

# **PART V: A VERY BRIEF INTRODUCTION TO DEEP UNSUPERVISED LEARNING**

\*elements taken from D. Kirkby lectures at KSPA

# TYPES OF MACHINE LEARNING

the machine is told what to look for

**SUPERVISED**

the machine is NOT told what to look for

**UN-SUPERVISED**

**TWO VERY BROAD TYPES OF MACHINE LEARNING ALGORITHMS**

# TYPES OF MACHINE LEARNING

the machine is told what to look for

**SUPERVISED**

LEARNS A MAP FROM  
X [FEATURES] TO Y  
[LABELS]

$$P(X|Y)$$

the machine is NOT told what to look for

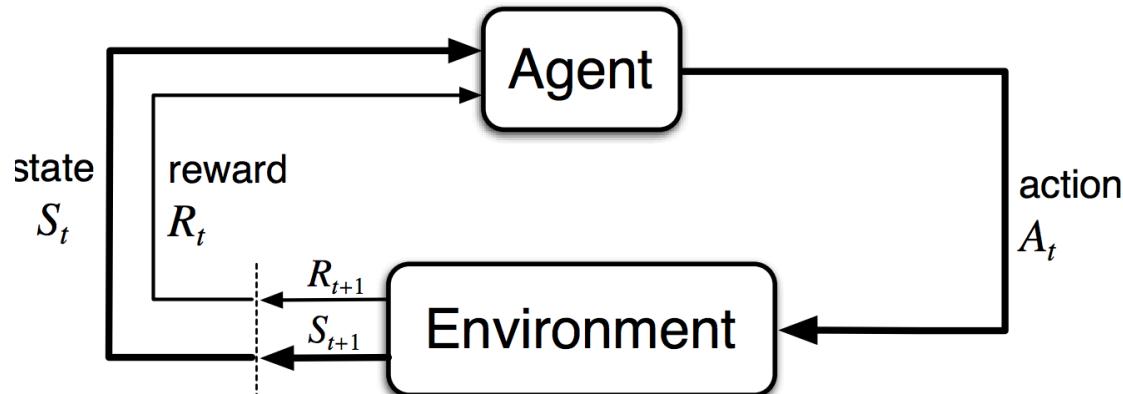
**UN-SUPERVISED**

NO LABELS - DISCOVER  
PATTERNS

$$P(X)$$

# ACTUALLY, THERE IS A THIRD AND FORTH TYPE THAT WE WILL NOT HAVE TIME TO COVER

## EINFORCMENT LEARNING:

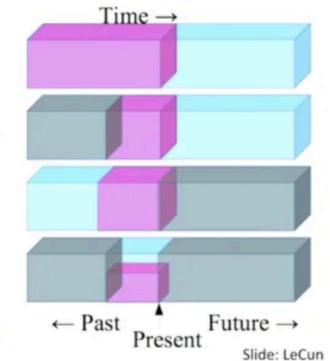


source: Sutton & Barto

TERATIVE LEARNING THROUGH TRIAL/ERROR;  
USED FOR ALPHAGO FOR EXAMPLE

## SELF-SUPERVISED LEARNING:

- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.
- ▶ Predict the **occluded** from the **visible**
- ▶ **Pretend there is a part of the input you don't know and predict that.**



That's also why more knowledge about the structure of the world can be learned through self-supervised learning than from the other two paradigms: the data is unlimited, and amount of feedback provided by each example is huge.

Y. LeCun

a self-supervised learning system attempts to predict parts of its inputs based on the other parts of its inputs

# IT IS AN ACTIVE FIELD OF RESEARCH BECAUSE LABELLING IS ONE OF THE MAIN BOTTLENECKS OF MACHINE LEARNING

- ▶ “Pure” Reinforcement Learning (**cherry**)
  - ▶ The machine predicts a scalar reward given once in a while.
  - ▶ **A few bits for some samples**
- ▶ Supervised Learning (**icing**)
  - ▶ The machine predicts a category or a few numbers for each input
  - ▶ Predicting human-supplied data
  - ▶ **10→10,000 bits per sample**
- ▶ Self-Supervised Learning (**cake génoise**)
  - ▶ The machine predicts any part of its input for any observed part.
  - ▶ Predicts future frames in videos
  - ▶ **Millions of bits per sample**



IN THIS LAST PART WE ARE GOING TO BRIEFLY INTRODUCE CURRENT TECHNIQUES OF UNSUPERVISED LEARNING WITH NEURAL NETWORKS

THERE ARE 3 MAJOR APPLICATIONS TO ASTRONOMY (THAT I CAN THINK OF):

- **DIMENSIONALITY REDUCTION:** HOW CAN I REPRESENT MY COMPLEX DATA MORE EFFICIENTLY TO GET NEW INSIGHTS INTO ITS STRUCTURE?
- **GENERATIVE MODELING:** HOW CAN I INTERPOLATE / EXTRAPOLATE A (SMALL, SPARSE) DATASET TO GENERATE NEW DATA SAMPLED FROM THE SAME (UNKNOWN) DISTRIBUTION?
- **PROBABILISTIC MODELING:** WHAT IS THE PROBABILITY THAT A NEW OBSERVATION IS DRAWN FROM THE SAME (UNKNOWN) DISTRIBUTION AS SOME REFERENCE (SMALL, SPARSE) DATASET?

# DIMENSIONALITY REDUCTION

NEED OF REPRESENTING THE DATA IN A SMALLER  
DIMENSIONALITY SPACE

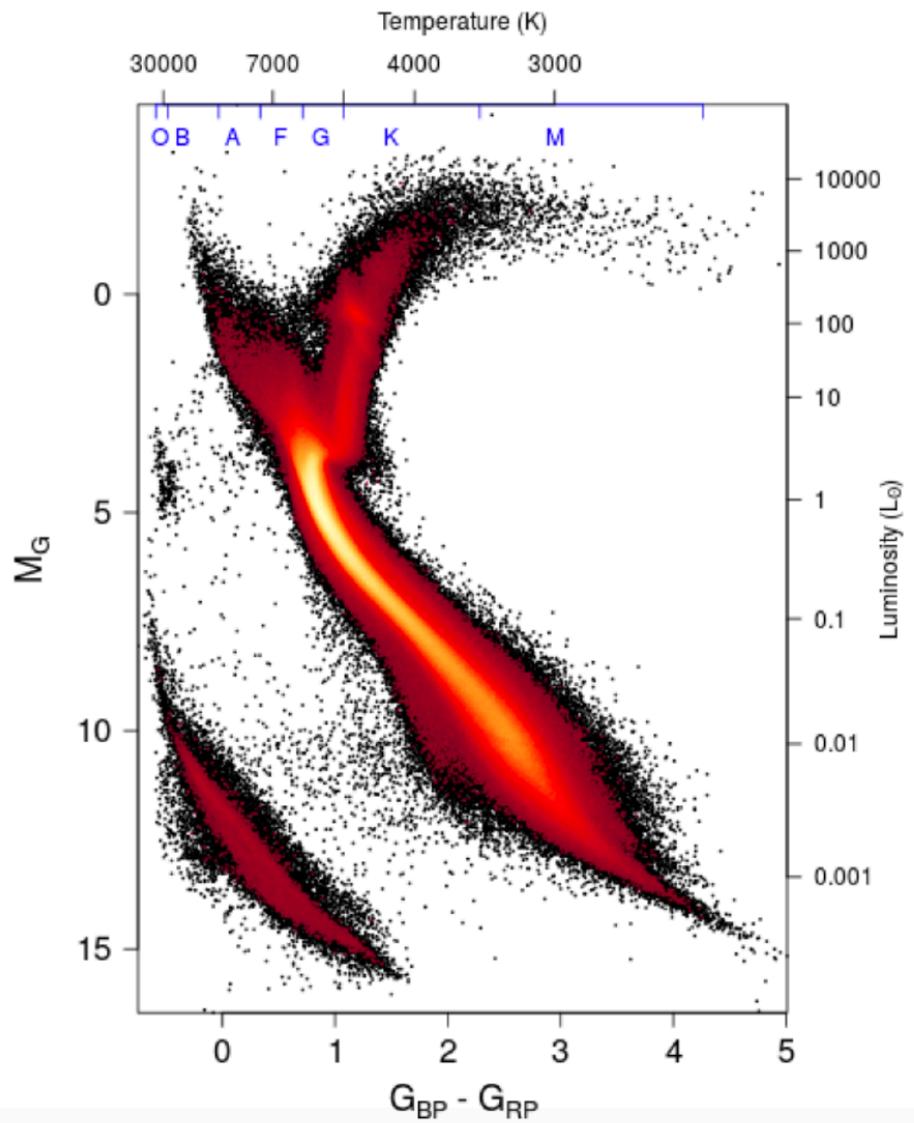
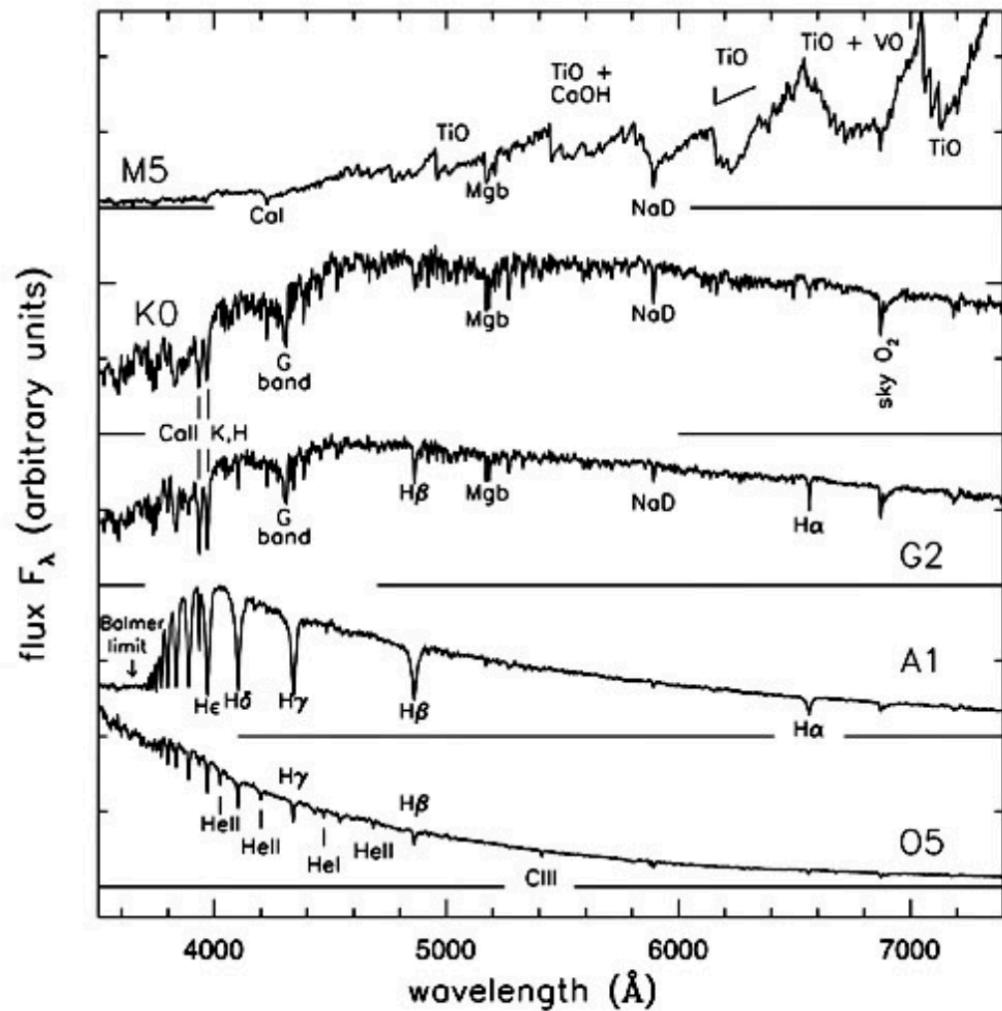
$$G_{\mu\nu} = \frac{8\pi G}{c^4} T_{\mu\nu}$$



@D. BARON

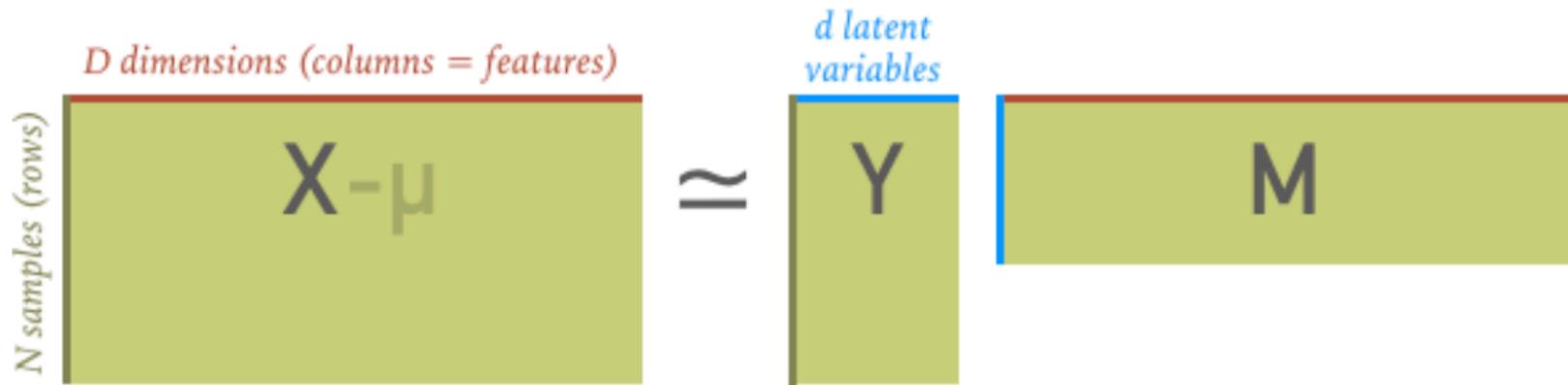
MOST OF THE SCIENTIFIC WORK IS ABOUT  
COMPRESSION OF NATURE IN SOME  
EXPLANATORY EQUATIONS

From: Gaia Collaboration et al. 2018



@D. BARON

**THE STELLAR MAIN SEQUENCE IS A  
DIMENSIONALITY REDUCTION**

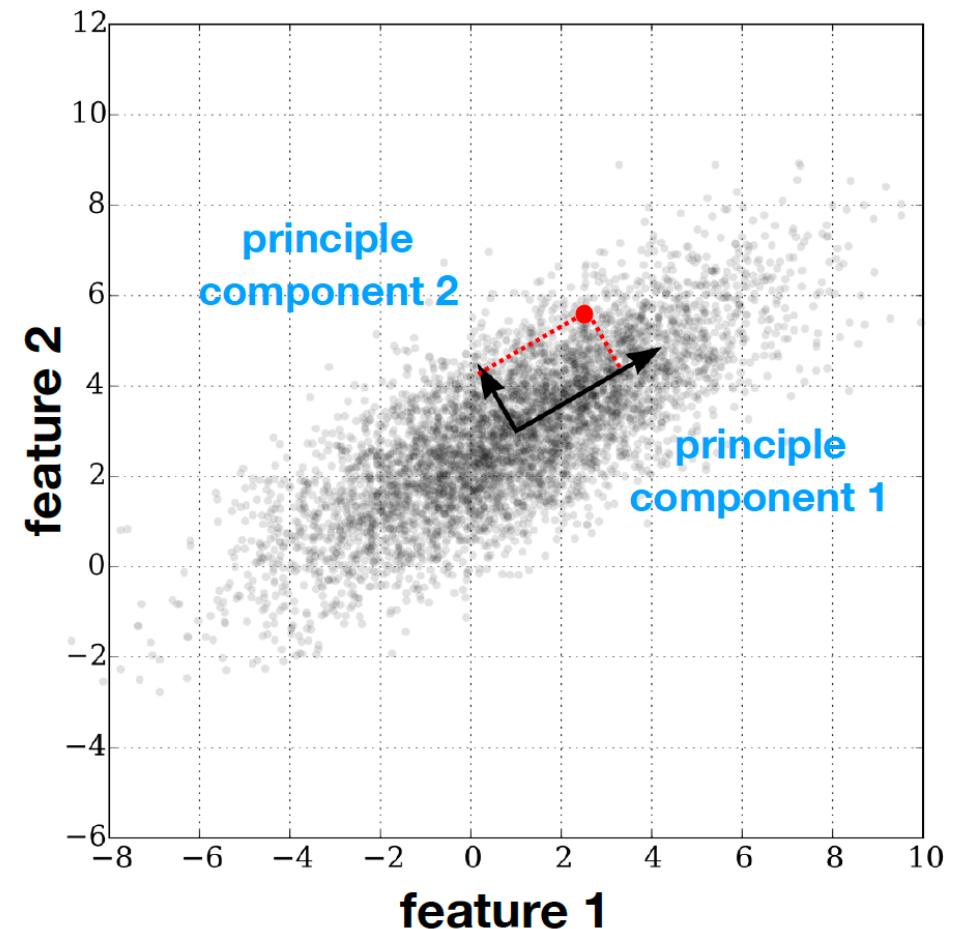


CLASSICAL METHODS FOR DIMENSIONALITY REDUCTION  
SEEK A LINEAR DECOMPOSITION THAT BEST EXPLAINS  
THE OBSERVATIONS  $X$  IN TERMS OF **LATENT VARIABLES  $Y$**

# PRINCIPAL COMPONENT ANALYSIS

PCA CONVERT A SET OF  
(CORRELATED) VARIABLES INTO A  
SET  
OF VALUES LINEARLY  
UNCORRELATED

1. THE FIRST PRINCIPLE COMPONENT (“PROTOTYPE”), HAS THE LARGEST POSSIBLE VARIANCE
2. THE FOLLOWING COMPONENTS HAVE THE LARGEST VARIANCES WITH THE ADDITIONAL CONSTRAINT THAT THEY ARE ORTHOGONAL TO THE PRECEDING COMPONENTS

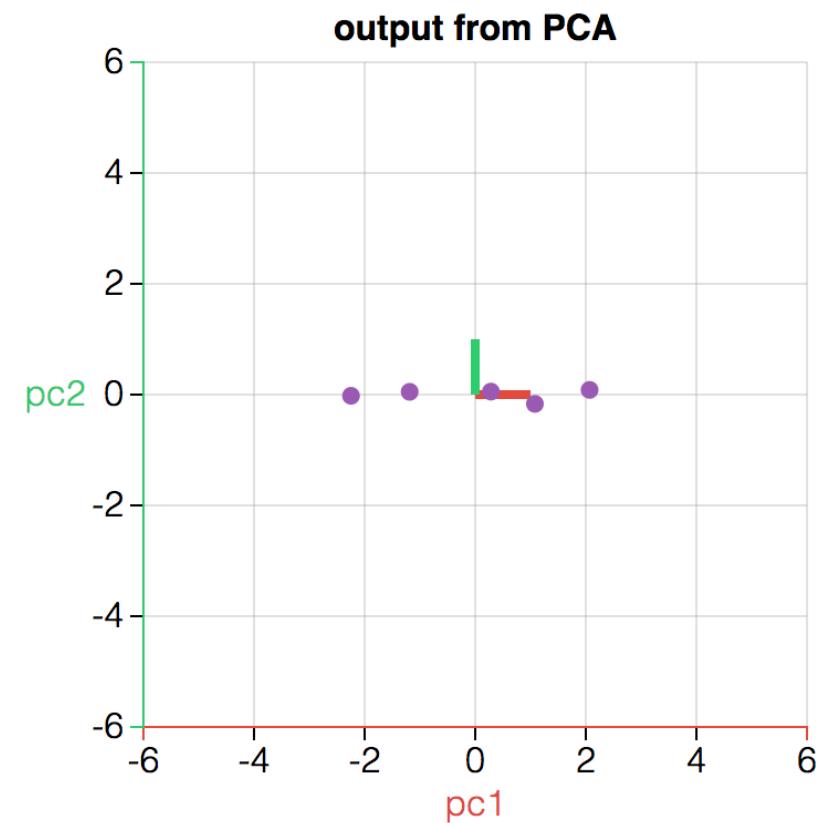
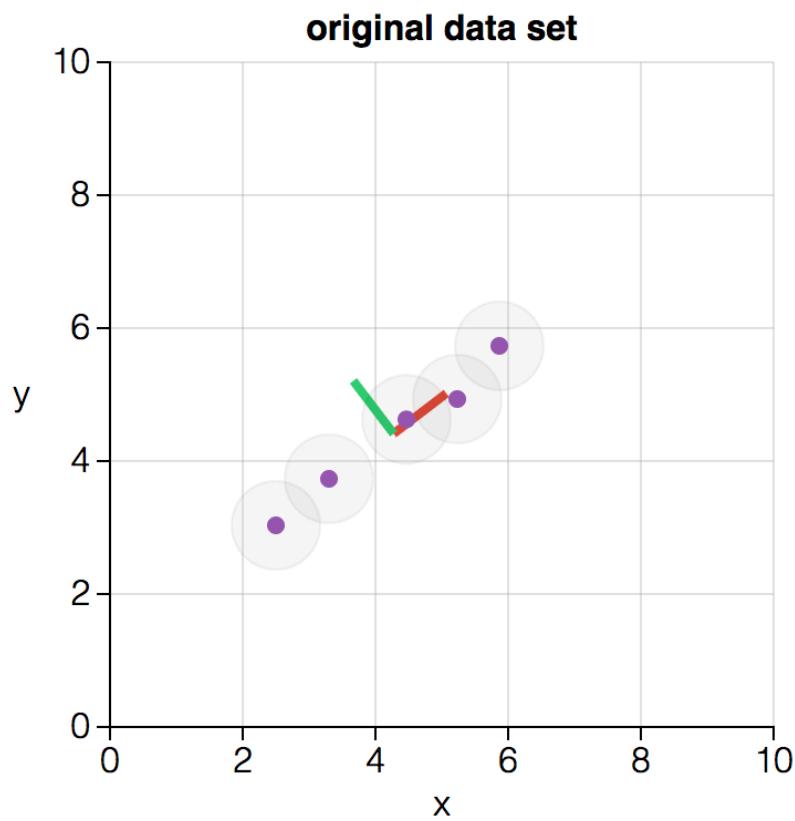


# PRINCIPAL COMPONENT ANALYSIS

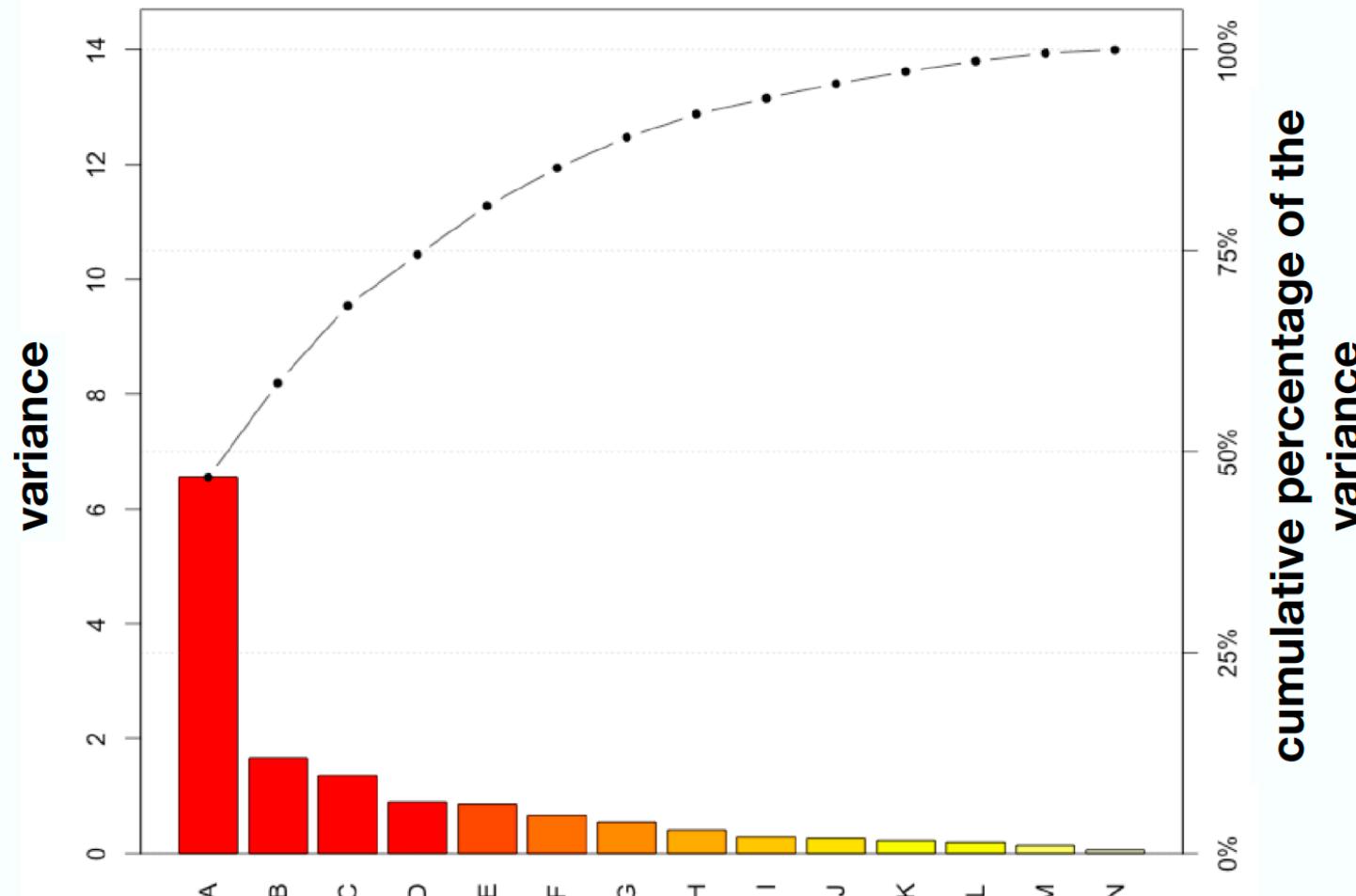
1. COMPUTE THE COVARIANCE MATRIX OF YOUR DATA
2. COMPUTE THE EIGENVALUES AND EIGENVECTORS OF THE COVARAINCE MATRIX
3. THE EIGENVECTORS PROVIDE THE DIRECTION OF MAXIMUM VARIANCE AND THE EIGENVALUES THE 'IMPORTANCE' OF THAT PARTICULAR FEATURE

**REMEMBER, THE COVARIANCE MATRIX IS NOTHING ELSE THAN THAT:**

$$K_{\mathbf{XX}} = \begin{bmatrix} E[(X_1 - E[X_1])(X_1 - E[X_1])] & E[(X_1 - E[X_1])(X_2 - E[X_2])] & \cdots & E[(X_1 - E[X_1])(X_n - E[X_n])] \\ E[(X_2 - E[X_2])(X_1 - E[X_1])] & E[(X_2 - E[X_2])(X_2 - E[X_2])] & \cdots & E[(X_2 - E[X_2])(X_n - E[X_n])] \\ \vdots & \vdots & \ddots & \vdots \\ E[(X_n - E[X_n])(X_1 - E[X_1])] & E[(X_n - E[X_n])(X_2 - E[X_2])] & \cdots & E[(X_n - E[X_n])(X_n - E[X_n])] \end{bmatrix}$$

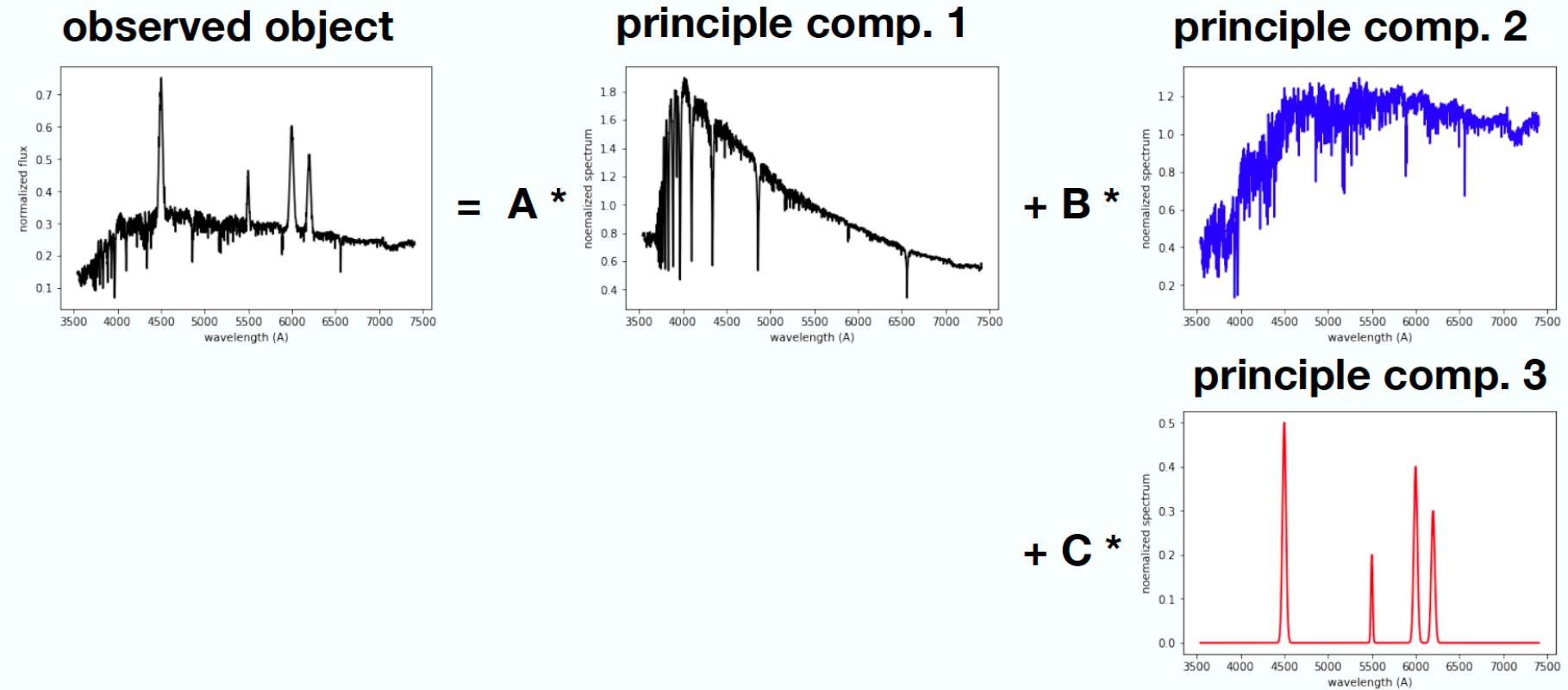


# PRINCIPAL COMPONENT ANALYSIS



IT RESULTS IN DATA COMPRESSION,  
BY REPRESENTING EACH OBJECT AS A PROJECTION OF THE FIRST PRINCIPLE COMPONENTS

# PRINCIPAL COMPONENT ANALYSIS



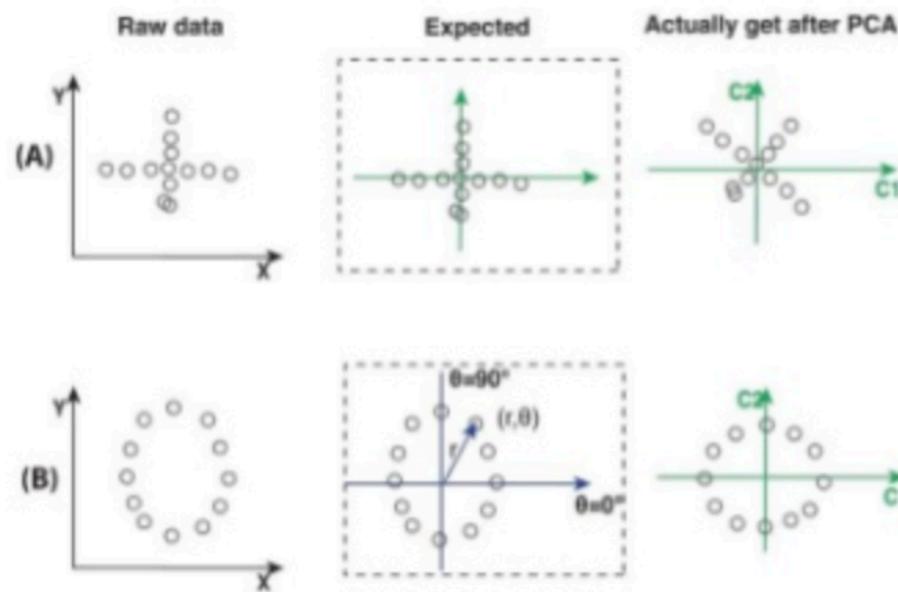
OBJECTS ARE DIVIDED INTO MAJOR  
PROTOTYPES AND ALL OBJECTS CAN BE  
OBTAINED AS LINEAR COMBINATIONS

@D. BARON

# LIMITATIONS OF PCA

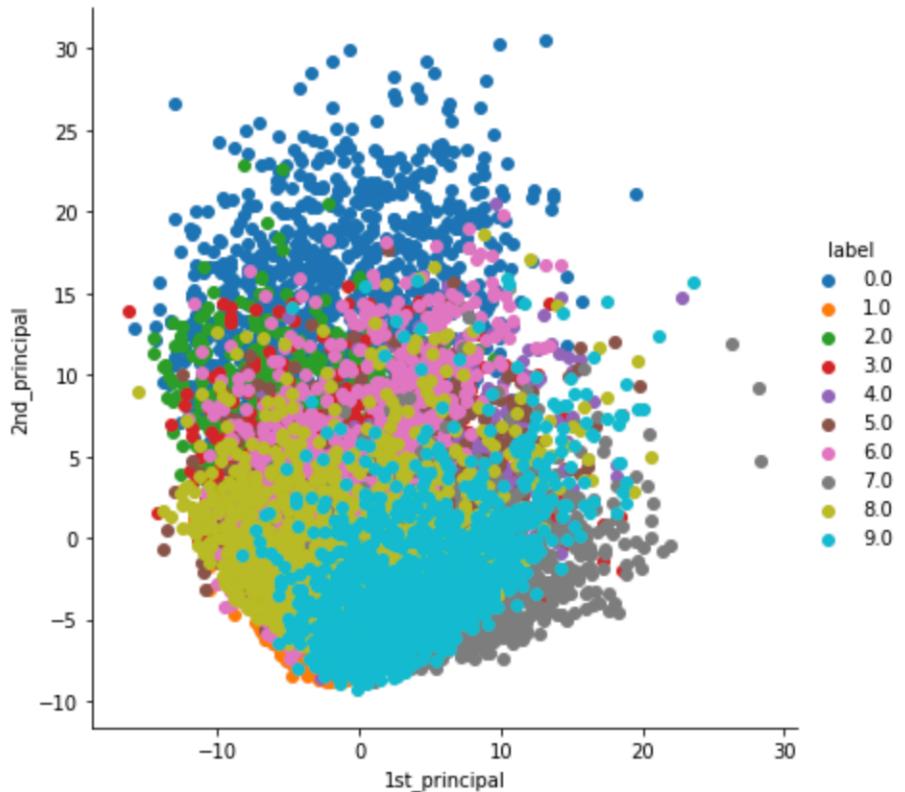
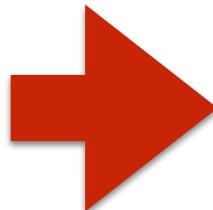
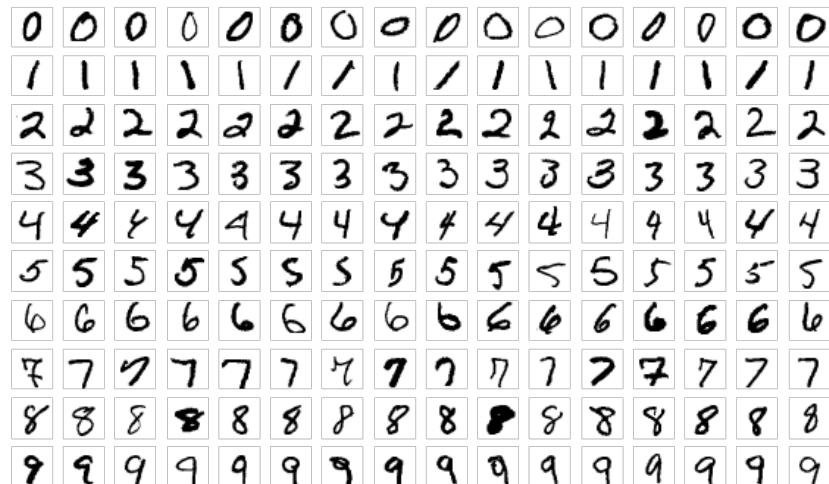
PCA APPLY LINEAR TRANSFORMATIONS

SINCE WE USE THE COVARIANCE MATRIX, IT ASSUMES THAT THE DATA FOLLOWS A **MULTIDIMENSIONAL NORMAL DISTRIBUTION**



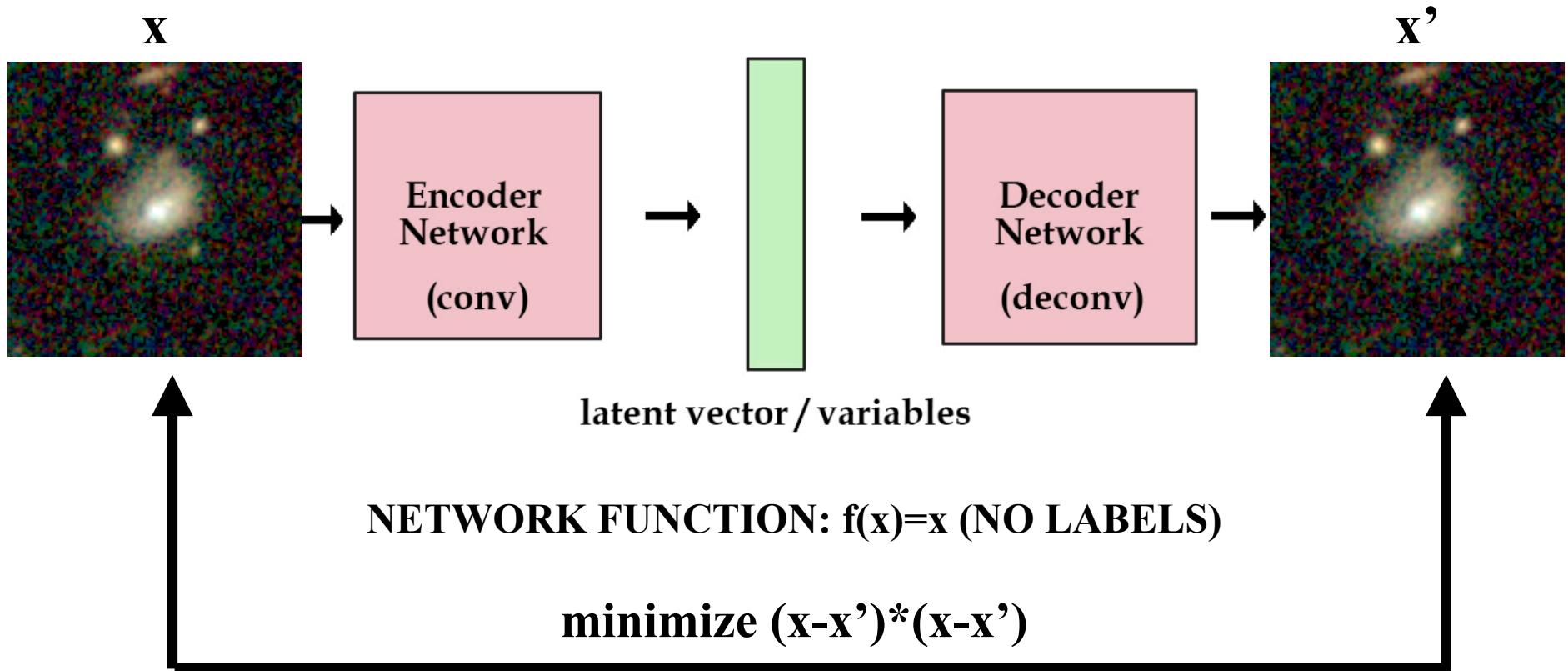
# LIMITATIONS OF PCA

AND DATA IS NOT ALWAYS GAUSSIAN....



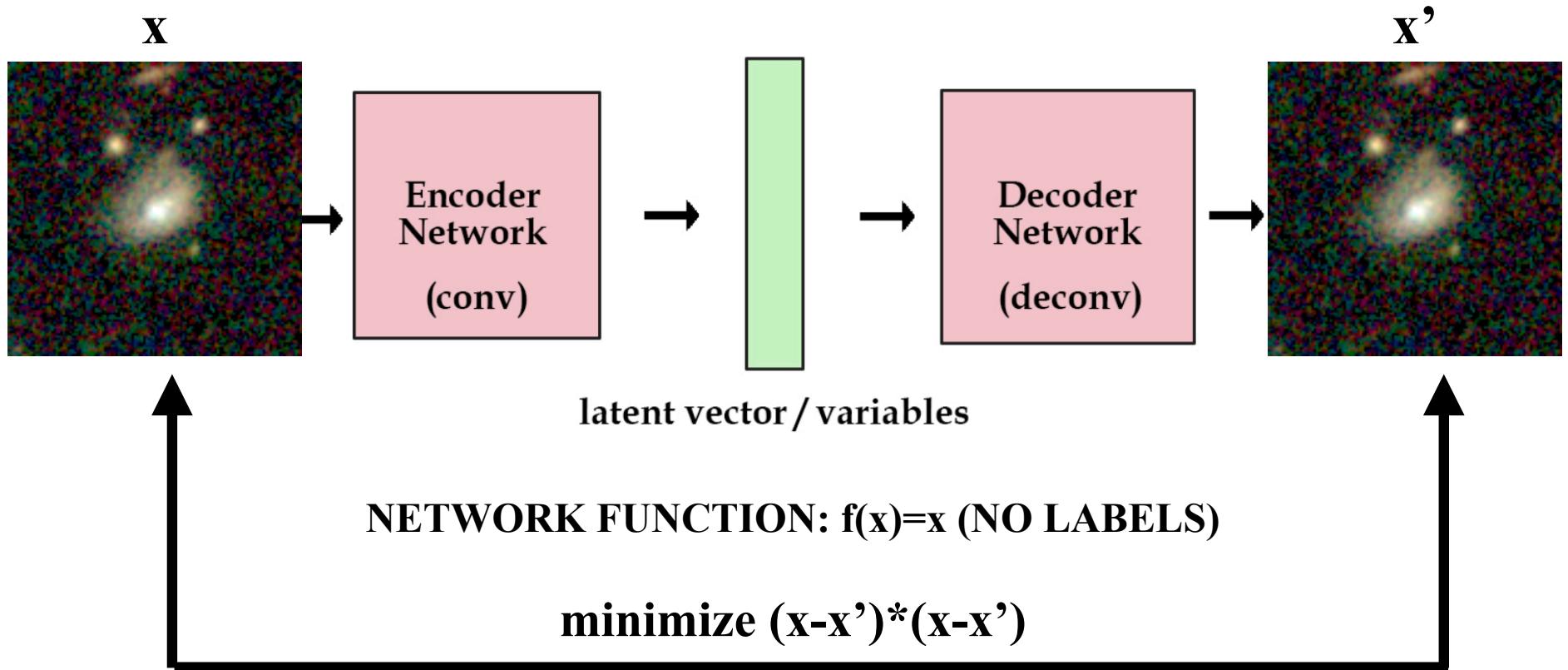
**CAN WE GENERALIZE THAT?**

# CONVOLUTIONAL AUTO-ENCODER



AN AUTO-ENCODER IS ANY NETWORK WITH IDENTICAL INPUT AND OUTPUT

# CONVOLUTIONAL AUTO-ENCODER



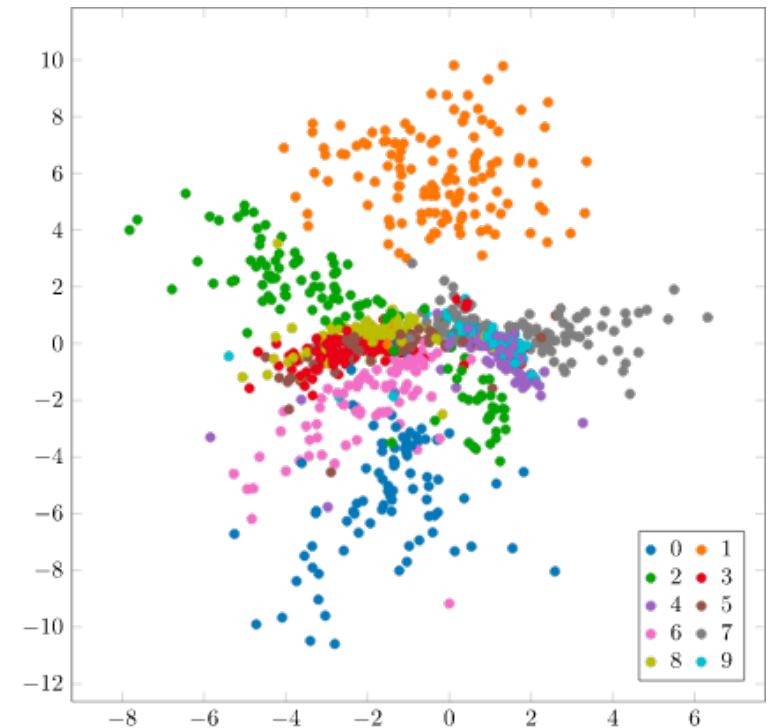
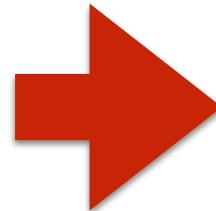
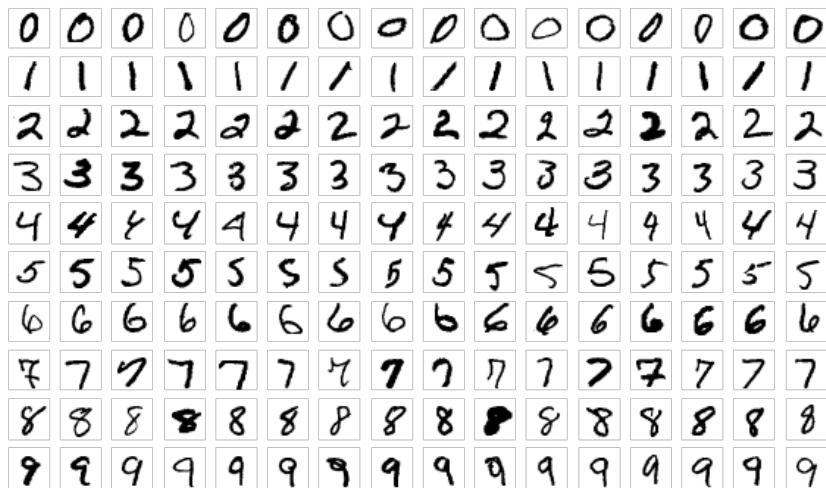
BY REDUCING THE DIMENSIONALITY IN THE LATENT SPACE WE FORCE THE NETWORK TO LEARN A REPRESENTATION OF THE INPUT DATA IN A LOWER DIMENSIONALITY SPACE

- \* NO NEED TO BE CONVOLUTIONAL - ANY NEURAL NETWORK WITH A BOTTLENECK WILL DO THE JOB

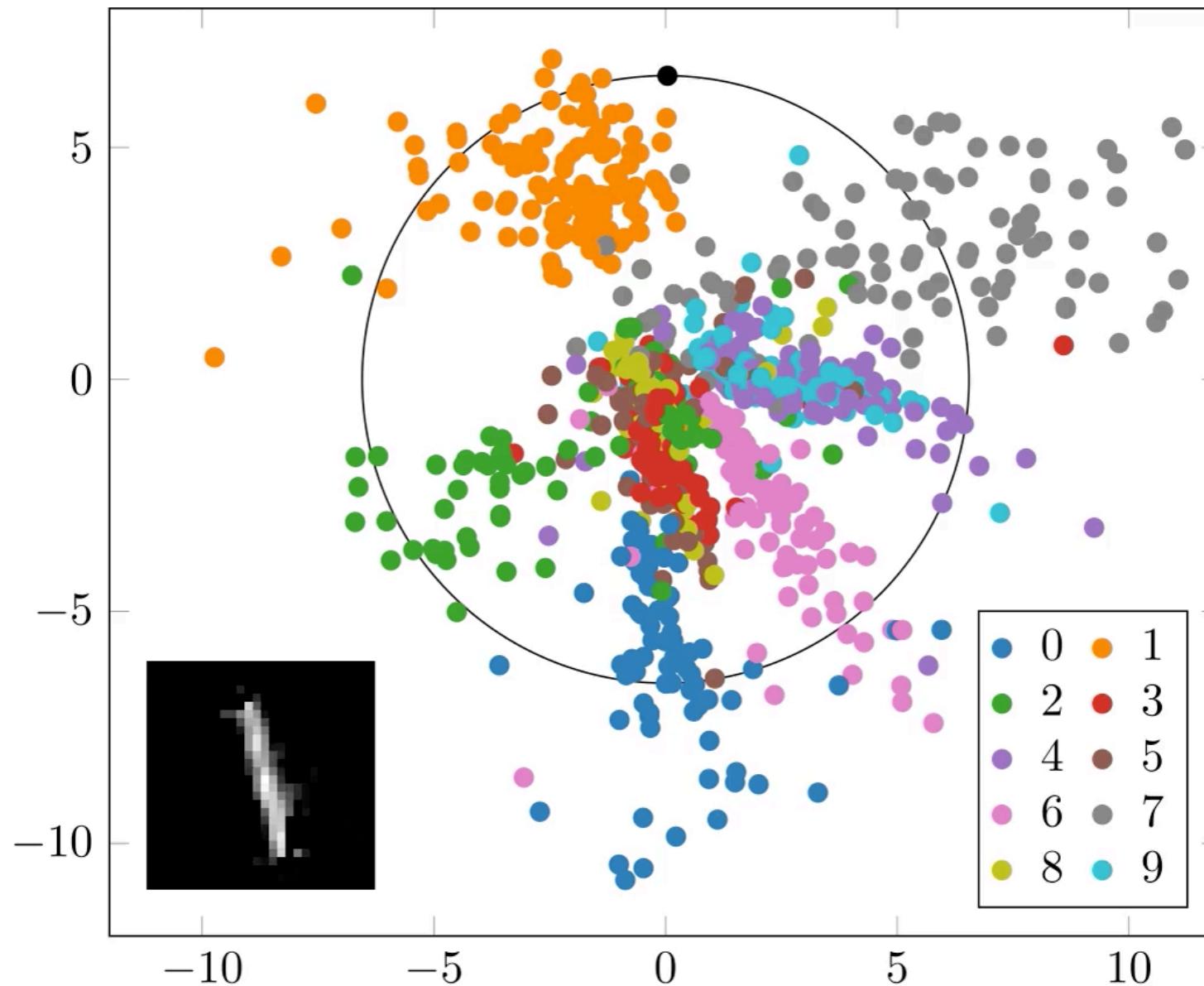
- \* **QUESTION:** WHAT WOULD HAPPEN IF WE SET AN AUTOENCODER WITH NO ACTIVATION FUNCTIONS?

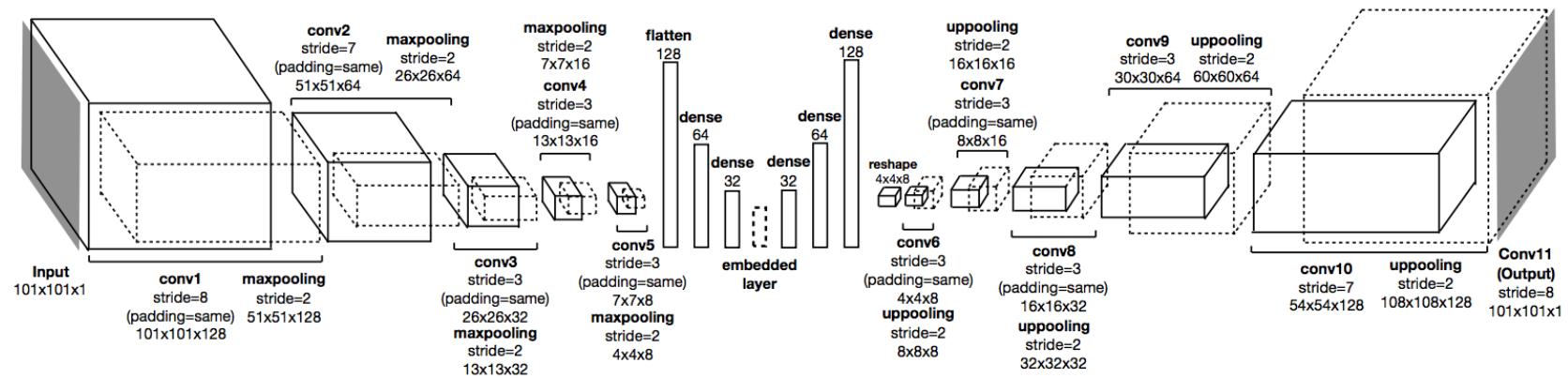
SEE EXAMPLE FROM TUTORIALS FOR STAR-GALAXY  
SEPARATION

# AUTOENCODER REPRESENTATION OF MNIST

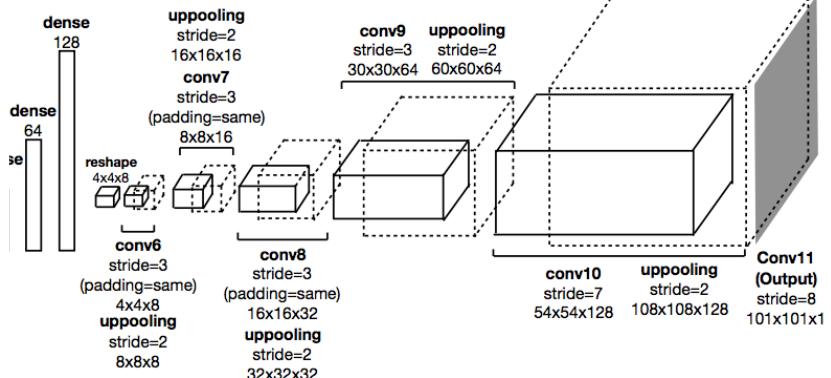
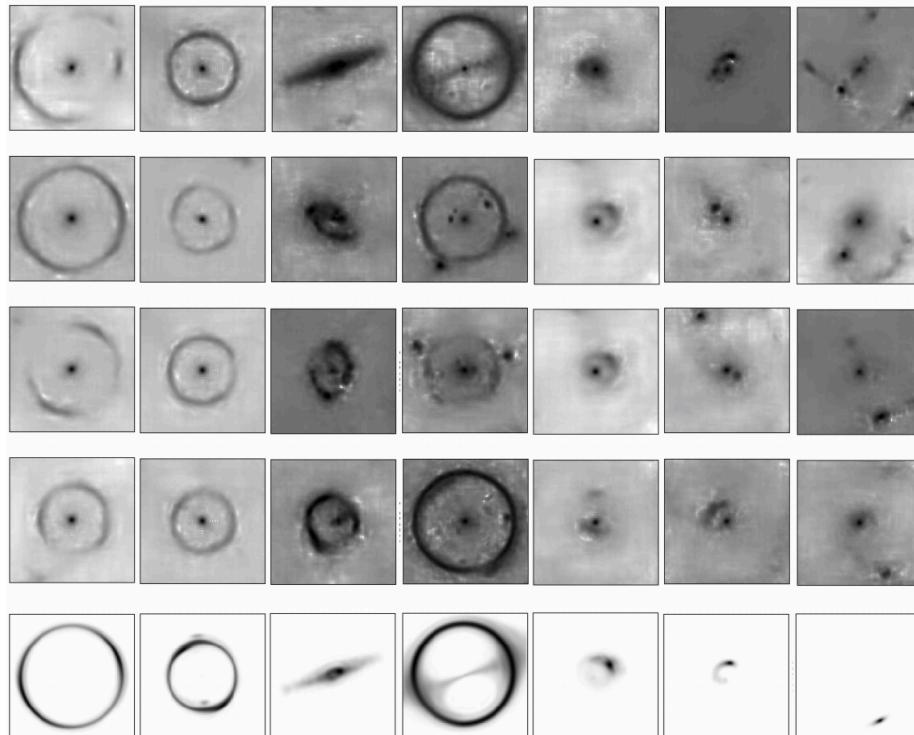


# AUTOENCODER REPRESENTATION OF MNIST

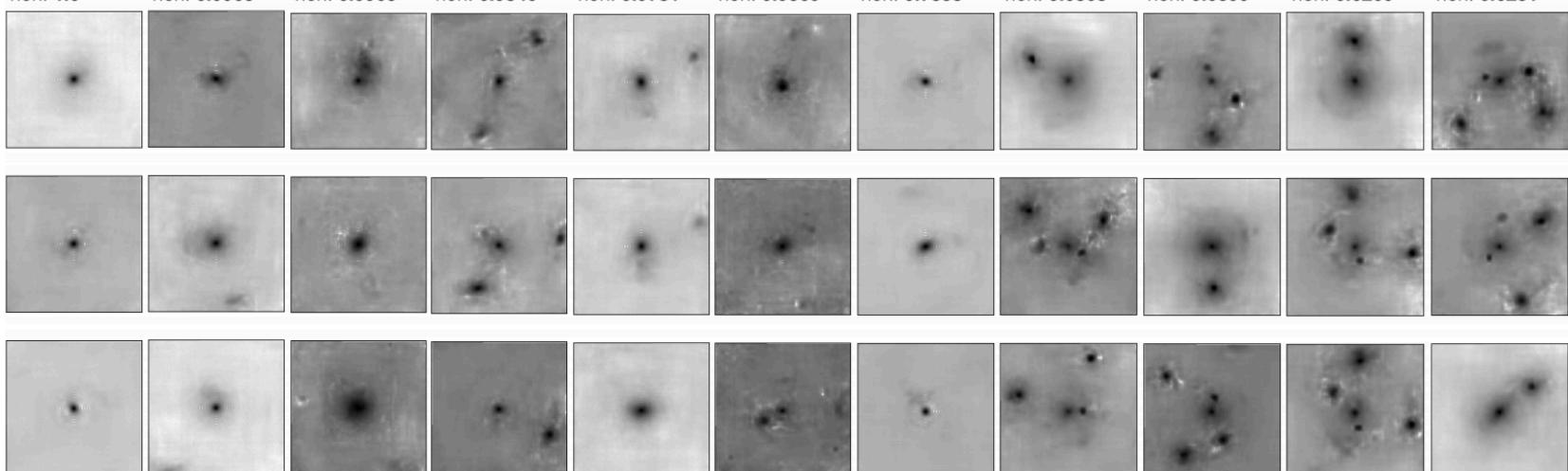




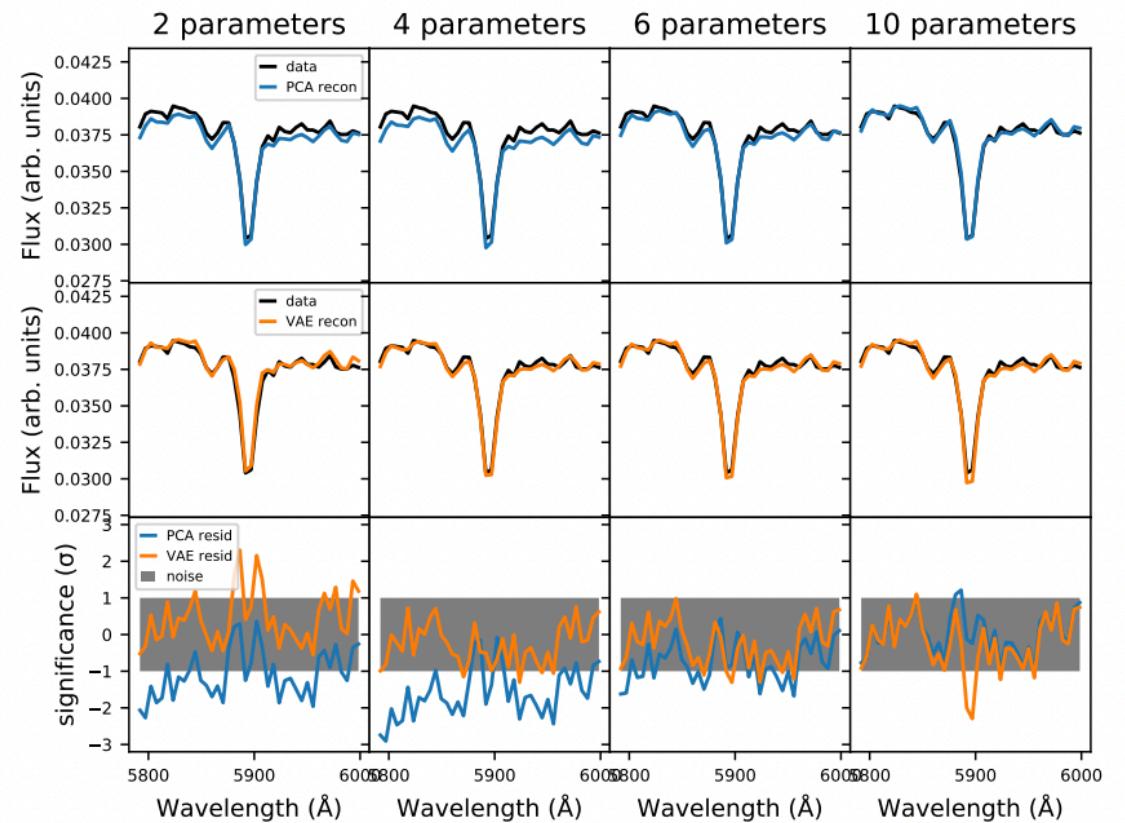
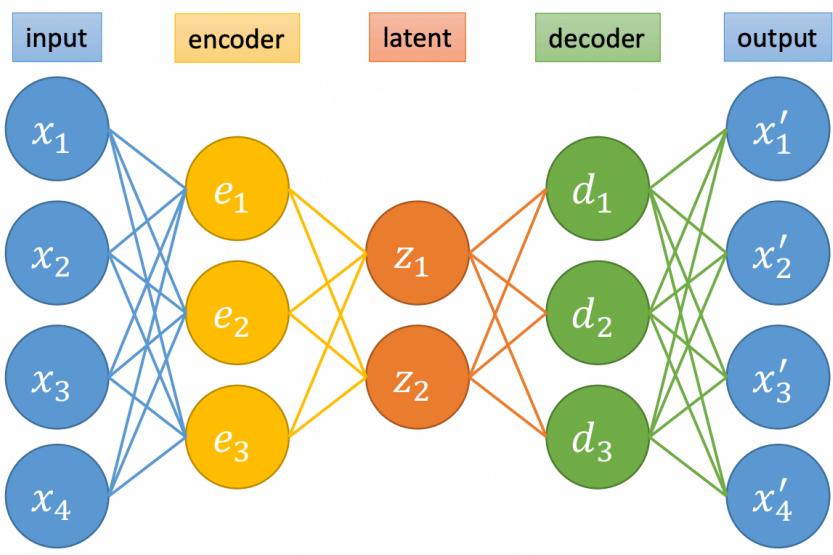
<b>Cluster 17:</b> lensing: 0.9873 non: 0.0127 F_len: 0.0914	<b>Cluster 21:</b> lensing: 0.9448 non: 0.0552 F_len: 0.0695	<b>Cluster 1:</b> lensing: 0.9159 non: 0.0841 F_len: 0.0731	<b>Cluster 6:</b> lensing: 0.8997 non: 0.1003 F_len: 0.0729	<b>Cluster 2:</b> lensing: 0.803 non: 0.197 F_len: 0.1945	<b>Cluster 20:</b> lensing: 0.6206 non: 0.3794 F_len: 0.0575	<b>Cluster 5:</b> lensing: 0.6170 non: 0.3830 F_len: 0.0734
---	---	--	--	--	---	--



<b>Cluster 14:</b> lensing: 0.0 non: 1.0	<b>Cluster 3:</b> lensing: 0.0037 non: 0.9963	<b>Cluster 8:</b> lensing: 0.0037 non: 0.9963	<b>Cluster 22:</b> lensing: 0.0154 non: 0.9846	<b>Cluster 16:</b> lensing: 0.0219 non: 0.9781	<b>Cluster 4:</b> lensing: 0.0431 non: 0.9569	<b>Cluster 0:</b> lensing: 0.2642 non: 0.7358	<b>Cluster 18:</b> lensing: 0.3132 non: 0.6868	<b>Cluster 19:</b> lensing: 0.3601 non: 0.6399	<b>Cluster 13:</b> lensing: 0.3731 non: 0.6269	<b>Cluster 9:</b> lensing: 0.3769 non: 0.6231
--	---	---	--	--	---	---	--	--	--	---



Cheng+20



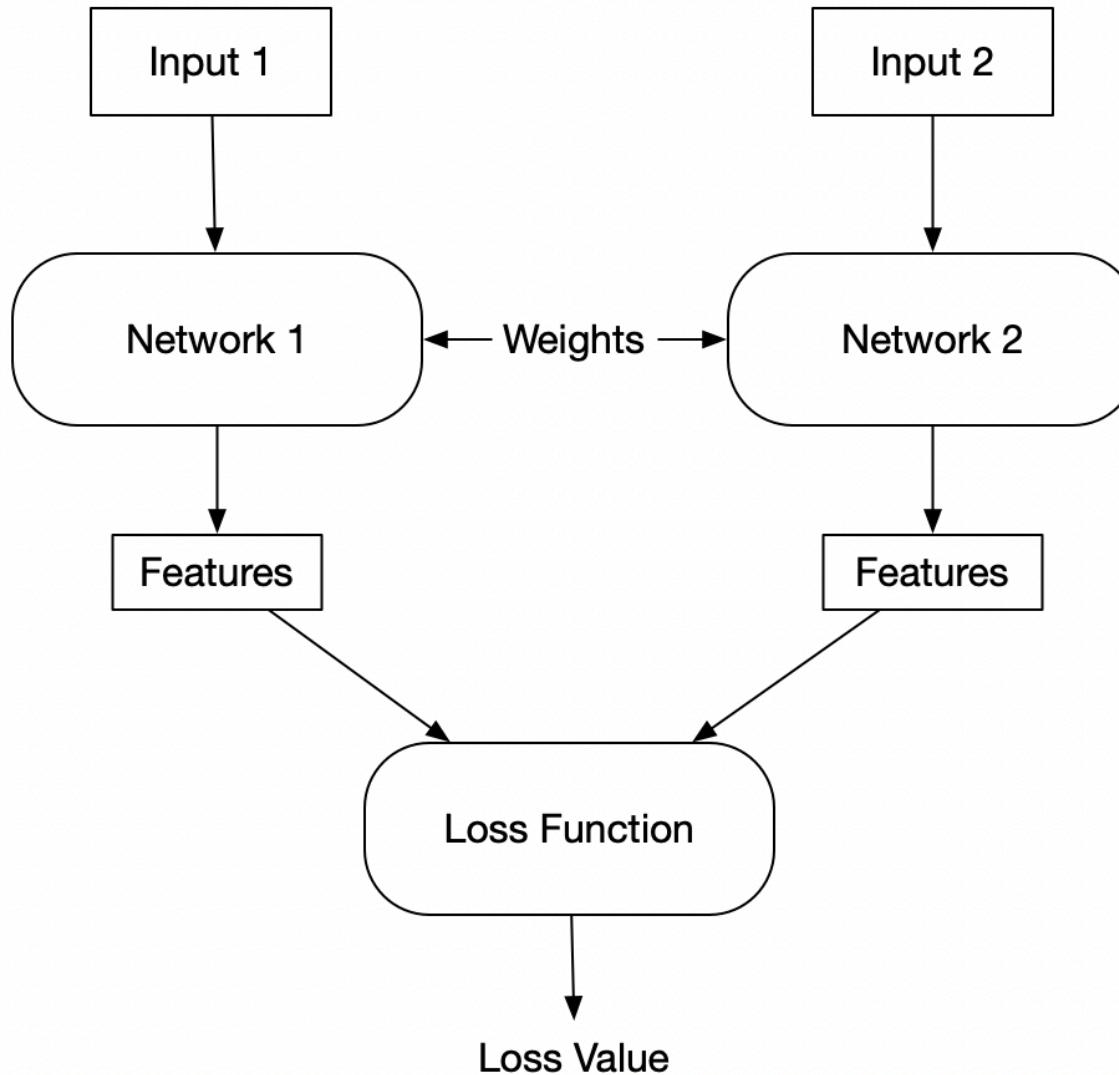
# **SELF-SUPERVISED CONTRASTIVE LEARNING**

Humans tend to identify objects without remembering all the details, by creating some abstract representations which are used to identify new objects of similar time.

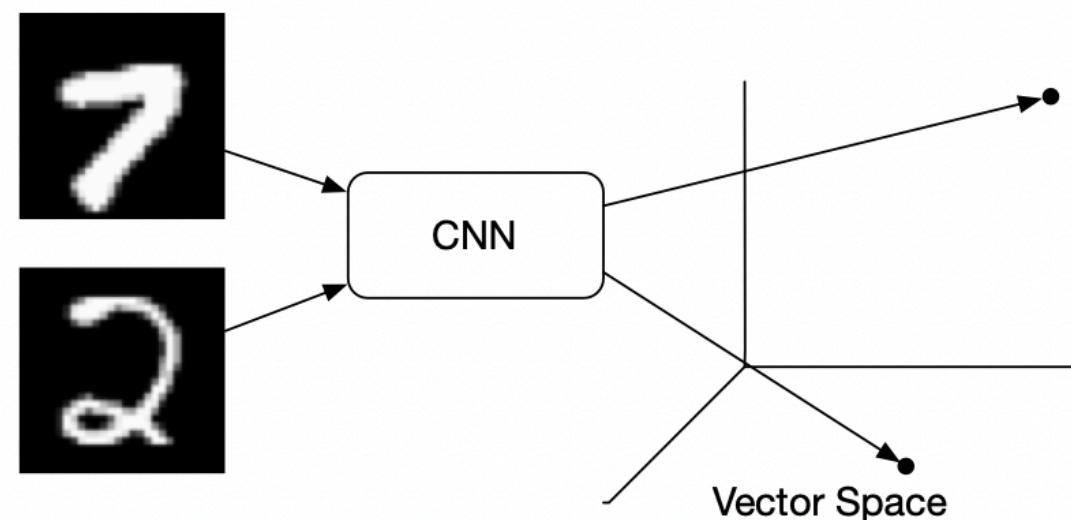
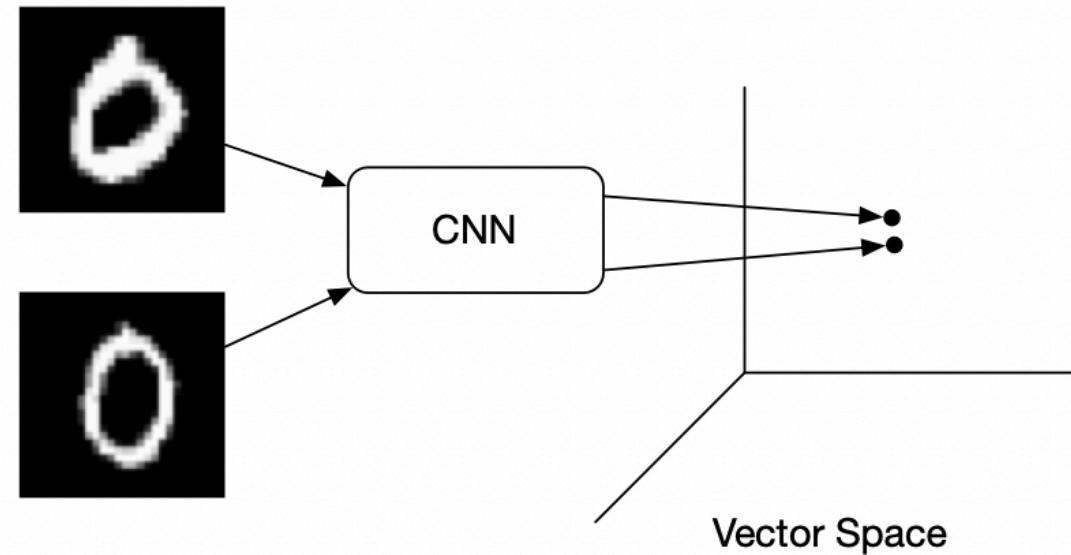
The main goal of self-supervised contrastive learning is to create and generalize these representations.

Therefore, it aims at building some general representations that can be used for other “downstream tasks”

Siamese networks are two networks sharing weights and common loss function



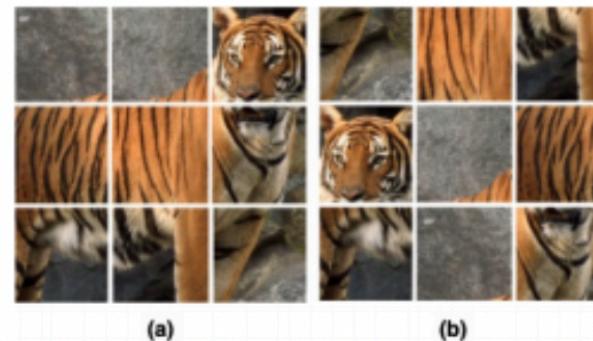
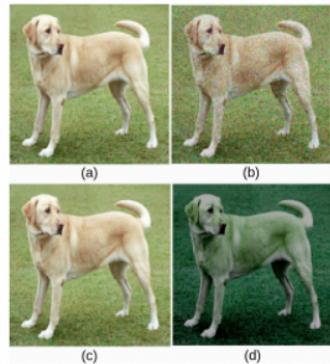
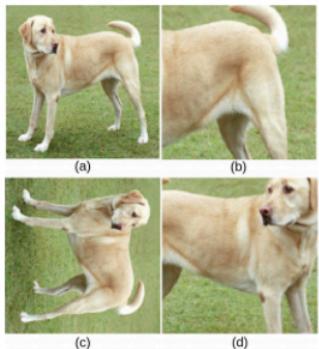
In a supervised setting, similar objects are used to be together by simply using a crossentropy loss



In most cases, labels are not available.

How do you recognise similar objects?

What would be the loss function?

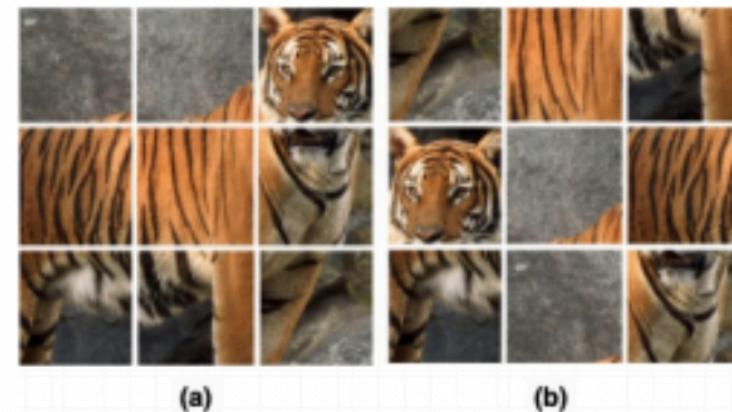
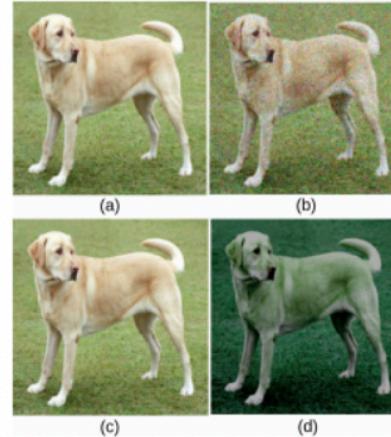
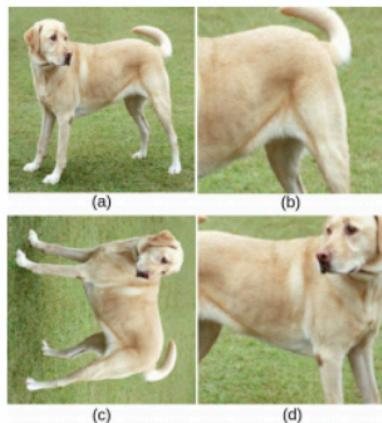


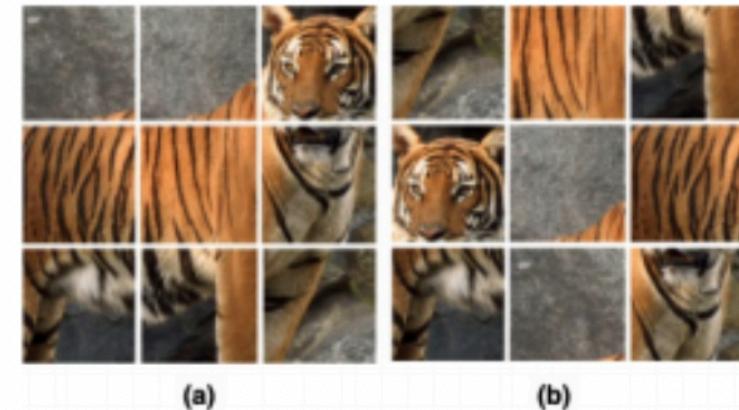
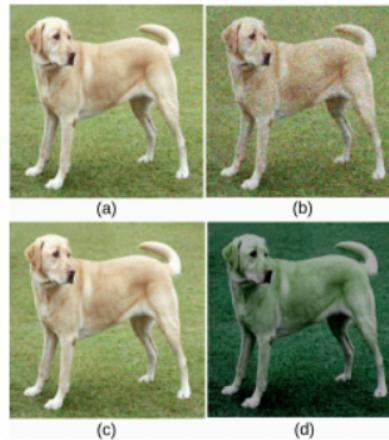
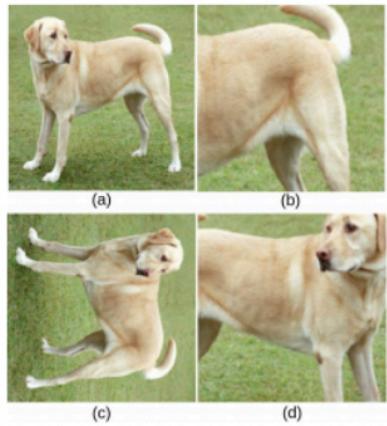
In most cases, labels are not available.

How do you recognise similar objects?

What would be the loss function?

We create “pretext tasks” from a unique image:





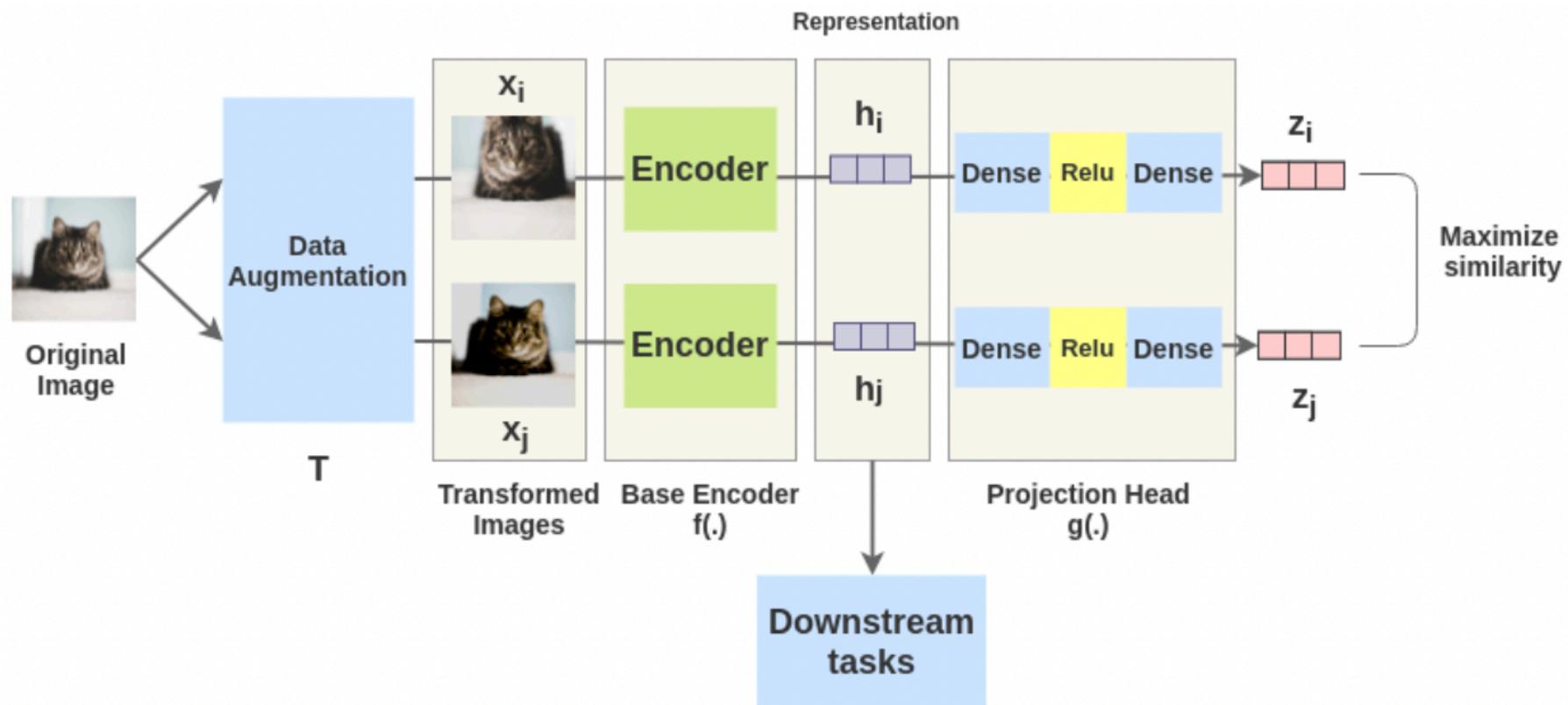
## Color Augmentation

## Image Rotation

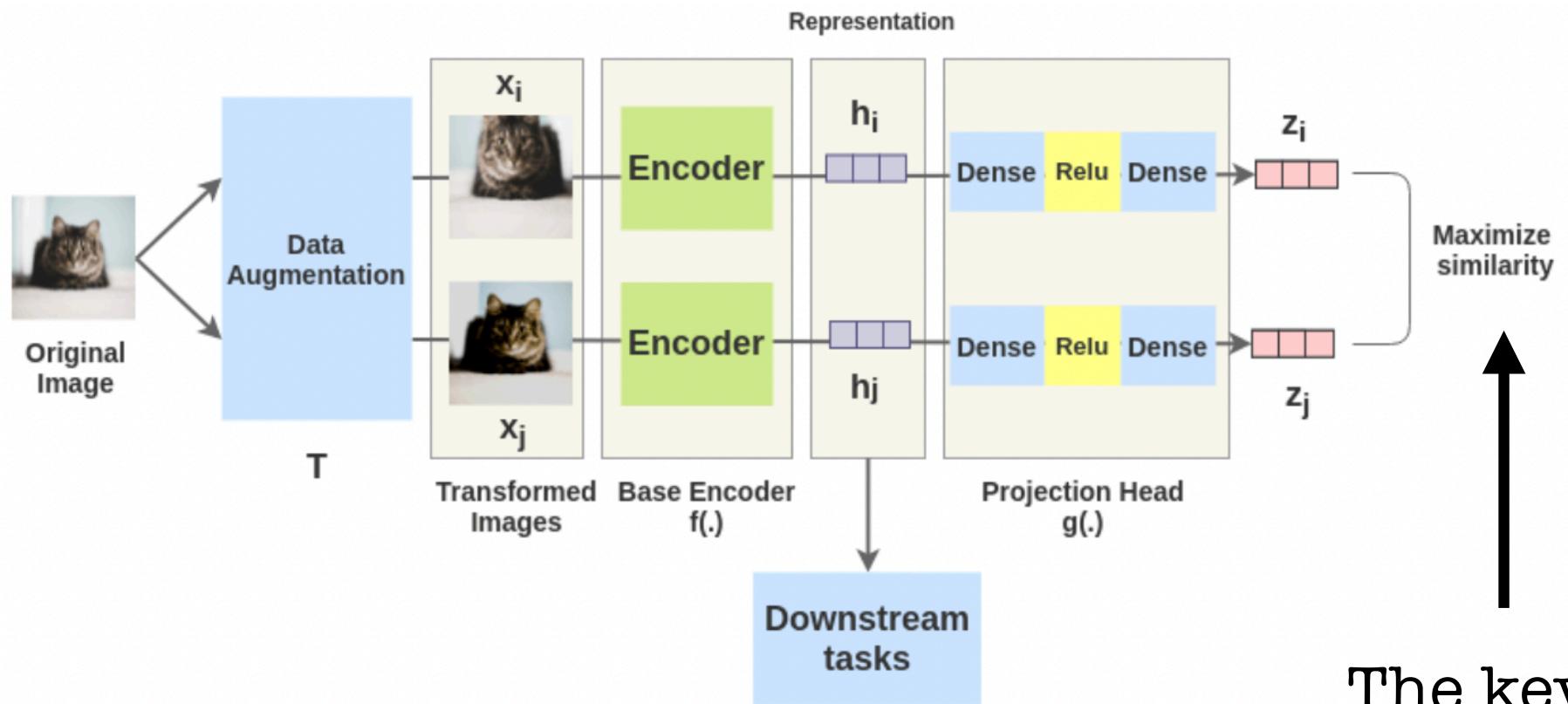
## Image Cropping

## Any geometrical transformation ...

The augmented versions of the images are passed through siamese networks and projected into a latent variable  $z$



The augmented versions of the images are passed through siamese networks and projected into a latent variable  $z$



The key is  
the loss  
function

Similarly between  
two representations of positive pairs

The contrastive loss:

$$l_{i,j} = - \log \frac{\exp(\langle z_i, z_j \rangle / h)}{\sum_{k=1, k \neq i}^{2N} \exp(\langle z_i, z_k \rangle / h)},$$

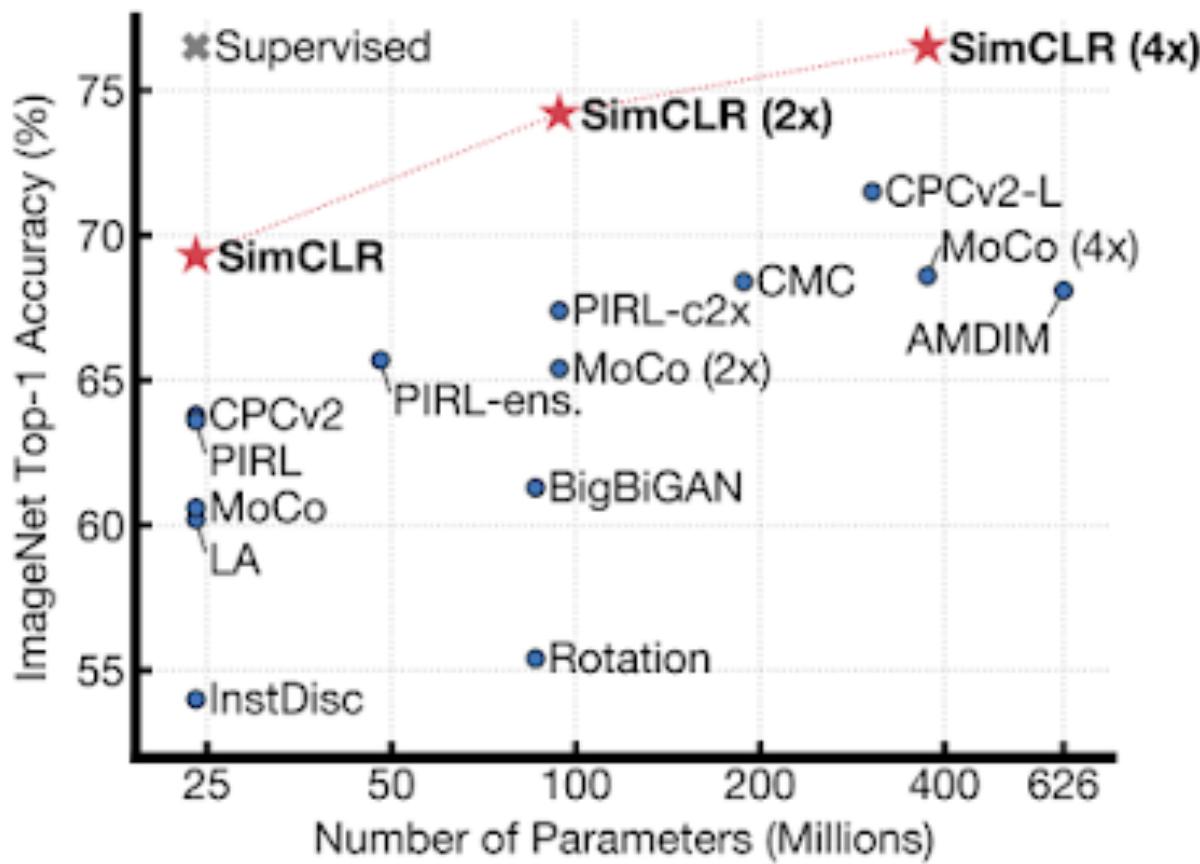
Sum of all similarities between  
negative pairs

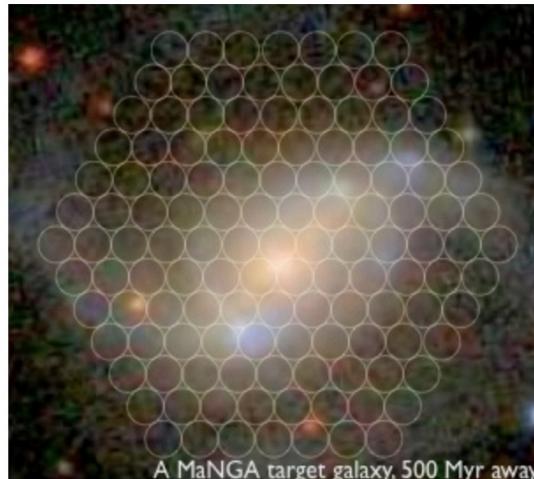
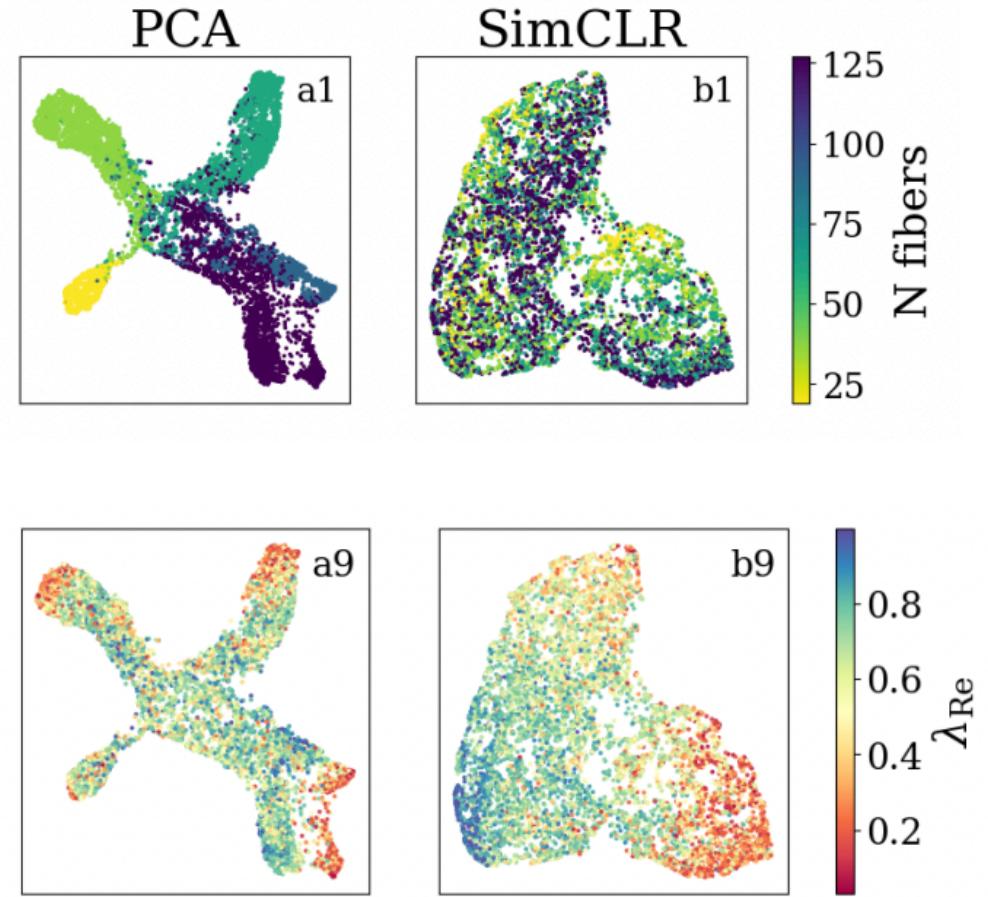
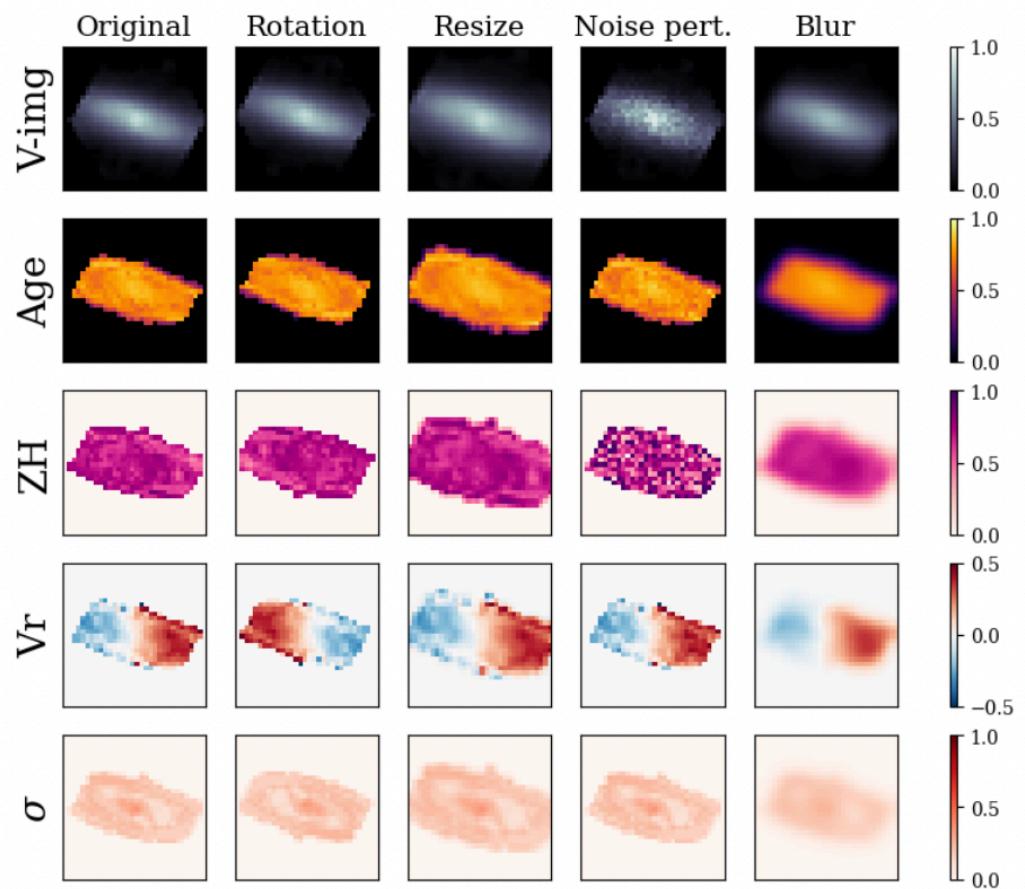
# LEARNING REPRESENTATIONS THROUGH **CONTRASTIVE** LEARNING

## **CONTRASTIVE LOSS:**

$$l_{i,j} = -\log \frac{\exp(\langle z_i, z_j \rangle / h)}{\sum_{k=1, k \neq i}^{2N} \exp(\langle z_i, z_k \rangle / h)},$$

# SELF-SUPERVISED LEARNING REACHES COMPARABLE ACCURACY TO FULLY SUPERVISED APPROACHES...





# Contrastive learning representation of Manga galaxies

Sarmiento+21