

DEEP LEARNING SPECIALIZATION
COURSE 5/5
SEQUENCE MODELS

D-5-1

(1/7)

→ EXAMPLES OF SEQUENCE DATA

* SPEECH RECOGNITION, MUSIC GENERATION, SENTIMENT CLASSIFICATION, DNA SEQUENCE ANALYSIS, MACHINE TRANSLATION, VIDEO ACTIVITY RECOGNITION, NAME ENTITY RECOGNITION

→ NOTATION

* MOTIVATION EXAMPLE: NAME-ENTITY RECOGNITION

$X: HARRY \mid POTTER \mid AND \mid HERMIONE \mid GRANGER \mid INVENTED \mid A \mid NEW \mid SPELL$

$x^{(1)} \quad x^{(2)} \quad \dots \quad x^{(t)} \quad \dots \quad x^{(9)}$

$T_x = 9 \rightarrow$

$y: 1 \mid 1 \mid 0 \mid 1 \mid 1 \mid 0 \mid 0 \mid 0 \mid 0$

$y^{(1)} \quad y^{(2)} \quad \dots \quad y^{(t)} \quad \dots \quad y^{(9)}$

$T_y = 9 \rightarrow$

$X^{(i)t}$ - t-th ELEMENT FROM i-th EXAMPLE

$T_x^{(i)}$ - SIZE (NO. OF ELEMENTS)

$y^{(i)t}$ - t-th ELEMENT FROM i-th OUTPUT

$T_y^{(i)}$ - SIZE (NO. OF ELEMENTS)



* VOCABULARY (DICTIONARY)

X: HARRY | POTTER | AND | HERMIONE | GRANGER | INVENTED | A | NEW | SPELL.

X⁽¹⁾ X⁽²⁾ X⁽³⁾ X⁽⁴⁾ X⁽⁵⁾ X⁽⁶⁾ X⁽⁷⁾ X⁽⁸⁾ X⁽⁹⁾

VOCABULARY

A	1	0	0	0	1	1
AARON	2	0	0	0	0	0
i	:	:	:	:	:	:
AND	367	0	0	1	367	0
:	:	:	:	:
HARRY	4075	1	4075	0	0	0
:	:	:	:	:	:	:
POTTER	6830	0	1	6830	0	0
:	:	:	:	:	:	:
ZULU	10000	0	0	0	0	0

↑
10000 WORD
DICTIONARY

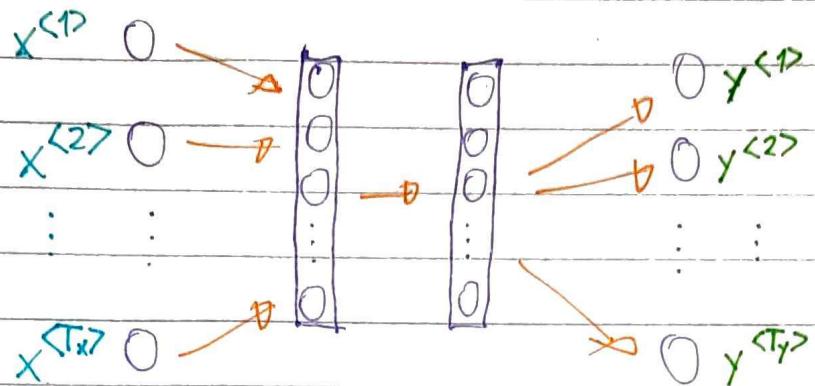
ONE-HOT ENCODING TO DETECT
WORDS

* WHAT IF A WORD IS NOT FOUND ON THE VOCABULARY?

- CREATE NEW TOKEN FOR UNKNOWN WORDS. IF A WORD ISN'T ON DICTIONARY, RETURN THIS TOKEN.



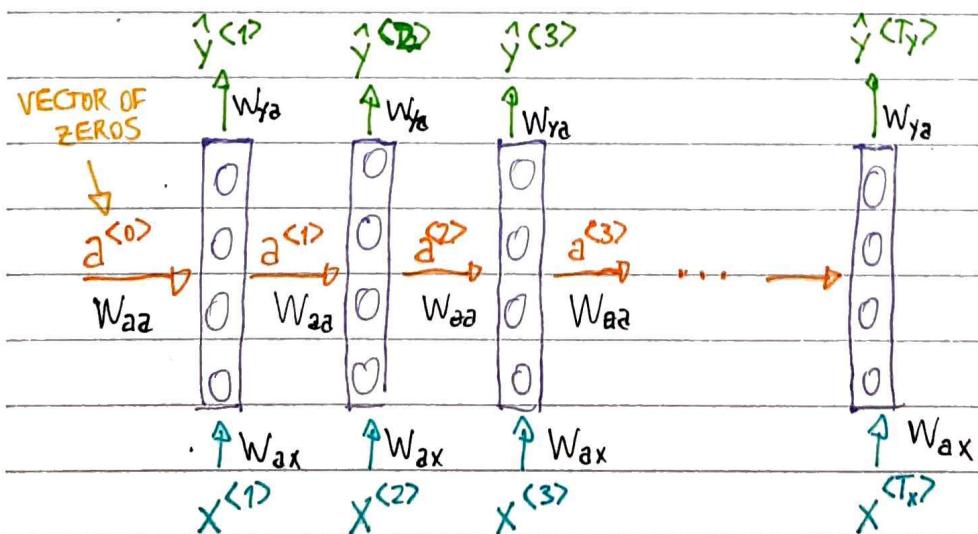
→ WHY NOT A STANDARD NETWORK?



* PROBLEMS:

- INPUTS AND OUTPUTS CAN HAVE DIFFERENT LENGTHS IN DIFFERENT EXAMPLES.
- DOESN'T SHARE FEATURES LEARNED ACROSS DIFFERENT POSITIONS OF TEXT.

→ RECURRENT NEURAL NETWORKS (RNN)



* DISADVANTAGE: RNN's ONLY USE EARLIER INFORMATION

* TO USE LATER INFORMATION: BIDIRECTIONAL RNN (BRNN)



→ FORWARD PROPAGATION ON RNN'S

$$a^{<0>} = \vec{0}$$

$$a^{<t>} = g_1(W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>} + b_a) \leftarrow g_1: \text{TANH/ReLU}$$

$$\hat{y}^{<t>} = g_2(W_{ya} \cdot a^{<t>} + b_y) \leftarrow g_2: \text{SOFTMAX/SIGMOID/ETC.}$$

* SIMPLIFYING NOTATION:

$$W_a = [W_{aa} : W_{ax}]$$

$$[a^{<t-1>} : x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ \vdots \\ x^{<t>} \end{bmatrix} \quad a^{<t>} = g_1(W_a[a^{<t-1>} : x^{<t>}] + b_a)$$

$$W_{ya} = W_y \rightarrow \hat{y}^{<t>} = g_2(W_y \cdot a^{<t>} + b_y)$$

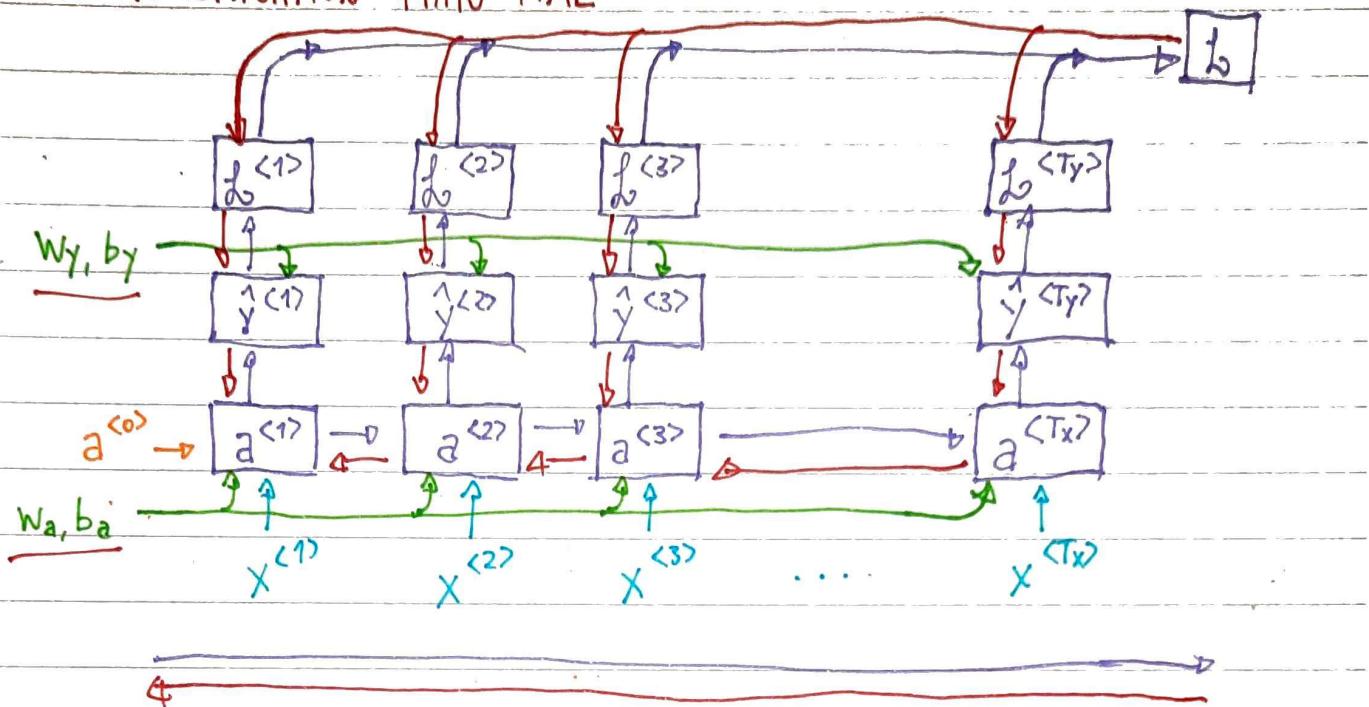
* FOR W_a :

- STACK MATRICES HORIZONTALLY. THEY HAVE THE SAME NUMBER OF ROWS, BUT NUMBER OF COLUMNS CAN BE DIFFERENT.

* FOR $[a^{<t-1>} : x^{<t>}]$:

- STACK MATRICES VERTICALLY.

→ BACKPROPAGATION THRU TIME



* LOSS FUNCTION: CROSS-ENTROPY

$$l^{(t)}(\hat{y}^{(t)}, y^{(t)}) = -y^{(t)} \log \hat{y}^{(t)} - (1-y^{(t)}) \log (1-\hat{y}^{(t)})$$

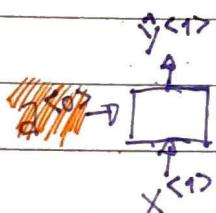
$$L(\hat{y}, y) = \sum_{t=1}^{T_x} l^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

- BACKPROPAGATE OVERALL LOSS TO PER-STEP LOSSES AND TO THE REST OF THE CHAIN

- MOST SIGNIFICANT RECURSIVE CALCULATION: ACTIVATIONS.
BACKPROP OCCURS "FROM RIGHT TO LEFT" (LATER INFORMATION TO EARLIER INFORMATION)

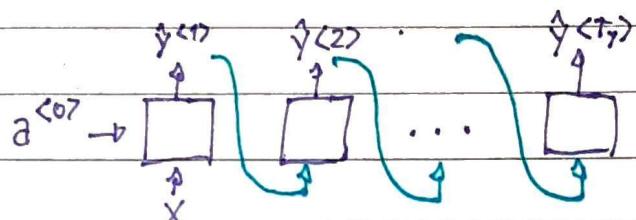
→ EXAMPLES OF RNN ARCHITECTURES

* ONE - TO - ONE



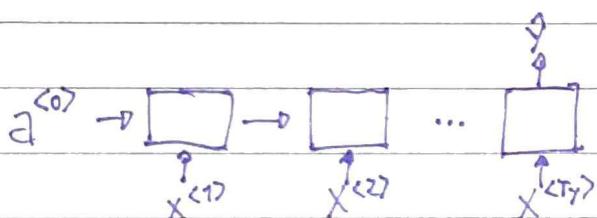
JUST A STANDARD N.N.
NOTHING SPECIAL.

* ONE - TO - MANY



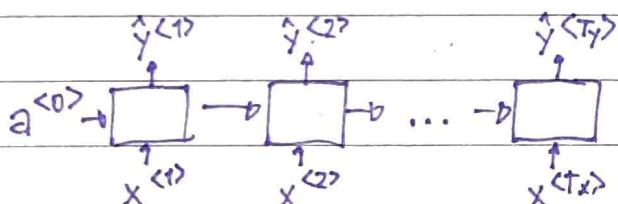
ONE INPUT, MANY OUTPUTS.

* MANY - TO - ONE



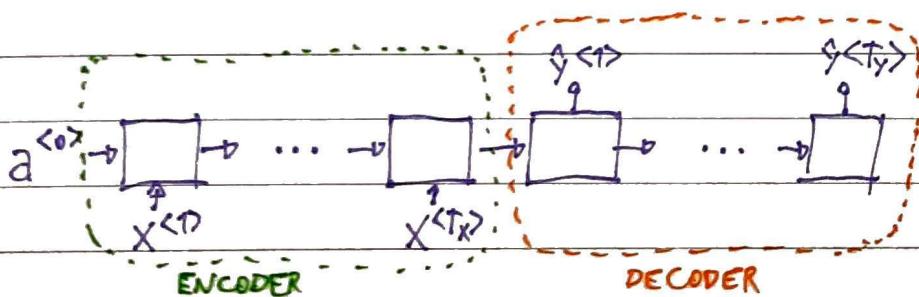
MANY INPUTS, ONE OUTPUT.

* MANY - TO - MANY ($T_x = T_y$)



INPUT AND OUTPUT HAVE
THE SAME LENGTH.

* MANY - TO - MANY ($T_x \neq T_y$) (ENCODER/DECODER)



ENCODER/DECODER STRUCTURE. INPUT AND OUTPUT HAVE DIFFERENT LENGTHS.



→ LANGUAGE MODELLING

* CALCULATE PROBABILITY OF A SENTENCE TO AVOID ERRORS/AMBIGUITY

* EXAMPLE: SPEECH RECOGNITION

- THE APPLE AND PAIR SALAD.
- THE APPLE AND PEAR SALAD.

$$P(\text{THE APPLE AND PAIR SALAD}) = 3.2 \times 10^{-13}$$

$$P(\text{THE APPLE AND PEAR SALAD}) = 5.7 \times 10^{-10}$$

$$P(\text{SENTENCE}) = ?$$

$$P(y^{<1>} , y^{<2>} , \dots , y^{<T>})$$

→ LANGUAGE MODELLING WITH AN RNN

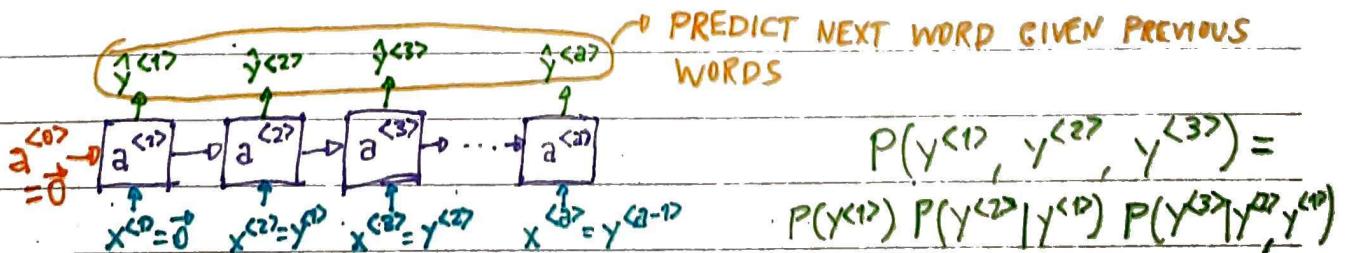
* TRAINING SET: LARGE CORPUS OF ENGLISH TEXT

* TOKENIZE WORDS BASED ON A VOCABULARY

THE EGYPTIAN UNK IS A BREAD OF CAT. <EOS>
 ↑ ↑ ↑ ↑ ↑ ↑ ↑ ↑ ↑ ↑
 y<1> y<2> y<3> y<4> y<5> y<6> y<7> y<8> y<9>

- <EOS>: END OF SENTENCE TOKEN
- <UNK>: UNKNOWN WORD TOKEN

* MODEL



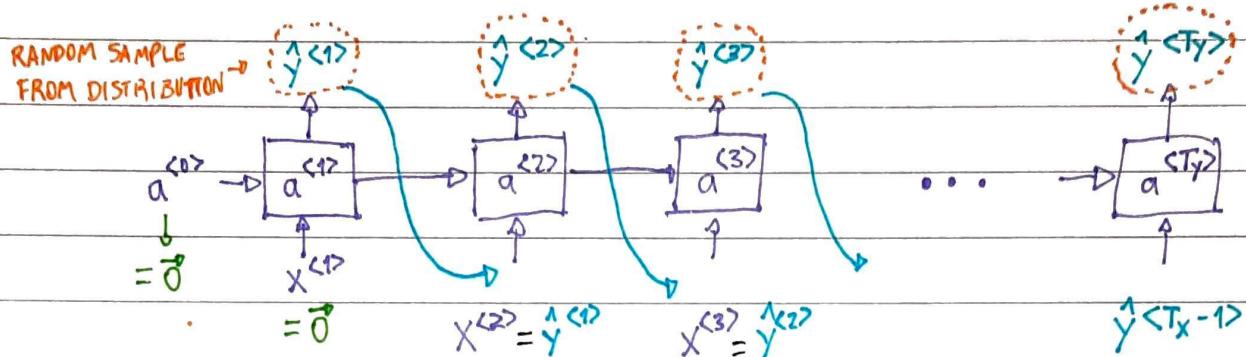
$$J_0(y^{<1>} , y^{<2>}) = - \sum_i y_i^{<t>} \log \hat{y}_i^{<t>} \quad \Big| \quad J_0 = \sum_t J_0^{<t>}(y^{<t>} , \hat{y}^{<t>})$$

→ SAMPLING A SEQUENCE FROM A TRAINED RNN

* TRAINED MODEL REPRESENTS A DISTRIBUTION OF PROBABILITY FOR SEQUENCES

* SAMPLING: TAKE RANDOM SAMPLES FROM THE PROBABILITY DISTRIBUTION

- NUMPY: np.random.choice



* END OF PREDICTION WHEN (EOS) IS RETURNED

* IF <EOS> TOKEN ISN'T ON VOCABULARY, LIMIT N° OF WORDS

* REJECTING <UNK> TOKENS OR NOT INCLUDING IT ON VOCABULARY

→ CHARACTER-LEVEL LANGUAGE MODEL

VOCABULARY = [a, b, c, ..., z, ←, →, ., ;, :, 0, 1, ..., 9, A, B, ..., Z]

* EVERY OUTPUT IS A CHARACTER

* LONGER SEQUENCES

* MORE COMPUTATIONALLY EXPENSIVE TO TRAIN

* USE W. SPECIALIZED APPLICATIONS

→ SEQUENCE GENERATION

* DEPENDS ON THE DATA YOU USE TO TRAIN THE MODEL

* EXAMPLES: TRAIN A LANGUAGE MODEL WITH NEWS TEXTS
OR SHAKESPEARE TEXTS



→ VANISHING GRADIENTS ON RNN'S

* EXAMPLE:

THE CAT, WHICH ALREADY ATE ... , WAS FULL.
THE CATS, WHICH ALREADY ATE ... , WERE FULL.

- RNN NEEDS PREVIOUS INFORMATION TO CHOOSE BETWEEN SINGULAR OR PLURAL NOUN

* RNN'S TEND NOT TO BE VERY GOOD AT CAPTURING LONG-RANGE DEPENDENCIES: VANISHING GRADIENTS ARE HARDER TO SOLVE

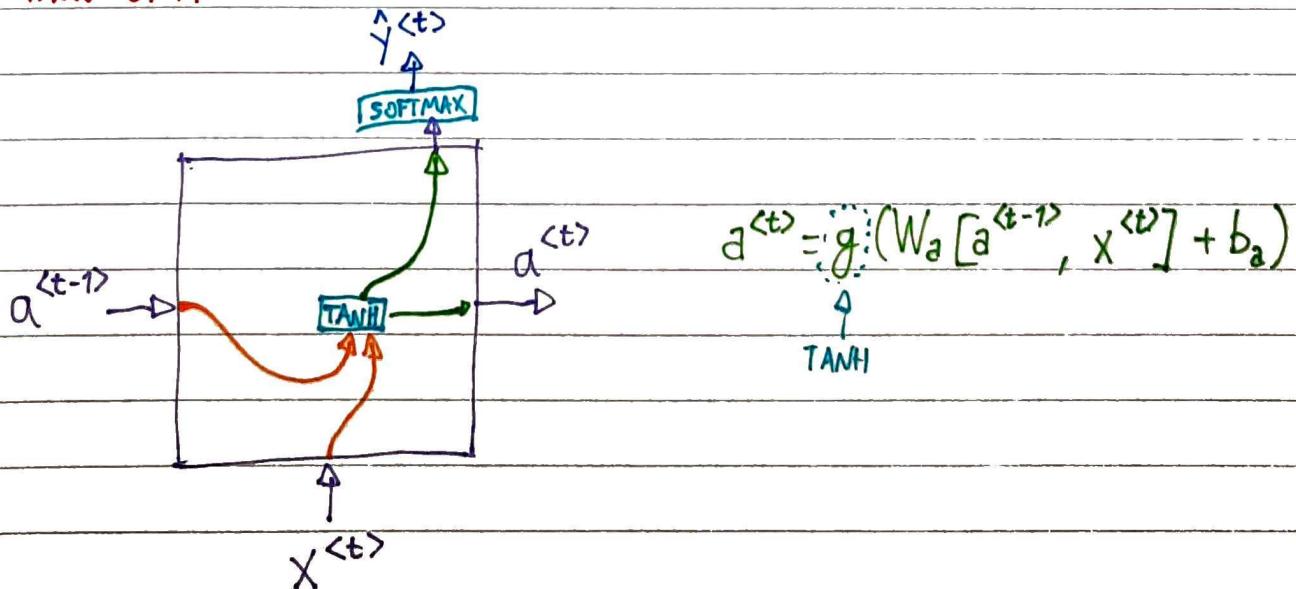
* EXPLODING GRADIENTS ARE EASIER TO DETECT AND SOLVE

- NaN: NUMERICAL OVERFLOW
- SOLUTION: GRADIENT CLIPPING

* PROPOSED SOLUTIONS

- GATED RECURRENT UNIT (GRU)
- LONG SHORT TERM MEMORY (LSTM)

→ RNN UNIT



→ GATED RECURRENT UNIT (GRU) SIMPLIFIED

* MEMORY CELL $C \Rightarrow C^{(t)} = a^{(t)}$

* C CAN BE A VECTOR

* CANDIDATE TO REPLACE $C^{(t)}$: $\tilde{C}^{(t)} = \tanh(W_c[C^{(t-1)}, x^{(t)}] + b_c)$

* GAMMA VALUE Γ_u BETWEEN 0 AND 1
 Γ_u
↑ UPDATE

$$\Gamma_u = \sigma(W_u[C^{(t-1)}, x^{(t)}] + b_u)$$

↑ SIGMOID.

* $C^{(t)}$, $\tilde{C}^{(t)}$ AND Γ_u HAVE THE SAME DIMENSION

* GAMMA IS A "GATE". IT DECIDES WHETHER OR NOT
 UPDATE THE MEMORY CELL

$$C^{(t)} = \Gamma_u * \tilde{C}^{(t)} + (1 - \Gamma_u) * C^{(t-1)}$$

ELEMENT-WISE ELEMENT-WISE

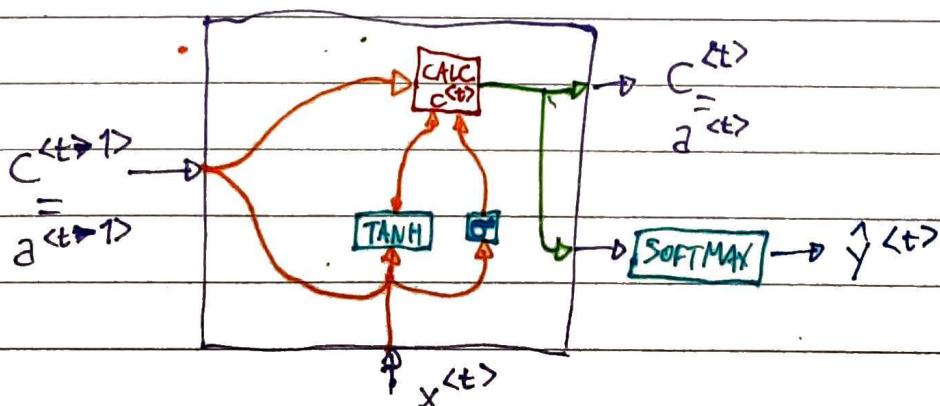
* EXAMPLE: C DETERMINES SINGULAR OR PLURAL NOUN

$$\Gamma_u = 1$$

$$C^{(t)} = 1 \quad \Gamma_u = 0 \quad \Gamma_u = 0 \quad \Gamma_v = 0, \dots \quad \Gamma_y = 1$$

THE CAT, WHICH ALREADY ATE, WAS: FULL

* SIGMOID KEEPS Γ_u CLOSE TO 0 OR CLOSE TO 1





→ FULL GATED RECURRENT UNIT (GRU)

* ADD RELEVANCE GATE Γ_r TO CALCULATE RELEVANCE OF C^{t-1}
IN THE CALCULUS OF \tilde{C}^t

$$\tilde{C}^t = \tanh (W_c [\Gamma_r * C^{t-1}, X^t] + b_c)$$

$$\Gamma_r = \sigma (W_r [C^{t-1}, X^t] + b_r)$$

$$\Gamma_u = \sigma (W_u [C^{t-1}, X^t] + b_u)$$

$$C^t = \Gamma_u * \tilde{C}^t + (1 - \Gamma_u) * C^{t-1}$$

* WHY USING RELEVANCE GATE TOO?

- MOST COMMONLY USED
- RELIABLE AND ROBUST

* MAYBE IN ACADEMIC LITERATURE, MEMORY CELLS, CANDIDATES
AND GATES MAY HAVE A DIFFERENT NOTATION

→ LONG-SHORT-TERM MEMORY (LSTM)

* SIMILAR TO THE IDEA OF GRU'S, BUT WITH THREE GATES
THAT WORK SLIGHTLY DIFFERENT

- Γ_u (UPDATE), Γ_f (FORGET), Γ_o (OUTPUT)

$$\tilde{c}^{(t)} = \tanh(W_c [a^{(t-1)}, x^{(t)}] + b_c)$$

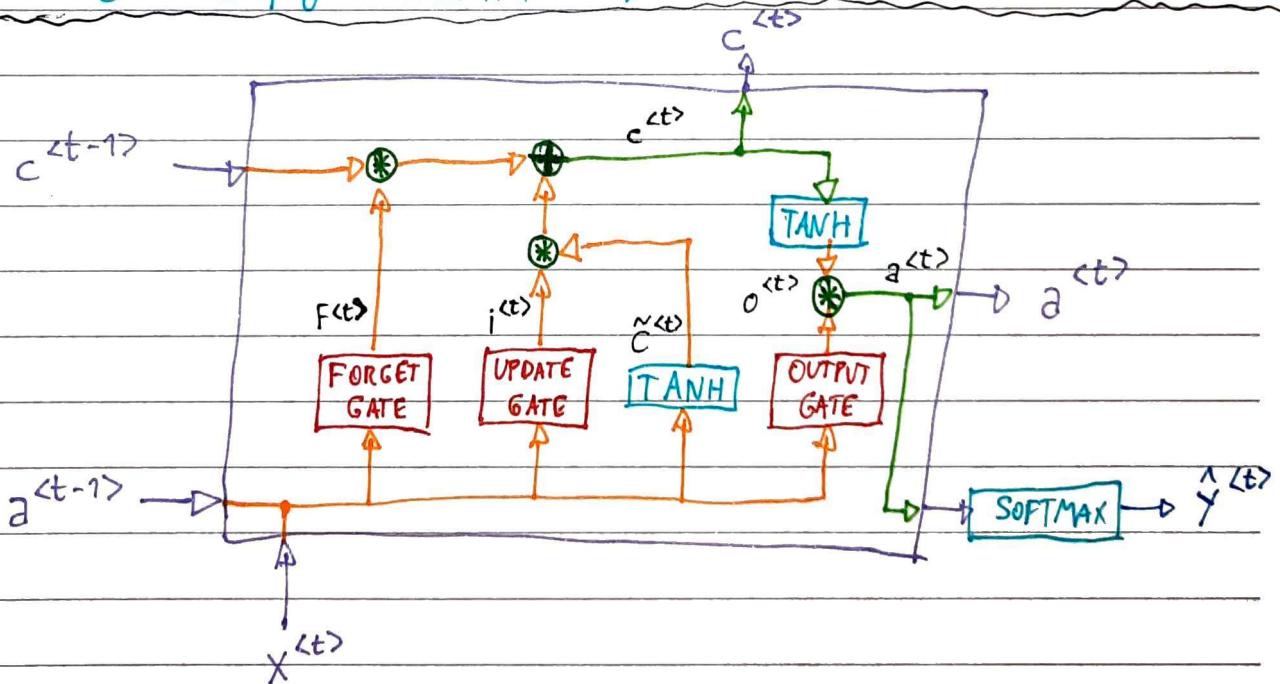
$$\Gamma_u = \sigma(W_u [a^{(t-1)}, x^{(t)}] + b_u)$$

$$\Gamma_f = \sigma(W_f [a^{(t-1)}, x^{(t)}] + b_f)$$

$$\Gamma_o = \sigma(W_o [a^{(t-1)}, x^{(t)}] + b_o)$$

$$c^{(t)} = \Gamma_u * \tilde{c}^{(t)} + \Gamma_f * c^{(t-1)}$$

$$a^{(t)} = \Gamma_o * \tanh(c^{(t)})$$



* VARIATION: "PEEPHOLE CONNECTION"

- Γ_o BEING DEPENDENT ON $c^{(t-1)}$ TOO
- CAN BE APPLIED TO Γ_u AND Γ_f TOO

→ GETTING INFORMATION FROM THE FUTURE (BIDIRECTIONAL RNN)

* IMAGINE TWO PHRASES WITH DIFFERENT CONTEXTS:

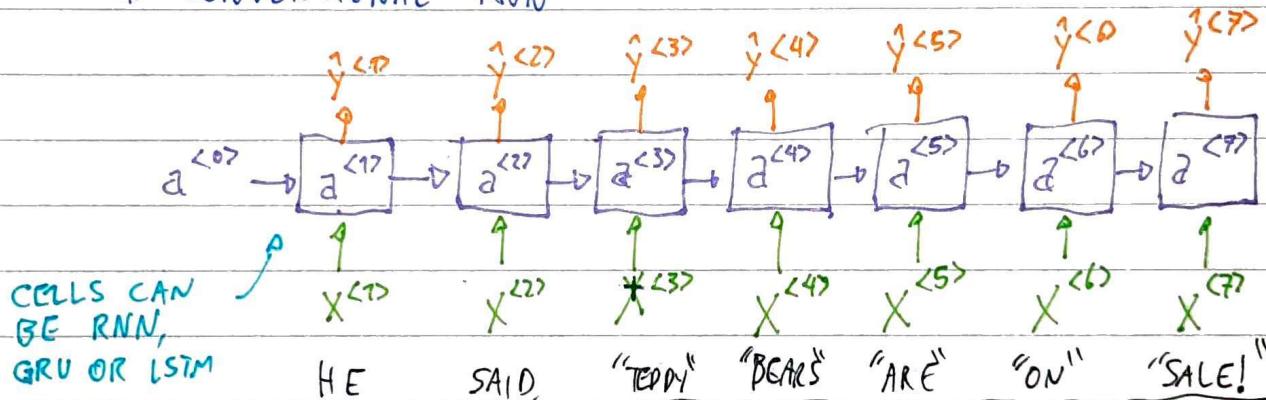
- HE SAID, "TEDDY BEARS ARE ON SALE"!
- HE SAID, "TEDDY ROOSEVELT WAS A GREAT PRESIDENT!"

↗ HOW TO PREDICT THE NEXT WORD

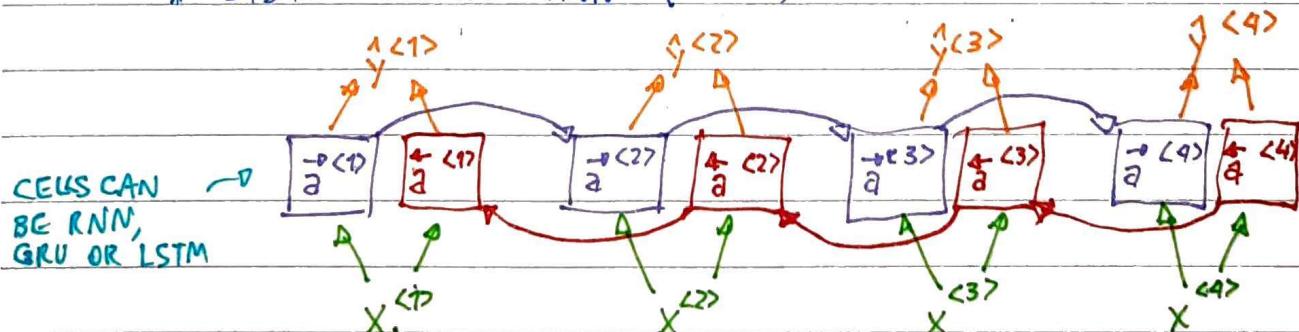
↖ WITH BETTER CORRELATION WITH CONTEXT?

GATHERING INFORMATION
FROM FUTURE!

* CONVENTIONAL RNN



* BIDIRECTIONAL RNN (BRNN)



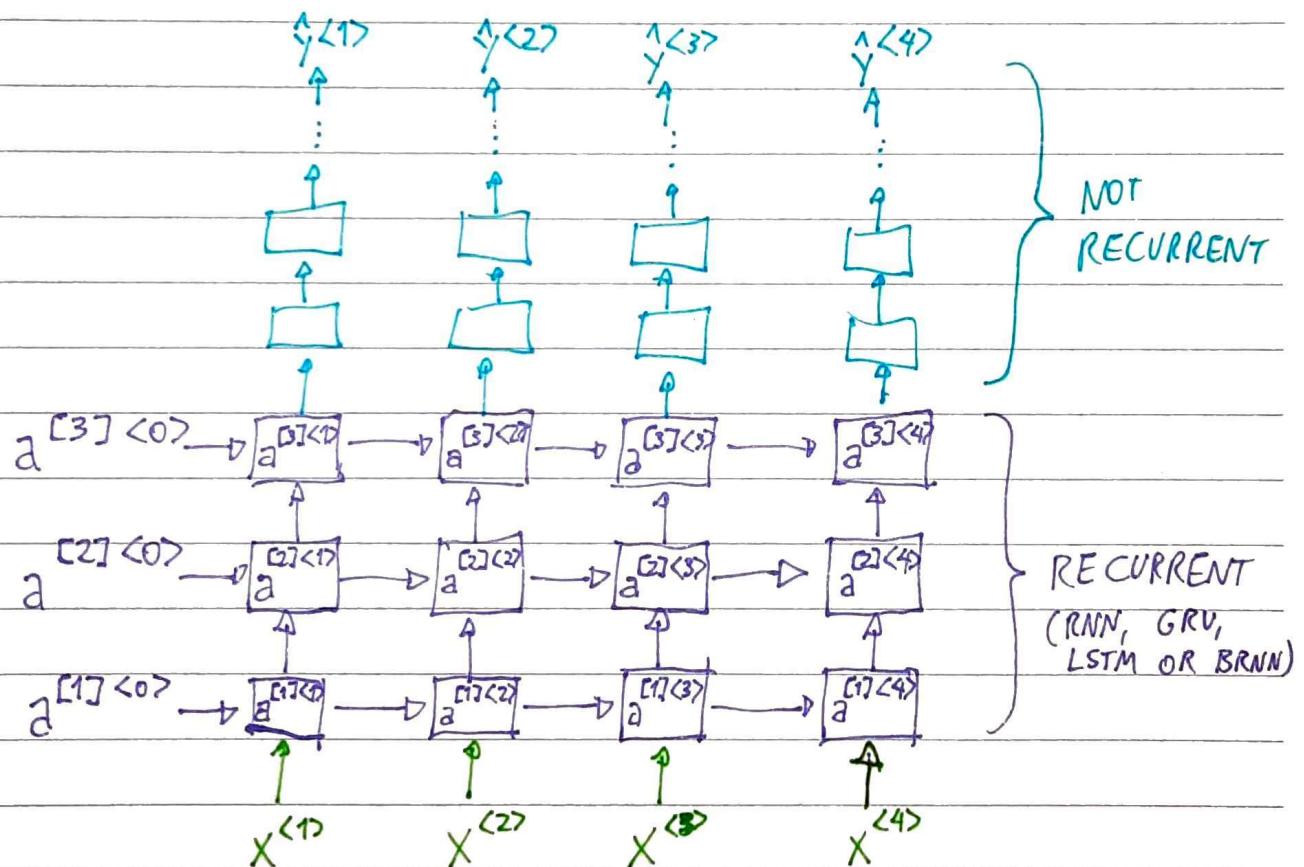
• ACYCLIC GRAPH

$$\hat{y}^{(t)} = g(W_y [\vec{a}^{(t)}, \overleftarrow{a}^{(t)}] + b_y)$$

→ DEEP RECURRENT NEURAL NETWORKS (DEEP RNN'S)

* STACKING RNN, GRU OR LSTM CELLS IN SEQUENCE,
FOLLOWED BY OTHER KINDS OF LAYERS (NOT NECESSARILY
RECURRENT)

* RECURRENT LAYERS ARE COMPUTATIONALLY EXPENSIVE TO
TRAIN, SO A FEW RECURRENT LAYERS ARE TYPICALLY USED
ON A DEEP NETWORK



* NOTATION: $\underset{\text{ACTIVATION LAYER}}{a} \underset{\text{TIMESTAMP}}{^q} \underset{\text{LAYER}}{^{[l]}} \underset{\text{LAYER}}{^{[t]}}$

$$a^{[l] < t >} = g(W_a^{[l]} [a^{[l] < t-1 >}, a^{[l-1] < t >}] + b_a^{[l]})$$

→ WORD REPRESENTATION

* ONE-HOT VOCABULARY REPRESENTATION

- WEAKNESS: DOESN'T ALLOW SEMANTIC RELATIONSHIP BETWEEN WORDS

I WANT A GLASS OF ORANGE JUICE

I WANT A GLASS OF APPLE ?????

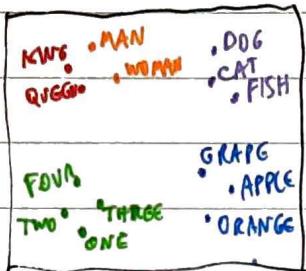


APPLE AND ORANGE HAVE A CORRELATION (BOTH ARE FRUITS),
BUT ONE-HOT ENCODING DOESN'T ALLOW THE MODEL TO
LEARN THIS SEMANTIC RELATIONSHIP. IN OTHER WORDS,
THE MODEL WON'T GENERALIZE TO PREDICT THAT THE NEXT
WORD MIGHT BE "JUICE".

→ FEATURIZED WORD REPRESENTATION: WORD EMBEDDING

	(5391) MAN	(9853) WOMAN	(4914) KING	(7157) QUEEN	(456) APPLE	(6257) 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
FOOD	0,03	0,01	0,02	0,01	0,95	0,97
:	:	:	:	:	:	:
VERB	:	:	:	:	:	:
NOUN	:	:	:	:	:	:
:	:	:	:	:	:	:
	e ₅₃₉₁	e ₉₈₅₃				

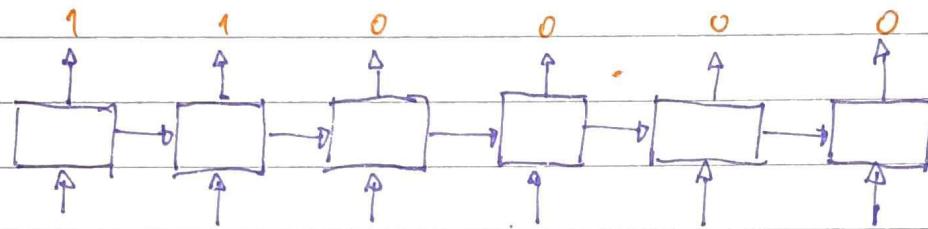
→ VISUALIZING WORD EMBEDDINGS



300 DIMENSIONS → 2-D EMBEDDING

ALGORITHM t-SNE

→ NAMED ENTITY RECOGNITION EXAMPLE



Sally Johnson is an orange Farmer. ↗ AN ORANGE FARMER IS A PERSON

Robert Lin is an apple farmer. ↗ APPLE/ORANGE ARE SIMILAR; ROBERT LIN IS A PERSON

Robert Lin is a durian cultivator. ↗ WITH EMBEDDINGS, IT'S POSSIBLE TO IDENTIFY EVEN UNCOMMON WORDS.

- TRAINING WITH 1B ~ 100B WORDS

- NAME RECOGNITION DATASET: 700K WORDS (TRANSFER LEARNING)

→ TRANSFER LEARNING AND WORD EMBEDDING

1. LEARN WORD EMBEDDINGS FROM LARGE TEXT CORPUS. (1~100B WORDS)

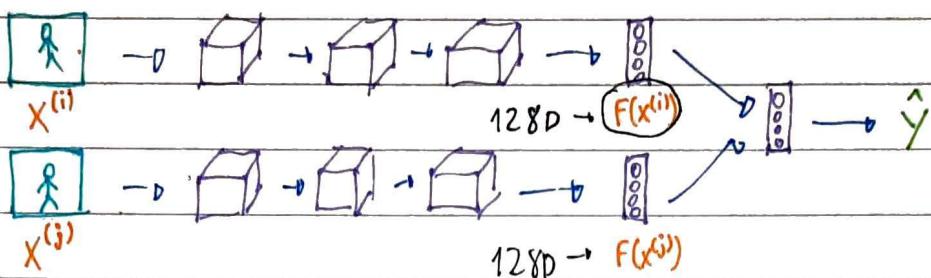
(OR DOWNLOAD PRE-TRAINED EMBEDDING ONLINE)

2. TRANSFER EMBEDDING TO NEW TASK WITH SMALLER TRAINING SET
(SAY, 100K WORDS)

3. CONTINUE TO FINETUNE THE WORD EMBEDDINGS WITH NEW DATA.

OBS: NOT TOO USEFUL FOR LANGUAGE MODELING AND MACHINE TRANSLATION

→ RELATION TO FACE ENCODING (OR EMBEDDING)



* $F(X^{(i)})$ AND $F(X^{(j)})$ ARE ENCODINGS (OR EMBEDDINGS)

* DIFFERENCE: FACE RECOGNITION COMPUTES ENCODINGS FOR ANY IMAGE INPUT; NLP HAS A STRICT VOCABULARY

→ ANALOGIES

	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	0,09	0,01	0,02	0,01	0.95	0.97
	e_{MAN}	e_{WOMAN}	e_{KING}	e_{QUEEN}		

MAN → WOMAN

AS

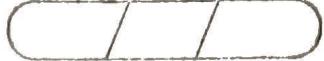
KING → ??? (QUEEN)

$$e_{\text{MAN}} - e_{\text{WOMAN}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

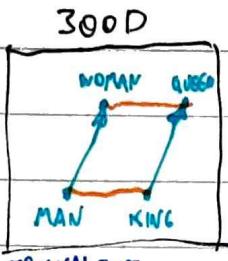
BOTH ARE → CLOSE

(ALMOST EQUAL)

$$e_{\text{KING}} - e_{\text{QUEEN}} \approx \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$



→ ANALOGIES USING WORD VECTORS



ATTENTION: t-SNE "BREAKS" LINEARITY AND WE CAN'T RELY ON THESE "PARALLELGRAM" DISTANCES ON A t-SNE REPRESENTATION.

$$e_{\text{MAN}} - e_{\text{WOMAN}} \approx e_{\text{KING}} - e_{\text{W}}$$

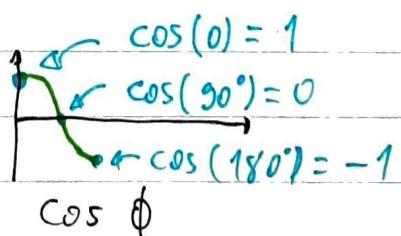
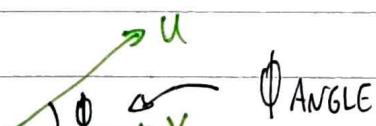
FIND WORD w: $\operatorname{argmax}_w [\text{SIMILARITY}(e_w, e_{\text{KING}} - e_{\text{MAN}} + e_{\text{WOMAN}})]$

- APPROX. 30 ~ 75% OF ACCURACY

→ COSINE SIMILARITY

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

$$\text{COSINE ANGLE}$$



→ EXAMPLES OF ANALOGIES

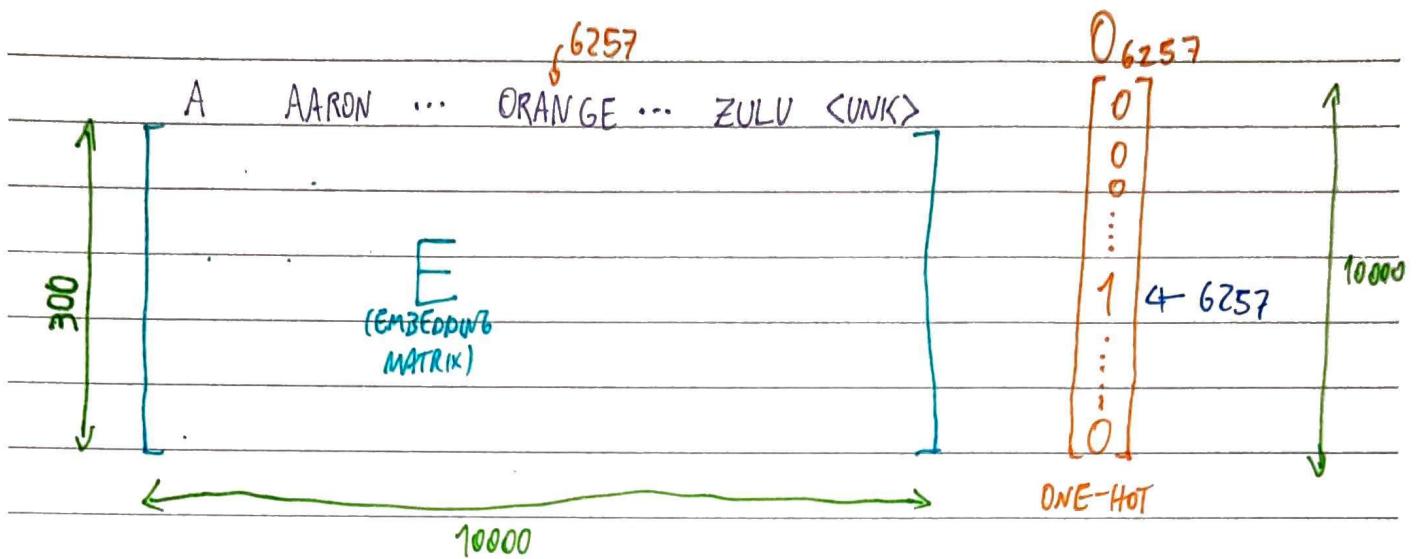
* MAN: WOMAN AS BOY: GIRL

* OTTAWA: CANADA AS NAIROBI: KENYA

* BIG: BIGGER AS TALL: TALLER

* YEN: JAPAN AS RUBLE: RUSSIA

→ EMBEDDING MATRIX



$$E_{(300, 10000)} \cdot O_{6257 (10000, 1)} = e_{6257} (300, 1)$$

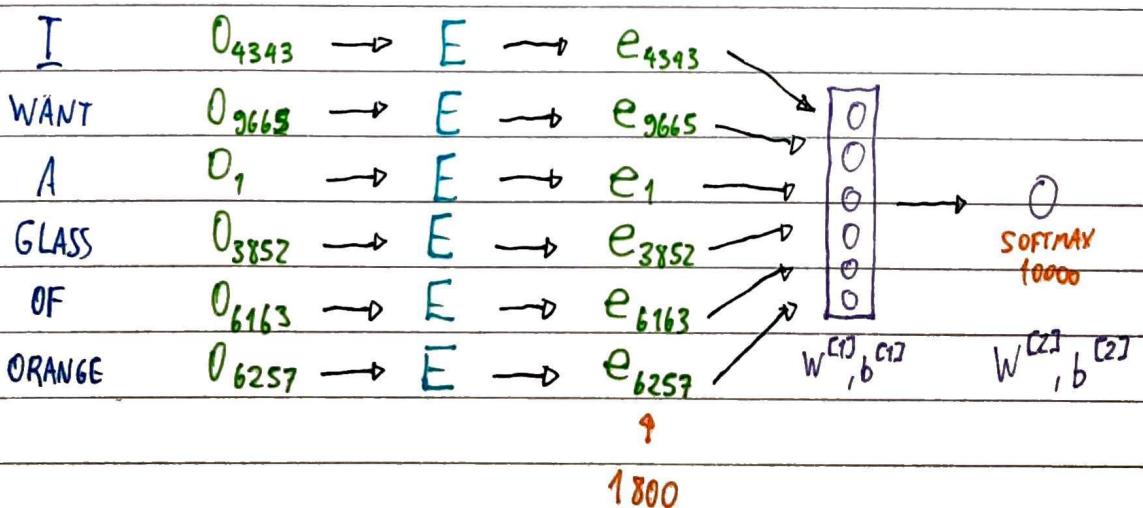
EMBEDDING
MATRIX

ONE-HOT ENCODING
FOR "ORANGE"

EMBEDDING VECTOR
FOR "ORANGE"

→ NEURAL LANGUAGE MODEL

I WANT A GLASS OF ORANGE JUICE.
4343 9665 1 3852 6163 6257





→ FIXED HISTORY ON NEURAL LANGUAGE MODELS

- * INPUT SIZES ALWAYS FIXED
- * USE ALWAYS THE LAST N WORDS TO PREDICT THE NEXT

→ OTHER CONTEXT/TARGET PAIRS

I WANT A GLASS OF ORANGE JUICE TO GO ALONG WITH MY CEREAL.

CONTEXT: → LAST 4 WORDS.

→ 4 WORDS ON LEFT & RIGHT

→ LAST 1 WORD → ORANGE

→ NEARBY 1 WORD → GLASS ...

A GLASS OF ORANGE TO GO ALONG WITH

→ SKIP-GRAMS (WORD2VEC)

* RANDOMLY PICK A WORD TO BE THE CONTEXT WORD

* RANDOMLY PICK ANOTHER WORD WITHIN SOME WINDOW ($\pm N$ WORDS) TO BE THE TARGET WORD

I WANT A GLASS OF JUICE TO GO ALONG WITH MY CEREAL.



CONTEXT	ORANGE	ORANGE	ORANGE
TARGET	JUICE	GLASS	MY

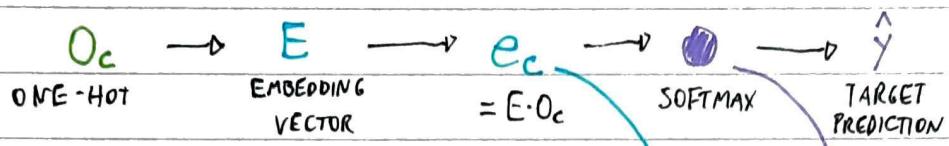
* LEARN GOOD WORD EMBEDDINGS



→ MODEL (VOCAB SIZE = 10000)

* EXAMPLE: LEARN MAPPING FROM CONTEXT WORD TO TARGET WORD

CONTEXT C ("ORANGE")₆₂₅₇ → TARGET t ("JUICE")₄₈₃₄



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

$\left(\theta_t: \text{PARAMETER ASSOCIATED WITH TARGET } t \right)$

$$\text{LOSS FOR SOFTMAX: } L(\hat{y}, y) = - \sum_{i=1}^{10000} y_i \log \hat{y}_i$$

\hat{y} : VECTOR WITH SOFTMAX PROBABILITIES

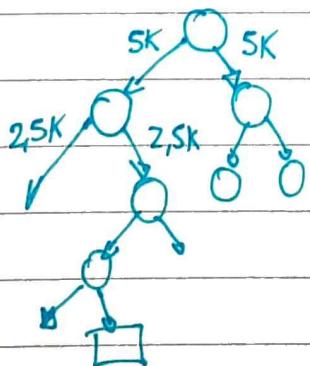
y : ONE-HOT VECTOR FOR TARGET WORD

→ PROBLEMS WITH SOFTMAX CLASSIFICATION

* COMPUTATIONALLY EXPENSIVE

* DENOMINATOR OF PROBABILITY WITH HUGE AMOUNT OF TERMS

* HIERARCHICAL SOFTMAX: (LOG COMPLEXITY RATHER THAN LINEAR)



IN PRACTICE, A BALANCED / PERFECTLY SYMMETRICAL TREE ISN'T USED; COMMON WORDS ARE POSITIONED ON TOP, AND LESS COMMON WORDS ARE PLACED DEEPER IN THE TREE.

* HOW TO SAMPLE CONTEXT C?

THE, OF, A, AND, TO... ← COMMON

ORANGE, APPLE, DURIAN.. ← LESS COMMON

DON'T TAKE UNIFORM SAMPLES;
TECHNIQUES TO AVOID OVERCALCULUS
OF COMMON WORDS

→ DEFINING A NEW LEARNING PROBLEM

* CHOOSE A PAIR OF WORDS THAT HAVE A CONTEXT - TARGET RELATIONSHIP

* RANDOMLY SAMPLE OTHER TARGET EXAMPLES WITHOUT RELATIONSHIP WITH THE CONTEXT WORD

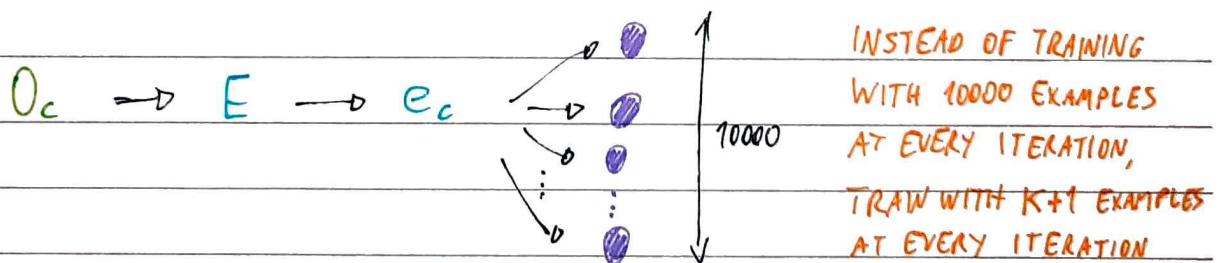
CONTEXT	WORD	POSITIVE EXAMPLE	TARGET ? } Y
ORANGE	JUICE	1	
ORANGE	KING	0	
ORANGE	BOOK	0	
ORANGE	THE	0	
ORANGE	OF	0	

SMALL DATASETS: K=5~20 → K
LARGE DATASETS: K=2~5

NEGATIVE SAMPLES

* LOGISTIC REGRESSION MODEL (LOWER COMPUTATIONAL COST)

$$P(Y=1 | c, t) = \sigma(\theta_t^T e_c)$$



→ SELECTING NEGATIVE EXAMPLES

* EMPIRICALLY TAKEN HEURISTIC VALUE

$$P(w_i) = \frac{f(w_i)^{\frac{3}{4}}}{\sum_{j=1}^{10000} f(w_j)^{\frac{3}{4}}}$$

F(w_i): OBSERVED FREQUENCY OF THE WORD w_i IN THE TRAINING SET TEXT CORPUS

→ GLOVE (GLOBAL VECTORS FOR WORD REPRESENTATION)

$X_{ij} = \# \text{ TIMES } j \text{ APPEARS IN CONTEXT OF } i$

* MODEL:

$$\text{MINIMIZE} \quad \sum_{i=1}^{10000} \sum_{j=1}^{10000} f(x_{ij}) (\theta_i^T e_j + b_i + b_j - \log x_{ij})^2$$

WEIGHTING TERM "θ_i^Te_j"

* IF $\log X_{ij} = 0 \rightarrow f(X_{ij}) = 0$ ("0 log 0" = 0)

* WEIGHTING TERM HELPS TO BALANCE THE AMOUNT OF COMPUTATION

FOR COMMON AND UNCOMMON WORDS

* IN GLOVE, θ_i AND e_j ARE SYMMETRIC $\rightarrow e_w^{(FINAL)} = \frac{(e_w + \theta_w)}{2}$

→ NOTE ON THE FEATURIZATION VIEW OF WORD EMBEDDINGS

* IN PREVIOUS EXAMPLES, EMBEDDING FEATURES WERE REPRESENTED AS "GENDER", "ROYAL", "AGE", "FOOD", ETC.

* IN REAL APPLICATIONS, EMBEDDINGS MAY NOT BE EASILY INTERPRETABLE

* EMBEDDINGS MAY BE COMBINATIONS OF FEATURES

* IT'S VERY DIFFICULT TO LOOK AT INDIVIDUAL COMPONENTS AND ASSIGN HUMAN INTERPRETATION

* PARALLELOGRAM MAP STILL WORKS FOR REAL WORD EMBEDDINGS

→ SENTIMENT CLASSIFICATION PROBLEM

X

Y

THE DESSERT IS EXCELLENT.

SERVICE WAS QUITE SLOW.

GOOD FOR A QUICK MEAL, BUT NOTHING SPECIAL.

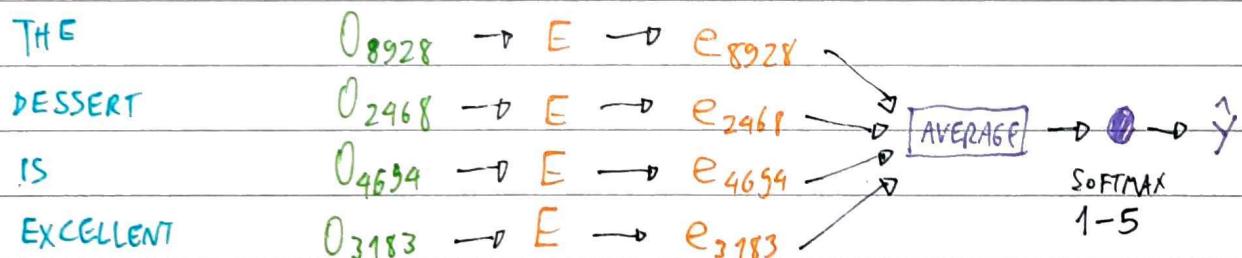
COMPLETELY LACKING IN GOOD TASTE, SERVICE AND AMBIENCE.



→ SIMPLE SENTIMENT CLASSIFICATION MODEL

THE DESSERT IS EXCELLENT

8928 2468 4694 3183



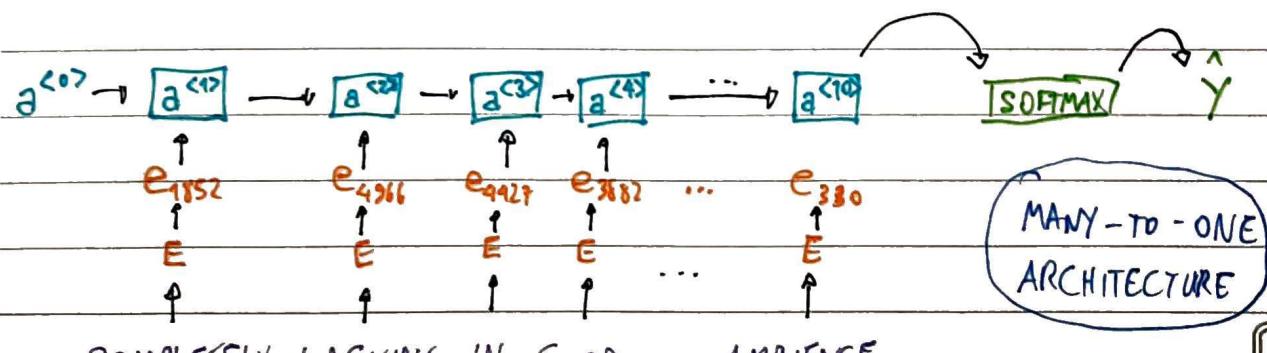
* AVERAGE MAY NOT WORK WELL IN SOME CASES.

COMPLETELY LACKING IN GOOD TASTE, GOOD SERVICE, AND GOOD AMBIENCE.

NEGATIVE REVIEW ↗

↑
THE WORD "GOOD" APPEARS THREE TIMES

→ RNN FOR SENTIMENT CLASSIFICATION



→ THE PROBLEM OF BIAS IN WORD EMBEDDINGS

* MAN: WOMAN AS KING: QUEEN ✓

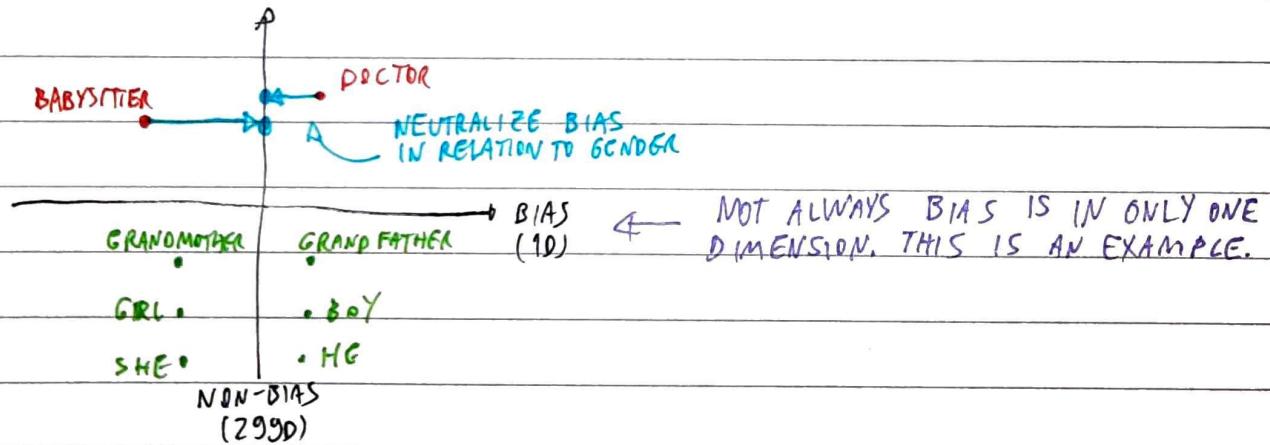
* MAN: COMPUTER_PROGRAMMER AS WOMAN: HOMEMAKER X

* 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.

→ ADDRESSING BIAS IN WORD EMBEDDINGS.

* EXAMPLE: BIAS IN GENDER DIRECTION



* STEPS FOR DEBIASING:

1. IDENTIFY BIAS DIRECTION. $\left\{ \begin{array}{l} e_{he} - e_{she} \\ e_{male} - e_{female} \end{array} \right. \xrightarrow{\text{AVERAGE}}$

2. NEUTRALIZE: FOR EVERY WORD THAT IS NOT DEFINITIONAL, PROJECT TO GET RID OF BIAS.

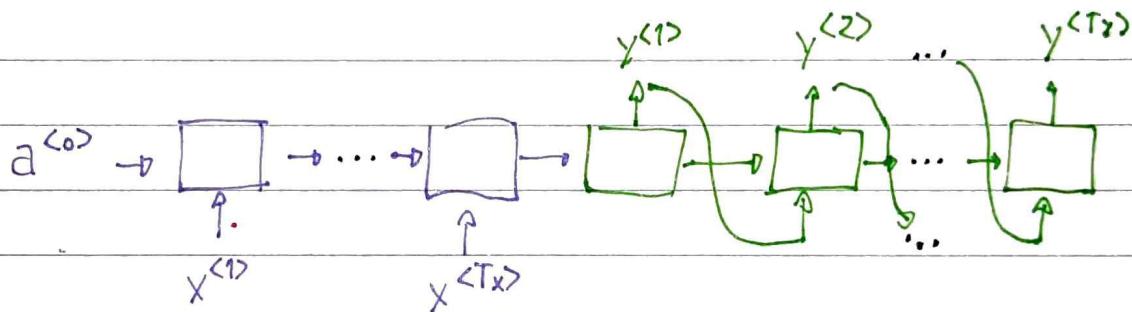
3. EQUALIZE PAIRS. $\left\{ \begin{array}{l} \text{GRANDMOTHER} - \text{GRANDFATHER} \\ \text{GIRL} - \text{BOY} \\ \vdots \end{array} \right.$

→ BASIC SEQUENCE MODELS

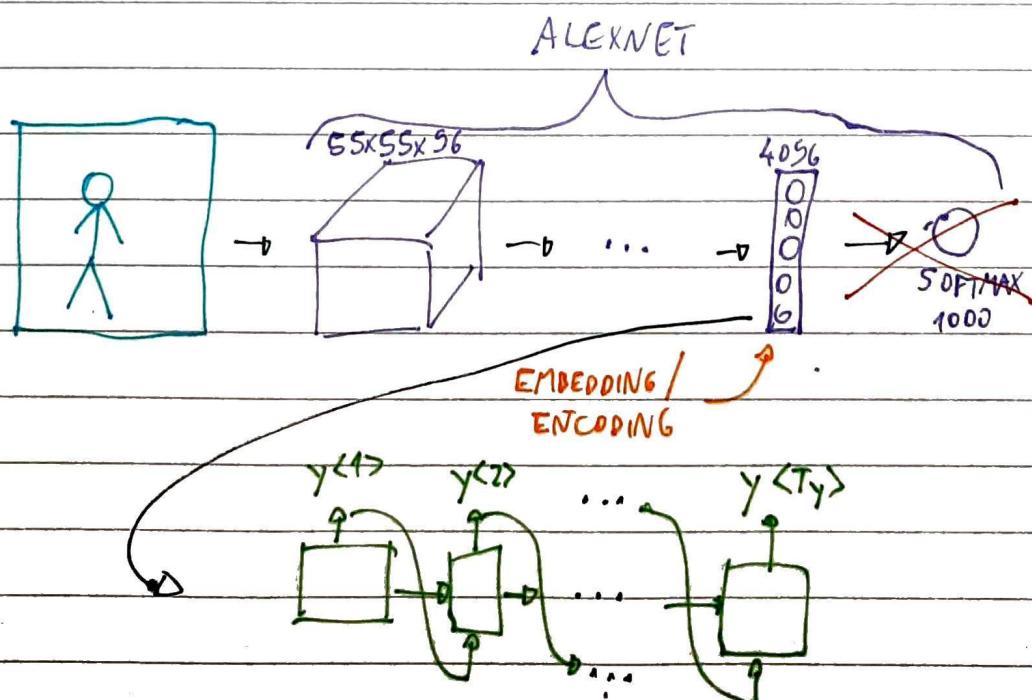
* SEQUENCE - TO - SEQUENCE MODEL (EXAMPLE: MACHINE TRANSLATION)



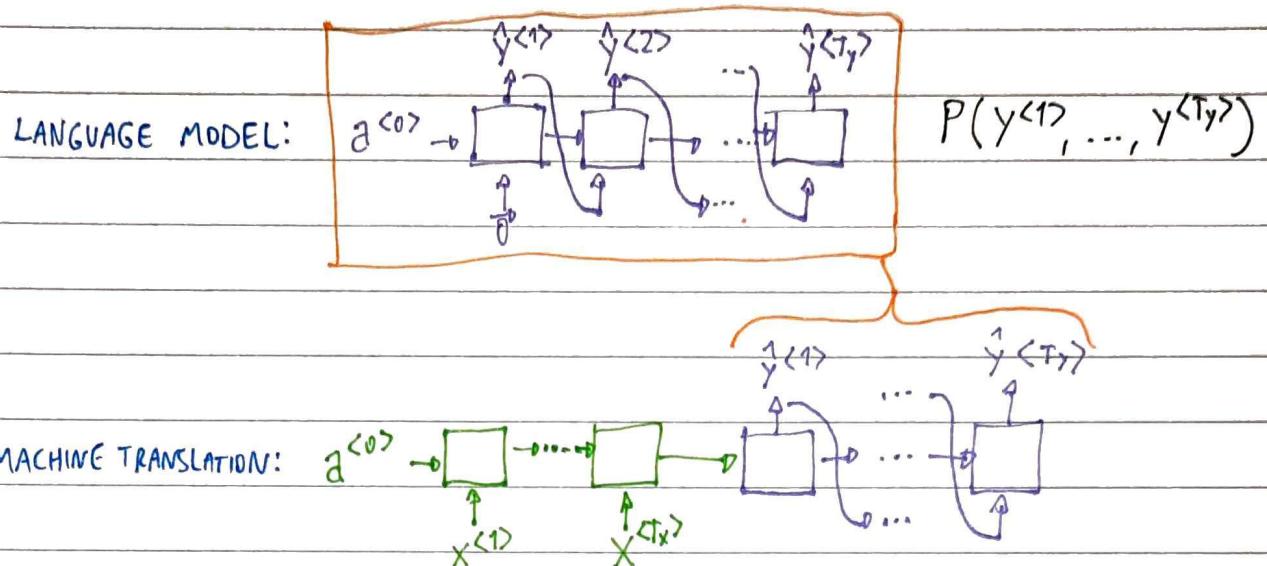
ENCODER - DECODER ARCHITECTURE



* IMAGE CAPTIONING (EXAMPLE WITH ALEXNET)



→ MACHINE TRANSLATION AS BUILDING A CONDITIONAL LANGUAGE MODEL



"CONDITIONAL LANGUAGE MODEL": $P(y^{<1>}, \dots, y^{<Ty>} | x^{<1>}, \dots, x^{<Tx>})$

→ FINDING THE MOST LIKELY TRANSLATION

JANE VISITE L'AFRIQUE EN SEPTEMBRE.

- JANE IS VISITING AFRICA IN SEPTEMBER. - **BEST TRANSLATION**
- JANE IS GOING TO BE VISITING AFRICA IN SEPTEMBER. - **GOOD TRANSLATION, BUT MORE VERBOSER**
- IN SEPTEMBER, JANE WILL VISIT AFRICA. - **ACCEPTABLE TRANSLATION**
- HER AFRICAN FRIEND WELCOMED JANE IN SEPTEMBER. - **WRONG TRANSLATION**

$$\underset{y^{<1>}, \dots, y^{<Ty>}}{\operatorname{argmax}} P(y^{<1>}, \dots, y^{<Ty>} | x)$$

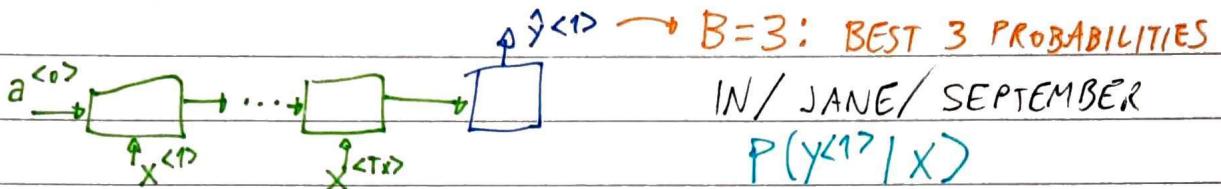
→ WHY NOT GREEDY SEARCH?

- * BEST INDIVIDUAL WORD PROBABILITIES DOESN'T ALWAYS LEAD TO THE BEST TRANSLATIONS
- * THEY COULD LEAD TO ACCEPTABLE TRANSLATIONS, BUT WITH MORE VERBOSITY
- * TOTAL POSSIBLE COMBINATIONS OF WORDS EXPONENTIALLY LARGER;
COMPUTATIONALLY EXPENSIVE

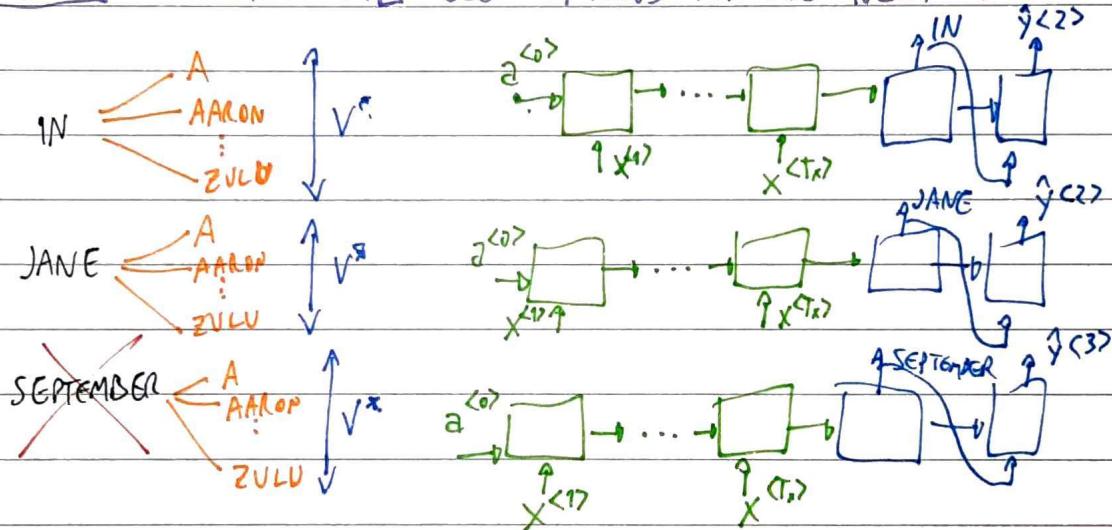
→ BEAM SEARCH ALGORITHM

* EXAMPLE: TRANSLATE "JANE VISITE L'AFRIQUE EN SEPTEMBRE"

STEP 1: PICK B (BEAM WIDTH) POSSIBLE FIRST WORDS.



STEP 2: COMBINE THE BEST OPTIONS FOR THE NEXT WORD



* V: VOCABULARY SIZE

BEST 3 PROBABILITIES: IN SEPTEMBER / JANE IS / JANE VISITS
 $P(y^{<1>}, y^{<2>} | x) = P(y^{<1>} | x) \cdot P(y^{<2>} | x, y^{<1>})$

STEP 3 AND NEXTS: REPEAT STEP 2 UNTIL <EOS>

IN SEPTEMBER JANE / JANE IS VISITING / JANE VISITS AFRICA

JANE VISITS AFRICA IN SEPTEMBER. <EOS>

* IF B=1: BEAM SEARCH = GREEDY SEARCH

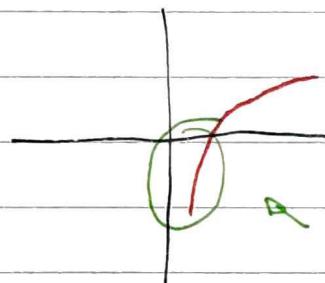
→ LENGTH NORMALIZATION

* MULTIPLY A LOT OF SMALL NUMBERS CAN CAUSE UNDERFLOW
 (NUMBERS SMALL ENOUGH TO BE STORED WITHOUT PRECISION)

$$\underset{Y}{\operatorname{argmax}} \prod_{t=1}^{T_y} P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

LOG:
CHANGES \prod TO \sum

$$\underset{Y}{\operatorname{argmax}} \sum_{t=1}^{T_y} \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$



PROBABILITIES ARE SMALL: SUM CAN RETRIEVE A
 LARGE NEGATIVE NUMBER

$$\underset{Y}{\operatorname{argmax}} \frac{1}{T_y} \sum_{t=1}^{T_y} \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)})$$

$\alpha = 0$: NO NORMALIZATION.
 $\alpha = 1$: FULL NORMALIZATION.

→ BEAM SEARCH DISCUSSION

- * LARGE B: BETTER RESULT, SLOWER
- * SMALL B: WORSE RESULT, FASTER
- * PRODUCTION: $B \approx 10$
- * ACADEMIC RESEARCH: $B \approx 1000 \sim 3000$

UNLIKE EXACT SEARCH ALGORITHMS LIKE
 BREADTH/DEPTH FIRST SEARCH, BEAM SEARCH
 RUNS FASTER, BUT IS NOT GUARANTEED TO
 FIND EXACT MAXIMUM FOR $\underset{Y}{\operatorname{argmax}} P(Y|X)$.

→ ERROR ANALYSIS ON BEAM SEARCH

* HOW TO FIGURE OUT IF RNN OR BEAM SEARCH IS THE MAJOR CAUSE OF ERRORS?

* EXAMPLE: "JANE VISITE L'AFRIQUE EN SEPTEMBRE"

→ HUMAN: JANE VISITS AFRICA IN SEPTEMBER. (y^*)

→ ALGORITHM: JANE VISITED AFRICA LAST SEPTEMBER (\hat{y})

COMPUTE $P(y^*|x)$ AND $P(\hat{y}|x)$ WITH RNN

CASE 1: $P(y^*|x) > P(\hat{y}|x)$

BEAM SEARCH CHOSE \hat{y} . BUT $P(y^*|x)$ ATTAINS HIGHER THAN $P(\hat{y}|x)$.

CONCLUSION: BEAM SEARCH IS AT FAULT.

CASE 2: $P(y^*|x) \leq P(\hat{y}|x)$

y^* IS A BETTER TRANSLATION THAN \hat{y} . BUT RNN PREDICTED $P(y^*|x) \leq P(\hat{y}|x)$.

CONCLUSION: RNN MODEL IS AT FAULT.

* ERROR ANALYSIS PROCESS

HUMAN	ALGORITHM	$P(y^* x)$	$P(\hat{y} x)$	AT FAULT?
ENGLISH 1-H	ENGLISH 1-A	2×10^{-10}	1×10^{-10}	BEAM
ENGLISH 2-H	ENGLISH 2-A	2×10^{-10}	6×10^{-10}	RNN
ENGLISH 3-H	ENGLISH 3-A	3×10^{-9}	2×10^{-11}	BEAM
:	:	:	:	:

FIGURES OUT WHAT FRACTION OF ERRORS ARE

"DUE TO" BEAM SEARCH VS. RNN MODEL.

→ BLEU SCORE (BILINGUAL EVALUATION UNDERSTUDY)

* EVALUATING MACHINE TRANSLATION

FRENCH: LE CHAT EST SUR LE TAPIS.

REFERENCE 1: THE CAT IS ON THE MAT.

REFERENCE 2: THERE IS A CAT ON THE MAT.

MACHINE TRANSLATION: THE THE THE THE THE THE THE

- FIRST PRECISION ATTEMPT: COUNT HOW MANY ^{TIME} PREDICTED WORDS APPEAR IN THE REFERENCES.

$$\text{PRECISION: } \frac{7}{7}$$

INCORRECT!

MODIFIED PRECISION: "CREDIT" THE WORDS WITH THE MAXIMUM NUMBER OF TIMES THEY APPEAR IN THE REFERENCES.

$$\text{PRECISION: } \frac{2}{7} \xrightarrow{\text{MORE REASONABLE.}} \text{COUNT}_{\text{clip}}(\text{"THE"})$$

→ COUNT ("THE")

→ BLEU SCORE ON BIGRAMS

* EXAMPLE: THE CAT THE CAT ON THE MAT.

	COUNT	COUNT _{clip}	
THE CAT (2x)	2	1	
CAT THE	1	0	
CAT ON	1	1	
ON THE	1	1	
THE MAT	1	1	

PRECISION = $\frac{4}{6}$



→ BLEU SCORE ON UNIGRAMS

* EXAMPLE:

REFERENCE 1: THE CAT IS ON THE MAT.

REFERENCE 2: THERE IS A CAT ON THE MAT.

MACHINE TRANSLATION: THE CAT THE CAT ON THE MAT. (\hat{y})

$$P_1 = \frac{\sum_{\text{UNIGRAM } e \in \hat{y}} \text{COUNT}_{\text{clip}}(\text{UNIGRAM})}{\sum_{\text{UNIGRAM } e} \text{COUNT}(\text{UNIGRAM})}$$

→ BLEU SCORE ON N-GRAMS

$$P_N = \frac{\sum_{N\text{-GRAM } e \in \hat{y}} \text{COUNT}_{\text{clip}}(N\text{-GRAM})}{\sum_{N\text{-GRAM } e} \text{COUNT}(N\text{-GRAM})}$$

* IF MT = REFERENCE: $P_1, P_2, \dots, P_N = 1.0$

→ BLEU DETAILS

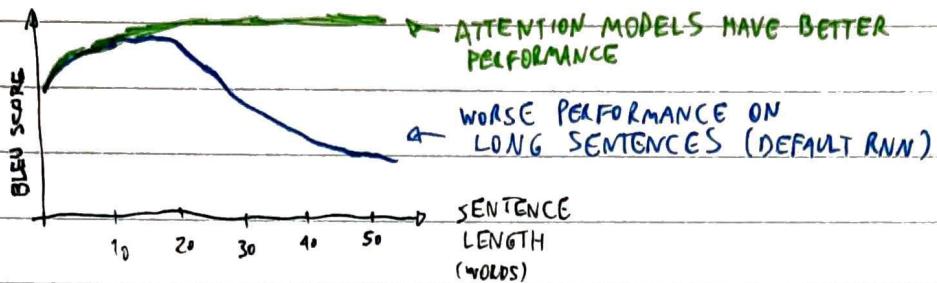
P_N - BLEU SCORE ON N-GRAMS ONLY

COMBINED BLEU SCORE: $BP \cdot \exp\left(\frac{1}{4} \sum_{n=1}^4 P_n\right)$

BP - BREVITY PENALTY

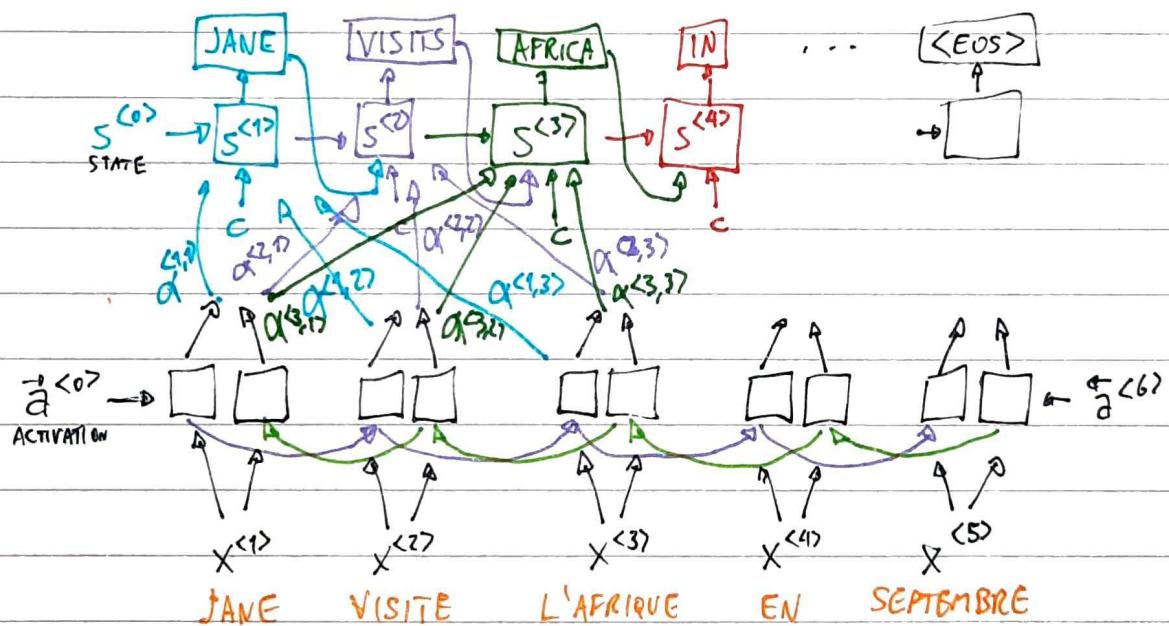
$$BP = \begin{cases} 1, & \text{IF } \text{MT-OUTPUT-LENGTH} > \text{REFERENCE-OUTPUT-LENGTH} \\ \exp\left(1 - \text{REFERENCE-OUTPUT-LENGTH}/\text{MT-OUTPUT-LENGTH}\right), & \text{OTHERWISE} \end{cases}$$

→ THE PROBLEM OF LONG SEQUENCES



ATTENTION MODELS TRY TO TRANSLATE PART BY PART: BETTER PERFORMANCE

→ ATTENTION MODEL INTUITION



* α : ATTENTION WEIGHT.

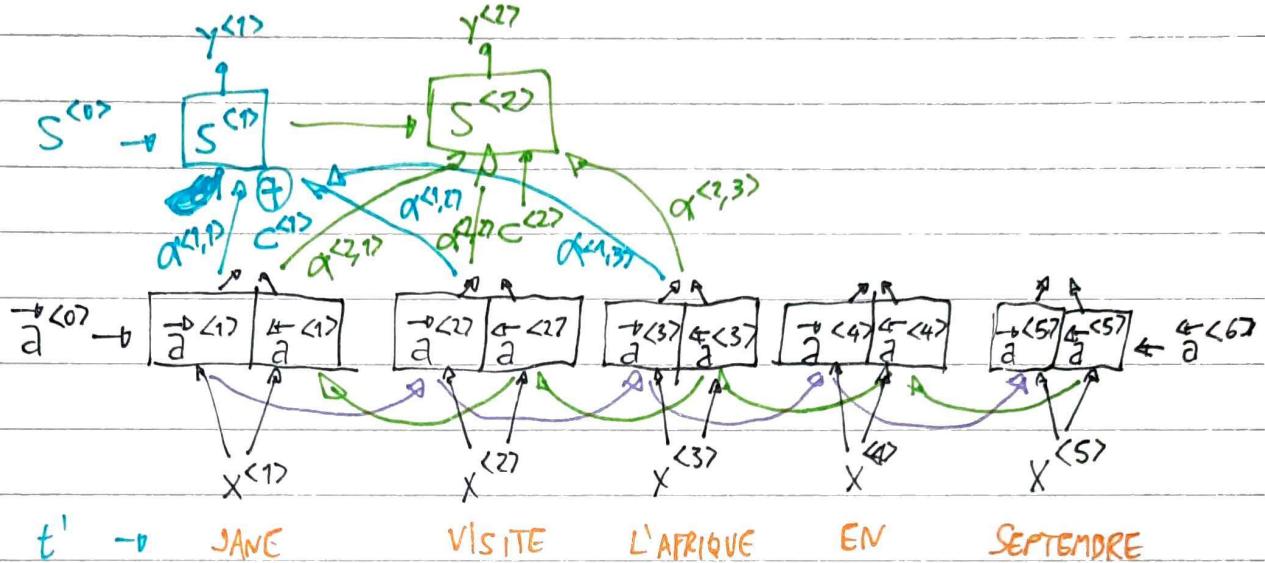
* $\alpha^{(j,w)}$: ATTENTION WEIGHT FROM WORD w TO GENERATE j .

FRENCH

ENGLISH

→ ATTENTION MODEL

* CONSIDERING $\vec{a}^{<t>} = (\vec{a}^{<t>}, \vec{a}^{<t>})$; t' FOR FRENCH SENTENCE Timestep; t FOR ENGLISH SENTENCE Timestep



$$\sum_{t'} \alpha^{<1,t'} = 1$$

$$c^{<1>} = \sum_{t'} \alpha^{<1,t'} a^{<t'>}$$

$$c^{<2>} = \sum_{t'} \alpha^{<2,t'} a^{<t'>}$$

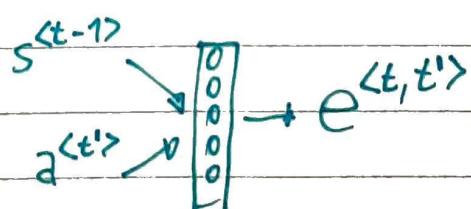
* $\alpha^{<t,t'>}$: AMOUNT OF "ATTENTION" $y^{<t>}$ SHOULD PAY TO $a^{<t'>}$.

→ COMPUTING ATTENTION $\alpha^{<t,t'>}$

$e^{<t,t'>}$: TRAINED PARAMETER

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

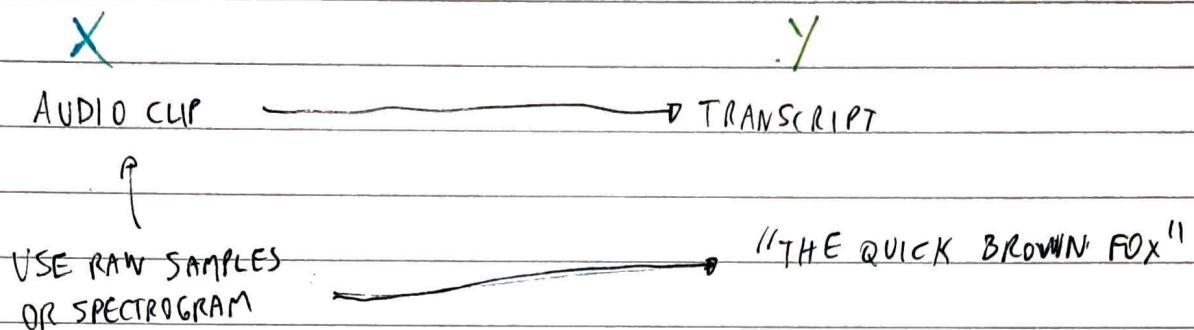
SOFTMAX!



* QUADRATIC COMPLEXITY ORDER (BUT IT'S ACTUALLY ACCEPTABLE)

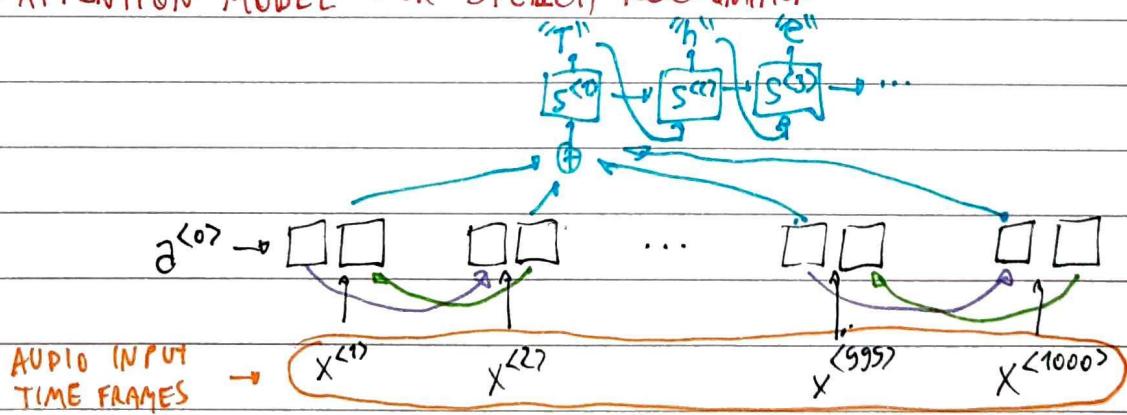
* EXAMPLE OF USE OF ATTENTION MODELS: DATE NORMALIZATION

→ SPEECH RECOGNITION PROBLEM



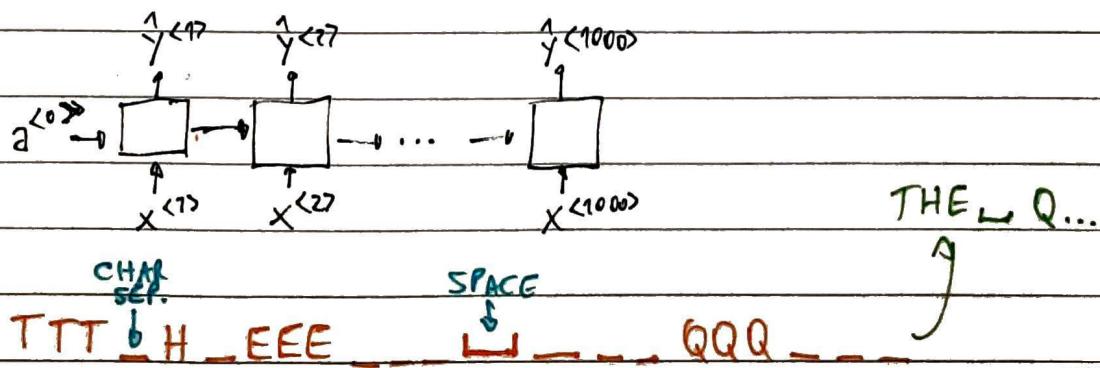
- * PAST APPROACHES USED PHONEMES; WITH END-TO-END DEEP LEARNING, THIS ISN'T NECESSARY ANYMORE
- * TRAINING ON LARGE DATASET: ~300 HOURS OF AUDIO
- * PRODUCTION SYSTEMS USUALLY TRAIN WITH 10000 ~ 100000 HOURS

→ ATTENTION MODEL FOR SPEECH RECOGNITION



→ CTC (CONNECTIONIST TEMPORAL CLASSIFICATION) COST FOR SPEECH RECOGNITION

* EXAMPLE: "THE QUICK BROWN FOX" → (19 CHARACTERS)

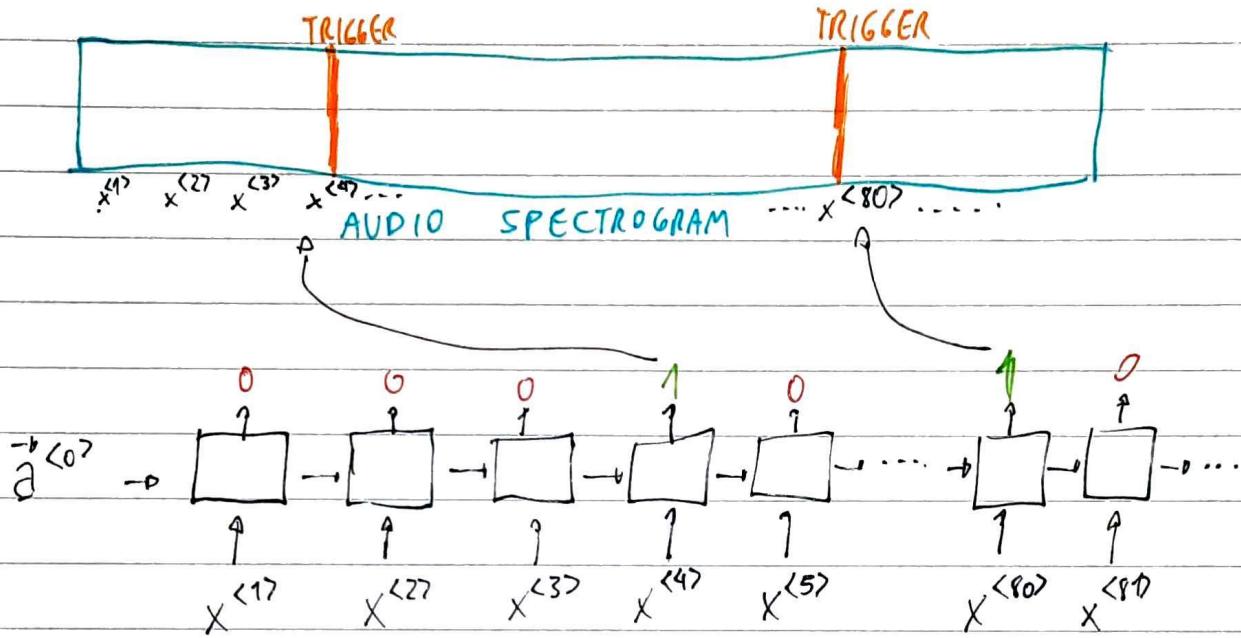




→ TRIGGER WORD DETECTION

* EXAMPLE: AMAZON ECHO (ALEXA), GOOGLE HOME (OK GOOGLE)

* ALGORITHM:



* PROBLEMS:

- UNBALANCED DATASET
- MORE 0'S THAN 1'S

* FIXING (HACKING) THE PROBLEM:

- EXTEND 1'S FOR A LITTLE BIT

0 0 0 1 1 1 1 0 0
 $X^{(1)}$ $X^{(2)}$ $X^{(3)}$ $X^{(4)}$ $X^{(5)}$ $X^{(6)}$ $X^{(7)}$ $X^{(8)}$ $X^{(9)}$...

→ TRANSFORMER NETWORK MOTIVATION

RNN

GRU

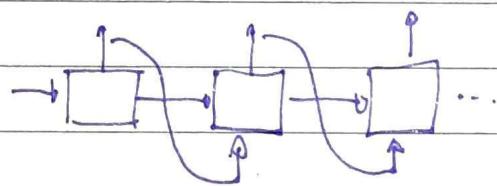
LSTM

INCREASED COMPLEXITY
SEQUENTIAL (ONE TOKEN AT A TIME)

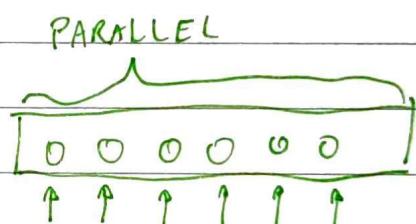
* HOW TO MAKE IT MORE EFFICIENT? IS THERE A WAY TO MAKE
A PARALLEL INPUT PROCESSING?

→ TRANSFORMER NETWORK INTUITION

* COMBINE ATTENTION WITH CNN STYLE OF PROCESSING



RNN



CNN

* ATTENTION + CNN

SELF-ATTENTION: COMPUTE N REPRESENTATIONS FOR
A SENTENCE WITH N WORDS IN PARALLEL

MULTI-HEAD ATTENTION: BASICALLY A "FOR LOOP"
OVER SELF-ATTENTION. MULTIPLE VERSIONS OF REPRESENTATIONS

→ SELF-ATTENTION INTUITION

$A(q, K, V)$ = ATTENTION-BASED VECTOR REPRESENTATION OF A WORD
 ↳ CALCULATE FOR EACH WORD $A^{(1)}, \dots, A^{(N)}$

RNN ATTENTION:

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

$A^{(1)}, A^{(2)}, A^{(3)}, A^{(4)}, A^{(5)}$
 JANE VISITE L'AFRIQUE EN SEPTEMBRE.

TRANSFORMERS ATTENTION

$$A(q, K, V) = \sum_i \frac{\exp(q \cdot K^{(i)})}{\sum_j \exp(q \cdot K^{(j)})} \cdot V^{(i)}$$

→ CALCULATING $A(q, K, V)$

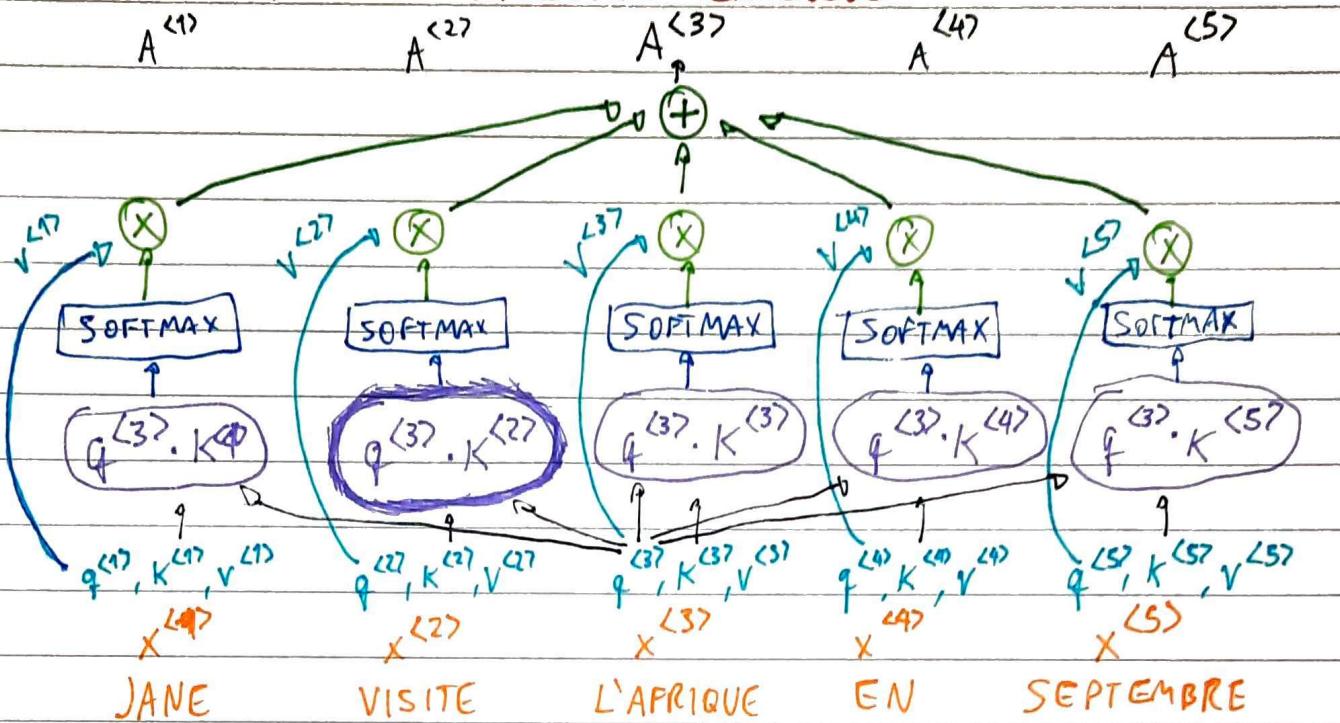
QUERY (q)	KEY (K)	VALUE (V)
$q^{(1)}$	$K^{(1)}$	$V^{(1)}$
$q^{(2)}$	$K^{(2)}$	$V^{(2)}$
$q^{(3)}$	$K^{(3)}$	$V^{(3)}$
$q^{(4)}$	$K^{(4)}$	$V^{(4)}$
$q^{(5)}$	$K^{(5)}$	$V^{(5)}$

5 WORDS.

- * ANALOGY WITH DATABASES: "QUERY" MIGHT BE A QUESTION;
 "KEY-VALUE" SHOULD BE LIKE "ANSWERS"
- * q, K AND V CALCULATED WITH TRAINED PARAMETERS

$$q^{(N)} = W^q \cdot x^{(N)} \quad | \quad K^{(N)} = W^K \cdot x^{(N)} \quad | \quad V^{(N)} = W^V \cdot x^{(N)}$$

→ EXAMPLE OF SELF-ATTENTION CALCULUS

* CALCULATING $A^{(3)}$:* $q^{(3)}$ MIGHT REPRESENT A QUESTION LIKE "WHATS HAPPENING THERE?"* MAYBE $K^{(2)}$ SHOULD BE THE BEST ANSWER.

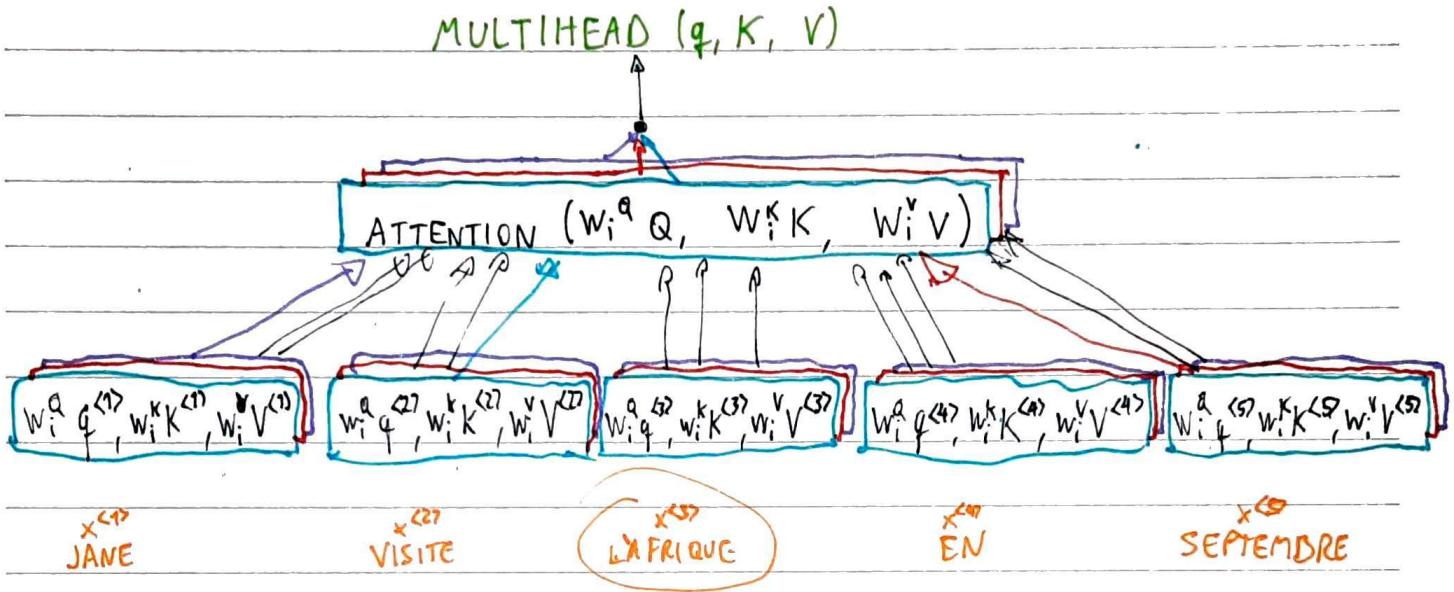
"THERE IS HAPPENING A VISIT".

 $q^{(3)}.K^{(2)}$ SHOULD HAVE THE GREATEST CONTRIBUTION.

$$\text{ATTENTION}(q, k, v) = \text{SOFTMAX}\left(\frac{qk^T}{\sqrt{d_k}}\right)v$$

→ MULTI-HEAD ATTENTION

* "HEAD": EACH SELF-ATTENTION CALCULATED



* EXAMPLES OF WHAT EACH HEAD SHOULD REPRESENT

HEAD 1: WHAT'S HAPPENING?

$$W_1^q, W_1^K, W_1^V$$

BEST ANSWER: $K^{(2)}$

HEAD 2: WHEN?

$$W_2^q, W_2^K, W_2^V$$

BEST ANSWER: $K^{(5)}$

HEAD 3: WHO?

$$W_3^q, W_3^K, W_3^V$$

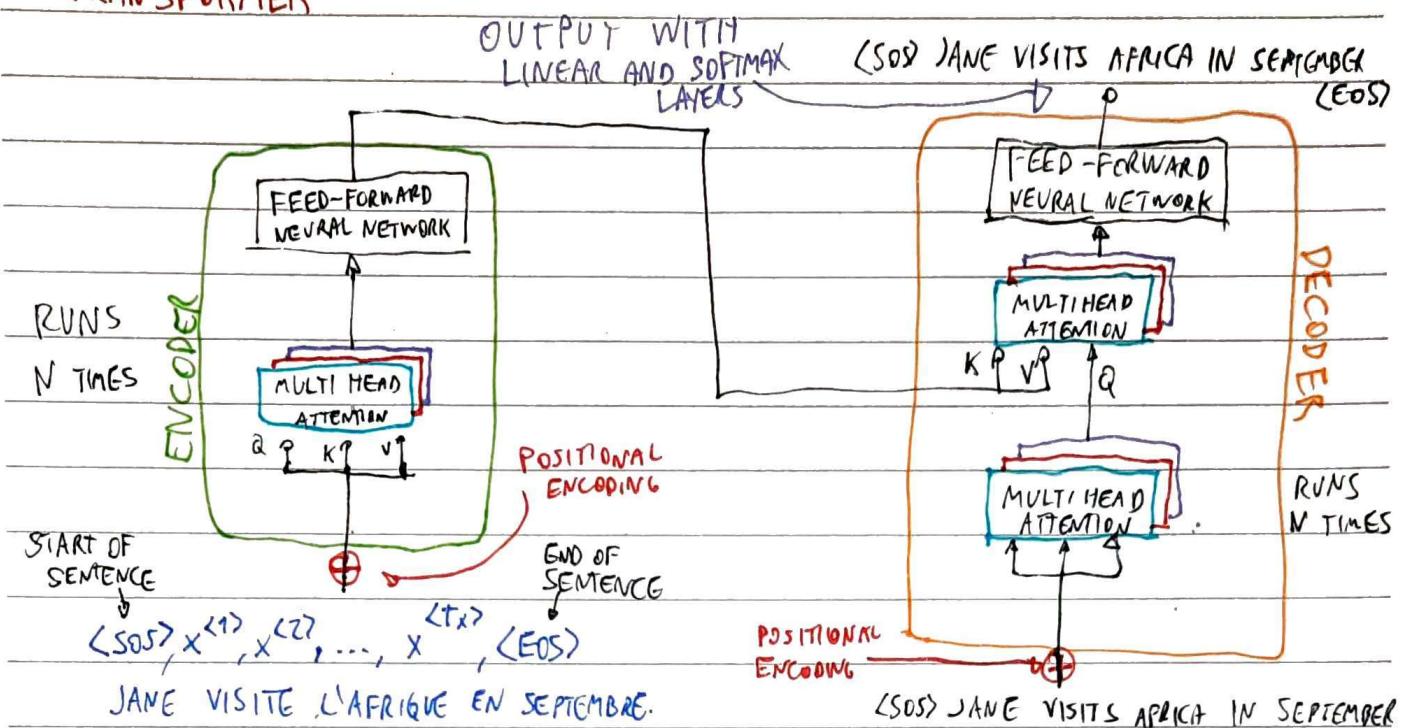
BEST ANSWER: $K^{(1)}$

MULTIHEAD (q, K, V) = CONCAT (HEAD₁, HEAD₂, ..., HEAD_h) W_o

HEAD_i = ATTENTION ($W_i^q q, W_i^K K, W_i^V V$)

- h: NUMBER OF HEADS

→ TRANSFORMER



→ TRANSFORMER DETAILS

* POSITIONAL ENCODING

$$\text{PE}(\text{pos}, z_i) = \sin\left(\frac{\text{pos}}{10000^{\frac{2i}{d}}}\right)$$

i: FEATURE OF THE EMBEDDING VECTOR

POS: POSITION OF WORD ON THE SENTENCE

$$\text{PE}(\text{pos}, z_{i+1}) = \cos\left(\frac{\text{pos}}{10000^{\frac{2i}{d}}}\right)$$

- POSITIONAL ENCODING CREATES A UNIQUE POSITION VECTOR FOR EACH WORD, WHICH ENCODES THE POSITION OF THE WORD
- POSITIONAL ENCODINGS COULD ALSO BE PASSED THRU RESIDUAL CONNECTIONS

* ADD & NORM

- "BATCHNORM-LIKE" LAYER

* MASKED MULTI-HEAD ATTENTION

- TECHNIQUE FOR TRAINING