

RECURRENT NEURAL NETWORKS

SEQUENCE PROBLEMS

IN	OUT	PURPOSE
Mr. Brown	THE QUICK BROWN FOX JUMPED...	SPEECH RECOGNITION
∅	♪ ♪ ♪ ♪ ♪	MUSIC GENERATION
THERE IS NOTHING TO LIKE IN THIS MOVIE	★ ★ ★ ★	SENTIMENT CLASSIFICATION
AGCCCCCTGTG AGGAACCTAG	AGCCCCCTGTG AGGAACCTAG	DNA SEQUENCE ANALYSIS
Voulez-vous chanter avec moi?	Do you want to sing with me?	MACHINE TRANSLATION
🏃‍♂️ 🏃‍♀️ 🏃	RUNNING	VIDEO ACTIVITY RECOGNITION
Yesterday Harry Potter met Hermione Granger	Yesterday Harry Potter met Hermione Granger	NAME ENTITY RECOGNITION

NAME ENTITY RECOGNITION

$x = \text{HARRY POTTER AND HERMIONE}$ $T_x = 9$
 $x^{<1>} x^{<2>} \dots$ (9 words)

GRANGER INVENTED A NEW SPELL

$$y = \begin{matrix} 1 & 1 & 0 & 1 & T_y = T_x \\ y^{<1>} & y^{<2>} & \dots & y^{<T>} & \end{matrix}$$

EXAMPLE OF A PROBLEM WHERE
EVERY $x^{<i>}$ HAS AN OUTPUT $y^{<i>}$

HOW DO WE REPRESENT WORDS?

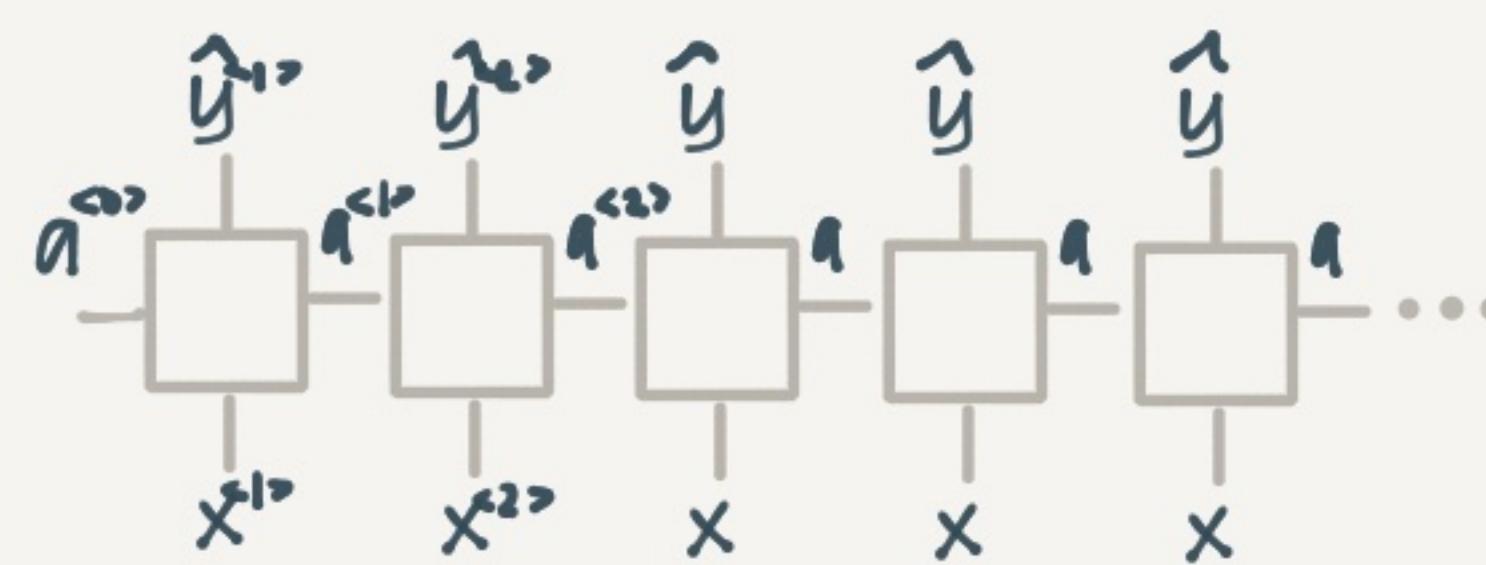
CREATE A VOCABULARY (EG 10K MOST COMMON WORDS IN YOUR TEXTS • OR DOWNLOAD EXISTING)

aaron	1	EACH WORD IS A ONE-HOT.
and	2	VECTOR
Harry	367	$\underline{\text{HARRY}} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$
Potter	4075	
Zulu	6830	
	10000	

WE COULD USE A STANDARD NETWORK BUT...

- (A) INPUT & OUTPUTS CAN HAVE DIFFERENT LENGTHS IN DIFF EXAMPLES
- (B) WE DON'T SHARE FEATURES LEARNED ACROSS DIFFERENT POSITIONS

RECURRENT NEURAL NET (RNN)



PREVIOUS RESULTS ARE PASSED IN AS INPUTS SO WE GET CONTEXT.

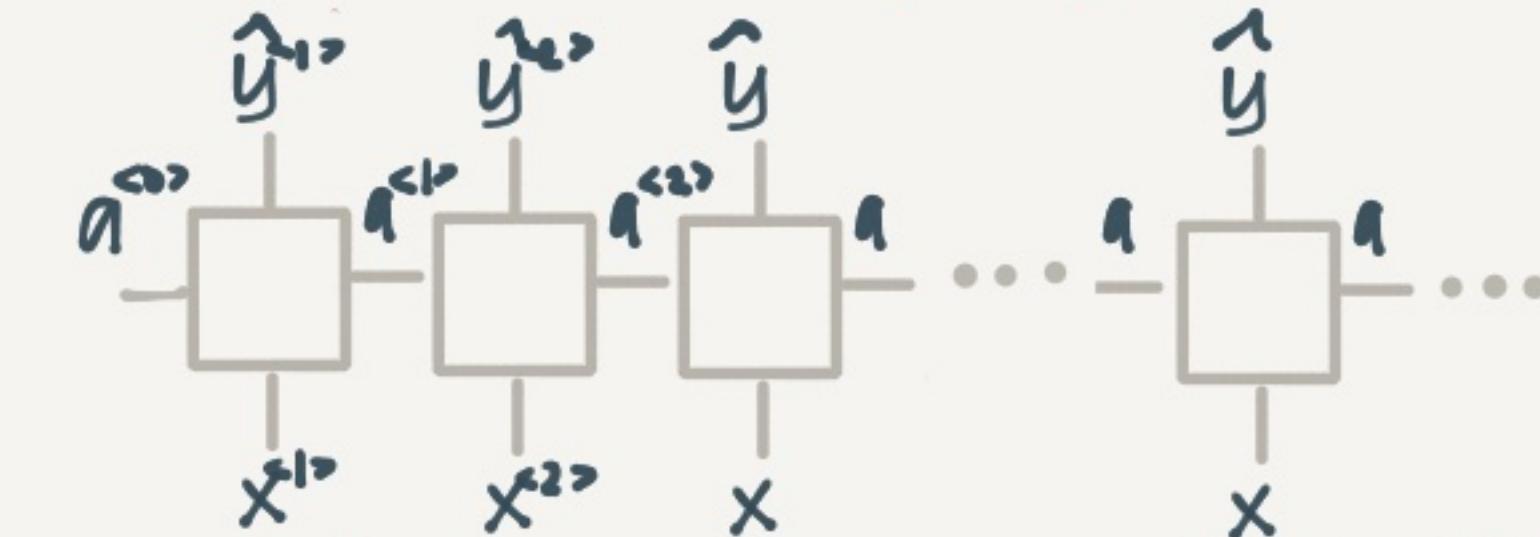
$$\begin{aligned} q^{<1>} &= g_1(W_1[a^{<0>}; x^{<1>}] + b_1) && \text{TANH / RELU} \\ \hat{y}^{<1>} &= g_2(W_{21}q^{<1>} + b_2) && \text{SIGMOID} \end{aligned}$$

THE SAME W & b ARE USED IN ALL TIME STEPS

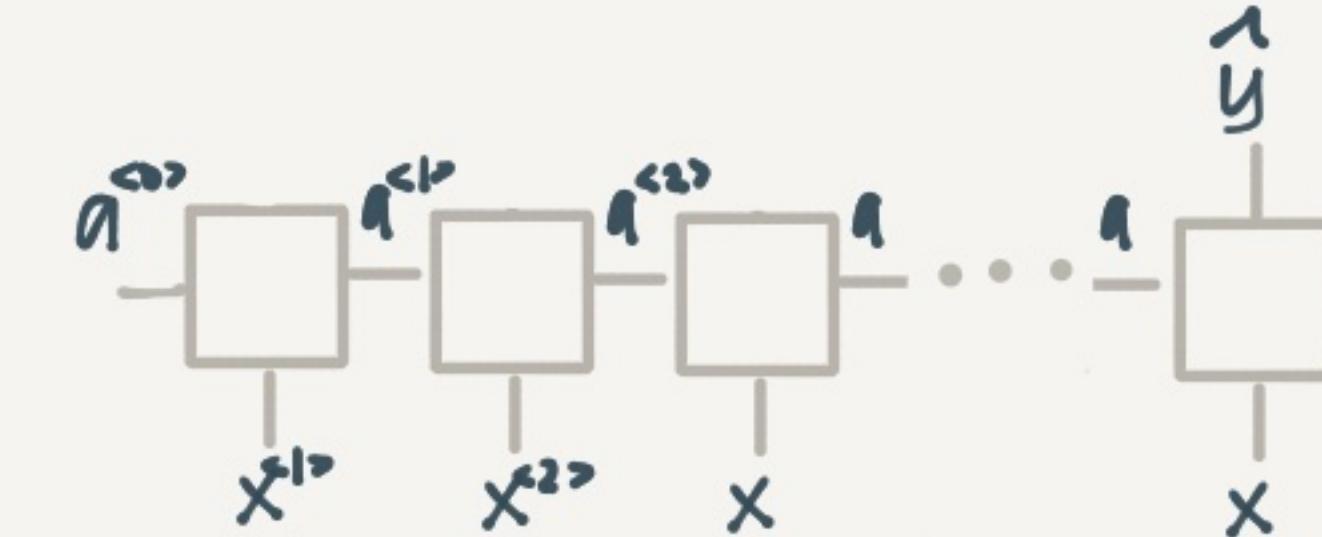
THE LOSS WE OPTIMIZE IS THE SUM OF $\mathcal{L}(\hat{y}, y)$ FROM 1-T

DIFFERENT TYPES OF RNN

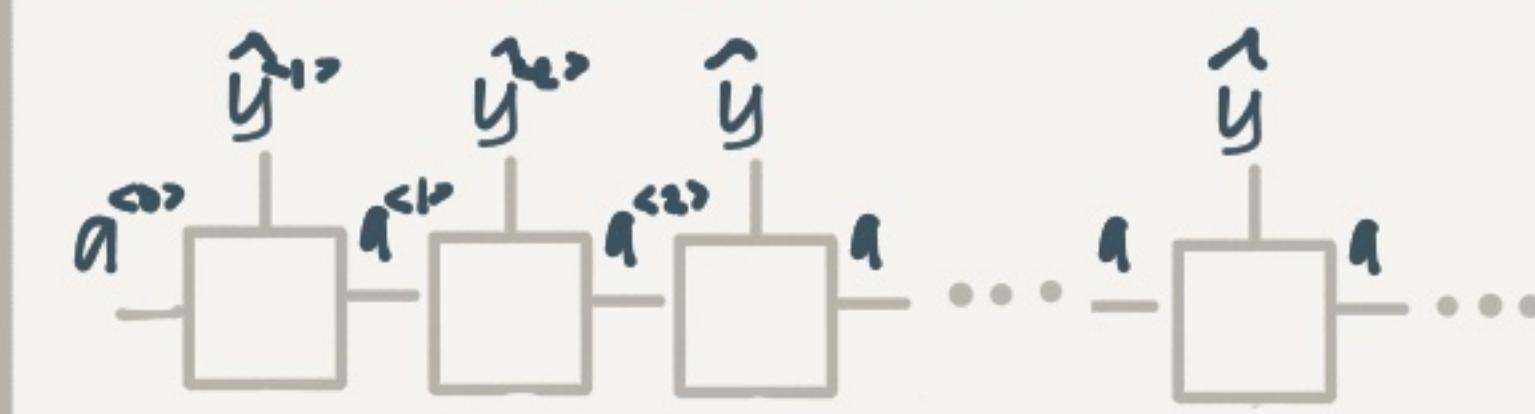
MANY-TO-MANY $T_x = T_y$



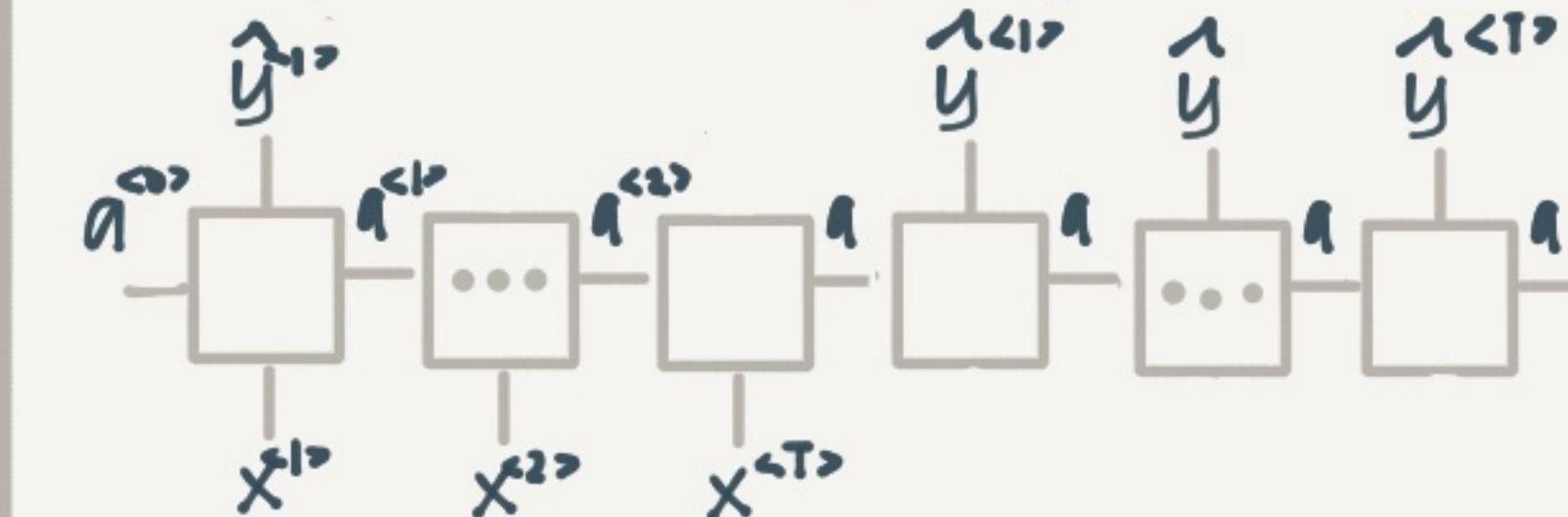
MANY-TO-ONE EX. SENTIMENT ANALYSIS



ONE-TO-MANY • MUSIC GENERATION



MANY-TO-MANY $T_x \neq T_y$ TRANSLATION



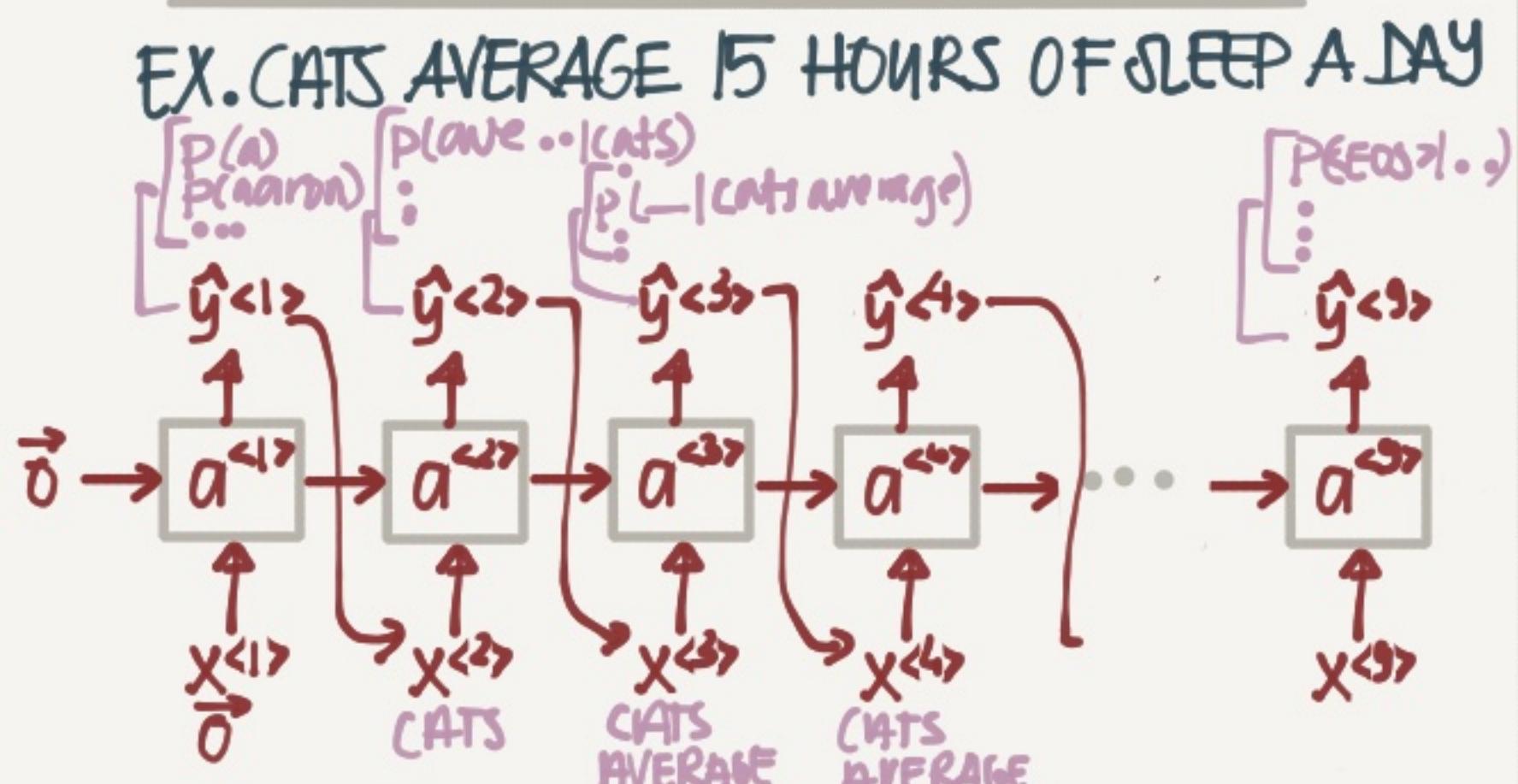
MORE ON RNNs

LANGUAGE MODELLING

HOW DO YOU KNOW IF SOMEONE SAID

THE APPLE AND PAIR SALAD OR
THE APPLE AND PEAR SALAD?

THE PURPOSE OF A LANG. MODEL IS TO
CALCULATE THE PROBABILITIES



SO GIVEN: CATS AVERAGE IS. WHAT IS THE PROB.
THE NEXT WORD IS HOURS?

SAMPLING SENTENCES

1. TRAIN ON ALL HARRY POTTER BOOKS.
2. RANDOMLY SELECT A WORD (ON OF THE TOP WORDS)
(EX. THE)
3. PASS THIS INTO THE NEXT TIMESTAMP
AND SAMPLE A NEW WORD
4. REPEAT UNTIL X WORDS OR YOU
REACHED <EOS>

CAN DO AT
CHARACTER LEVEL
AS WELL

VANISHING GRADIENTS

THE CAT, WHO ALREADY ATE APPLES AND ORANGES
AND A FEW MORE THINGS BUT ~~BU~~ WAS FULL
THE CATS ~~W~~ ALREADY ATE ...
... ~~W~~ WERE FULL

NEED TO REMEMBER
SING/PLURAL FOR A LONG
TIME

SINCE LONG SENTENCE \Rightarrow DEEP RNN
WE GET THE VANISHING GRADIENTS PROB WE
HAVE IN STANDARD NNs - I.E. THE GRADIENTS
FOR CAT/CATS HAVE LITTLE OR NO EFFECT
ON WAS/WERE.

NOTE: SOMETIMES YOU SEE EXPLODING GRAD
(AS OVERFLOW NAN) BUT THIS IS EASILY FIXED
WITH GRADIENT CLIPPING

GATED RECURRENT UNIT GRU
HELPS RECALL IF CAT WAS SING.
OR PLURAL



THE GRU ACTS AS A MEMORY
— AT EVERY Timestep IT
CALCULATES A NEW \tilde{c} TO STORE
AND A GATE Γ_u DECIDES TO
UPDATE c TO \tilde{c} OR NOT

YAY! YOU ARE NOW
YOUR OWN J.K. ROWLING

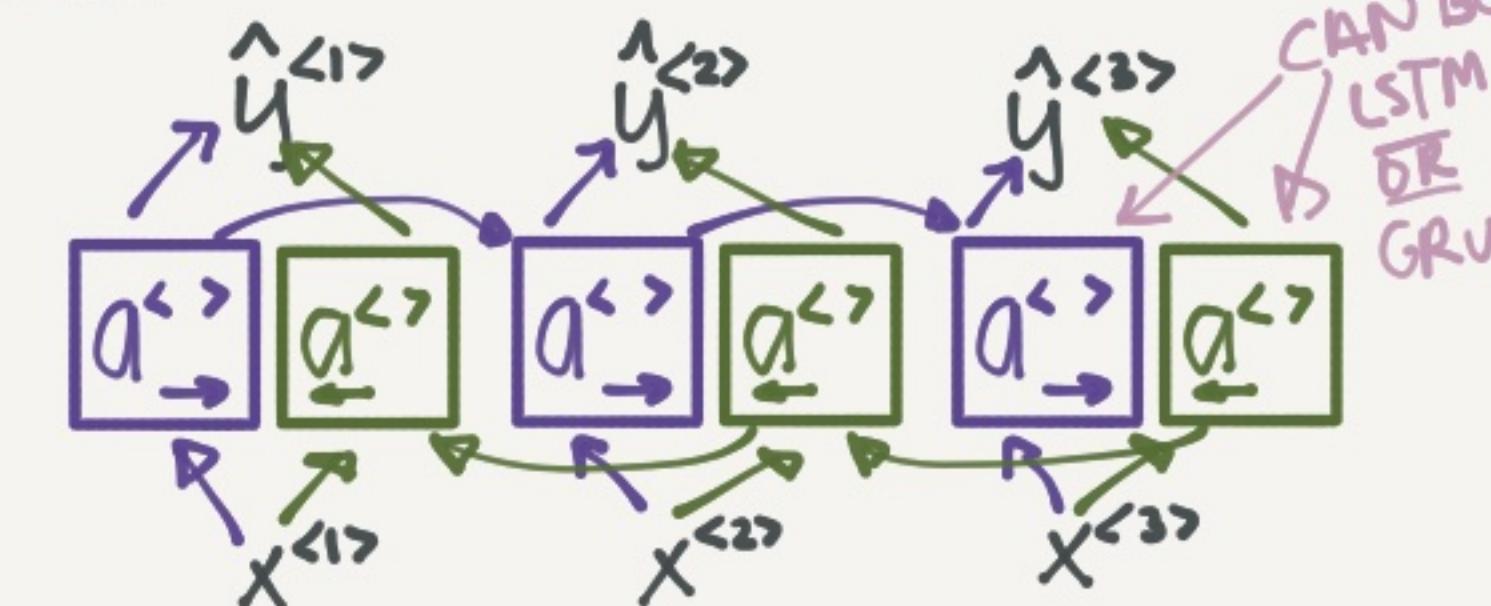
LONG SHORT TERM MEMORY (LSTM)

THE LSTM IS A VARIATION ON
THE SAME THEME AS GRU
BUT WITH AN ADDITIONAL Γ_f
FORGET GATE

BI-DIRECTIONAL RNNs (BRNN)

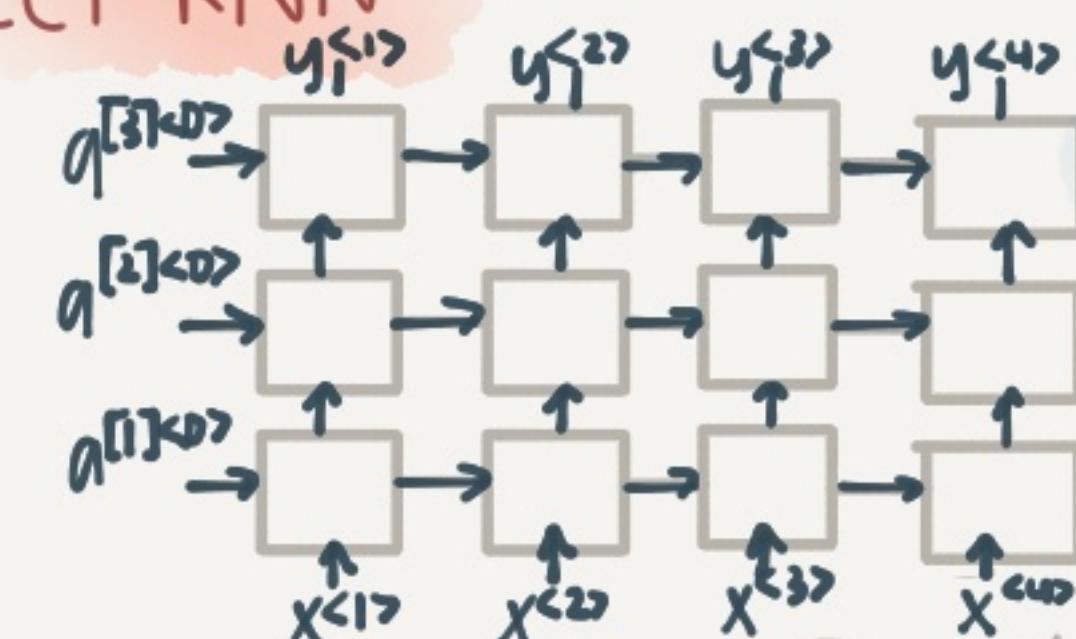
HE SAID, 'TEDDY BEARS ARE ON SALE'
HE SAID, 'TEDDY ROOSEVELT WAS A
GREAT PRESIDENT'

PROBLEM: WITHOUT LOOKING FORWARD WE
CAN'T SAY IF TEDDY IS A TOY OR A NAME



ONE DISADVANTAGE IS THAT YOU NEED
THE FULL SENTENCE BEFORE YOU BEGIN-
SO NOT SUITABLE FOR LIVE SPEECH RECO

DEEP RNN



SINCE THEY
ARE ALREADY
MEMORY,
DEEP THEY
USUALLY
DON'T HAVE
A LOT OF
LAYERS

NLP & WORD EMBEDDINGS

MAN IS TO WOMAN AS
KING IS TO QUEEN

PROBLEM: THE ONE-HOT REPR \mathbf{q}_apple OF
APPLE HAS NO INFO ABOUT ITS RELATIONSHIP
TO $\mathbf{o}_{\text{orange}}$

I WANT A GLASS OF ORANGE —
I WANT A GLASS OF APPLE —

SOLUTION: CREATE A MATRIX OF
FEATURES TO DESCRIBE THE WORDS

WORD EMBEDDINGS

	MAN	WOMAN	KING	QUEEN	APPLE	ORANGE
5391	1	-0.95	0.97	0.00	0.01	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
:						
e_{5391}						

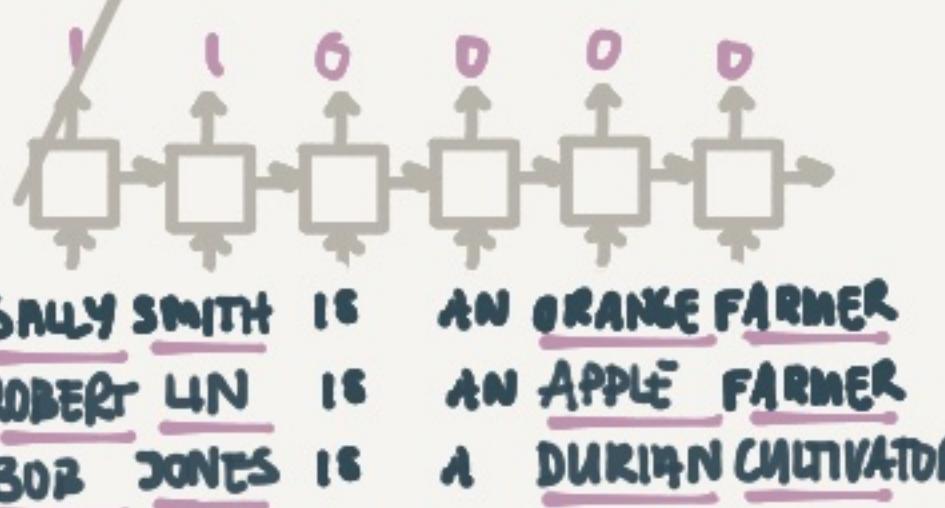
IN REALITY • THE FEATURES ARE
LEARNED & NOT AS STRAIGHTFWD
AS GENDER/AGE

• MAN	• woman	• dog
• king	• cat	• fish
• queen	• apple	• orange
• four	• grape	
• three	• orange	
• one		
• two		

t-SNE
VISUAL
REPRESENT
OF 300D
WORD
EMBEDDINGS

USING WORD EMBEDDINGS

EX. NAME/ENTITY RECOGN



WITH WORD EMBEDDINGS WE
UNDERSTAND THAT AN ORANGE
FARMER IS A PERSON \Rightarrow SALLY
SMITH = NAME

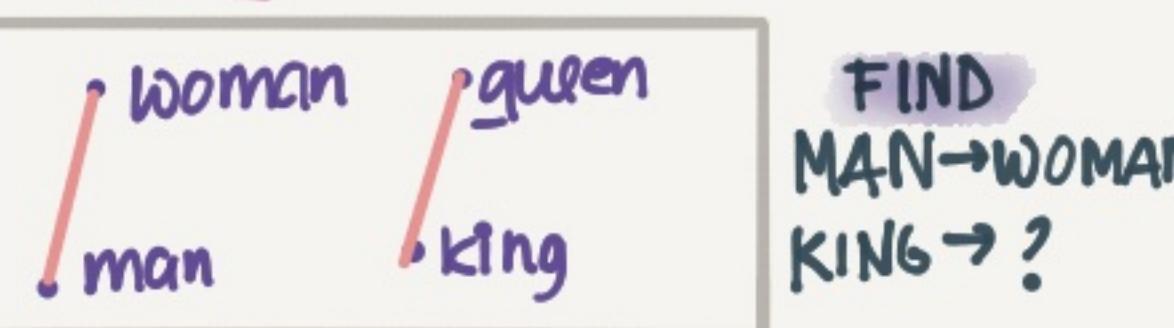
- APPLE ~ ORANGE \Rightarrow PERSON
- USING WORD EMBEDDINGS TRAINED
ON LOTS OF TEXT WE ALSO GET EMB
FOR MORE UNCOMMON WORDS
(DURIAN, CULTIVATOR)

EX. MAN IS TO WOMAN AS
KING IS TO ?

$E = \text{EMBEDDING MATRIX}$

	MAN	WOMAN	KING	QUEEN
e_{man}	5391	9853	9914	7157
GENDER	-1	1	-0.95	0.97
ROYAL	0.01	0.02	0.93	0.95
AGE	0.03	0.02	0.7	0.69
FOOD	0.04	0.01	0.02	0.01
...				
e_{5391}				

$$\begin{bmatrix} e_{\text{man}} - e_{\text{woman}} \\ e_{\text{king}} - e_{\text{queen}} \end{bmatrix} \xrightarrow{\text{EVERY SIMILAR}} \begin{bmatrix} -2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$



$$\text{FIND}(w) : \arg \max \text{sim}(e_w, e_{\text{king}} - e_{\text{man}} + e_{\text{woman}})$$

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

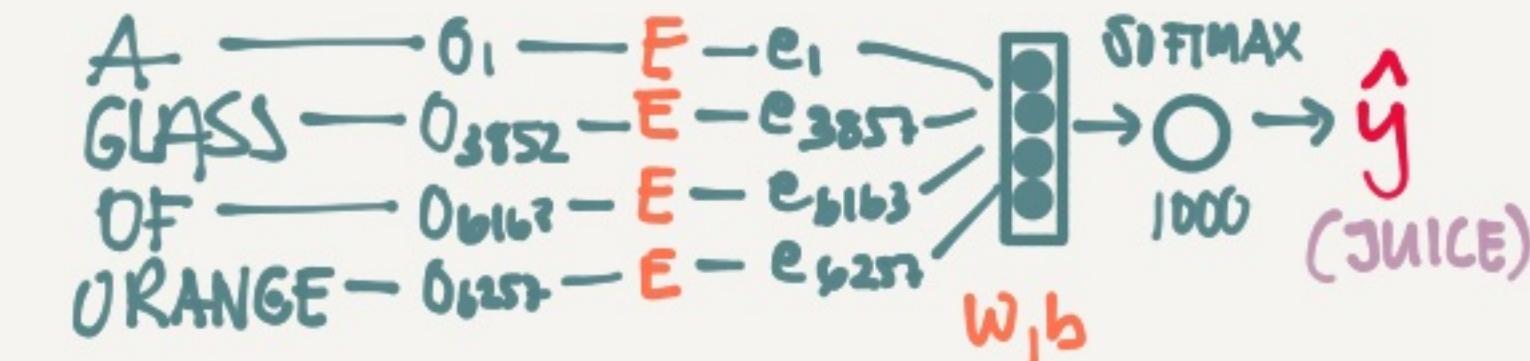
COSINE SIMILARITY

LEARNING WORD EMBEDDINGS

HOW DO WE LEARN THE EMBEDDING MATRIX E ?

IDEA1: USING A NEURAL LANG MODEL

I WANT A GLASS OF ORANGE \hat{y}



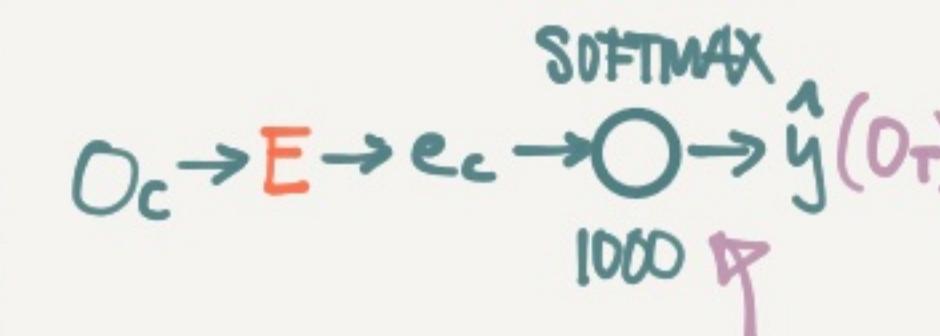
WE CAN USE DIFFERENT CONTEXTS THAN THE LAST 4 WORDS

- LAST 4 WORDS
 - 4 WORDS LEFT+RIGHT
 - LAST 1 WORD
 - NEARBY 1 WORD
- SKIPGRAM**
RANDOM WITHIN EX 5 WORDS

IDEA2: SKIP-GRAMS WORD2VEC

I WANT A GLASS OF ORANGE JUICE TO GO ALONG WITH MY CEREAL
PICK RANDOM CONTEXT/TARGET PAIRS (WITHIN EX 5 WORDS)

CONTEXT	TARGET
ORANGE	JUICE
ORANGE	GLASS
ORANGE	MY
...	...



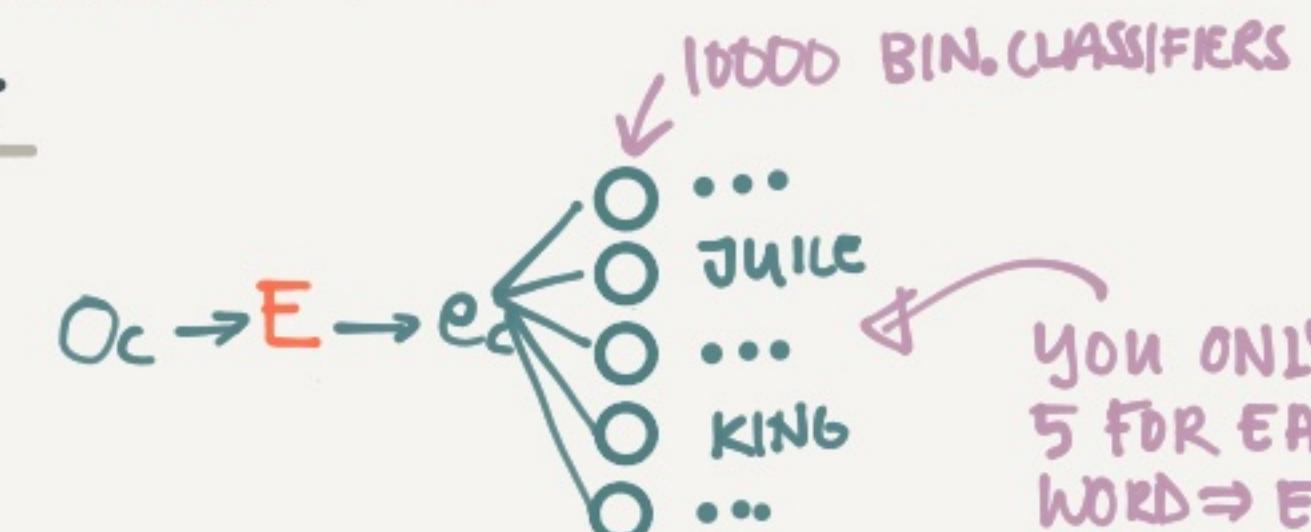
NOTE: WHILE THIS
SIMPLE NN PREDICTS O_t
OUR REAL GOAL IS TO
LEARN E

THIS IS VERY COMPUTATIONALLY EXPENSIVE
BUT WE CAN OPTIMIZE BY USING A HIERARCHICAL
SOFTMAX CLASSIFIER

IDEA: NEGATIVE SAMPLING

1. PICK A CONTEXT/TARGET PAIR AS A POSITIVE EXAMPLE
2. PICK A FEW NEG EXAMPLES CONTEXT + RANDOM

CONTEXT	WORD	TARGET
ORANGE	JUICE	1
ORANGE	KING	0
FRANGE	BOOK	0
ORANGE	THE	0
ORANGE	OF	0



NOTE: SOMETIMES BY
CHANCE YOU PICK A
POS PAIR • BUT IT DOESN'T
MATTER

YOU ONLY TRAIN
5 FOR EACH CONTEXT
WORD \Rightarrow EFFICIENT
TO TRAIN

WORD EMBEDDINGS

CONTINUED...

Glove WORD VECTORS

$x_{ij} = \# \text{TIMES WORD } i \text{ APPEARS IN THE CONTEXT OF } j$

TARGET CONTEXT
(HOW RELATED THEY ARE)

$$\text{MINIMIZE } \sum_{i=1}^{10k} \sum_{j=1}^{10k} f(x_{ij})(\theta_i^T e_j + b_i + b_j - \log x_{ij})^2$$

IF NO CONTEXT
(ALSO HELPS WEIGHING VERY FREQ WORDS (THE, OF...) & VERY INFREQUENT (PURPLE))

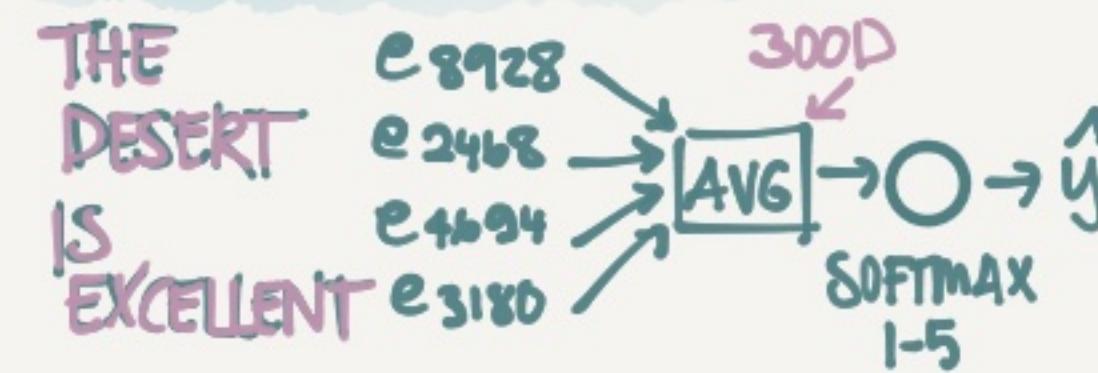
EVERYTHING LED UP TO THIS VERY SIMPLE ALGORITHM

SENTIMENT CLASSIFICATION

X	Y
THE DESSERT IS EXCELLENT	★★★☆
SERVICE WAS QUITE SLOW	★☆
GOOD FOR A QUICK MEAL BUT NOTHING SPECIAL	★☆☆
COMPLETELY LACKING IN GOOD TASTE, GOOD SERVICE AND GOOD AMBIENCE	*

PROBLEM: YOU MAY NOT HAVE A LARGE DATASET
BUT YOU CAN USE AN EMBEDDING MATRIX E THAT IS ALREADY PRE-TRAINED

IDEA: SIMPLE CLASSIFICATION

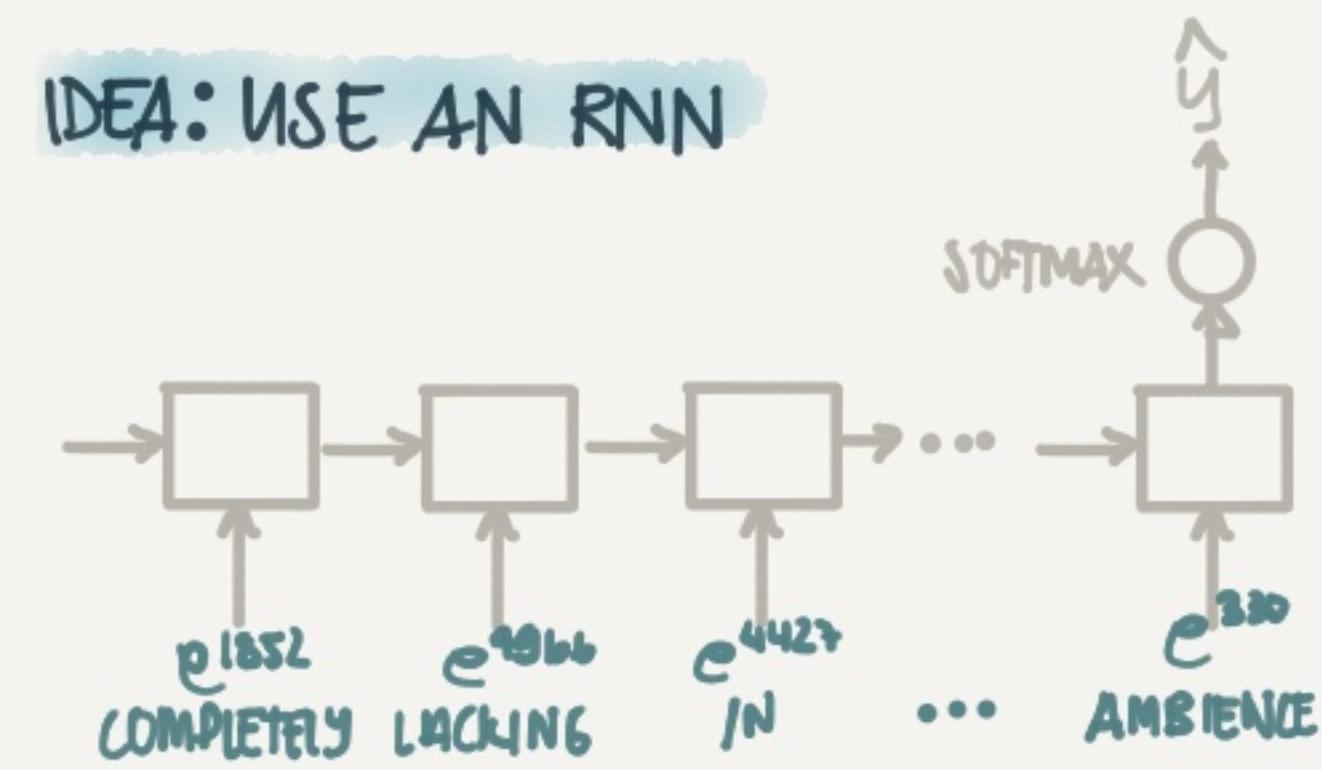


WORKS WELL FOR SHORT SENTENCES BUT DOESN'T TAKE ORDER INTO ACCOUNT

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

THIS MAY BE SEEN AS A ~~++~~ REVIEW

IDEA: USE AN RNN



THIS CAN NOW TAKE INTO ACCOUNT THAT COMPLETELY LACKING NEGATES THE WORD GOOD

ELIMINATING BIAS IN WORD EMBEDDINGS

MAN IS TO COMPUTER PROGRAMMER AS WOMAN IS TO HOME MAKER

SOMETIMES THE TEXT CONTAINS ♂ ALBOS LEARN A GENDER, RACE, AGE... BIAS WE DON'T WANT OUR MODELS TO HAVE. EX. HIRING BASED ON GENDER, SENTENCING BASED ON RACE ETC.

ADDRESSING BIAS

1. IDENTIFY BIAS DIRECTION

♂ he → e_he
♀ male → e_female

2. NEUTRALIZE

FOR EVERY WORD THAT IS NOT DEFINITIONAL (GIRL, BOY, HE, SHE...) PROJECT TO GET RID OF BIAS



3. EQUALIZE PAIRS

THE ONLY DIFF BETWEEN EX GIRL/BOY SHOULD BE GENDER

HOW DO YOU KNOW WHICH WORDS TO NEUTRALIZE?

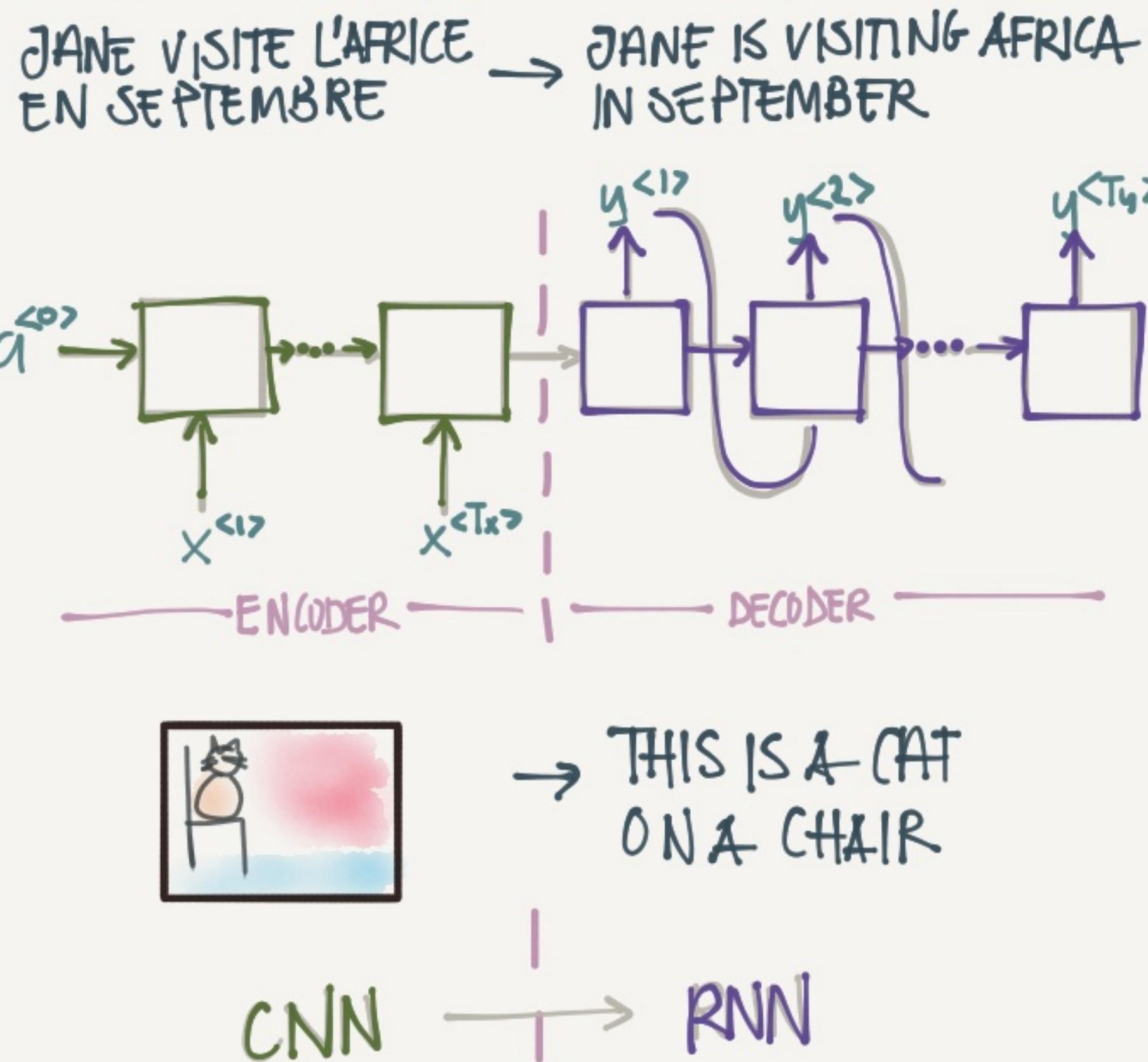
DOCTOR, BEARD, SEWING MACHINE?

A: BY TRAINING A CLASSIFIER TO FIND OUT IF A WORD IS DEFINITIONAL

TURNED OUT THE # OF PAIRS IS FAIRLY SMALL SO YOU CAN EVEN HAND PICK THEM

SEQUENCE TO SEQUENCE

BASIC MODELS



HOW DO YOU PICK THE MOST LIKELY SENTENCE?

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

WE DON'T WANT A RANDOMLY GENERATED SENTENCE (WE WOULD SOMETIMES GET A GOOD, SOMETIMES BAD)
INSTEAD WE WANT TO MAXIMIZE

$$\text{ARG MAX } P(y^{<1>} | \dots | y^{<Ty>} | x)$$

IDEA: USE GREEDY SEARCH

1. PICK THE WORD WITH THE BEST PROBABILITY
2. REPEAT UNTIL DEAD

WITH THIS WE COULD GET

- JANE IS GOING TO BE VISITING AFRICA THIS SEPTEMBER

INSTEAD OF

- JANE IS VISITING AFRICA THIS SEPTEMBER

SOLUTION

OPTIMISE THE PROB OF THE WHOLE SENTENCE INSTEAD

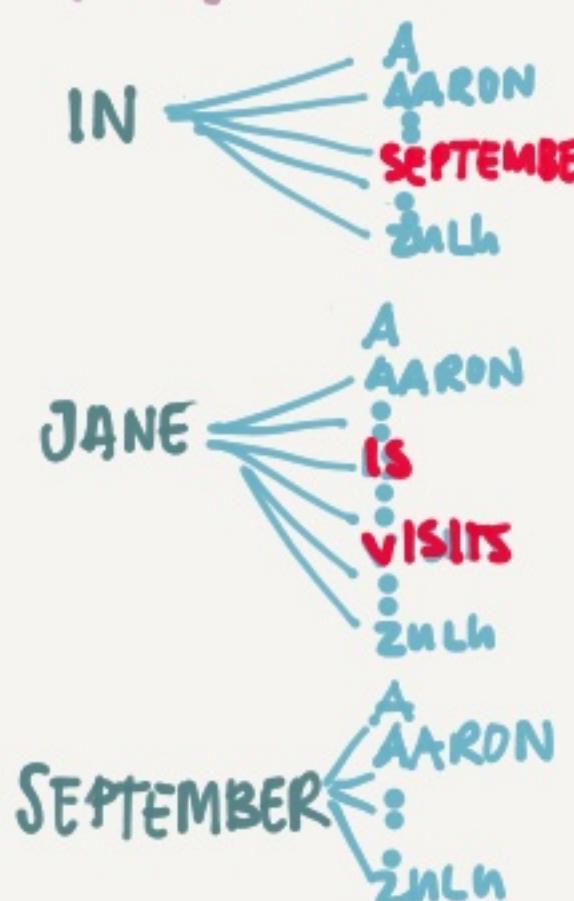
BEAM SEARCH

1. PICK THE FIRST WORD

PICK THE B (EX 3) BEST ALTERNATIVES
(IN, JANE, SEPTEMBER)

2. FOR EACH B WORDS PICK THE NEXT WORD AND EVALUATE THE PAIRS TO END UP w B PAIRS

$$P(y^{<1>} | y^{<2>} | x) = P(y^{<1>} | x) P(y^{<2>} | x, y^{<1>})$$



(IN SEPTEMBER, JANE IS, JANE VISITS)

3. REPEAT TIL DONE

$$\text{ARG MAX } \prod_{t=1}^{Ty} P(y^{<t>} | x, y^{<1>} | \dots | y^{<t-1>})$$

OVERFLOWS

PROBLEM: MULTIPLYING PROBABILITIES ($0 < p \ll 1$) RESULTS IN A VERY SMALL NUMBER

PROBLEM II: IF WE OPTIMIZE FOR THE MULT WE WILL PREFER SHORT SENTENCES. SINCE EACH WORD WILL REDUCE PROB

INSTEAD WE CAN OPTIMIZE FOR THIS

$$\frac{1}{Ty} \alpha \sum_{t=1}^{Ty} \log(P^{<t>} | x, y^{<1>} | \dots | y^{<t-1>})$$

HOW DO WE PICK B?

LARGE B: BETTER RESULT, SLOWER
SMALL B: WORSE RESULT, BETTER

IN PROD YOU MIGHT SEE B=10.
100 IS PROBABLY A BIT TOO HIGH -
BUT ITS DOMAIN DEPENDENT

ERROR ANALYSIS IN BEAM S.

HUMAN: JANE VISITS AFRICA IN SEPT... y^*
ALSO: JANE VISITED AFRICA LAST SEPTEMBER \hat{y}

HOW DO WE KNOW IF ITS OUR RNN OR OUR BEAM SEARCH WE SHOULD WORK ON?

$$\text{LET THE RNN GIVE } P_y^* = P(y^* | x) \text{ & } P_{\hat{y}} = P(\hat{y} | x)$$

IF $P_y^* > P_{\hat{y}}$:

BEAM PICKED THE WRONG ONE
TRY A HIGHER B

ELSE:

THE RNN PICKED THE WRONG PROBS - SO FOCUS ON THE RNN

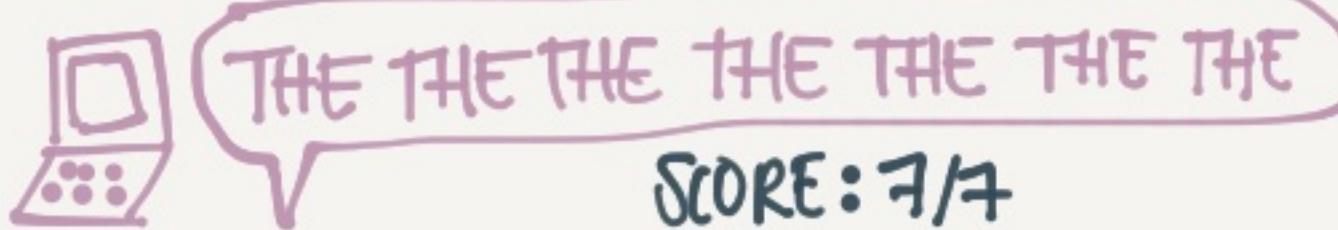
SEQUENCE TO SEQUENCE

FRENCH: LE CHAT EST SUR LE TAPIS
 HUMAN1: THE CAT IS ON THE MAT
HUMAN2: THERE IS A CAT ON THE MAT

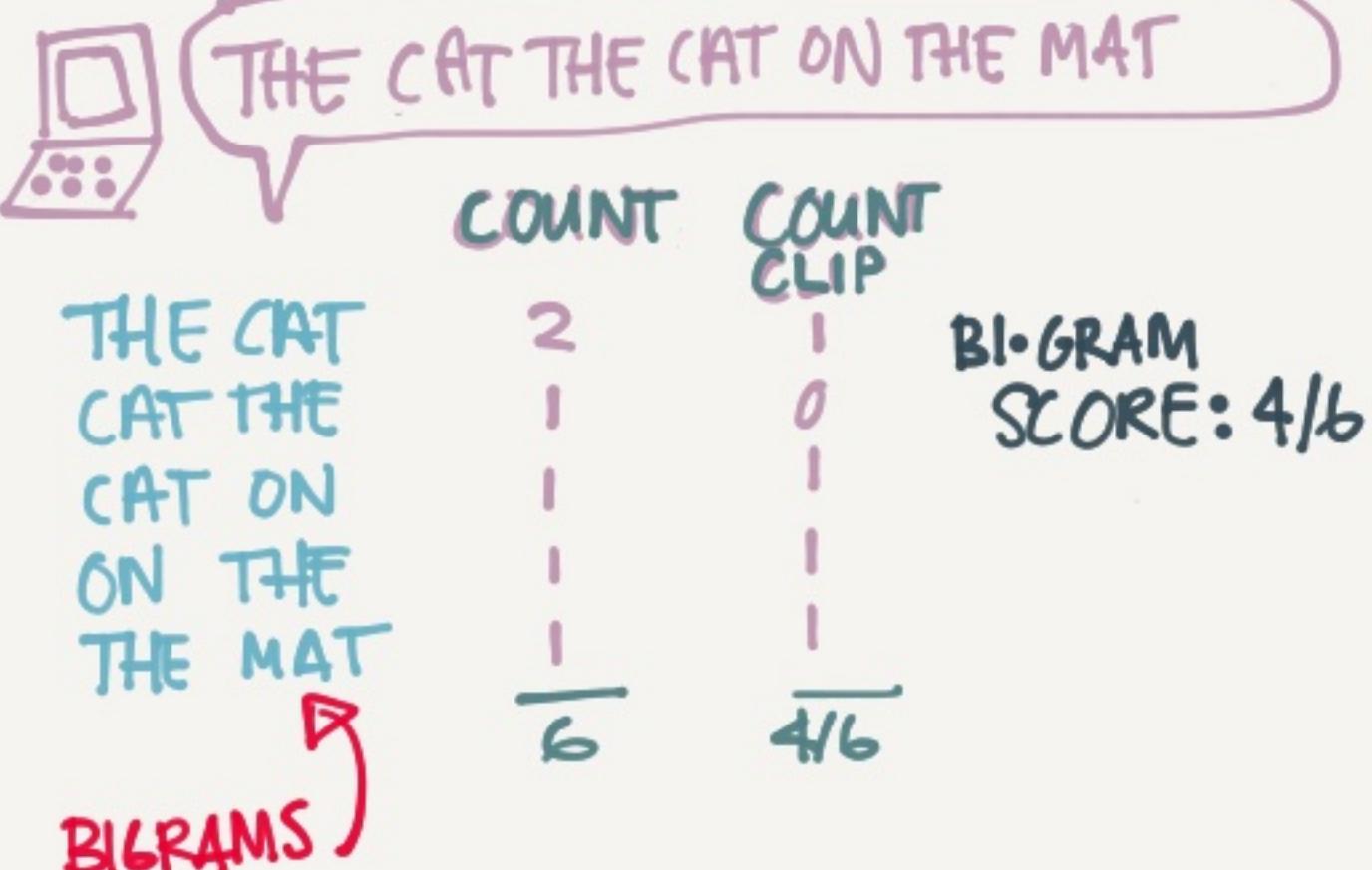
How do you evaluate the machine translation when multiple are right?

BLEU SCORE

IDEA: CHECK IF THE WORDS ^{MR} APPEAR IN THE REAL TRANSLATION



IDEA: ONLY GIVE CREDIT FOR A WORD THE MAX # TIMES IT APPEARS IN A TARGET SENTENCE
 SCORE: 2/7 ^{COUNT CLIP}



COMBINED BLEU SCORE

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

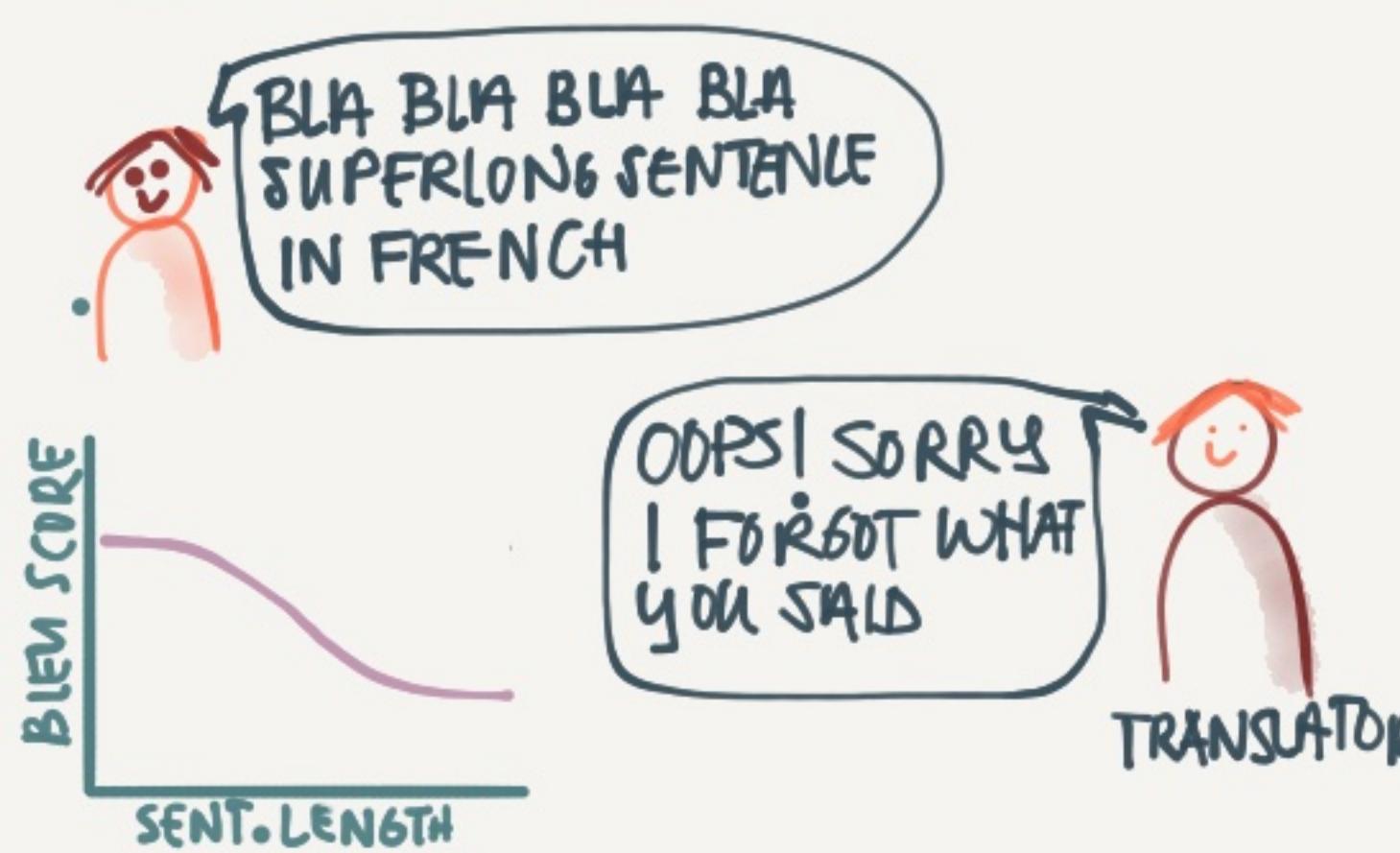
P_1 = SCORE SINGLE WORD
 P_2 = SCORE BIGRAMS

...
 BP = BREVITY PENALTY

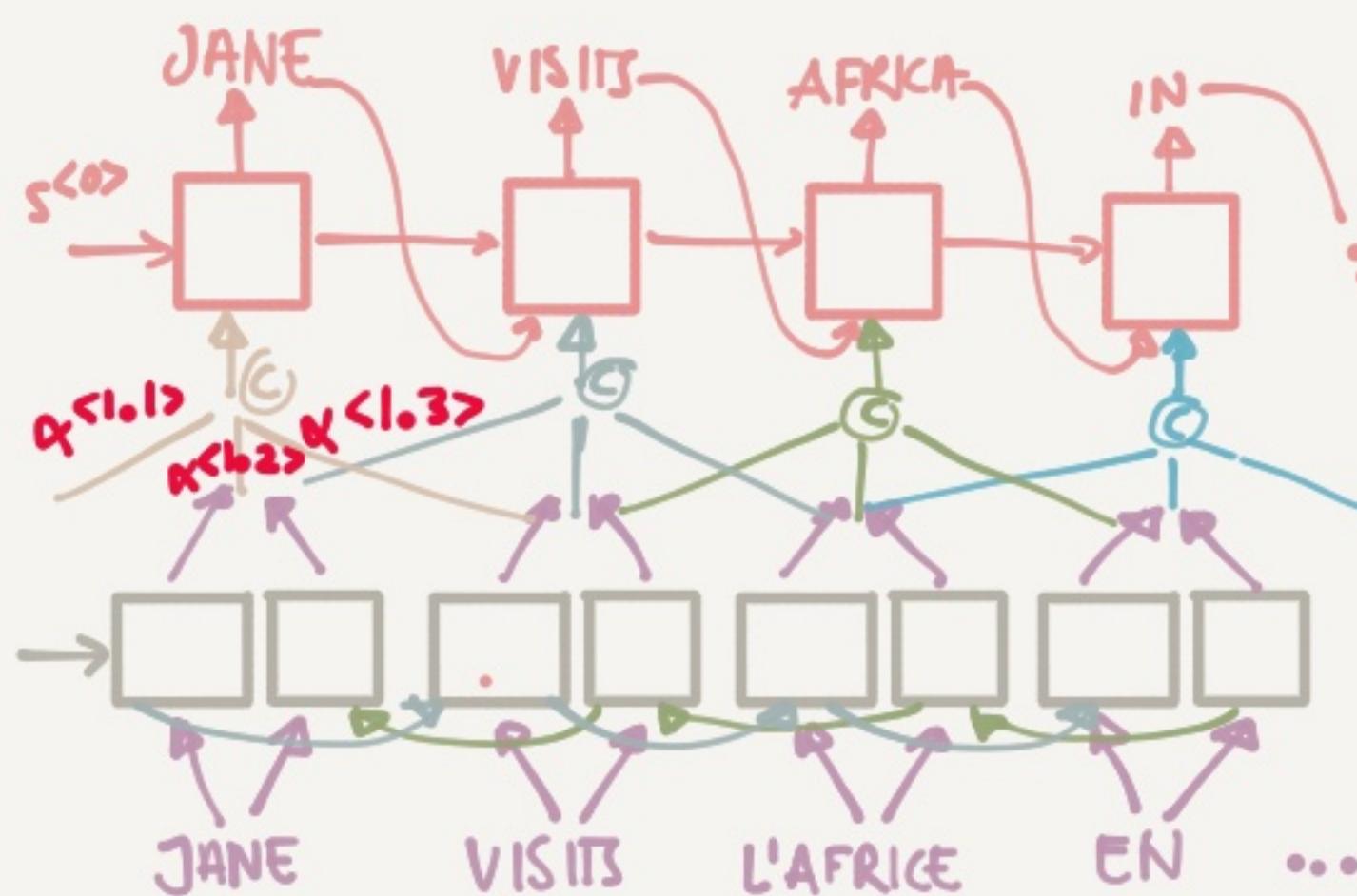
PENALIZES
SENTENCES
SHORTER
THAN THE
TARGET

A USEFUL SINGLE NUMBER
EVAL METRIC

ATTENTION MODEL



SOLUTION TRANSLATE A LITTLE AT A TIME USING ONLY PARTS OF THE SENTENCE AS CONTEXT



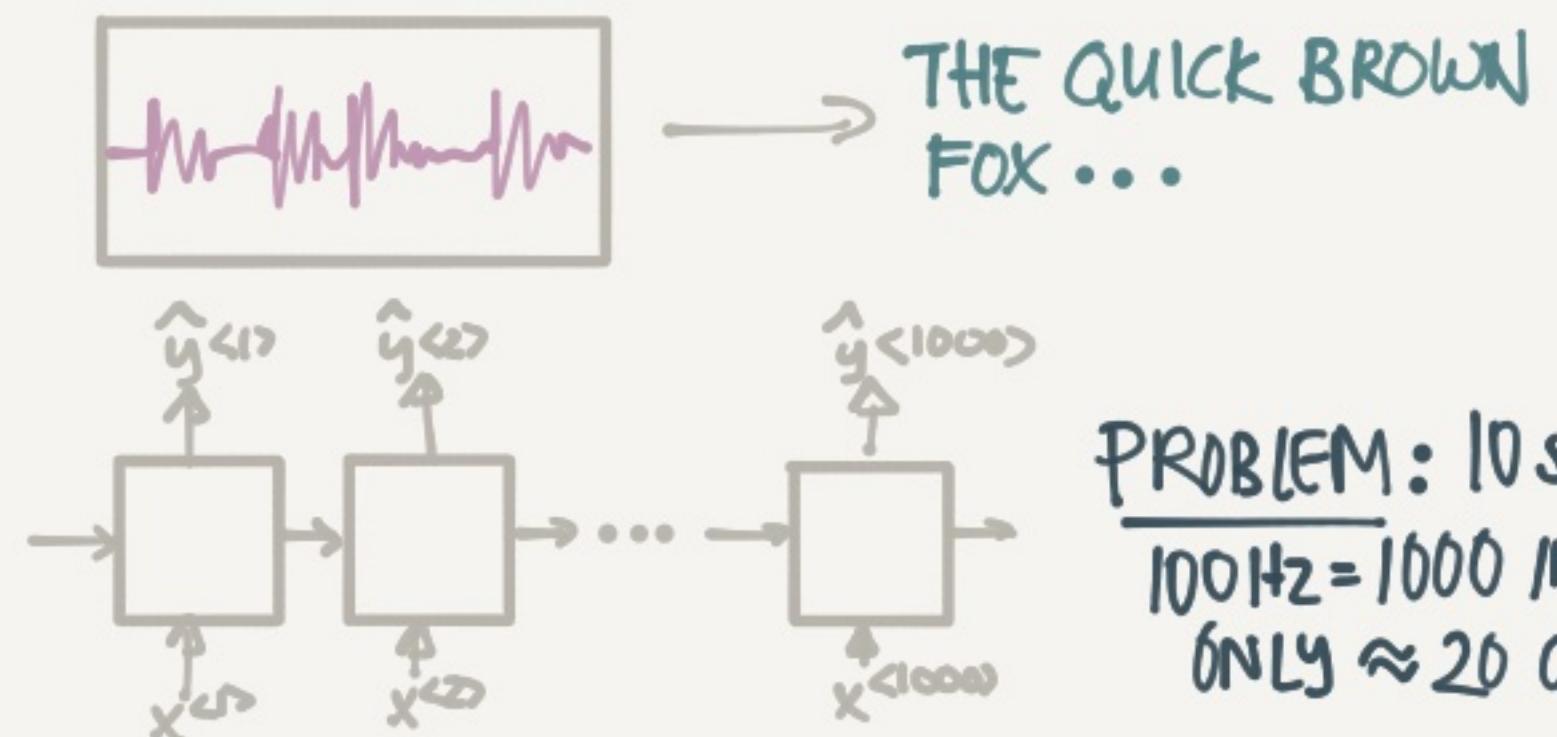
$$\alpha^{<t,t>} = \text{HOW MUCH ATTENTION } y^{<t>} \text{ SHOULD PAY TO } a^{<t>}$$

$$C^{<2>} = \sum_{t'} \alpha^{<2,t>} \cdot a^{<t>} \quad \sum_t \alpha^{<1,t>} = 1$$

α IS CALCULATED USING A SMALL NEURAL NETWORK

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

SPEECH RECOGNITION



SOLUTION USE CTC COST (CONNECTION TEMPORAL CLASSIFICATION)

t t - h _ e e e - - - l u - - - q q q - - - o

COLLAPSE REPEATED CHARS NOT SEP BY BLANK

TRIGGER WORD DETECTION

