Lecture Notes for

# Neural Networks
# and Machine Learning

Transformers and Vision
Transformers

# Logistics and Agenda

- Logistics
  - Paper presentations
- Agenda
  - Transformers
- Next Time:
  - Vision Transformers
  - Paper Presentation
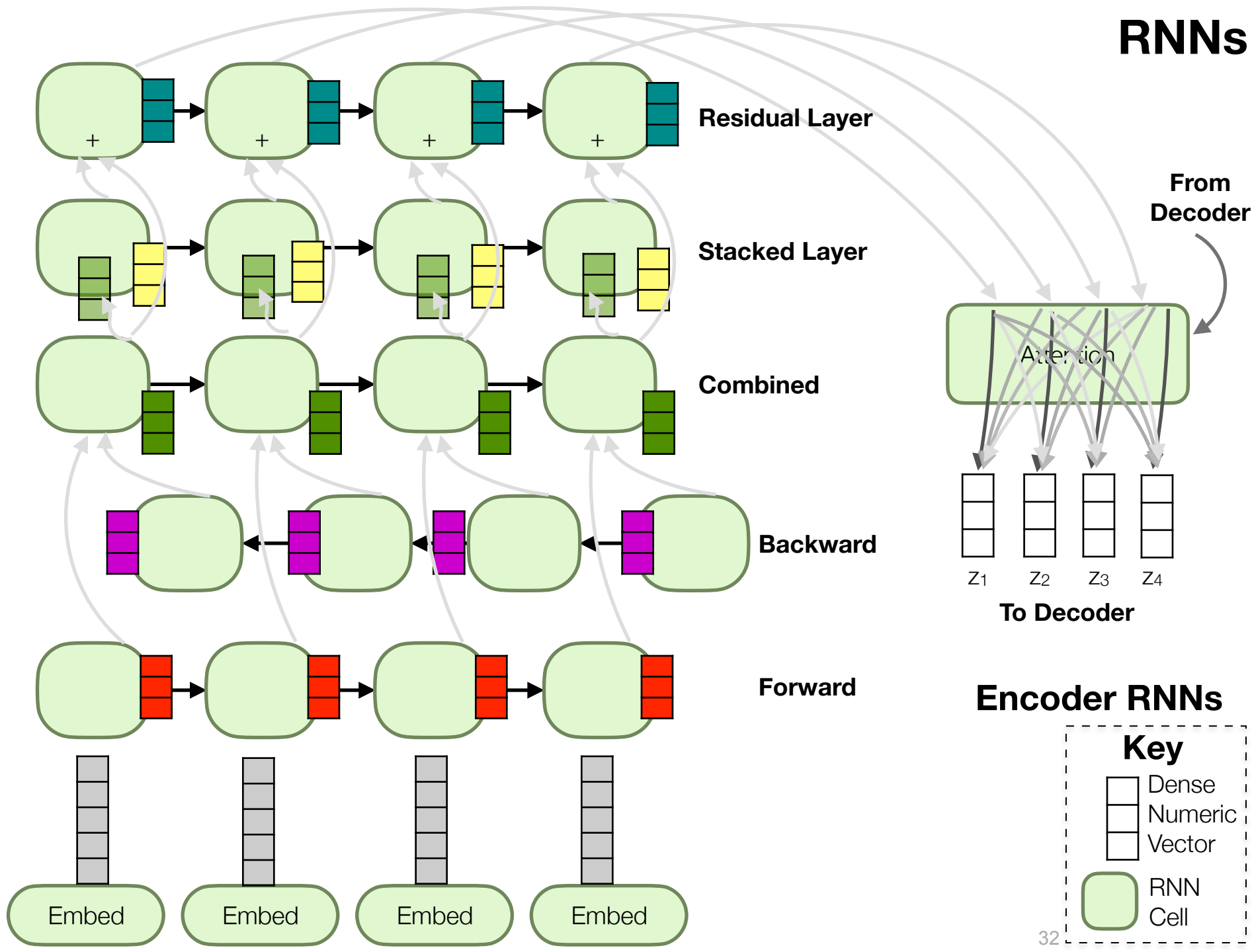  - Self-supervised learning and other consistency losses

# Transformers

**Dr Simone Stumpf** @DrSimoneS... · 13h ⋯
God grant me the confidence of an average machine learning expert.

# RNNs



**Residual Layer**

**Stacked Layer**

**Combined**

**Backward**

**Forward**

**From Decoder**

Attention

$Z_1$ $Z_2$ $Z_3$ $Z_4$

**To Decoder**

**Encoder RNNs**

Embed    Embed    Embed    Embed
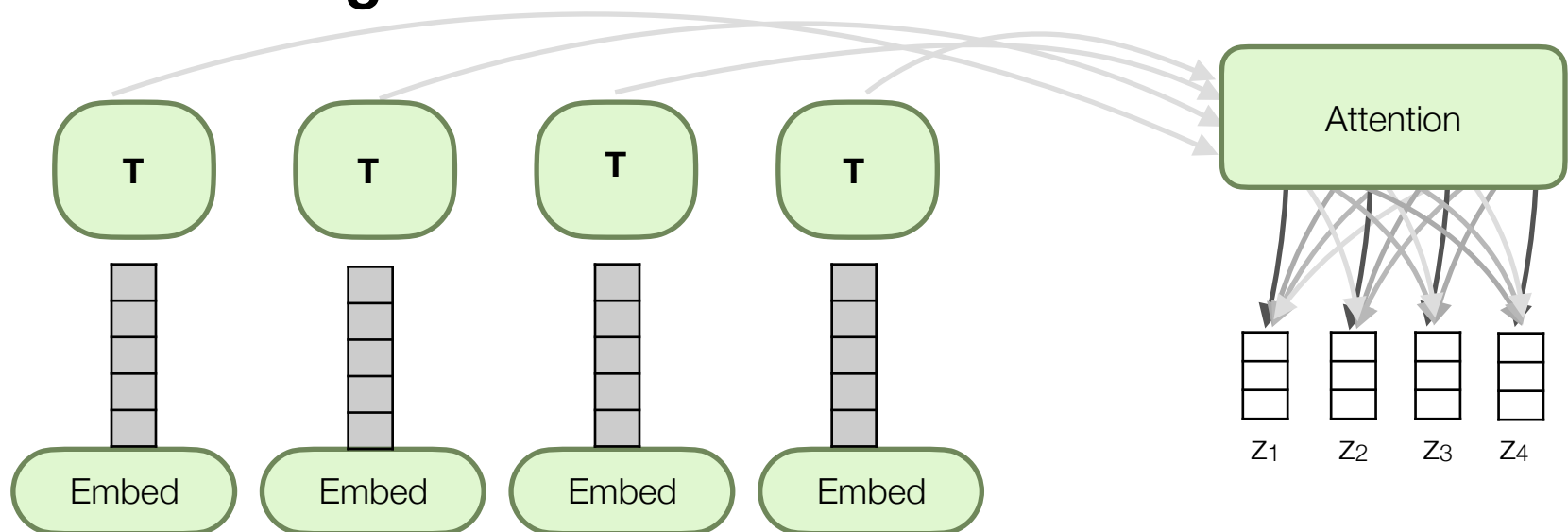
**Key**

Dense
Numeric
Vector

RNN
Cell

32

# Transformers Intuition

- Recurrent networks track state using an "updatable" state vector, but this takes lots of processing to across sequence

- Attention mechanism (in RNNs) already takes a weighted sum of state vectors to generate new token in a decoder

- … so why not just use attention on a transformation of the embedding vectors? **Do away with the recurrent state vector all together?**

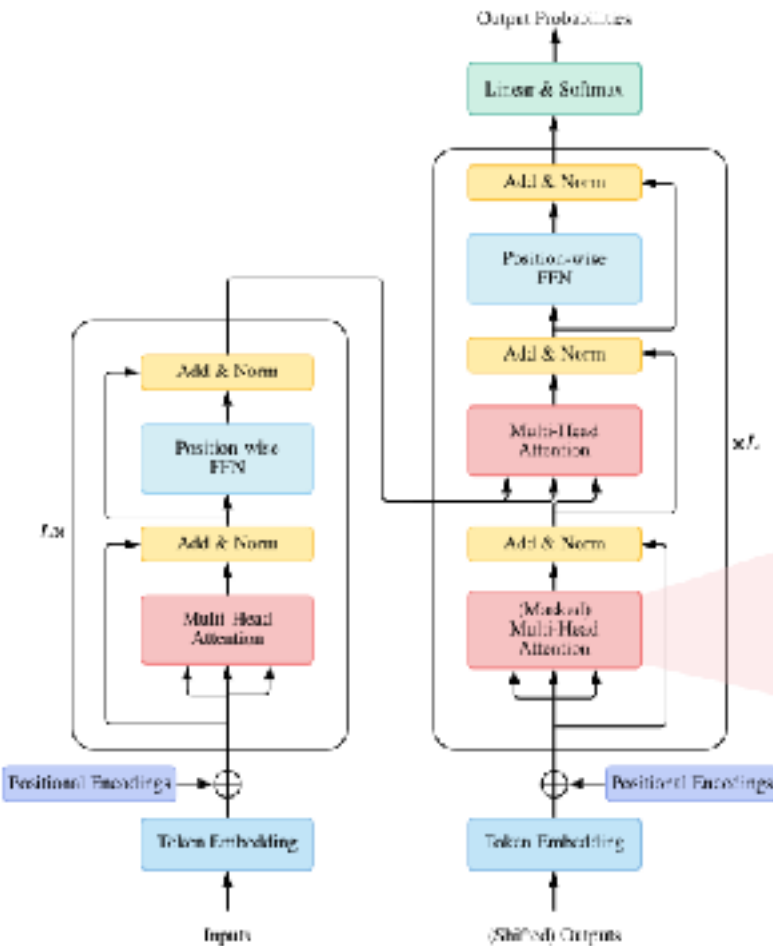# Attention is All You Need

- **Continued Motivation**:
  - RNNs are not inherently parallelized or efficient at remembering based on state vector
  - CNNs are not resilient to long-term word relationships, limited by filter size
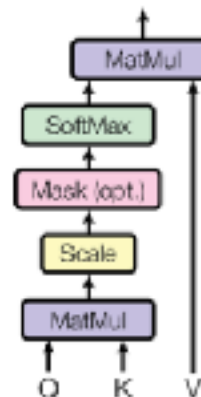- **Transformer Solution:**
  - Build attention into model from the **beginning**
  - Compare all words to each other through **self-headed** attention
  - Define a notion of "**position**"in the sequence
  - *Should be resilient to long term relationships and be highly parallelized for GPU computing!!*
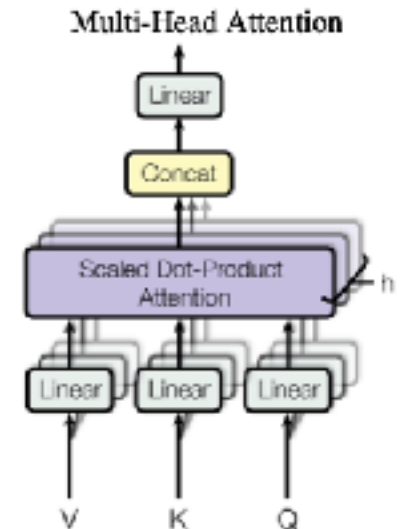
# Transformer



Scaled Dot-Product Attention

for each word

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

Multi-Head Attention

more than one
Q,K,V use in document

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, ..., \text{head}_h)W^O$$

$$\text{where head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

https://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf

# Transformer: in more detail



Input: Thinking Machines

Embedding: $X_1$ $X_2$

Outputs of Matrix Multiplications:

Queries: $q_1$ $q_2$

Keys: $k_1$ $k_2$

Values: $v_1$ $v_2$

Learned Matrices: $W^Q$, $W^K$, $W^V$

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

# Transformer: in more detail

Input

Embedding

Queries

Keys

Values

Calc. q, k, v for each word

**Thinking**  **Machines**

$x_1$    $x_2$

$q_1$    $q_2$

$k_1$    $k_2$

$v_1$    $v_2$

Straight forward to do this operation in matrix form:

X    $W^Q$    Q

Thinking Machines    $\times$    =    $\leftarrow d_k \rightarrow$

X    $W^K$    K

Thinking Machines    $\times$    =    $\leftarrow d_k \rightarrow$

X    $W^V$    V

Thinking Machines    $\times$    =    $\leftarrow d_v \rightarrow$

Score

$q_1 \cdot k_1 = 112$    $q_1 \cdot k_2 = 96$

Divide by 8 ($\sqrt{d_k}$)

in visual, $d_k = 3$

14    12

Softmax

0.88    0.12

Softmax
X
Value

Calc weights for $z_1$

$v_1$    $v_2$

weighted sum for all words in document

Q    $K^T$    V

$\mathrm{softmax}\left(\dfrac{\phantom{XX} \times \phantom{XX}}{\sqrt{d_k}}\right)$

Z

= $\begin{matrix} z_1 \\ z_2 \end{matrix}$

Sum

$z_1$    $z_2$

attention for word 1    attention for word 2

Size of W matrices:
$W^V$: |Embed Size| x $d_v$
$W^{Q,K}$: |Embed Size| x $d_k$

Size of Q,K,V:
|Seq Len| x $d_v$

**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

# Transformer: Multi-headed Attention



one row for each word

Length of row determined by $W^o$

$x\ W^o = Z$

# Transformer: Positional Encoding



- Objective: add notion of position to embedding
- Attempt in paper: add sin/cos to embedding
- But could be anything that encodes position
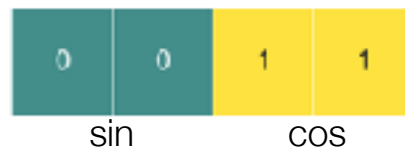
$p$: in sequence

$d\_m$: 0.5*max dim of embed

$i$ = position in embed

$$PE_{(p,i \in 0...d_m-1)} = \sin(p/10000^{i/d_m})$$

$$PE_{(p,i \in d_m...2d_m)} = \cos(p/10000^{(i-d_m)/d_m})$$

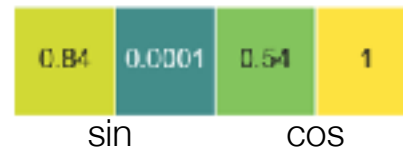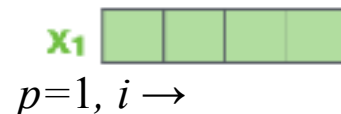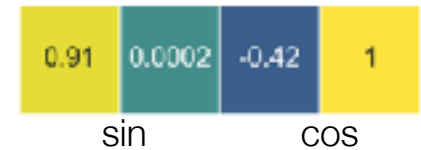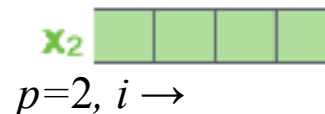Now use the new embeddings, with position, into transformer architecture
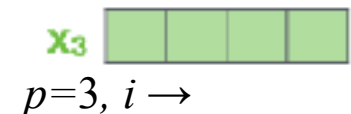


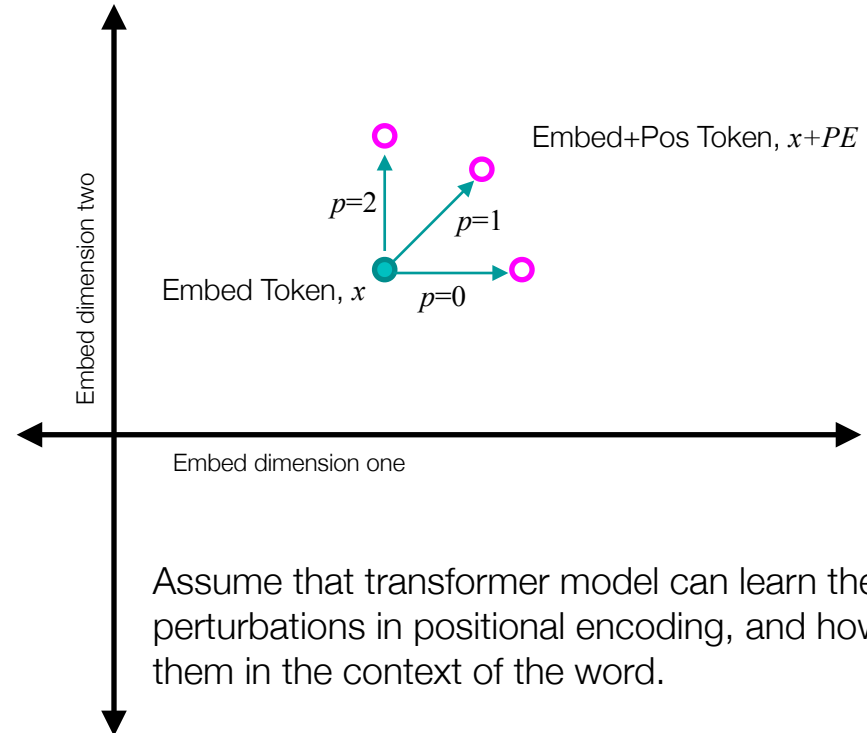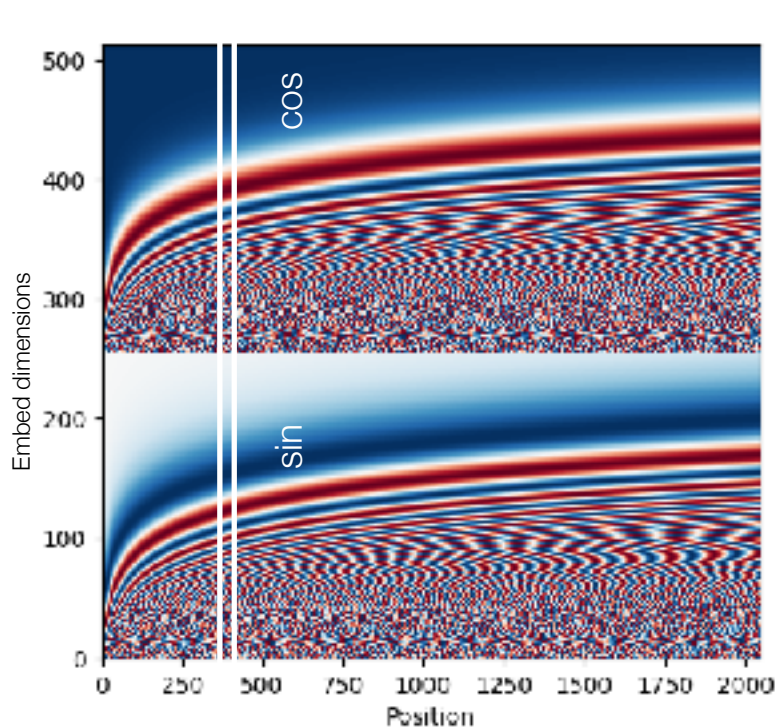$p=1, i \rightarrow$

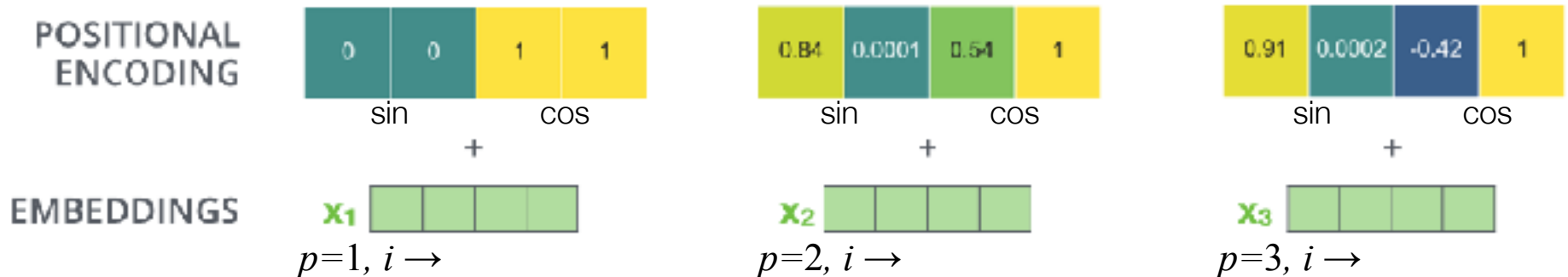$p=2, i \rightarrow$

$p=3, i \rightarrow$

**Hypothesis**: Now the word proximity is encoded in the embedding matrix, with other pertinent information. Well, it does help… so it could be true that this is a good way to do it.
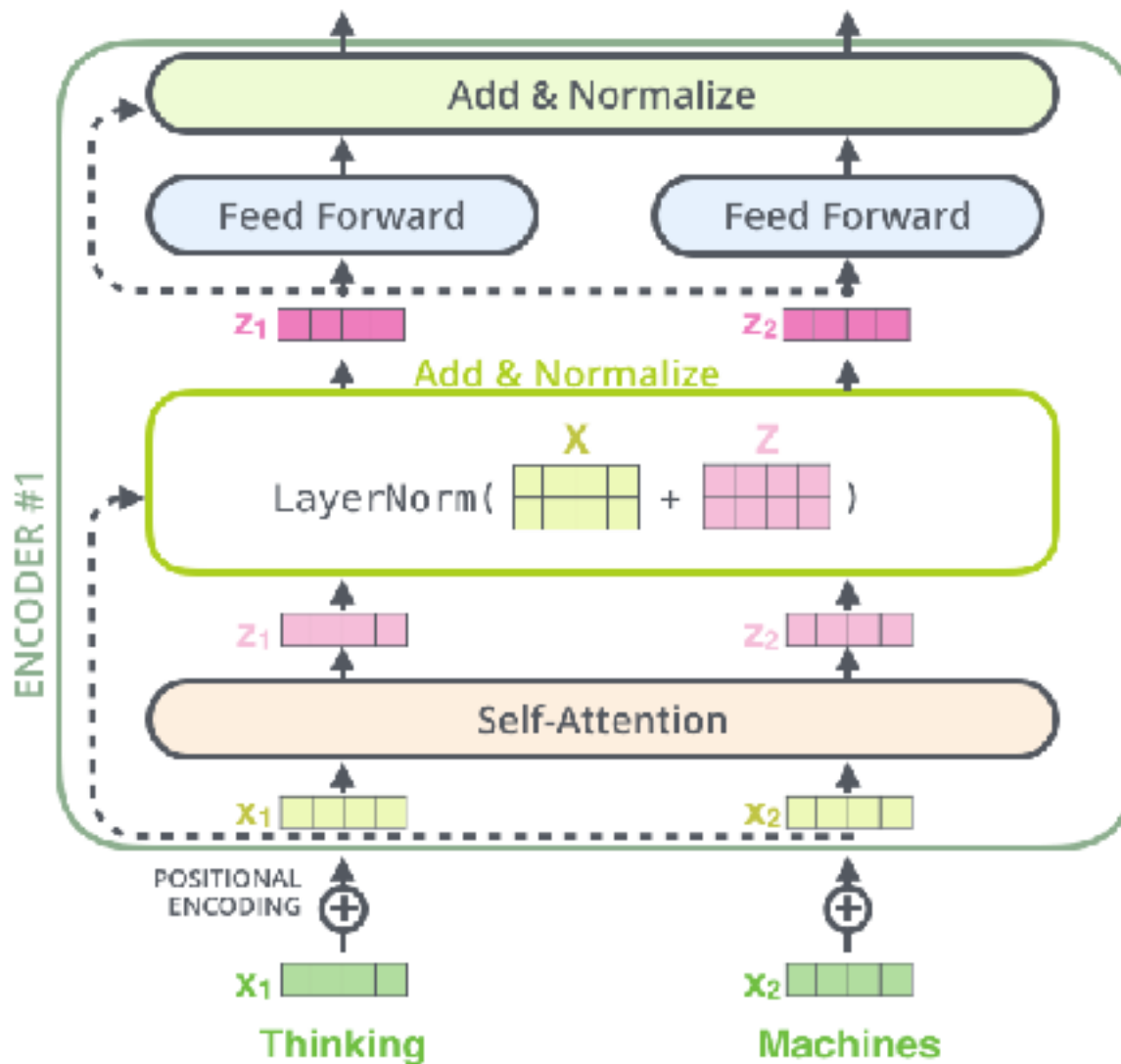
**Excellent Blog on Transformers:** http://jalammar.github.io/illustrated-transformer/

# Positional Intuition, Geometrically



Embed+Pos Token, $x+PE$

$p=2$

$p=1$

Embed Token, $x$

$p=0$

Embed dimension two

Embed dimension one

Assume that transformer model can learn the small perturbations in positional encoding, and how to use them in the context of the word.

**POSITIONAL ENCODING**

| 0 | 0 | 1 | 1 |
|---|---|---|---|

sin          cos

+

**EMBEDDINGS**

$X_1$

$p=1, i \rightarrow$

| 0.84 | 0.0001 | 0.54 | 1 |
|---|---|---|---|

sin          cos

+

$X_2$

$p=2, i \rightarrow$

| 0.91 | 0.0002 | -0.42 | 1 |
|---|---|---|---|

sin          cos

+
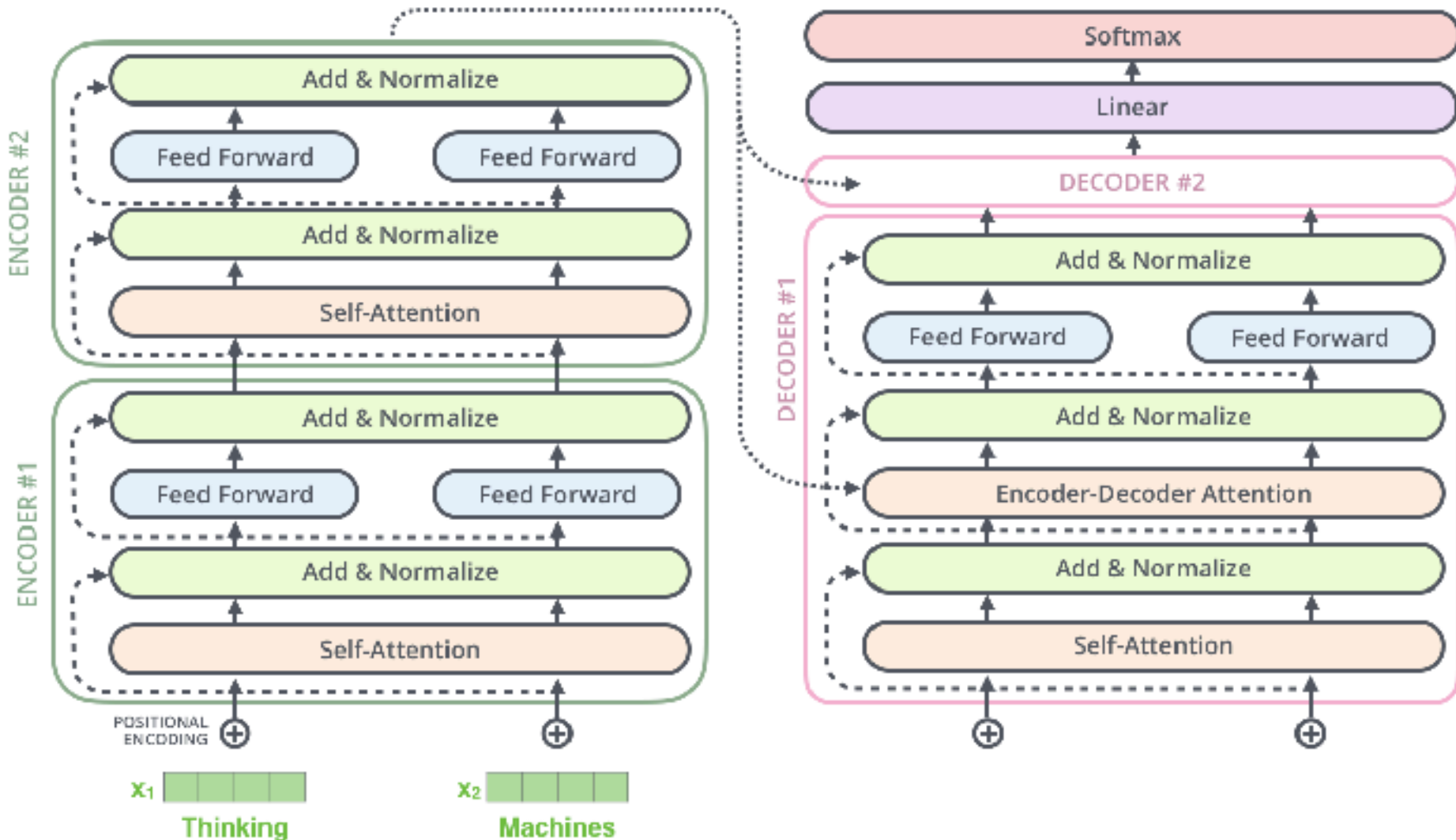
$X_3$

$p=3, i \rightarrow$

# Transformer: Residual Connections



LN: prevents vanishing gradients from softmax in attention

# Transformer: Putting it all together

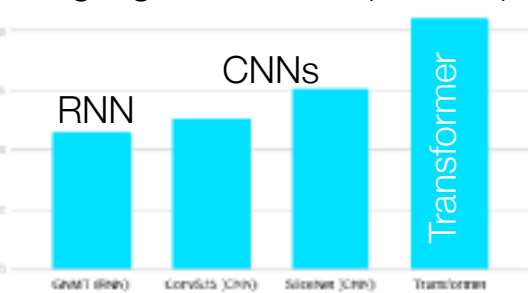# Transformer: Putting it all together

# Results

Language Translation (German)          Language Translation (French)



https://ai.googleblog.com/2017/08/transformer-novel-neural-network.html

**Implementations**:
- Many open source Keras or Tensorflow Implementations Exist
  - https://www.tensorflow.org/text/tutorials/transformer
- Many people like PyTorch for this
- HuggingFace has asme great implementations for free
  - https://huggingface.co/docs/transformers/index

## Text Classification GLUE

| Task | Dataset Variant | Best Model |
|---|---|---|
| **Text Classification** | GLUE | deberta-v3-small |
| **Sentiment Analysis** | SST-2 Binary classification | T5-11B |
| **Semantic Textual Similarity** | STS Benchmark | StructBERT+RoBERTa ensemble |
| **Natural Language Inference** | MultiNLI | T5-11B |
| **Natural Language Inference** | RTE | PaLM 540B |

General Language Understanding Evaluation (**GLUE**) benchmark is a collection of nine natural language understanding tasks, including single-sentence tasks CoLA and SST-2, similarity and paraphrasing tasks …

## IMDb Classifications

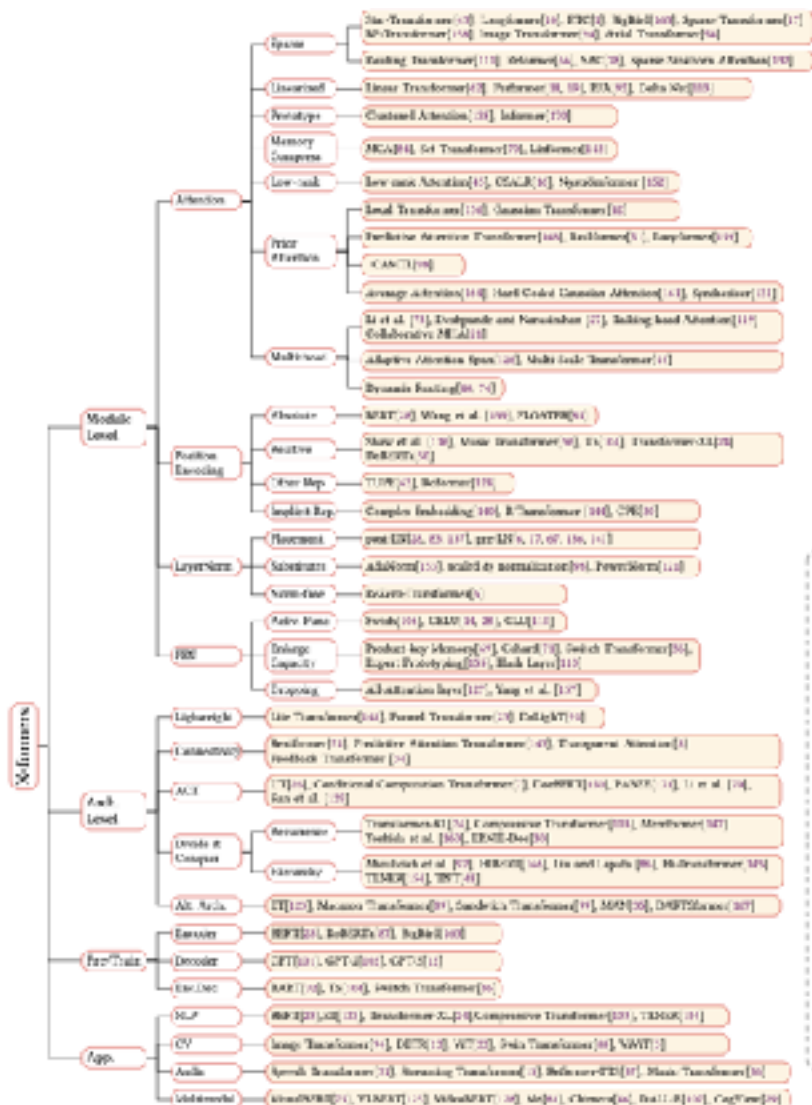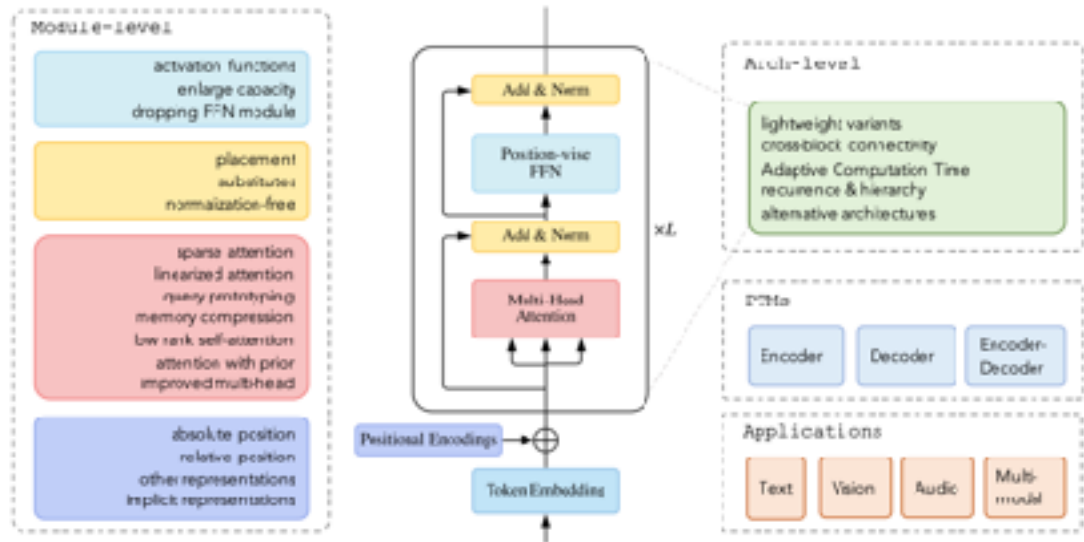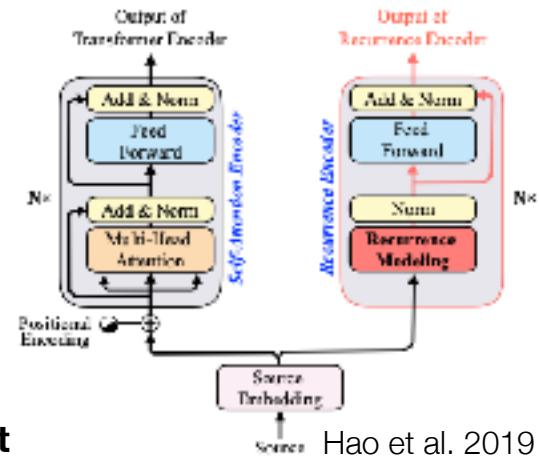| Task | Dataset Variant | Best Model |
|---|---|---|
| **Text Classification** | IMDb | ERNIE-Doc-Large |
| **Sentiment Analysis** | IMDb | XLNet |
| **Sentiment Analysis** | User and product information | MA-BERT |
| **SQL Parsing** | IMDb | Seq2Seq with copying |
| **Node Clustering** | IMDb | MAGNN |
| **Graph Similarity** | IMDb | SimGNN |
| **Link Prediction** | IMDb | Event2vec |

# MELVA Results (my lab)

- Measuring English Language Vocabulary Acquisition
- Or results from my lab:
  - Can students use science terms in a sentence?
  - Collect and transcribe student verbal responses regarding a scientific term
- Combine transcribed sentence and "good example"
- Collected about 6000 sentences
- Put through a language model (recurrent or transformer)
- Transfer learn based upon LM output
  - Without transformer LM: ~78%
  - With transformer LM: ~84%

# Lots of Transformer Variants



Architecture Tuning Matters Some
Pre-Training Matters
Sparse Attention Matters
Positional Encoding Doesn't
Recurrence Might… ?
**X-formers are NOT just for Text**

Hao et al. 2019

Lin et al "Survey of X-formers, 2021, https://arxiv.org/pdf/2106.04554.pdf

Lecture Notes for

# Neural Networks and Machine Learning

Transformers

**Next Time:**
SSL, Multi-Modal and Multi-Task
**Reading:** Keras F-API