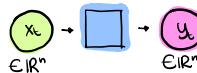
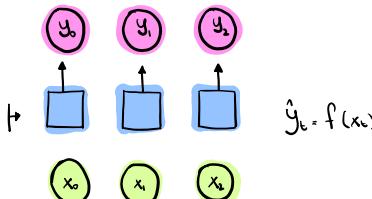


THE PERCEPTRON RUSHED

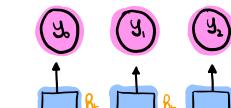


WE CAN REPLICATE THE SAME PROCESS FOR EACH TIME STEP



WE NEED A WAY TO RELATE THE NETWORK COMPUTATION AT A PARTICULAR TIME TO ITS HISTORY!

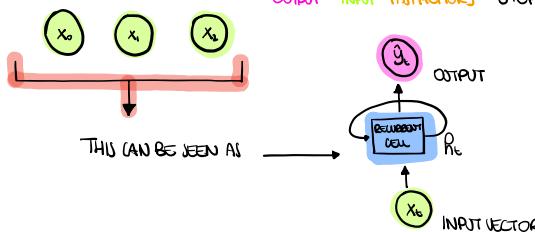
WE CAN LINK THE INFORMATION AND THE COMPUTATION OF THE NETWORK AT DIFFERENT TIME STEPS TO EACH OTHER!



$$\hat{y}_t = f(x_t, h_{t-1})$$

OUTPUT INPUT PAST MEMORY STEP.

NOW THE OUTPUT IS DEPENDENT NOT ONLY ON THE INPUT AT A PARTICULAR TIME STEP BUT ALSO ON THE STATE OF THE PREVIOUS



RNNs → RECURRENT NEURAL NETWORKS

↳ RNNs HAVE A STATE h_t , THAT IS UPDATED AT EACH Timestep AS A SEQUENCE IS PROCESSED

THIS IS DONE BY APPLY A RECURRENT RELATION AT EVERY TIME STEP TO PROCESS A SEQUENCE

$$h_t = f_w(x_t, h_{t-1})$$

CELL STATE WITH WEIGHTS
INPUT
OLD STATE
 w

KEY POINT → THE SAME FUNCTION AND SET OF PARAMETERS ARE USED AT EVERY TIME STEP.

IN PSEUDO CODE:

`my_rnn = RNN()``hidden_state = [0, 0, 0, 0] # THE FIRST HIDDEN STATE``Sentence = ["I", "love", "nebully"]``for word in sentence: # THE RECURRENT RELATION IS CAPTURED BY THIS LOOP``prediction, hidden_state = my_rnn(word, hidden_state)``next_word_prediction = prediction # TOKEN PREDICTION`THE INTERNAL COMPUTATION OF THE RNN INCLUDES BOTH THIS INTERNAL STATE OUTPUT AS WELL AS ULTIMATELY TRYING TO OUTPUT THE PREDICTION WE ARE INTERESTED IN (THE OUTPUT VECTOR \hat{y}_t)GIVEN AN INPUT VECTOR x_t WE PASS THAT IN TO COMPUTE THE RNN INTERNAL STATE COMPUTATION ↴

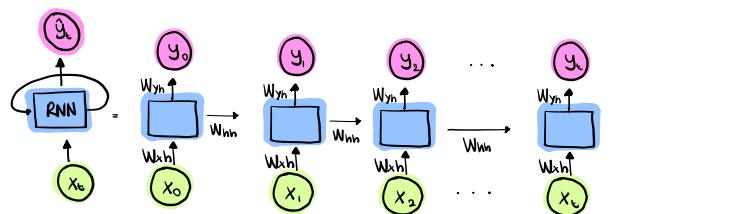
$$\rightarrow h_t = \text{tanh}(W_{hh}^T h_{t-1} + W_{hx}^T x_t) \text{ UPDATE HIDDEN STATE}$$

BREAKING THIS FUNCTION DOWN, WHAT IT'S DOING IT'S JUST A STANDARD NN OPERATION. IN THIS CASE WE MULTIPLY BOTH THE PAST HIDDEN STATE BY A WEIGHT MATRIX W AS WELL AS THE CURRENT INPUT x_t BY ANOTHER WEIGHT MATRIX, AND THEN WE ADD THEM TOGETHER AND APPLY A NON-LINEARITY.

$$\hat{y}_t = W_{hy}^T h_t$$

↳ THE FINAL STEP IS TO GENERATE THE OUTPUT FOR A GIVEN TIME STEP. THIS IS DONE BY TAKING THAT INTERNAL STATE AND SIMPLY MODIFYING IT, FOLLOWING A MULTIPLICATION BY ANOTHER WEIGHT MATRIX AND THEN USING THIS AS OUR GENERATED OUTPUT.

RNN REPRESENTED AS A COMPUTATIONAL GRAPH UNROLLED ACROSS TIME.



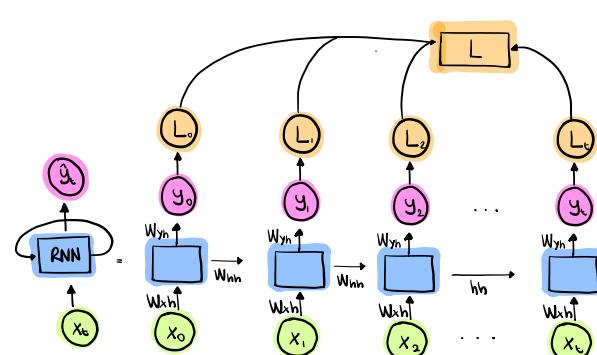
WE RE-USE THE SAME WEIGHT MATRICES AT EVERY TIME STEP.

BUT HOW DO WE TRAIN THIS NETWORK? WE NEED A LOSS!!

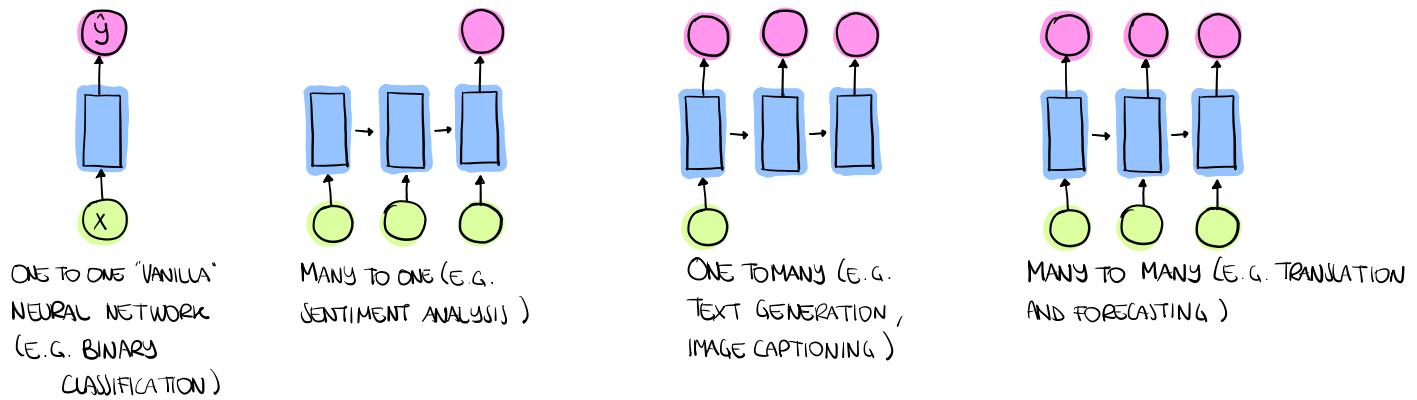
WE CAN GENERATE AN INDIVIDUAL LOSS FOR EACH OF THESE INDIVIDUAL TIME STEPS

ACCORDING TO WHAT THE OUTPUT OF THE TIME STEP IS. AND WE CAN GENERATE A

TOTAL SUM LOSS BY TAKING THESE STEPS AND SUM THEM TOGETHER.



RNNs FOR SEQUENCE MODELING



SEQUENCE MODELING - DESIGN CRITERIA

WE NEED TO:

- 1) HANDLE VARIABLE-LENGTH SEQUENCES → FEEDFORWARD NEURAL NETWORKS ARE NOT ABLE TO HANDLE VARIABLE LENGTH BECAUSE THEY HAVE INPUT OF FIXED DIMENSIONALITY. BUT BECAUSE WITH RNNs WE ARE UNROLLING ACROSS TIME WE ARE ABLE TO HANDLE THIS.
- 2) TRACK LONG-TERM DEPENDENCIES → WE NEED INFORMATION FROM THE DISTANT PAST TO ACCURATELY PREDICT THE NEXT WORD.
- 3) MAINTAIN INFORMATION ABOUT ORDER
- 4) SHARE PARAMETERS ACROSS THE SEQUENCE TO BE ABLE TO KEEP TRACK OF THESE DEPENDENCIES

REPRESENTING THE LANGUAGE TO A NN

NEURAL NETWORKS REQUIRE NUMERICAL INPUTS, WE NEED A WAY TO REPRESENT THE WORDS NUMERICALLY → EMBEDDING → TRANSFORM INDEXES INTO A VECTOR OF FIXED SIZE

① VOCABULARY

CORPUS OF WORDS

WALK
AM MY
FOR THIS I
TOOK CAT
MORNING

② INDEXING

WORD TO INDEX

A → 1
CAT → 2
...
WALK → N

③ EMBEDDING

THESE INDEXES CAN BE MAPPED TO AN EMBEDDING WHICH IS JUST A FIXED LENGTH VECTOR.

ONE-HOT EMBEDDING
"CAT" = [0, 1, 0, 0, ..., 0]
↑
i-th INDEX

IDEA: TAKING AN INPUT OF WORDS THAT ARE GOING TO BE INDEXED NUMERICALLY, WE CAN LEARN AN EMBEDDING OF THOSE WORDS IN ONE LOWER DIMENSIONAL SPACE.

USING MACHINE LEARNING WE CAN MAP THE

LEARNED EMBEDDING
RUN WALK
DAY SUN
HAPPY SAD
DOG CAT

ANOTHER OPTION TO GENERATE AN EMBEDDING IS TO USE SOME SORT OF MACHINE LEARNING MODEL TO LEARN AN EMBEDDING

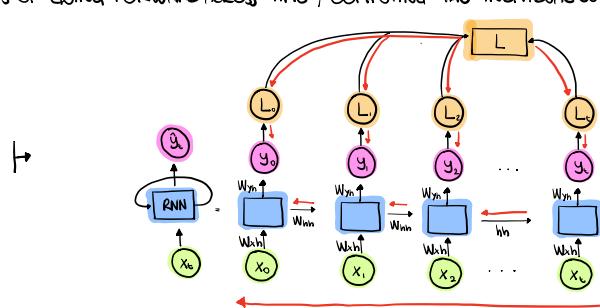
MEANING OF THE WORDS TO AN ENCODING THAT IS MORE INFORMATIVE → WORDS THAT HAVE SIMILAR MEANING WILL HAVE SIMILAR EMBEDDINGS.

BACKPROPAGATION THROUGH TIME (BPTT)

IN THE RNNs THE FORWARD PASS THROUGH THE NETWORK CONSISTS OF GOING FORWARD ACROSS TIME, COMPUTING THE INDIVIDUAL LOSS VALUES AT THE INDIVIDUAL TIME STEPS AND THEN SUMMING THEM TOGETHER.

↓

TO BACKPROPAGATE, INSTEAD OF BACKPROPAGATING ERRORS THROUGH A SINGLE FEEDFORWARD NETWORK, WE HAVE TO BACKPROPAGATE ERRORS INDIVIDUALLY ACROSS EACH TIME STEP AND ACROSS EACH TIME STEP.



BETWEEN EACH TIME STEP WE NEED TO PERFORM THESE INDIVIDUAL MATRIX MULTIPLICATION → COMPUTING THE GRADIENT WITH RESPECT TO θ INVOLVES MANY FACTORS OF $W \cdot b$, AND REPEATED GRADIENT COMPUTATION. WHY THIS CAN BE PROBLEMATIC:

- MANY VALUES (WEIGHTS OR GRADIENTS) $> 1 \rightarrow$ EXPLODING GRADIENT, THE GRADIENTS ARE BECOMING EXTREMELY LARGE DUE TO THIS REPEATED MULTIPLICATION OPERATION. POSSIBLE SOLUTION: GRADIENT CLIPPING / SCALE BIG GRADIENTS.
- MANY VALUES $< 1 \rightarrow$ VANISHING GRADIENT: MULTIPLY MANY SMALL VALUES TOGETHER → ERRORS DUE TO FURTHER BACK TIME STEPS HAVE SHOWER AND FAUER GRADIENTS → BIAS PARAMETERS TO CAPTURE SHORT TERM DEPENDENCIES

POSSIBLE SOLUTIONS TO THE VANISHING GRADIENT PROBLEM:

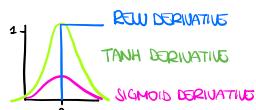


THE MODEL WILL IGNORE THE LONG TERM DEPENDENCIES THAT MAY EXIST.

T

① ACTIVATION FUNCTIONS

USING RELU PREVENTS f' FROM SHRINKING THE GRADIENTS WHEN $x > 0$



② PARAMETER INITIALIZATION

INITIALIZE WEIGHTS TO IDENTITY MATRIX AND INITIALIZE BIASES TO ZERO. THIS HELPS PREVENT THE WEIGHTS SHRINKING TO ZERO.

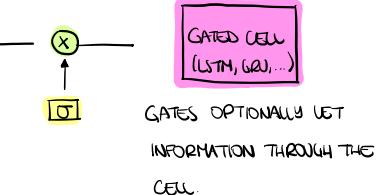
$$I_n = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 \end{pmatrix}$$

POINTWISE
MULTIPLICATION
SIGMOID LAYER

③ GATED CELL

USE GATES TO SELECTIVELY ADD OR REMOVE INFORMATION WITHIN EACH RECURRENT UNIT

WITH:



LSTMS (LONG SHORT TERM MEMORY)

NETWORK RELY ON GATED CELL TO TRACK INFORMATION THROUGH OUT MANY TIME STEPS. GATED LSTM CELLS CONTROL INFORMATION FLOW:

FORGET

FORGET GATE LETS RID OF IRRELEVANT INFORMATION

STORE

STORE RELEVANT INFORMATION FROM CURRENT INPUT

UPDATE

SELECTIVELY UPDATE CELL STATE

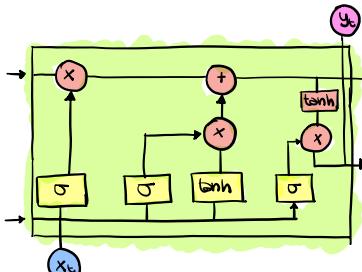
OUTPUT

OUTPUT GATE RETURNS A FILTERED VERSION OF THE CELL STATE

THANKS TO THIS GATED STRUCTURE THE BACK PROPAGATION THROUGH TIME ALGORITHM BECOMES MUCH MORE STABLE! THE VANISHING GRADIENT PROBLEM IS MITIGATED BY HAVING FEWER REPEATED MATRIX MULTIPLICATIONS.



BPTT WITH PARTIALLY UNINTERRUPTED GRADIENT FLOW



LSTM CELLS ARE ABLE TO TRACK INFORMATION THROUGHOUT MANY TIME STEPS

LIMITATIONS OF RNNs

① ENCODING BOTTLENECK

WE TAKE A LOT OF CONTENT (E.G. A VERY LONG BODY OF TEXT) AND CONDENSE IT INTO A REPRESENTATION THAT CAN BE PREDICTED ON.



INFORMATION CAN BE LOST IN THIS ENCODING OPERATION

② SLOW, NO PARALLELIZATION

RECURRENT NEURONS AND RECURRENT MODELS ARE NOT EFFICIENT AT ALL.



THEY REQUIRE TO PROCESS INFORMATION SEQUENTIALLY → THEY ARE VERY INEFFICIENT ON MODERN GPU HARDWARE

③ NOT LONG MEMORY

RECURRENT MODELS DON'T HAVE AN HIGH MEMORY CAPACITY, WHILE THEY CAN HANDLE SEQUENCES IN THE ORDER OF 100 OR 1000 TIME STEPS, THEY DON'T REALLY SCALE WELL TO SEQUENCES THAT ARE OF LENGTH OF THOUSANDS OR TEN THOUSANDS OF TIME STEPS.

THE MAIN LIMITATION OF RNN'S IS THAT THEY PROCESS THOSE TIME STEPS INDIVIDUALLY DUE TO THE RECURRENT RELATION.

WE WANT TO ELIMINATE THE RECURRENT OPERATION. WE CAN'T SIMPLY TAKE THE INPUT, SQUASH IT AND PAVE IT INTO A DENIC NETWORK BECAUSE:

• IT'S NOT SCALABLE FOR VERY LONG SEQUENCES

• THERE IS NO ORDER

• NO LONG MEMORY

• THERE IS NO NOTION OF WHAT POINT IN OUR SEQUENCE IS IMPORTANT.

ATTENTION IS ALL YOU NEED ← MECHANISM BEHIND TRANSFORMER

IDEA: IDENTIFY AND ATTEND TO WHAT'S IMPORTANT



SUPER-ATTENTION → THE ABILITY TO TAKE AN INPUT AND ATTEND TO THE MOST IMPORTANT

PARTS OF THAT INPUT: 1) IDENTIFY WHICH PARTS TO ATTEND TO

2) EXTRACT THE FEATURES WITH HIGH ATTENTION

UNDERSTANDING ATTENTION WITH SEARCH

WHEN WE SEARCH SOMETHING ON THE INTERNET, FOR EXAMPLE A VIDEO ON YOUTUBE WE INPUT SOME **QUERY (Q)** INTO THE SEARCH BAR. FOR EVERY VIDEO IN THE DB WE ARE GOING TO EXTRACT SOME KEY INFORMATION WHICH WE CALL THE **KEY (K)**. THE OVERLAPS BETWEEN THE QUERY AND THE KEYS IN THE DATABASE ARE GOING TO BE COMPUTED, AND AS WE DO THIS, AT EACH CHECK WE MAKE WE'LL ASK HOW SIMILAR IS THAT KEY TO OUR QUERY.

THIS IS THE IDEA OF COMPUTING AN ATTENTION MASK.

THE NEXT AND FINAL STEP IS TO EXTRACT THE INFORMATION THAT WE CARE ABOUT BASED ON THIS COMPUTATION, AND WE CALL THIS THE **VALUE (V)**.

LEARNING SELF-ATTENTION WITH NEURAL NETWORKS

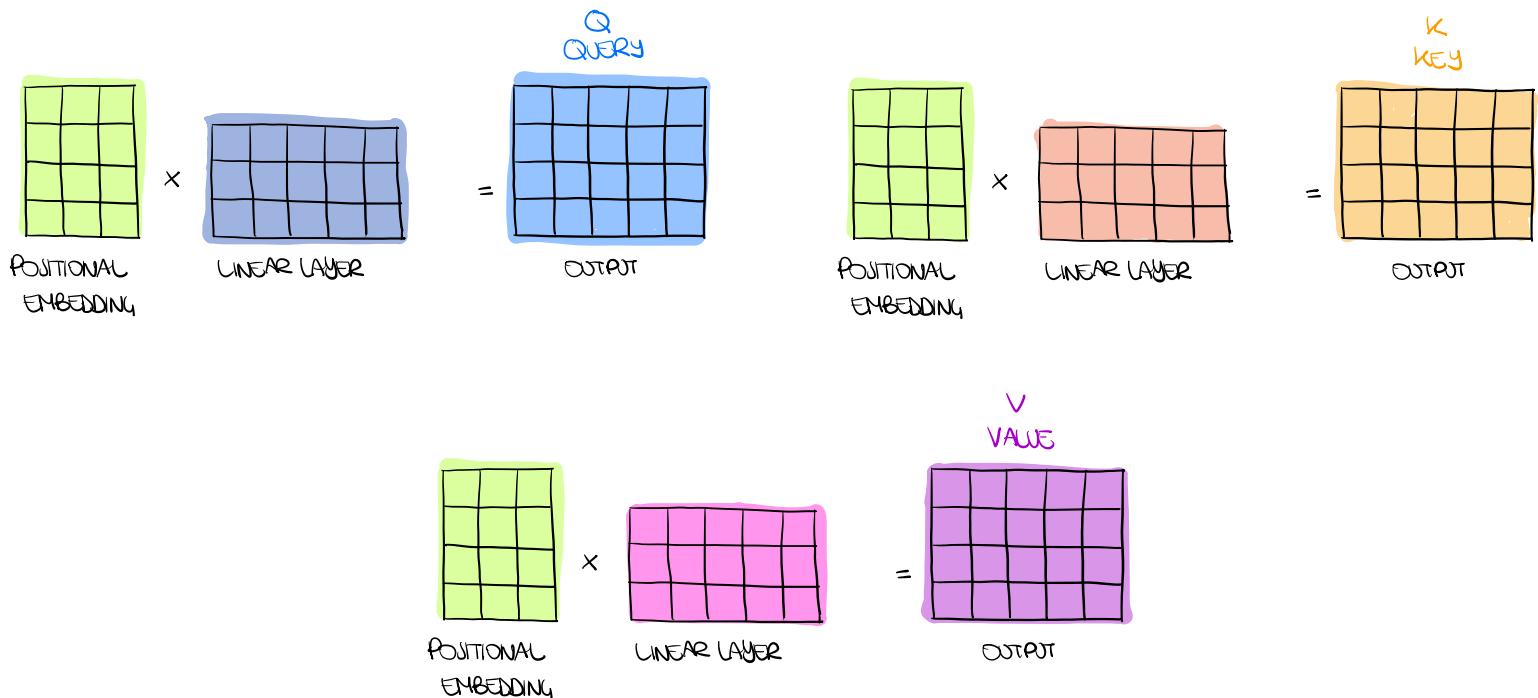
WE DON'T WANT TO PROCESS THE INFORMATION TIME STEP BY TIME STEP \Rightarrow WE NEED A WAY TO ENCODE THE POSITION INFORMATION TO UNDERSTAND ORDER.

↳ DATA IS FED ALL AT ONCE!

WE USE AN EMBEDDING THAT IS GOING TO INCORPORATE SOME NOTION OF POSITION.

↳ COMPUTE THE WORD EMBEDDING, TAKE SOME METRICS THAT CAPTURES POSITION INFORMATION WITHIN THAT SEQUENCE AND THEN COMBINE THEM TOGETHER.

AFTER WE COMPUTE THE POSITIONAL EMBEDDING, WE NEED TO FIGURE OUT HOW TO EXTRACT THE QUERY, KEY AND VALUE.



THE NEXT STEP IS TO COMPUTE THE ATTENTION USING THESE THREE MATRICES, FIGURING OUT HOW MUCH ATTENTION TO PAY TO AND WHERE

↳ WE FOCUS ON THE **SIMILARITY** BETWEEN OUR QUERY AND KEY.

↓

ATTENTION SCORES \leftarrow COMPUTE PAIRWISE SIMILARITY BETWEEN EACH QUERY AND VALUE: $\frac{QK^T}{\text{SCALING}}$



THE SIMILARITY METRIC \leftarrow IS PASSED IN A SOFTMAX FUNCTION THAT SQUASH EVERYTHING BETWEEN 0 AND 1

THE FINAL STEP IS TO USE THE WEIGHTING MATRIX TO EXTRACT FEATURES WITH HIGH ATTENTION.

$$\text{SOFTMAX} \left(\frac{QK^T}{\text{SCALING}} \right) V = A(Q, K, V)$$

WE CAN DO THIS MULTIPLE TIMES \Rightarrow WE CAN HAVE **MULTIPLE INDIVIDUAL ATTENTION HEADS** FOR ATTENDING TO DIFFERENT THINGS.