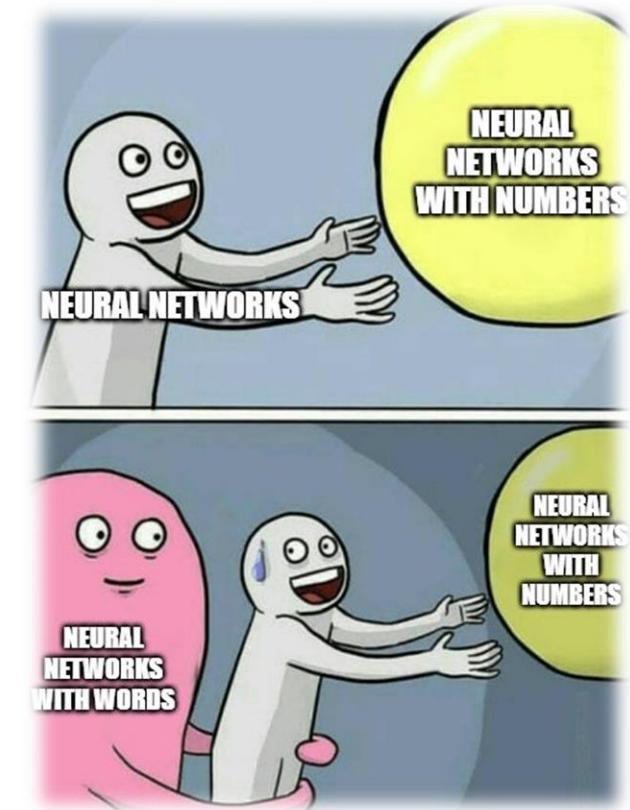# Embeddings

Alberto TONDA, Ph.D. (Senior permanent researcher, DR)

*UMR 518 MIA-PS, INRAE, AgroParisTech, Université Paris-Saclay*
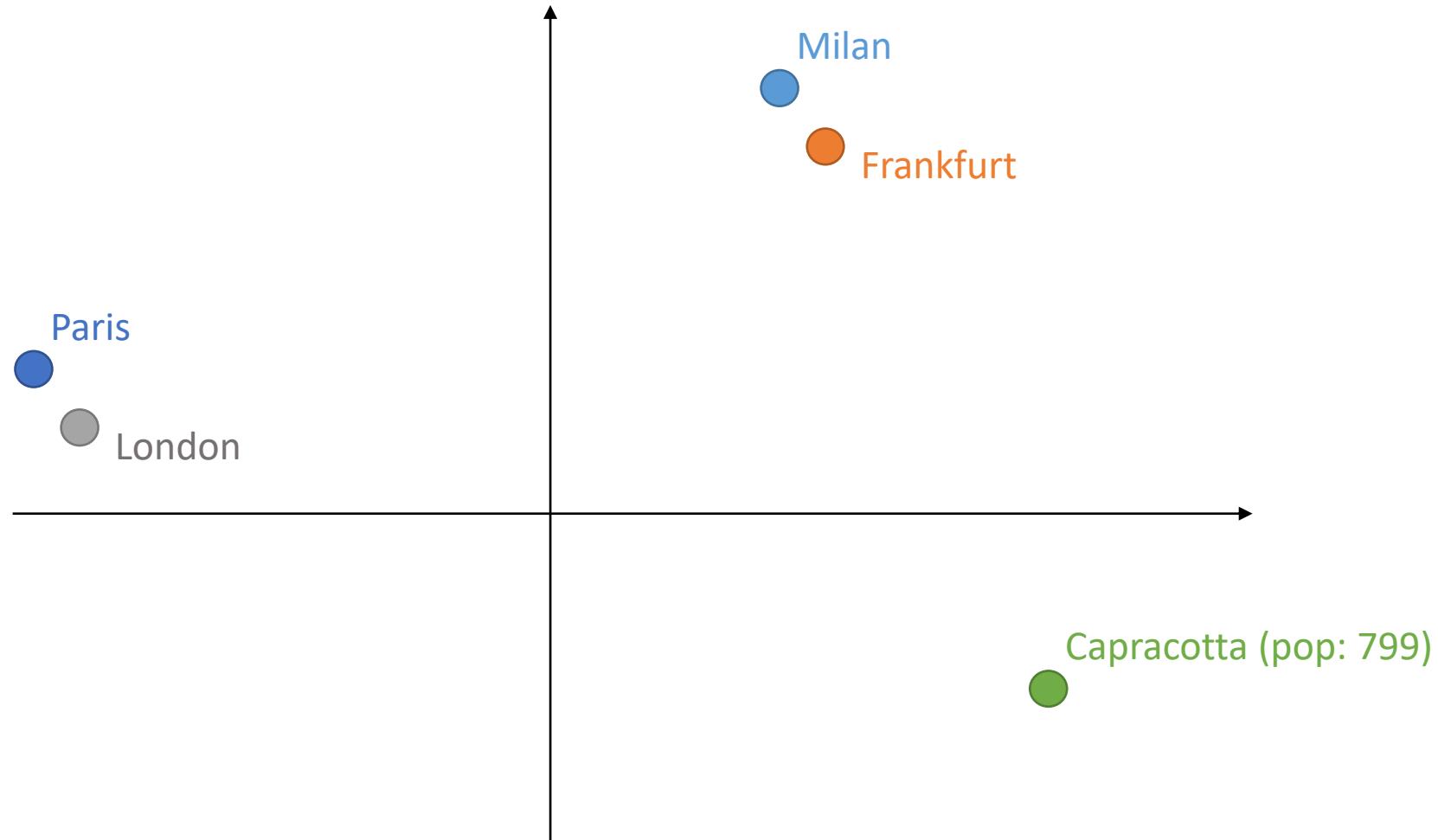*UAR 3611, Institut des Systèmes Complexes de Paris Île-de-France*

# Outline

- Embeddings
- Building embeddings
- Embeddings from trained deep networks
- Embeddings with Autoencoders
- Embeddings of Vocabularies

EMBEDDINGS
Alberto TONDA, Team EKINOCS, UMR 518 MIA-PS, INRAE, Université Paris-Saclay

# Embeddings

- Curiously, an embedding is not necessarily related to DL

- An embedding is a **vectorial space**, $x \in \mathbb{R}^d$ with semantics
  - Distances and relative positions between points have a *meaning*
  - And displacements/transformations might also have meaning
  - Vocabulary clash, "vector"/"point"


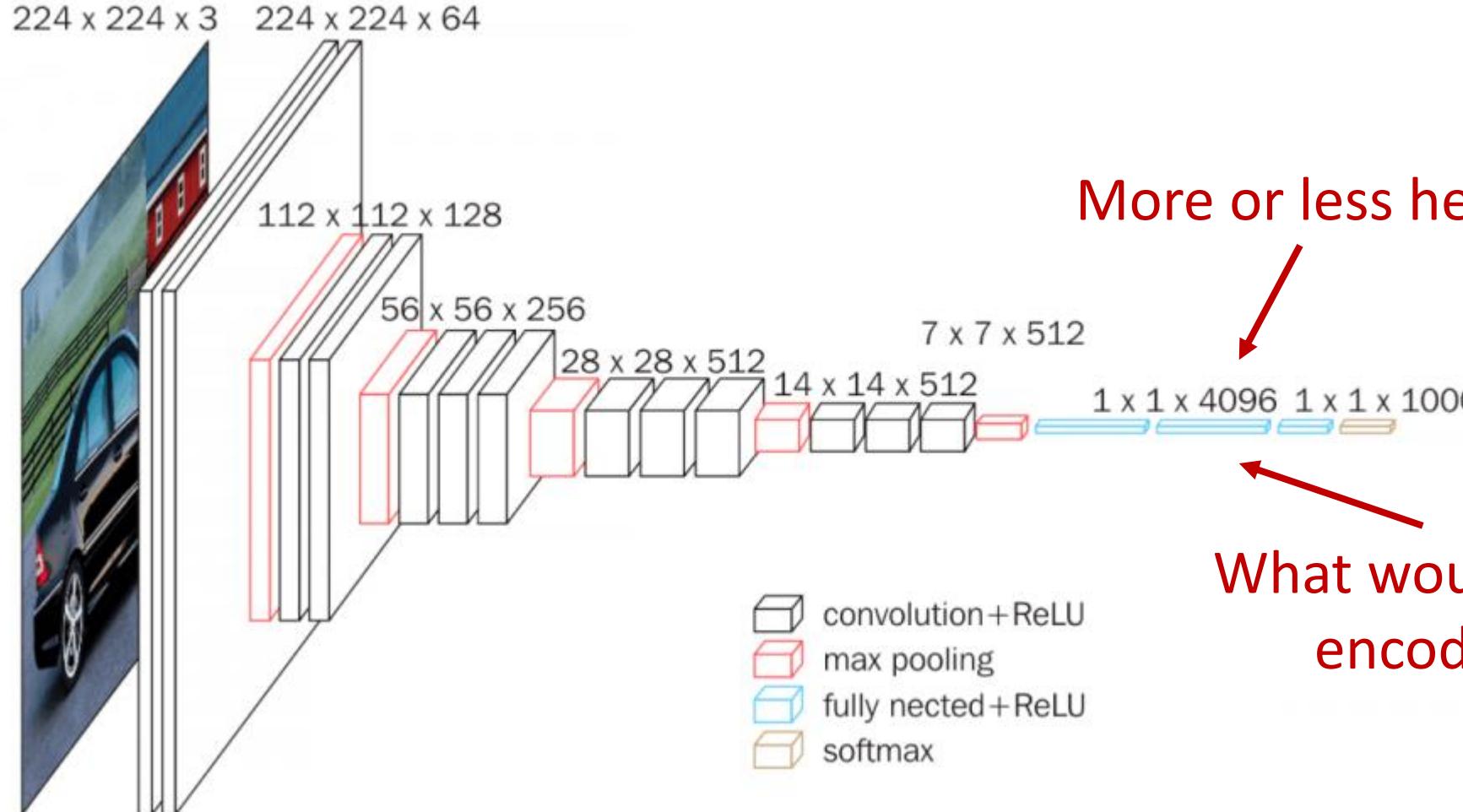- Let's build an embedding! Meaningful dimensions for **cities**?

# Embedding for cities

# Building embeddings

- Building embeddings is pretty *hard*
    - Already when we have meaningful features
    - But for **relational data**? Pixels, words, videos, sound, …?


- Supervised task in DL architecture *might* create embedding
    - Output of a (deep) module can be interpreted as vectorial space
    - Encoding **meaning related to supervised task**

Alberto TONDA, Team EKINOCS, UMR 518 MIA-PS, INRAE, Université Paris-Saclay

# Embeddings from trained deep networks



224 x 224 x 3    224 x 224 x 64

112 x 112 x 128

56 x 56 x 256

28 x 28 x 512

14 x 14 x 512

7 x 7 x 512

1 x 1 x 4096    1 x 1 x 1000

More or less here

What would this encode?

convolution+ReLU
max pooling
fully nected+ReLU
softmax

Alberto TONDA, Team EKINOCS, UMR 518 MIA-PS, INRAE, Université Paris-Saclay
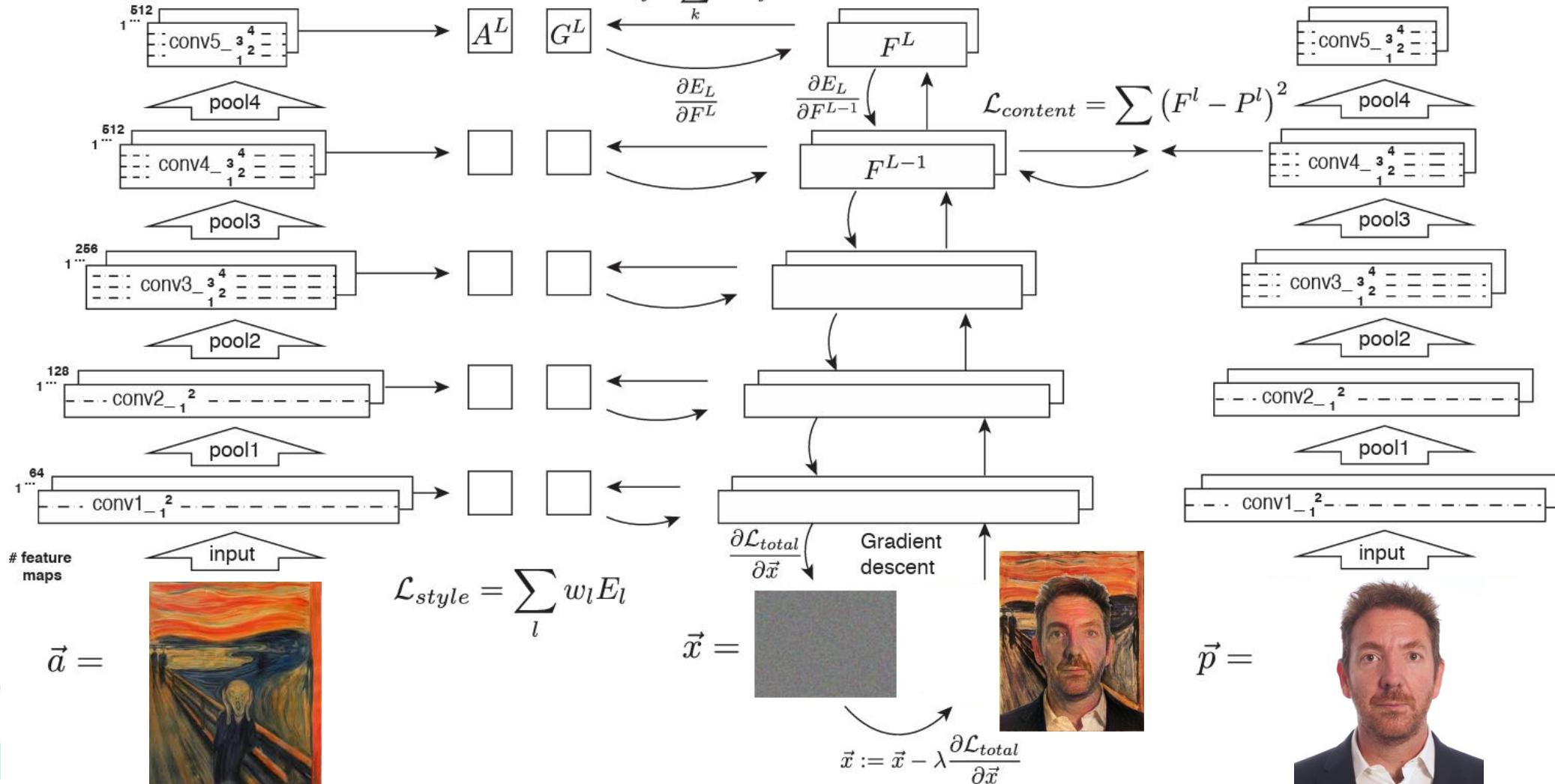
# Embeddings from trained deep networks

- Taking VGG-16, output of deep modules
  - Might encode samples of same class close together
  - If we are lucky, samples of *visually similar classes* close together
  - However, encoded meaning related to **visual aspects**, only
  - Photo of bear and tiger very different, both ferocious predators
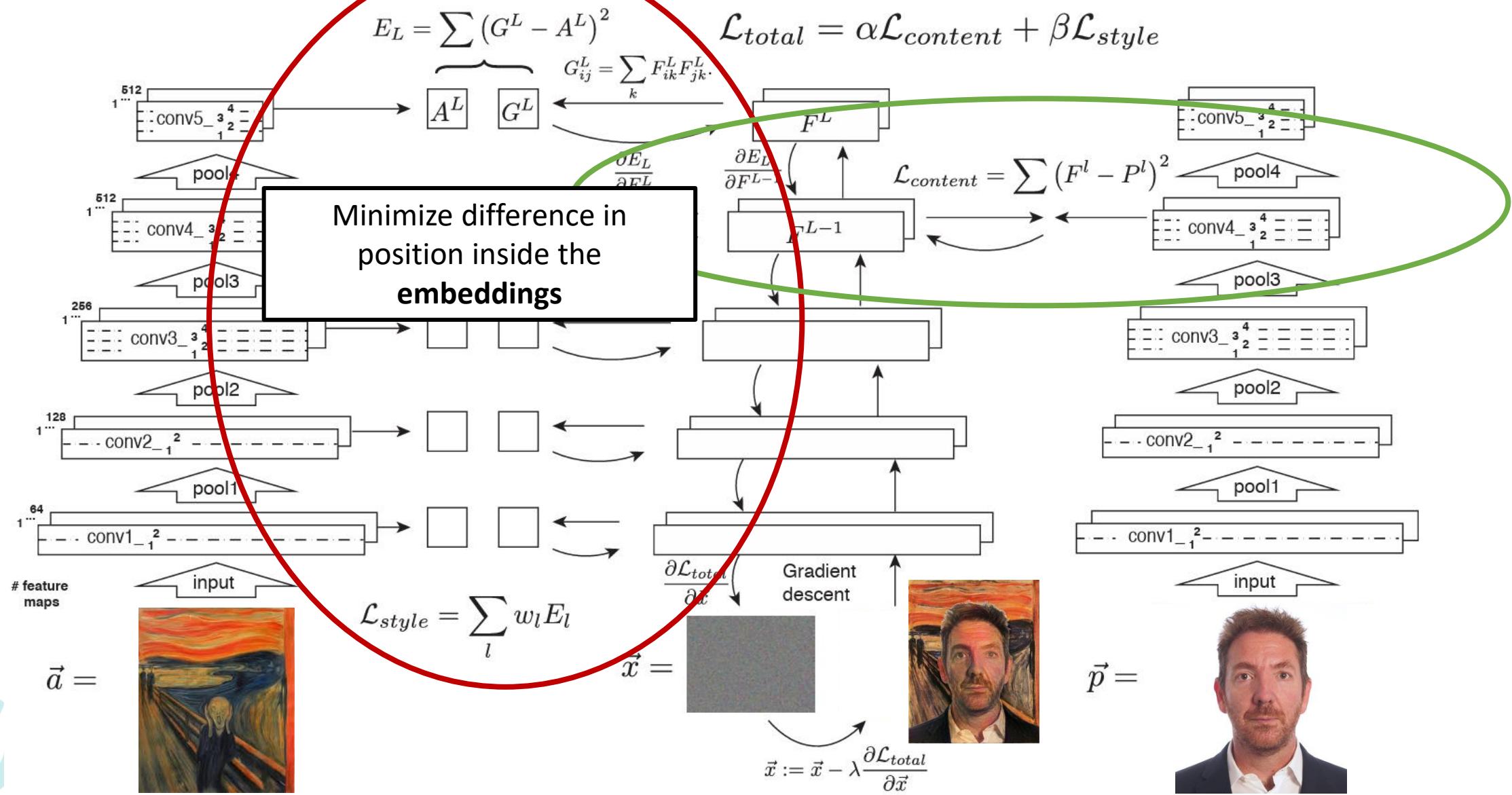
# Embeddings from trained deep networks



$$E_L = \sum \left( G^L - A^L \right)^2 \qquad \mathcal{L}_{total} = \alpha \mathcal{L}_{content} + \beta \mathcal{L}_{style}$$

$$G_{ij}^L = \sum_k F_{ik}^L F_{jk}^L.$$

$$\frac{\partial E_L}{\partial F^L} \qquad \frac{\partial E_L}{\partial F^{L-1}} \qquad \mathcal{L}_{content} = \sum \left( F^l - P^l \right)^2$$

$$\mathcal{L}_{style} = \sum_l w_l E_l$$

$$\frac{\partial \mathcal{L}_{total}}{\partial \vec{x}} \quad \text{Gradient descent}$$

$$\vec{x} := \vec{x} - \lambda \frac{\partial \mathcal{L}_{total}}{\partial \vec{x}}$$

$$E_L = \sum \left(G^L - A^L\right)^2 \qquad \mathcal{L}_{total} = \alpha \mathcal{L}_{content} + \beta \mathcal{L}_{style}$$

$$G_{ij}^L = \sum_k F_{ik}^L F_{jk}^L$$

$$\mathcal{L}_{content} = \sum \left(F^l - P^l\right)^2$$

Minimize difference in position inside the **embeddings**

$$\mathcal{L}_{style} = \sum_l w_l E_l$$

$$\vec{a} = \qquad \vec{x} = \qquad \text{Gradient descent} \qquad \vec{p} =$$

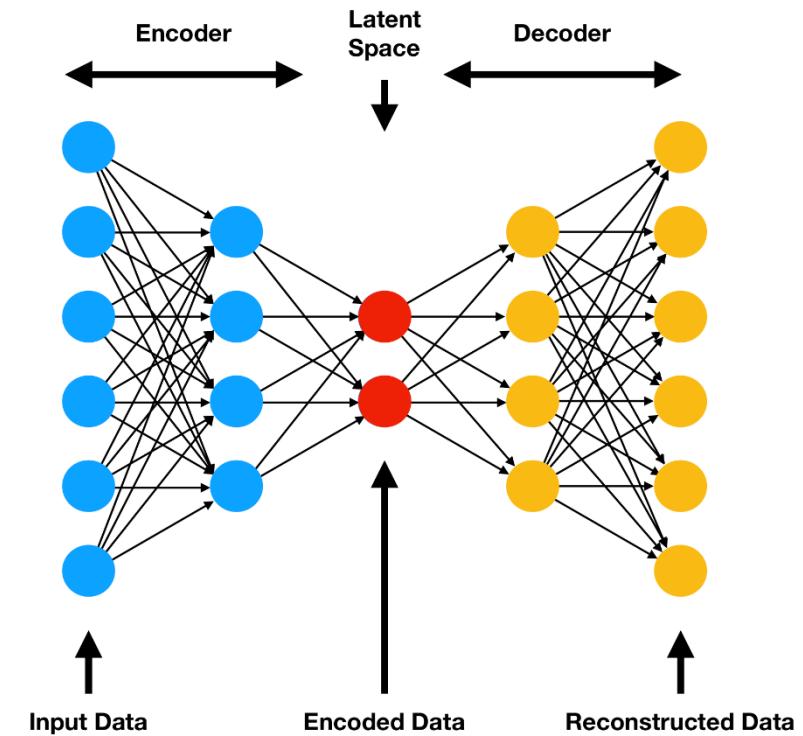$$\vec{x} := \vec{x} - \lambda \frac{\partial \mathcal{L}_{total}}{\partial \vec{x}}$$

# Embeddings, no tasks?
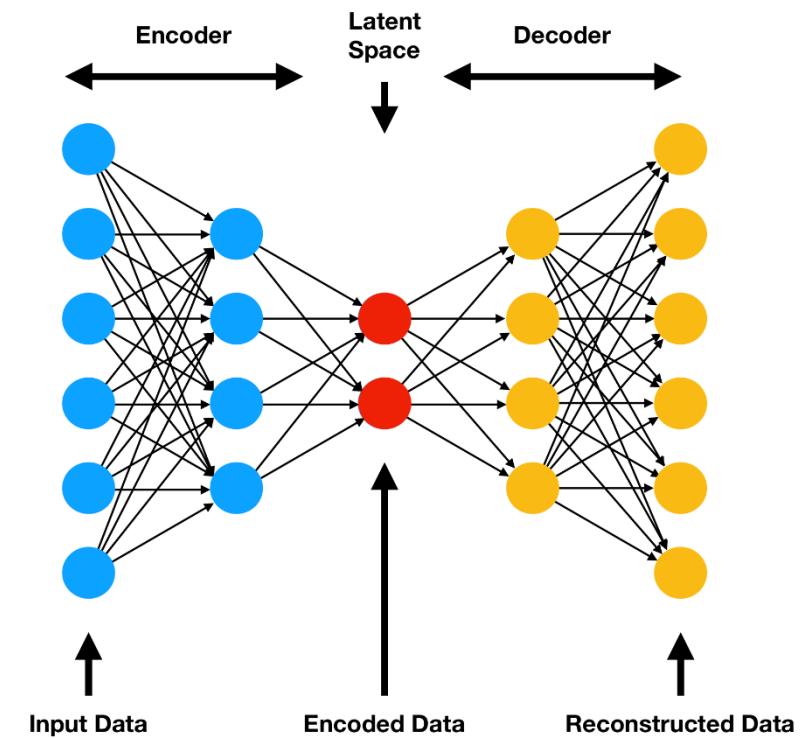
- What if do not have a task, or a ground truth?

# Embeddings with Autoencoders

- Deep learning architecture
  - Series of modules with less and less output dimensions
  - A bit like the "funnel" in CNNs
  - Module in the middle, lowest dimension
  - Second part, "funnel out"
  - Back to the original dimension of input
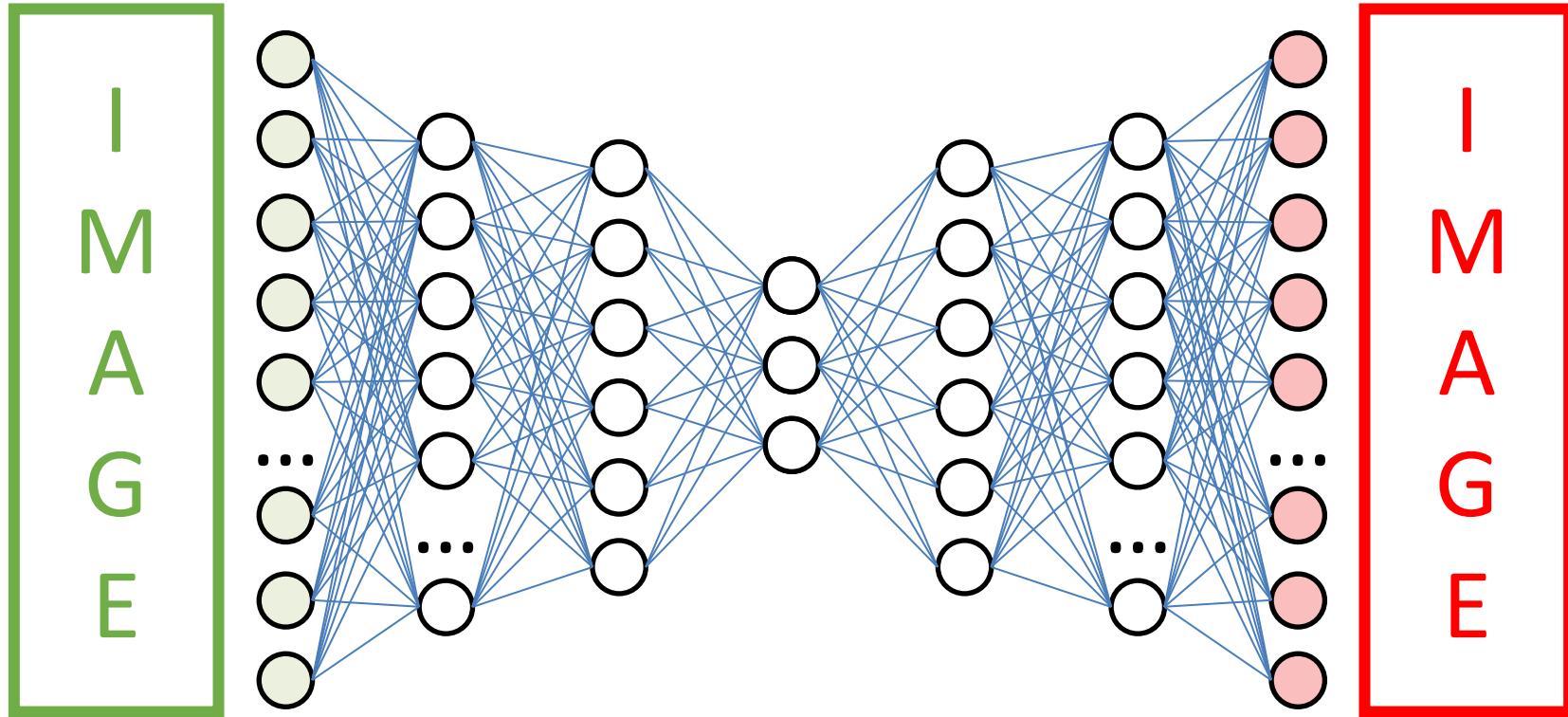- Loss function: difference input-output
  - What does it mean?
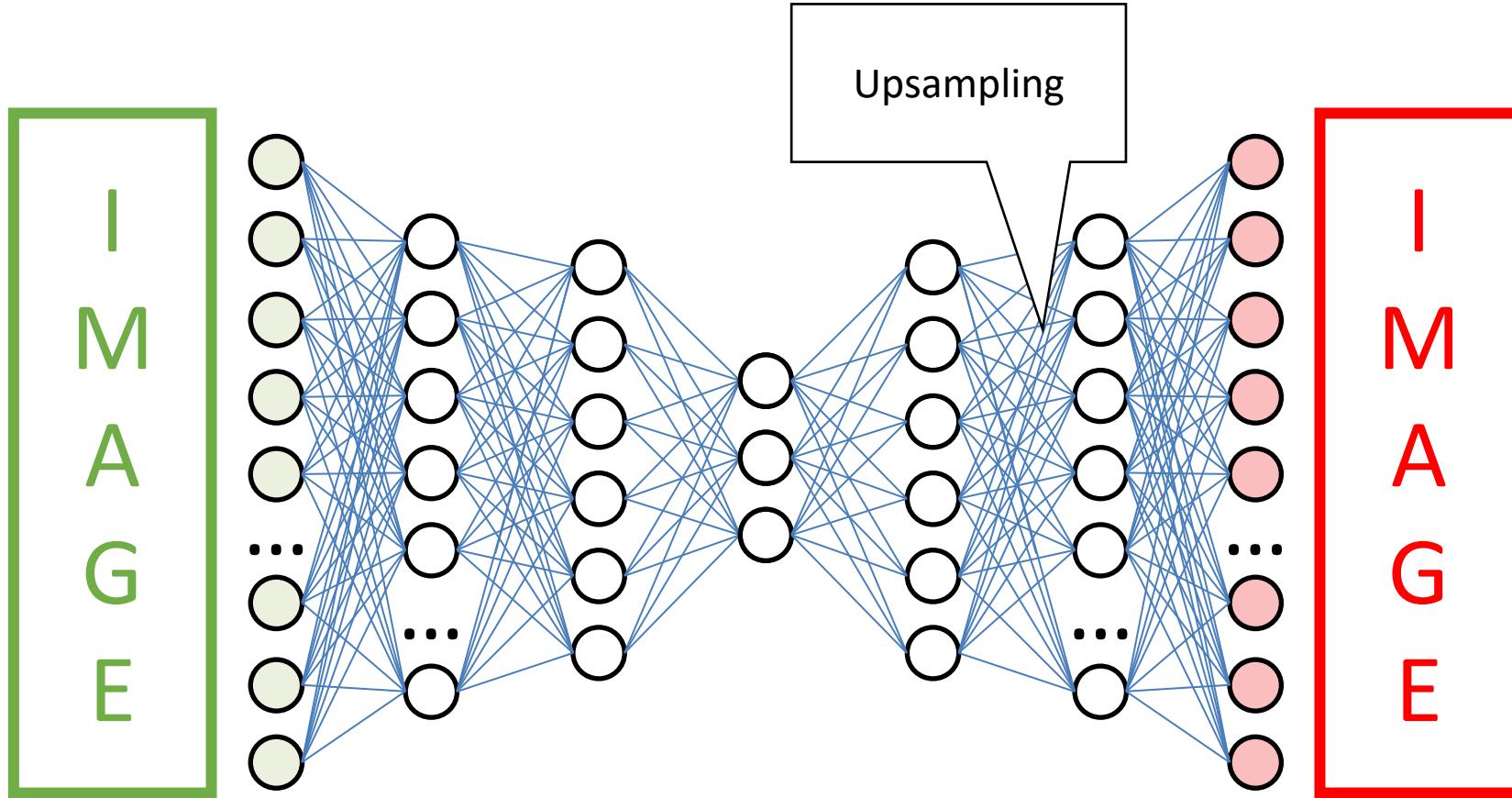
# Embeddings with Autoencoders

- Deep learning architecture
  - Series of modules with less and less output dimensions
  - A bit like the "funnel" in CNNs
  - Module in the middle, lowest dimension
  - Second part, "funnel out"
  - Back to the original dimension of input
- Loss function: difference input-output
  - Dimensionality reduction
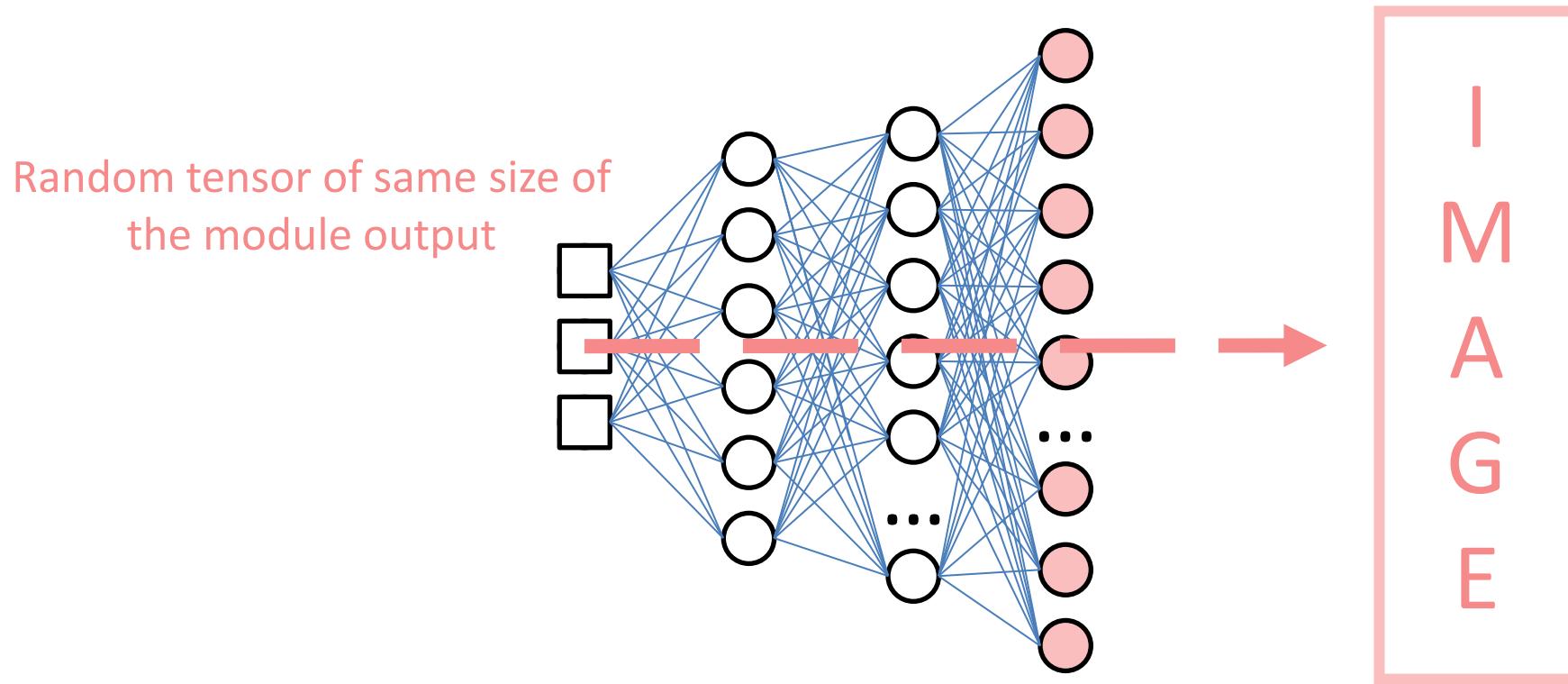  - **Meaningful**, to be able to reconstruct

# Embeddings with Autoencoders

EMBEDDINGS
Alberto TONDA, Team EKINOCS, UMR 518 MIA-PS, INRAE, Université Paris-Saclay

# Embeddings with Autoencoders

EMBEDDINGS
Alberto TONDA, Team EKINOCS, UMR 518 MIA-PS, INRAE, Université Paris-Saclay

# Embeddings with Autoencoders



Random tensor of same size of the module output

I M A G E

# Embeddings with Autoencoders

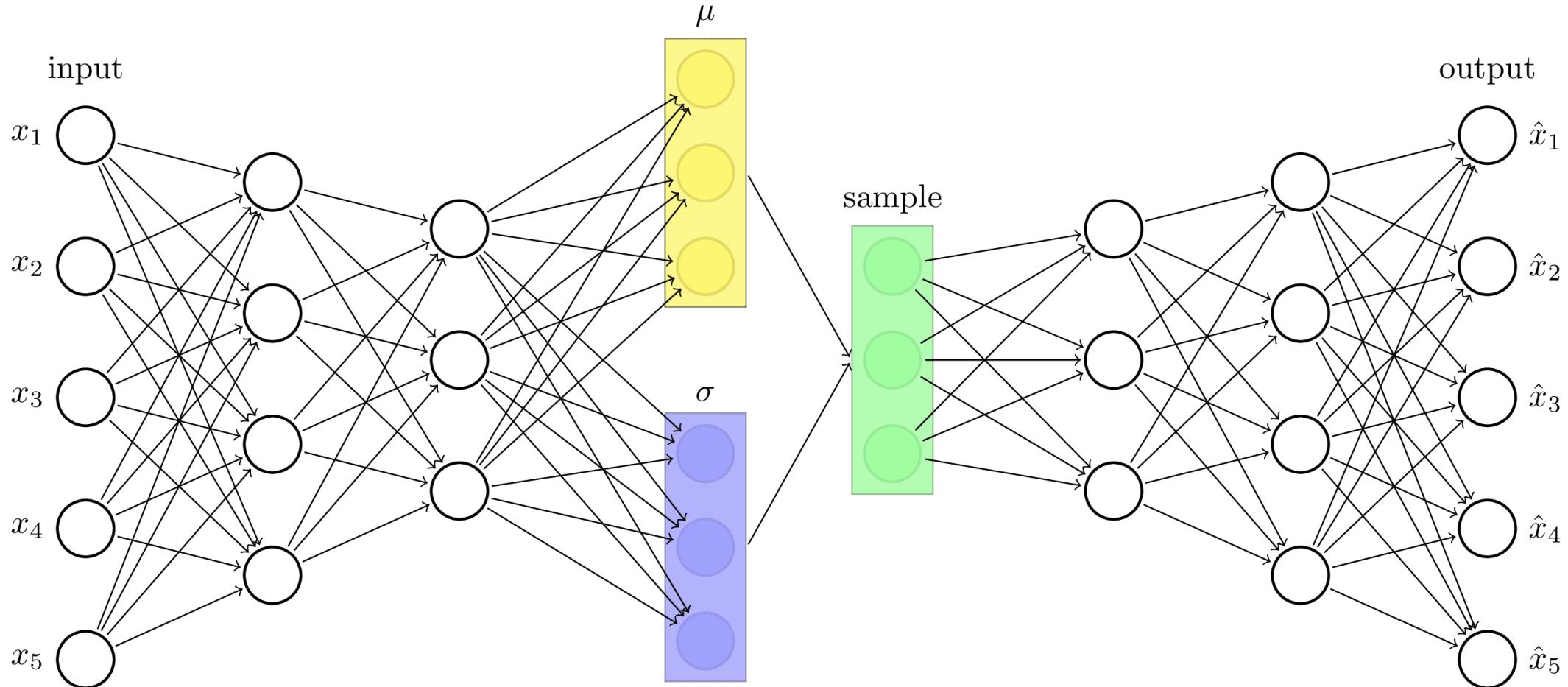- Two different points (vectors) in the same embedding...

# Embeddings with Autoencoders

- The latent space is brittle, full of "holes"

- Can we force this latent space to behave better?

- And if so, what parts of the NN would **you** act on?

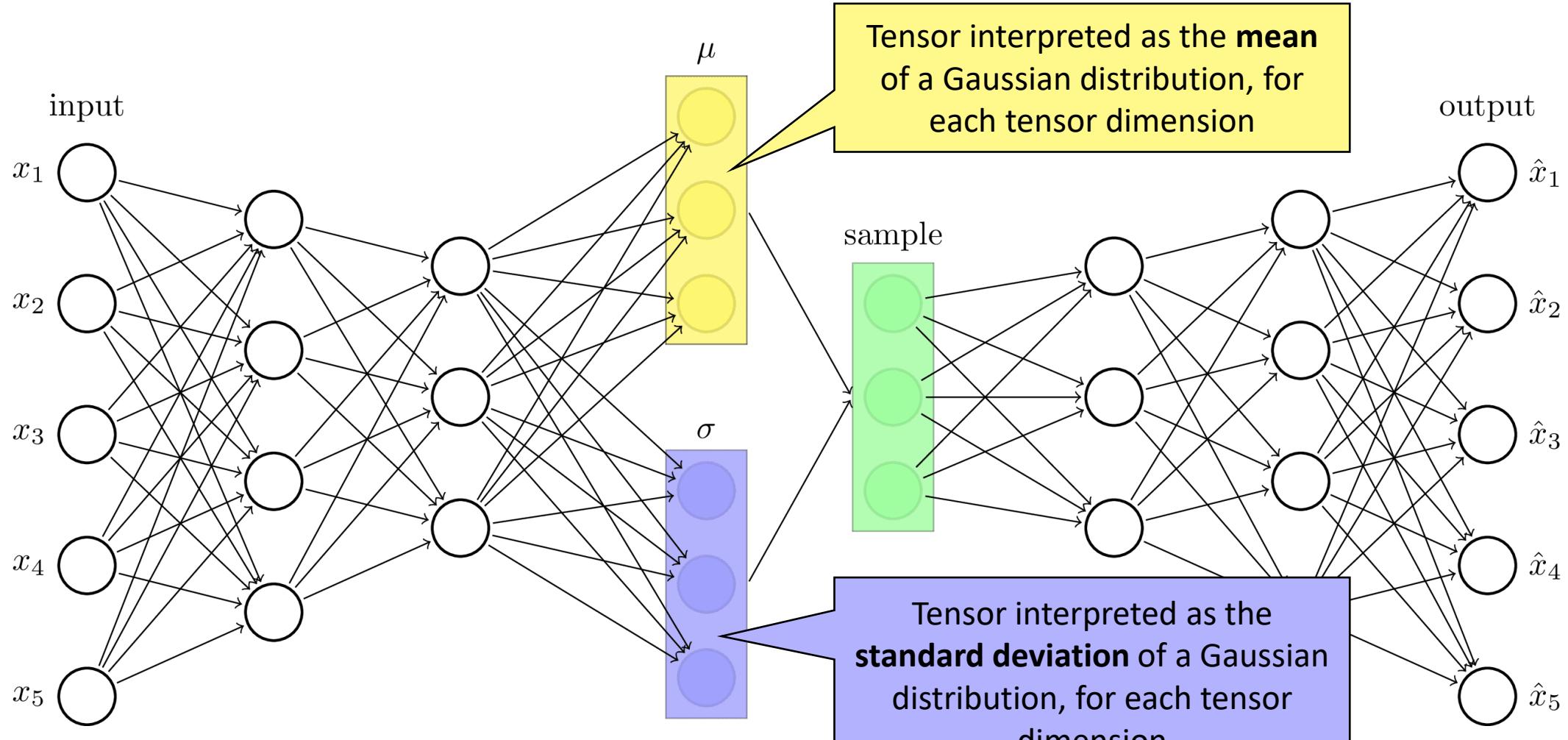Alberto TONDA, Team EKINOCS, UMR 518 MIA-PS, INRAE, Université Paris-Saclay

# Variational Autoencoders

- Each image does **not** correspond to a single point

- But rather to a (Gaussian) distribution around a point

- This makes it possible to

    - Express *uncertainty* (small or large variance)
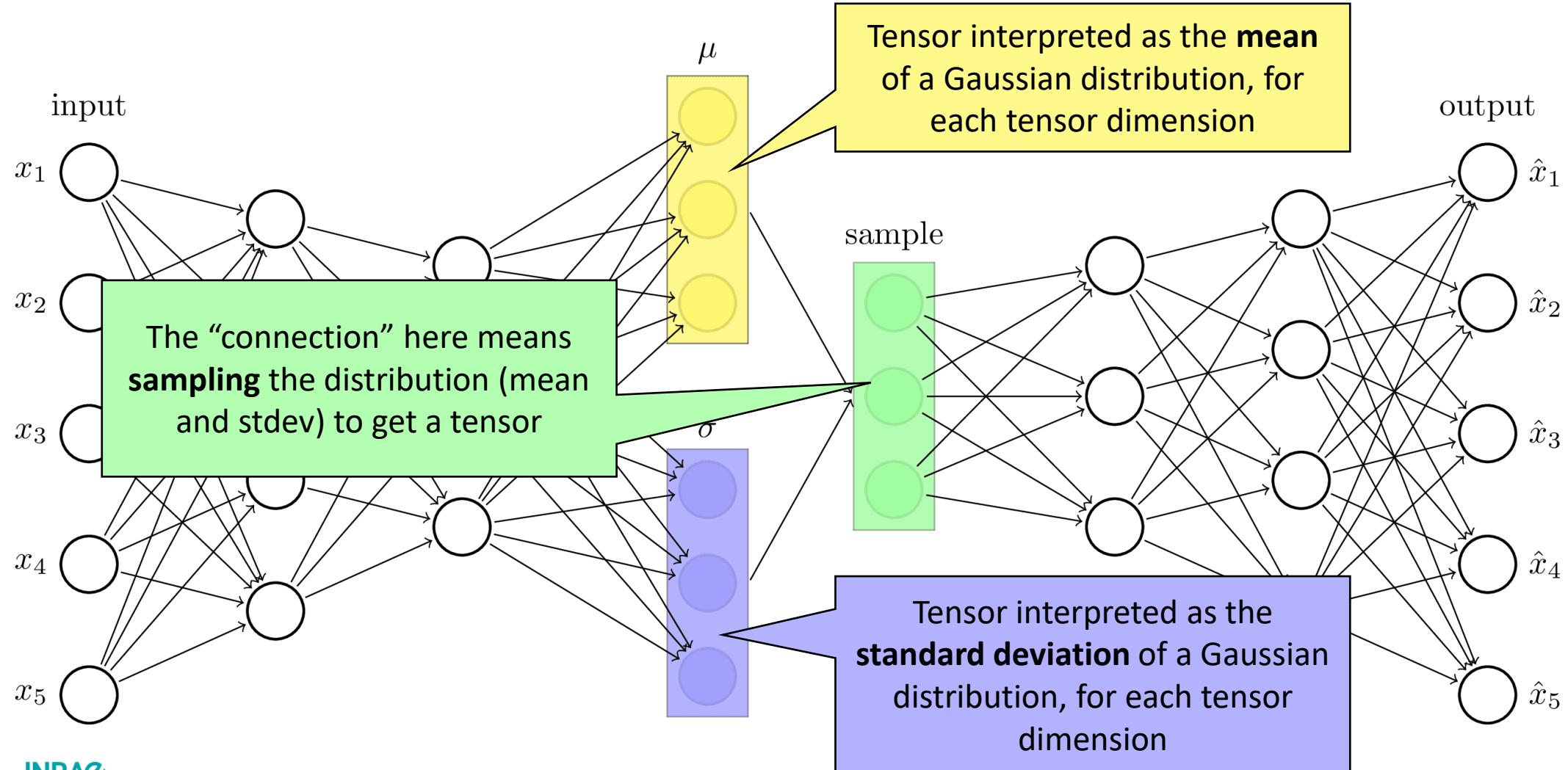    - **Enforce structure** of the latent space (less brittle)

# Variational Autoencoders

# Variational Autoencoders



input

$x_1$
$x_2$
$x_3$
$x_4$
$x_5$

$\mu$

Tensor interpreted as the **mean** of a Gaussian distribution, for each tensor dimension

$\sigma$

sample

Tensor interpreted as the **standard deviation** of a Gaussian distribution, for each tensor dimension

output

$\hat{x}_1$
$\hat{x}_2$
$\hat{x}_3$
$\hat{x}_4$
$\hat{x}_5$

Alberto TONDA, Team EKINOCS, UMR 518 MIA-PS, INRAE, Université Paris-Saclay

# Variational Autoencoders



input

$x_1$
$x_2$
$x_3$
$x_4$
$x_5$

$\mu$

Tensor interpreted as the **mean** of a Gaussian distribution, for each tensor dimension

sample

The "connection" here means **sampling** the distribution (mean and stdev) to get a tensor

$\sigma$

Tensor interpreted as the **standard deviation** of a Gaussian distribution, for each tensor dimension

output

$\hat{x}_1$
$\hat{x}_2$
$\hat{x}_3$
$\hat{x}_4$
$\hat{x}_5$

Alberto TONDA, Team EKINOCS, UMR 518 MIA-PS, INRAE, Université Paris-Saclay

# Variational Autoencoders

EMBEDDINGS

Alberto TONDA, Team EKINOCS, UMR 518 MIA-PS, INRAE, Université Paris-Saclay

# Variational Autoencoders

EMBEDDINGS
Alberto TONDA, Team EKINOCS, UMR 518 MIA-PS, INRAE, Université Paris-Saclay
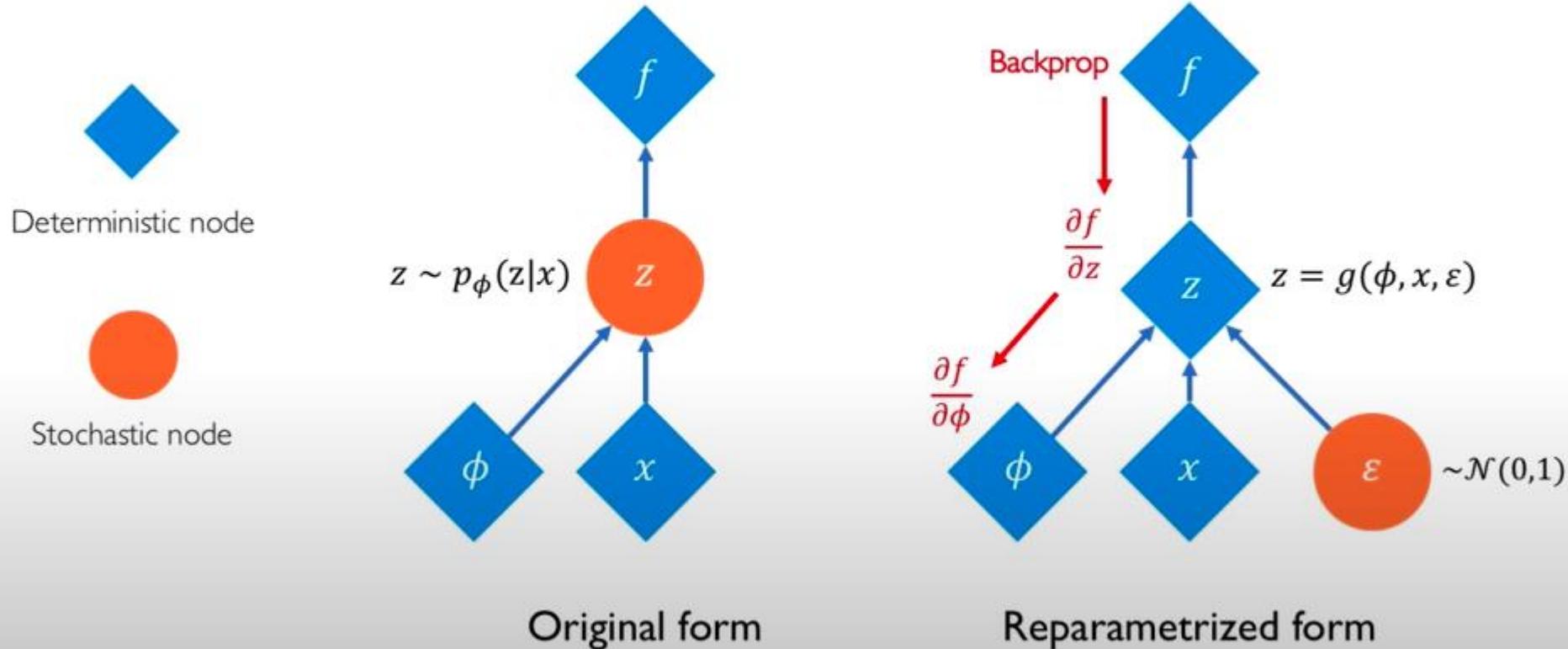
$$z \in \mathcal{N}(\mu, \sigma^2) \longrightarrow z = \mu + \sigma \odot \epsilon, \text{ where } \epsilon \in \mathcal{N}(0, 1)$$



input

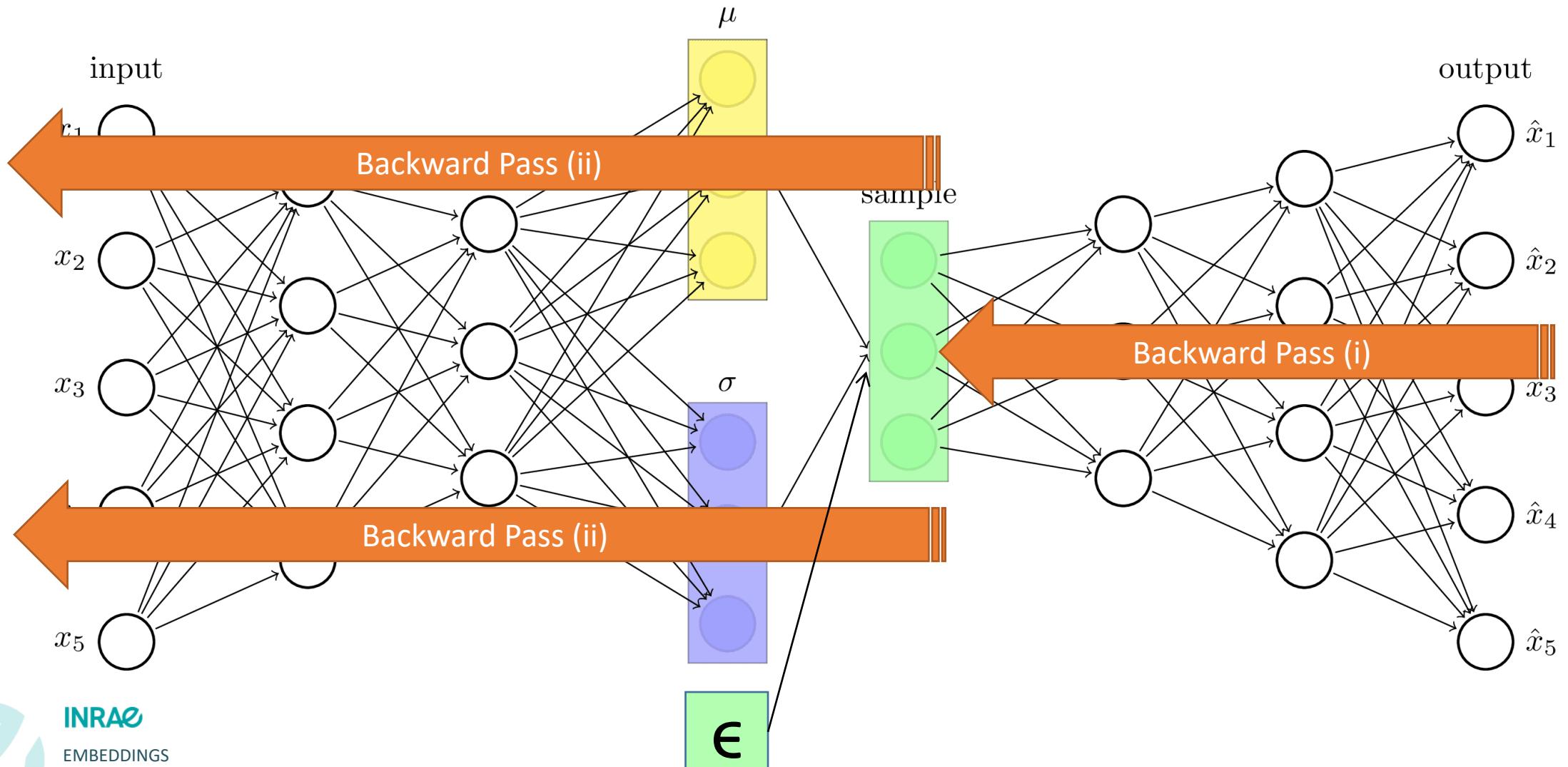$\mu$

However, there is another way of writing the sampling of a Gaussian distribution! (**reparametrization trick**)

sample

output

$x_1$

$\hat{x}_1$

$x_2$

$\hat{x}_2$

$\sigma$

Backward Pass (i)

$x_3$

$\hat{x}_3$

$x_4$

$\hat{x}_4$

$x_5$

$\hat{x}_5$

$\epsilon$

# Variational Autoencoders



Reparametrizing the sampling layer

Deterministic node

Stochastic node

$z \sim p_\phi(z|x)$

$z = g(\phi, x, \varepsilon)$

Backprop

$\frac{\partial f}{\partial z}$

$\frac{\partial f}{\partial \phi}$

$\sim \mathcal{N}(0,1)$

Original form

Reparametrized form

6.S191 Introduction to Deep Learning
introtodeeplearning.com    @MITDeepLearning

Kingma+ *ICLR* 2014. 1/28/20

Alberto TONDA, Team EKINOCS, UMR 518 MIA-PS, INRAE, Université Paris-Saclay

# Variational Autoencoders

# Variational Autoencoders

- Loss function in two parts

- The first one checks correctness of the output
    - For example, pixel-wise MSE or binary cross-entropy (BCE)
    - BCE assumes that pixels are "on" or "off" {0,1}

- The second one tries to keep distribution close to $\mathcal{N}(0, I)$
    - In other words, for each dimension Gaussian has $\mu = 0$
    - And covariance matrix is an identity matrix $I$
    - Which means that all latent dimensions are independent

- They can be weighted with a hyperparameter (β-VAE)

# Variational Autoencoders

- Evaluating the difference between two distributions is hard

- Kullback–Leibler divergence

$$D_{\mathrm{KL}}\left(P \parallel Q\right) = \sum_{x \in \mathcal{X}} P(x) \log \frac{P(x)}{Q(x)}.$$

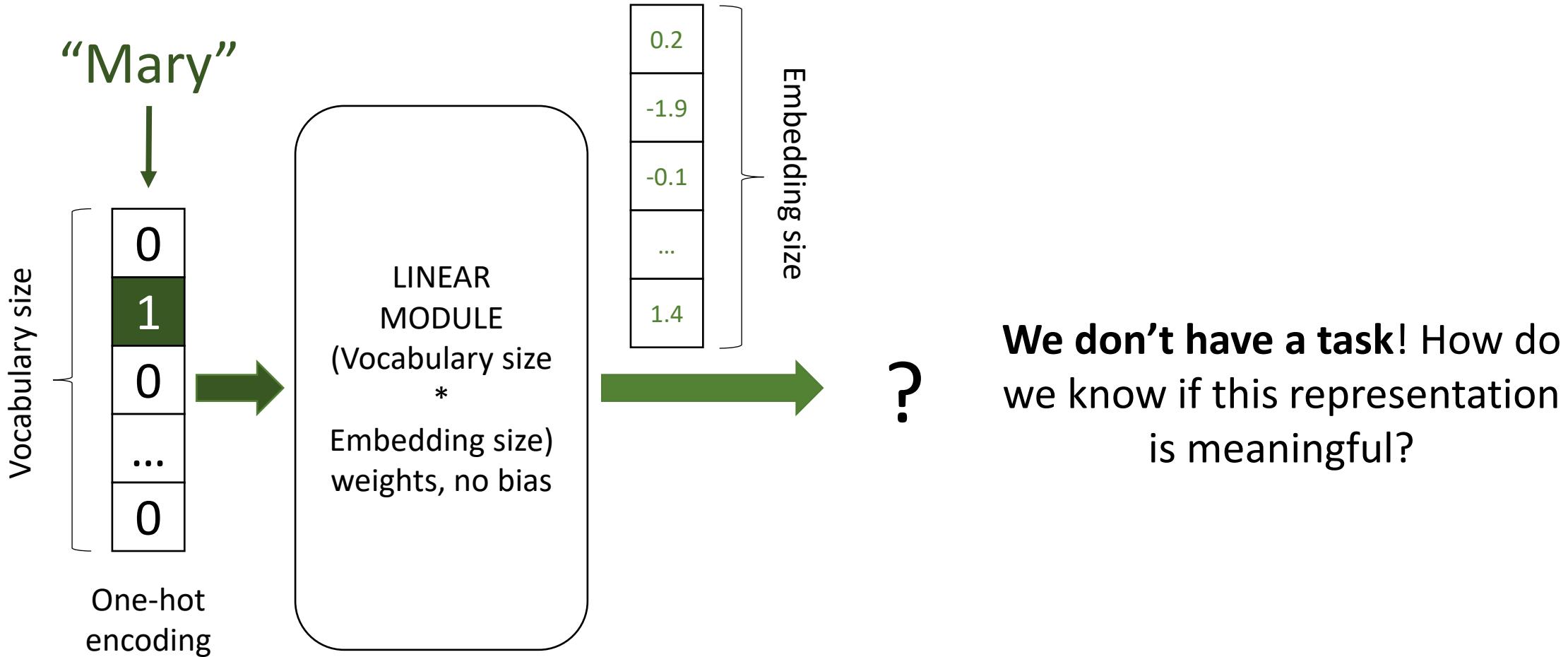- However, if both distributions are Gaussians, becomes easier

EMBEDDINGS
Alberto TONDA, Team EKINOCS, UMR 518 MIA-PS, INRAE, Université Paris-Saclay

# Embeddings of Vocabularies

- Vocabulary
  - Set of finite size, discrete elements (element: "token")
  - Classic example: words, characters, bitstrings
  - But also, e.g. possible chessboard configurations in chess

- **Words**, in particular, are difficult to manipulate for AI
  - Map to symbols, but not 1:1
  - "Chicken" vs "chicken"! "Mole" vs "mole" vs "mole"! *Context*?!?
  - Discrete: an integer per word? Not really expressive...

# Vocabulary embedding: Word2Vec

- Revolutionary idea (grandfather of ChatGPT)
  - Arbitrary number of dimensions for the embedding
  - (n_vocabulary, n_embedding)

- Simple neural network architecture
  - Input: one-hot encoding of a vocabulary
  - Output/Loss: task connected to meaning
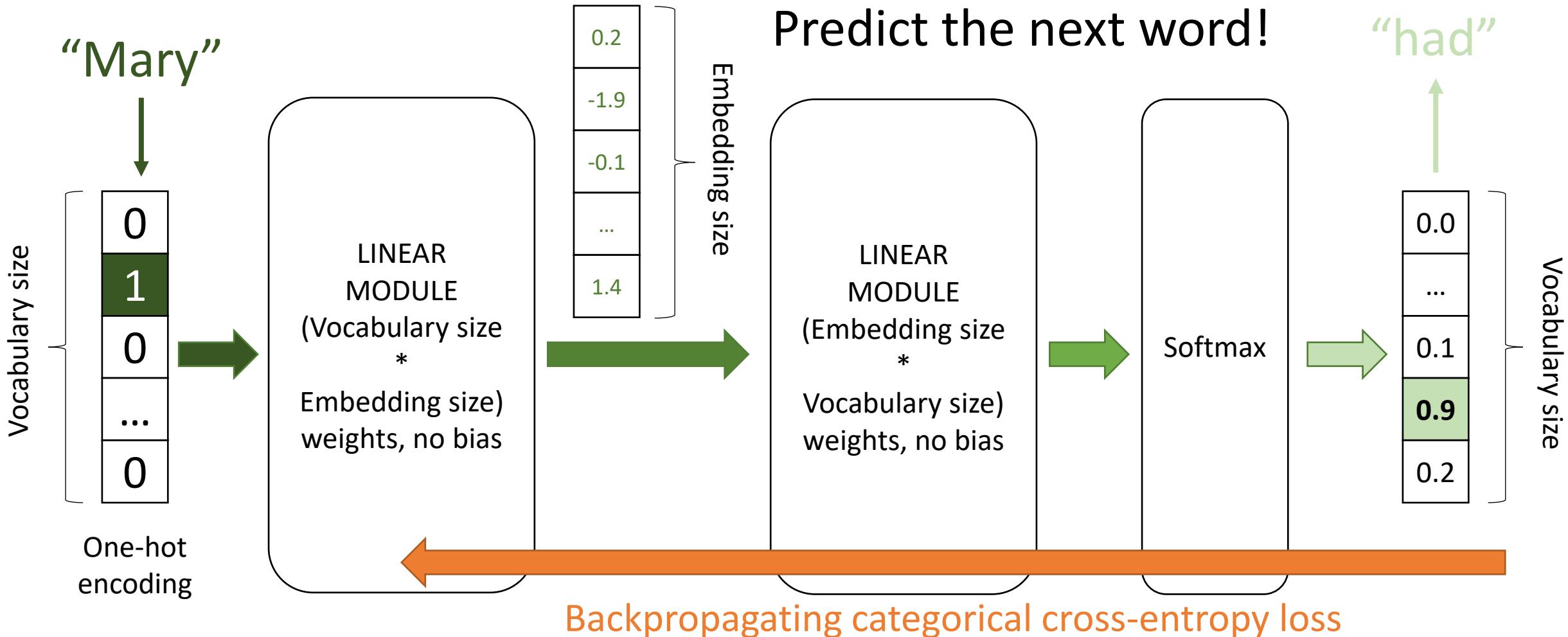  - Use **learned weights** as **word embeddings**

# Vocabulary embedding: Word2Vec

"Mary"

Vocabulary size

| |
|---|
| 0 |
| 1 |
| 0 |
| ... |
| 0 |

One-hot encoding

LINEAR MODULE
(Vocabulary size
*
Embedding size)
weights, no bias

| |
|---|
| 0.2 |
| -1.9 |
| -0.1 |
| ... |
| 1.4 |

Embedding size

?

**We don't have a task**! How do we know if this representation is meaningful?

# Vocabulary embedding: Word2Vec

- We already have lots of organized **relational data**
  - Text! Words are not randomly organized: sentences, paragraphs
  - Use text to generate meaningful training samples
  - Classic **unsupervised learning** tricks

- Slide a window on the text
  - For example, predict the next word
  - "Mary had a little lamb"
  - (Mary, had), (had, a), (a, little), (little, lamb)

# Vocabulary embedding: Word2Vec

"Mary"

Predict the next word!

"had"

One-hot encoding

Vocabulary size

| 0 |
|---|
| **1** |
| 0 |
| ... |
| 0 |

LINEAR MODULE
(Vocabulary size
*
Embedding size)
weights, no bias

Embedding size

| 0.2 |
|-----|
| -1.9 |
| -0.1 |
| ... |
| 1.4 |

LINEAR MODULE
(Embedding size
*
Vocabulary size)
weights, no bias

Softmax

Vocabulary size

| 0.0 |
|-----|
| ... |
| 0.1 |
| **0.9** |
| 0.2 |

Backpropagating categorical cross-entropy loss

# Vocabulary embeddings: Word2Vec

- Continuous Bag of Words (CBoW)
  - Predict the word between two other words
  - "Mary had a little lamb"
  - (Mary, had), (a, had); (had, a), (little, a); (a, little), (lamb, little)

# Vocabulary embeddings: Word2Vec

- Another possibility, skip-gram
  - "Context window" centered on a word, with $n$ words on each side
  - Create training samples from central word to each side word
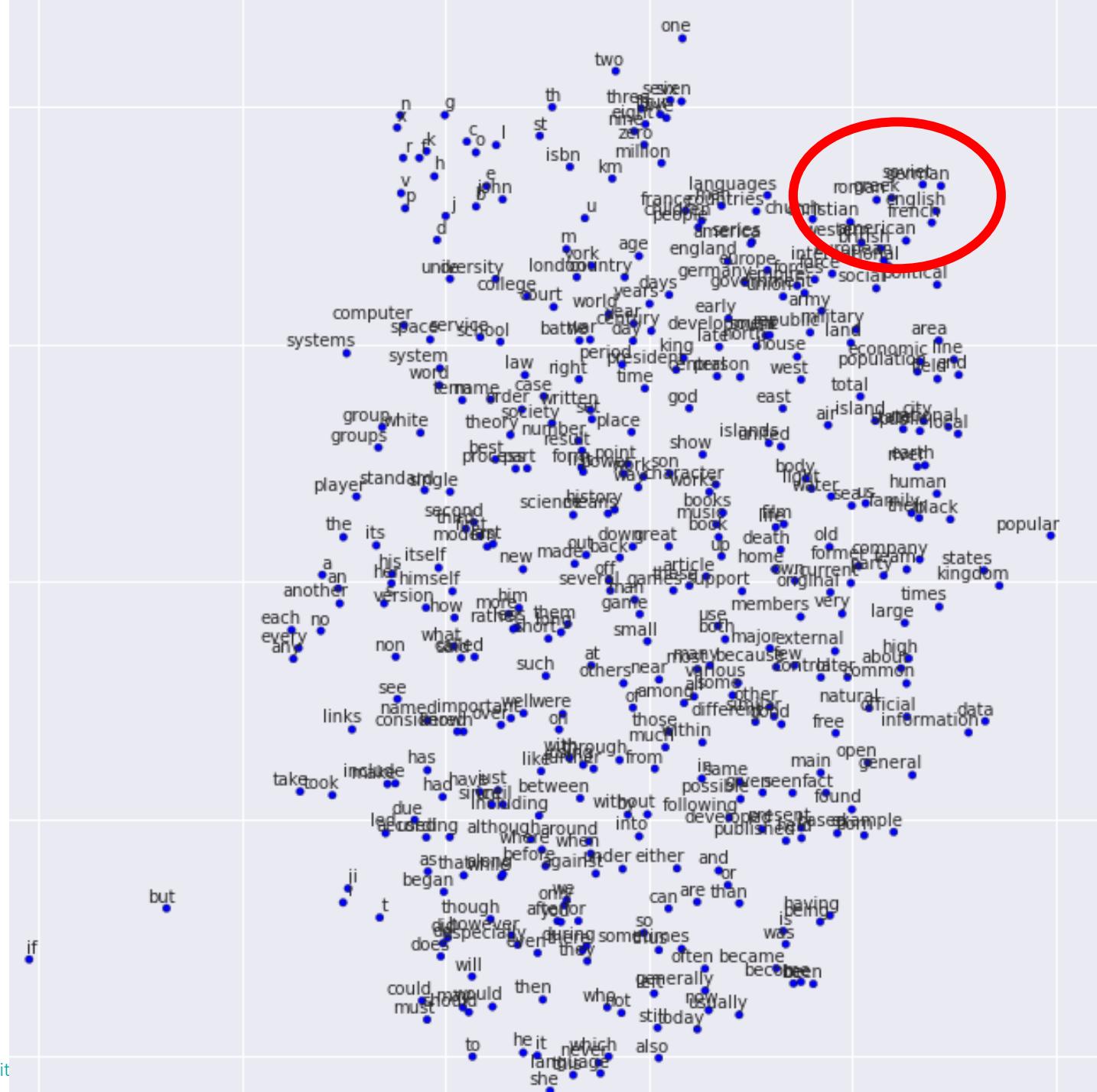
The [ wide road **shimmered** in the ] hot sun.

shimmered, wide
shimmered, road
shimmered, in
shimmered, the

- Skip-gram also uses **negative sampling**
  - For each positive sample (shimmered, wide)
  - Get *random words* that **don't appear** in context window
  - Loss function to get softmax outputs as close as possible to zero
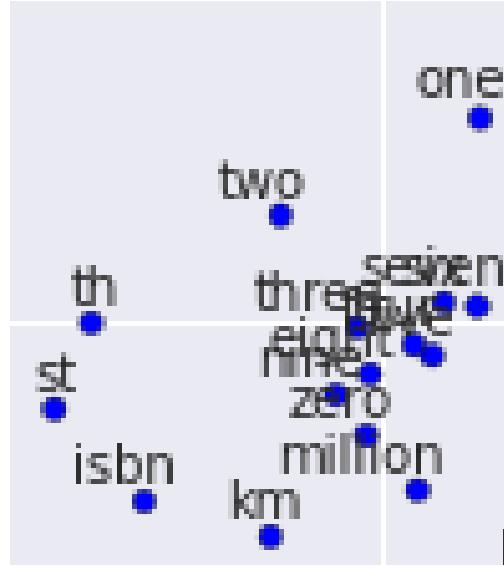
# Word2Vec

# Word2Vec

# Word2Vec



- "French", "British", "American"…
  - Adjectives for nationality!
  - Nearby, you have "languages", "countries"
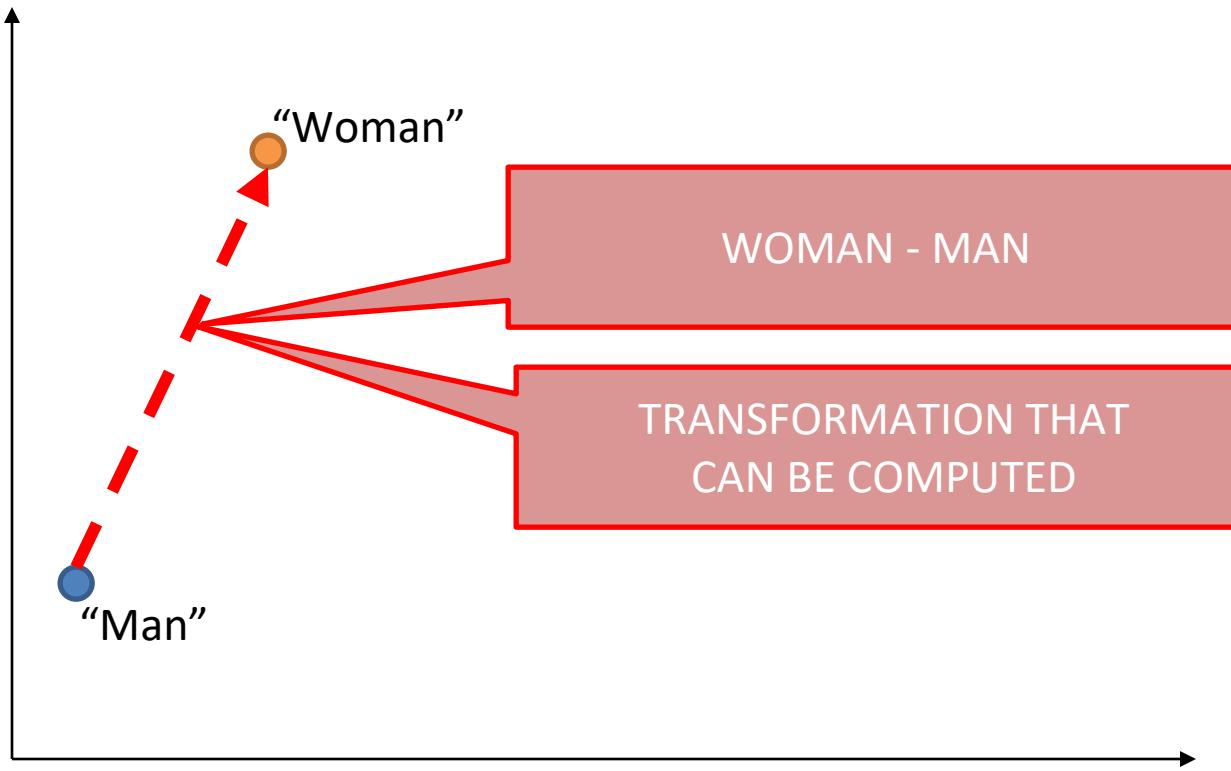  - Also, "England", "Europe", "International", …

# Word2Vec

# Word2Vec



- "one", "two", "zero", "seven", "million"…
  - Numbers, quantities; nearby, some units of measurement
  - Also "th", and "st", as in 9-th, 1-st
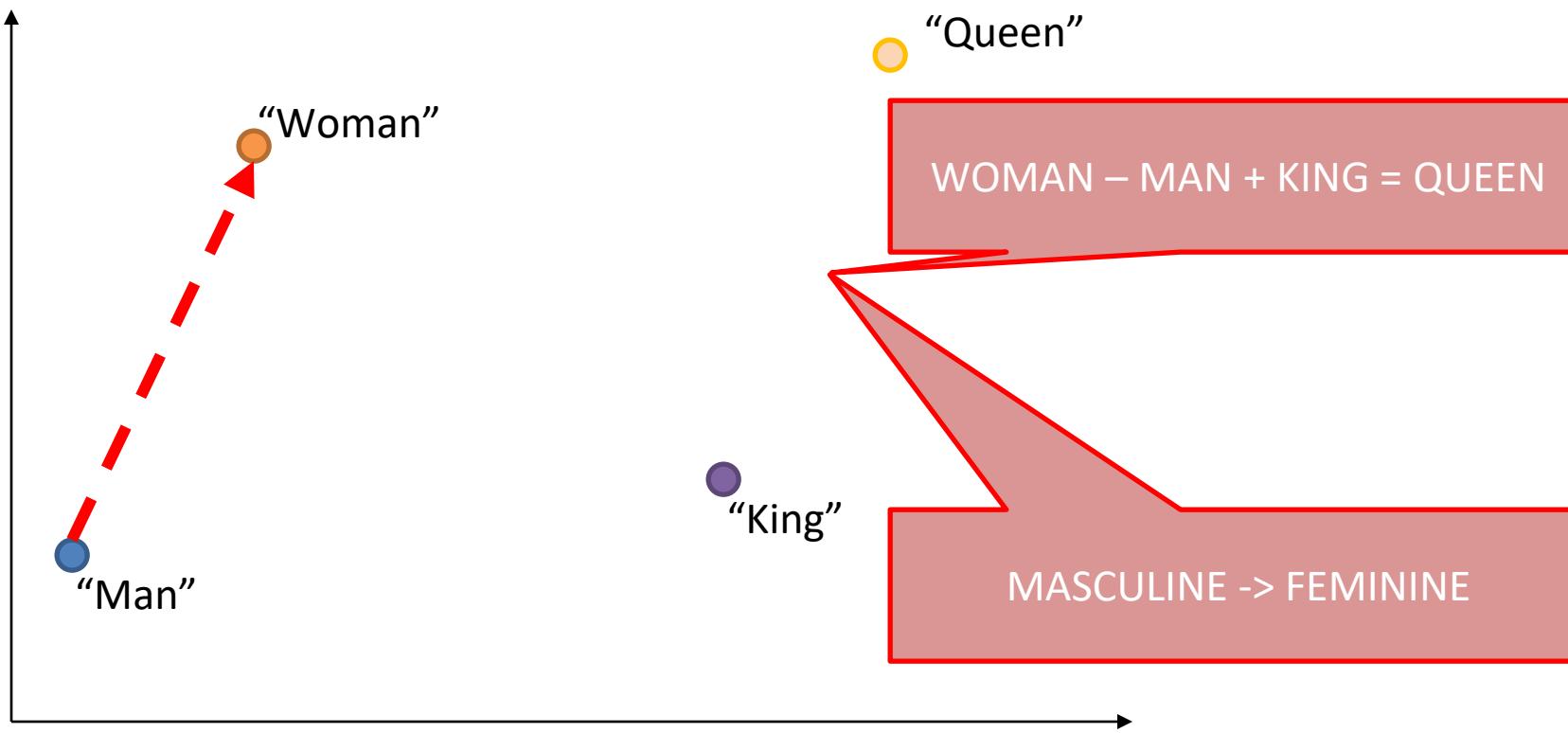  - ISBN (guess usually appears nearby numbers!)

# Word2Vec

- What is happening here?
  - Algorithm has **no semantic** info (**no meaning**)
  - But words with **similar meaning** are **close**
- Just by looking at the position of words in text
  - Similar *use* -> same *positions* with respect to other words
  - Word2Vec captures *some* aspects of meaning
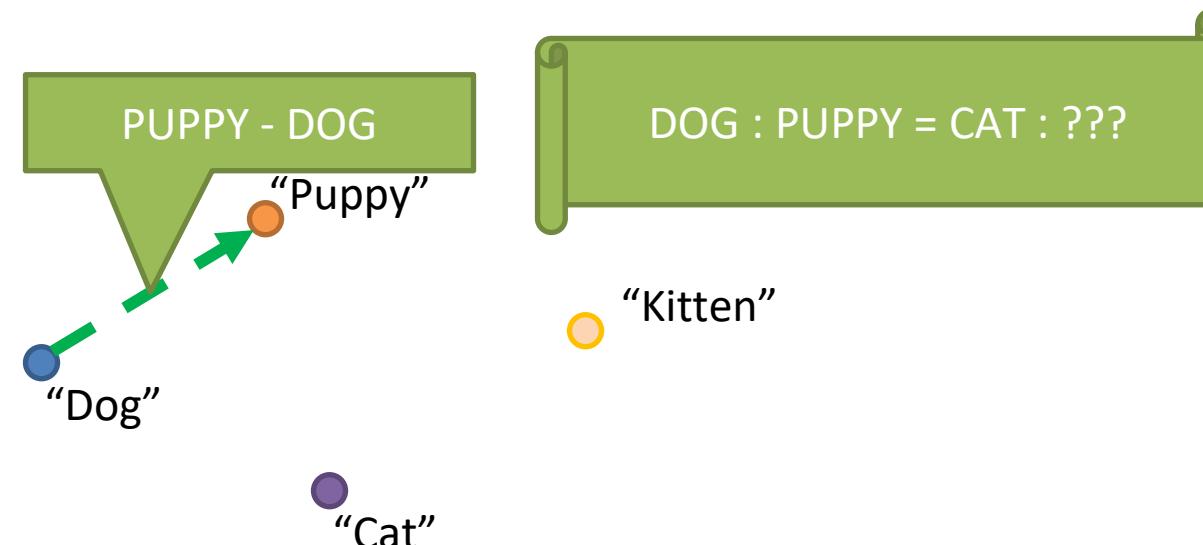- Can we do something else with Word2Vec?

# Word2Vec



"Woman"

"Man"

WOMAN - MAN

TRANSFORMATION THAT CAN BE COMPUTED

# Word2Vec

"Queen"

"Woman"

WOMAN – MAN + KING = QUEEN

"King"

MASCULINE -> FEMININE

"Man"

# Word2Vec

- Let's think about this for a second
  - No meaning inserted into the algorithm
  - Words with similar meaning cluster together
  - *Mathematical* transformations provide *meaningful* results

PUPPY - DOG

DOG : PUPPY = CAT : ???

"Puppy"

"Dog"

"Kitten"

"Cat"

# Word2Vec

- What does all this mean?
  - We are not really sure
  - Obtain *semantics* just from word positions
  - Maybe semantics is not as hard of a problem as it looks
- And Word2Vec is an incredibly *simple* architecture!
  - More complex architectures might get even better results
  - And they did (Transformers, ChatGPT, etc.)

EMBEDDINGS
Alberto TONDA, Team EKINOCS, UMR 518 MIA-PS, INRAE, Université Paris-Saclay

# Creating a Vocabulary

- Creating a Vocabulary is not straightforward
  - Elements in Vocabulary are called "tokens"
  - "Tokenization" is the process to create good tokens


- It's another optimization problem
  - "What is the set of tokens that gives me the best result?"
  - E.g., discard or keep less frequent tokens?

# Creating a Vocabulary

- Example: text from natural language?

# Creating a Vocabulary

- Example: text from natural language?

- Seems pretty intuitive, 1 word = 1 token; but nuances
  - Special tokens, <SOS> (Start of Sequence) and <EOS> (End)
  - Actually more efficient to use **parts of a word**
    - "Cod" might be better than "Code": "Cod-es", "Cod-ing", "En-cod-ing"
    - Kind like mapping to word roots, even if not exactly

- What happens if input is "**Grzegorz Brzęczyszczykiewicz**"?

# Creating a Vocabulary

- In practice, solved by declaring **character-level tokens**
  - Tokenization attempts on largest possible matching tokens
  - If it fails, fall back on smaller and smaller tokens

- Test GPT-4 tokenization, platform.openai.com/tokenizer

# Creating a Vocabulary

EMBEDDINGS

Alberto TONDA, Team EKINOCS, UMR 518 MIA-PS, INRAE, Université Paris-Saclay

# Questions?

Bibliography
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). *Efficient estimation of word representations in vector space*. arXiv preprint arXiv:1301.3781.

Images and videos: unless otherwise stated, I stole them from the Internet. I hope they are not copyrighted, or that their use falls under the Fair Use clause, and if not, I am sorry. Please don't sue me.