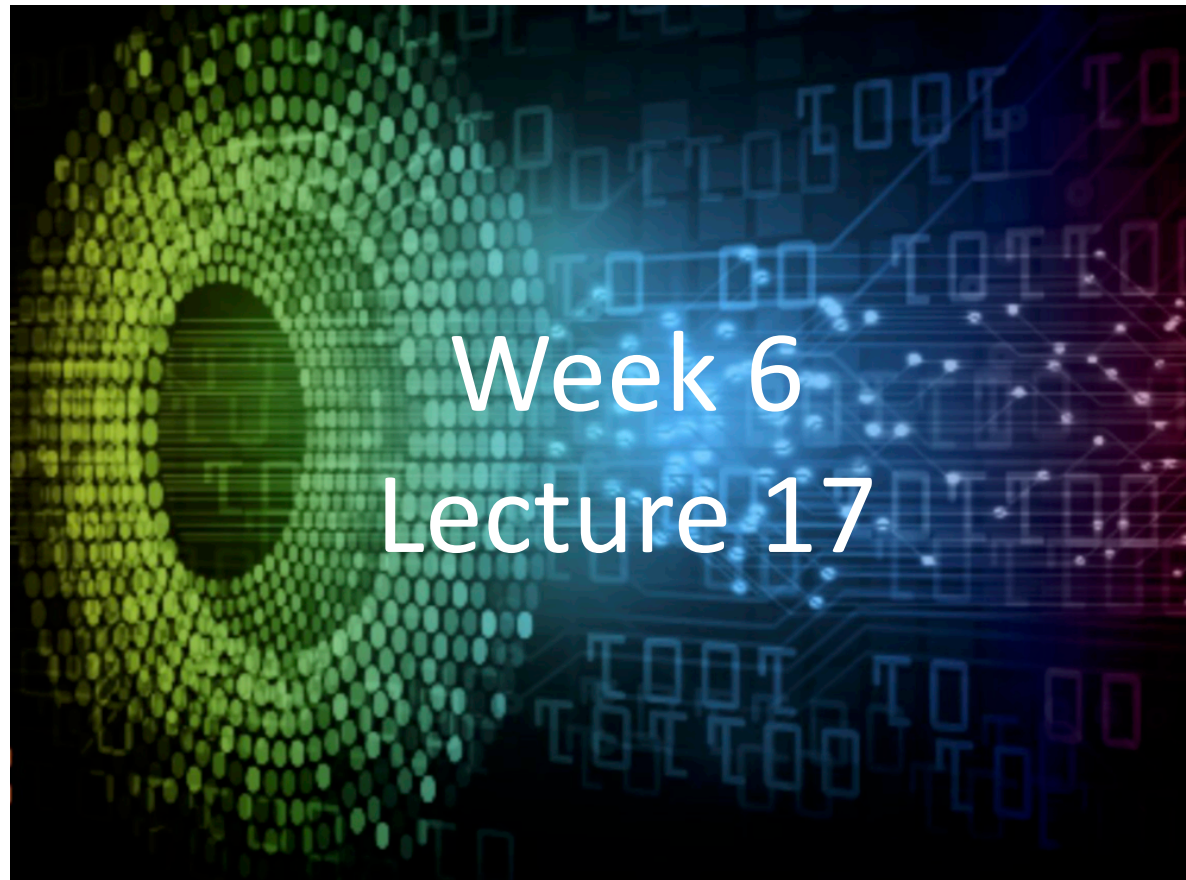# Introduction to Deep Learning Applications and Theory



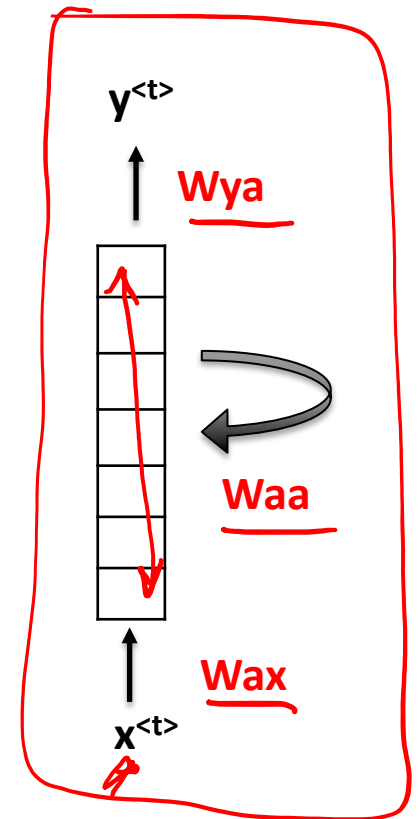Week 6
Lecture 17

# ECE 596 / AMATH 563

# Previous Lecture:
# Recurrent Neural Networks (RNNs)

- Input-Output

- Definition

- Notation
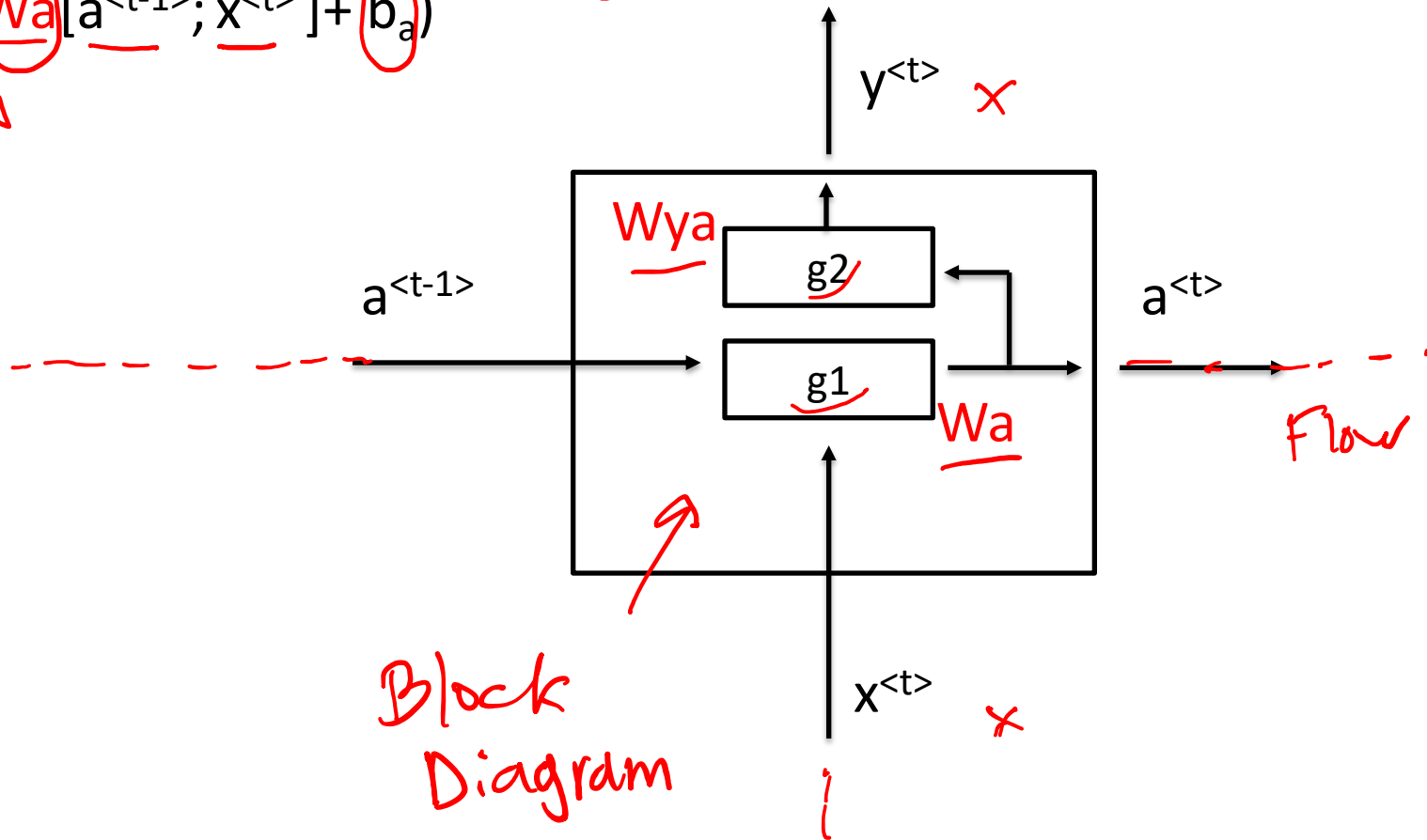
- Backward Propagation

- Diminishing/ Exploding Gradients

# RNN

$$y^{<t>} = g2(W_{ya} a^{<t>} + b_y)$$

$$a^{<t>} = g1(W_a [a^{<t-1>}; x^{<t>}] + b_a)$$

Classification

$g2 \sim tanh, \dots, \mathbb{G}$

$g1 \sim tanh$



Block Diagram

Flow

# Vanishing/Exploding Gradients

"I grew up in France and moved to the United States, therefore I speak _____"

$\mathbf{a^{<0>}}$

$\mathbf{x^{<1>}}$   $\mathbf{x^{<5>}}$

$\mathbf{x^{<14>}}$

$\mathbf{y^{<1>}}$

one-hot

$L(y, \hat{y}^{<1>})$

# Vanishing Gradients

$$a^{<t>} = \tanh(\text{Wa}[a^{<t-1>}; x^{<t>}] + b_a)$$
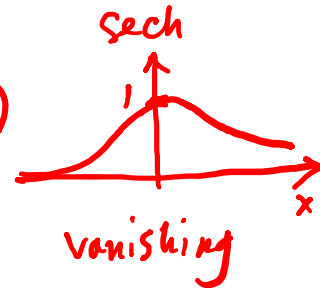
① > 1  blow up

< 1  vanish

$$\frac{\partial L}{\partial W_a} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a^{<Tx>}} \cdot \left( \prod_{t=2}^{Tx} \frac{\partial a^{<t>}}{\partial a^{<t-1>}} \right) \cdot \frac{\partial a^{<1>}}{\partial W_a}$$

B. P. T

②

sech

②

vanishing

①

$$\frac{\partial a^{<t>}}{\partial a^{<t-1>}} = \text{sech}^2 \left( W_a a^{<t-1>} + W_x x^{<t>} \right) \cdot W_a$$

② ①

$$W_a \left[ a^{<t-1>}, x^{<t>} \right]$$
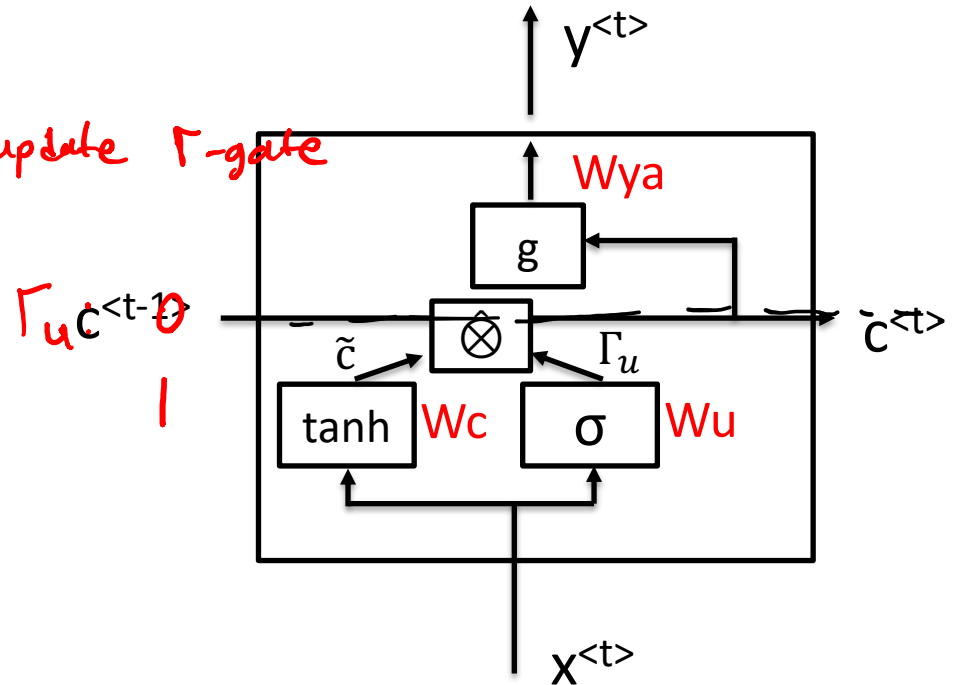
# Gated Recurrent Unit (GRU)

$$\tilde{c}^{<t>} = \tanh(Wc[c^{<t-1>}; x^{<t>}]+ b_c)$$

$$\Gamma_u = \sigma(Wu[c^{<t-1>}; x^{<t>}]+ b_u)$$

u-update  $\Gamma$-gate

$$c^{<t>} = \Gamma_u \tilde{c}^{<t>} + (1 - \Gamma_u)c^{<t-1>}$$

$$h^{<t>} == c^{<t>} == a^{<t>}$$

$\Gamma_u c^{<t-1>}$

$y^{<t>}$

Wya

g

$\tilde{c}$  $\otimes$  $\Gamma_u$

$\dot{c}^{<t>}$

tanh  Wc  $\sigma$  Wu

$x^{<t>}$

Cho, K., Van Merriënboer, B., Bahdanau, D., & Bengio, Y. (2014).
On the properties of neural machine translation: Encoder-decoder approaches.

Chung, Junyoung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. (2014).
Empirical evaluation of gated recurrent neural networks on sequence modeling.

# Vanishing Gradients

$$\frac{\partial L}{\partial W_a} = \frac{\partial L}{\partial \hat{y}} \ \frac{\partial \hat{y}}{\partial c^{<Tx>}} \cdot \left( \prod_{t=2}^{Tx} \frac{\partial c^{<t>}}{\partial c^{<t-1>}} \right) \cdot \frac{\partial c^{<1>}}{\partial W_c}$$

$a = c$

$$\frac{\partial c^{<t>}}{\partial c^{<t-1>}} = \Gamma_u{}' \tanh\left( W_c c^{<t-1>} + W_x x^{<t>} \right)$$

← ① $\Gamma'$

$$\Gamma_u \operatorname{sech}^2\left( W_c c^{<t-1>} + W_x x^{<t>} \right) \cdot W_c +$$

← ②

$$\Gamma_u{}' \ c^{<t-1>} +$$

← ③ ← $\Gamma'$

$$(1 - \Gamma_u)$$

④

# Full GRU

$$\tilde{c}^{<t>} = \tanh(Wc[\Gamma_r c^{<t-1>}; x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(Wu[c^{<t-1>}; x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(Wr[c^{<t-1>}; x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u \tilde{c}^{<t>} + (1 - \Gamma_u)c^{<t-1>}$$

# Long Short Term Memory

$$\tilde{c}^{<t>} = \tanh(Wc[a^{<t-1>}; x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(Wu[a^{<t-1>}; x^{<t>}] + b_u)$$

$$\Gamma_f = \sigma(Wf[a^{<t-1>}; x^{<t>}] + b_f)$$

$$\Gamma_o = \sigma(Wo[a^{<t-1>}; x^{<t>}] + b_0)$$

$$c^{<t>} = \Gamma_u \tilde{c}^{<t>} + \Gamma_f c^{<t-1>}$$

$$a^{<t>} = \Gamma_o \tanh c^{<t>}$$

Sepp Hochreiter

Jurgen Schmidhuber

Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural computation, 9*(8), 1735-1780.

# LSTM

$$\tilde{c}^{<t>} = \tanh(Wc[a^{<t-1>}; x^{<t>}] + b_c)$$
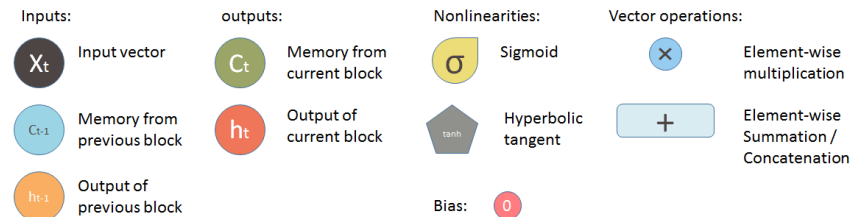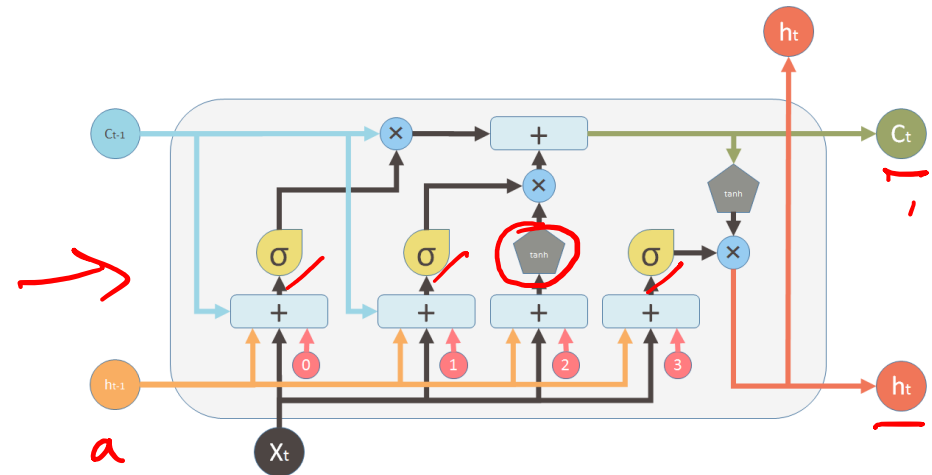
$$\Gamma_u = \sigma(Wu[a^{<t-1>}; x^{<t>}] + b_u)$$

$$\Gamma_f = \sigma(Wf[a^{<t-1>}; x^{<t>}] + b_f)$$

$$\Gamma_o = \sigma(Wo[a^{<t-1>}; x^{<t>}] + b_0)$$

$$c^{<t>} = \Gamma_u \tilde{c}^{<t>} + \Gamma_f c^{<t-1>}$$

$$a^{<t>} = \Gamma_o \tanh c^{<t>}$$



| Inputs: | | outputs: | | Nonlinearities: | | Vector operations: | |
|---|---|---|---|---|---|---|---|
| $X_t$ | Input vector | $C_t$ | Memory from current block | $\sigma$ | Sigmoid | $\times$ | Element-wise multiplication |
| $C_{t-1}$ | Memory from previous block | $h_t$ | Output of current block | tanh | Hyperbolic tangent | $+$ | Element-wise Summation / Concatenation |
| $h_{t-1}$ | Output of previous block | | | Bias: | 0 | | |

Hochreiter, S., & Schmidhuber, J. (1997).  Long short-term memory. *Neural computation*, *9*(8), 1735-1780.

# Bi-Directional RNN



Forward

$\hat{y}^{<1>}$   $\hat{y}^{<2>}$   $\hat{y}^{<3>}$

$a^{<1>}$   $a^{<2>}$

$x^{<1>}$   $x^{<2>}$   $x^{<Tx>}$

He said, Tesla is a unit of magnetic field strength

He said, Tesla is an electric automotive and sustainable energy company
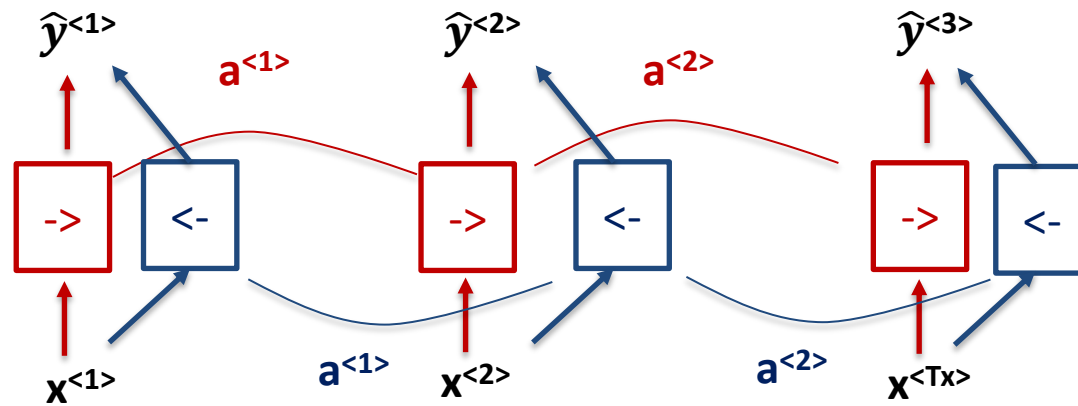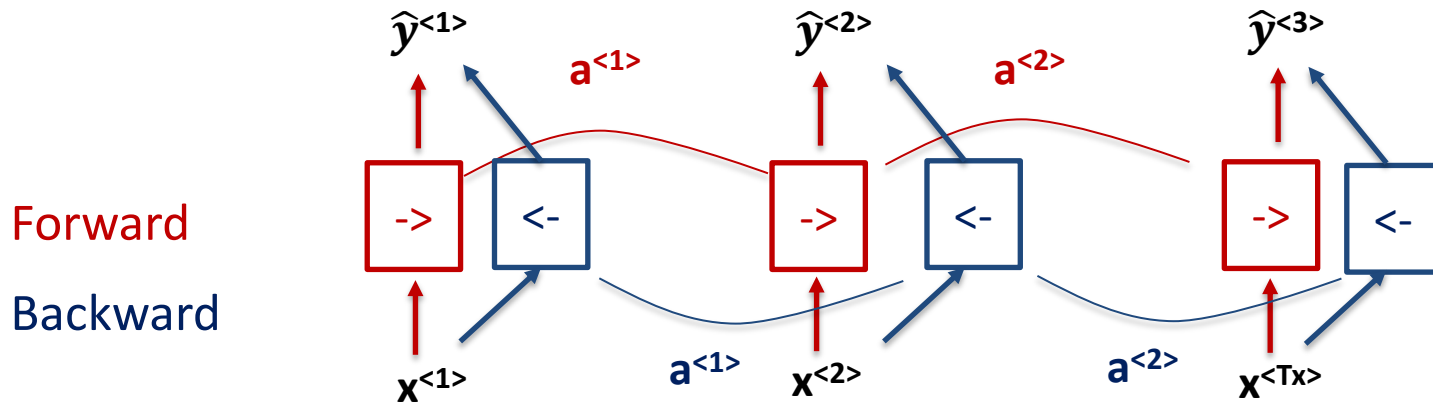
# Bi-Directional RNN

# Bi-Directional RNN



$$\hat{y}^{<t>} = g(Wy\,[a^{F<t>}; a^{B<t>}] + by)$$

# Deep RNN

# Deep RNN



$$a^{[2]<3>} = g($$
$$Wa^{[2]}\left[a^{[2]<2>}; a^{[1]<3>}\right] + b_a{}^{[2]})$$

# Natural Language Processing

## Vocabulary/Dictionary

| | |
|---|---|
| a | 1 |
| and | 2 |
| … | |
| Harry | 400 |
| … | |
| Jordan | |
| Michael | |
| … | |
| Potter | |
| … | |
| Table | |
| … | |
| | 10,000 |

$\xrightarrow{\text{one-hot}}$

Harry  Potter

| $x^{<1>}$ | $x^{<2>}$ | | |
|---|---|---|---|

$(O_{400})$

$\begin{matrix} 0 \\ 0 \\ … \\ 1 \\ … \\ 0 \\ 0 \\ … \\ … \\ 0 \end{matrix}$

$\begin{matrix} 0 \\ 0 \\ … \\ 0 \\ … \\ 0 \\ 1 \\ … \\ … \\ 0 \end{matrix}$ $(O_{7568})$

# Word Embeddings

I would really like to drive a Tesla _____ .

I would really like to drive a Porsche _____.

Embedding vector

| | **Man** | **Woman** | **King** | **Queen** | **Orange** | **Apple** | **Tesla** | **Porsche** |
|---|---|---|---|---|---|---|---|---|
| Gender | -1 | 1 | -0.95 | 0.97 | 0.01 | 0 | -0.5 | 0.03 |
| Royal | 0.01 | 0.02 | 0.94 | 0.95 | -0.02 | 0.04 | -0.01 | 0.1 |
| Food | | | | | | | | |
| Size | | | | | | | | |
| Engine | 0.07 | -0.01 | 0.03 | 0.02 | 0.001 | 0 | 0.95 | 0.98 |

Visualization 300 -> 2D: t-SNE

# Named Entity Recognition Example

$\hat{y}^{<1>} = P_1$   $\hat{y}^{<2>}$   $\hat{y}^{<END>}$

a^{<0>}

...

x^{<1>}   y^{<1>}   y^{<Ty>}

Sally   Johnson   is   a   Porsche   driver.

Magic   Jordan   is   a   Tesla   motorist.

# Transfer Learning

- Learn word embeddings from large text corpus (1-100B words). Or download.

- Transfer to new task with smaller training set (~100k).

- Continue to tune the embeddings with the new data.

# Properties of Embedding Vectors

|        | Man   | Woman | King  | Queen | Orange | Apple | Tesla | Porsche |
|--------|-------|-------|-------|-------|--------|-------|-------|---------|
| Gender | -1    | 1     | -0.95 | 0.97  | 0.01   | 0     | -0.5  | 0.03    |
| Royal  | 0.01  | 0.02  | 0.94  | 0.95  | -0.02  | 0.04  | -0.01 | 0.1     |
| Food   |       |       |       |       |        |       |       |         |
| Size   |       |       |       |       |        |       |       |         |
| Engine | 0.07  | -0.01 | 0.03  | 0.02  | 0.001  | 0     | 0.95  | 0.98    |

Can define distances (similarity):

$e_{Man} - e_{Woman} \sim= [-2;0;0;..;0]$

$e_{King} - e_{Queen} \sim= [-2;0;0;..;0]$

$e_{Man} - e_{Woman} \sim= e_{King} - e_?$

# Similarity

We can define similarity in the space of embedding vectors (Full space)

$$\text{argmax}_w \ \text{sim}(e_w, v)$$

# Embedding Matrix

E=

A and ... Harry ... Jordan Michael ... Potter ...Table ...

300

10000

$$\overrightarrow{e_j} = E \cdot \overrightarrow{o_j}$$