

• What if  $Z^{(l+2)}$  and  $a^{(l)}$  does not have the same dimension?

- First, many convolutions in ResNets are "same conv" (preserves dimension)
- if dimensions are different, you can add an extra matrix  $W_s$  to change the dimension of  $a^{(l)}$  →
  - can be learned
  - or be fixed

→ See the paper (2015) by He et al.

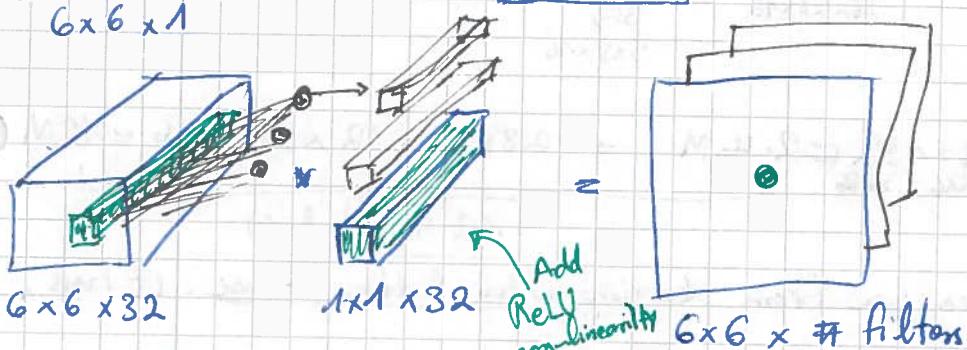
Networks in Networks :  $1 \times 1$  convolutions.

{ • Adds non-linearity  
• Helps decrease # channels if needed.

$$\begin{bmatrix} 1 & 2 & 3 & 4 & 5 \\ 6 & - & - & - & - \\ - & - & - & - & - \\ - & - & - & - & - \end{bmatrix} \quad * \quad \begin{bmatrix} 2 & 4 & 6 & 8 & 10 \end{bmatrix} = \begin{bmatrix} \text{...} \end{bmatrix}$$

$6 \times 6 \times 1$

• nothing interesting for a  $(6 \times 6 \times 1)$  one ~~channel~~ channel image



• now it becomes interesting  
• we can now think of it as single neuron taking as input 32 numbers from different channels, and add a ReLU non-linearity to it

• It becomes even more interesting if you have more filters.

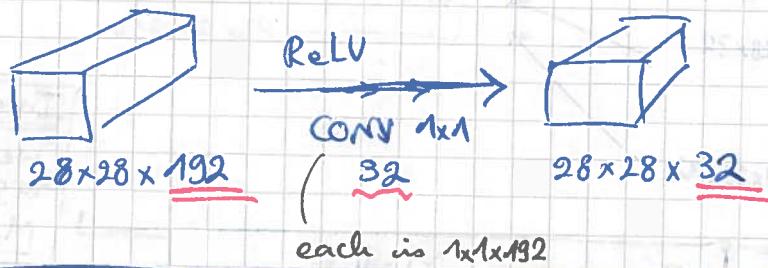
• Called "Network in Network"

• Can carry out a pretty non-trivial computation on the input volume.

• It is like a Fully connected layer at a specific position  $(h, w)$  of your input volume.

• Is basis for many interesting innovations including the Inception.

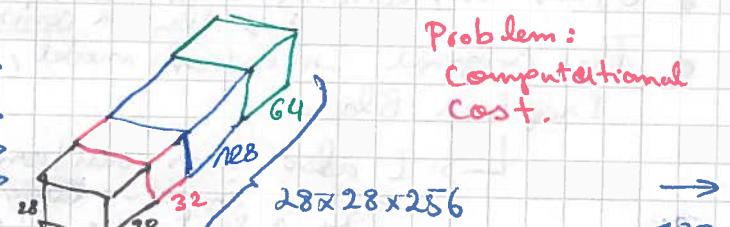
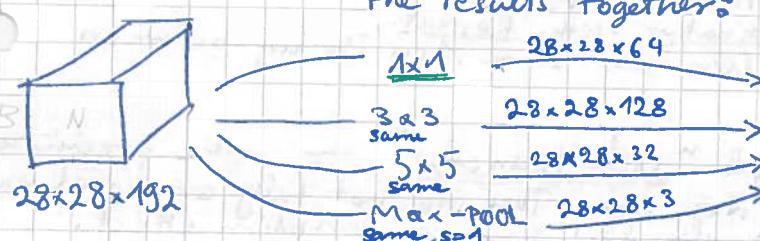
• It can also be used to "shrink number of channels"



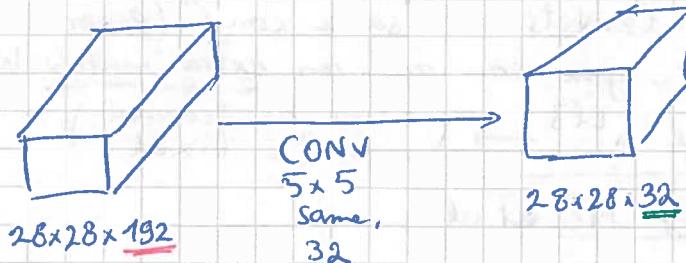
Networks  
In  
Network

Motivation for Inception: (2014)

Inception layer: Apply all the possible different filters in one layer and stack up the results together:



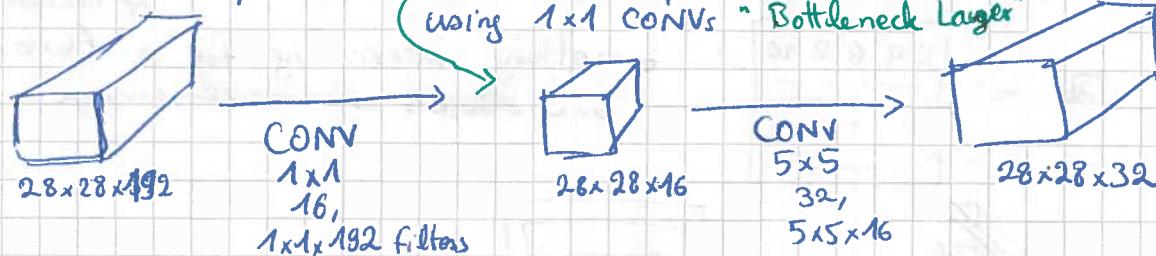
The problem of computational cost: Let's focus on the  $5 \times 5$  CONV:



- 32 filters of size  $5 \times 5 \times 192$
- out put is  $28 \times 28 \times 32$  numbers for each we need to do  $5 \times 5 \times 192$  multiplications:  
⇒ total multiplications:

$$28 \times 28 \times 32 \times 5 \times 5 \times 192 = 120 \text{ M}$$

So, let's try this: Shrink the representation first using  $1 \times 1$  CONVs "Bottleneck Layer"

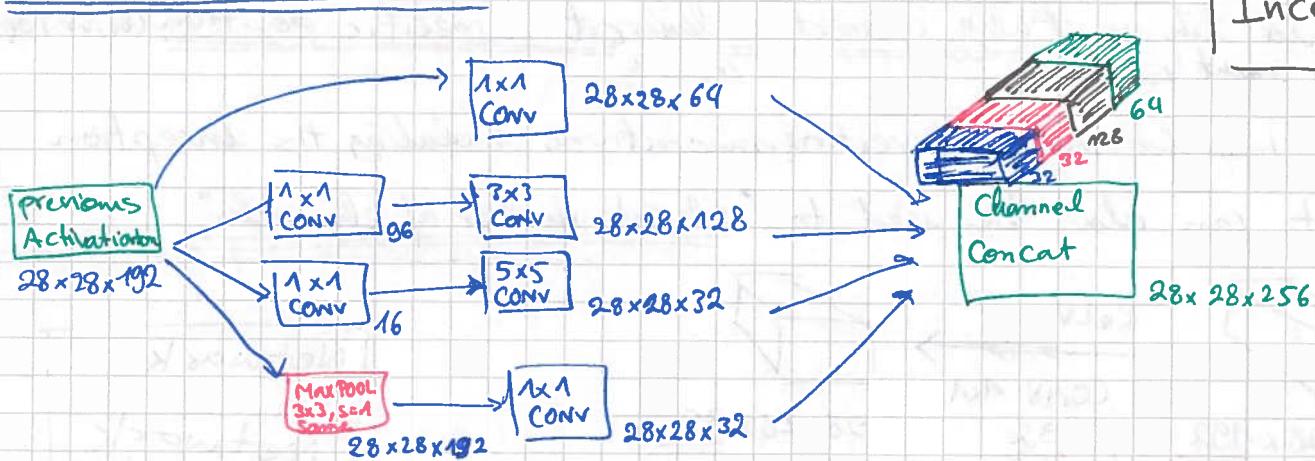


# total Multiplications :  $28 \times 28 \times 192 \times 16 \approx 2.4 \text{ M} + 28 \times 28 \times 32 \times 5 \times 5 \times 16 \approx 10 \text{ M}$

$12.4 \text{ M}$

- Inception module frees you from deciding what filter to use. It says, let's do them all.
- To reduce the comp. cost, use the "Bottleneck layer" to shrink down the representation.
- Bottleneck layers do not hurt the performance as long as used reasonably.

## Inception module :



Inception

One Inception Module

Inception  
"We need to go deep"  
Meme

- An Inception Network is stacking up multiple of these modules together. (see Inception in Video) → O-O-O-O-O
- Other variations exist, Inception V2, V3, ↗ Also a combination with ResNet.
- The original Inception model, has also some Max Pooling between Inception Blocks.  
↳ I also draw additional "side branches" to make predictions earlier in the Network. They are just fully connected layers with a softmax. ↗ Inception effect.

Google Le Net

# Practical Advices for using ConvNets

- Start with: Using Open-source Implementations: find Author's implementations to get going faster.
- ↳ Pre trained networks.
  - ↳ Use transfer learning.

## Transfer Learning:

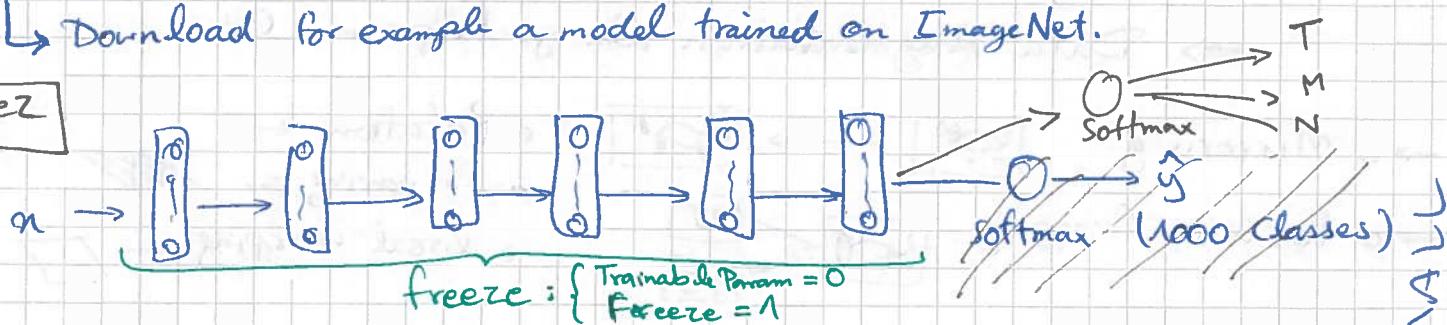
### Transfer Learning

- Download pre-trained weights.  
↳ Data sets: ImageNet, COCO, Pascal VOC

- Imagine you want to train a classifier for you two cats neither  
But you don't have a large dataset of their images: Tiger Misty

↳ Download for example a model trained on ImageNet.

### Freeze



- remove the last softmax unit in original network and train your own Softmax unit with 3 outputs.  
↳ freeze the weights in previous layers so that they don't get retrained.  
↳ All DL Frameworks support this.

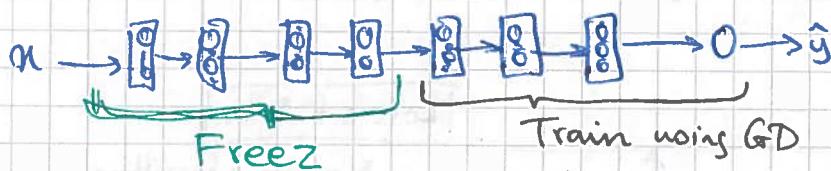
- Another Option: Consider previous layers a frozen fixed function that takes an input  $x$  and maps it to some activation

### Fixed Function

↳ Save that activation as a feature vector  
↳ train a new shallow softmax model to make the prediction.

- If you have a larger training set

↳ Freeze fewer layers, and train the rest.



↳ or even remove them and add your own hidden units.

The more data you have, the more layers you can retain

- If you have a lot of data use the pre trained weights as initialization and retrain the whole network.

Medium

High

Low

Because dataset on the internet are so big and existing networks have spent weeks on training, Almost Always.

Better off, when you start with a pre-trained Model!

- Transfer Learning, always the first option, unless you have a exceptionally very large dataset and huge computational budget.

## Data Augmentation:

- CV involves learning pretty decent complicated functions.  
↳ In order to learn these very complicated functions, we ~~can't~~ always benefit from more data. We can't have enough data.  
⇒ Data Augmentation always helps in CV

→ o Mirroring 

o Rotation ↵

→ o Random Cropping 

o Shearing ↵

o local warping ↵

→ o Color shifting : R+20, G-20, B+20 (Distort RGB channels)

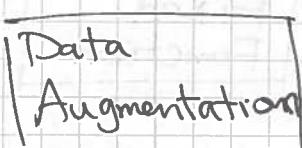
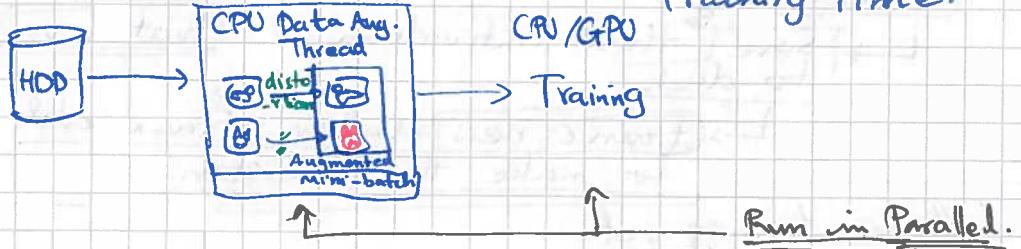
↳ In practice the values for R, G, B are drawn from a probability dist.

↳ Makes the learning more Robust to changing of color in images

→ PCA color augmentation : Change the dominant color channels more.

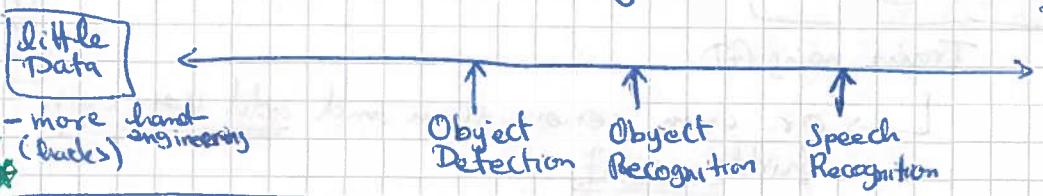
↳ AlexNet paper.

o Online Data Augmentation : Do data augmentation right at the training time.



## State of Computer Vision : Some special aspects of DL in CV

- Data vs. hand-engineering :



## Two Sources of Knowledge

↳ labeled data : (x,y) pairs

↳ Hand engineered features/network architecture/other components.

Lots of Data

- simpler algorithms  
- less hand engineering

Transfer learning helps a lot when you have less data

## Tips for doing well on benchmarks / winning competitions:

- doing well on a benchmark normally increases the chance of the paper getting accepted.
  - But the tricks are most of the time not very useful for a real system in production environment.
- ① Ensembling :
- Train [Needs much memory] →  $y$
  - several networks independently and average their output.
  - 3-15 networks are typical. Gains of up to 2%.
- ② Multi-Crop at test time :
- Run classifier on multiple versions of the test images and average the results. (e.g. 10-crop)
  - 4 corners, 1 center = 10 crops
  - flip

## Course 4 - CNNs - Week 2 - Programming 1 - Keras

- ③ Keras : high level NN API capable of running on top of several lower level framework like TF and CNTK

```
def model(input_shape):
```

```
X_input = Input(input_shape) # input placeholder.
```

```
X = ZeroPadding2D((3,3))(X_input)
```

```
# CONV → BN → ReLU Block applied to X
```

```
X = Conv2D(32, (7,7), strides=(1,1), name='conv0')(X)
```

```
X = BatchNormalization(axis=3, name='bn0')(X)
```

```
X = Activation('relu')(X)
```

```
# MAXPOOL
```

```
X = MaxPooling2D((2,2), name='max_pool')(X)
```

```
# Flatten X + Fully Connected
```

```
X = Flatten()(X)
```

```
X = Dense(1, activation='sigmoid', name='fc')(X)
```

```
# Create Model instance to train/test
```

```
model = Model(inputs=X_input, outputs=X, name='HappyModel')
```

```
return model
```



Very cool :)

1. Create Model

2. Compile model : `model.compile(optimizer="--", loss="--", metrics=["accuracy"])`

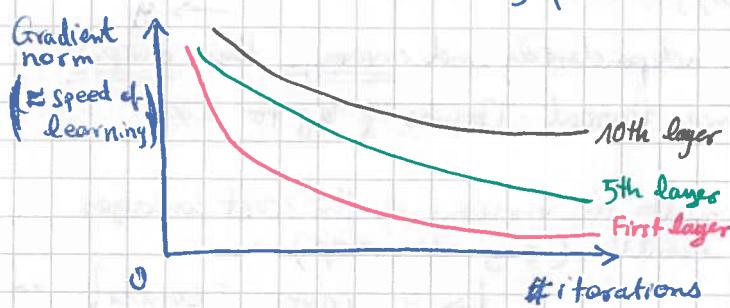
3. Train the model : `model.fit(x=..., y=..., epochs=..., batch_size=...)`

4. Test the model : `model.evaluate(x=..., y=...)`

`model.summary()`, `model.plot()`

## Course 4 - CNNs - Week 2 - Programming 2 - ResNets

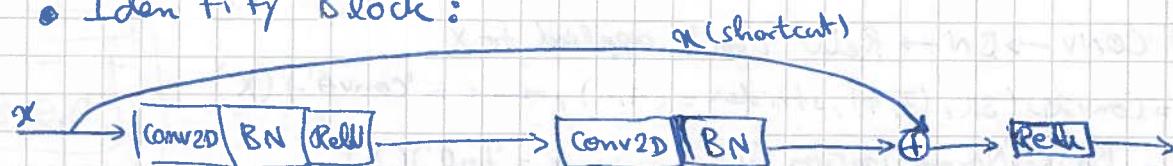
- Very deep networks are hard to train due to vanishing or exploding gradients.
- During training we might actually see the magnitude (or norm) of the gradients for earlier layers decreases to zero very rapidly as the training proceeds.



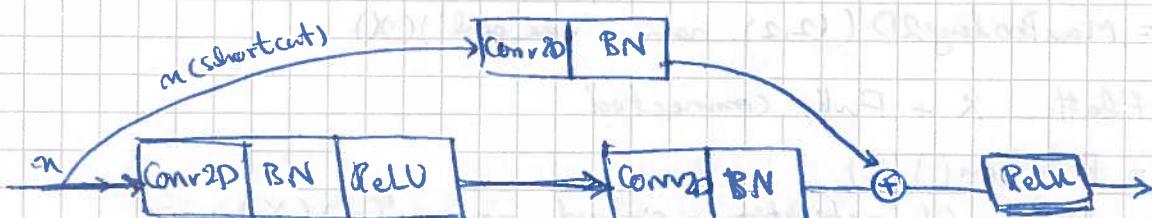
Speed of learning decreases very rapidly for the earlier layers as the network trains.

- The ResNet "shortcut" or "skip connection" allows the gradient to be directly backpropagated to earlier layers.
- Because of the skip connection they can also very easily learn the Identity function.
- Depending on whether input/output dimensions are the same, two main types of blocks are used in a ResNet:

- Identity Block:



- Convolutional Block : when input and output dimensions are different



the conv2d block in shortcut path resizes the input  $n$  to a different dimension.

## Course 5 - RNNs - Week 2

1.4.4 in slide  
2.4 pages

Natural language processing and Word embeddings



- NLP with DL: An important combination.
- Word vector representations, Embedding layers
- Example Applications: sentiment analysis, named entity recognition and machine translation.

## Word Representations:

- DL has revolutionized NLP
- Word embedding: A way of representing words that lets your algorithm automatically learn analogies.

So far: We had a dictionary (vocabulary) with one-hot word representations:

$$V = [a, aaron, \dots, zulu, \text{<UNK>}] \quad |V| = 10,000$$

Man (5391)	Woman (9853)	King (4914)	Queen (7175)	Apple (456)	Orange (6257)
$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$
$\xrightarrow{\quad}$ 0 <sub>5391</sub>	$\xrightarrow{\quad}$ 0 <sub>9853</sub>				

I want a glass  
of orange juice.

I want a glass  
of apple juice.

- Problems with dictionary and one-hot representation:

- no generalization based on the concepts.  
for example, the relationship between orange and apple is no different from relationship between orange and King.
- Because the inner product of any two different one-hot vectors is 0, and euclidean distance between only pairs of these vectors is also the same.

!  $\Rightarrow$  This representation does not allow the algorithm to know that somehow "apple and orange" are more similar than "King and orange".

$\Rightarrow$  Featurized Representation: word embedding

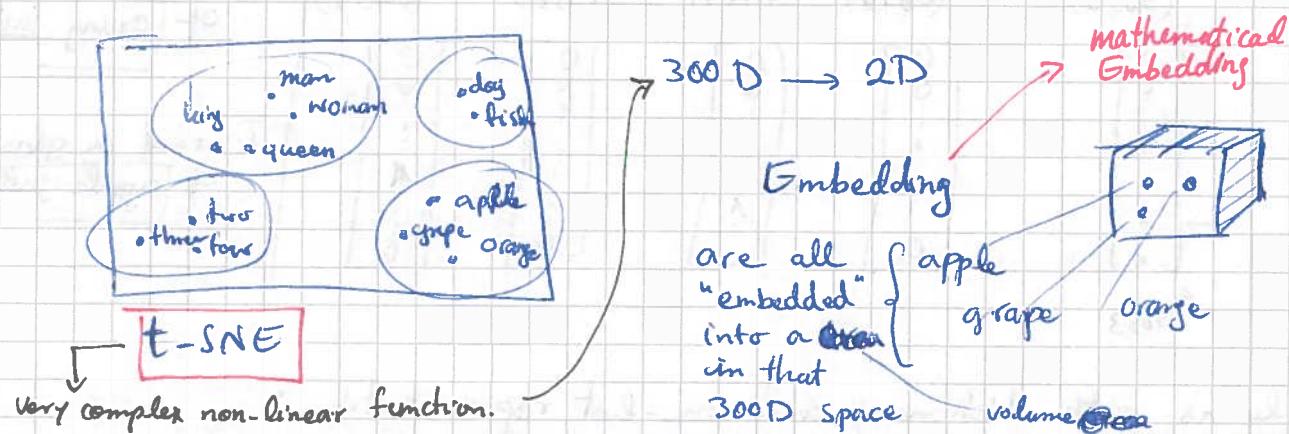
	Man (5391)	Woman (9853)	King (4914)	Queen (7175)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.04	0.01	0.02	0.01	0.95	0.97
Size	1	1	1	1	1	1
Cost	1	1	1	1	1	1
alive	1	1	1	1	1	1
Verb	1	1	1	1	1	1
e.g. 300 features	1	1	1	1	1	1
BRUNNEN	1	1	1	1	1	1
e <sub>5391</sub>	e <sub>9853</sub>					

NOW: becomes easy:

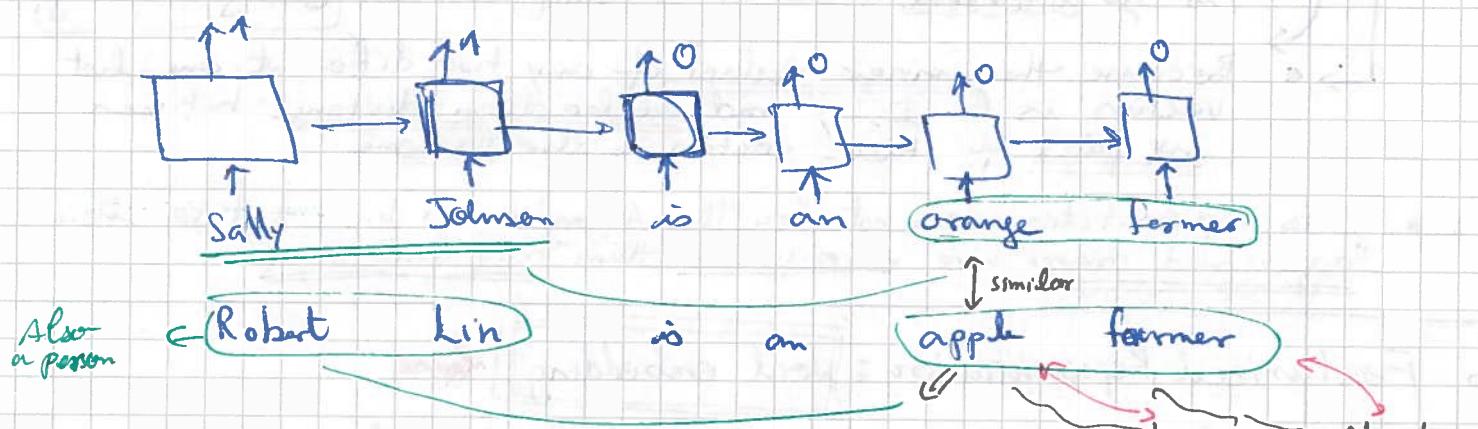
- I want a glass of orange juice.
- I want a glass of apple juice.

- We want to learn word embeddings; That is, learn high dimensional feature vectors that give a better representation than one-hot representation
- Of course the features of a WE won't be as clear interpretation like "gender", "age", "royal". What they exactly represent might be difficult to interpret, but still they capture very well the similarities like between apple and orange.

## Visualizing Word Embeddings:



## Using WEs: Sample: Named Entity Recognition



- This even works for words we haven't seen in our training set as long as we have WEs that capture the similarities of those words with words we already know.

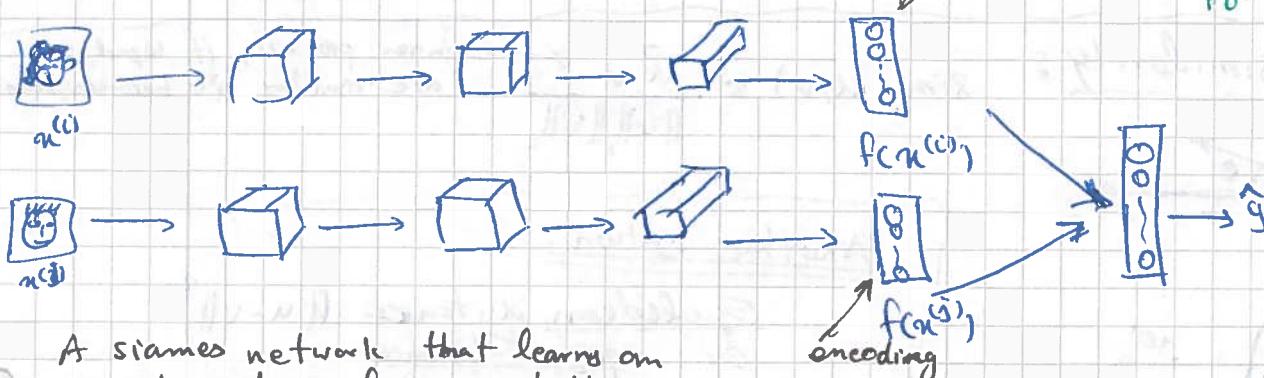
- Transfer learning for WEs learned over very large datasets ( $10^{8-10}$  words) into ~~other~~ task over much smaller datasets (100k words)

- Learn word embeddings from large text corpus (1-100B words)  
(Or download pre-trained embeddings online)
- Transfer embeddings to new tasks with smaller training set (say, 100K words)
- Optional: Continue to fine tune (adjust) the word embeddings with new data. normally only if you have enough data

Word Embeddings make a big difference when the task you trying to carry out, has a relatively small training set.

- Named Entity Recognition
- Text summarization
- Coreference resolution
- Parsing

### Relation to Face Encoding:



A siamese network that learns an encoding of a face, and then compares the two encodings to see if they belong to the same person.

### Difference between face encoding and word embedding:

- In face encoding we want to be able to compute the encoding for any new image, even those we have not seen before,
- In WE we have a fixed vocabulary of say 10K and we learn only for each of them a feature vector representing it.

### Properties of WEs:

- WEs help with analogy reasoning

Mikolov et al., 2013  
"Linguistic similarities in continuous space word representations"

	Man	Woman	King	Queen	Apple	Orange
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.70	0.69	0.03	-0.02
Food	0.69	0.01	0.02	0.01	0.95	0.97

### Question (Analogy Reasoning):

Man  $\xrightarrow{\text{to}}$  Woman is like King  $\xrightarrow{\text{to}}$  ? (what?)

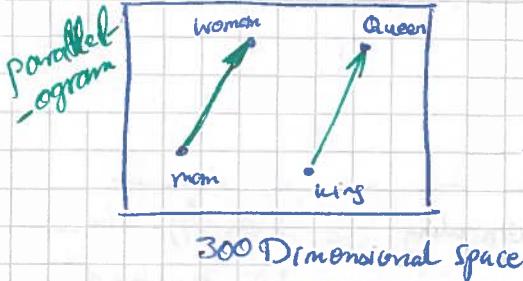
$$e_{\text{man}} = e_{\text{woman}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \end{bmatrix}$$

$$e_{\text{king}} - e_{\text{queen}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \end{bmatrix}$$

$$\Rightarrow e_{\text{man}} - e_{\text{woman}} \approx e_{\text{king}} - e_{\text{queen}} \Rightarrow \text{Queen}$$

difference at man and woman is like difference at king and queen.

## ← Analogies using Vector

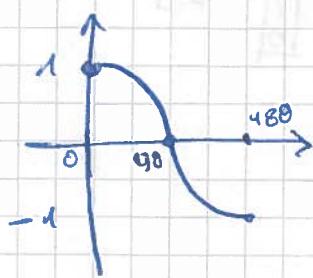


Cosine similarity:



$$\text{sim}(u, v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$$

inner product, if  $u$  and  $v$  are similar  $u^T v$  becomes very large



Another option:

Euclidean distance  $\|u - v\|^2$   
or squared distance

Cosine more often, because it normalizes better on length of  $\vec{u}$  and  $\vec{v}$

- Generality of analogy relationships that WEs can learn is remarkable:

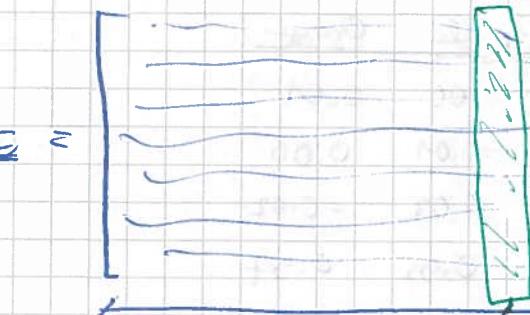
! Man: Woman as Boy: Girl || Ottawa: Canada as Nairobi: Kenya  
Big: Bigger as Tall: Taller || Yen: Japan as Ruble: Russia.

## Embedding Matrix:

our goal is to learn the embedding matrix E

Vocabulary = [a aaron

----- orange ----- rulu cunkr];  $M = 10,000$



- $E$  is the  $(10,000)$  dim embedding matrix.
- Each column of  $E$  is the embedding for the corresponding word from the dictionary

$$(300, 10k) \quad (10k, 1) \quad (300, 1)$$

$$E \cdot o_{6257} = e_{6257}$$

one-hot representation of Orange

$$o_{6257} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{bmatrix} \quad \leftarrow 6257$$

$$E \cdot o_j = e_j$$

Embedding matrix      One-hot represent. of word  $j$       Embedding of word  $j$

# Learning Word Embeddings : Neural Language Model (Example)

I want a glass of orange

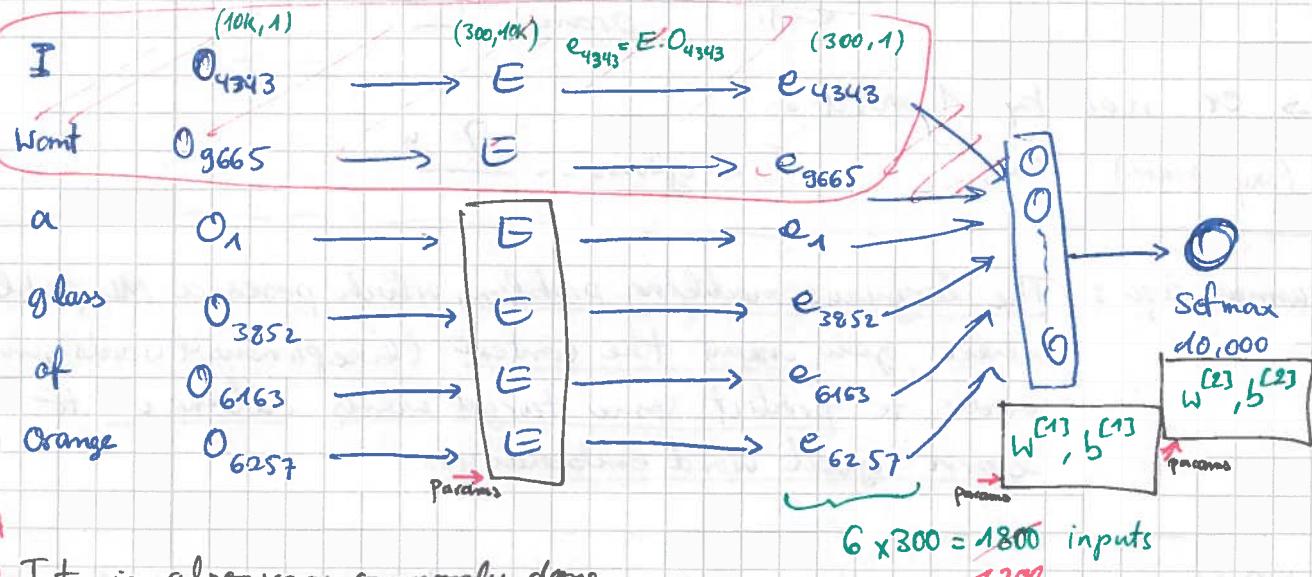
9343	9665	1	3852	6163	6257
------	------	---	------	------	------

Fixed historical window

We want to predict this.

Bengio et al.  
2003  
"A neural probabilistic language model"

- Building a Language Model is a reasonable way to learn a set of embeddings



- It is also very commonly done to have a fixed historical window

↳ => you can deal with different length sentences, because input sizes are always fixed

- Params of this model are: Matrix E (always the same for all) the words

$$\begin{aligned} &\rightarrow W^{[1]}, b^{[1]} \\ &\rightarrow W^{[2]}, b^{[2]} \end{aligned}$$

Now we can use Backprop and GD

to maximize the likelihood of the training set

try to predict what is the next word given previous words in the sequence.

- Interestingly this Algorithm will learn pretty decent word embeddings

↳ Because it is in the algorithm's incentive to learn pretty similar embeddings for say "orange" and "apple" for the example sentence above.

Other context/target pairs:

> cereal

I want a glass of orange juice to go along with my

Context

target

Context: 4 words before

↳ we can experiment with the context, given the task we have →

- ← for example, if the task is to learn a language model, the context could be "the last 4 words"
- But the context could be: "4 words on the left and right"  
e.g. "a glass of orange ? to go along with."
- Or a simpler context: "last 1 word"  
e.g. "Orange ?"
- Or nearby 1 word:  
(skip gram) e.g. "...glass... ?"

to summarize: The language modeling problem, which poses a ML problem where you input the context (like previous 4 words), in order to predict some target words allow us to learn good word embeddings.

## Word2Vec :

Compared to a language model, the research has shown that even simpler context → target mappings, e.g. Only one word will also result in very good WE results.

### Skip Gram:

I want a glass of orange juice to go along with my cereal.

Mikolov et al. '13  
"Efficient estimation of word representations in vector space"

Context  
orange  
orange  
orange

target  
juice  
glass  
my

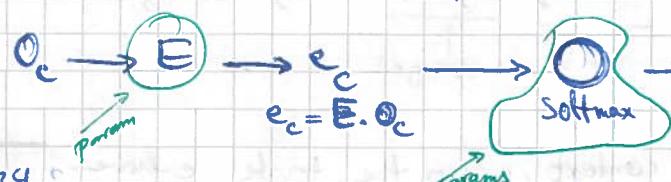
- Create a set of randomly picked context-target pairs
- Select the context word first randomly
- then randomly pick another word within a given window.

- We set up a supervised learning problem that given the context word, what could be a randomly chosen word predict within a say ±10 word window.

Vocab size = 10k → can even be much larger > 1M

Not an easy problem; because the set of possible words in the vocabulary is large to choose from.  
But the purpose of this problem is in fact to learn the WEs.

Learn a mapping:  
context  $c$  ("orange")  $\xrightarrow{6257}$  target  $t$  ("juice")  $\xrightarrow{4824}$



$$\text{Softmax: } p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10000} e^{\theta_j^T e_c}}$$

$\theta_t$ : parameter associated with output  $t$

← The loss function for the Softmax will be the usual negative log likelihood

$$L(\hat{y}, y) = - \sum_{i=1}^{10000} y_i \log \hat{y}_i \quad \rightsquigarrow y \text{ as one-hot vector}$$

0
1
3
0

4834 ←

Parameters:

- Matrix  $E$ : all the embedding vectors  $e_c$
- $\Theta_t$  parameters of the softmax.

→  $\hat{y}$  will be a 10,000 dim vector output from softmax unit, with probabilities for each of the words in the vocab being the target.

← This model is called **Skip Gram**-model. It can learn pretty descent Embeddings

- But it has some problems:

- Computational speed

$$p(t|c) = \frac{e^{\Theta_t^T e_c}}{\sum_{j=1}^{10k} e^{\Theta_j^T e_c}}$$

Sum over the whole vocab.

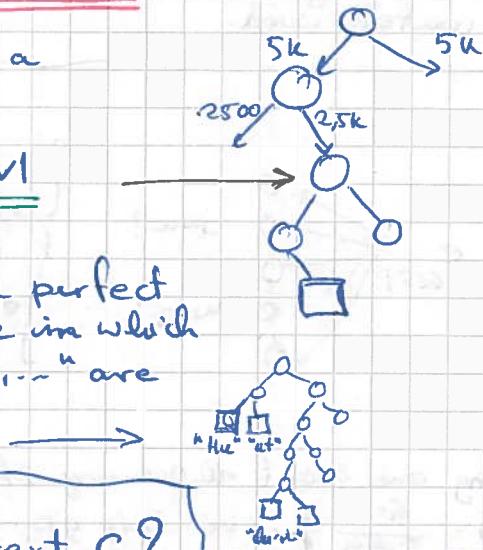
Very slow!

⇒ use a Hierarchical Softmax:

- each interior node is a binary classifier

⇒ Cost will be order  $\log V$  rather than  $V$

- in practice, don't use a perfect balanced tree, but a tree in which common words like "the, at, ..." are more at top



How to sample the context  $C$ ?

↳ Uniformly? → there are very frequent words "the, of, a, and, to, ..."

⇒ don't sample uniformly; use some heuristics that let you sample also less frequent words fairly.

Negative Sampling : Skip Gram can find very useful embeddings, but the down side of it is the computational cost.  
↳ too slow.

"I want a glass of orange juice to go along with my cereal."

Question : Given a pair of words, is it a context-target pair?

Context	word	Target?	
Orange	juice	1	pick target like before "randomly within a given window"
Orange	king	0	Pick the non-targets simply randomly from the dictionary.
Orange	book	0	

- to summarize
- How to generate the dataset
- Pick a context word, and a target within a window.  
That gives us the "positive example" → label 1
  - Then for k times, pick random words from the dictionary to create "negative examples" → label 0

The supervised Learning problem : Input: the word pair → Output. 1 or 0

- How to choose k?  $k = 5-20$  for smaller DSets.  
 $k = 2-5$  for larger DSets.

Model :

$$P(y | c, t) = \sigma(\theta_t^T e_c)$$

estimate the probability  
that  $y = 1$

X		Y	
context	word	target?	
orange	juice	1	
orange	king	0	
orange	book	0	
k orange	the	0	
orange	of	0	
c	t	y	

- Similar parameters as softmax before:

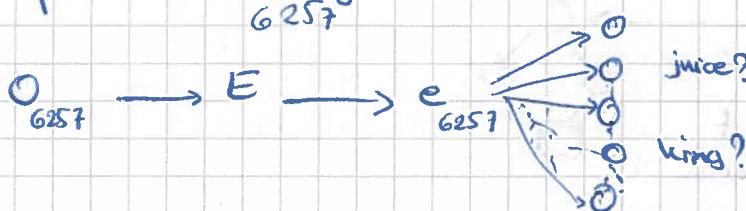
→  $\theta_t$  one param vector for each possible target word

→  $e_c$  the embedding vector for each possible context word.

- $k$  to 1: the ratio of negative to positive examples.

Input word: orange

6257



• 10,000 binary logistic regression classifiers

- Instead of training all 10,000 of them, at each iteration we only train  $k+1$  of them.

- Instead of having one giant 10,000 way softmax, which is very expensive to train ( $P(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_j e^{\theta_t^T e_{c_j}}}$ , the skip gram model), we now have 10,000 binary classification problems, which are really cheap to compute and on every iteration, we are going to train only  $k+1$  of them.
- This technique is called **Negative Sampling** because for each "positive example" we have  $k$  "negative examples"

- How to sample the negative examples?

one extreme

- Sample based on  $w_i$  frequency of each word

$$P(w_i) \propto f(w_i)$$

- you end up having much higher representation of frequent words: "at, if, the, an"

In Between

- works best
- empirical

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10k} f(w_j)^{3/4}}$$

Other extreme

- Sample uniformly with probability for each word  $\frac{1}{|V|}$

# GLOVe : (Global Vectors for word Representations)

Pennington et. al. 2014  
"GLOVe: Global Vectors for word representations"

I want a glass of orange juice to go along with my cereal.

$c, t$

$$X_{ij} = \# \text{ times } i \text{ appears in context of } j \quad \begin{matrix} \uparrow \\ t \end{matrix} \quad \begin{matrix} \uparrow \\ c \end{matrix} \quad \begin{matrix} \text{counted co-appearances,} \\ \text{within a given window} \end{matrix}$$

depending on the choice of what is context and what is target,  $X_{ij}$  can be symmetric:  $X_{ij} = X_{ji}$

Model :

minimize  
the difference  
 $(\theta_i^T e_j - \log X_{ij})$

$$\sum_{i=1}^{10k} \sum_{j=1}^{10k} f(X_{ij}) (\underbrace{\theta_i^T e_j + b_i + b_j}_{\text{weighting term}} - \log X_{ij})^2$$

bias  
 $b_t$      $b_c$

how often they appear together.

how related are words  $t$  and  $c$ ,

which is measured by

- learn two vectors ( $\theta_i$  and  $e_j$ ) whose inner product is a good predictor for how often these two words  $i, j$  appear together.

- $f(X_{ij})$ , the weighting term,  $f(X_{ij}) = 0$  if  $X_{ij} = 0$  (because  $\log 0$  not defined) and

also weighting factor can give more weight to less frequent words, and less weight to more frequent (like stop word) words.

both initialized randomly  $\theta$  and  $e$

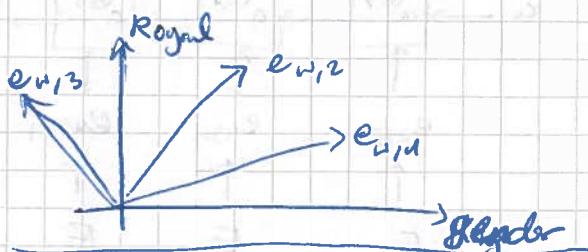
- Interestingly, looking at math the role of  $\theta_i$  and  $e_j$  are now pretty symmetric.

$$\Rightarrow e_w^{(\text{final})} = \frac{\theta_w + \theta_w}{2}$$

We can calculate the final embedding as the average over learned  $\theta$  and  $e$

A note on featurization view of word embeddings:

	Man 5291	Woman 9277	King 4464	Queen 2457
Gender	-1	1	0.95	0.97
Royal	0.01	0.02	0.93	0.95
Age	0.03	0.02	0.70	0.69
Food	0.09	0.01	0.02	0.01

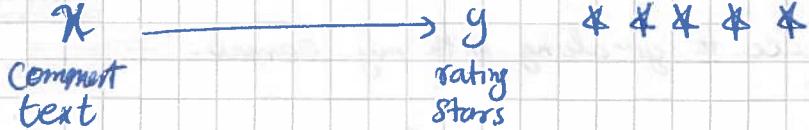


- the components of the WE vector are not interpretable.
- They might not even be orthogonal to each other.
- The parallelogram math still works.

$$\text{minimize} \sum_{i=1}^{10k} \sum_{j=1}^{10k} f(X_{ij}) (\underbrace{\theta_i^T e_j + b_i + b_j - \log X_{ij}}_2)$$

$$(A\theta_i)^T (A^T e_j) = \theta_i^T A^T A e_j \\ = \theta_i^T e_j$$

## Sentiment Classification :



- One of challenges of Sentiment Analysis is you may not have a huge labeled Data set.  
 e.g. (10k - 100k words)

"The desert is excellent"

8928 2468 4694 3180

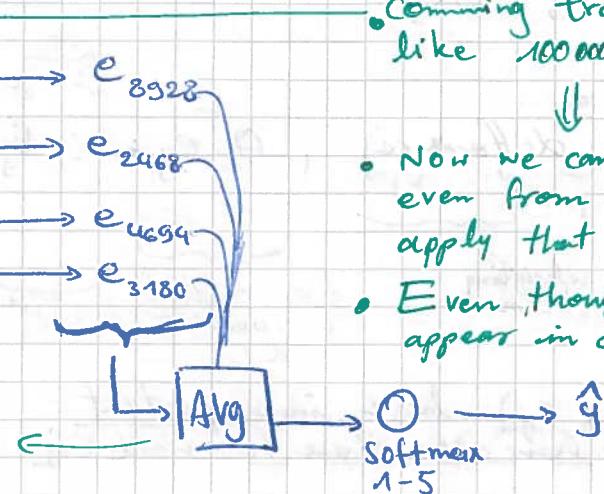
★★★☆ (4 stars)

• Coming trained on a very large data set like 100000M words.  $10^9$  words



- Now we can take a lot of knowledge even from infrequent words and apply that to our problem.
- Even though those word do not even appear in our training set.

Avg allows working with reviews that are short or long.



• This Algorithm works very well. It sums (averages) the meanings in our words.

• A problem is that, it ignores **word order**. e.g.

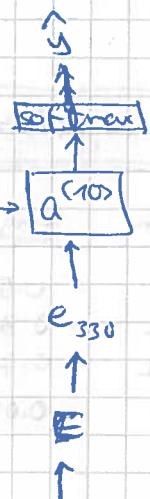
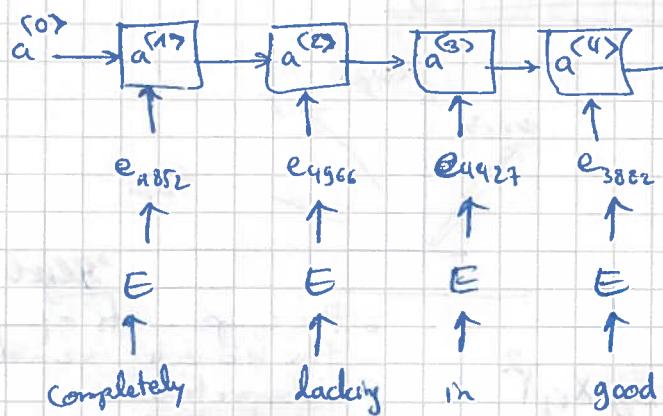
"Completely lacking in **good** taste, **good** service, and **good** ambience."

One star  
 : (★)

But by simply averaging the words, our algorithm would think this is a very "good" review.

⇒ Instead of just summing over all of the WEs, we can instead ~~use~~ use a RNN for sentiment classification.

## RNN for sentiment classification:



{ actually the one-hot representation of the words }

(Many-to-one RNN)

# Debiasing word Embeddings

Bodukbasi et.al. '16  
"Man is to computer as woman is to homemaker"  
Debiasing word embeddings

Man : Woman as King : Queen

Man : Computer-programmer as Woman : Homemaker X Gender bias!

Father : Doctor as Mother : Nurse X !

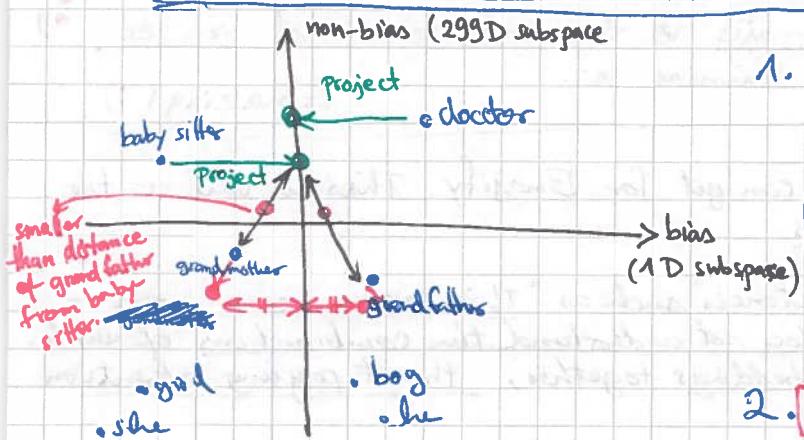
Word embeddings can reflect gender, ethnicity, age, sexual orientation, and other biases of the text used to train the model.

- Learning algorithms are increasingly making very important decisions:

- College Admissions
- Jobs ~~are gendered~~
- Loan applications

- Even sentencing guidelines in the criminal justice systems. !

## Addressing bias in Word Embeddings :



Example here for gender.  
Same procedure for any other bias.

In original paper the bias direction can be multidim. and instead of Avg, SVD is used to find the bias direction

### 1. Identify the bias direction:

{  
e-he - e-she  
e-male - e-female  
}  
Averaging  $\rightarrow$  This gives the bias direction

### 2. Neutralize:

for every word that is not definitional (has gender in its definition, such as girl, boy, grandma, grandpa) project to get rid of bias.

→ eliminate their components that are in bias direction (contribute to creating bias)

### 3. Equalize pairs.

for pairs like, grandfather - grandmother  
girl - boy  
woman - man

that have gender,  
make sure they have "the same distance" from neutralized word.

For words that have gender in their definition, we want the only difference in their embedding to be the gender.

Make sure that words like grandmother and grandfather has "exactly the same" distance from

words that should be gender neutral like babysitter or doctor.

How to choose words that should be Neutralized? Authors of the paper trained a classifier,

most words in English language are gender neutral for that.

accordingly the set of pairs we want to neutralize is relatively small.

# Course 5 - RNNs - Week 2 - Programming 1 - Operations on Word Vectors

- Load pre-trained word vectors and measure similarity using cosine similarity
- Use WEs to solve word analogy problems such as Man to Woman is King to ?
- Modify WEs to reduce their gender bias

Word Analogy problem

Debiasing:

It is really :)  
amazing! See Notebook

- Using a pre-trained set of word vectors from the internet is often a good way to start a NLP application.

$$\begin{aligned} \text{"I"} & e_i = \begin{pmatrix} 0.12 \\ 1 \\ 0.32 \end{pmatrix} \\ \text{"love"} & e_{\text{love}} = \begin{pmatrix} 0.75 \\ 1 \\ 0.11 \end{pmatrix} \\ \text{"you"} & e_{\text{you}} = \begin{pmatrix} 0.54 \\ 1 \\ 0.13 \end{pmatrix} \end{aligned}$$

(50, 1)

Word embedding matrices pre-trained using GloVe

Emojify with softmax

Avg

Avg

0.47

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

0.19

1

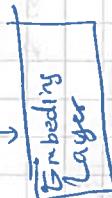
</div

- In Keras, the embeddings are represented as a layer that takes integers (word indices) and outputs dense embedding vectors. This layer can be initialized using existing pre-trained embeddings or be trained.

"I love you"  
"I will be ready"  
word-to-index  
↓

$(2, \text{max\_len}) \left[ \begin{array}{c} [185457, 226278, 394475, 0, 0] \\ [185457, 387696, 71090, 302259, 0] \end{array} \right]$

zero padding



$[e_{185457}, e_{226278}, e_{394475}, 0, 0]$   
 $[e_{185457}, e_{387696}, e_{71090}, e_{302259}, 0]$ ,  
 $(2, \text{max\_len}, 50)$

length of GloVe Emb. Vectors

### Things to remember

- If you have an NLP task where the training set is small, using word embeddings can help your algorithm significantly.
- Training sequence models in Keras (and most other DL Frameworks) requires few important details:
  - To use mini-batches, all the sequences must be padded, so that all the examples in a mini-batch have the same length.
  - The Embeddings() layer can be initialized with pre-trained values. You can further train those values (if you have a large enough Data set) or leave them fixed.
  - LSTM() has a flag called "return\_sequences" to decide if you want to return every hidden state or only the last one.
  - We can use Dropout() right after LSTM() to regularize the network.

## Course 5 - Sequence Models - Week 3 - Sequence Models and Attention Mechanism

### Week 3

### Sequence Models and Attention Mechanism

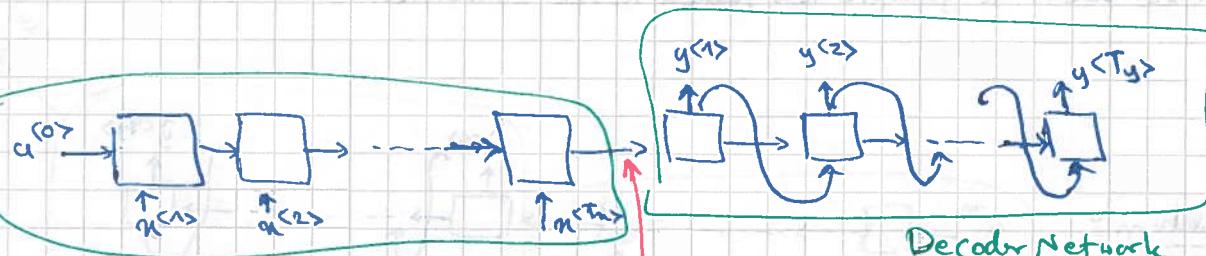
video 14:42m  
prob. 2 h 30 m

#### Sequence to Sequence Models : Basic Models

↳ Machine Translation, Speech Recognition, etc.  
Input sequence:  
Jane visite l'Afrique en Septembre

→ Jane is visiting Africa in September. ← Output sequence  
 $y^{(1)}, y^{(2)}, y^{(3)}, y^{(4)}, y^{(5)}, y^{(6)}$

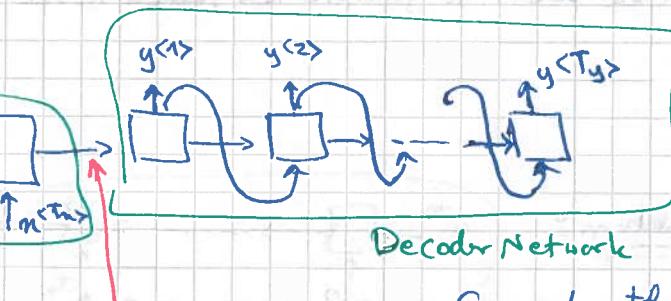
Sutskever et al. '14  
"Sequence to sequence learning with neural networks"  
Cho et al. '14  
"Learning phrase representations using RNN encoder-decoder for statistical machine translation"



#### Encoder Network

- get the input sequence one word at a time and encodes it

BRUNN  
• Can be any kind of RNN like LSTM, GRU.



#### Decoder Network

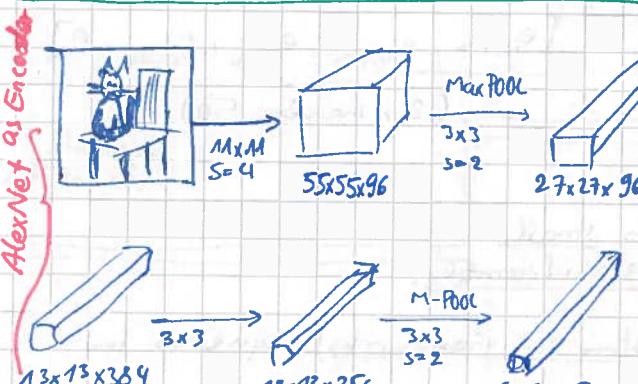
- Generates the output sequence
- Similar to the way we generated text after learning a language model.

Given enough pairs of french-english training examples, this model works pretty decently :)

→  
100%

← Image Captioning → An architecture very similar to the encoder-decoder for machine translation can be used for automatic Image Captioning.

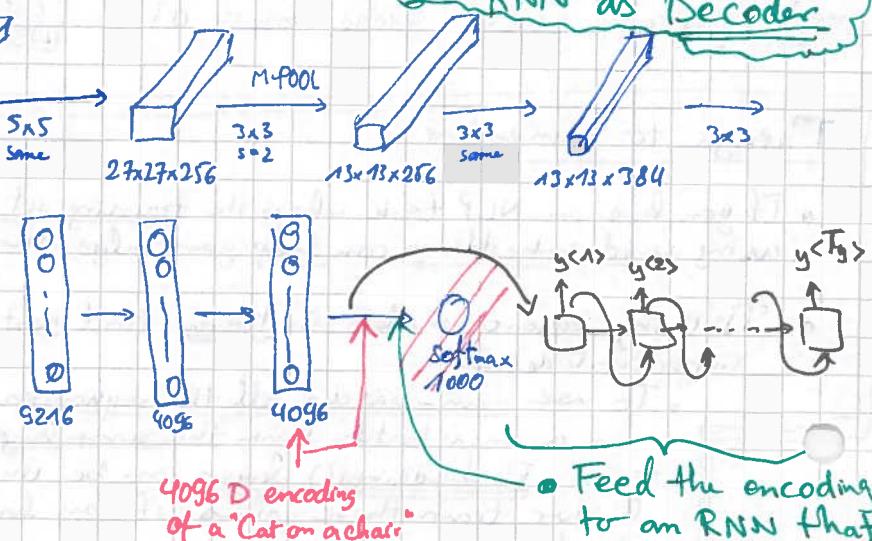
- Mao et al., '14, "Deep captioning with multi-modal recurrent neural networks"
- Vinyals et al., '14, "Show and tell: Neural Image Caption generator"
- Karpathy and Fei Fei, '15, "Deep visual semantic alignments for generating image descriptions"



$y^{(1)} \ y^{(2)} \ y^{(3)} \ y^{(4)} \ y^{(5)} \ y^{(6)}$   
A cat sitting on a chair.

AlexNet as Encoder  
RNN as Decoder

Interesting, this model works pretty well.



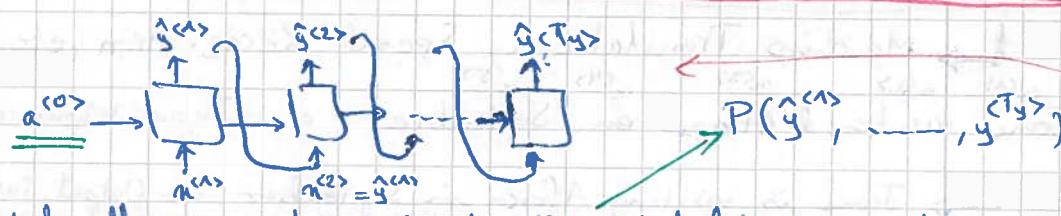
- Feed the encoding to an RNN that generates an output sequence of words.

- One important difference of sequence-to-sequence models and the models we used for synthesizing novel text using a language model is:
  - Here you don't want a randomly chosen translation. Here you want the most likely translation; We don't want a randomly chosen caption, but repeat the best and most likely.

## Picking the Most Likely Sentence:

- We can think of machine translation as "conditional language model".

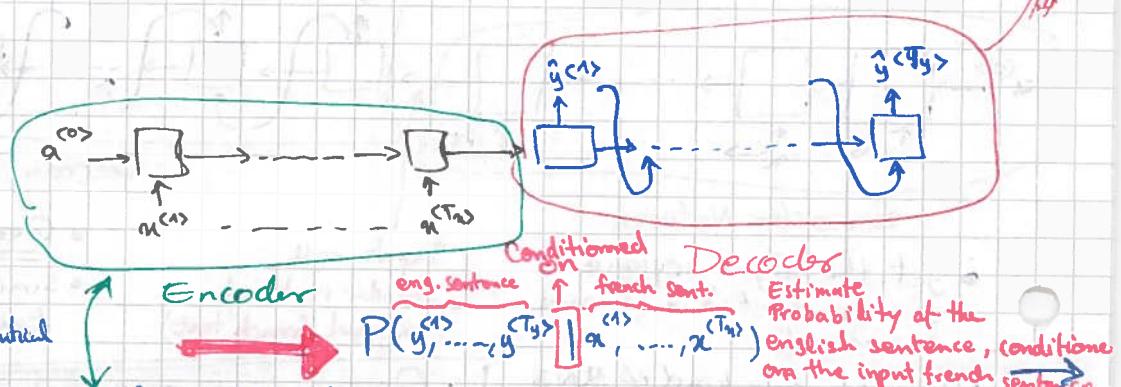
Language Model :



- A language model allows us to estimate the probability of a sentence
- We could also use it to generate new ~~sente~~ novel sentences.

Identical

Machine Translation :



- The decoder of a machine translation model is similar/identical to a language model.

Conditioned Language Model.

- Machine translation model is very similar to the language model, except that instead of starting off with the vector of all zeros, it starts off with the context vector z<0>.

## Finding the most likely translation:

- Given the french sentence  $x$ , the model should tell us what is the conditional probability  $P(\hat{y}^{(1)}, \dots, \hat{y}^{(T_y)} | x)$  of different english translations.
- We don't want to sample words for the english sentence randomly  
Because we could then end up with very stupid translations, like:  
→ "Her african friend welcomed Jane in September."
- What we want is to find an english sentence  $y$ , that maximizes that conditional probability:

$$\arg \max_{y^{(1)}, \dots, y^{(T_y)}} P(y^{(1)}, \dots, y^{(T_y)} | x)$$

### Why not a greedy Search?

- ↳ • just pick the most likely word at each time-step, according to the conditional language model.
- This does not work. We want the joint probability of all words in the whole sentence to be maximized given the french sentence.

Example:

→ Jane is visiting Africa in September. ← Better translation, higher  $P(y|x)$

→ Jane is going to be visiting Africa in September. ← Less good, lower  $P(y|x)$

• But greedy search would generate the second one! !

↳ Because  $P(\text{Jane is going} | x) > P(\text{Jane is visiting} | x)$

only based on maximum probability of the next word

:C

So we need something that considers the joint conditional probability of the whole sentence

↳ But the space of possible sentences is huge; exponentially large:  $|V|^L \approx 10,000$

⇒ We need an approximate search algorithm

↳ Tries to pick the sentence  $y$  that maximizes  $P(y|x)$

• Although it is not guaranteed to find  $y$ .

Space of possible sentences with  $\approx 10,000$   
10 words ↑

impossible to enumerate them all

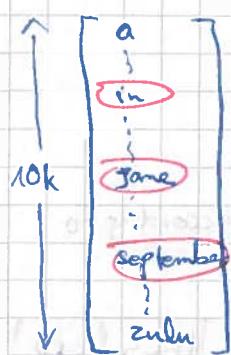


## ← Beam Search Algorithm :

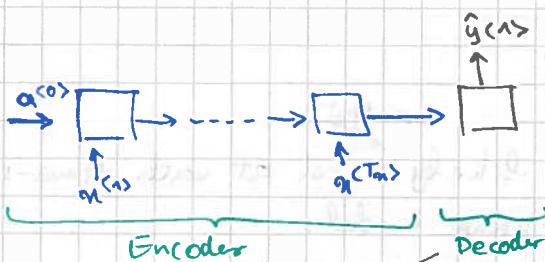
- Given a french sentence as input, we don't want to output a random english sentence. We want to output the best, the most likely sentence.
- Also given an audio signal, we want to output the most likely transcription
- Beam Search** is a widely algorithm to achieve this.

Beam Search : Instead of picking only the one most likely word at each step, pick the  $B$  (beam width) top most likely words, and build up a tree of search by picking  $B$  words at each step. (Greedy Search)

### Step 1



Beam Width,  $B = 3 \Rightarrow$  Pick the three most likely words at first step.



This network calculates

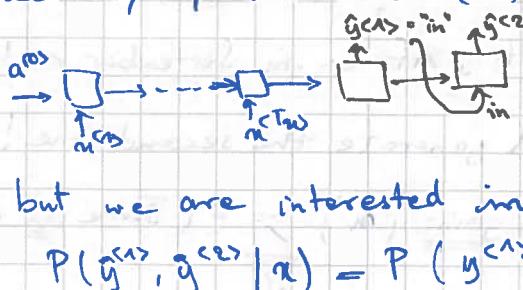
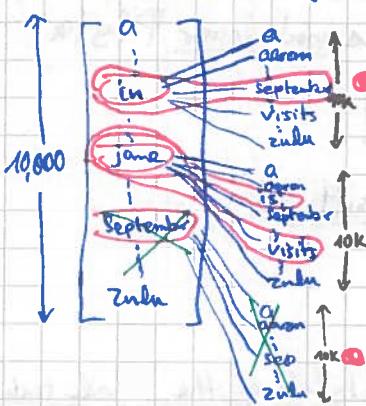
$$\rightarrow P(y^{(1)} | x)$$

Let's say evaluating  $P(y^{(1)} | x)$  outputs the words "in", "Jane", and "September" as the most likely three words.

A softmax output with 10,000 possibilities, out of which we pick the three most likely.

### Step 2

: Pick the second ~~most likely~~ words, based on the probability of the whole (2) combinations  $P(y^{(1)}, y^{(2)} | x)$



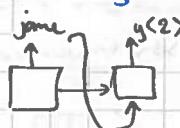
This network evaluates:

$$P(y^{(2)} | x, "in")$$

but we are interested in

$$P(y^{(1)}, y^{(2)} | x) = P(y^{(1)} | x) \cdot P(y^{(2)} | x, y^{(1)})$$

We do the same thing for the second choice of the most likely first word i.e. "Jane".

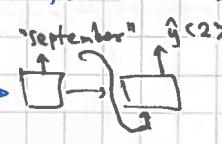


This network computes:

$$P(y^{(2)} | x, "Jane")$$

$$\text{we want } P(y^{(1)}, y^{(2)} | x) = P(y^{(1)} | x) \cdot P(y^{(2)} | x, y^{(1)})$$

• The same thing for the third ~~most likely~~ first word.



... - - -



⇒ We can now reject "September" for the choice of the first word. X

Now we evaluate all the 30,000 possibilities for the ~~first~~ combination of the first and second word and pick the most likely top 3 again.

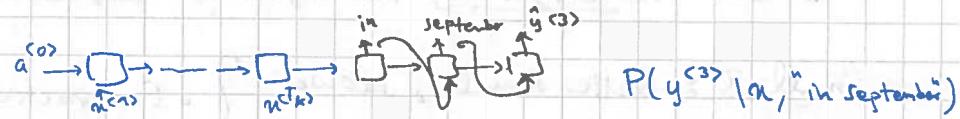
"combinations!"

"in september"  
"Jane is"  
144 "Jane visits"

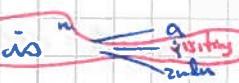
- At every step we instantiate  $B=3$  copies of the network, and evaluate all the  $30,000 P(y^{(1)}, y^{(2)} | n)$  probabilities using them.

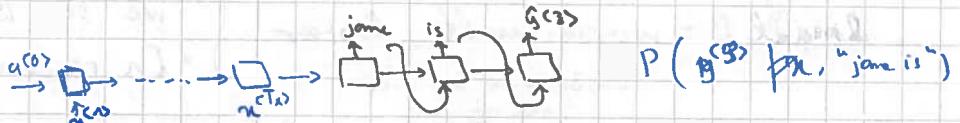
Step 3: so far we have saved 3 possibilities for the most likely combinations for the first 2 words.

"in September" 



$$P(y^{(3)} | n, "in September")$$

"jane is" 



$$P(y^{(3)} | n, "jane is")$$

"jane visits" 



$$P(y^{(3)} | n, "jane visits")$$

$$P(y^{(1)}, y^{(2)}, y^{(3)} | n) = P(y^{(1)}, y^{(2)} | n) \cdot P(y^{(3)} | n, y^{(1)}, y^{(2)})$$

! Note:  $B=1$  would be equivalent to "Greedy Search"

## Beam Search Refinements:

Length Normalization: What Beam search does is maximizing:

$$\arg \max_y \prod_{t=1}^{T_y} P(y^{(t)} | n, y^{(1)}, \dots, y^{(t-1)})$$

But all these  $P(\dots)$  are numbers  $< 1$  (small numbers) and multiplying all of them together will end up in a very very small number and numerical underflow. :(

$\Rightarrow$  Reformulate the problem: by adding a log before the product. [Because log is a strictly monotonic function, we know that maximizing the  $\log(f(n))$  is the same as maximizing  $f(n)$ .]

$$\Rightarrow \arg \max_y \sum_{t=1}^{T_y} \underbrace{\log(P(y^{(t)} | n, y^{(1)}, \dots, y^{(t-1)}))}_{\text{sum of log probabilities}} !$$

Normalization for Beam search

$$\approx \log P \leq 0$$

Another problem is that, both these objective functions ( $\prod$  and  $\sum \log$ ) prefer shorter sentences (Because longer sentences make  $\prod$  smaller, or  $\sum \log$  a larger negative number)

$\Rightarrow$  normalize the  $\sum \log$  sum by number of words  $T_y$  in the sentence, to reduce the penalty for outputting longer sentences.

A Hyperparameter to be tuned.

$$\arg \max_y \frac{1}{T_y^\alpha} \sum_{t=1}^{T_y} \log P(y^{(t)} | n, y^{(1)}, \dots, y^{(t-1)})$$

In practice use  $\alpha$  to make the normalization softer.

$\alpha=1$ , completely normalize by length

$\alpha=0$ , no normalization

$\alpha=0.7$  often used.

How to choose the final sentence?

Normalized log probability objective

Let's say we ran BeamSearch for 30 steps and  $B=3$

We see a lot of sentences with  $T_y = 1, 2, \dots$ , up to 30

We pick the best performing sentences and score them again this score. So, the best score wins.

## How to choose the Beam width B?

- The larger  $B$ , the more possibilities we are considering and so higher the chance of finding a good sentence.
  - But the larger  $B$ , the higher the computational cost, and memory need.
  - Larger  $B$ : Better Results, slower
  - Smaller  $B$ : worse results, faster
    - Because we keep track of less possibilities
- In practice, in production systems we see  $B \approx 10$
- In research we see even  $B = 100 \text{ to } 3000$
- 1, 3, 10, 100, 1000, 3000

Note: Unlike exact search algorithms like BFS or DFS, Beam Search runs faster but is not guaranteed to find exact maximum for  $\arg\max_y P(y|u)$ .

Error Analysis in Beam Search: Error analysis for finding best params and tuning BS

- BS is an approximate search alg; Also called a heuristic search app.
- Error Analysis also helps us in case of low performance to figure out if it is the beam search that causes the problem, or it is the RNN model that needs more work.

Example:

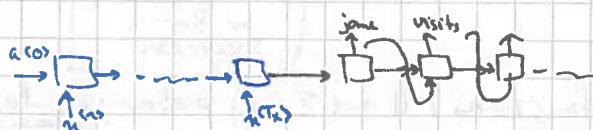
Input: Jane visite d'Afrique en Septembre.

Human: Jane visits Africa in September ( $\hat{y}^*$ ) ← better

Algorithm: Jane visited Africa last September ( $\hat{y}$ ) ← worse ← ?

Error Analysis for RNN with Beam Search

- Our model has two main components { RNN (encoder and decoder) } Which is responsible RNN or BS ???
  - Should we collect more training data?
  - Should we increase the beam width?



RNN computes  $P(y|u)$

⇒ let RNN compute  $P(\hat{y}^*|u) > ?$  and  $P(\hat{y}|u) < ?$

depending on the result we can decide what to do next

Human: Jane visits Africa in September ( $\hat{y}^*$ )  
Alg.: Jane visited Africa last September ( $\hat{y}$ )

let RNN compute  $P(\hat{y}^*|u)$  and  $P(\hat{y}|u)$

Case 1:  $P(\hat{y}^*|u) > P(\hat{y}|u)$

⇒ { Beam search chose  $\hat{y}$ . But RNN gave a higher probability to  $\hat{y}^*$   
⇒ Beam search is a fault. Because it fails to find a value of  $y$  that maximizes  $\arg\max_y P(y|u)$ .

Case 2:  $P(\hat{y}^*|u) < P(\hat{y}|u)$

⇒ {  $\hat{y}^*$  is a better translation than  $\hat{y}$ . But RNN predicted  $P(\hat{y}^*|u) < P(\hat{y}|u)$  →

Note: In case we use normalization, we should evaluate the optimization objective instead of the probabilities.

## Error Analysis Process:

- Go through the development set, and find the mistakes the algorithm did.

Human	Algorithm	$P(y^*(x))$	$P(\hat{y}(x))$	At fault?
Jane visits Africa in September.	Jane visited Africa in September.	$2 \times 10^{-10}$	$1 \times 10^{-10}$	Beam Search RNN BS R ?

=> We figure out what fraction of errors are due to BS and what due to RNN.

↳ We can attribute the error to either the search algorithm or to the objective function for RNN that generates the objective function that BS is supposed to be maximizing.

- If RNN is at fault, then use other network improvement techniques (Course 3) to improve it (change architecture, add more data, add regularization, ...)
- Such error analysis is very useful, whenever we have an optimization algorithm (e.g. BS) that wants to optimize the objective function output by a learning algorithm.

## Bleu Score:

### Bleu Score

Papineni et al.'02; "Bleu: A method for automatic evaluation of machine translation"

- A difficulty in evaluating machine translation systems is that there are multiple correct answers possible. Unlike CV and image recognition, where only one correct answer exists and hence accuracy is a good measure. How do you measure accuracy if there're multiple correct answers?

=> The Bleu Score

### Example:

French: le chat est sur le tapis.  
 Both translations → Reference 1: The cat is on the mat.  
 good and correct. → Reference 2: There is a cat on the mat.

Bleu Score automatically measures a score that tells how good is a piece of machine translation.

### Bleu: Bilingual evaluation underly.

↳ (means substitute for an actor in theatre)

Intuition: every translation that is near to the given reference translations, gets a high Bleu score.

- Intuition: Look at machine generated output and see if types of words it generates appear in at least one of the human generated references.

Imagine: MT output: the the the the the the.

precision: 1 if a word from output appears in the references.

Lets define  $\frac{7}{7} \leftarrow$  each word appears also in references  
 $\frac{7}{7} \leftarrow$  # words in output

### Credit & Modified Precision

=> Precision is not a useful measure here.

Modified Precision: define credit as the max number of times a word appears in references.

↳ "the" got credit 2 => Modified Precision =  $\frac{2}{7} \leftarrow$  max count("the") in refs.

Bleu Score in Bigrams : **Bigram**: pairs of words appearing next to each other.

- We don't wanna look at words only in isolation - Let's look at pairs of words as well.

Ref 1: The cat is on the mat.

Ref 2: There is a cat on the mat.

MT output: The cat the cat on the mat.

Bigrams	Count	Count		Modified Bigram Precession ( $P_2$ )	4
		clip	clip		
the cat	2	1	1		
cat the	1	0	0		
Cat on	1	1	1		
on the	1	1	1		
the mat	1	1	1		
	6	4	4		

$$P_n = \frac{\sum_{\substack{\text{unigrams} \in Y \\ \text{paper}}} \text{Count}_{\text{clip}}(\text{unigram})}{\sum_{\substack{\text{unigrams} \in Y}} \text{Count}(\text{unigram})}$$

for n-grams

$$P_n = \frac{\sum_{\substack{\text{n-grams} \in Y}} \text{Count}_{\text{clip}}(\text{n-gram})}{\sum_{\substack{\text{n-grams} \in Y}} \text{Count}(\text{n-gram})}$$

- Modified Precision Scores allow to measure the degree of how much is the MT similar (in other words "Overlaps") with the reference translations.

↳ If MT output is exactly equal to one of the references, then all the values  $P_1$  to  $P_n$  will be equal to 1.0 ( $P_1 = P_2 = \dots = P_n = 1.0$ )

Bleu Details :  $P_n$  = Bleu score on n-grams only

To compute the one Bleu score, by convention, compute  $P_1$  to  $P_4$ :

**Combined Bleu Score** :  $BP \exp \left( \frac{1}{4} \sum_{n=1}^4 P_n \right)$

**BP** : Bravity Penalty:

is an adjustment factor that penalizes too short translations.

(Because) It is easier to get high precision if we output very short translations  
(Because) most of the words we output appear in the references. But we don't want translations that are too short.

$$BP = \begin{cases} 1 & \text{if } \text{MT\_output\_length} > \text{reference\_output\_length} \\ \exp(1 - \text{MT\_output\_length} / \text{reference\_output\_length}) & \end{cases} \quad \text{see paper for details.}$$

- Refer to Course 3 on the importance of having one single real numbers evaluation metric. The Bleu score provided the field of MT with such a metric and it was revolutionary and accelerated the advancement in this field a lot.

- It is also used to evaluate image captioning systems, and generally any kind of system that generates text.

truth:  
good  
one  
have

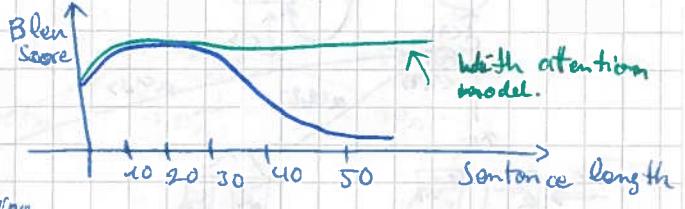
# Attention Model Intuition: one of the most influential ideas in DL that makes the encoder-decoder Architectures work much better

## The problem of very long sequences:

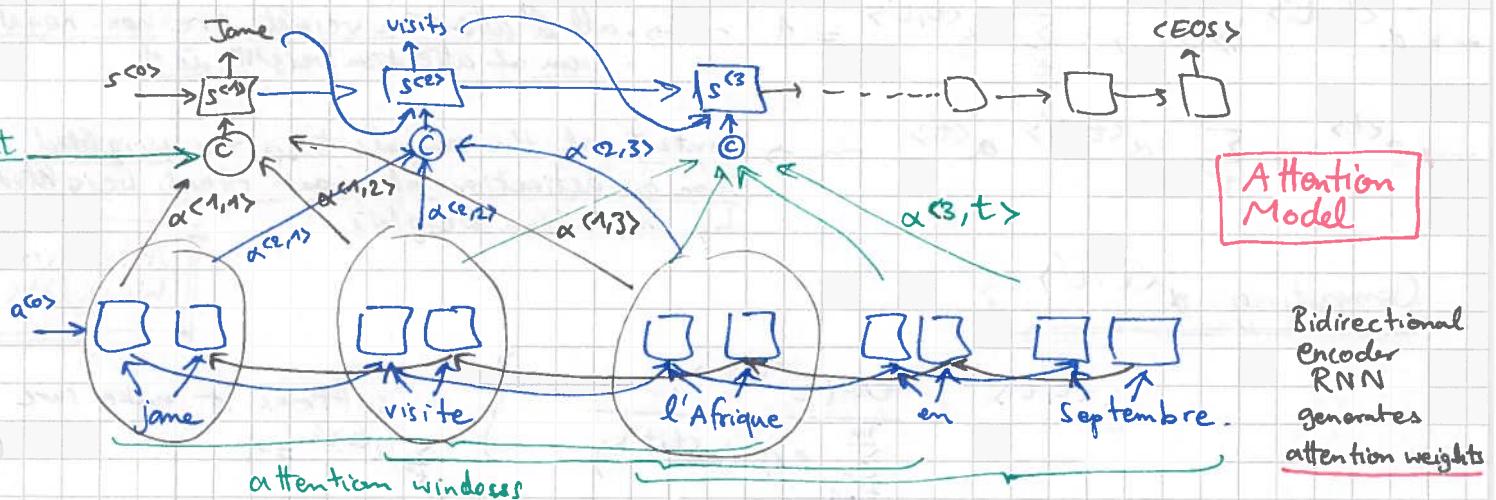
Bahdanau et al '14  
"Natural machine translation by jointly learning to align and translate"

Very seminal

- the NN cannot memorize the whole sentence in the encoder  $\rightarrow$  
- Performance degrades when sentences become too long.
- A human translator would not memorize the whole sentence first before translating. He would translate a little part of sentence first, then another little part, then another little part and so on. Attention Model does something similar.
- Attention Model (AM) originally developed ~~for~~ for MT, but is also very successfully applied to other ~~fields~~ domains. Very influential. Also applied to image captioning.



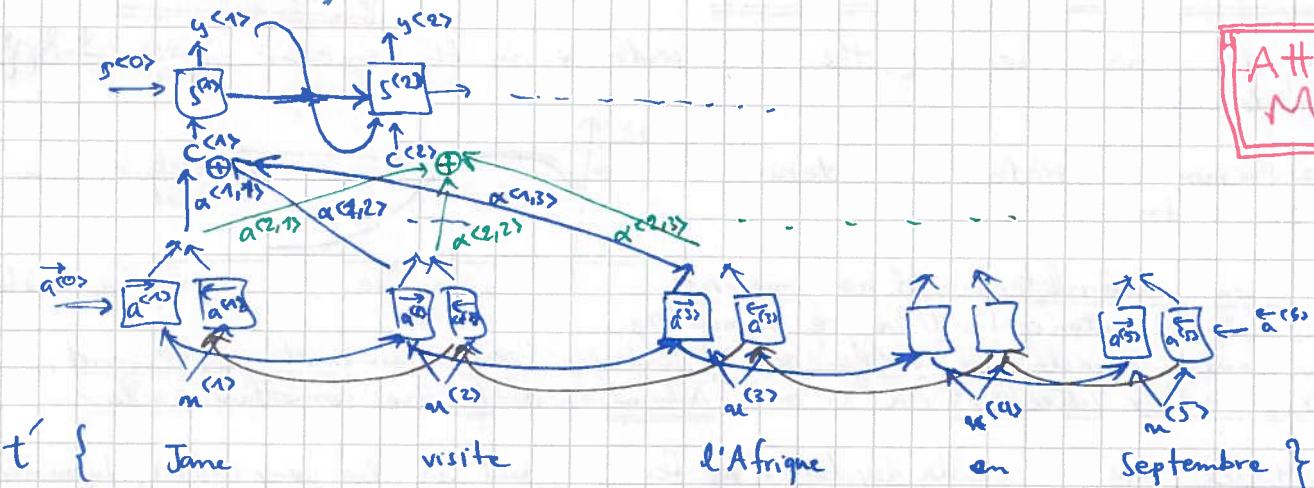
- Use the concept of "Attention Weights" to denote how much should you be paying attention to input words, and define a Window of attention in the input sequence.



- $\alpha^{<1,2>}$ : How much attention should we pay to input word 2 in order to generate output word 1
- The amount of attention to pay at generating the 3rd output to the input French word at time t depends on the activations of bidirectional RNN at time t, and on the state of the previous step.  $\alpha^{<3,t>}$  depends on  $\left\{ \begin{array}{c} s^{<2>} \\ \alpha^{(t)}, \bar{\alpha}^{(t)} \end{array} \right\}$
- The RNN marches forward generating one word at the time governed by  $\alpha^{<t,t>}$  at each step, that tell how much attention to pay to the  $t^{\text{th}}$  input word.

## Attention Model:

- Attention model allows the NN to only pay attention to parts of the input sentence when generating the output sequence (Much as a human translator would do)



Attention Model

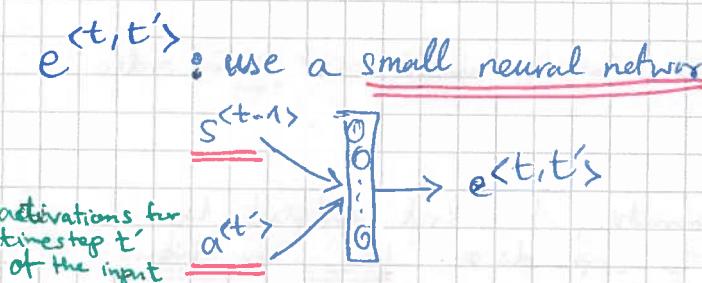
- $a^{(t')} = (\vec{a}^{(t')}, \overleftarrow{a}^{(t')})$  activations of the BiRNN at timestep  $t'$  of input
- $\alpha^{(t,t')}$ : amount of attention  $y^{(t)}$  should pay to  $a^{(t')}$
- $\alpha^{(t,t')} \geq 0$ ,  $\sum_t \alpha^{(t,t')} = 1$  → all attention weights are non-negative  
• sum of attention weights is 1.
- $C^{(t)} = \sum_{t'} \alpha^{(t,t')} a^{(t')}$  → Context of the output  $t$  is the weighted sum of activations of input steps, weighted by attention weights.

Computing  $\alpha^{(t,t')}$ :

$$\alpha^{(t,t')} = \frac{\exp(e^{(t,t')})}{\sum_{t=1}^T \exp(e^{(t,t')})}$$

} a softmax to make sure  
 $\sum_t \alpha^{(t,t')} = 1$

Attention Weights



use a small neural network to compute it based on  $a^{(t')}$  and  $s^{(t-1)}$   
→ a very small NN usually 1 hidden layer NN.

Intuition behind a small NN:

- Based on our architecture, we know  $s^{(t)}$  depends on  $s^{(t-1)}$  and  $a^{(t')}$
- But we don't know what is the function that computes each attention weight based on these values.
- So we use a small NN to compute this function.

- One downside of this algorithm is that it takes quadratic cost  
total #  $\alpha^{(t,t')}$  =  $T_x \times T_y$

Examples: Date Normalization?

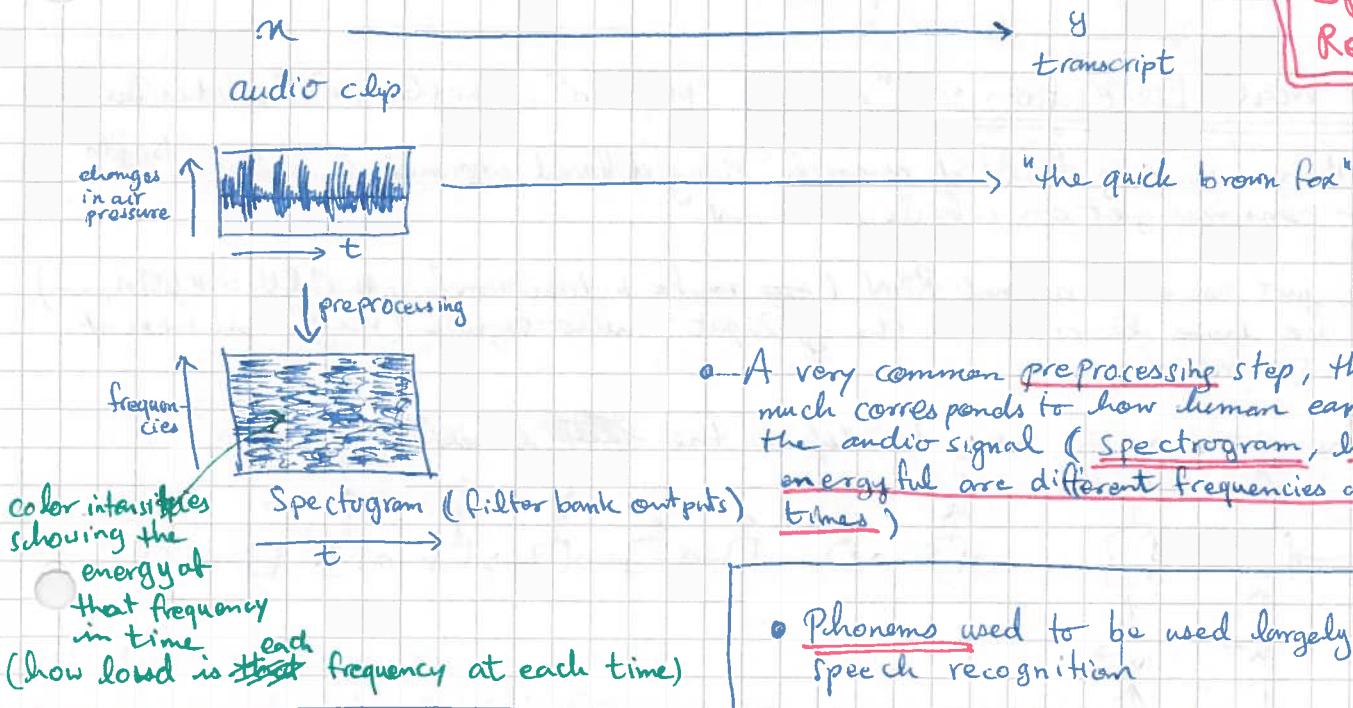
July 20th 1969 → 1969-07-20

23 April 1564 → 1564-04-23

MO.  $f_i: \{1, \dots, n\} \rightarrow \{1, \dots, n\}$  • Also, visualizing attention weights shows attention weights are high

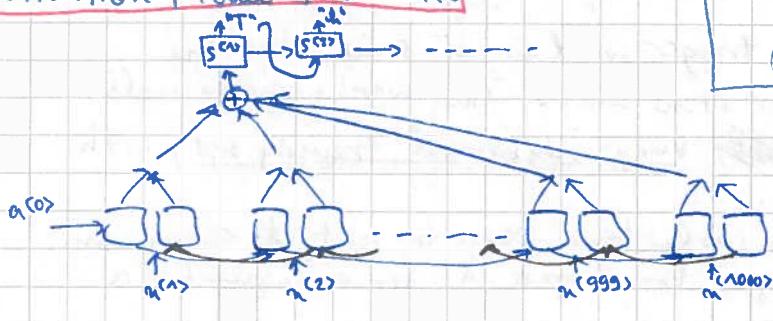
## Speech Recognition :

- The advances in sequence-to-sequence modelling has dramatically improved the Speech Recognition systems.



- A very common preprocessing step, that very much corresponds to how human ear processes the audio signal (spectrogram, how loud or energetic are different frequencies at different times)

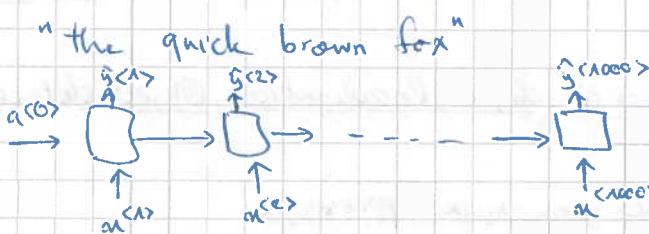
## Attention Model for SR:



- Phonemes used to be used largely for speech recognition
- But if we have enough data, now an E2E DL would perform better than hand-engineered phonemes as input.

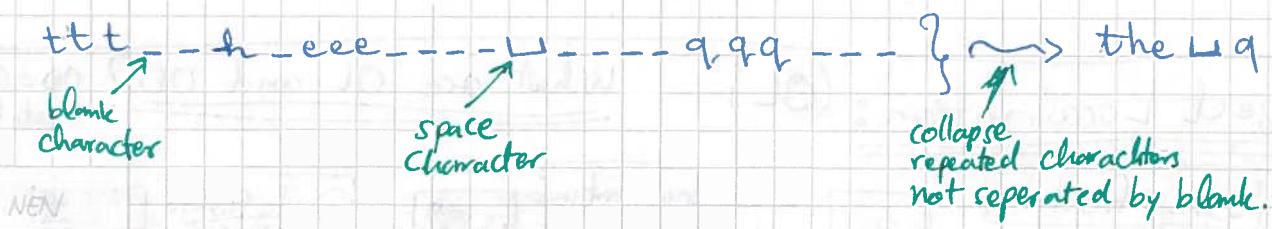
- Horizontal axis is single time frames of the input signal
- Attention model outputs the transcript.

## CTC cost for SR: (connectionist temporal classification)



- RNN with equal number of input/output.
- In practice a bidirectional RNN with LSTM or GRU and usually deeper

- Usually in SR, number of input time-steps is much bigger than the number of output time-steps. For example 10s of audio with 100 Hz  $\rightsquigarrow$  1000 input features.
- The CTC allows the NN to output a sequence like:

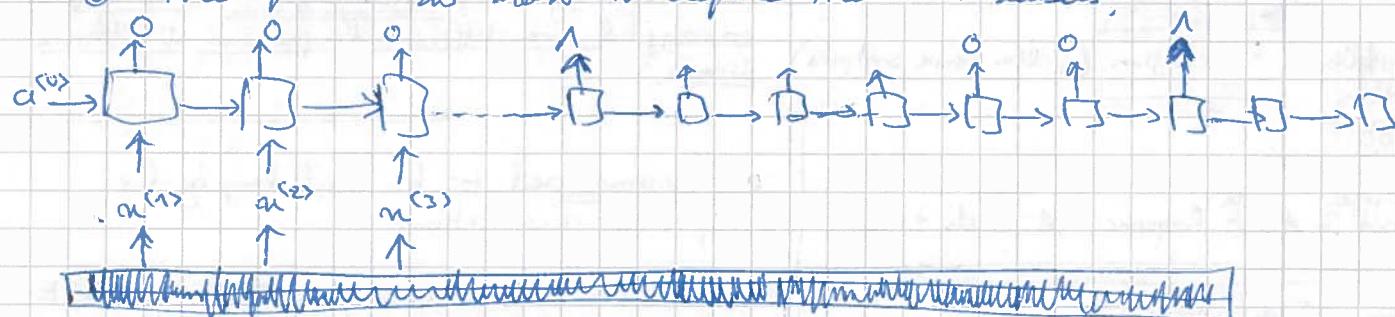


- Building an effective and production scale SR system is a pretty significant effort and requires a lot of data.
- Building a "Trigger word detection" System is much easier and requires much less amount of data.

## Trigger Word Detection: "Alexa"; "Hey Siri"; "Okay Google"; "xiadunihao"

↑ Baidu.

- Still an evolving field of research. Many different approaches. No consensus yet on which is the best.
- We just use a normal RNN (e.g. maybe bidirectional, with GRU, with LSTM, ... ) as we have learned, with the input audio signal (maybe in form of an spectro gram).
- The question is how to define the ~~labels~~ labels?



- We can use output 0 when trigger word is not said, set the output to one 1 as soon as it is said. This works pretty well, but has the problem of a ~~very imbalanced training set~~, with a lot more zeros than 1's.

→ Instead of having one single time-step to output 1, we can have several time steps to output 1, for example for a defined period of time.

## Course 4 - CNNs - Week 3 - Object Detection

1h34m video  
2h prog.

### Week 3

#### Learning Objectives:

- Understand the challenges of object localization, Object detection and Landmark Finding.
- Understand and implement non-max suppression
- Understand and implement intersection over union
- Understand how to label a dataset for object detection task.
- Learn the vocabulary of Object detection (landmark, anchor, bounding box, grid, ...)

## Object Localization: (OL)

What are OL and OD? OD = Object detection

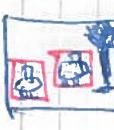
Classification  
1 Object



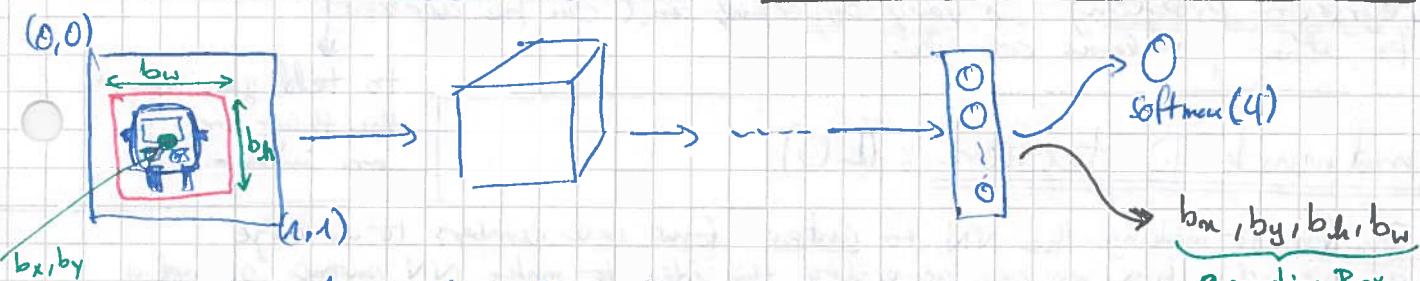
Classification with localization  
1 object and its location



Detection:  
multiple Objects and their location.



# Classification with Localization: At most only 1 object in image



- The standard pipeline for classification, say with 4 classes: 1. pedestrian, 2. car, 3. motorcycle, 4. background.
- In order to do OL (Object Localization) we extend this pipeline with few more output units that output a bounding box.
- Now the training set need to include the bounding box. In the example: e.g.  $b_x = 0.5, b_y = 0.6, b_h = 0.3, b_w = 0.4$ . These are not coordinates, but relative ratios.

## Defining the Target label $y$ :

1 - pedestrian  
2 - Car  
3 - Motorcycle  
4 - Background  $\leftarrow p_c = 0$

Example:

$$x = \begin{array}{|c|} \hline \text{[Image]} \\ \hline \end{array}, y = \begin{bmatrix} 1 \\ b_x \\ b_y \\ b_h \\ b_w \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

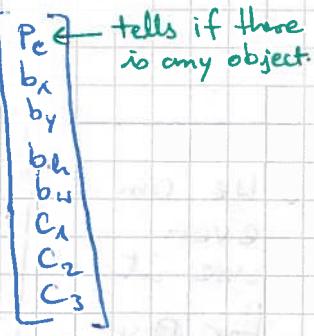
- We define output  $y$  as follows:

The training set needs to have bounding boxes now  $(b_x, b_y, b_h, b_w)$

$$x = \begin{array}{|c|} \hline \text{[Image]} \\ \hline \end{array}, y = \begin{bmatrix} ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \\ ? \end{bmatrix} \quad ? = \text{"don't care"}$$

## Object Localization

Output:  $b_x, b_y, b_h, b_w$ , Class label (1-4)



## Loss Function:

$$L(\hat{y}, y) = \begin{cases} (\hat{y}_1 - y_1)^2 + (\hat{y}_2 - y_2)^2 \\ + \dots + (\hat{y}_8 - y_8)^2 \end{cases} \quad \text{if } y = 1$$

In this case only  $\hat{y}_1$  matters, because all other elements of  $\hat{y}$  are "don't care"

if  $y = 0$

$$y = \begin{bmatrix} p_c \\ b_x \\ b_y \\ b_h \\ b_w \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}$$

logistic reg.  
loss  
squared error  
log likelihood  
of softmax

- Note that if  $y_1 = 1$ , we penalize deviations in all 8 components; whereas if  $y_1 = 0$  we only care about  $\hat{y}_1$  (the  $p_c$ ), and don't care about the rest.

- In this example we used squared error for all 8 components of  $y$ , just for simplification. In practice we normally use a log likelihood loss for  $C_1, C_2, C_3$  of the Softmax output, squared error for bounding box, and for  $p_c$  something like Logistic Regression loss.

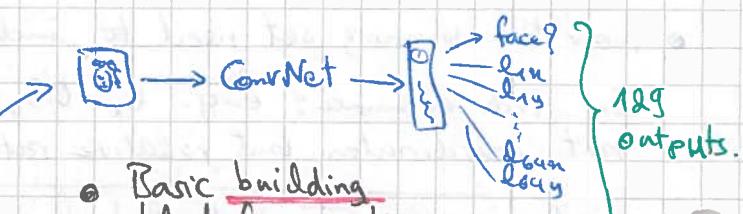
- The idea of having the NN output a bunch of real numbers (like a regression problem) is very powerful and can be applied to other problems as well.

to tell you where the things are in an image

## Landmark Detection : (LD)

- Similar to making the NN to output four real numbers to localize a bounding box, we can generalize this idea to make NN output  $n$  many coordinates of important points in image (the Landmarks)
- Make NN automatically output  $(x, y)$  coordinates of important points in image

### Example : Landmarks of a face



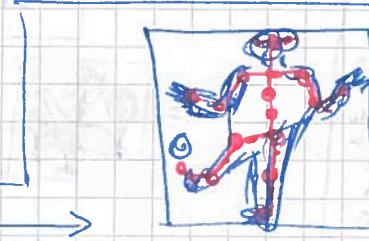
- Basic building block for emotion recognition and AR filters like Snapchat.

- We define let's say 64 keypoints in face

$l_{1x}, l_{1y}$   
 $l_{2x}, l_{2y}$   
 $\vdots$   
 $l_{64x}, l_{64y}$

- An train the NN to output the  $(x, y)$  coordinates of those points.

- We can even use it for Pose-detection



Amazing :)

**Landmark Detection**

## Object Detection : Sliding Windows

- First train a ConvNet with a data set of closely cropped images of car



- Then crop the input image using a sliding window, and feed every time the small cropped portion of the input to the trained classifier and hope that it will predict the cars.

- Change the window size and repeat the process.

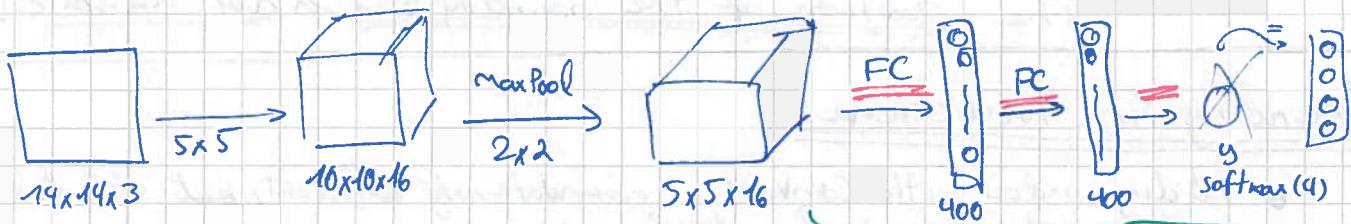
- Problem :** Computational cost. For a single image we have to repeat the classification task several times. Specifically in order to have good performance we have to choose a small stride

**Sliding Windows**

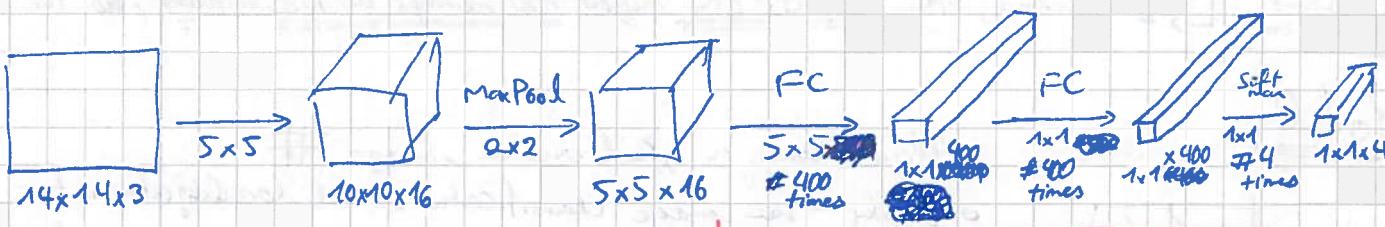
Solution  
A green checkmark icon.

## Convolutional Implementation of Sliding Windows :

## Turning FC Layers into Convolutional Layers:

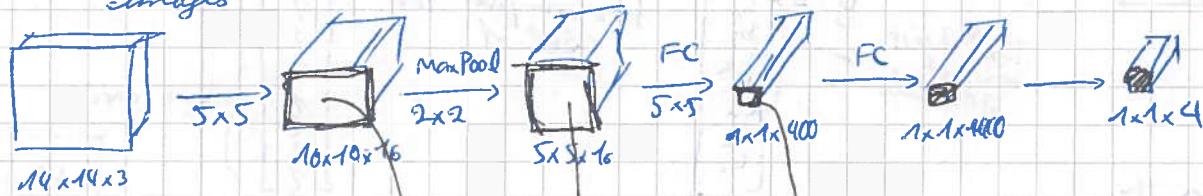


- We can turn the fully connected layers above into convolutional layers by using equivalent convolutional filters:



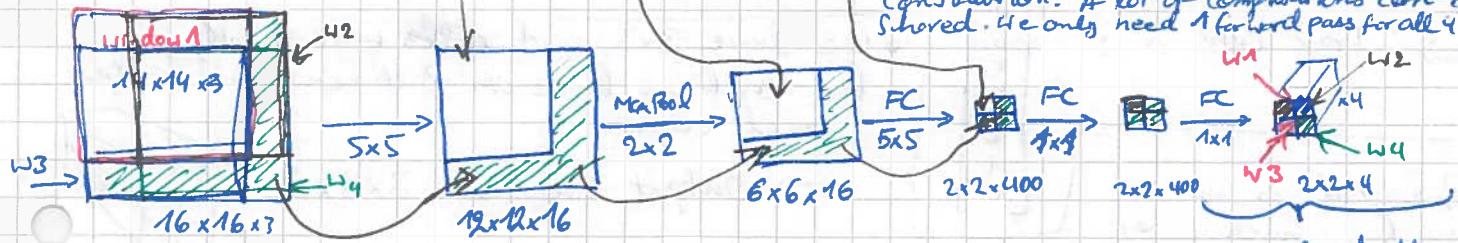
## Convolutional Inception of Sliding Windows:

- Let's assume we have trained our ConvNet for  $14 \times 14 \times 3$  images



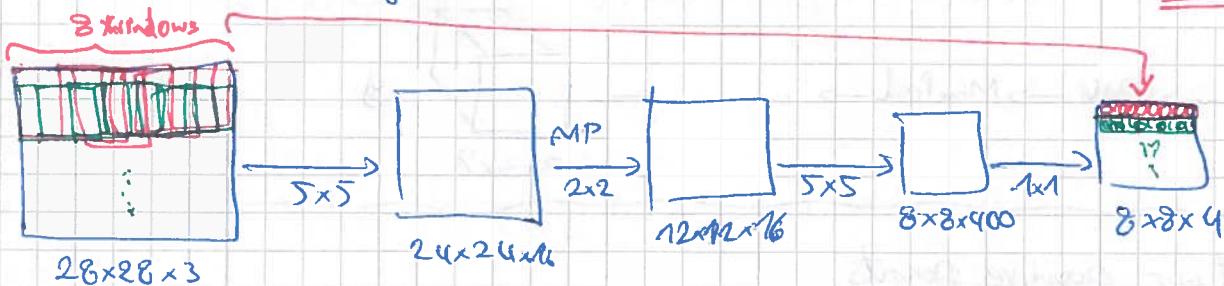
- Now we input a  $16 \times 16 \times 3$  Test image

Sermanet et al '14  
"SupFeat: Integrated recognition, localization, and detection using convolutional networks"



## For Another Example for a 28x28

- Another example:  $28 \times 28 \times 3$  input:



- Using Convolutional filters, we can do all the predictions of all the windows of an input image at the same time in one single forward Pass.

- Sliding Window with ConvNets has still a problem:  
The position of the bounding Box is not accurate!

## Bounding Box Prediction:

- Sliding window with ConvNet is computationally efficient, but still has the problem of inaccurate Bounding Boxes:

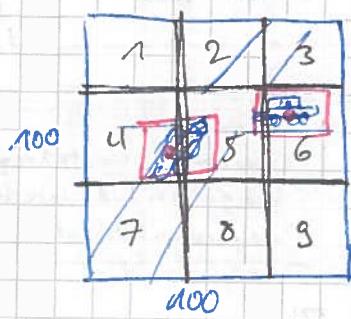
Redmon et.al '15

\* You only look once:

Unified real-time  
Object detection

- first, depending on the window size and stride, often, none of the Boxes really matches up with the ~~area~~ object)
- Second, the objects BB might not always be rectangular, as the windows

## YOLO:

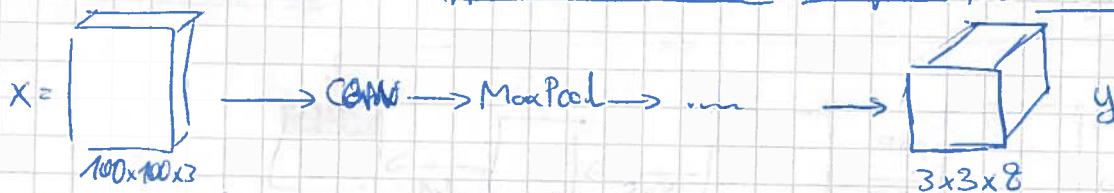


- YOLO considers the center point of the BB and assigns the object to the ~~BB~~ that contains the center point

⇒ For cell 3 :

- The grid is in practice much finer than 3x3. ⇒ It is not probable that a cell contains two objects.

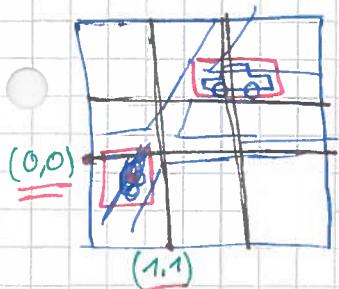
- Now we implement a typical ConvNet to output the desired output volume



Yolo has two positive points:

- Yolo is very much like the object localization algorithm, and outputs the BB coordinates explicitly ⇒ It can output BBs of any aspect ratio and is not restricted to a forced window size.  
↳ with very precise
- This is a convolutional implementation. The result for all grid cells is computed in one single forward pass.

How do we encode the BBs ( $b_x, b_y, b_w, b_h$ ) in a cell grid?



•  $b_x, b_y, b_w, b_h$  are defined relative to the grid cell.

• example :

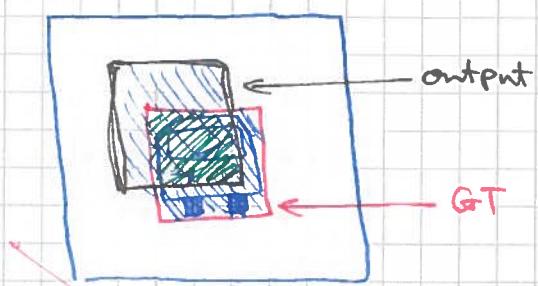
1
0.3
0.6
0.8
0.7
0
1
0

}  $b_x$  and  $b_y$  must be  $0 \leq b_x, b_y \leq 1$   
} can be  $\geq 1$ , for the case the  
BB crosses the cell boundaries. }!

How to measure the performance of an Object detection Algorithm?

Intersection over Union :

IoU



IOU = 1  
⇒ Perfect Match

IoU = Intersection over union of  
the two BBs

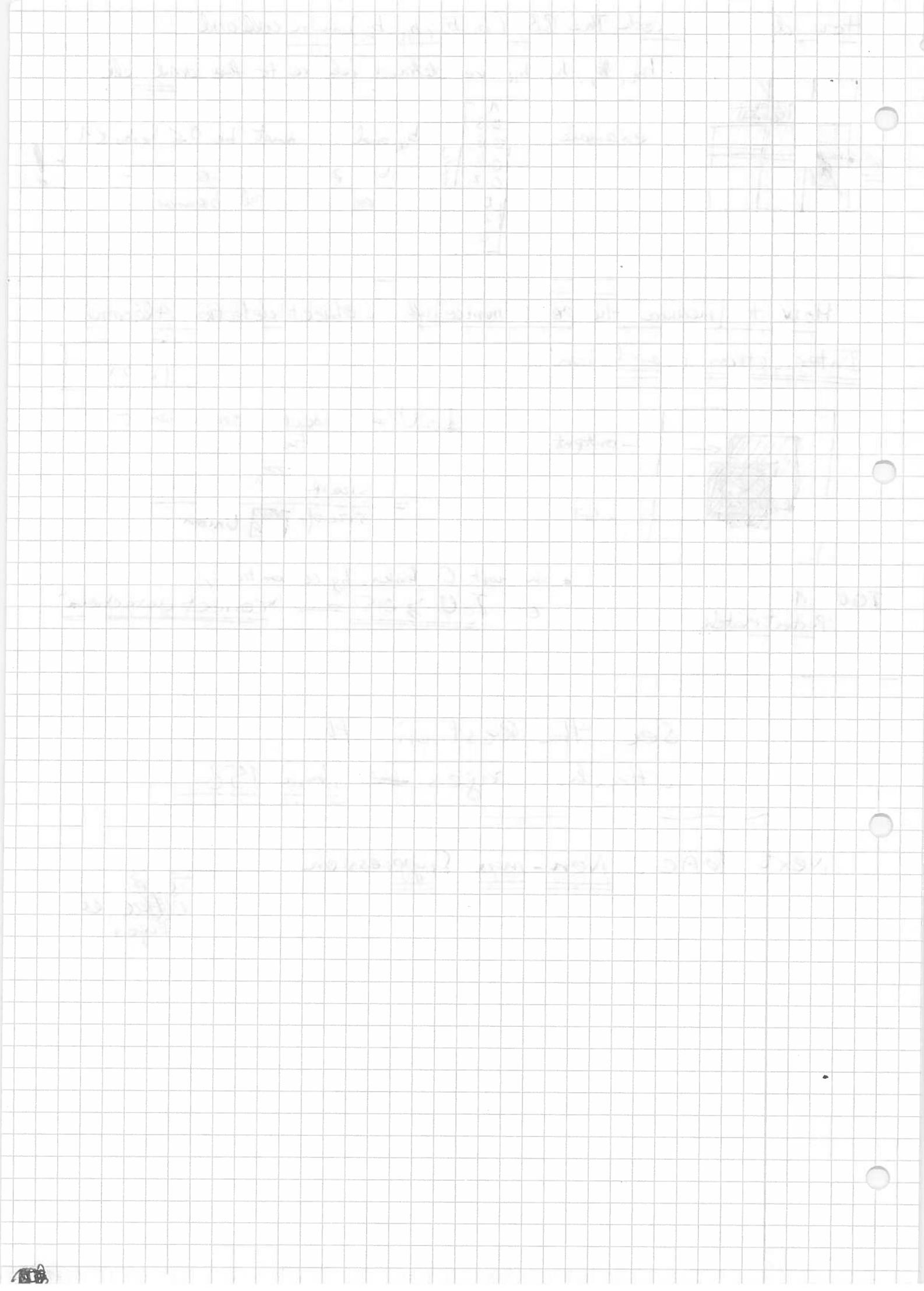
$$= \frac{\text{size of intersection}}{\text{size of Union}}$$

- in most CV tasks, by convention,  
if  $\underline{\text{IoU}} \geq 0.5 \Rightarrow \underline{\text{"Correct detection"}}$

See the Rest in the  
attached pages  $\Rightarrow$  from 158

Next topic : Non-max Suppression

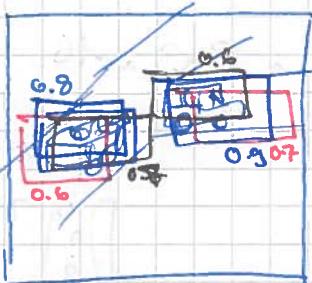
$\Rightarrow$   
see  
attached  
Pages



## Non-Max Suppression:

- So far our OD algorithm have the problem of multiple detection of the same object.

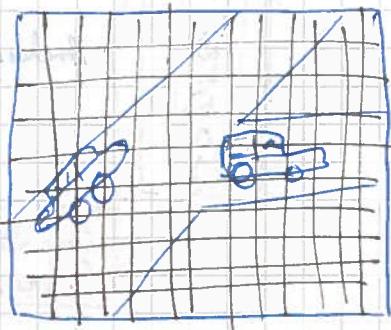
Example :



- Multiple cells might think that the center of the object is contained in them.  
I.e., multiple cells might have a P<sub>c</sub> high.

- ⇒ Take the BB with the highest probability, and suppress all others with IoU with it high
- Supress all overlapping rectangles.

Details :



19x19 grid

Each output prediction is

$$\begin{bmatrix} P_c \\ b_x \\ b_y \\ b_h \\ b_w \end{bmatrix}$$

• Assume we have only 1 class!

Non-Max Suppression

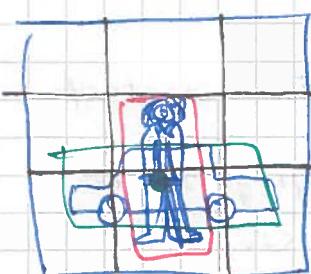
- Discard all boxes with  $P_c < 0.6$   
(Get rid of all low probability output boxes)
- While there are any remaining BBs:
  - pick the Box with the largest  $P_c$  and output that as prediction.
  - Discard any remaining BB with  $IoU \geq 0.5$  with the output box.

- For multiple classes, run separate Non-max suppression once for each of the classes

## Anchor Boxes:

- ↑ So far we have the problem that each of the grid cells can detect only one object.

## Overlapping Objects:



- Previously: each cell only one  $\Rightarrow y =$  object

Anchor box 1:



• now each cell multiple objects each corresponding to one anchor box

Anchor box 2:



Anchor box 1

$$y =$$

$$\begin{bmatrix} P_{c1} \\ b_{x1} \\ b_{y1} \\ b_{h1} \\ b_{w1} \\ C_1 \\ C_2 \\ P_{c2} \\ b_{x2} \\ b_{y2} \\ b_{h2} \\ b_{w2} \\ C_3 \\ C_4 \end{bmatrix}$$

Anchor box 2

## ← Anchor Box Algorithm

### Previously:

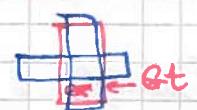
Each object in training image is assigned to grid cell that contains that object's midpoint

→ Output  $y$ :  $3 \times 3 \times 8$

With two anchor boxes:

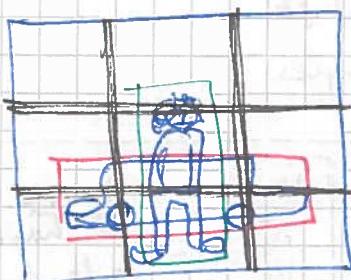
Each object in the training set is assigned to grid cell that contains objects midpoint and anchor box for the grid cell with highest IoU

⇒  $(\text{grid}, \text{cell}, \text{anchor box})$



⇒ Output:  $3 \times 3 \times 16 \approx (3 \times 3 \times 2 \times 8)$

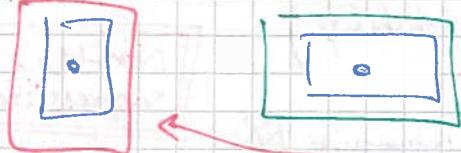
### Example



$y =$

Pc bx by bw bh Cn C2 Cs Pc bx by bw bh Cn C1 C2 Cs	=	<table border="1"> <tr><td>1</td><td>bx</td><td>by</td><td>bw</td><td>bh</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>1</td><td>bx</td><td>by</td><td>bw</td><td>bh</td><td>1</td><td>0</td><td>0</td></tr> <tr><td>1</td><td>bx</td><td>by</td><td>bw</td><td>bh</td><td>0</td><td>1</td><td>0</td></tr> </table>	1	bx	by	bw	bh	0	0	0	1	bx	by	bw	bh	1	0	0	1	bx	by	bw	bh	0	1	0
1	bx	by	bw	bh	0	0	0																			
1	bx	by	bw	bh	1	0	0																			
1	bx	by	bw	bh	0	1	0																			

Anchor box 1: Anchor box 2:



Car only?

0	0	0	0	0	0	0	0
---	---	---	---	---	---	---	---

Anchor box 1

1	bx	by	bw	bh	0	1	0
---	----	----	----	----	---	---	---

Anchor box 2

Anchor Boxes

! Note: Two cases are not handled well with this algorithm:

- If we have, say 2 Anchorboxes, but three overlapping objects.
- If we have two objects, that correspond both to the same anchor box shape

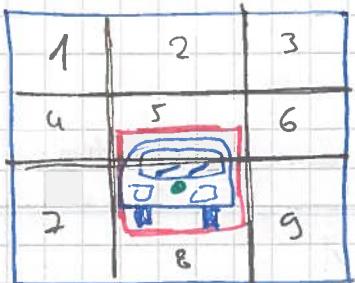
• Generally, in practice, with a fine grid (say  $19 \times 19$ ) the chance of two objects overlapping is not that high.

• Anchor boxes allow your algorithm to specialize better. Particularly if your dataset has tall objects and wide objects.

### How to choose Anchor boxes?

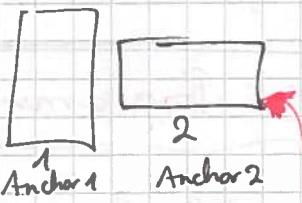
- Normally chosen by hand, 5 or 10 shapes, that represent a variety of shapes of objects.
- Even better is to use a k-Means algorithm to group together different shapes of objects and choose anchor box shapes to cover a wide range of object groups.

## Yolo Algorithm :



### 3 Classes

- 1 - Pedestrian
- 2 - car
- 3 - Motor cycle



$y$  is  $3 \times 3 \times 2 \times 8$

grid      #anchors       $5 + \# \text{ classes}$

$P_c, b_x, b_y, b_w, b_h$

$$y =$$

$P_c$	0
$b_x$	?
$b_y$	?
$b_w$	?
$b_h$	?
$C_1$	?
$C_2$	?
$C_3$	?

$$y_{\text{cell } 1} =$$

0
2
?
?
?
?
?
?

$$\leftarrow P_{C_1}$$

0
?
?
?
?

$$\leftarrow P_{C_2}$$

$$y_{\text{cell } 8} =$$

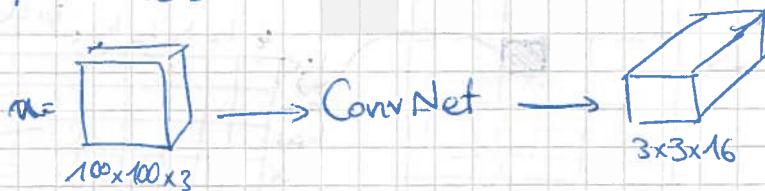
0
1
?
?
?

$$\leftarrow P_{C_3}$$

Anchor 1 has higher  
prob. than Anchor 2.

so with the GT.

- Now lets train a ConvNet :



- Making Predictions as you expect

- Outputting the non-max suppressed outputs :

- for each grid cell, get 2 predicted bounding boxes (for two anchors)
- Then get rid of low probability predictions.
- for each of the classes, independently run the non-max suppression.

## Region Proposal :

- quite influential in CV

Girshick et al. '13;  
"Rich feature hierarchies for accurate object detection and semantic segmentation"

- Many areas of an input image are often empty. There is nothing interesting in them.

$\Rightarrow$  Run the convnet only on selected regions

## Region Proposal

### R-CNN

- Use a segmentation algorithm to find proposed regions (blobs)
- Run the classification on the proposed regions one at a time

} R-CNN

$$\rightarrow$$

## ← • Fast R-CNN:

- Propose regions. Use convolution implementation of sliding windows to classify all the proposed regions.

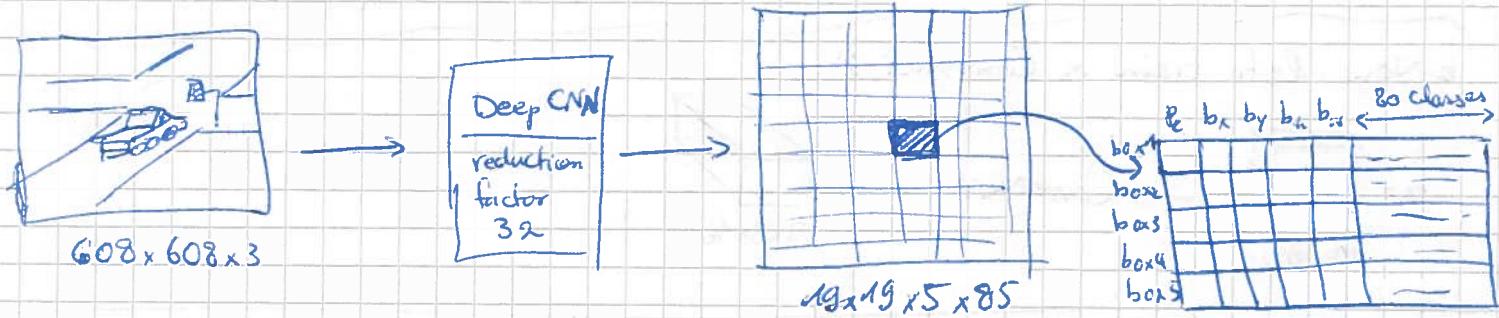
## • Faster R-CNN:

- Use conv nets to propose regions instead of a classical algorithm.

## Course 4 - CNNs - Week 3 - Programming 1

### Car Detection with Yolo: 80 classes, 5 Anchors

- Input : Batch of images of shape  $(m, 608, 608, 3)$
- Output : list of BBs with recognized classes. Each BB has is represented by 6 numbers ( $p_c, b_x, b_y, b_w, b_h, c$ ). Classes are expanded into a 80 dim one-hot vector  $\Rightarrow$  Each BB is a 85 dim vector.



- We flatten the last two dimension (for simplicity)  
 $\Rightarrow$  Output ~~BBs~~ shape  $19 \times 19 \times 425$

### Summary of Yolo :

- Input image  $(608, 608, 3)$
- The input image goes through a CNN, resulting in a  $(19, 19, 5, 85)$  dimensional output.
- flattening the last two dimensions of the output  $\rightarrow (19, 19, 425)$
- Each cell ~~over~~ in a  $19 \times 19$  grid over the input image gives 425 numbers.
- $425 = 5 \times 85$ . Each cell contains predictions for 5 boxes, corresponding to 5 Anchorboxes
- $85 = 5 + 80$ , 5 because of  $(p_c, b_x, b_y, b_w, b_h)$  and 80 then number of classes.
- We then ~~select~~ select only few boxes based on:
  - Score thresholding: throwaway boxes that have detected a class with a score less than the threshold
  - Non-max Suppression: Compute the intersection over union and avoid selecting overlapping boxes.
- This gives us Yolo's final output.

# Course 4 - CNNs - Week 4 - Special Applications

## Week 4

### Face Recognition & Neural Style transfer

#### What is Face Recognition?

- A Demo video from Baidu's HQ in China.
- Face Recognition      }  
 • liveness detection      } opens entrance gate

#### Face Verification vs. Face Recognition:

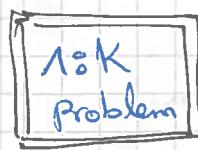
##### Verification:



- Input: Image, name/ID
- Output: Whether the input image is that of the claimed person.

We build a very accurate face verification system, and use that for face recognition.

##### Recognition:



- Has a Database of K persons
- gets an input image
- Outputs ID if the image is one of the K persons (or "not recognized")

↳ much harder than verification problem.

- ↳ Imagine you have a FV system with 99% accuracy
- Now you apply it to a dataset with 100 persons => you have 100 times 1% chance of making a mistake =>
- ! → • Your Face Verification System needs to be very very (99,9%) accurate to be used for Face Recognition

- One of the main problems that needs to be solved is "One-Shot Learning"

#### One Shot Learning:

- You need to be able to recognize a person given only one sample or one example of that person's face!
- Historically DL algorithms don't work well, if you only have one training example.



## One-Shot learning :



- Learning from one example to recognize the person again.



↳ Because we might have only one picture of the person in our database.



- ~~Not working solution~~



- It does not work:

- we don't have enough data to train it robustly.
- What if a new person is added to the DB?

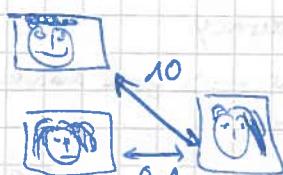
↳ Do we have to train the CNN over again?

⇒ Let's learn a "Similarity function"

$$d(\text{img}^1, \text{img}^2) = \text{degree of difference between two images}$$

- $d$  small for two images of one person.

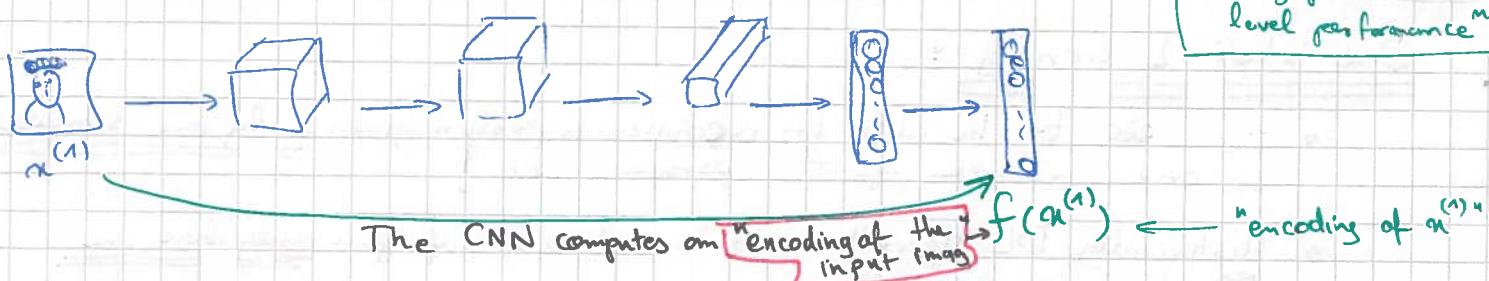
If  $d(\text{img}^1, \text{img}^2) \leq T \rightarrow \text{"same person"}$  } Verification.  
 $> T \rightarrow \text{"different"}$



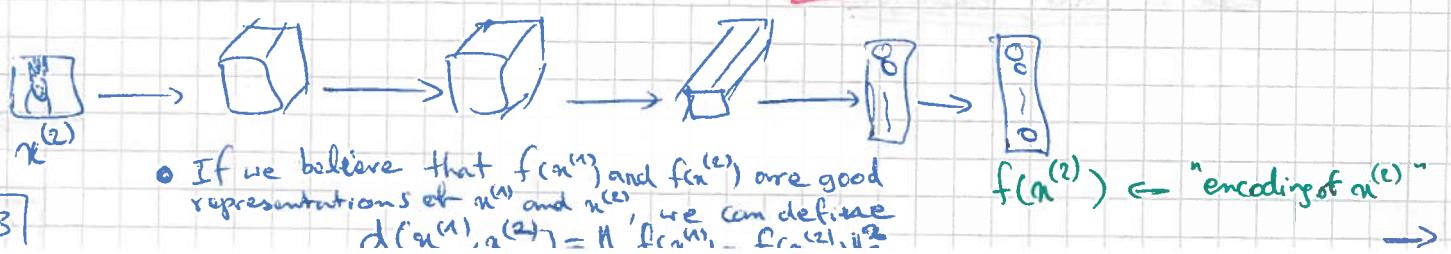
- Compare the new input image with all the images of the DB and compute the  $d$ .
- If now a new person is added to the DB, we don't need to learn anything new.

## Siamese Network :

- ~~A good way to calculate the  $d()$  function~~



Taigman et al '14,  
"DeepFace: closing the gap for human-level performance"



- ← • The idea of running two identical CNNs on two different inputs and then comparing them, is called a Siamese Network.

## ② How to train Siamese NN's?

↳ The two NNs in previous page have same parameters

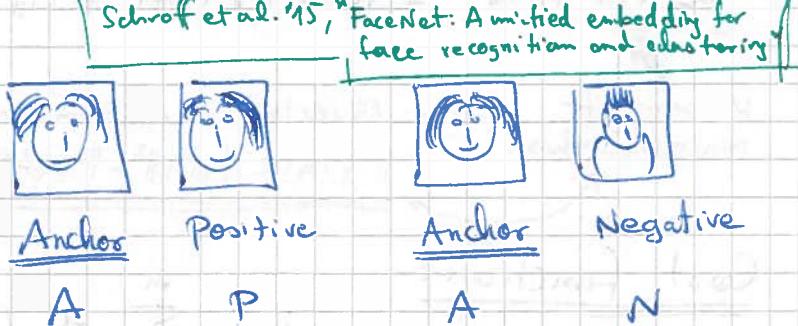
→ So we want to train a neural network, so that the encoding it computes results in a function  $d$  that tells us when two pictures are of the same person.

- Parameters of NN define an encoding  $f(x^{(i)})$  ↗ say a 128 dim vector.
- Learn Parameters so that:
  - If  $x^{(i)}, x^{(j)}$  are the same person,  $\|f(x^{(i)}) - f(x^{(j)})\|^2$  is small
  - If  $x^{(i)}, x^{(j)}$  are different persons,  $\|f(x^{(i)}) - f(x^{(j)})\|^2$  is large.

## Triplet Loss :

### Learning Objective :

- We want the encoding between Anchor and positive to be very similar
- Whereas the encoding of anchor and negative should be far from each other.



⇒ Term, Triplet loss, because at any time we are looking at three images (A, P, N)

### What we want:

$$\underbrace{\|f(A) - f(P)\|^2}_{d(A, P)} \leq \underbrace{\|f(A) - f(N)\|^2}_{d(A, N)}$$

$$\Rightarrow \|f(A) - f(p)\|^2 - \|f(A) - f(N)\|^2 \leq 0$$

! Problem :  $f(\text{img}) = \vec{0}$  will trivially satisfy this equation; or  $f(\text{img}_1) = f(\text{img}_2)$

→ To make sure our NN does not learn these trivial useless encodings modify the objective function:

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha \leq 0$$

$\alpha$  is now the margin

↓ Similar to margin in SVM. 164

- now let say the margin is  $\alpha = 0.2$   $\rightarrow$

~~$$\|f(A) - f(P)\|^2 \leq \|f(A) - f(N)\|^2 + \alpha$$~~

$$\Rightarrow |d(A, P) - d(A, N)| \geq 0.2$$

- Margin defines the gap between distances of the anchor and the Positive and the distance of the anchor and those negative.

- The margin pushes the  $(A, P)$  pair and the  $(A, N)$  pair, further away from each other.

### Loss function :

Given 3 images  $A, P, N$  :

$$l(A, P, N) = \max(\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha, 0)$$

We want to minimize this. | The Objective function was :

$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha \leq 0$$

### Cost function :

$$J = \sum_{i=1}^m l(A^{(i)}, P^{(i)}, N^{(i)})$$

- So, if we have a Training set of 10,000 pictures of 1000 persons, we generate triplets from the 10000 pictures and then train the algorithm to minimize that cost function.
- For the purpose of training the system, we do need pairs of  $A$  and  $P$ . Pairs of pictures of the same person.
- But after training, we can use the network for doing one-shot learning.

### Choosing the Triplets, $A, P, N$

- During training if  $A, P, N$  are chosen randomly,  $d(A, P) + \alpha \leq d(A, N)$  is easily satisfied.

Because if we choose  $A$  and  $N$  randomly, chances are high that they are very different and so the  $\|f(A) - f(P)\|^2 + \alpha \leq \|f(A) - f(N)\|^2$  is easily satisfied and the NN does not learn much from it.

⇒ We want to choose triplets that are "hard" to train on.

⇒ That is choosing triplets for which  $d(A, P) \approx d(A, N)$

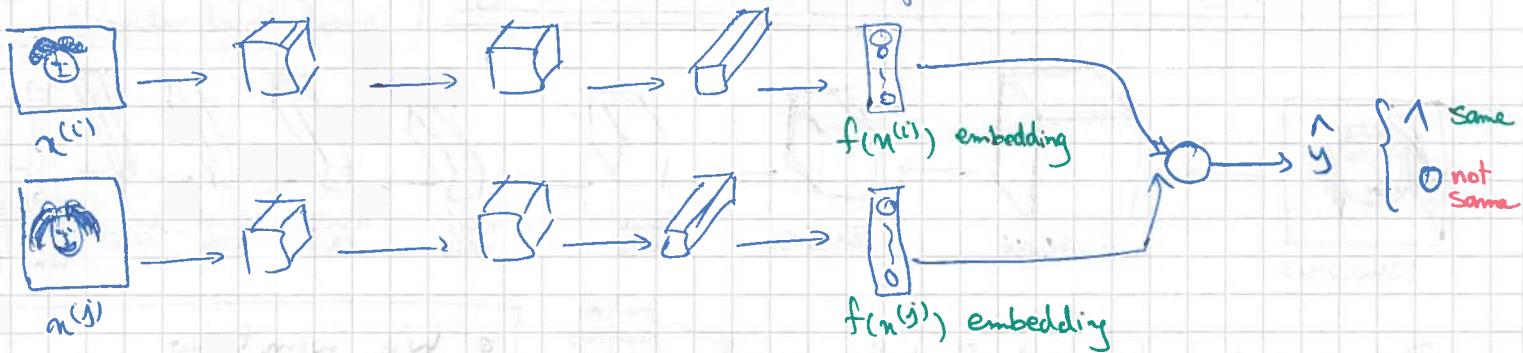
so that the NN learns to push down push up.

- If the problem is too easy, the GD does not have much to do, and the algorithm does not learn well.

## Face Verification and Binary Classification:

- Posing FR problem as a binary classification problem.

- This is an alternative to the triplet loss



- Feed the  $f(m^{(i)})$  and  $f(m^{(j)})$  embeddings into a logistic regression unit to make the prediction. The target output = 1 if same person, and = 0 if ~~different~~ different person.

$$\hat{y} = \sigma \left( \sum_{k=1}^{128} w_i |f(m^{(i)})_k - f(m^{(j)})_k| + b \right)$$

element-wise difference of embedding vectors

other formulas also possible

$$\begin{cases} X^2 \text{ similarity} \\ (f(m^{(i)})_k - f(m^{(j)})_k)^2 \\ \frac{f(m^{(i)})_k + f(m^{(j)})_k}{2} \end{cases}$$

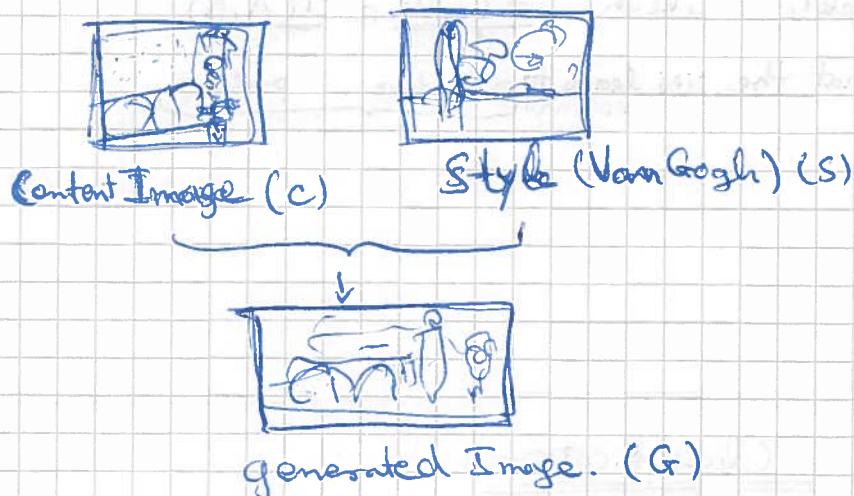
- The input is a pair of images → The output is 0 or 1

- We train a Siamese Network again.

- Computational Tip: For the images in the DB, we can have already the embedding computed, and saved.

Then we only need to compute the embedding for the new input image.

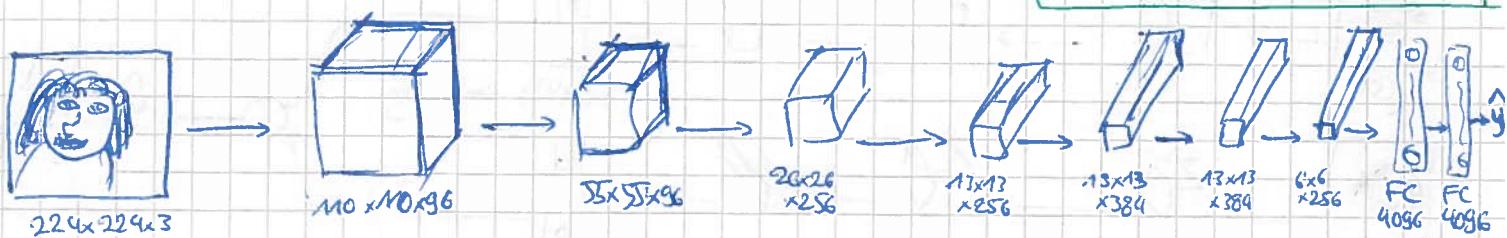
## Neural Style Transfer :



### Important:

- Neural Style transfer is a specific application of supervised learning in which
  - The network weights does not change. They don't get optimized.
- We use a pre-trained network, trained on a dataset of images. This network is capable of detecting features. ← We use the activations of this network (i.e. the features detected) to optimize our generate image

## Visualizing what a Deep network is learning



- Pick a unit in layer 1. Find the nine image patches that maximize the unit's activation.

↳ Pass your training set through the NN and find out what is the image that maximizes that particular unit's activation.

Important lecture, but hard to understand

- We want to visualize what the hidden ~~layer~~ units in different layers are computing

Earlier layers look at smaller patches of the input image and find more fine grained features like lines. → See video

Later layers look at larger patches of the input image and find more complex features, like faces and dogs

## Neural Style Transfer Cost function :

$J(G)$  : measure how good is a generate image G

↳ We separate  $J(G)$  into two components

$$J(G) = \alpha J_{\text{content}}(G, C) + \beta J_{\text{style}}(G, S)$$

How "Similar" is G to the input content image.

How "Similar" is G to the input Style image.

We also define two Hyperparameters  $\alpha$  and  $\beta$  to specify the relative weighting between the content and style costs.

Gatys et al '15; "A Neural Algorithm of artistic style"

## Finding the Generated Image G :

- 1. Initialize  $G$  randomly

$$G : 100 \times 100 \times 3$$

- 2. Use Gradient descent to minimize  $J(G)$

$$G := G - \frac{\partial}{\partial G} J(G) \quad \left. \right\} \rightarrow \text{updating the pixel values of the image } G$$

## Content Cost Function :

$$J(G) = \alpha J_{\text{content}}(C, G) + \beta J_{\text{style}}(S, G)$$

- Say we use ~~the~~ hidden layer  $l$  to compute content cost.

- if  $l$  is small, say first layer, then the generated image will be very much like at pixel level like the input content image.
- If  $l$  is very deep, then the layer figures out there is a dog in your input content image and says, well, let's put a dog somewhere in the output generated image.
- So we want  $l$  to be somewhere in between. Neither too shallow nor too deep.

- Use a pre-trained ConvNet. (E.g. VGG network)

→ We want to measure, given a content image and a generated image  $G$  How similar are their content.

- Let  $a^{[l](C)}$  and  $a^{[l](G)}$  be the activation of layer  $l$  on the images

- If  $a^{[l](C)}$  and  $a^{[l](G)}$  are similar both images have similar content.

$$J_{\text{content}}(C, G) = \| a^{[l](C)} - a^{[l](G)} \|_2$$

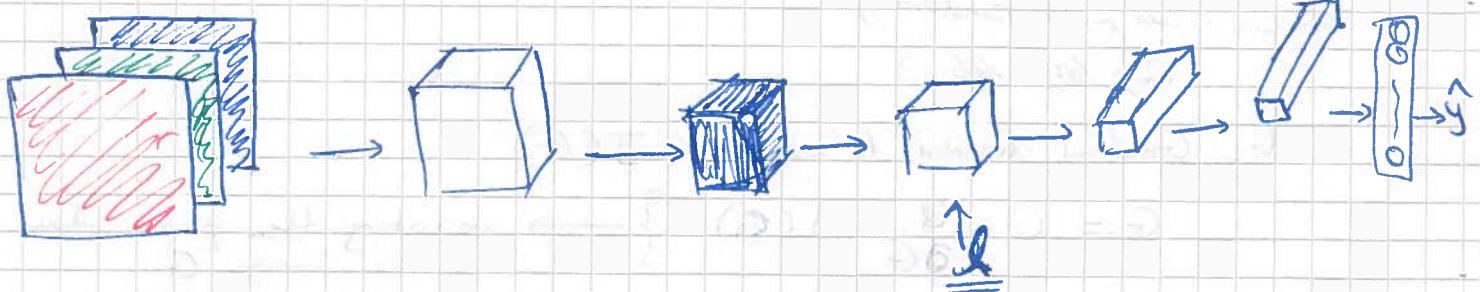
We can also have a normalization factor, say  $\frac{1}{2}$  there. But it does not matter

unroll both activations into vectors and calculate their difference

## Style Cost Function :



## Meaning of the "Style" of an Image :



- Say we are using layer l's activation to measure "Style"
- We define Style as correlation between activations across different channels.

Correlation between activations means, whenever part of the image that fires that activation is present, ~~is~~ another part of the image that fires the other activation is also present

### See Video

- The correlations tell us, which of the high level features/textures/forms/lines/shapes/ of the image occur together. (Like, whenever you have vertical narrow lines, you have an orange like color texture)

"Gram matrix"

Style Matrix : Captures all the correlations between activations across ~~all~~ different channels.

- Let  $a_{i,j,k}^{[l]}$  = activation at  $(i,j,k)$ .  $G^{[l]}$  is  $n_c^{[l]} \times n_c^{[l]}$

$\begin{matrix} H \\ W \\ \text{position} \\ (i,j) \end{matrix}$   $\begin{matrix} C \\ \text{channel} \\ k \end{matrix}$

Style matrix

Because we want to measure how correlated are each pair of channels (pairwise correlations between channels)

- $G_{kk'}$  measures how correlated are the activations in channels  $k$  and  $k'$ ;  $k, k' = 1, \dots, n_c$

$$G_{kk'}^{[l]} = \sum_{i=1}^{n_H^{[l]}} \sum_{j=1}^{n_W^{[l]}} a_{ijk}^{[l]} \cdot a_{ijk'}^{[l]}$$

Technically this is not correlation, but the unnormalized cross-covariance. But fulfills the same purpose

- We compute  $G^{[l]}$  for both the input style image and the generate image  $G$

$$\begin{aligned} J_{\text{style}}^{[l]}(S, G) &= \frac{1}{(n_H^{[l]} n_W^{[l]} n_C^{[l]})^2} \| G^{[l](S)} - G^{[l](G)} \|_F^2 \\ &= \frac{1}{n_H^{[l]} n_W^{[l]} n_C^{[l]}} \sum_i \sum_j \sum_k (G_{ijk}^{[l](S)} - G_{ijk}^{[l](G)})^2 \end{aligned}$$

i.e. correlated

- If both activations are large together, then  $G^{[l]}$  is large for these two channels, and
- If activations are not both large i.e. uncorrelated  $G^{[l]}$  will be small.

## Style Cost Function:

$$\| G^{[l](s)} - G^{[l](G)} \|_F^2$$

$$J_{\text{style}}^{[l]}(S, G) = \frac{1}{(2n_H^{[l]} n_W^{[l]} n_C^{[l]})^2} \sum_k \sum_{k'} (G_{kk'}^{[l](s)} - G_{kk'}^{[l](G)})^2$$

- It turns out that we get more visually pleasing results, if we use the style cost function from multiple different layers.

So the overall Style Cost Function can be defined as:

$$J_{\text{style}}(S, G) = \sum_l \lambda^{[l]} J_{\text{style}}^{[l]}(S, G)$$

↑  
weighting factor for style cost function  
of different layers.

- This allows us to take into account style cost of different layers in CNN
  - earlier ones, considering lower level features like edges
  - and later ones, considering some higher level features.

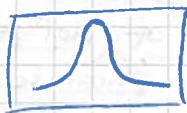
## Convolutions on 1D and 3D Data :

- Exactly as we apply convolutions on 2D data (like images), we can also apply a convolution on 1D or 3D data

Example : EKG



\*



$$\xrightarrow{\quad 14 \times 1 * 5 \times 1 \quad} 10 \times 16$$

16 filters

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----

1D signal  
say 14 \*  
size

\*

1	3	10	3	1
---	---	----	---	---

1D filters  
filter size 5

$$\xrightarrow{\quad 10 \times 16 * 5 \times 16 \quad} 6 \times 32$$

32 filters

All the same ideas as we saw in 2D data also apply here in 1D data

3D data :



\*



3D data Volume

$$14 \times 14 \times 14 * 5 \times 5 \times 5 \xrightarrow{\quad (\# \text{filters}) \quad} 10 \times 10 \times 10 \times 16$$

$$\xrightarrow{\quad (\# \text{filters}) \quad} 6 \times 6 \times 6 \times 32$$

- CT scan, CAD 3D models
- Video, motion data.

# Course 4 - CNNs - Week 4 - Programming 1

## Art Generation with Neural Style Transfer :

- In Neural Style Transfer, we optimize a cost function to get pixel values. (as opposed to other approaches, where we optimized a cost function to get parameter values)
  - NST can be considered as a transfer learning task. It uses another pre-trained network and builds on top of that.
  - We use VGG-19 (19 layers) trained on the very large ImageNet database.
    - ↳ This model has already learned to recognize a variety of low-level features (at earlier layers) and high-level features (at later layers)
  - We load the weights in a dictionary. Each layer.
- model["input"].assign(image) ← feed the image through network
- After that
- To get the activation of a layer : sess.run(model["conv4\_2"])
- NST in three steps :
    - ↳ Build the content cost  $J_{content}(C, G)$
    - Build the style cost  $J_{style}(S, G)$
    - Overall cost  $J(G) = \alpha J_{content}(C, G) + \beta J_{style}(S, G)$

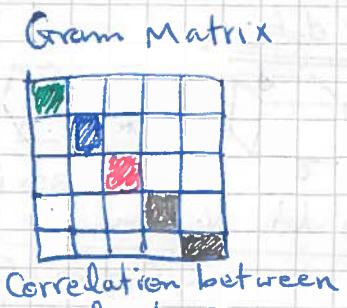
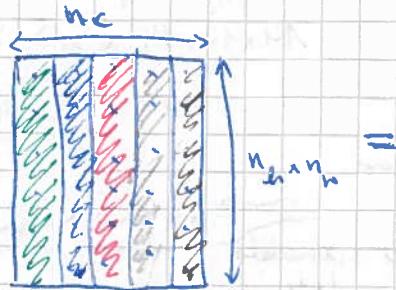
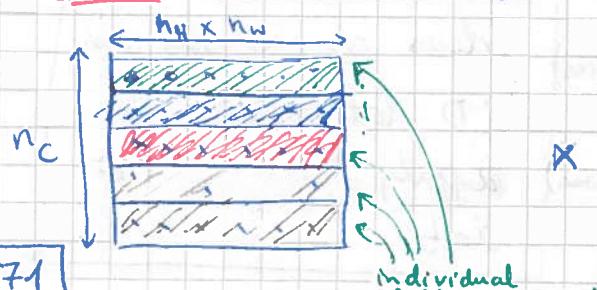
tf.transpose(), tf.reshape(), tf.reduce\_sum(), tf.square()

Gram Matrix : Gram matrix<sup>G</sup> of a set of vectors  $(v_1, \dots, v_n)$  is the matrix of dot products of vectors pairwise

$$G_{ij} = v_i^T v_j = \text{np.dot}(v_i, v_j)$$

In other words  $G_{ij}$  computes how similar are  $v_i$  and  $v_j$ . If  $v_i$  and  $v_j$  are similar, then their dot product would be high.

In NST we can compute the Style matrix by multiplying the "unrolled" filter matrix with its transpose



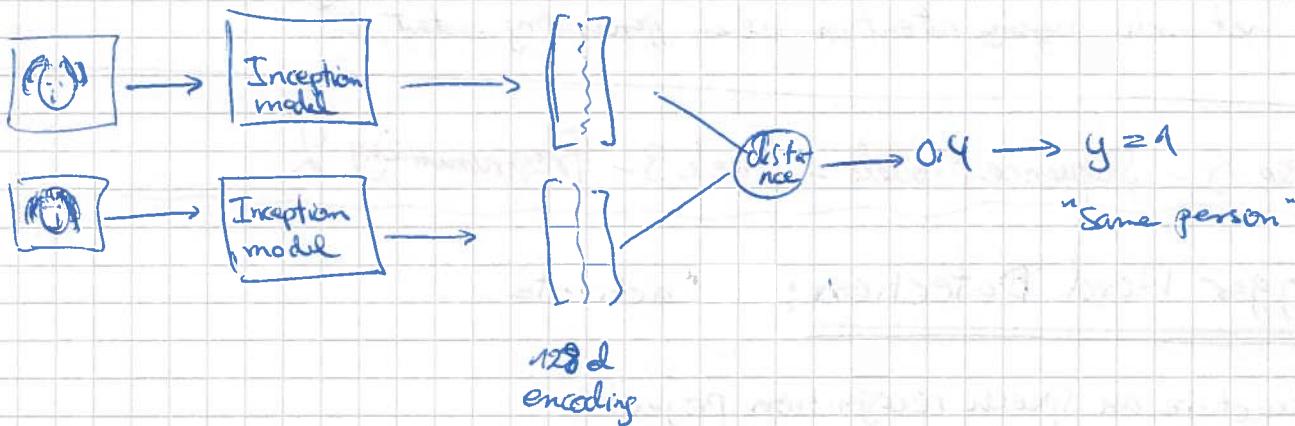
refer to  
notebook

- The resulting
- Gram matrix is a  $n_c \times n_c$  matrix. The  $G_{ij}$ 's measure how similar the activation of filter  $i$  is with the activation of filter  $j$
  - The diagonal of the  $G$  matrix ( $G_{ii}$ 's) is a measure of how active filter  $i$  is. For example, let's say filter  $i$  is ~~not~~ detecting vertical textures in the image; Then  $G_{ii}$  measures how common are vertical textures in the image as a whole. If  $G_{ii}$  is large, this means the image has a lot of vertical texture.
  - By capturing the prevalence of different types of features ( $G_{ii}$ ) as well as how much different features occur together ( $G_{ij}$ ) the style matrix  $G$  measures the style of an image

## Course 4 - CNNs - Week 4 - Programming 2

### Face Recognition

- Implement the Triplet loss function
- Use a pretrained model to map face images into 128-dimensional encodings
- Use encodings to perform face verification and face recognition.
- We use Inception model and channel first for input data.



### Triplet loss:

$$J = \sum_{i=1}^m \left[ \| f(A^{(i)}) - f(P^{(i)}) \|_2^2 - \| f(A^{(i)}) - f(N^{(i)}) \|_2^2 + \alpha \right]_+$$

$$[z]_+ = \max(z, 0).$$

## Cours 5 - Sequence models - Week 3 - Programming 1

### Neural Machine Translation (NMT) with Attention Model:

- Translating Human readable dates into machine readable dates :
  - "the 29th of August 1958" → 1958-08-29
  - "03/30/1968" → 1968-03-30
  - "24 JUNE 1987" → 1987-06-24
- Dataset 10000 pairs ("9 May 1998", "1998-05-09")
- human vocab : Python dict mapping all characters of human readable dates to int.
- machine vocab : Python dict = all characters of machine readable dates to int.

$T_x = 30$  max length of human readable date

$T_y = 10$  (YY-MM-DD)

- Machine translation models can map one sequence to another. They can be used for tasks like date format translation too
- We can visualize attention weights  $\alpha(t, t')$  to know to what input  $t'$  is the network paying attention when generating output  $t$ .

## Course 5 - Sequence Models - Week 3 - Programming 2

### Trigger Word Detection: "activate"

- Structure an speech recognition project
- Synthesize and process audio recordings to create train/dev datasets.
- Train a trigger word detection model.
- Data set → mixture of positive / negative words
  - different background noise
  - different accents.
- Synthesizing data also allows us to easily set the ground truth labels. →

- ← Input:  $5511 = T_x$  (spectrogram output)
- Output:  $1375 = T_y$
- Output  $\approx 1$  exactly right after someone finished saying "activate" time-step + 50 further timesteps.

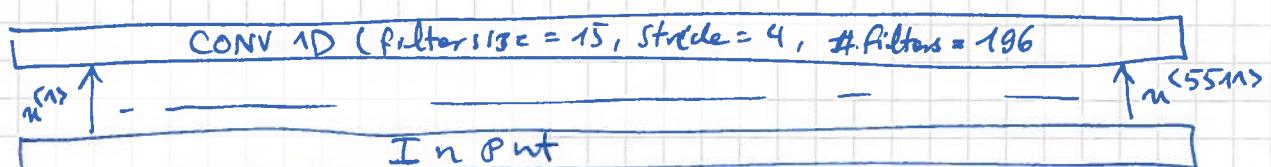
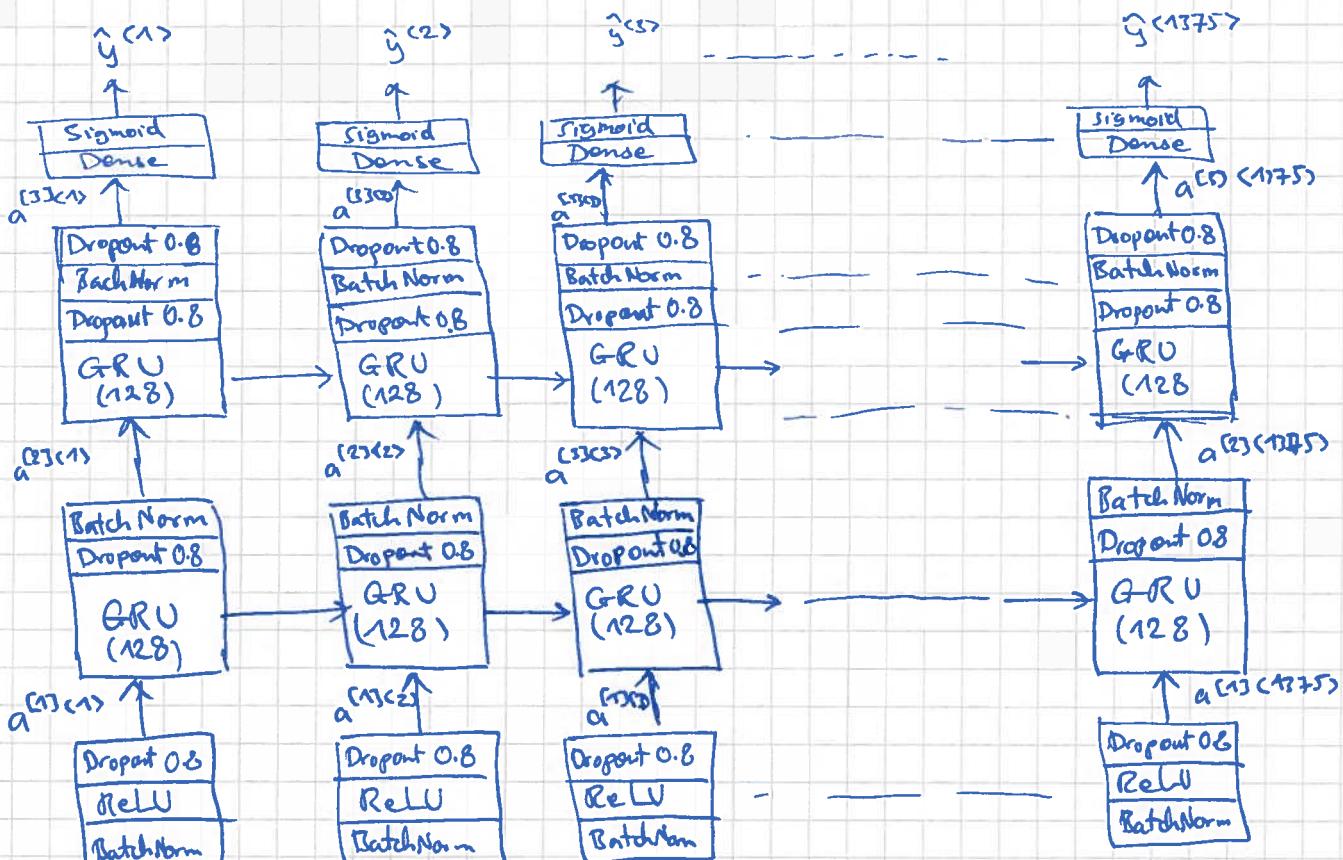
raw audio: 441000 samples  
10 sec.

### Synthesizing data:

- Pick a random 10 sec background noise audio clip.
- randomly insert 0-4 audio clips of "activate" into it
- randomly insert 0-2 audio clips of negative words into it

Model:

- 1D conv layers; GRU layers; dense layers.



- make input smaller ( $5511 \rightarrow 1375$ )  $\Rightarrow$  speed up the network
- detect first set of low-level features extract

Fantastic!

- Data Synthesis is an effective way to create a large training set for speech problems, specifically trigger word detection.
- Using a Spectrogram and optionally a 1D convolutional layer is a common pre-processing step prior to passing audio data to an RNN, GRU, or LSTM.
- An end-to-end deep learning approach can be used to build a very effective trigger word detection system.



