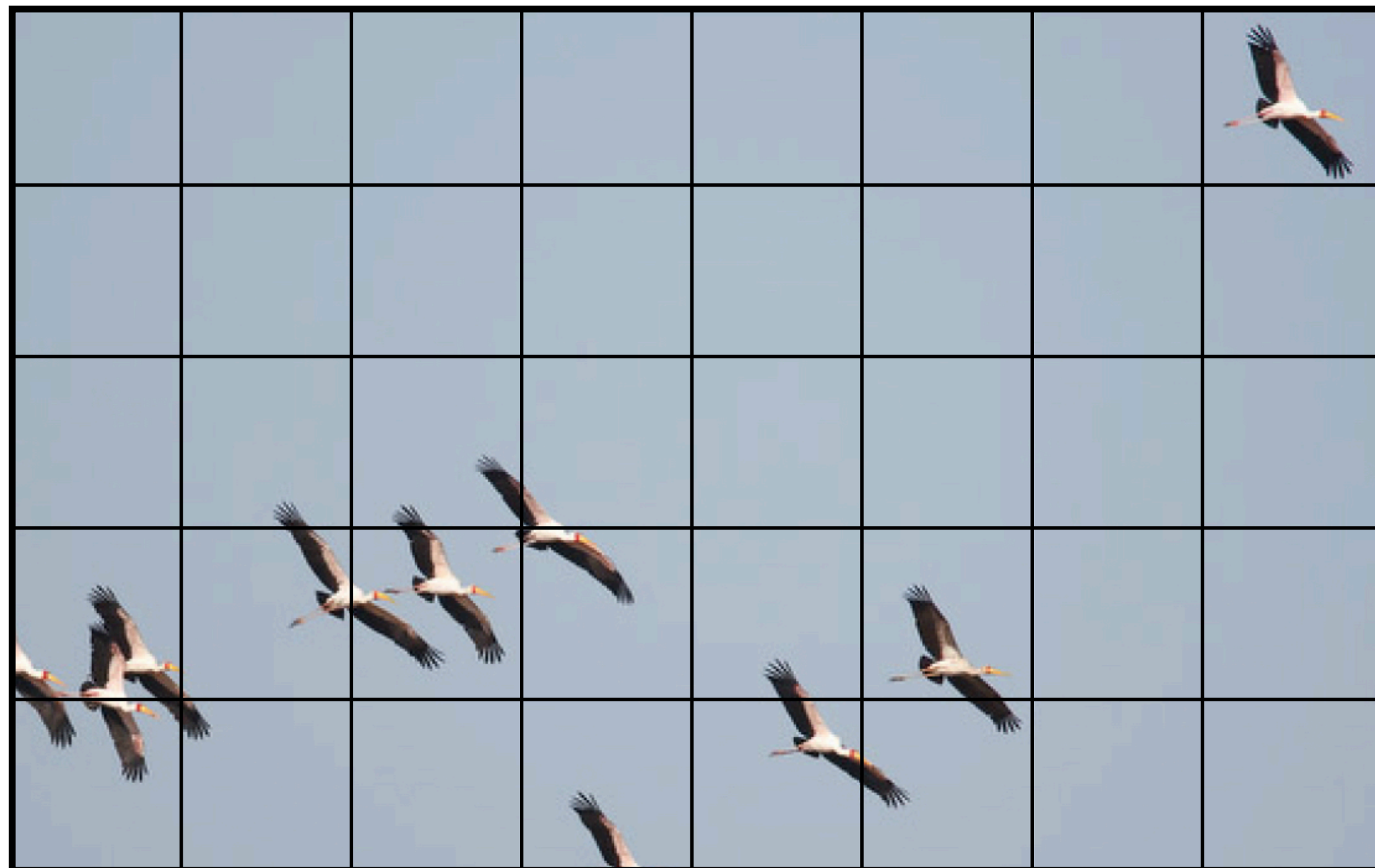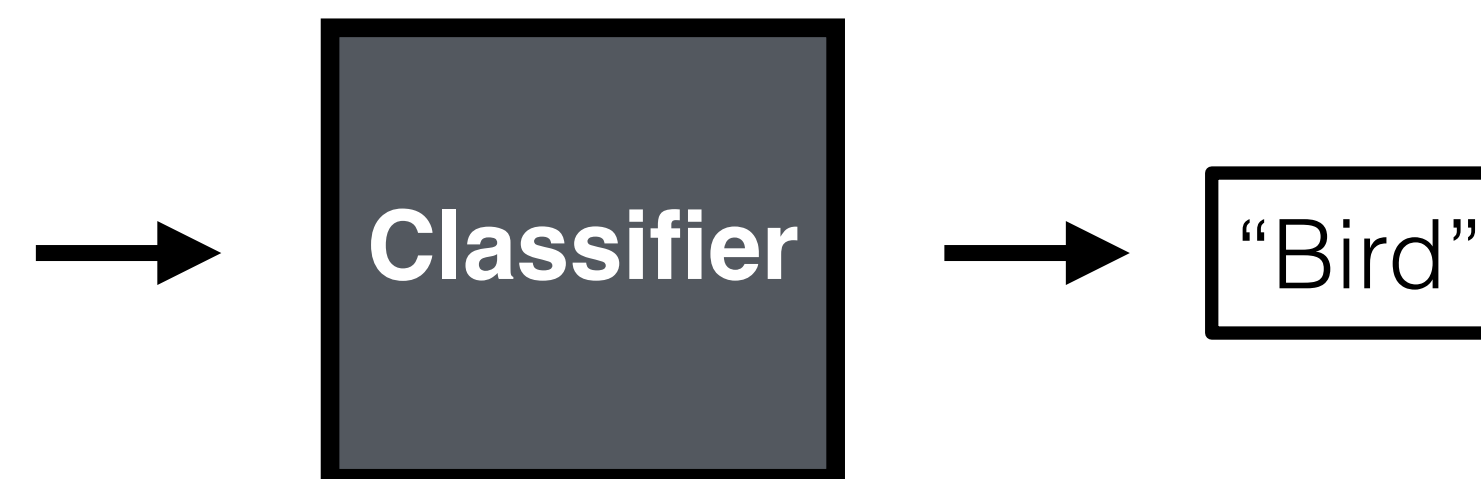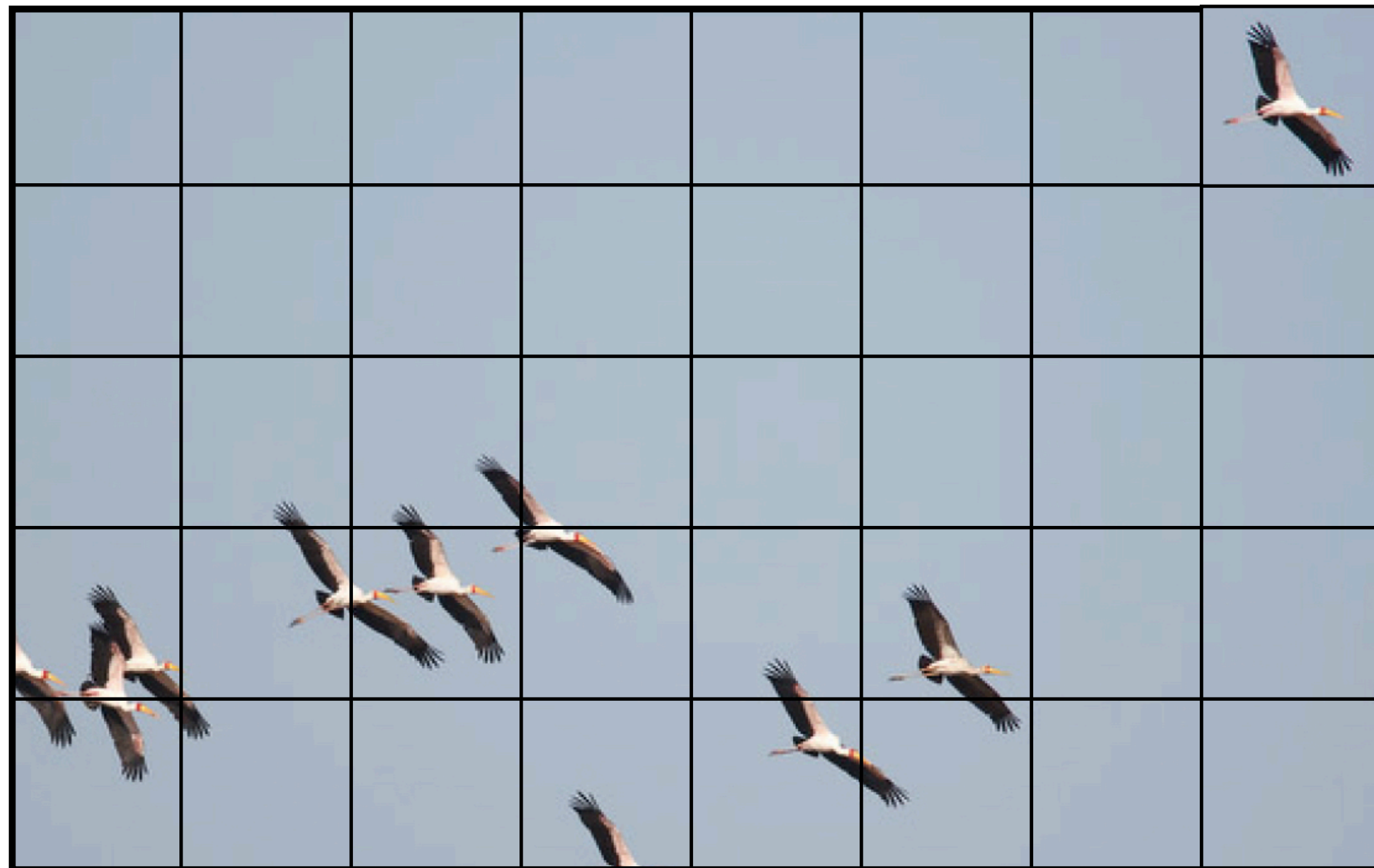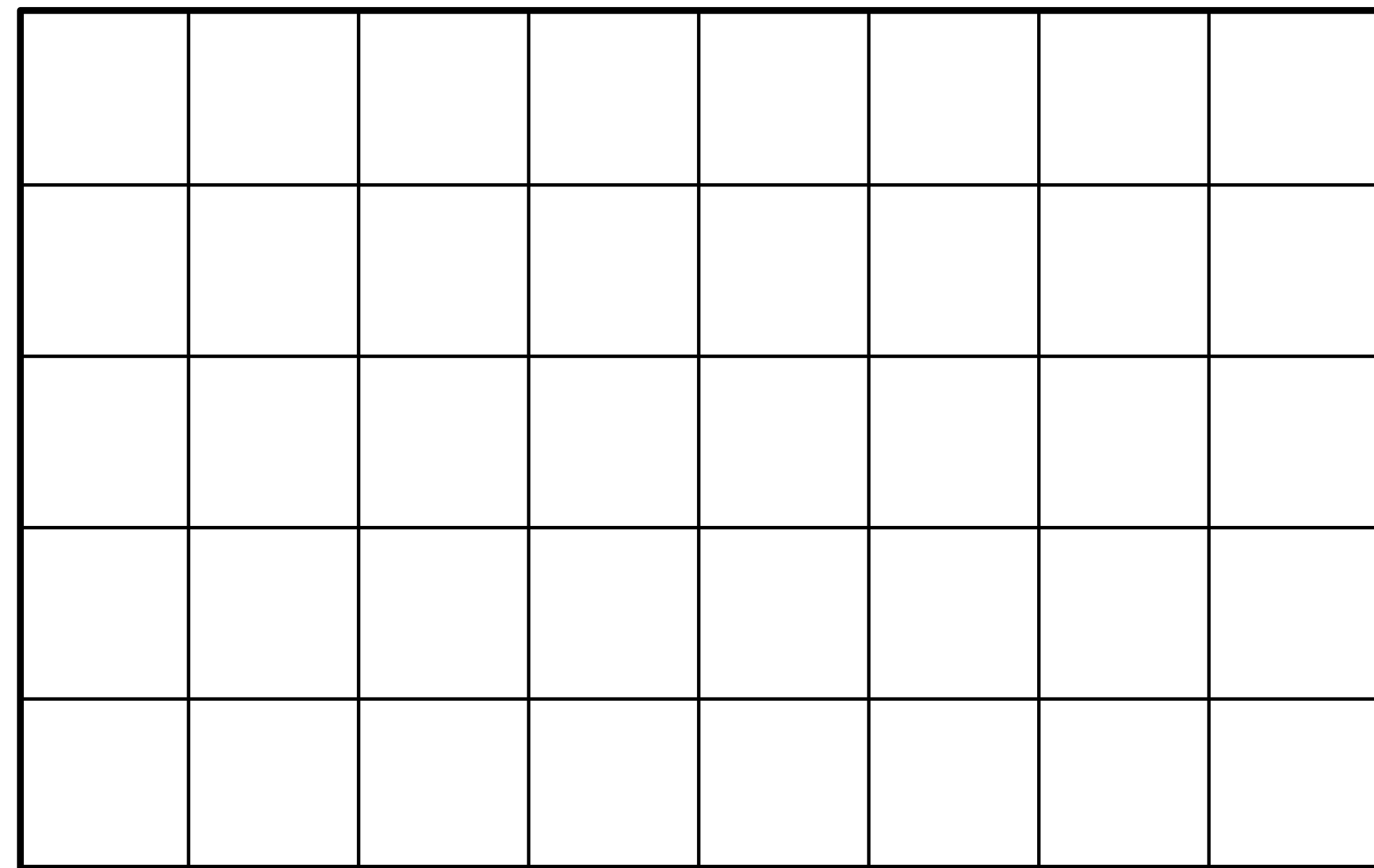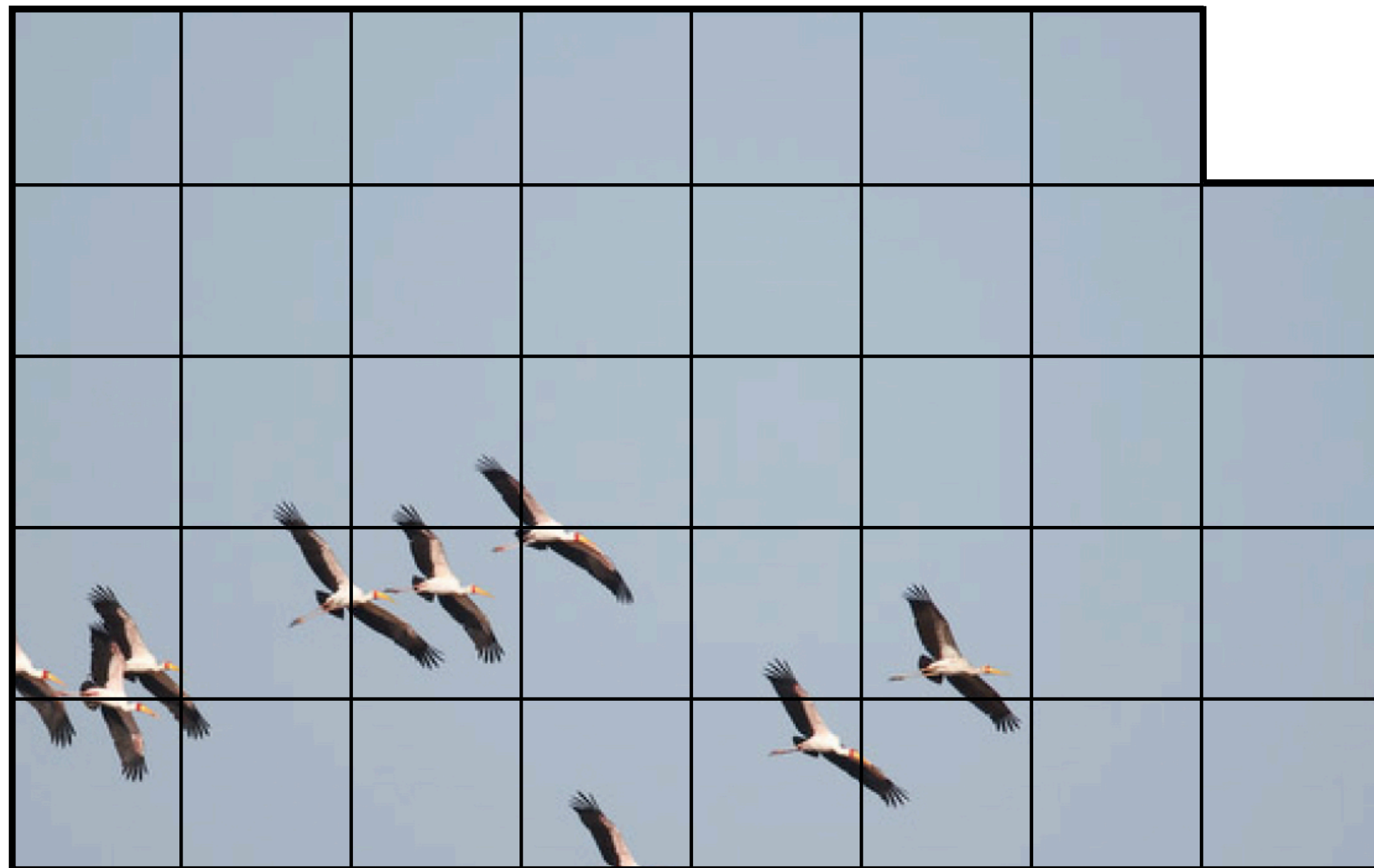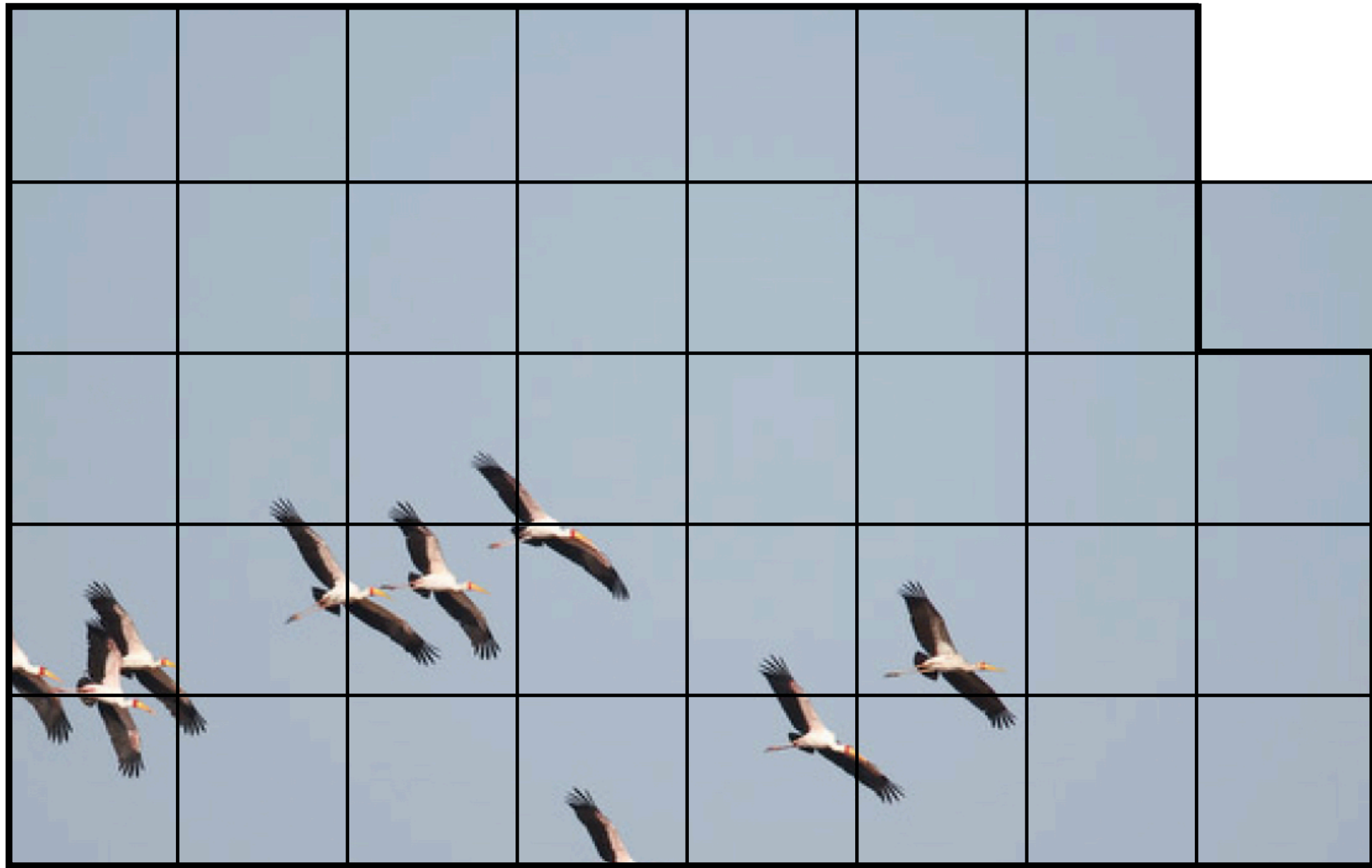
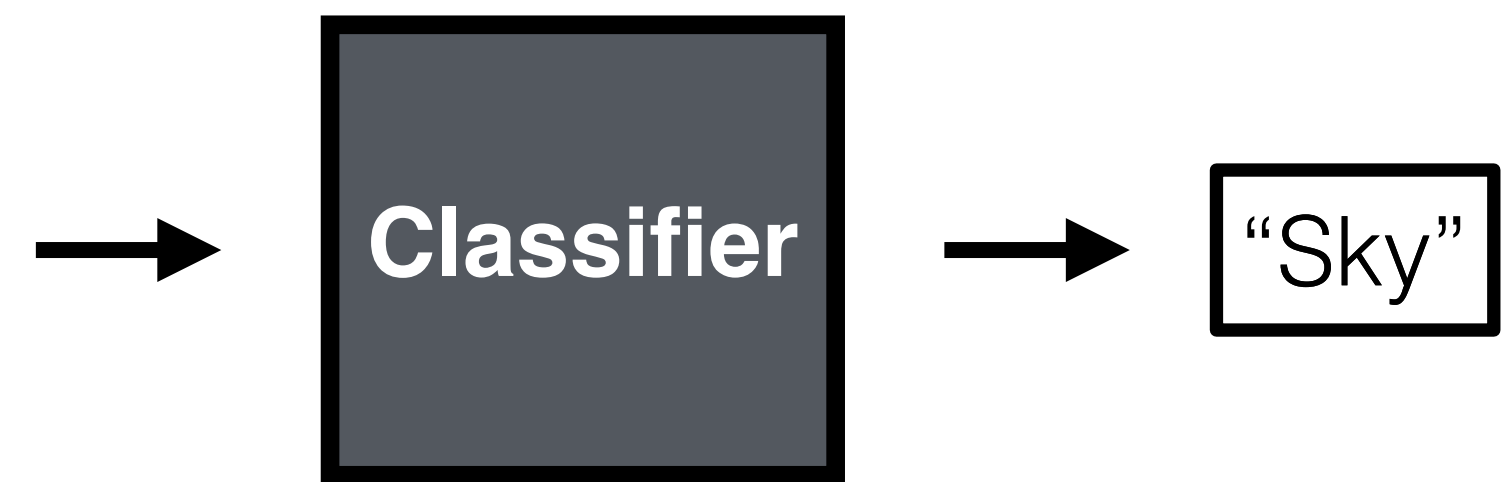# Image classification



image **x**

**Classifier**

"Fish"

label y

Photo credit: Fredo Durand

**Classifier** → "Bird"

Bird

Classifier → "Sky"

| Sky | Sky | Sky | Sky | Sky | Sky | Sky | Bird |
|------|------|------|------|------|------|------|------|
| Sky | Sky | Sky | Sky | Sky | Sky | Sky | Sky |
| Sky | Sky | Sky | Sky | Sky | Sky | Sky | Sky |
| Bird | Bird | Bird | Sky | Bird | Sky | Sky | Sky |
| Sky | Sky | Sky | Bird | Sky | Sky | Sky | Sky |

**Problem:**

What happens to objects that are bigger?

What if an object crosses multiple cells?

"Cell"-based approach is limited.

What can we do instead?

What's the object class of the center pixel?

$f$ → "Bird"

$f$ → "Bird"

$f$ → "Sky"

$f$ → "Sky"

Training data

$\mathbf{x}$      $y$

{ "Bird" ,

{ "Bird" ,

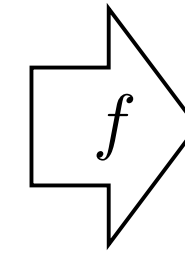{ "Sky" ,

⋮

What's the object class of the center pixel?

$f$ → "Bird"

$f$ → "Bird"

$f$ → "Sky"

$f$ → "Sky"

$f$

(Colors represent one-hot codes)

This problem is called **semantic segmentation**

What's the object class of the center pixel?

"Bird"

"Bird"

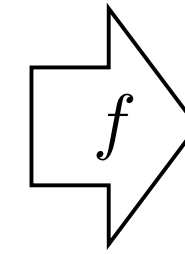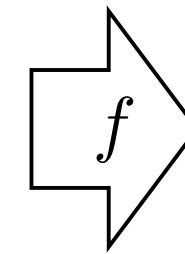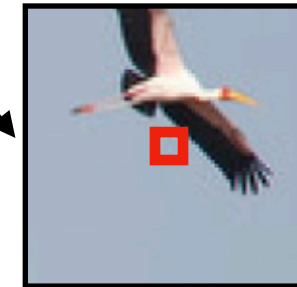"Sky"

"Sky"

Translation invariance: process each patch in the same way.

An *equivariant* mapping:

$$f(\texttt{translate}(x)) = \texttt{translate}(f(x))$$

**W** computes a weighted sum of all pixels in the patch



**W**

**W**

**W**

**W** is a convolutional kernel applied to the full image!

# Convolution



filter

# Fully-connected network

**Fully-connected (fc) layer**

# Locally connected network



Often, we assume output is a **local** function of input.

If we use the same weights (**weight sharing**) to compute each local function, we get a convolutional neural network.

# Convolutional neural network

**Conv layer**



$$\mathbf{z} = \mathbf{w} \circ \mathbf{x} + \mathbf{b}$$

$\mathbf{x}$

$\mathbf{z}$  $g(\mathbf{z})$

Often, we assume output is a **local** function of input.

If we use the same weights (**weight sharing**) to compute each local function, we get a convolutional neural network.

# Weight sharing

**Conv layer**



$$\mathbf{z} = \mathbf{w} \circ \mathbf{x} + \mathbf{b}$$

Often, we assume output is a **local** function of input.

If we use the same weights (**weight sharing**) to compute each local function, we get a convolutional neural network.

# Linear system: $\mathbf{y} = f(\mathbf{x})$

A linear function f can be written as a matrix multiplication:

$$\begin{bmatrix} y \end{bmatrix} = \begin{bmatrix} & & & \\ & & h\,[n,k] & \\ & & & \end{bmatrix} \begin{bmatrix} x \end{bmatrix}$$

n indexes rows,
k indexes columns

It can also be represented as a fully connected linear neural network



$$y\,[n] = \sum_{k=0}^{N-1} h\,[n,k]\, x\,[k]$$

$h\,[n,k]$ Is the strength of the connection between x[k] and y[n]

# Convolution

A linear shift invariant (LSI) function f can be written as a matrix multiplication:



$$h\left[n-k\right]$$ n indexes rows, k indexes columns

It can also be represented as a convolutional layer of neural net:



$$y\left[n\right] = \sum_{k=-1}^{1} h\left[k\right] x\left[n-k\right]$$

$$h\left[n-k\right]$$ Is the strength of the connection between x[k] and y[n]

**Toeplitz matrix**

$$\begin{pmatrix} a & b & c & d & e \\ f & a & b & c & d \\ g & f & a & b & c \\ h & g & f & a & b \\ i & h & g & f & a \end{pmatrix}$$

$\mathbf{y}$ $=$  $*$ $\mathbf{x}$

e.g., pixel image

- Constrained linear layer
- Fewer parameters —> easier to learn, less overfitting

$$\mathbf{y} = \quad * \quad \mathbf{x}$$

Conv layers can be applied to arbitrarily-sized inputs

# Five views on convolutional layers

1. Equivariant with translation    $f(\mathtt{translate}(x)) = \mathtt{translate}(f(x))$

2. Patch processing

3. Image filter

4. Parameter sharing

5. A way to process variable-sized tensors

# What if we have color?

(aka multiple input channels?)

# Multiple channel inputs

**Conv layer**



$$\mathbf{y} = \sum_c \mathbf{w}_c \circ \mathbf{x}_c$$

$$\mathbb{R}^{N \times C} \rightarrow \mathbb{R}^{N \times 1}$$

$\mathbf{x}$

$\mathbf{y}$

# Multiple channel *outputs*

**Conv layer**



$$\mathbf{y}_k = \sum_c \mathbf{w}_{k_c} \circ \mathbf{x}_c$$

$$\mathbb{R}^{N \times C} \to \mathbb{R}^{N \times K}$$

# Multiple channels

**Conv layer**



$$\mathbf{y}_k = \sum_c \mathbf{w}_{k_c} \circ \mathbf{x}_c$$

$$\mathbb{R}^{N \times C} \rightarrow \mathbb{R}^{N \times K}$$

Input features     A bank of 2 filters

2-dimensional output
**feature maps**

$F^1$

$\Sigma$

$F^2$

$\Sigma$

$$\mathbf{x}_l \in \mathbb{R}^{H \times W \times C_l} \quad \rightarrow \quad \mathbf{x}_{(l+1)} \in \mathbb{R}^{H \times W \times C_{(l+1)}}$$

[Figure modified from Andrea Vedaldi]

# Feature maps



conv1

relu1

conv2

relu2

- Each layer can be thought of as a set of C **feature maps** aka **channels**
- Each feature map is an NxM image

# Multiple channels: Example

$\mathbf{x}_l$

128

128

3

Filter Bank with
3x3 filters

$\mathbf{x}_{(l+1)}$

...

...

128

128

96

How many parameters does each *filter* have?

(a) 9      (b) 27      (c) 96      (d) 864

# Multiple channels: Example

$\mathbf{x}_l$

128

128

3

Filter Bank with
3x3 filters

$\mathbf{x}_{(l+1)}$

...

...

128

128

96

# How many filters are in the bank?

(a) 3        (b) 27        (c) 96        (d) can't say

# Filter sizes

When mapping from

$$\mathbf{x}_l \in \mathbb{R}^{H \times W \times C_l} \quad \rightarrow \quad \mathbf{x}_{(l+1)} \in \mathbb{R}^{H \times W \times C_{(l+1)}}$$
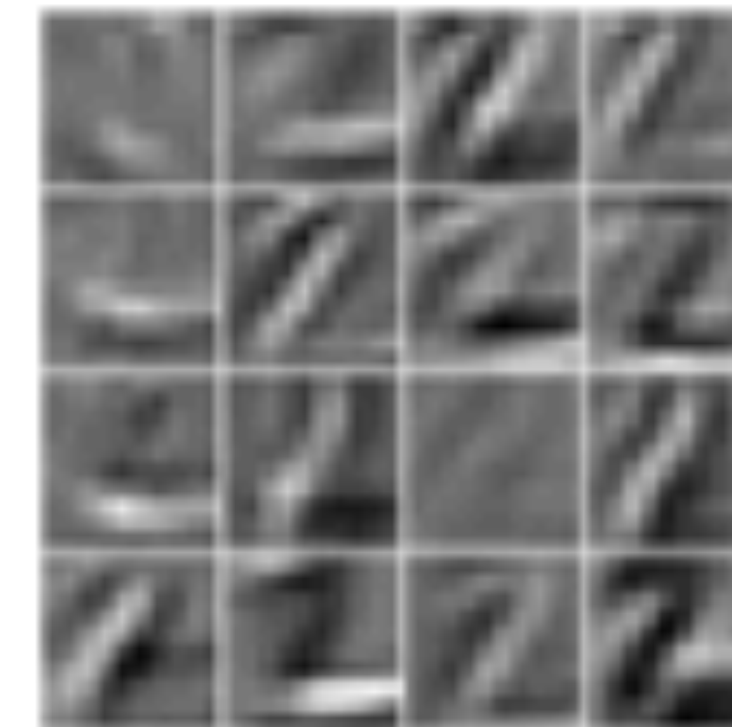
using an filter of spatial extent $\quad M \times N$

Number of parameters per filter: $\quad M \times N \times C_l$

Number of filters: $\quad C_{(l+1)}$

# Pooling and downsampling



We need translation and **scale** invariance

# Image pyramids

# Gaussian Pyramid



1/2

1/2

1/2

1/2

1/2

# Multiscale representations are great!



512   256   128   64   32   16   8

512   256   128   64   32   16   8

Gaussian Pyr

Laplacian Pyr

# How can we use multi-scale modeling in Convnets?

# Steerable Pyramid



Oriented filters

Blur

Downsampling

# Pooling



**Filter**

**Pool**

$\mathbf{x}$

$\mathbf{w}$

$b$

1

$\mathbf{z}$  $\mathbf{h}$

max

$\mathbf{y}$

**Max pooling**

$$y_j = \max_{j \in \mathcal{N}(j)} h_j$$

# Pooling



**Filter**  **Pool**

$\mathbf{x}$  $\mathbf{w}$  $b$  $1$  $\mathbf{z}$ $\mathbf{h}$  $\Sigma$  $\mathbf{y}$

**Max pooling**

$$y_j = \max_{j \in \mathcal{N}(j)} h_j$$

**Mean pooling**

$$y_j = \frac{1}{|\mathcal{N}|} \sum_{j \in \mathcal{N}(j)} h_j$$

# Pooling — Why?

Pooling across spatial locations achieves
stability w.r.t. small translations:

# Pooling — Why?

Pooling across spatial locations achieves
stability w.r.t. small translations:



large response
regardless of exact
position of edge

# Pooling — Why?

Pooling across spatial locations achieves
stability w.r.t. small translations:

# CNNs are stable w.r.t. diffeomorphisms



$$\approx$$

["Unreasonable effectiveness of Deep Features as a Perceptual Metric", Zhang et al. 2018]

# Pooling *across channels* — Why?

Pooling across feature channels (filter outputs)
can achieve other kinds of invariances:



large response
for any edge,
regardless of its
orientation

# Computation in a neural net

*Filter*

*ReLU*

*Pool*

*Classify*

"clown fish"

$$f(\mathbf{x}) = f_L(\ldots f_2(f_1(\mathbf{x})))$$

# Downsampling

**Filter**

**Pool and downsample**

$\mathbf{x}$

$\mathbf{w}$

$b$

1

$\mathbf{z}$ $g(\mathbf{z})$

# Downsampling

**Filter**                    **Downsample**



$\mathbf{x}$

$\mathbf{w}$

$b$

1

$\mathbf{z}$ $g(\mathbf{z})$

# Strided operations

**Conv layer**



**x**

**w**

Stride 2

$\mathbf{z}$   $g(\mathbf{z})$

**Strided operations** combine a given operation (convolution or pooling) and downsampling into a single operation.

# Computation in a neural net



$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

# Receptive fields

# Receptive fields



*Pool and downsample by 2*

*3x1 Filter*

*Pool and downsample by 2*

Sky

Sky

**Bird**

Sky

$RF = RF*2$

$RF = RF + floor(3/2)*2$

$RF = RF*2$

kernel size

downsample factor

# Effective Receptive Field

Contributing input units to a convolutional filter.

@jimmfleming // fomoro.com

**Input Features**

**7 // 2 Convolution**
Each filter sees 7 input units

strides continue...

**Convolutional Features**

**2 // 2 Max Pool**
Each filter sees 9 input units

**Max Pool Features**

**3 // 1 Convolution**
Each filter sees 17 input units

**Convolutional Features**

Features

Conv1D Filter

Padding or Stride

Receptive Field

[http://fomoro.com/tools/receptive-fields/index.html]

# Example: simple line classification



Class 0: vertical lines
Class 1: horizontal lines

# Example: simple line classification

| |
|---|
| *Filter* |
| *ReLu* |
| *Pool* |

$$\mathbf{z}_{1_i} = \mathbf{w}_i \circ \mathbf{x} + \mathbf{b}_i \qquad \text{with 2 learned kernels} \;\; \mathbf{w}_1, \mathbf{w}_2$$

$$\mathbf{h}_i = \max(\mathbf{z}_{1_i}, 0)$$

$$\mathbf{z}_{2_i} = \frac{1}{NM} \sum_{n,m} \mathbf{h}_i[n,m]$$

| |
|---|
| *fc* |

$$\mathbf{z}_3 = \mathbf{W}\mathbf{z}_2 + \mathbf{c}$$

| |
|---|
| *Softmax* |

$$\mathbf{y}_i = \frac{e^{-\tau \mathbf{z}_{3_i}}}{\sum_{k=1}^{K} e^{-\tau \mathbf{z}_{3_k}}}$$

parameters $\;\mathbf{w}_1, \mathbf{w}_2, \mathbf{b}_1, \mathbf{b}_2, \mathbf{W}, \mathbf{c}$

# Network training and evaluation

**1**  **0**  **0**  **0**  **1**  **1**

Training

**0**  **0**  **0**  **0**  **1**  **0**

Testing
3/10000 misclassified

**0**  **1**  **1**  **0**  **0**  **1**

# Network visualization

9x9 learned kernels:



fc layer learned weight:

$$\mathbf{W} = \begin{bmatrix} 2.83 & -2.36 \\ -0.60 & 1.14 \end{bmatrix}$$

# Out of domain generalization

Out of domain test samples are classified correctly
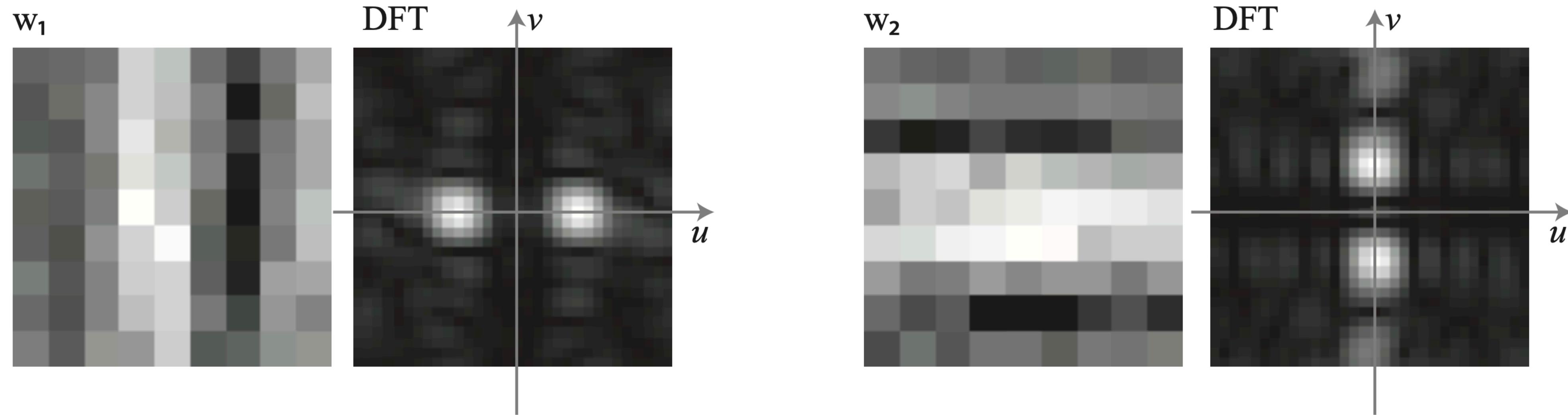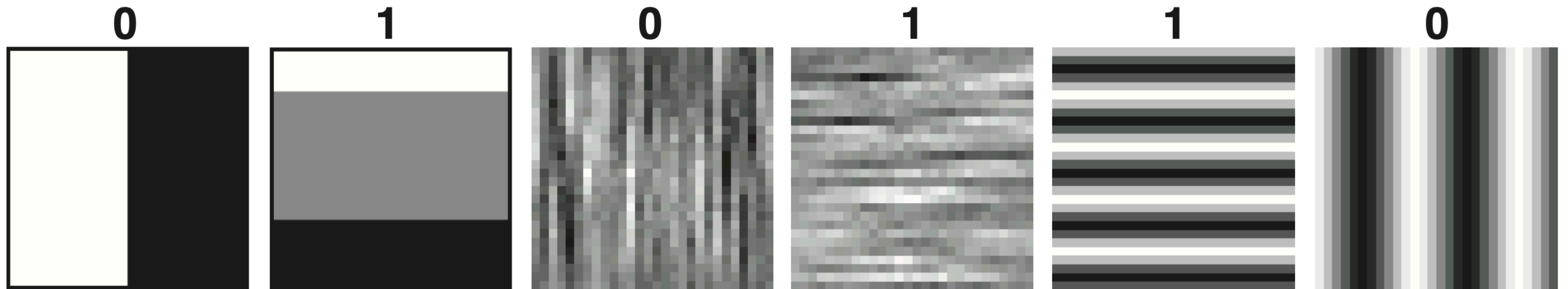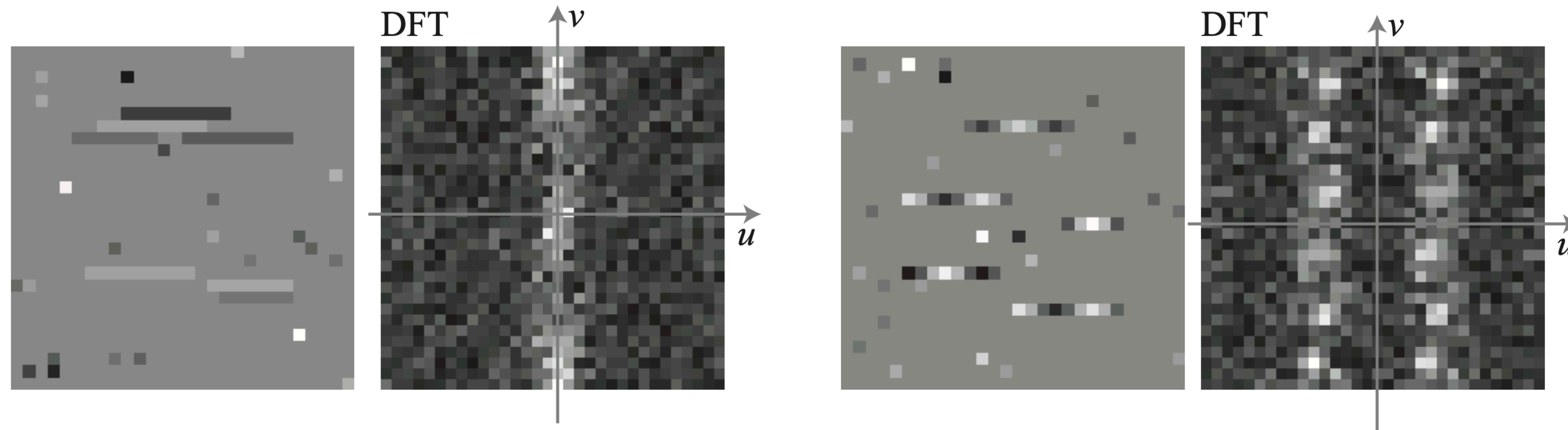
# Identifying vulnerability

- Modulation: multiplying image with sinusoidal wave moves spectral content horizontally



- All lines with sinusoidal texture are misclassified

# Analysis to run when training a large system
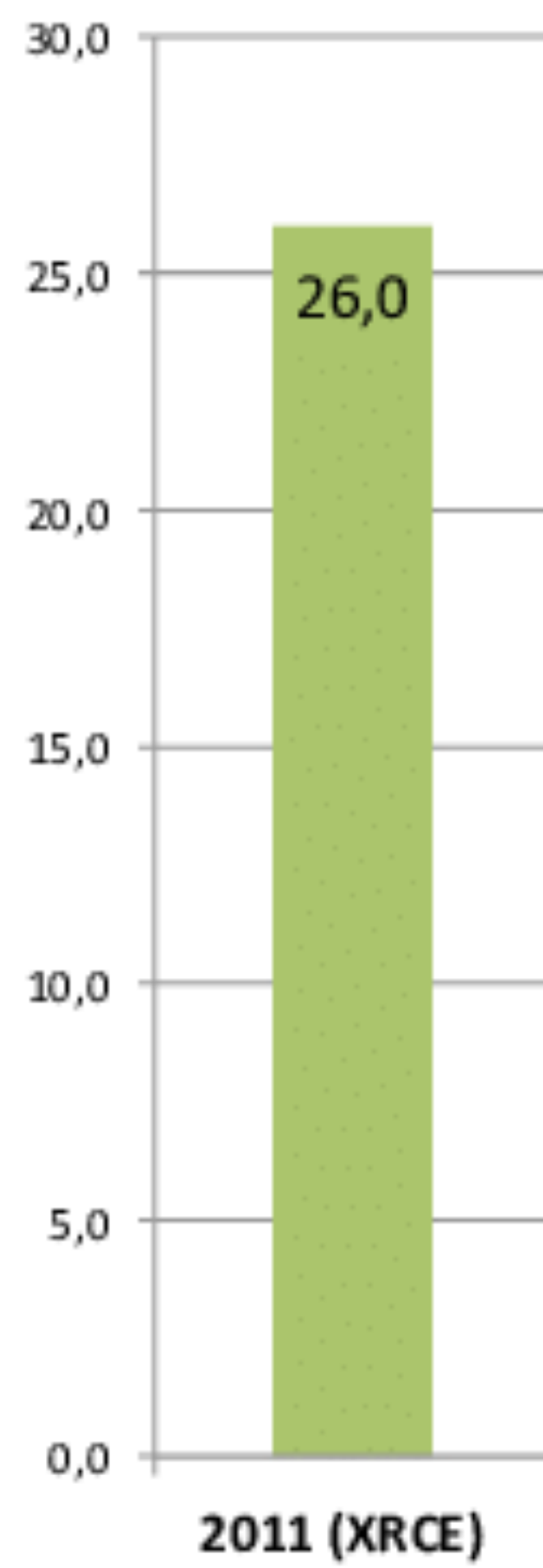
- Training and evaluation
- Visualize and understand the network
- Out of domain generalization
- Identifying vulnerability

# Some networks

… and what makes them work

# ImageNet Classification Error (Top 5)



26,0

2011 (XRCE)

ImageNet Classification Error (Top 5)

2012: AlexNet
5 conv. layers

11x11 conv, 96, /4, pool/2

5x5 conv, 256, pool/2

3x3 conv, 384

3x3 conv, 384

3x3 conv, 256, pool/2

fc, 4096

fc, 4096

fc, 1000

Error: 16.4%

*[Krizhevsky et al: ImageNet Classification with Deep Convolutional Neural Networks, NIPS 2012]*

# Alexnet — [Krizhevsky et al. NIPS 2012]

[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0
[27x27x96] MAX POOL1: 3x3 filters at stride 2
[27x27x96] NORM1: Normalization layer

11x11 conv, 96, /4, pool/2

5x5 conv, 256, pool/2

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2
[13x13x256] MAX POOL2: 3x3 filters at stride 2
[13x13x256] NORM2: Normalization layer

3x3 conv, 384

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

3x3 conv, 384

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

3x3 conv, 256, pool/2

[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1
[6x6x256] MAX POOL3: 3x3 filters at stride 2

fc, 4096

[4096] FC6: 4096 neurons

fc, 4096

[4096] FC7: 4096 neurons

fc, 1000

[1000] FC8: 1000 neurons (class scores)

11x11 conv, 96, /4, pool/2

5x5 conv, 256, pool/2

3x3 conv, 384

3x3 conv, 384

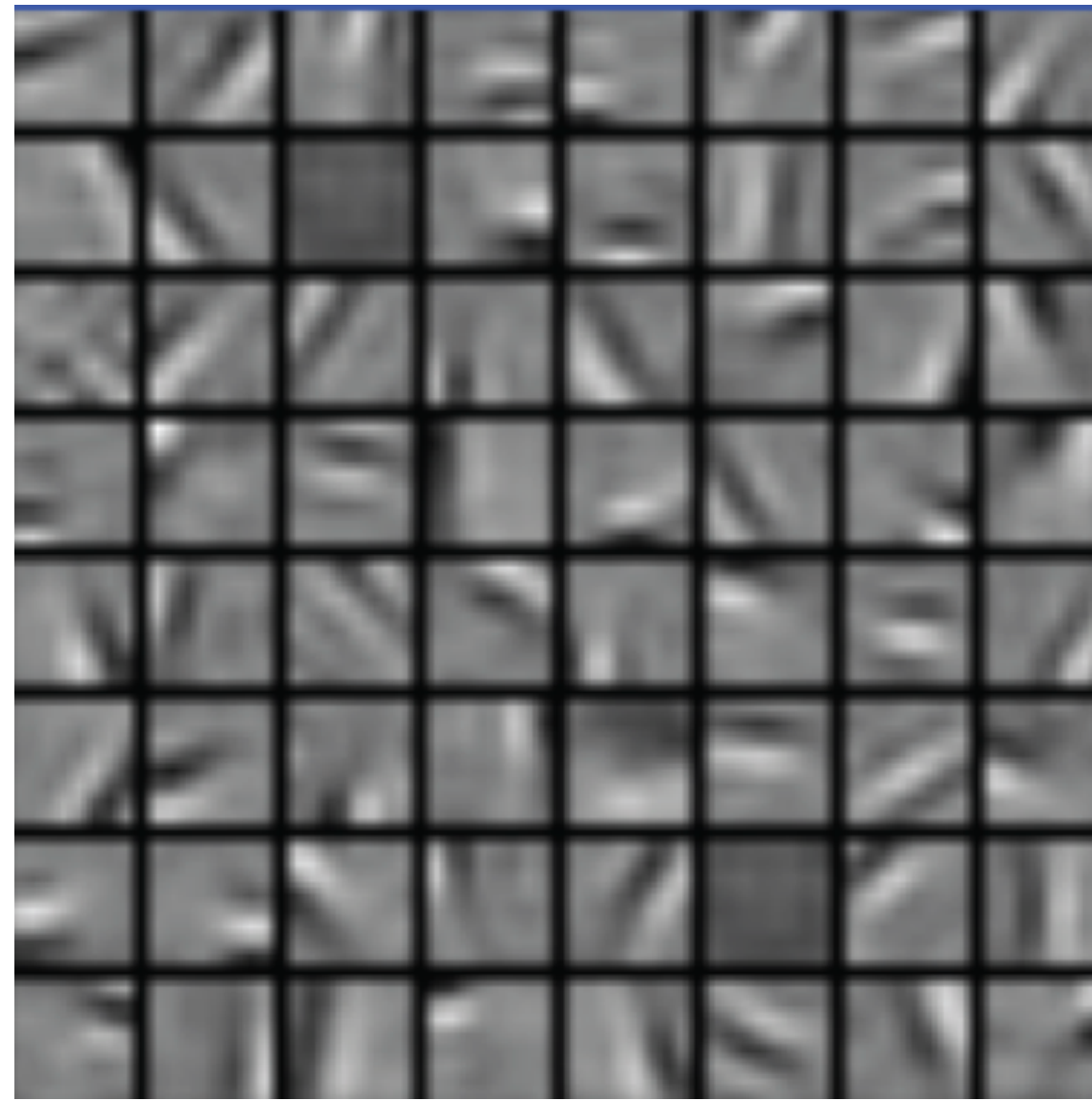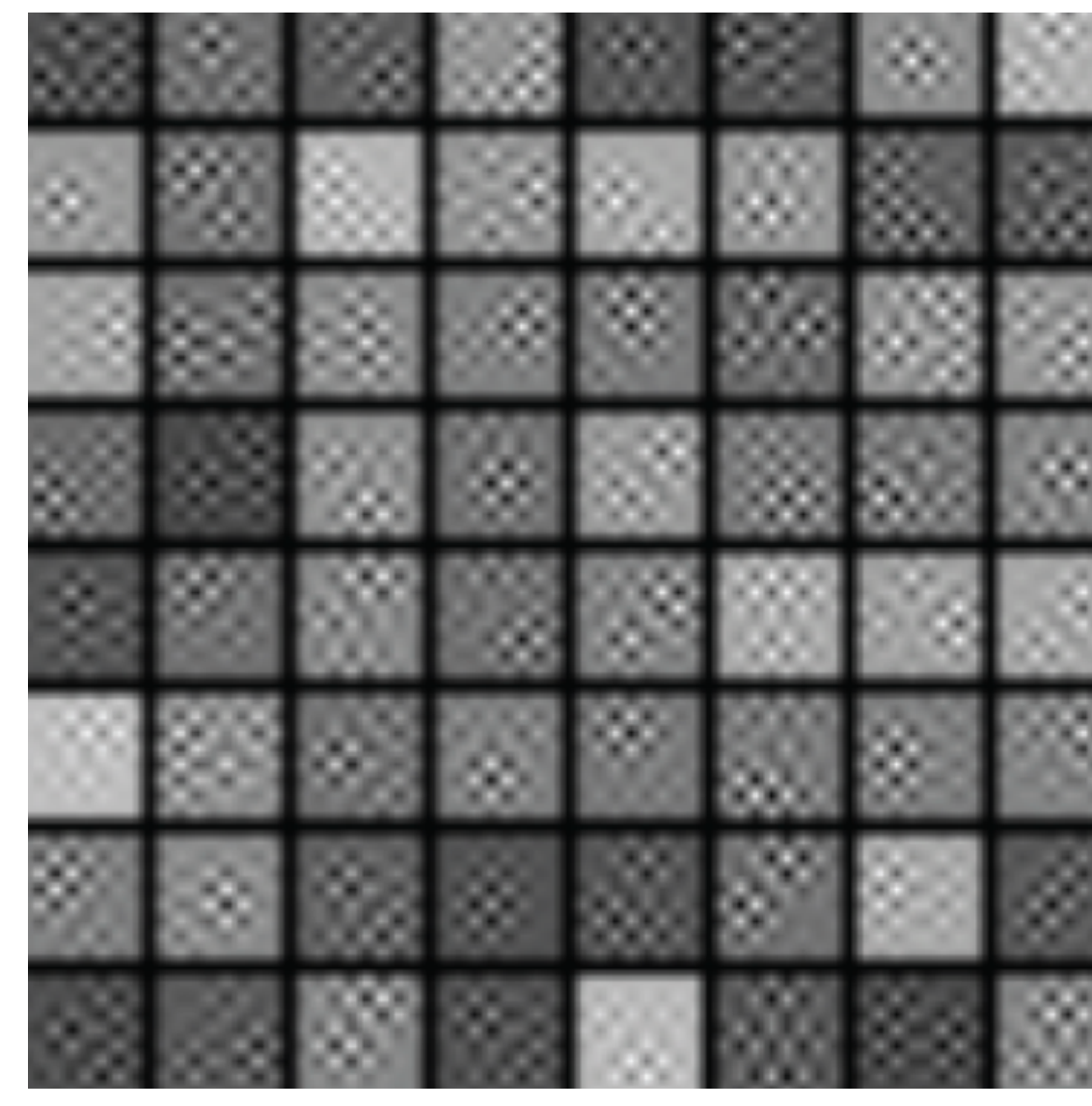3x3 conv, 256, pool/2

fc, 4096

fc, 4096

fc, 1000

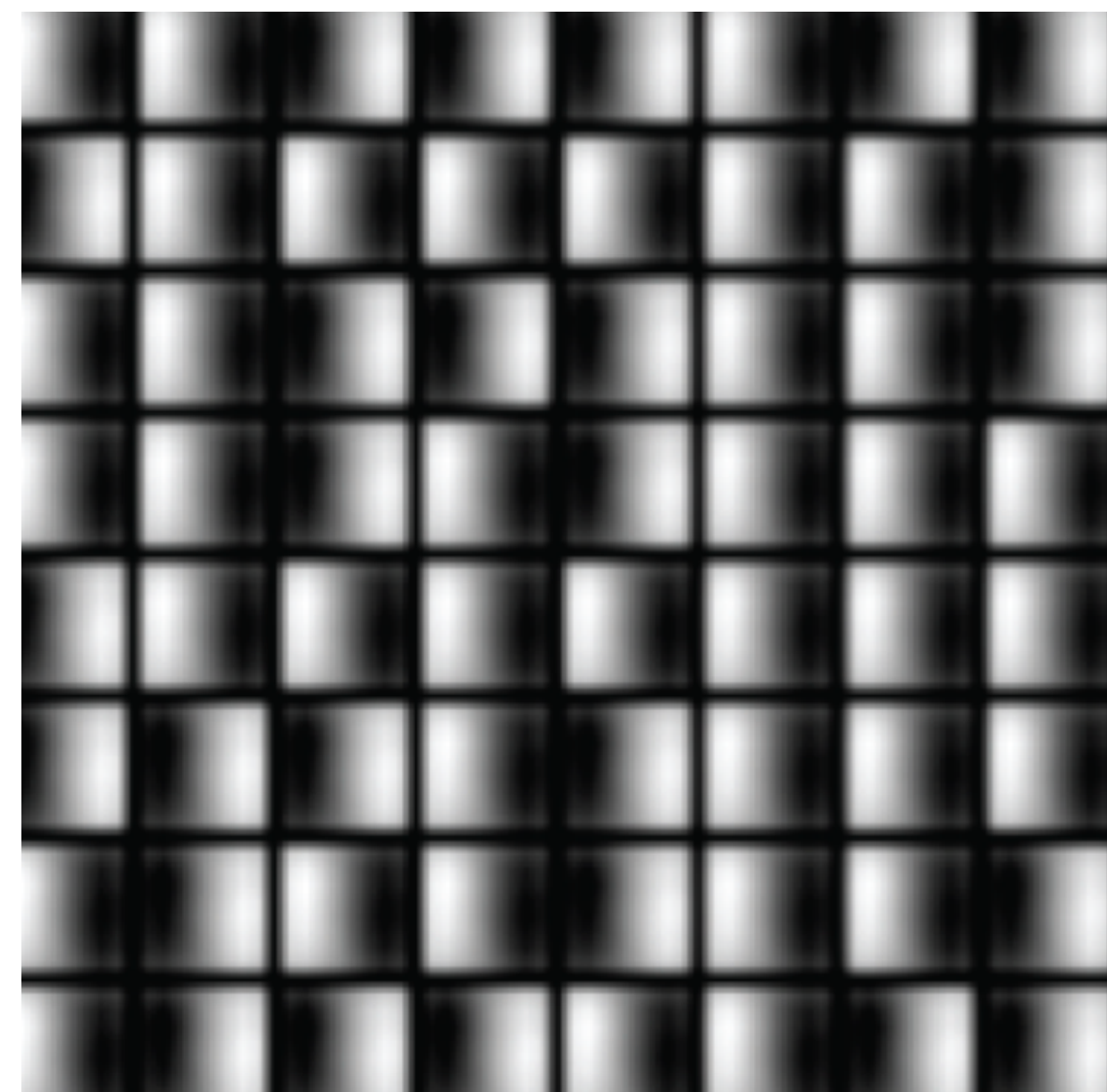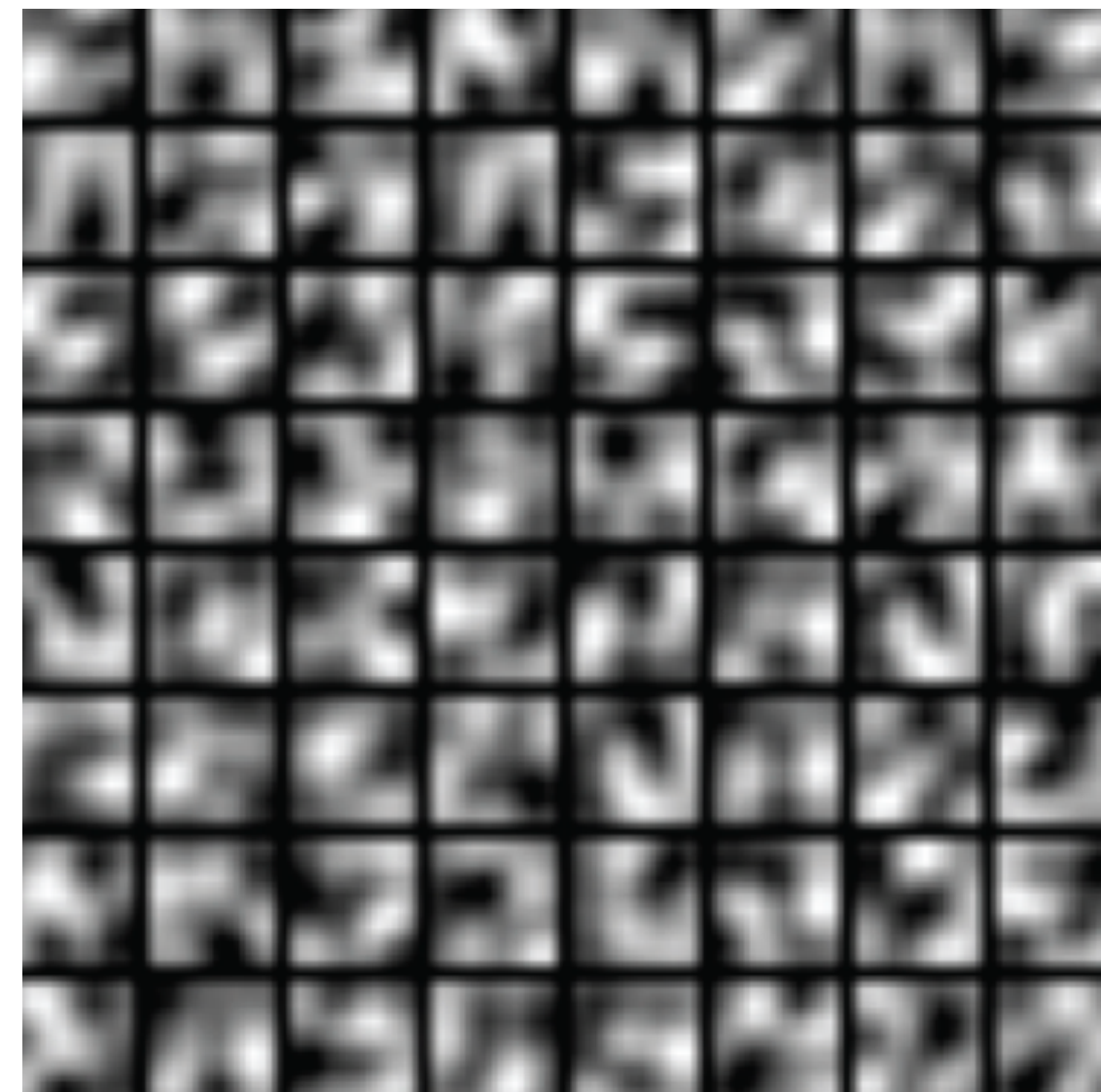# What filters are learned?

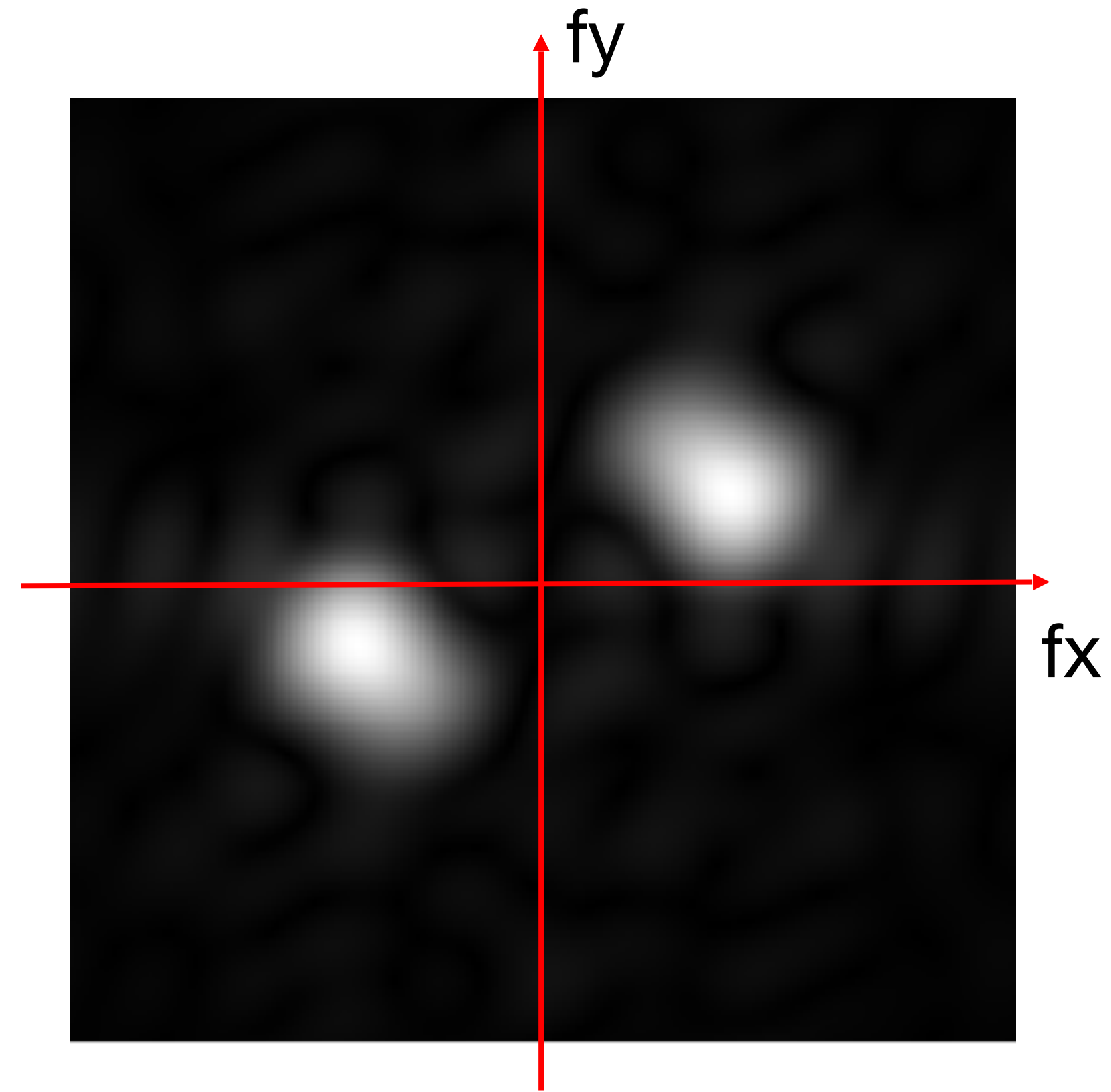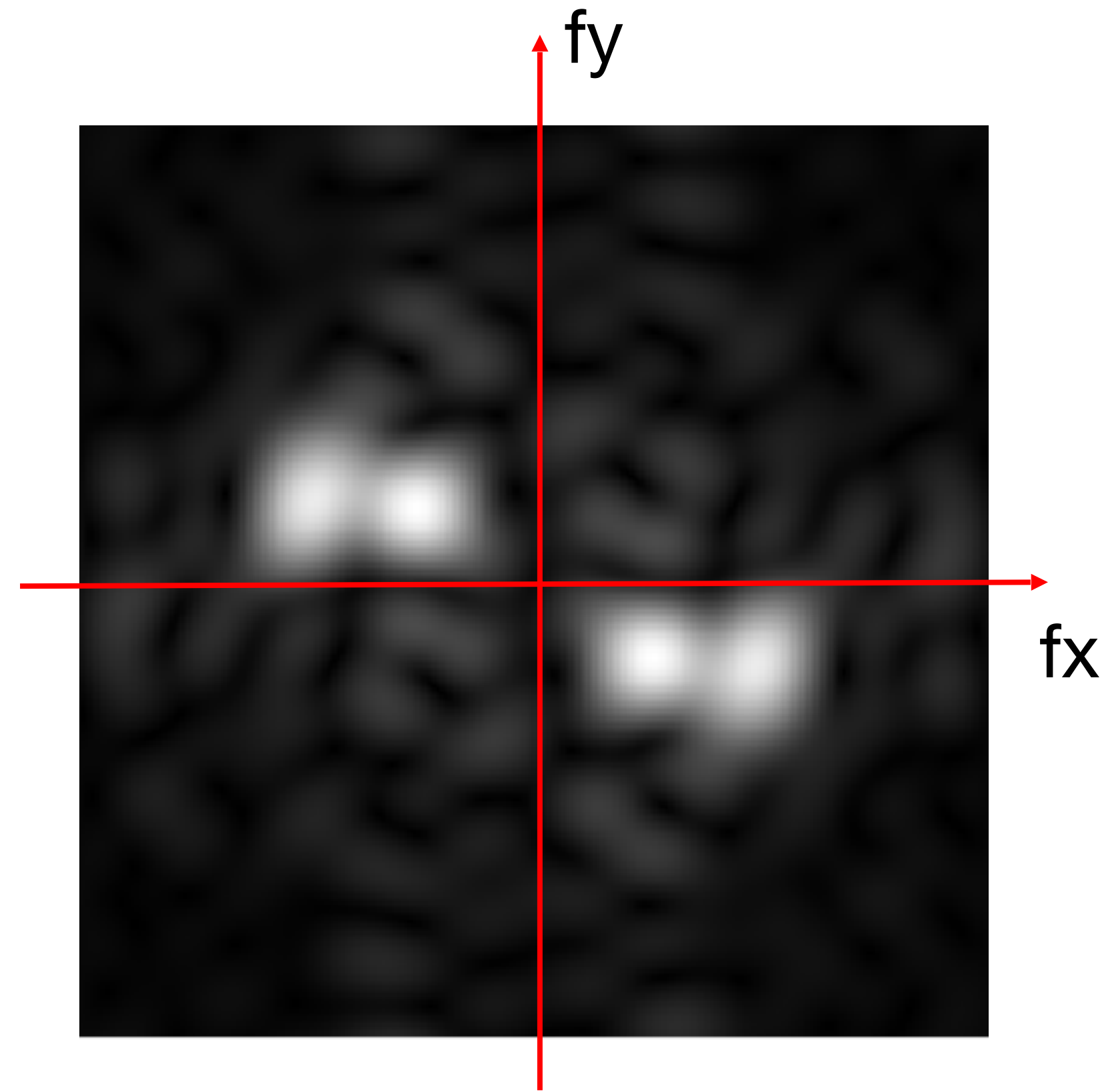# What filters are learned?



A

B

C

D

# Get to know your units



11x11 convolution kernel
(3 color channels)

# Get to know your units

# Get to know your units

# Get to know your units

# Get to know your units

# Get to know your units
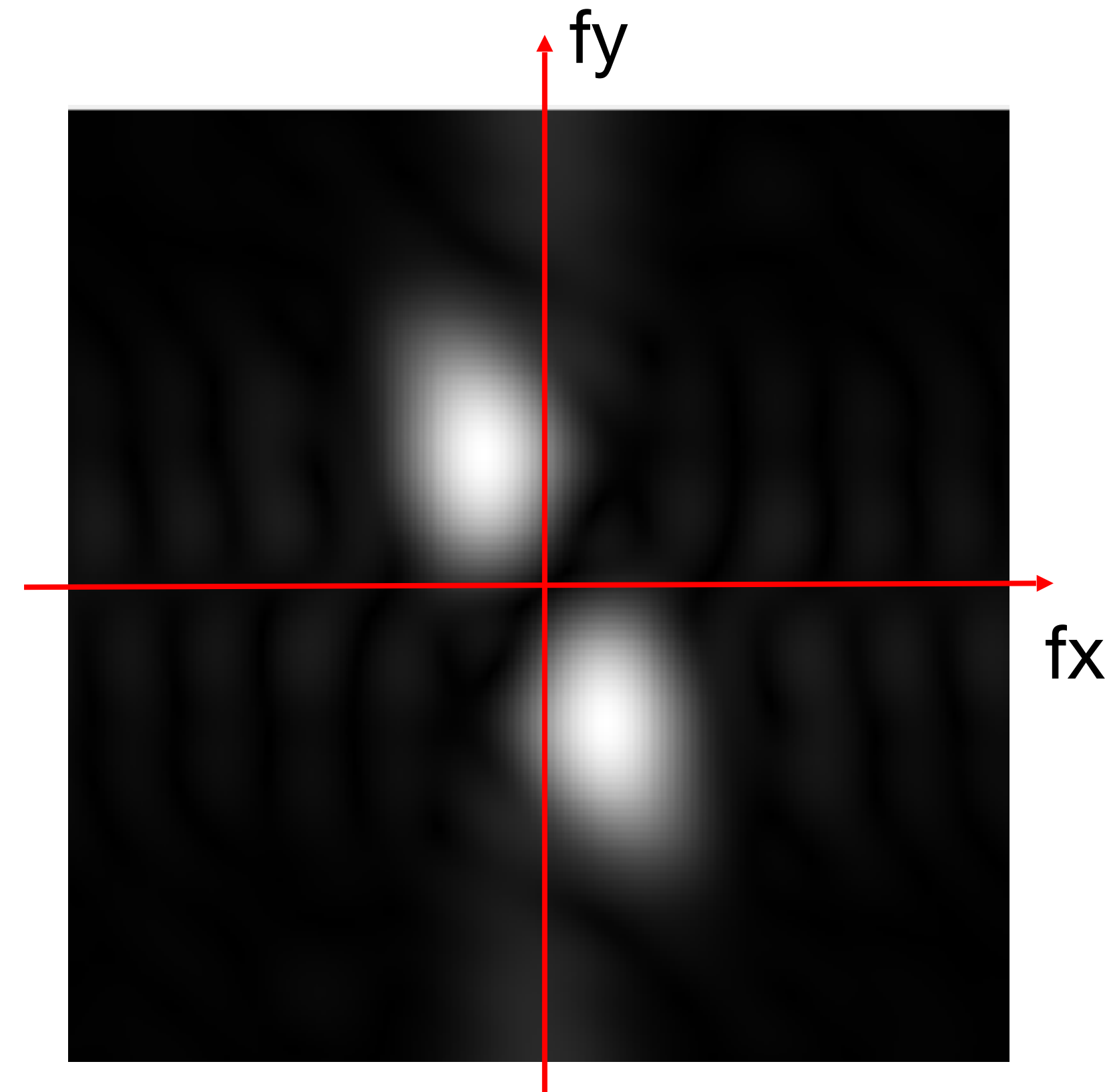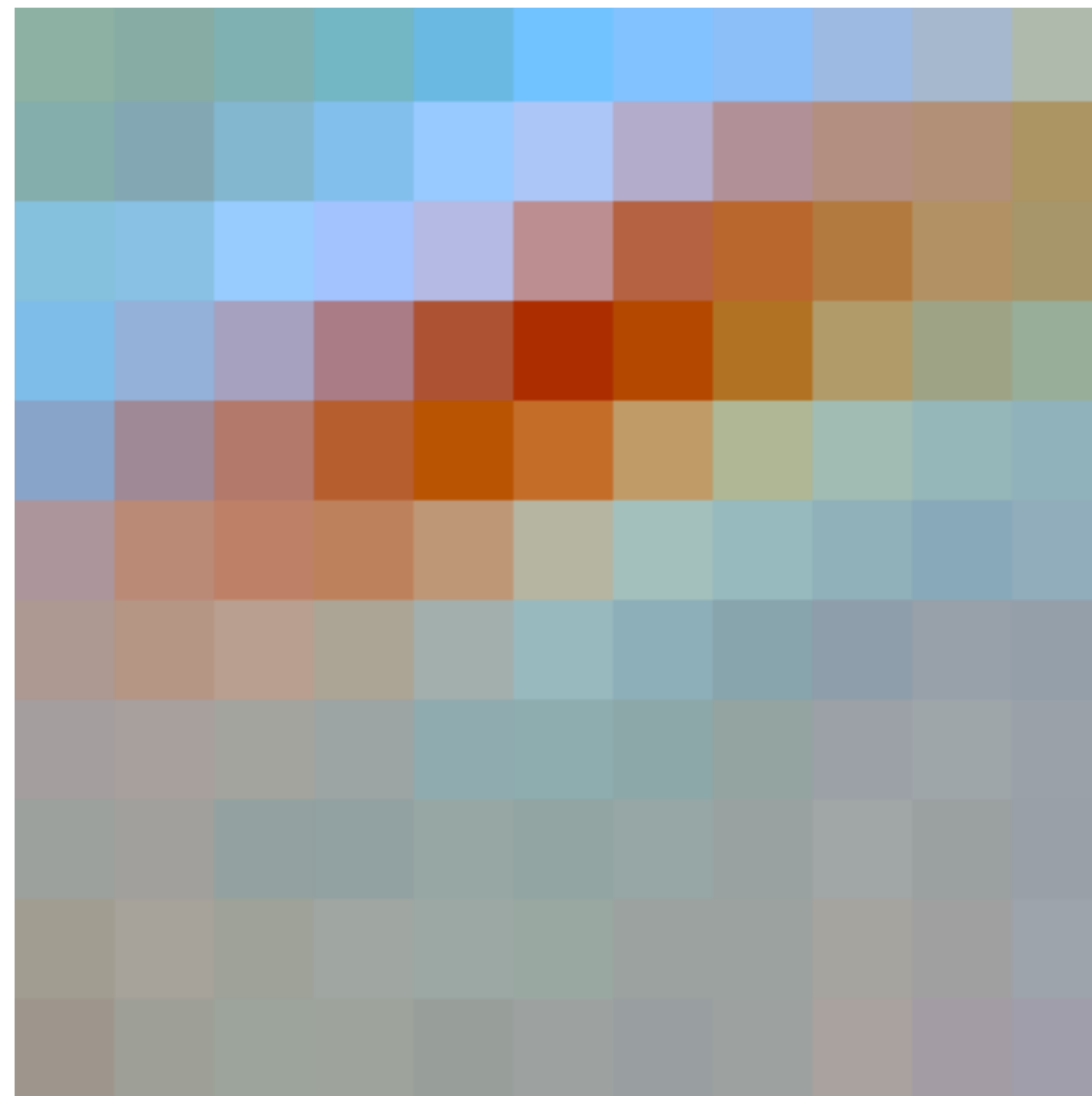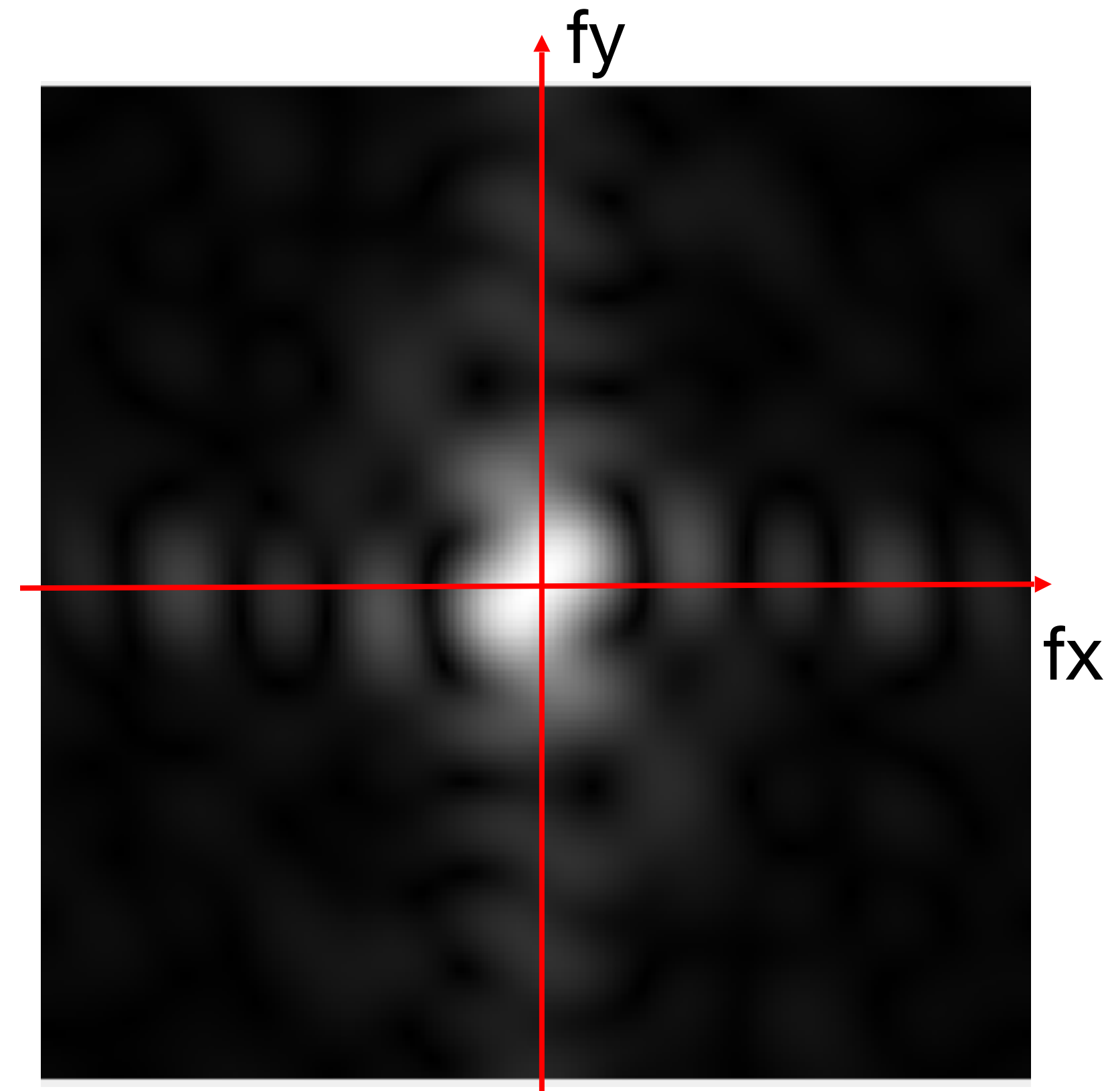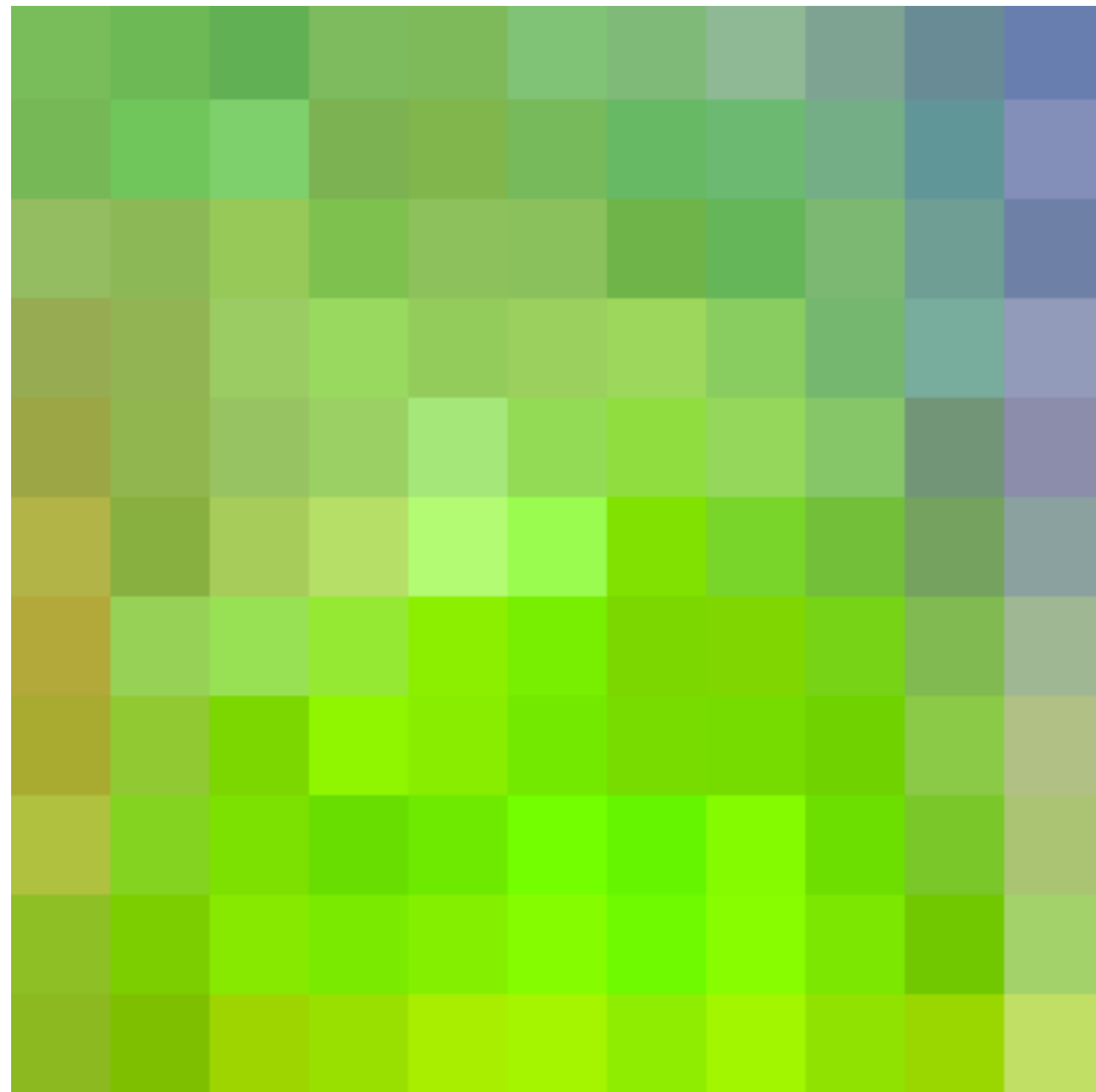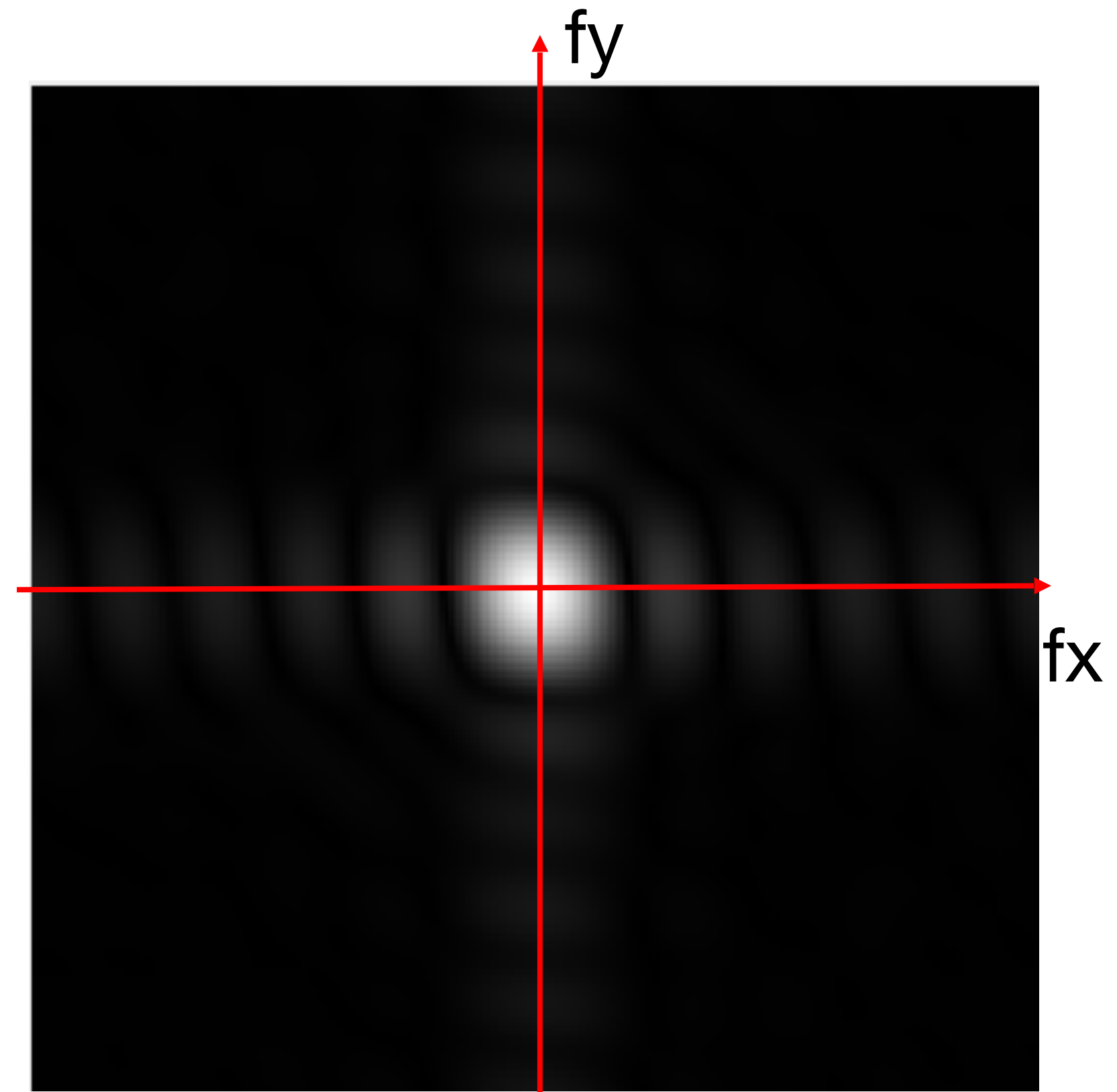
# Get to know your units
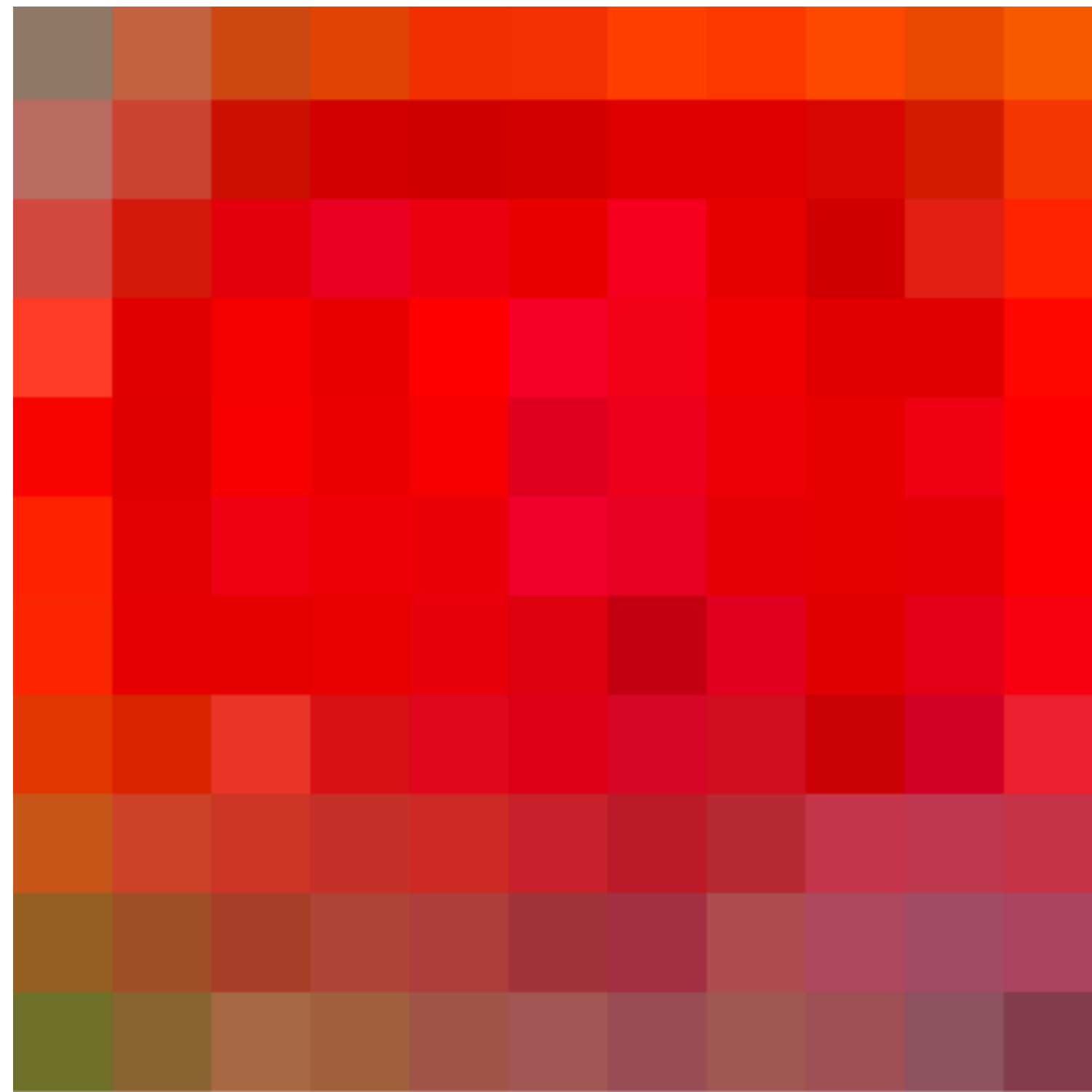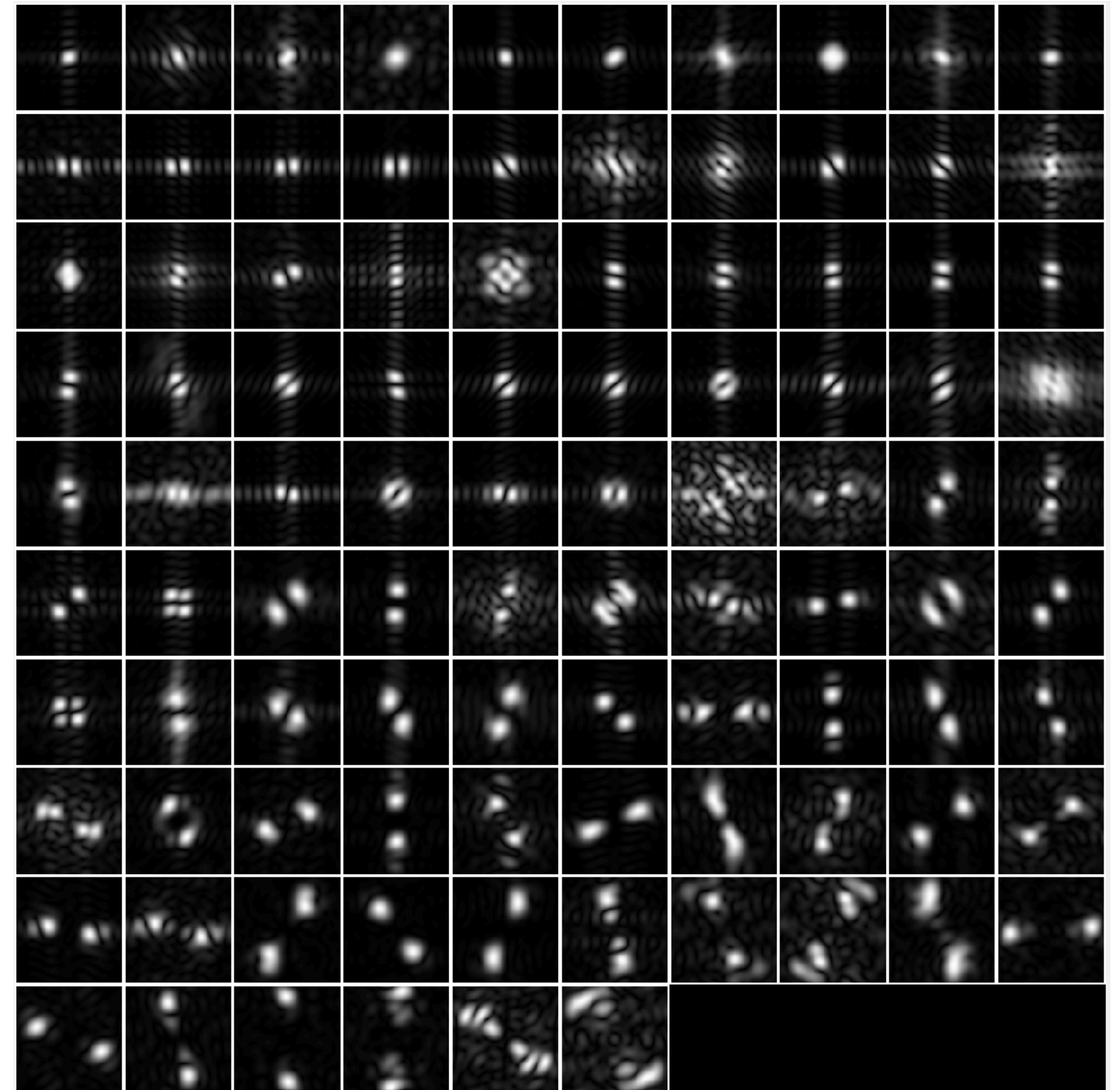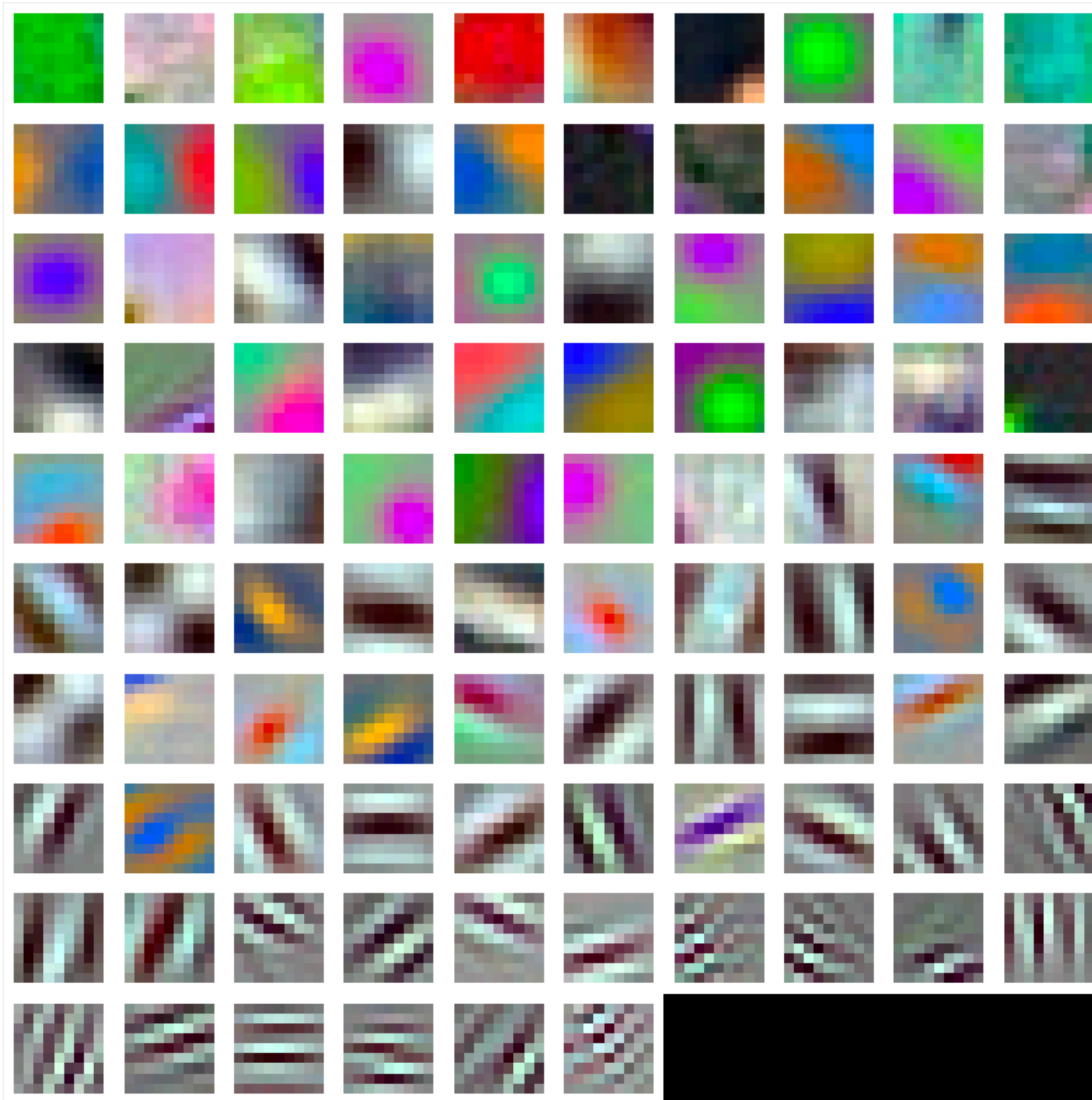
# Get to know your units
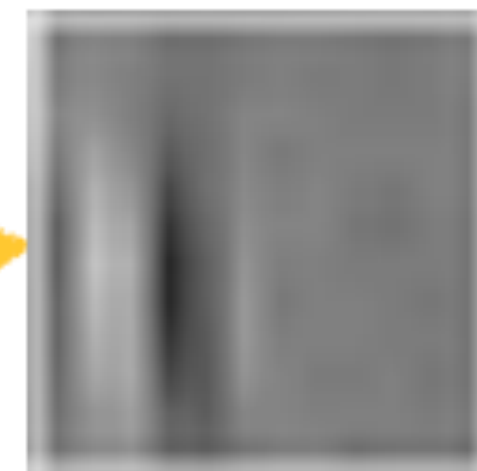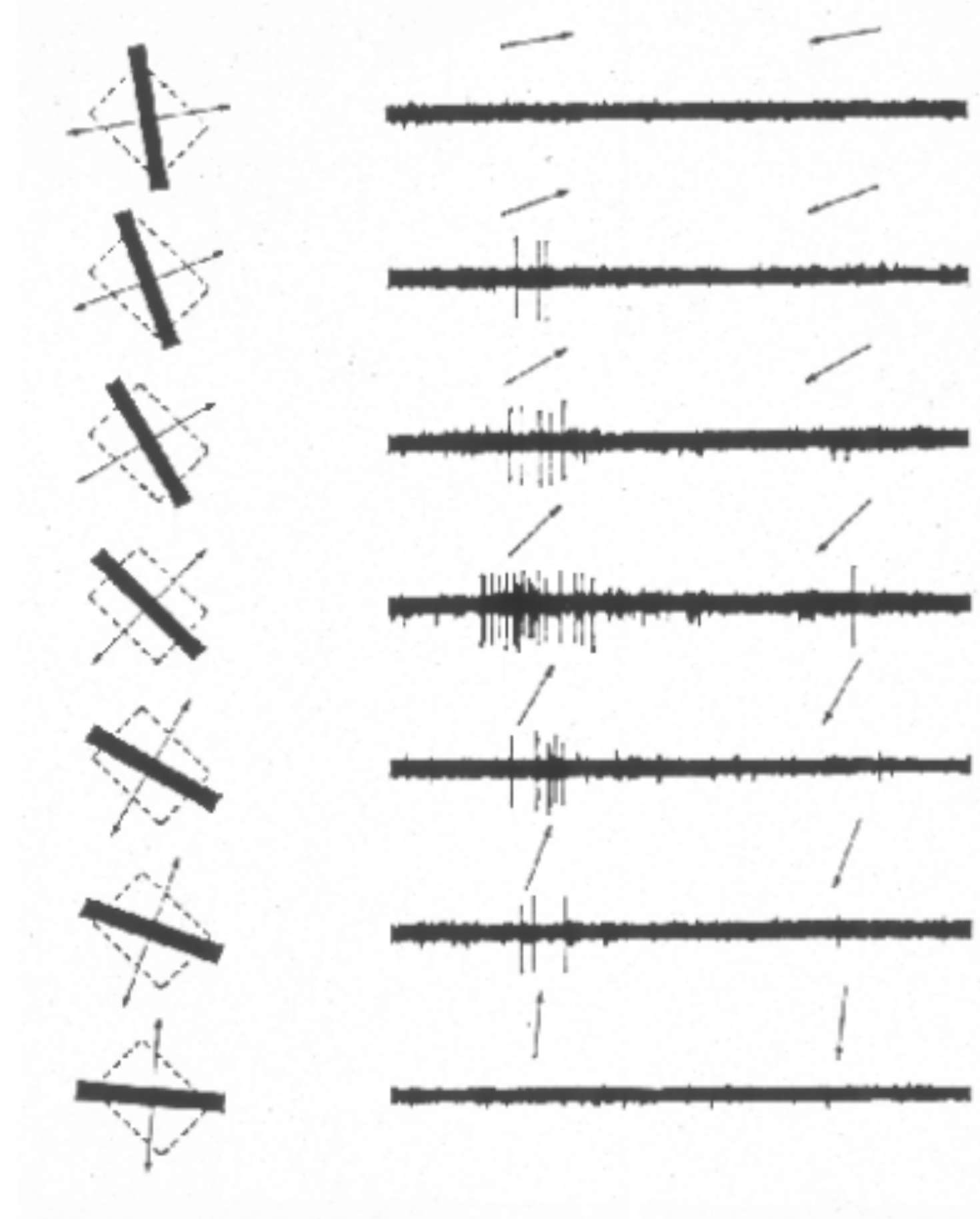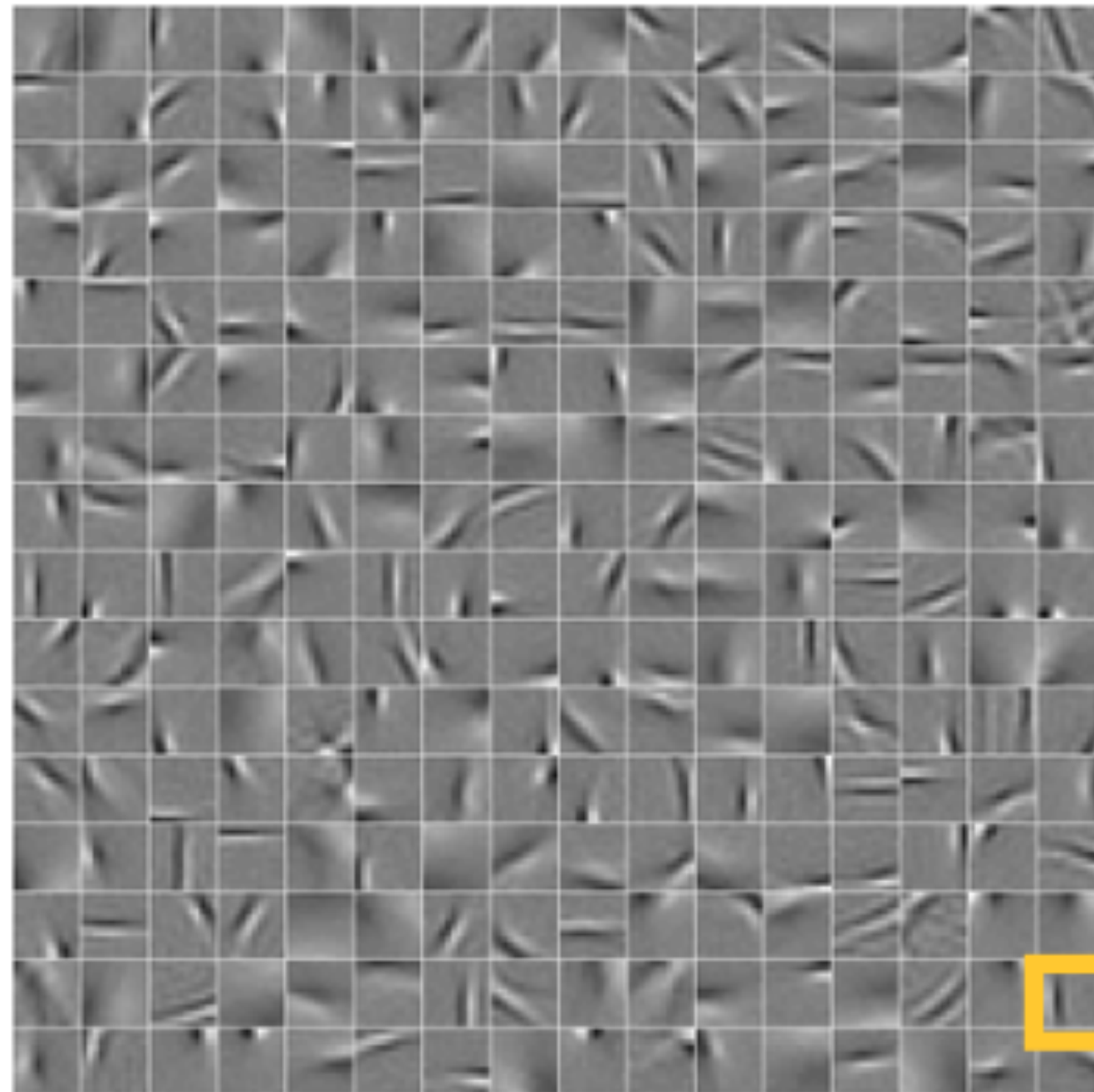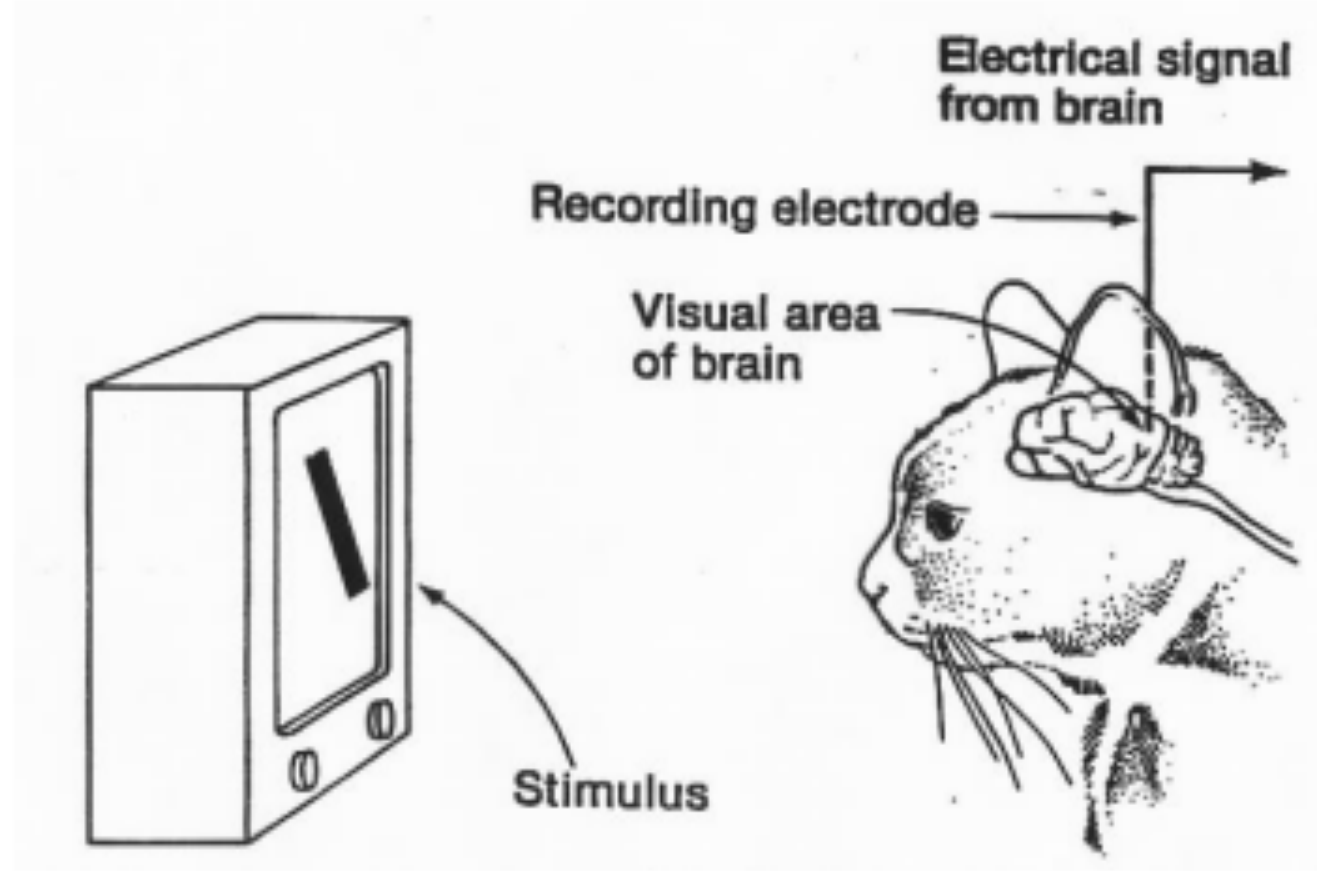


96 Units in conv1

# [Hubel and Wiesel 59]



oriented filter

[Slide from Andrea Vedaldi]

ImageNet Classification Error (Top 5)

## ImageNet Classification Error (Top 5)

2014: VGG
16 conv. layers

Error: 7.3%

*[Simonyan & Zisserman: Very Deep Convolutional Networks for Large-Scale Image Recognition, ICLR 2015]*

# VGG-Net [Simonyan & Zisserman, 2015]

2014: VGG
16 conv. layers

| |
|---|
| 3x3 conv, 64 |
| 3x3 conv, 64, pool/2 |
| 3x3 conv, 128 |
| 3x3 conv, 128, pool/2 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

Error: 7.3%

**Main developments**

- Small convolutional kernels: only 3x3

- Increased depth (5 -> 16/19 layers)

# Chaining convolutions

3x3 ○ 3x3 = 5x5

25 coefficients, but only
18 degrees of freedom

○ =

9 coefficients, but only
6 degrees of freedom.
Only separable filters... would this be enough?

ImageNet Classification Error (Top 5)

**ImageNet Classification Error (Top 5)**

2016: ResNet
>100 conv. layers

Error: 3.6%

*[He et al: Deep Residual Learning for Image Recognition, CVPR 2016]*

# ResNet [He et al, 2016]

2016: ResNet
>100 conv. layers

Error: 3.6%

**Main developments**

- Increased depth possible through residual blocks

# Residual Blocks

# Residual Blocks



$\mathcal{F}(x)$

weight layer

relu

weight layer

$x$

identity

$\mathcal{F}(x) + x$

relu

**Why do they work?**

- Gradients can propagate faster (via the identity mapping)

- Within each block, only small residuals have to be learned

# Make them bigger



ResNet
>100 conv. layers

GoogLeNet
22 conv. layers

VGG
16 conv. layers

AlexNet
5 conv. layers

**2012**          **2013**          **2014**          **2015**          **2016**

# Some debugging advice

# Other good things to know

- Check gradients numerically by finite differences
- Visualize hidden activations — should be uncorrelated and high variance



**Good training**: hidden units are sparse across samples and across features.

[Derived from slide by Marc'Aurelio Ranzato]

# Other good things to know

- Check gradients numerically by finite differences
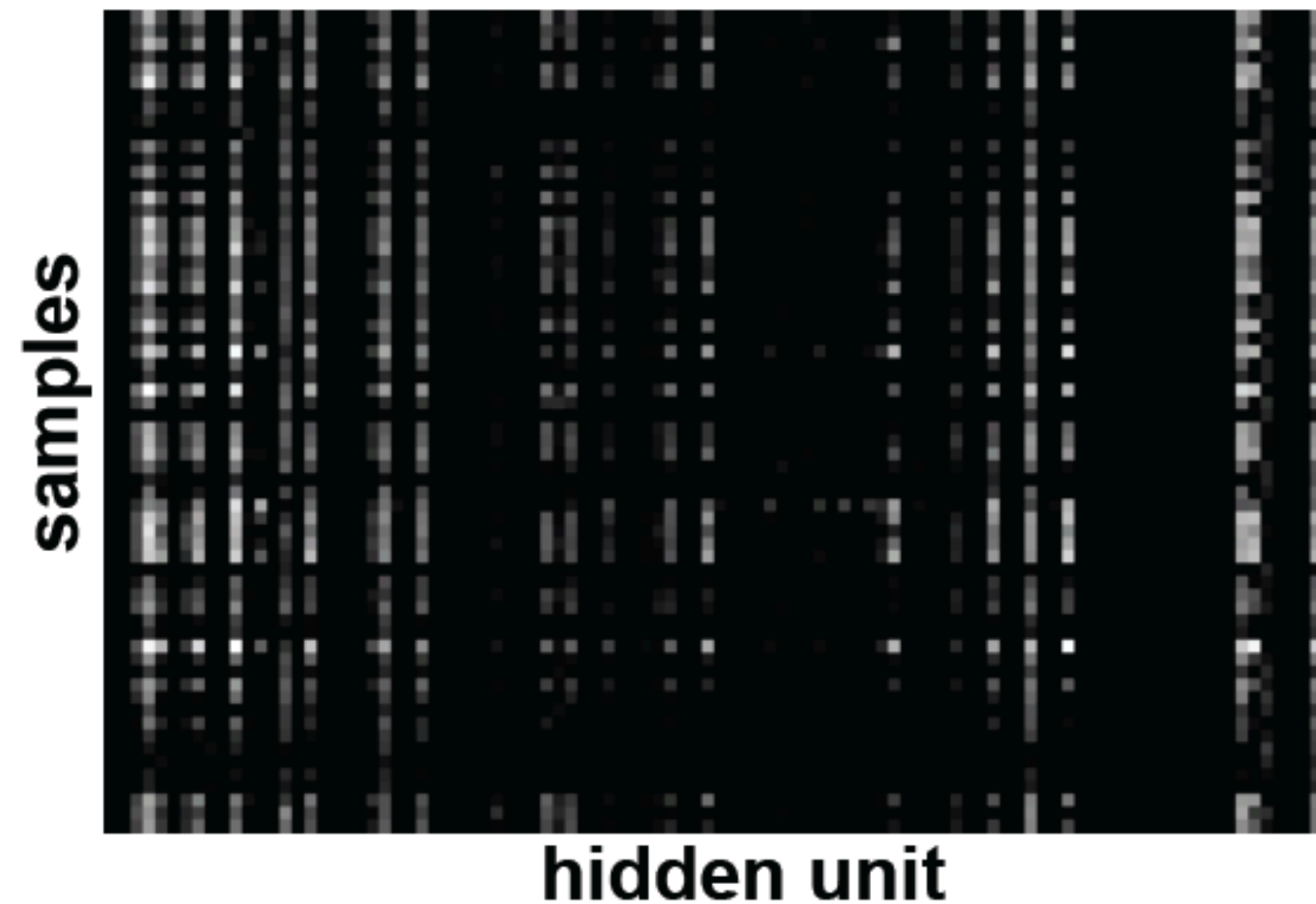- Visualize hidden activations — should be uncorrelated and high variance



**Bad training**: many hidden units ignore the input and/or exhibit strong correlations.

[Derived from slide by Marc'Aurelio Ranzato]

# Other good things to know

- Check gradients numerically by finite differences
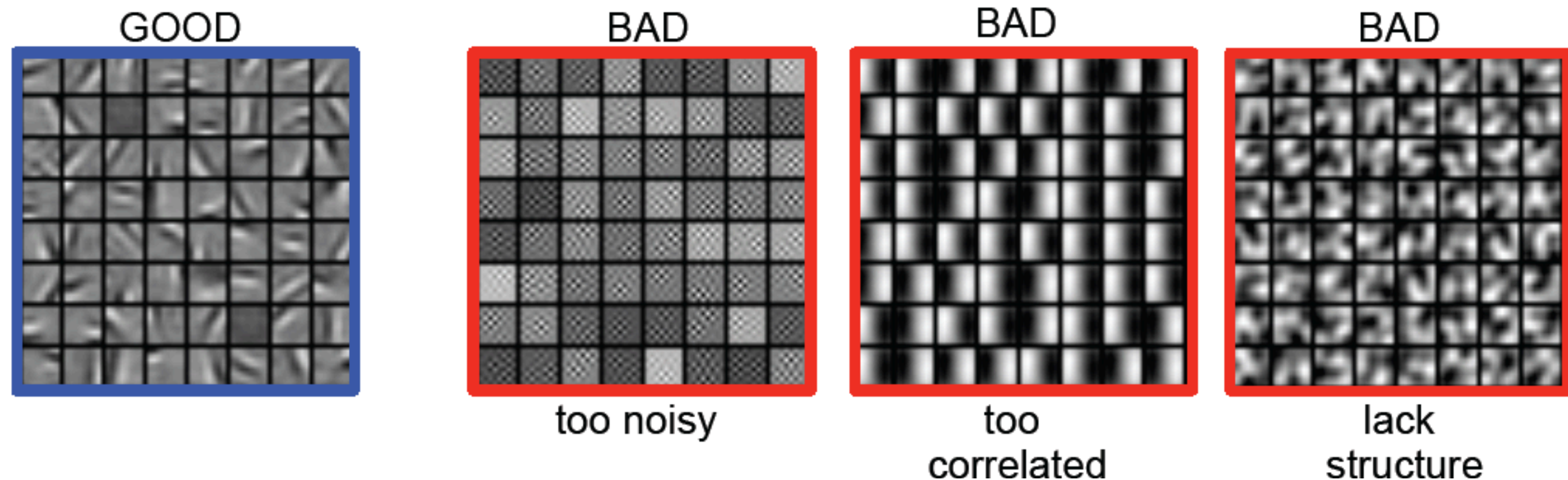- Visualize hidden activations — should be uncorrelated and high variance
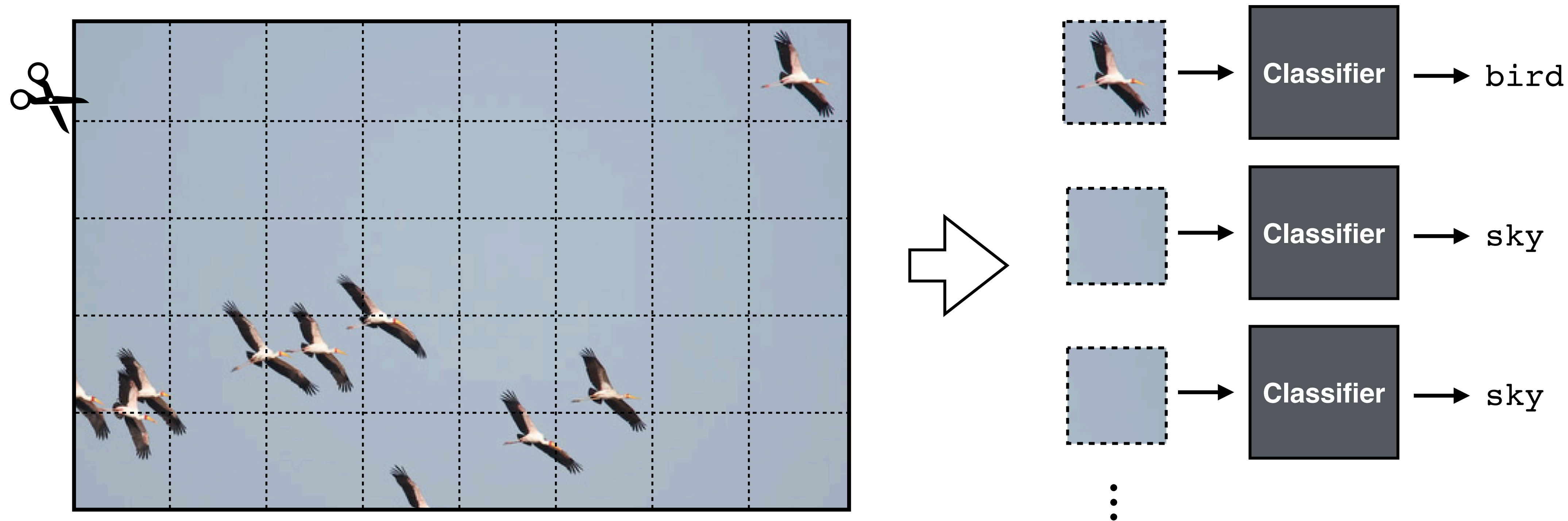- Visualize filters



**Good training**: learned filters exhibit structure and are uncorrelated.

[Derived from slide by Marc'Aurelio Ranzato]

# Transformers

Convnets in Disguise

Enduring principles:
1. Chop up signal into patches (divide and conquer)
2. Process each patch identically (and in parallel)

# 9. CNNs and Spatial Processing

- How to use deep nets for images

- New layer types: convolutional, pooling

- Feature maps and multichannel representations

- Popular architectures: Alexnet, VGG, Resnets

- Getting to know learned filters

- Unit visualization