# Comparing RNNs to Transformer

**RNNs**
(+) LSTMs work reasonably well for long sequences.
(-) Expects an ordered sequences of inputs
(-) Sequential computation: subsequent hidden states can only be computed after the previous ones are done.

**Transformer:**
(+) Good at long sequences. Each attention calculation looks at all inputs.
(+) Can operate over unordered sets or ordered sequences with positional encodings.
(+) Parallel computation: All alignment and attention scores for all inputs can be done in parallel.
(-) Requires a lot of memory: N x M alignment and attention scalers need to be calculated and stored for a single self-attention head. (but GPUs are getting bigger and better)

# Attention Is All You Need

**Ashish Vaswani**[*]
Google Brain
avaswani@google.com

**Noam Shazeer**[*]
Google Brain
noam@google.com

**Niki Parmar**[*]
Google Research
nikip@google.com

**Jakob Uszkoreit**[*]
Google Research
usz@google.com

**Llion Jones**[*]
Google Research
llion@google.com

**Aidan N. Gomez**[* †]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser**[*]
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin**[* ‡]
illia.polosukhin@gmail.com

"ImageNet Moment for Natural Language Processing"

**Pretraining**:
Download a lot of text from the internet

Train a giant Transformer model for language modeling

**Finetuning:**
Fine-tune the Transformer on your own NLP task

# On the Opportunities and Risks of Foundation Models
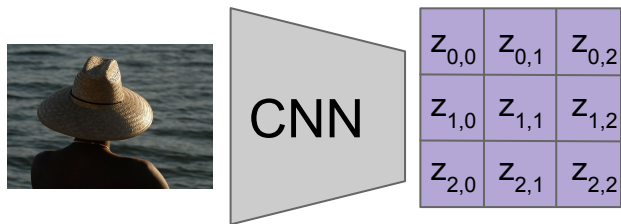
Rishi Bommasani*   Drew A. Hudson   Ehsan Adeli   Russ Altman   Simran Arora
Sydney von Arx   Michael S. Bernstein   Jeannette Bohg   Antoine Bosselut   Emma Brunskill
Erik Brynjolfsson   Shyamal Buch   Dallas Card   Rodrigo Castellon   Niladri Chatterji
Annie Chen   Kathleen Creel   Jared Quincy Davis   Dorottya Demszky   Chris Donahue
Moussa Doumbouya   Esin Durmus   Stefano Ermon   John Etchemendy   Kawin Ethayarajh
Li Fei-Fei   Chelsea Finn   Trevor Gale   Lauren Gillespie   Karan Goel   Noah Goodman
Shelby Grossman   Neel Guha   Tatsunori Hashimoto   Peter Henderson   John Hewitt
Daniel E. Ho   Jenny Hong   Kyle Hsu   Jing Huang   Thomas Icard   Saahil Jain
Dan Jurafsky   Pratyusha Kalluri   Siddharth Karamcheti   Geoff Keeling   Fereshte Khani
Omar Khattab   Pang Wei Koh   Mark Krass   Ranjay Krishna   Rohith Kuditipudi
Ananya Kumar   Faisal Ladhak   Mina Lee   Tony Lee   Jure Leskovec   Isabelle Levent
Xiang Lisa Li   Xuechen Li   Tengyu Ma   Ali Malik   Christopher D. Manning
Suvir Mirchandani   Eric Mitchell   Zanele Munyikwa   Suraj Nair   Avanika Narayan
Deepak Narayanan   Ben Newman   Allen Nie   Juan Carlos Niebles   Hamed Nilforoshan
Julian Nyarko   Giray Ogut   Laurel Orr   Isabel Papadimitriou   Joon Sung Park   Chris Piech
Eva Portelance   Christopher Potts   Aditi Raghunathan   Rob Reich   Hongyu Ren
Frieda Rong   Yusuf Roohani   Camilo Ruiz   Jack Ryan   Christopher Ré   Dorsa Sadigh
Shiori Sagawa   Keshav Santhanam   Andy Shih   Krishnan Srinivasan   Alex Tamkin
Rohan Taori   Armin W. Thomas   Florian Tramèr   Rose E. Wang   William Wang   Bohan Wu
Jiajun Wu   Yuhuai Wu   Sang Michael Xie   Michihiro Yasunaga   Jiaxuan You   Matei Zaharia
Michael Zhang   Tianyi Zhang   Xikun Zhang   Yuhui Zhang   Lucia Zheng   Kaitlyn Zhou
Percy Liang*[1]

Center for Research on Foundation Models (CRFM)
Stanford Institute for Human-Centered Artificial Intelligence (HAI)
Stanford University

# Image Captioning using Transformers

**Input**: Image **I**
**Output:** Sequence $\mathbf{y} = y_1, y_2,..., y_T$



Extract spatial features from a pretrained CNN

Features: H x W x D

# Image Captioning using Transformers

**Input**: Image $I$
**Output:** Sequence $y = y_1, y_2, ..., y_T$

**Encoder**: $c = T_W(z)$
where $z$ is spatial CNN features
$T_W(.)$ is the transformer encoder
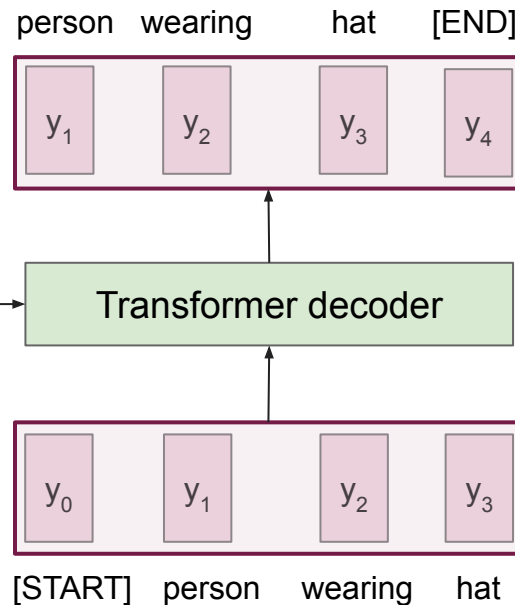


CNN

Extract spatial features from a pretrained CNN

Features:
$H \times W \times D$

$z_{0,0}$ $z_{0,1}$ $z_{0,2}$
$z_{1,0}$ $z_{1,1}$ $z_{1,2}$
$z_{2,0}$ $z_{2,1}$ $z_{2,2}$

$c_{0,0}$ $c_{0,1}$ $c_{0,2}$ ... $c_{2,2}$

Transformer encoder

$z_{0,0}$ $z_{0,1}$ $z_{0,2}$ ... $z_{2,2}$

# Image Captioning using Transformers

**Input**: Image **I**
**Output:** Sequence $\mathbf{y} = y_1, y_2,..., y_T$

**Decoder**: $y_t = T_\mathbf{D}(\mathbf{y}_{0:t-1}, \mathbf{c})$
where $T_\mathbf{D}(.)$ is the transformer decoder

**Encoder**: $\mathbf{c} = T_\mathbf{W}(\mathbf{z})$
where $\mathbf{z}$ is spatial CNN features
$T_\mathbf{W}(.)$ is the transformer encoder



Extract spatial features from a pretrained CNN

Features:
H x W x D

# The Transformer encoder block

Made up of N encoder blocks.

In vaswani et al. N = 6, $D_q$ = 512
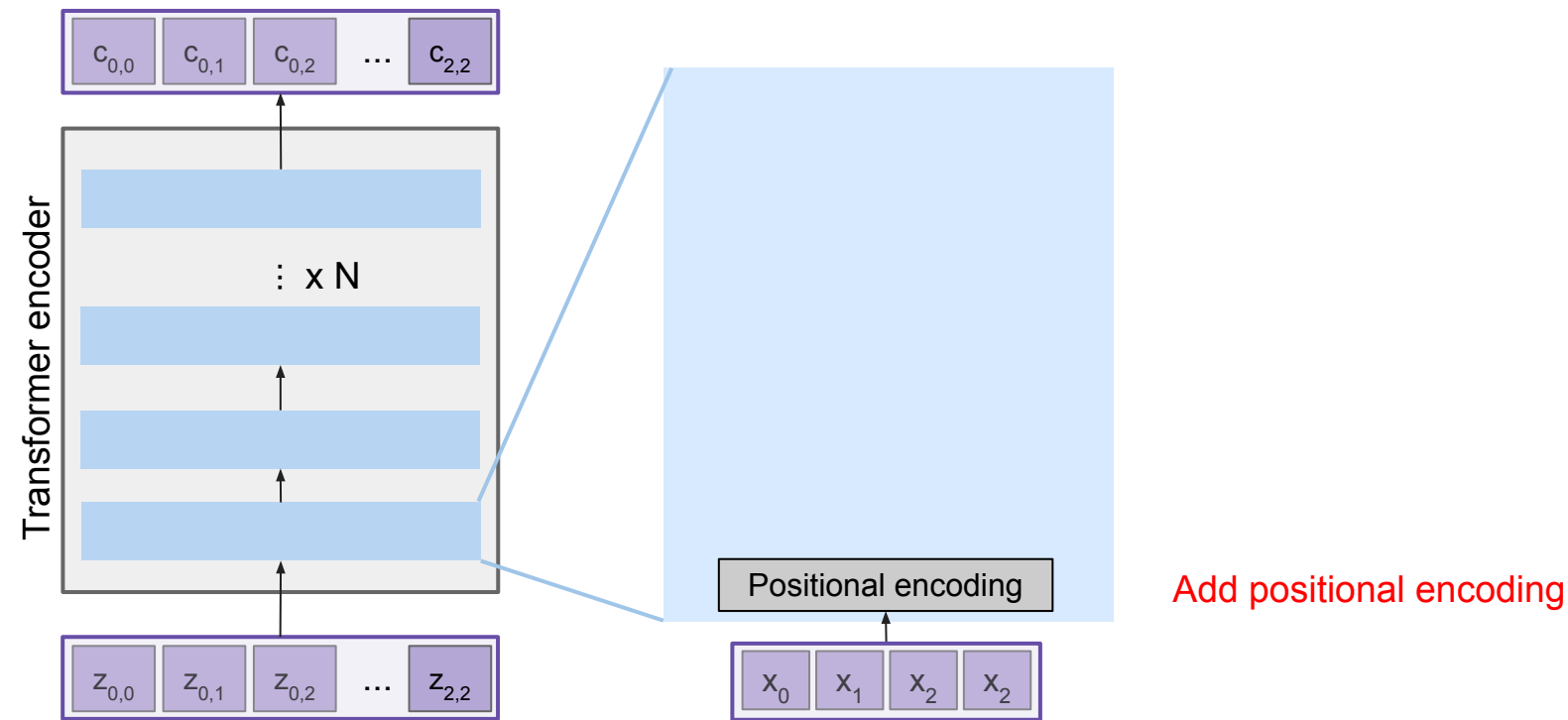
Vaswani et al, "Attention is all you need", NeurIPS 2017
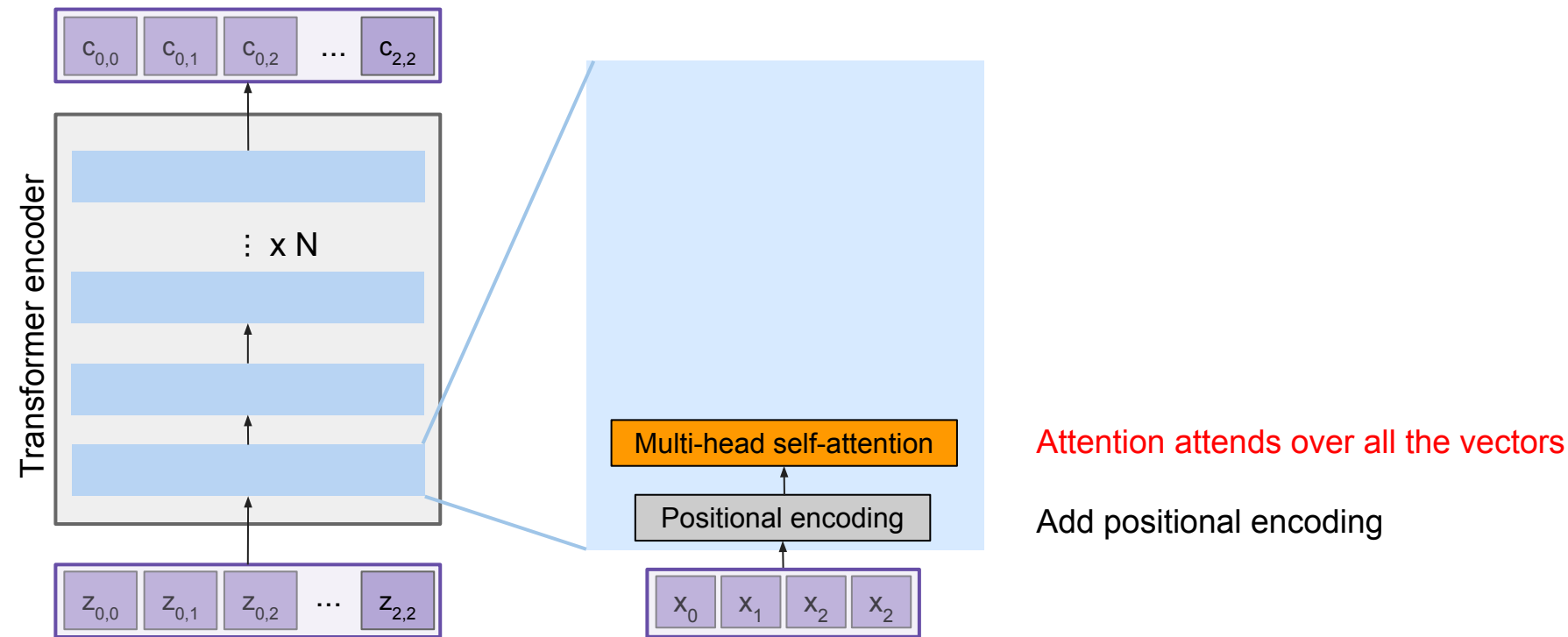
# The Transformer encoder block



$c_{0,0}$ $c_{0,1}$ $c_{0,2}$ ... $c_{2,2}$

Transformer encoder

: x N

$z_{0,0}$ $z_{0,1}$ $z_{0,2}$ ... $z_{2,2}$

Let's dive into one encoder block

$x_0$ $x_1$ $x_2$ $x_2$

Vaswani et al, "Attention is all you need", NeurIPS 2017

# The Transformer encoder block



Add positional encoding

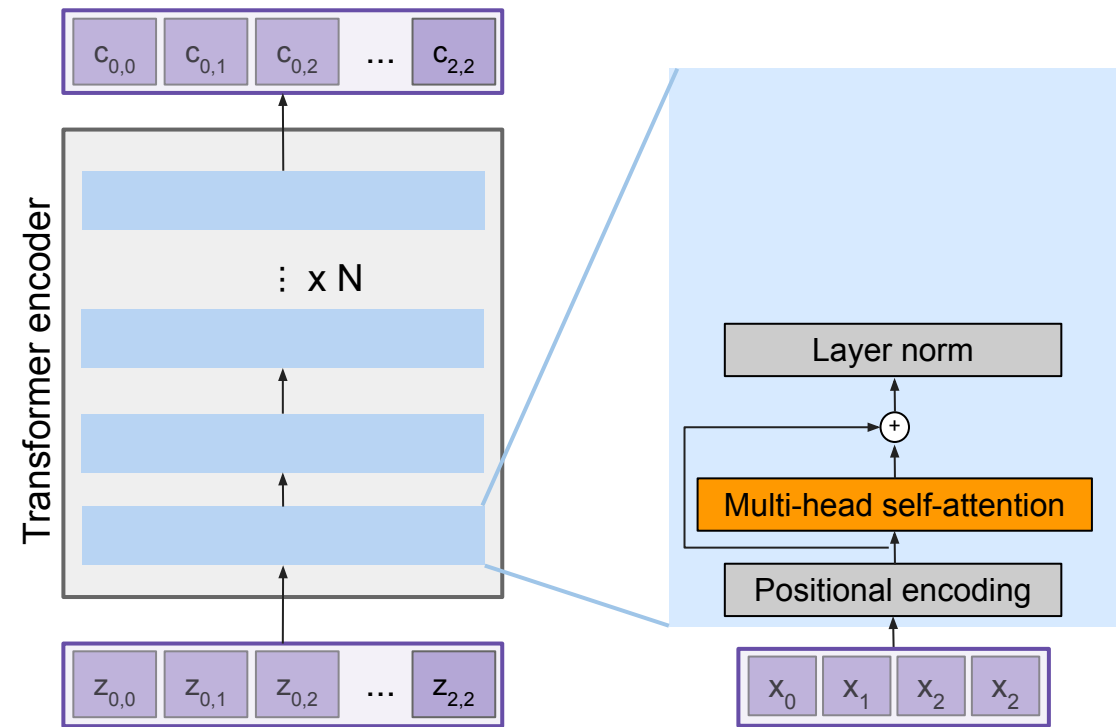Vaswani et al, "Attention is all you need", NeurIPS 2017

# The Transformer encoder block



Attention attends over all the vectors

Add positional encoding

Vaswani et al, "Attention is all you need", NeurIPS 2017
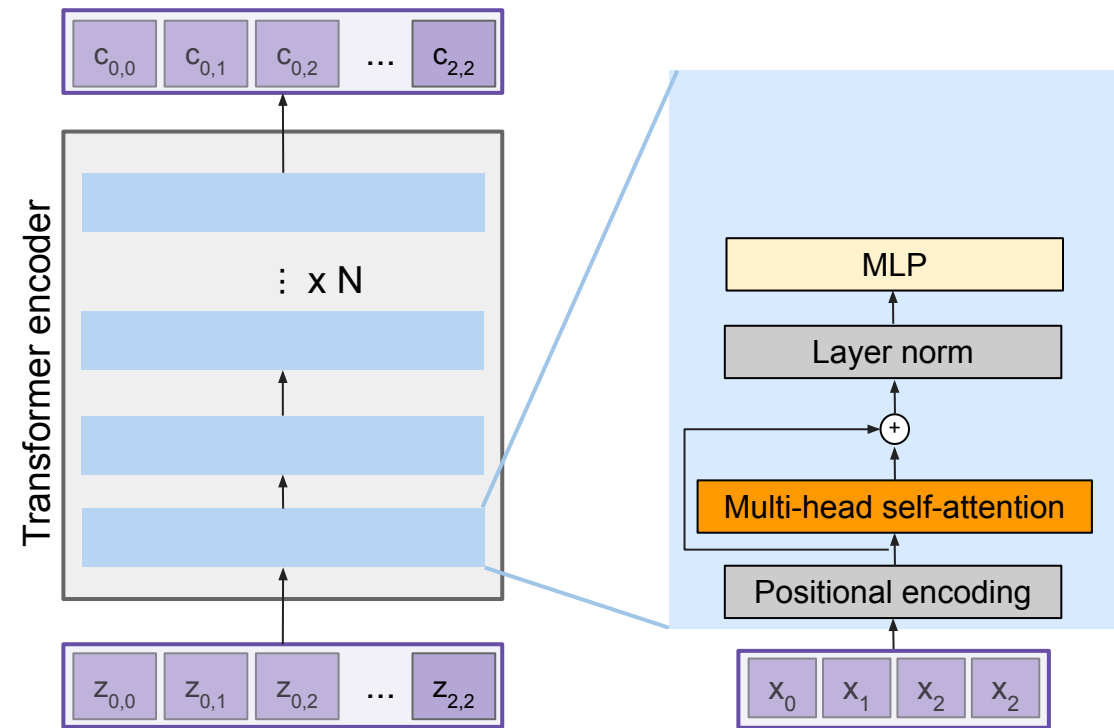
# The Transformer encoder block



Residual connection

Attention attends over all the vectors

Add positional encoding

Vaswani et al, "Attention is all you need", NeurIPS 2017

# The Transformer encoder block
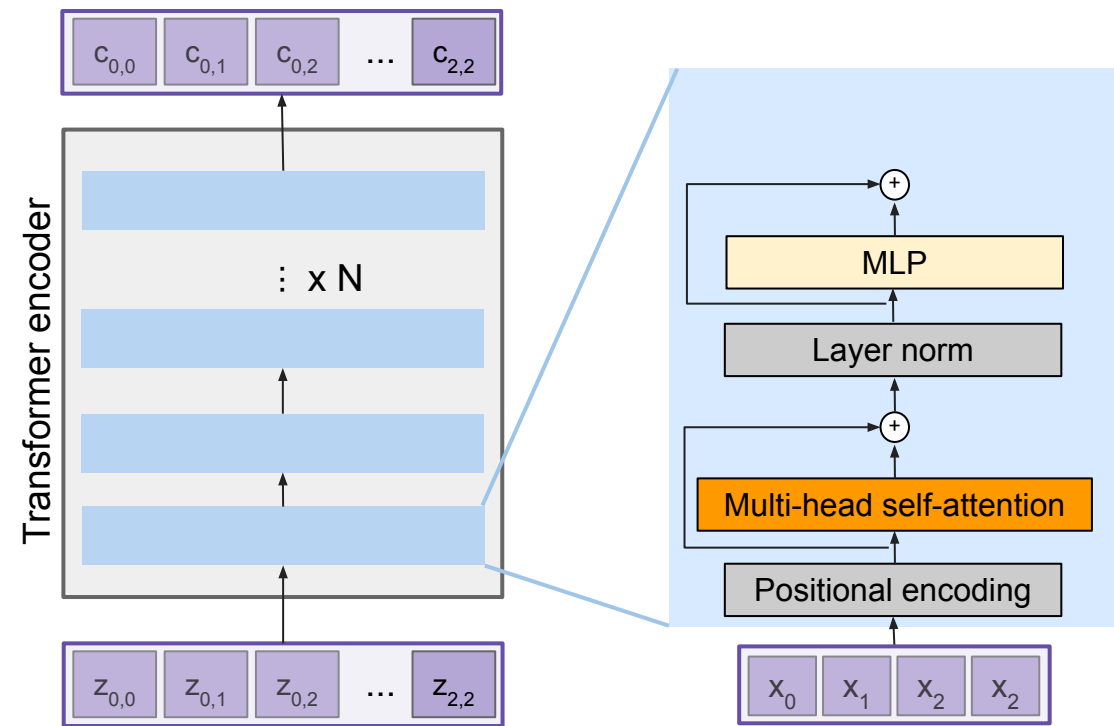


LayerNorm over each vector individually

Residual connection

Attention attends over all the vectors

Add positional encoding

Vaswani et al, "Attention is all you need", NeurIPS 2017

# The Transformer encoder block



MLP over each vector individually

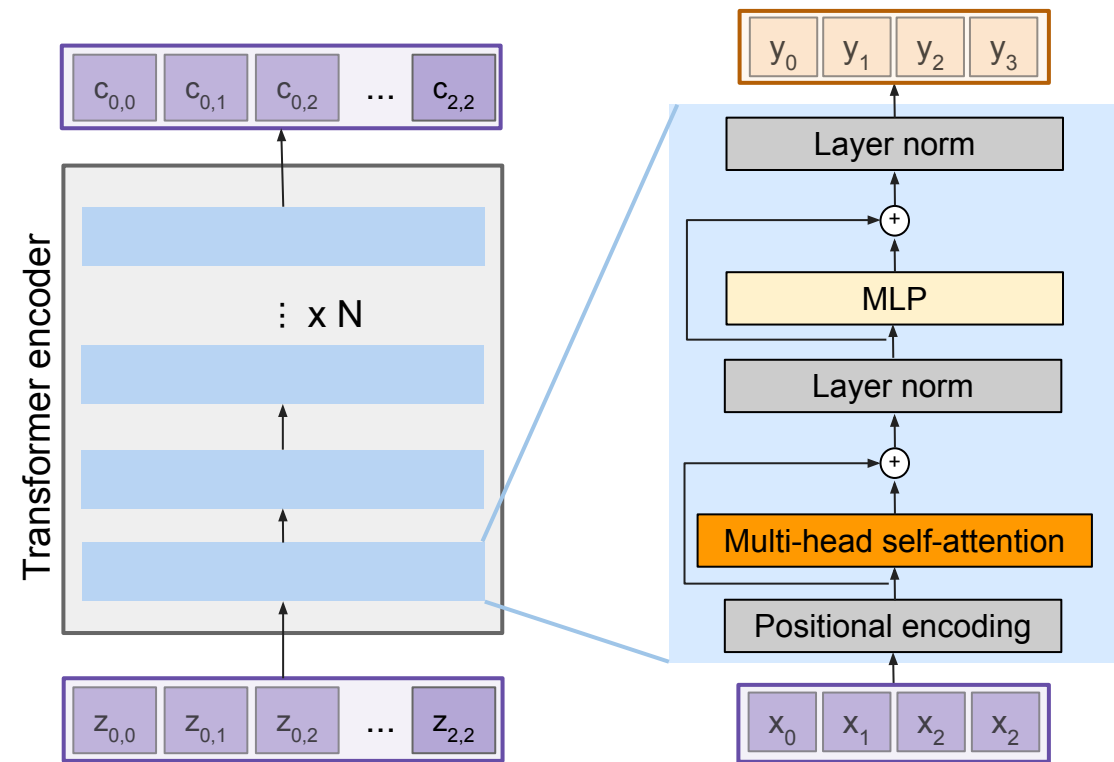LayerNorm over each vector individually

Residual connection

Attention attends over all the vectors

Add positional encoding

Vaswani et al, "Attention is all you need", NeurIPS 2017

# The Transformer encoder block



Residual connection

MLP over each vector individually

LayerNorm over each vector individually

Residual connection

Attention attends over all the vectors

Add positional encoding

Vaswani et al, "Attention is all you need", NeurIPS 2017

# The Transformer encoder block



**Transformer Encoder Block:**

**Inputs**: Set of vectors **x**
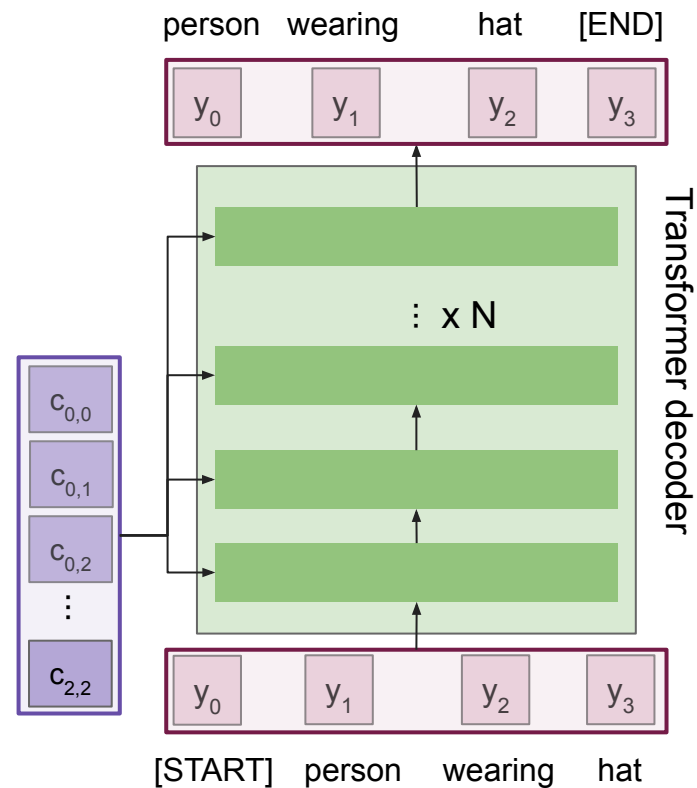**Outputs**: Set of vectors **y**

Self-attention is the only interaction between vectors.

Layer norm and MLP operate independently per vector.

Highly scalable, highly parallelizable, but high memory usage.

Vaswani et al, "Attention is all you need", NeurIPS 2017
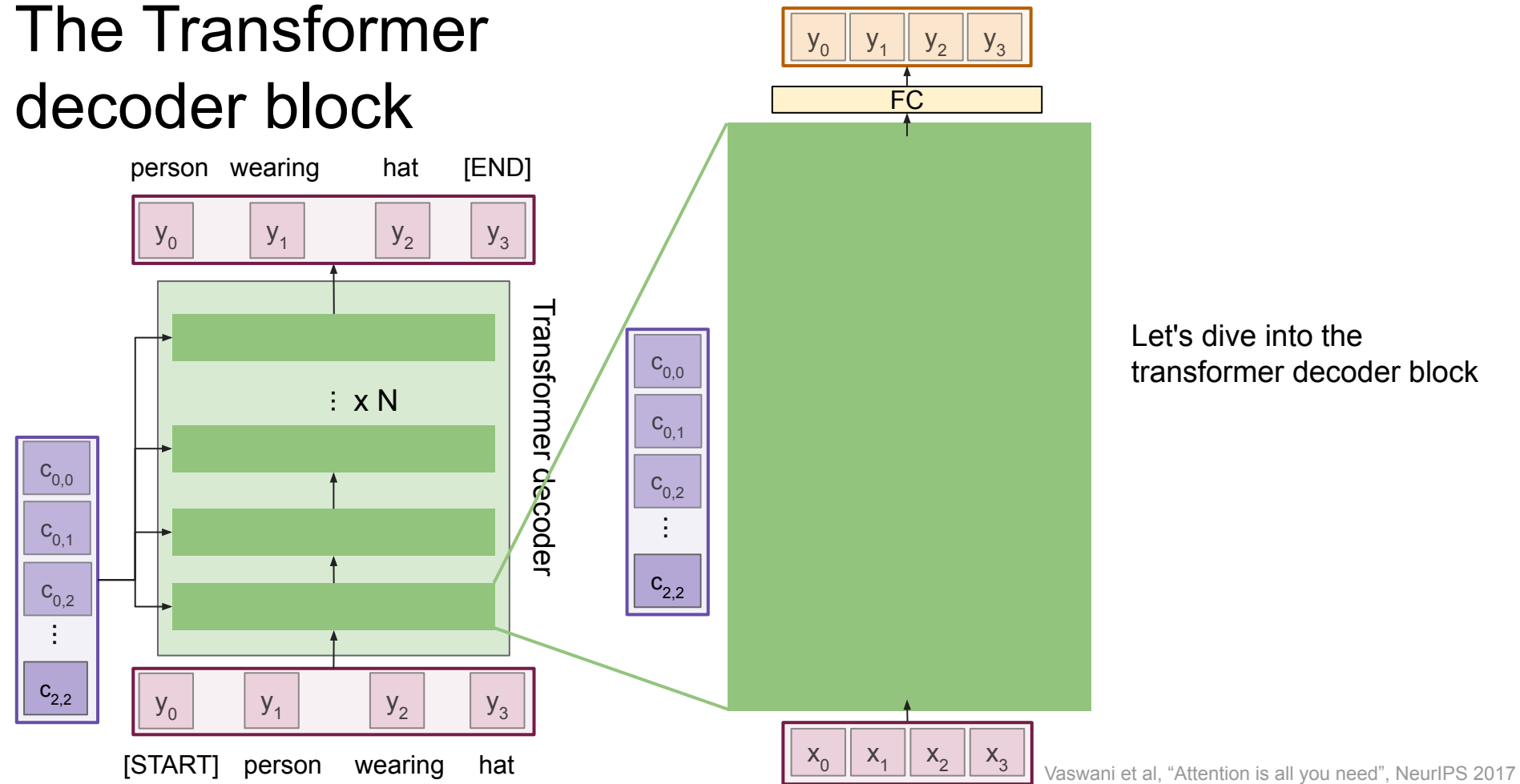
# The Transformer decoder block

person    wearing       hat      [END]

| $y_0$ | $y_1$ | $y_2$ | $y_3$ |

Transformer decoder

⋮ x N

$c_{0,0}$
$c_{0,1}$
$c_{0,2}$
⋮
$c_{2,2}$

| $y_0$ | $y_1$ | $y_2$ | $y_3$ |

[START]    person    wearing    hat

Made up of N decoder blocks.

In vaswani et al. N = 6, $D_q$ = 512

Vaswani et al, "Attention is all you need", NeurIPS 2017

# The Transformer decoder block

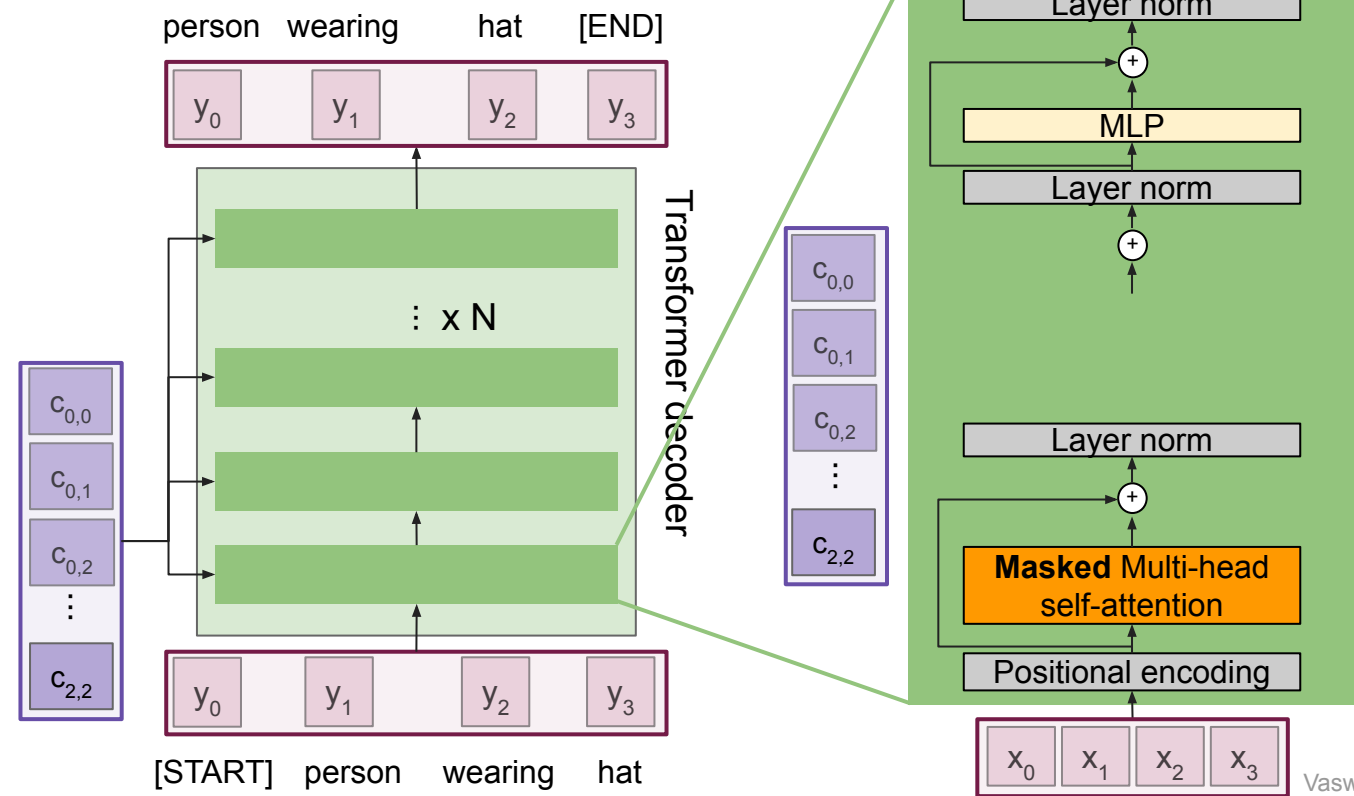person    wearing    hat    [END]



Let's dive into the transformer decoder block

Vaswani et al, "Attention is all you need", NeurIPS 2017

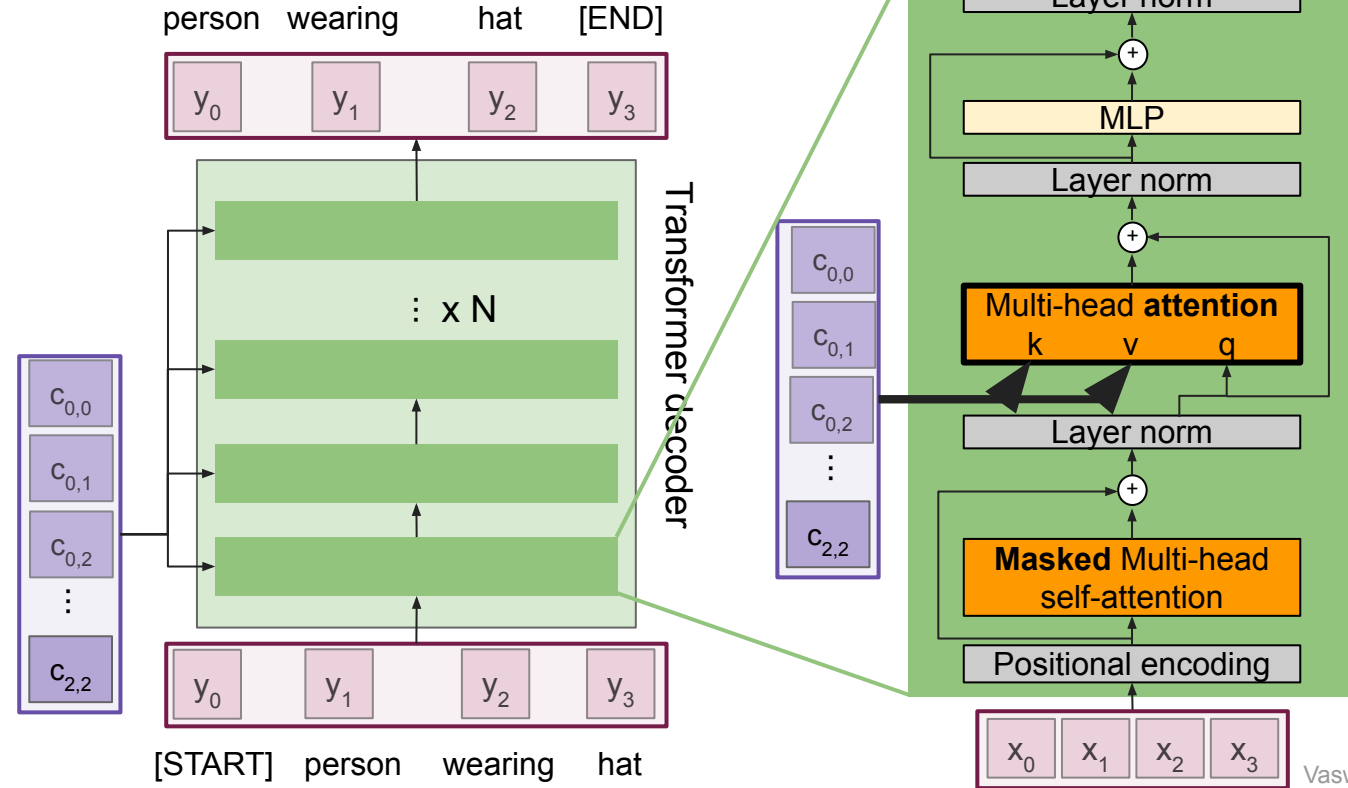# The Transformer Decoder block



Most of the network is the same the transformer encoder.

Vaswani et al, "Attention is all you need", NeurIPS 2017
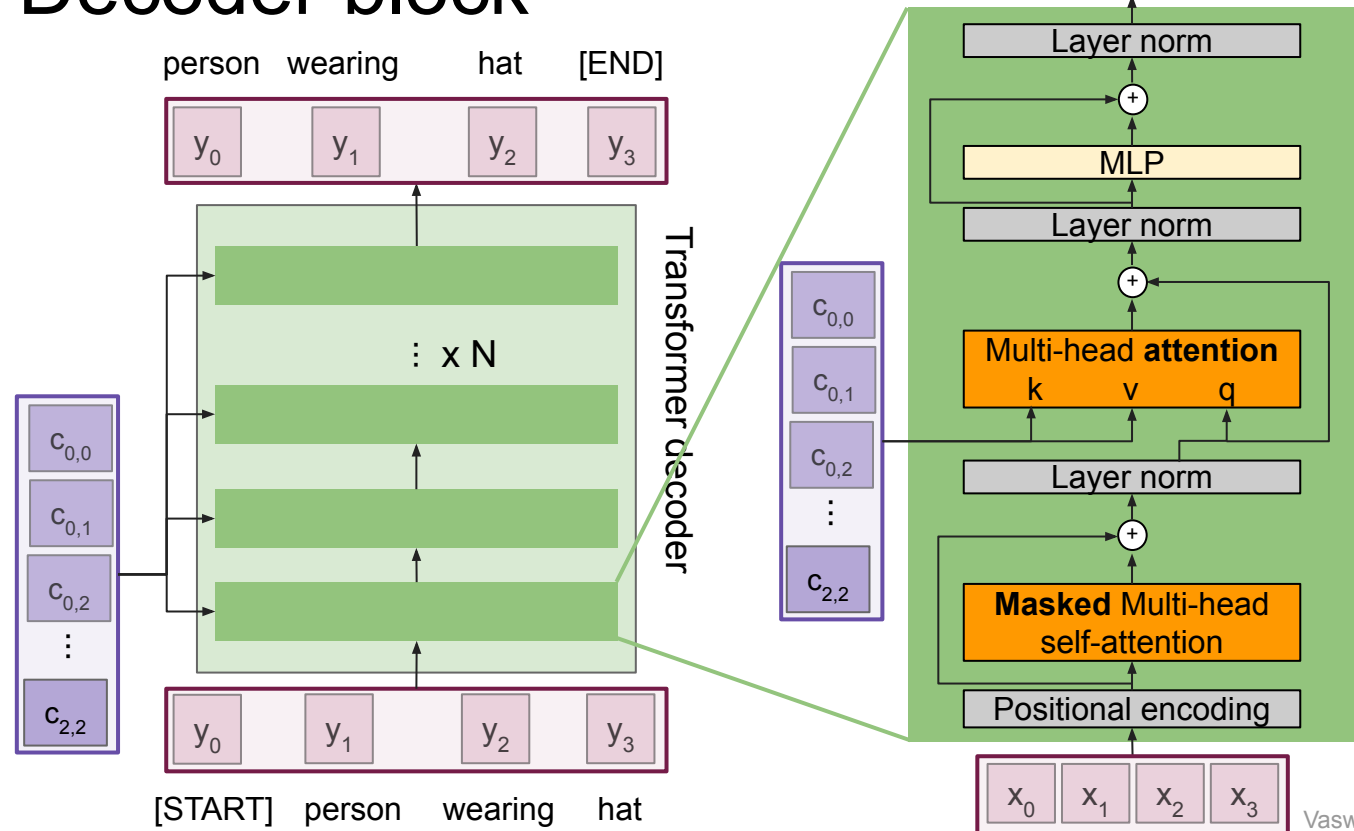
# The Transformer Decoder block



Multi-head attention block attends over the transformer encoder outputs.

For image captioning, this is how we inject image features into the decoder.

Vaswani et al, "Attention is all you need", NeurIPS 2017

# The Transformer Decoder block



**Transformer Decoder Block:**

**Inputs**: Set of vectors **x** and Set of context vectors **c**.
**Outputs**: Set of vectors **y**.

Masked Self-attention only interacts with past inputs.

Multi-head attention block is NOT self-attention. It attends over encoder outputs.

Highly scalable, highly parallelizable, but high memory usage.

Vaswani et al, "Attention is all you need", NeurIPS 2017

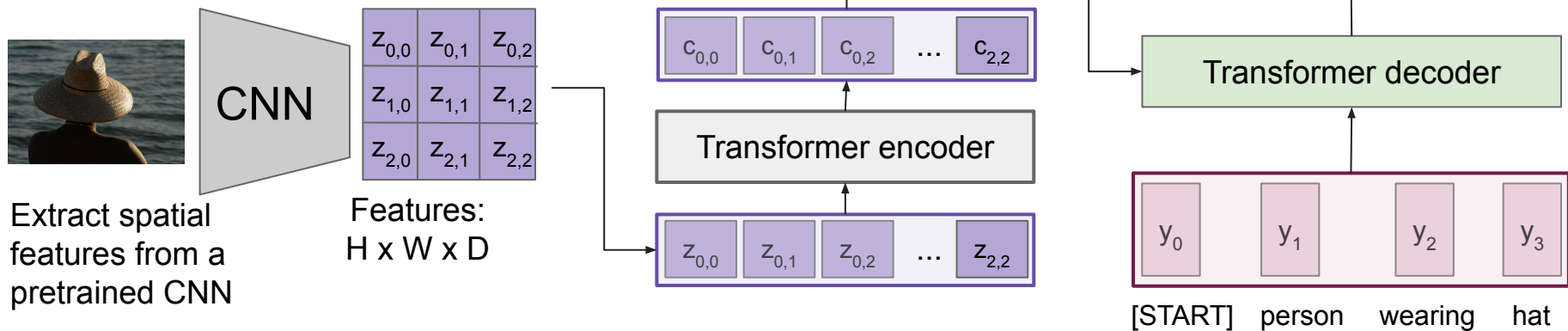# Image Captioning using transformers

- **No recurrence at all**

# Image Captioning using transformers

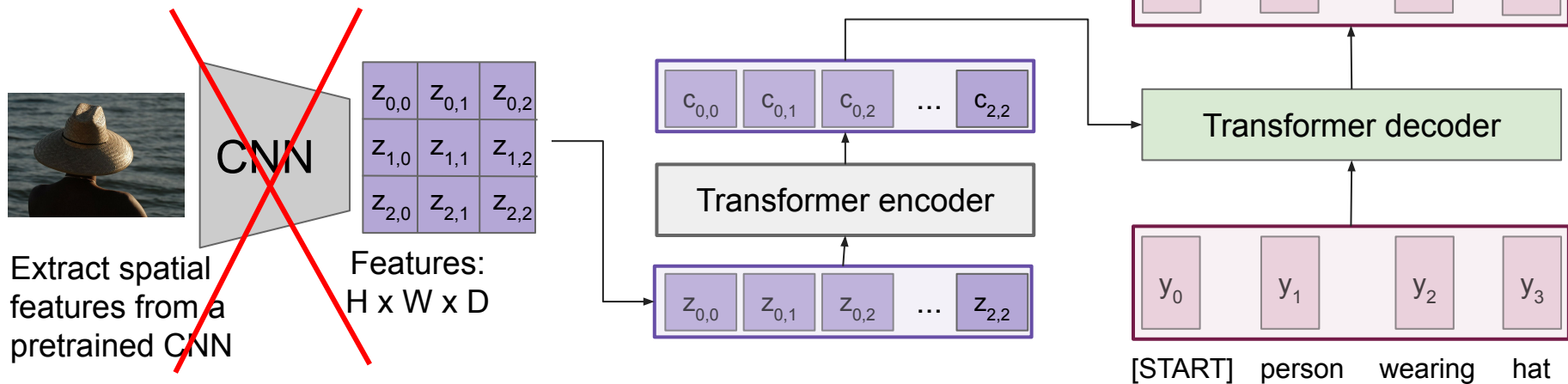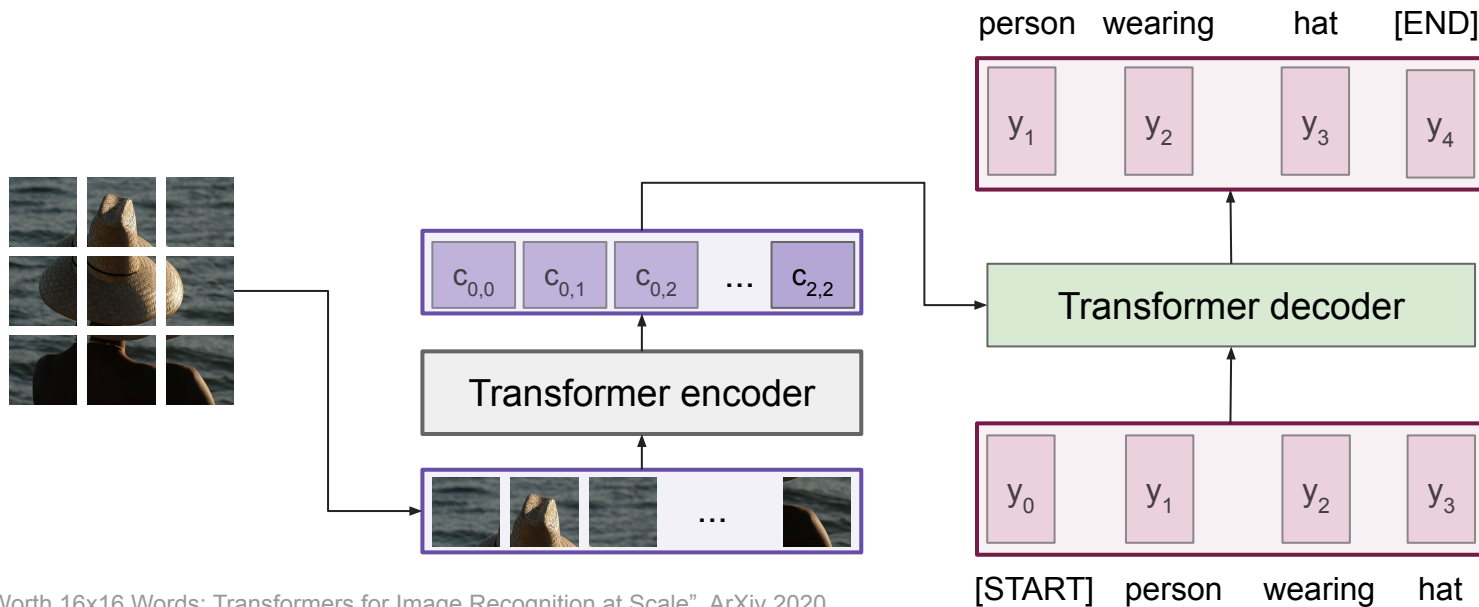- **Perhaps we don't need convolutions at all?**



Extract spatial features from a pretrained CNN

Features: H x W x D

# Image Captioning using ONLY transformers

- **Transformers from pixels to language**



person wearing hat [END]

Transformer decoder

Transformer encoder

[START] person wearing hat

# Summary

- Adding **attention** to RNNs allows them to "attend" to different parts of the input at every time step
- The **general attention layer** is a new type of layer that can be used to design new neural network architectures
- **Transformers** are a type of layer that uses **self-attention** and layer norm.
  - It is highly **scalable** and highly **parallelizable**
  - **Faster** training, **larger** models, **better** performance across vision and language tasks
  - They are quickly replacing RNNs, LSTMs, and may(?) even replace convolutions.