

LANGUAGE & AI: DEEP LEARNING



Chris Emmery
Department of Cognitive Science & AI
Tilburg University

[@cmry](https://twitter.com/_cmry) • [@_cmry](https://github.com/cmry) • [@cmry](https://github.com/cmry) • cmry.github.io



RECAP PREVIOUS LECTURES

- We looked at how language might be **noisy** as input.
- We discussed several way to **represent** and **model** language: count and prediction-based, and sequential.
- We looked at framing **prediction** tasks around those respective representations.

Today, we discuss Deep Learning: end-to-end representation learning and predictions in one black box.



RECURRENT MODELS





RELATION TO MATERIAL

- Representing word similarity (Weeks 1 and 4).
- Predicting outputs given input (Weeks 3 and 5).
- Sequential (or temporal) structure (Week 5).
- Semantic representations (Week 4)
- Context windows (Week 4).



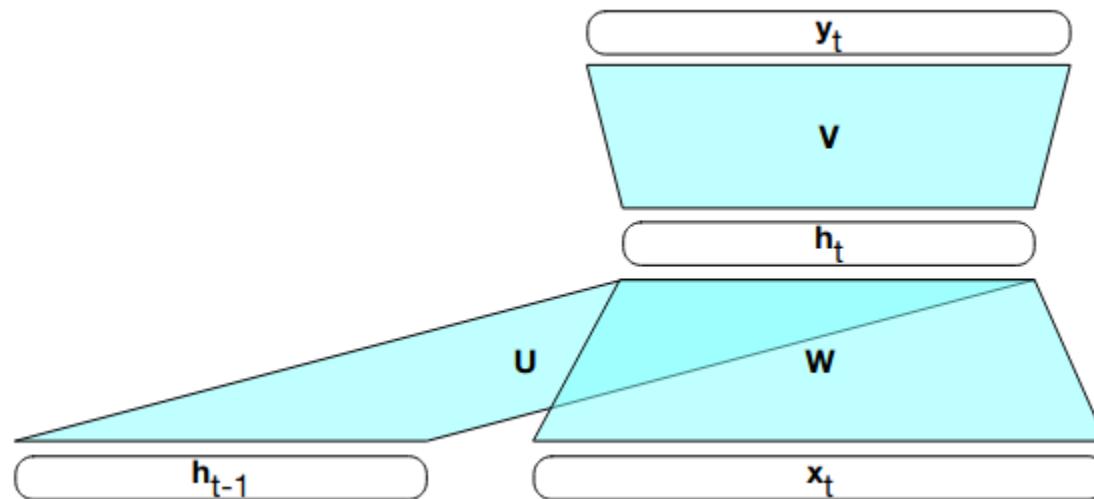
RECURRENT NEURAL NETWORKS (RNNs)

- Recurrent unit: encodes earlier steps in future decisions via 'lookback' for the hidden layers.



HOW DO THEY WORK?

Magic? No, just **U**.





FORMALLY

$$\mathbf{h} = \sigma(\mathbf{Wx} + \mathbf{b}) \quad \rightarrow$$

$$\mathbf{h}_t = g(\mathbf{Uh}_{t-1} + \mathbf{Wx}_t)$$

$$\mathbf{z} = \mathbf{Uh}$$

$$\hat{\mathbf{y}}_t = \text{softmax}(\mathbf{Vh}_t)$$

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{z})$$

- Dimensionality given d_{in} , d_h , and d_{out} :
- $\mathbf{W} \in \mathbb{R}^{d_h \times d_{in}}$,
- $\mathbf{U} \in \mathbb{R}^{d_h \times d_h}$, and
- $\mathbf{V} \in \mathbb{R}^{d_{out} \times d_h}$.

RNNS AS LMS

$$\mathbf{e} = [\mathbf{Ex}_{t-3}; \mathbf{Ex}_{t-2}; \mathbf{Ex}_{t-1}]$$

$$\mathbf{h} = \sigma(\mathbf{We} + \mathbf{b})$$

$$\mathbf{z} = \mathbf{Uh}$$

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{z})$$

$$\mathbf{e}_t = \mathbf{Ex}_t$$

$$\mathbf{h}_t = g(\mathbf{Uh}_{t-1} + \mathbf{We}_t)$$

$$\hat{\mathbf{y}}_t = \text{softmax}(\mathbf{Vh}_t)$$

Actual probabilities as dot product:

$$P(w_{t+1} = i \mid w_1, \dots, w_t) = \mathbf{y}_t[i]$$

$$P(w_{1:n}) = \prod_{i=1}^n \mathbf{y}_i [w_i]$$

Cross-entropy loss over predicted probabilities:

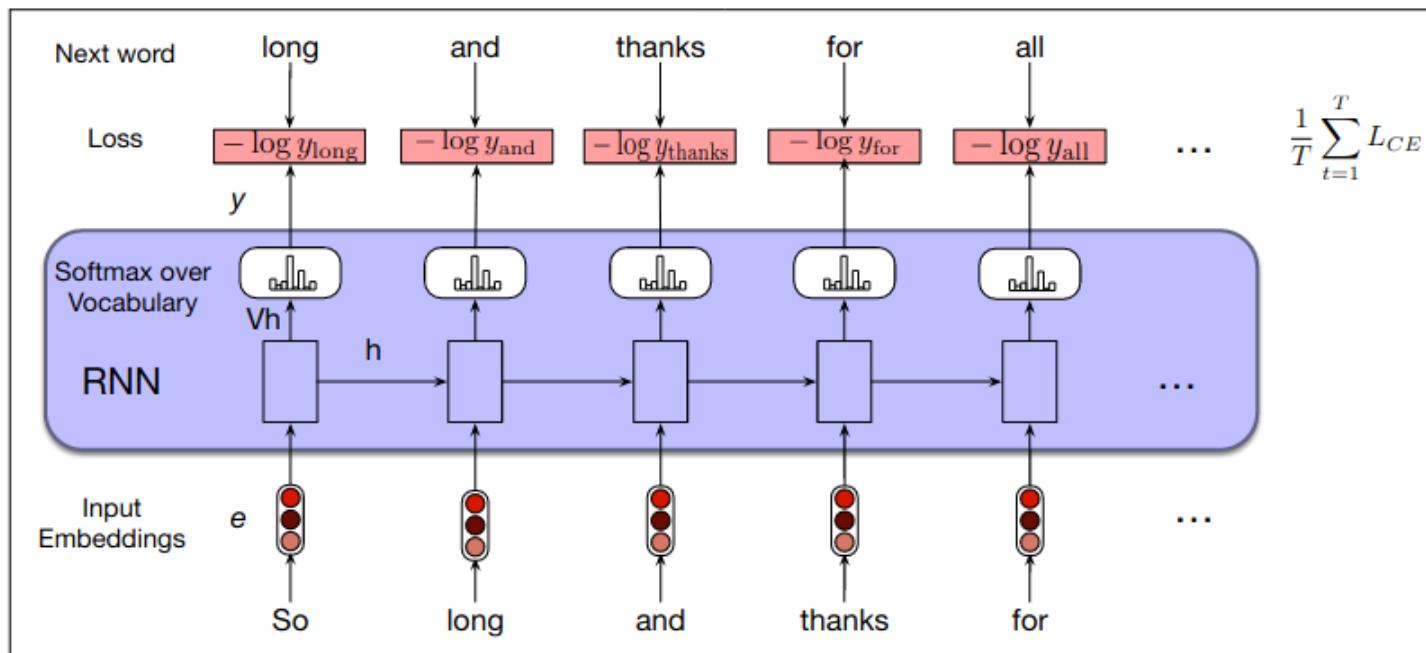
$$L_{CE} = - \sum_{w \in V} \mathbf{y}_t[w] \log \hat{\mathbf{y}}_t[w]$$

$$L_{CE} (\hat{\mathbf{y}}_t, \mathbf{y}_t) = - \log \hat{\mathbf{y}}_t [w_{t+1}]$$

LM Cross-entropy loss for = probability the model assigns to the correct next word (one-hot vector).

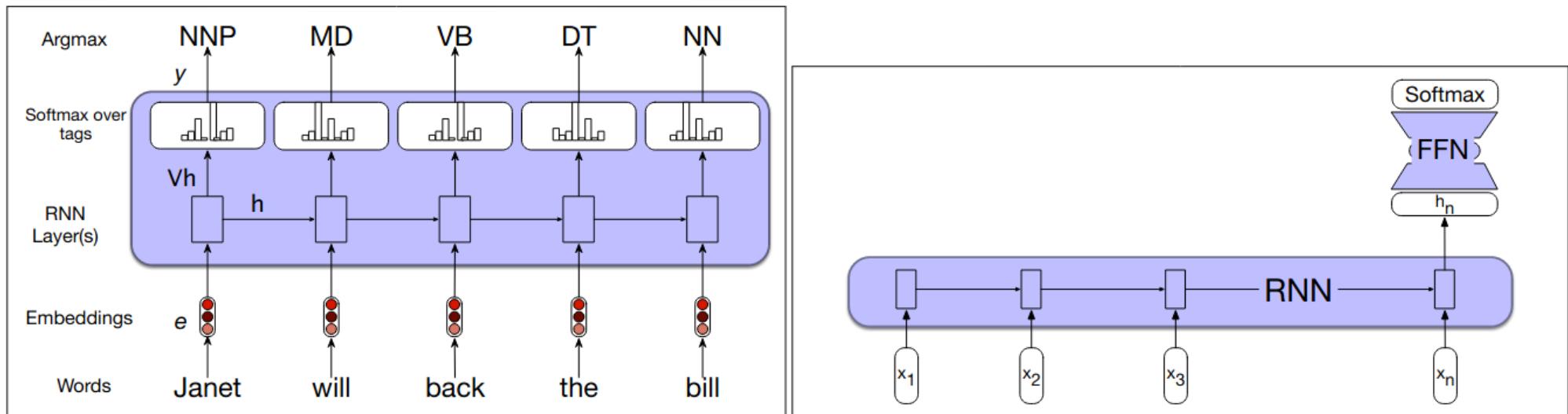


PUTTING IT TOGETHER





GENERALIZING RNN ARCHITECTURE TO OTHER TASKS





LET'S INCREASE THE COMPLEXITY!

- Like feed-forward layers: **stack** RNNs, use top as input.
- **Bidirectional** RNNs (every step / last): $\mathbf{h}_t = [\mathbf{h}_t^f; \mathbf{h}_t^b]$.
- !! Remaining issues:
 - Two tasks in one:
 - Carrying info forward.
 - Representation for current task.
 - Short context focus.
 - Vanishing gradients.



LONG SHORT-TERM MEMORY (LSTM) NETWORK

- Manage information required in the long term.
- Adds context vector and masking gates to the representation.



REMAINING LIMITATIONS

- Recurrence leads to information loss and training issues.
- Sequential nature blocks parallelism (as in e.g. CNNs).

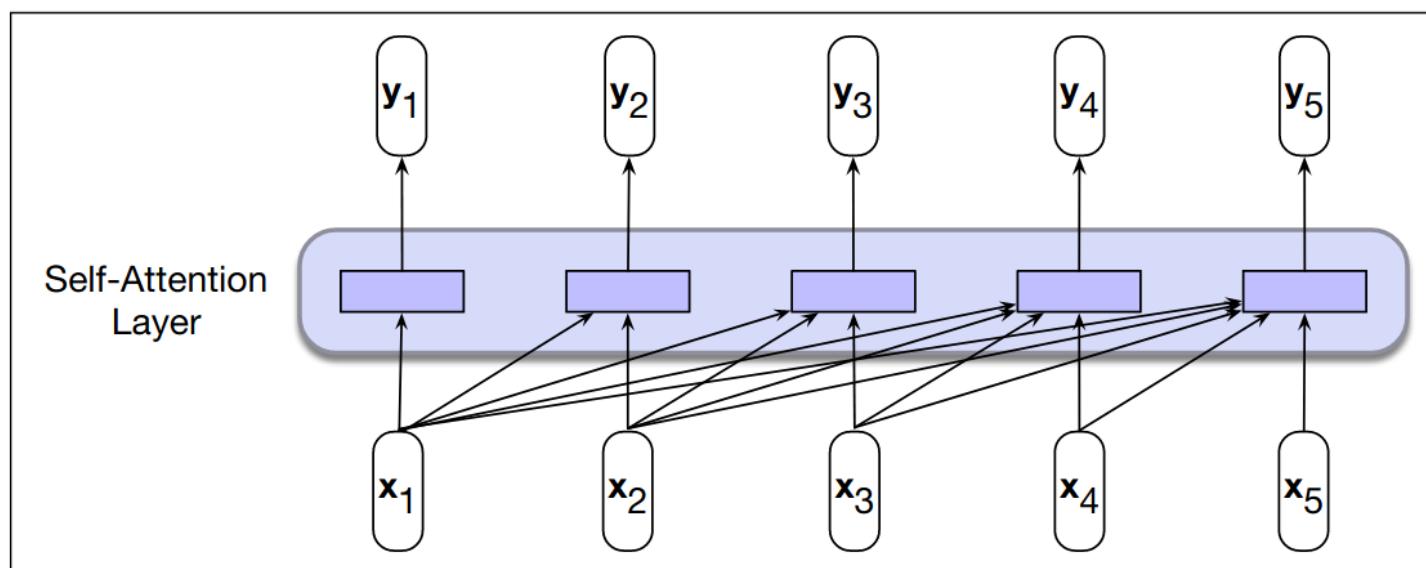


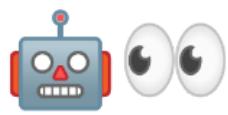
TRANSFORMERS



👀 ATTENTION

$$\begin{aligned}\alpha_{ij} &= \text{softmax}(\text{score}(\mathbf{x}_i, \mathbf{x}_j)) \quad \forall j \leq i & \text{score}(x_i, x_j) &= x_i \cdot x_j \\ &= \frac{\exp(\text{score}(\mathbf{x}_i, \mathbf{x}_j))}{\sum_{k=1}^i \exp(\text{score}(\mathbf{x}_i, \mathbf{x}_k))} \quad \forall j \leq i & \mathbf{y}_i &= \sum_{j \leq i} \alpha_{ij} \mathbf{x}_j\end{aligned}$$



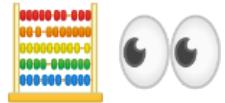


TRANSFORMER ATTENTION

- **query (\mathbf{q}_i)**: attention focus wrt preceding inputs -- \mathbf{W}^Q
- **key (\mathbf{k}_i)**: preceding input being wrt attention focus -- \mathbf{W}^K
- **value (\mathbf{v}_i)**: computes the output for attention focus -- \mathbf{W}^V

Project \mathbf{x}_i :

$\mathbf{q}_i = \mathbf{W}^Q \mathbf{x}_i; \mathbf{k}_i = \mathbf{W}^K \mathbf{x}_i; \mathbf{v}_i = \mathbf{W}^V \mathbf{x}_i$, where
dims: $1 \times d$. Later: $\mathbf{W}^Q \in \mathbb{R}^{d \times d}$, $\mathbf{W}^K \in \mathbb{R}^{d \times d}$,
and $\mathbf{W}^V \in \mathbb{R}^{d \times d}$.



MATRIX SELF-ATTENTION

$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{q}_i \cdot \mathbf{k}_j$$



$$\text{score}(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{q}_i \cdot \mathbf{k}_j}{\sqrt{d_k}}$$

$$\mathbf{y}_i = \sum_{j \leq i} \alpha_{ij} \mathbf{v}_j$$

Actually: $\mathbf{X} \in \mathbb{R}^{N \times d}$ (embeddings),
so: $\mathbf{Q} \in \mathbb{R}^{N \times d}$, $\mathbf{K} \in \mathbb{R}^{N \times d}$,
and $\mathbf{V} \in \mathbb{R}^{N \times d}$, giving:

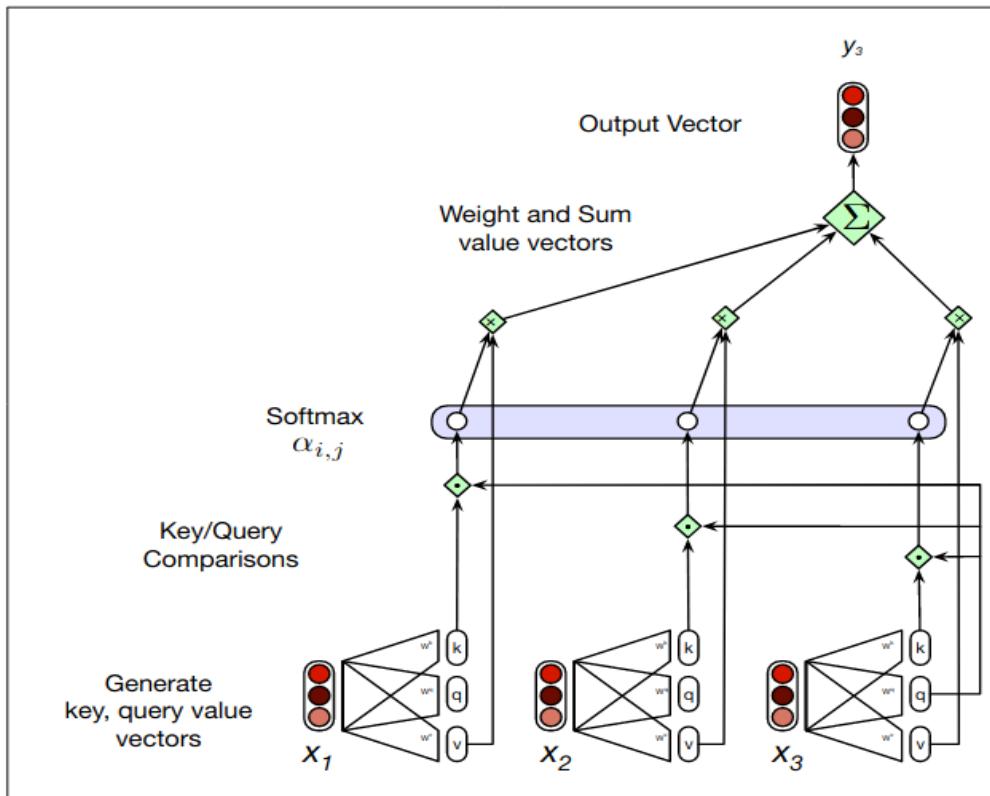
$$\mathbf{Q} = \mathbf{XW}^{\mathbf{Q}}; \mathbf{K} = \mathbf{XW}^{\mathbf{K}}; \mathbf{V} = \mathbf{XW}^{\mathbf{V}}$$

which gives:

$$\text{SelfAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{QK}^T}{\sqrt{d_k}}\right)\mathbf{V}$$



VISUALLY



N	q1•k1	$-\infty$	$-\infty$	$-\infty$	$-\infty$
q2•k1	q2•k2	$-\infty$	$-\infty$	$-\infty$	$-\infty$
q3•k1	q3•k2	q3•k3	$-\infty$	$-\infty$	
q4•k1	q4•k2	q4•k3	q4•k4	$-\infty$	
q5•k1	q5•k2	q5•k3	q5•k4	q5•k5	

Setting weights to -inf to hide future.



TRANSFORMER BLOCKS

$\mathbf{z} = \text{LayerNorm}(\mathbf{x} + \text{Self Attn}(\mathbf{x}))$

$\mathbf{y} = \text{LayerNorm}(\mathbf{z} + \text{FFNN}(\mathbf{z}))$

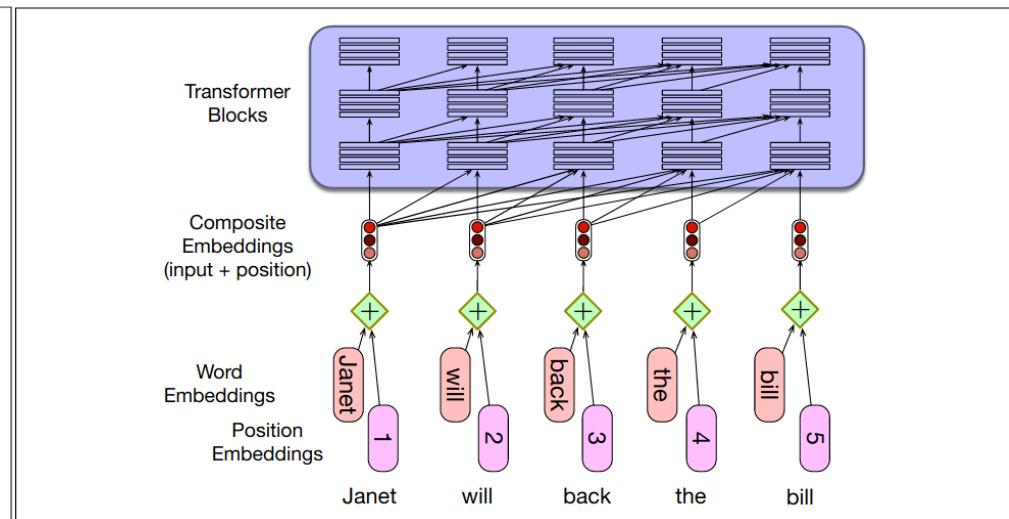
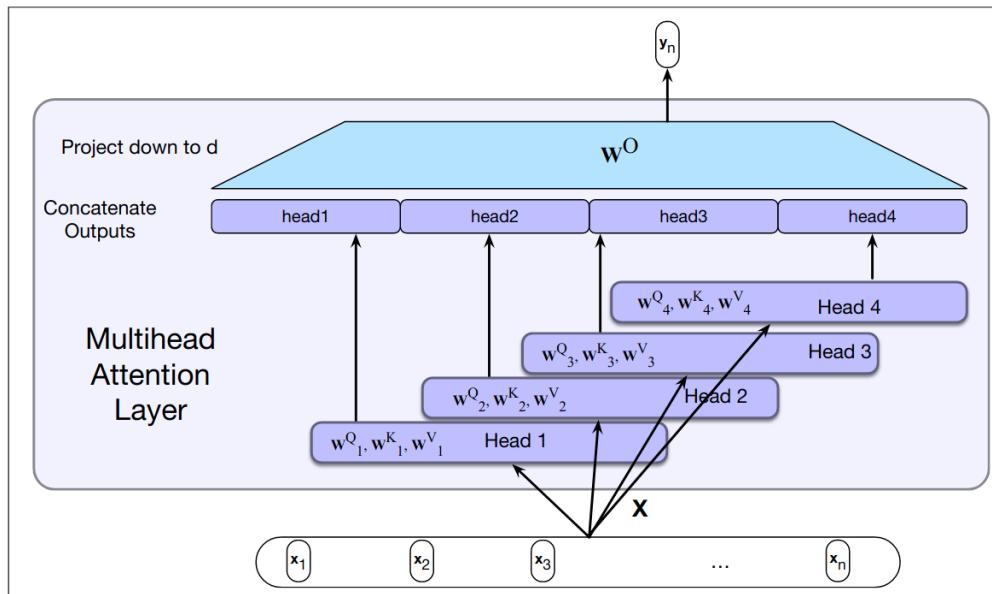
$$\hat{\mathbf{x}} = \frac{(\mathbf{x} - \mu)}{\sigma} \quad \text{LayerNorm} = \gamma \hat{\mathbf{x}} + \beta$$

$$\mu = \frac{1}{d_h} \sum_{i=1}^{d_h} x_i$$

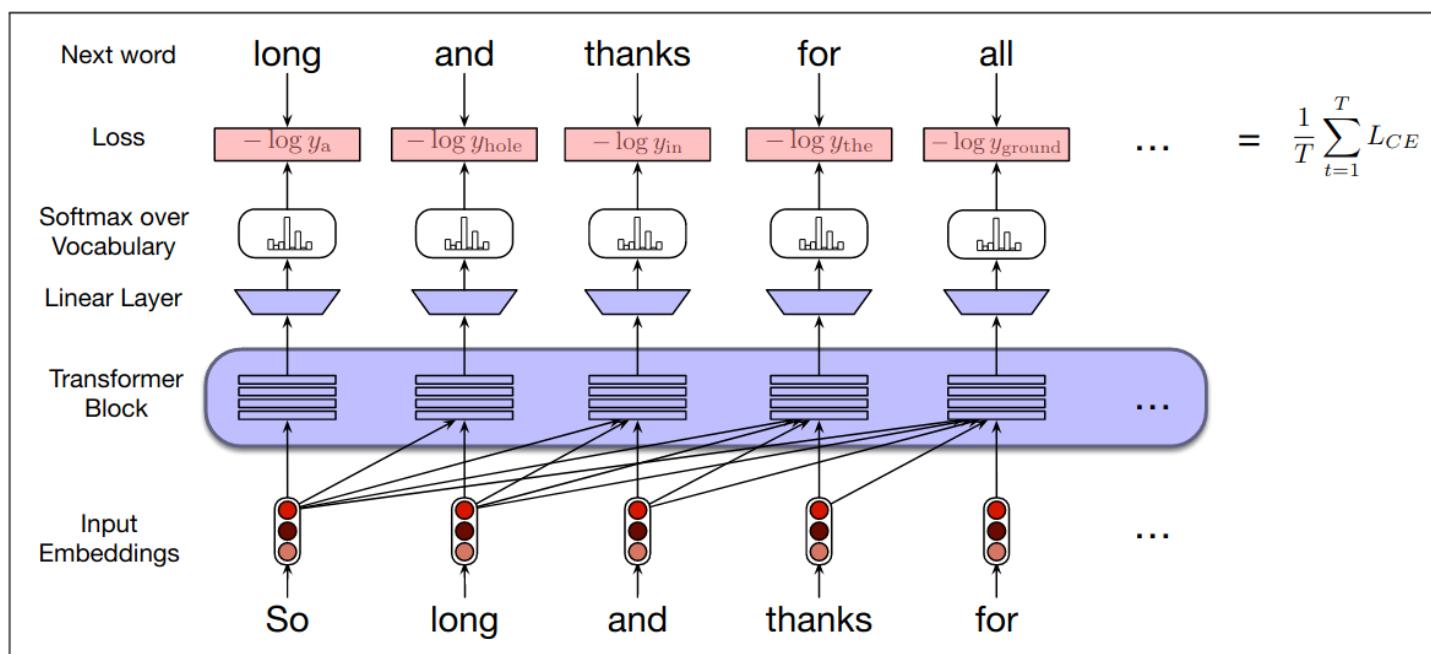
$$\sigma = \sqrt{\frac{1}{d_h} \sum_{i=1}^{d_h} (x_i - \mu)^2}$$

MORE TRICKS

$\text{MultiHeadAttn}(\mathbf{X}) = (\text{head}_1 \oplus \text{head}_2 \dots \oplus \text{head}_h) \mathbf{W}^O$
 $\mathbf{Q} = \mathbf{X}\mathbf{W}_i^Q; \mathbf{K} = \mathbf{X}\mathbf{W}_i^K; \mathbf{V} = \mathbf{X}\mathbf{W}_i^V$
 $\text{head}_i = \text{SelfAttention}(\mathbf{Q}, \mathbf{K}, \mathbf{V})$

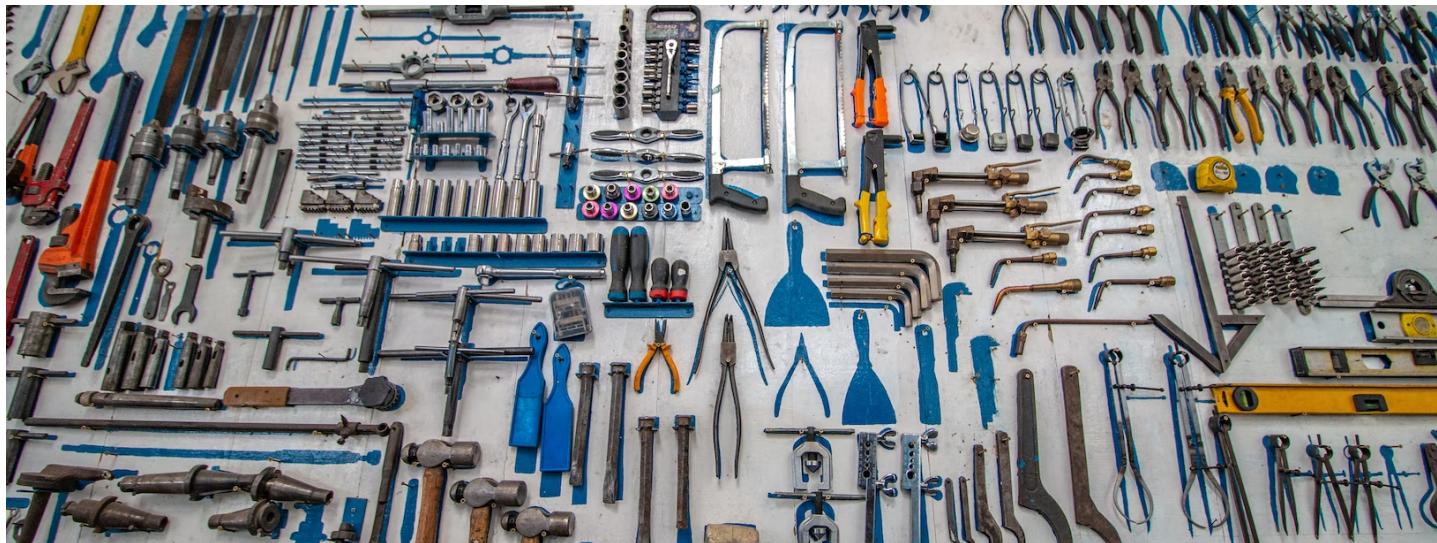


APPLIED





DEEP LEARNING LANDSCAPE





QUESTIONS

Post on the Discussion board and join class on Thursdays!