

DOTE 6635: Artificial Intelligence for Business Research

# Transformers

Renyu (Philip) Zhang

1

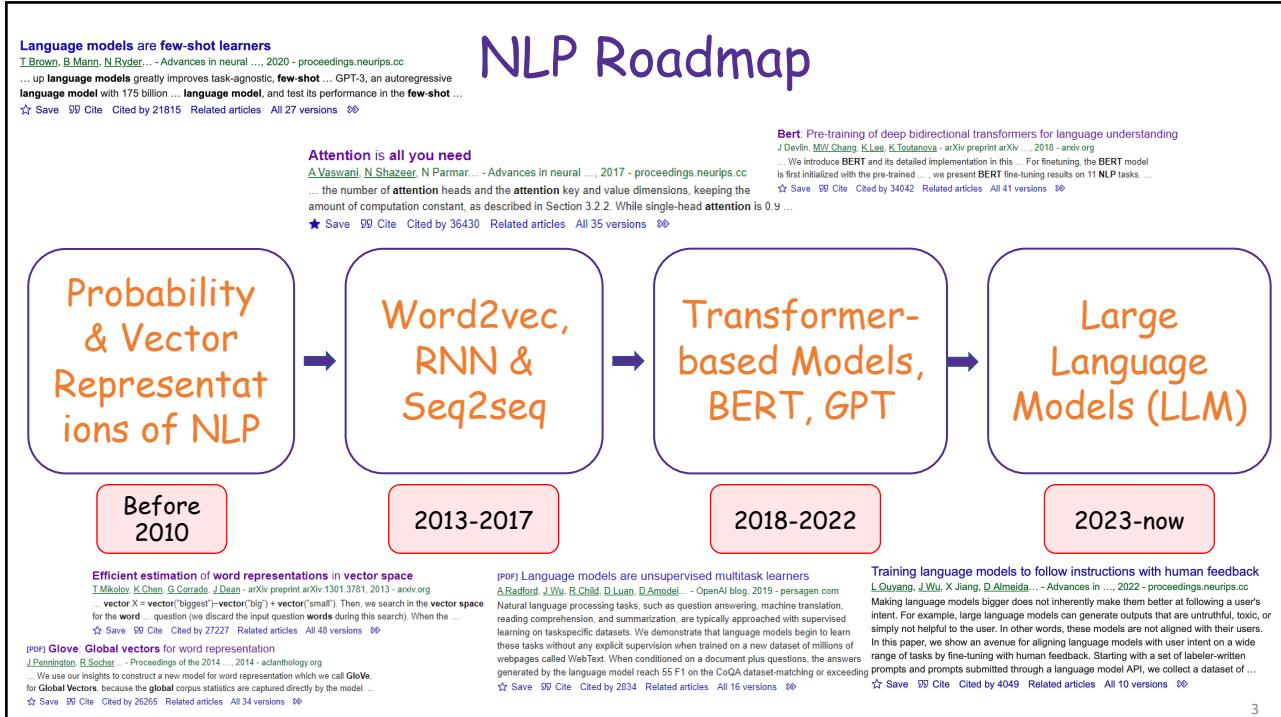
## Agenda

- Sequence-to-Sequence (Seq2seq) and Neural Machine Translation (NMT)
- Attention is All You Need
- Extensions, Adaptations, and Substitutes of Transformers

2

2

1



3

3

NOVEMBER  
27  
2024

## Announcing the NeurIPS 2024 Test of Time Paper Awards

COMMUNICATIONS CHAIRS 2024 / 2021 Conference

By Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub Tomczak, Cheng Zhang

We are honored to announce the Test of Time Paper Awards for NeurIPS 2024. This award is intended to recognize papers published 10 years ago at NeurIPS 2014 that have significantly shaped the research field since then, standing the test of time.

This year, we are making an exception to award two Test of Time papers given the undeniable influence of these two papers on the entire field. The awarded papers are:

- Generative Adversarial Nets  
Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, Yoshua Bengio
- Sequence to Sequence Learning with Neural Networks  
Ilya Sutskever, Oriol Vinyals, Quoc V. Le

Sequence to Sequence Learning with Neural Networks has been cited more than 27,000 times as of this blog post. With the current fast advances of large language models and foundation models in general, making a paradigm shift in AI and applications, the field has benefited from the foundation laid by this work. It is the cornerstone work that set the encoder-decoder architecture, inspiring later attention-based improvements leading to today's foundation model research.

Seq2Seq paper has passed the test of time!

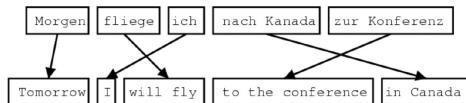
4

2

# Neural Machine Translation (NMT)

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>

- NMT is a way to do machine translation with a single end-to-end neural network: Sequence-to-sequence (**seq2seq**), which involves **2 RNNs**.
- Machine translation is **highly nontrivial** and once was a huge research field in CS and NLP.



1519年600名西班牙人在墨西哥登陆，去征服几百万人口的阿兹特克帝国，初次交锋他们损兵三分之二。

In 1519, six hundred Spaniards landed in Mexico to conquer the Aztec Empire **with a population of a few million**. They lost two thirds of their soldiers in the first clash.

[translate.google.com \(2009\)](#): 1519 600 Spaniards landed in Mexico, **millions of people to conquer the Aztec empire**, the first two-thirds of soldiers against their loss.

[translate.google.com \(2013\)](#): 1519 600 Spaniards landed in Mexico **to conquer the Aztec empire**, hundreds of millions of people, the initial confrontation loss of soldiers two-thirds.

[translate.google.com \(2015\)](#): 1519 600 Spaniards landed in Mexico, **millions of people to conquer the Aztec empire**, the first two-thirds of the loss of soldiers they clash.

5

5

# Seq2Seq for NMT

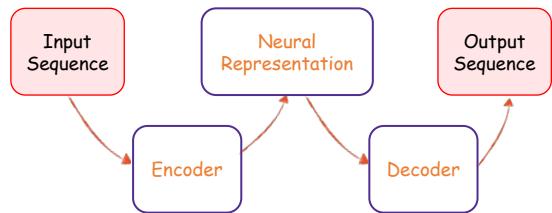
Stanford CS224N: <https://web.stanford.edu/class/cs224n/>

- Seq2seq is a Conditional Language Model:**
  - Predicting the next word of the target sentence  $y$  conditioned on the source sentence  $x$  and prior texts.

$$P(y|x) = P(y_1|x) P(y_2|y_1, x) P(y_3|y_1, y_2, x) \dots P(y_T|y_1, \dots, y_{T-1}, x)$$

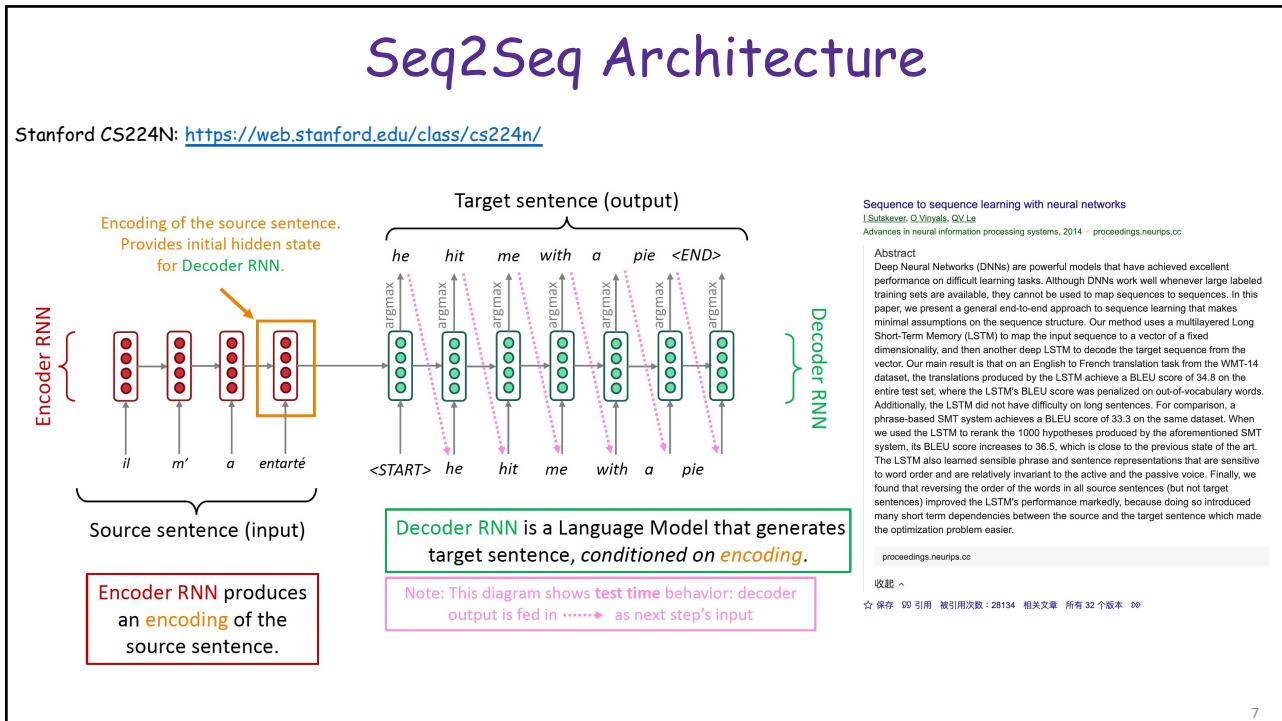
Probability of next target word, given target words so far and source sentence  $x$

- Encoder-decoder architecture:** Encoder takes input and produces a **neural representation**; Decoder produces output based on that **neural representation**.
  - Seq2seq:** both input and output are **sequences**.
  - Summarization:** Long text  $\rightarrow$  short text
  - Dialogue:** previous utterances  $\rightarrow$  next utterance
  - Parsing:** Input text  $\rightarrow$  output parse as a sequence
  - Code generation:** Natural language  $\rightarrow$  Python code

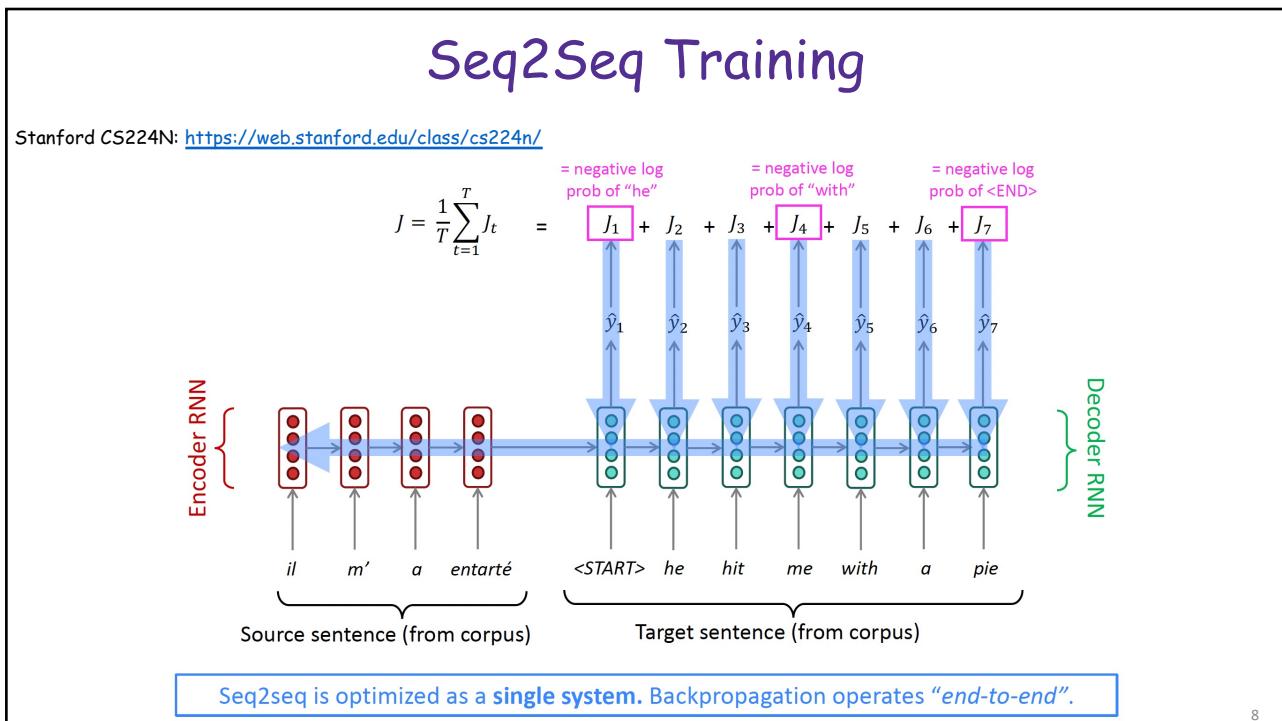


6

6



7



8

## NMT: The First Major Success of DL-NLP

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>

- NMT transformed from a mere research attempt in 2014 to the leading standard in 2016.
- 2014: The Seq2seq paper.
- 2016: Adopted by Google Translate.
- 2018: Adopted by everyone.



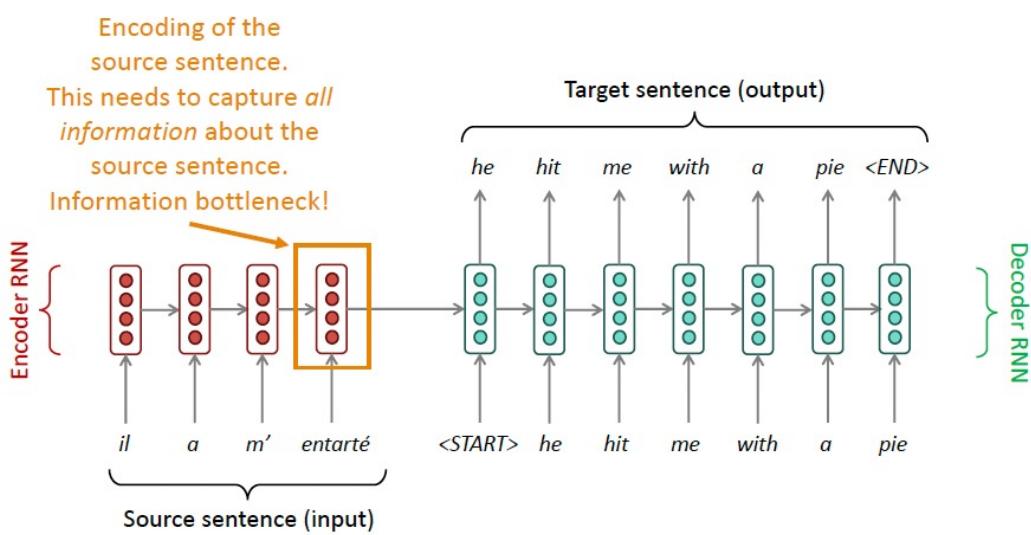
- The original statistical machine translation (SMT) system, built by hundreds of engineers over many years, soon outperformed by NMT trained by small groups of engineers in a few months.

9

9

## Information Bottleneck in RNN

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>



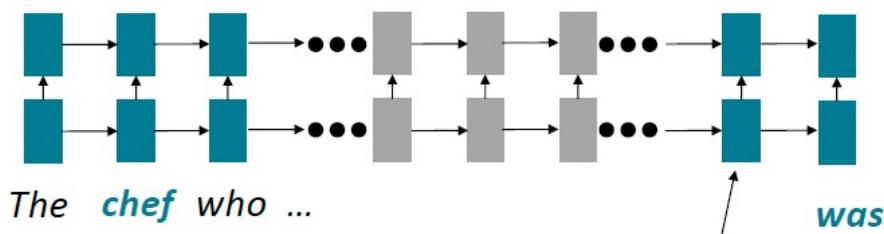
10

10

## Issue with RNN: Linear Interaction Distance

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>

- Human languages are intrinsically NOT linearly ordered.



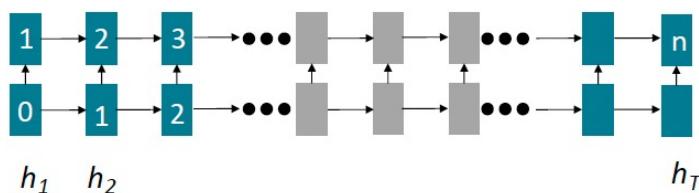
11

11

## Issue with RNN: Non-parallelizability

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>

- Forward and backward passes both have  $O(\text{sequence length})$  unparallelizable operations.
- GPUs can perform independent small computations quickly in a large scale.
- Future hidden states cannot be computed (in full) before past RNN hidden states have been computed.
- Cannot scale with a very large dataset.



Numbers indicate min # of steps before a state can be computed

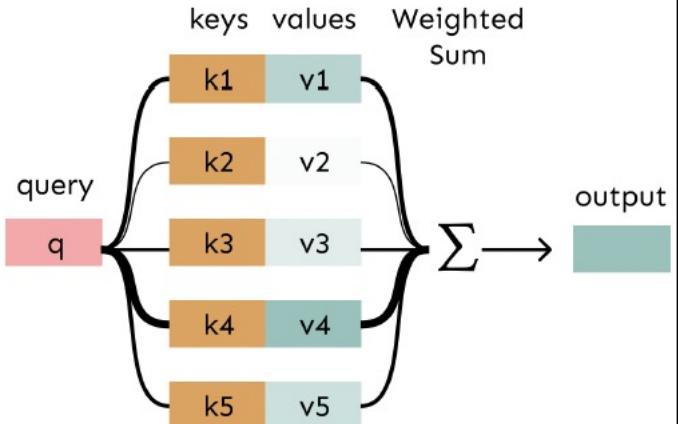
12

12

## Attention as a Very General DL Technique

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>

- **Attention:** Given a set of vector values and a vector of query, attention is a technique to compute a weighted sum of the values dependent on the query.
  - The weighted sum is a selective summary of the information contained in the values, where the query determines which values to focus on.
  - A fixed-size representation of an arbitrary set of representations (values), dependent on some other representation (query).
- In seq2seq + attention, each decoder hidden state (query) attends to all the encoder hidden states (values)..



13

13

## A Family of Attention Models

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>

Name	Alignment score function	Citation
Content-base attention	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \text{cosine}[\mathbf{s}_t, \mathbf{h}_i]$	<a href="#">Graves2014</a>
Additive(*)	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{v}_a^\top \tanh(\mathbf{W}_a[\mathbf{s}_t; \mathbf{h}_i])$ $\alpha_{t,i} = \text{softmax}(\mathbf{W}_a \mathbf{s}_t)$	<a href="#">Bahdanau2015</a>
Location-Based	Note: This simplifies the softmax alignment to only depend on the target position.	<a href="#">Luong2015</a>
General	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{s}_t^\top \mathbf{W}_a \mathbf{h}_i$ where $\mathbf{W}_a$ is a trainable weight matrix in the attention layer.	<a href="#">Luong2015</a>
Dot-Product	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \mathbf{s}_t^\top \mathbf{h}_i$	<a href="#">Luong2015</a>
Scaled Dot-Product(^)	$\text{score}(\mathbf{s}_t, \mathbf{h}_i) = \frac{\mathbf{s}_t^\top \mathbf{h}_i}{\sqrt{n}}$ Note: very similar to the dot-product attention except for a scaling factor; where n is the dimension of the source hidden state.	<a href="#">Vaswani2017</a>

14

14

## Agenda

- Sequence-to-Sequence (Seq2seq) and Neural Machine Translation (NMT)
- Attention is All You Need
- Extensions, Adaptations, and Substitutes of Transformers

15

15

## Attention is All You Need

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>  
<https://peterbloem.nl/blog/transfomers>

- Transformer: No RNN architecture, just attention mechanism.
- Self-attention: To generate  $y_t$ , we need to pay attention to  $y_{\leq t}$ .

$$\begin{aligned} w'_{ij} &= \frac{q_i^T k_j}{\sqrt{k}} \\ q_i &= W_q x_i & k_i &= W_k x_i & v_i &= W_v x_i \\ w'_{ij} &= q_i^T k_j & \text{Value} & & & \\ w_{ij} &= \text{softmax}(w'_{ij}) & & & \text{Why does it work?} & \\ y_i &= \sum_j w_{ij} v_j. & & & & \end{aligned}$$

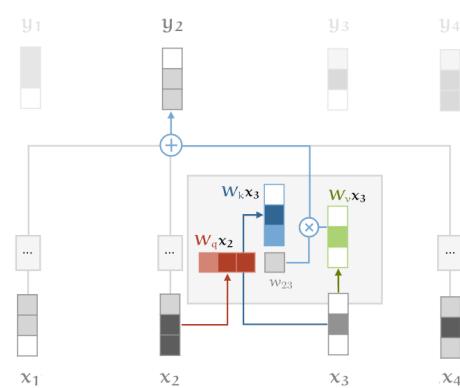


Illustration of the self-attention with key, query and value transformations.

**Attention is all you need**  
[A Vaswani, N Shazeer, N Parmar... - Advances in neural ...](#), 2017 - proceedings.neurips.cc  
 ... to attend to all positions in the decoder up to and including that position. We need to prevent  
 ... We implement this inside of scaled dot-product attention by masking out (setting to  $-\infty$ ) ...  
 ☆ Save ⌂ Cite Cited by 150580 Related articles All 91 versions ➔

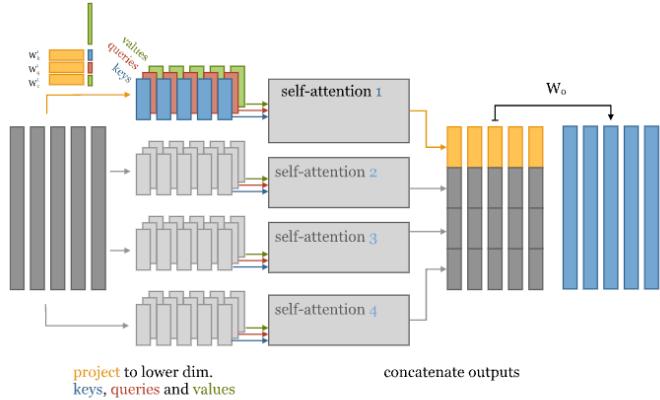
16

16

## Multi-head Attention

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>  
<https://peterbloem.nl/blog/transfomers>

- Multi-head attention is a way to **speed up the training procedure**.
- Instead of using a large matrix to compute all attentions, we can **compute multiple attention matrices and concatenate the final vectors**.
- Allows for **parallel computing**: Deploy attention mechanisms to multiple computing cores in parallel and sum them up at the end.
- Input dim = 256, 8 attention heads, each with 32 dimensions.



The basic idea of multi-head self-attention with 4 heads. To get our **keys**, **queries** and **values**, we project the input down to vector sequences of smaller dimension.

17

17

## Position Encoding

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>  
<https://peterbloem.nl/blog/transfomers>

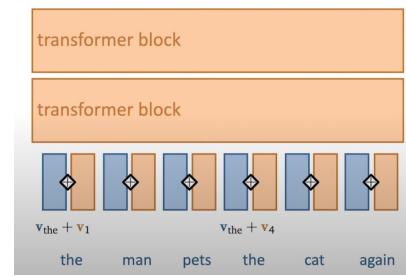
- **Position embeddings**: Position vectors which are learned.
- **Position encoding**: The function from position to vector.
- The final input of the model is the **sum of word embeddings and position embeddings**.

word embeddings:

$v_{\text{the}}$ ,  $v_{\text{man}}$ ,  $v_{\text{pets}}$ ,  $v_{\text{cat}}$ ,  $v_{\text{again}}$

position embeddings:

$v_1$ ,  $v_2$ ,  $v_3$ ,  $v_4$ ,  $v_5$ , ...



18

18

## Auto-Regression

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>  
<https://peterbloem.nl/blog/transfomers>

- Self-supervised learning for transformers.
- To use self-attention in decoders, we need to mask the future.
- Inefficient implementation: Change the set of keys and queries to include only past words.
- Parallelizable implementation: Mask out attention to future words by setting the weight to -inf.

$$w'_{ij} = \begin{cases} q_i^T k_j, j \leq i \\ -\infty, j > i \end{cases}$$

[START]	The	chef	who
[START]	−∞	−∞	−∞
The		−∞	−∞
chef			−∞
who			

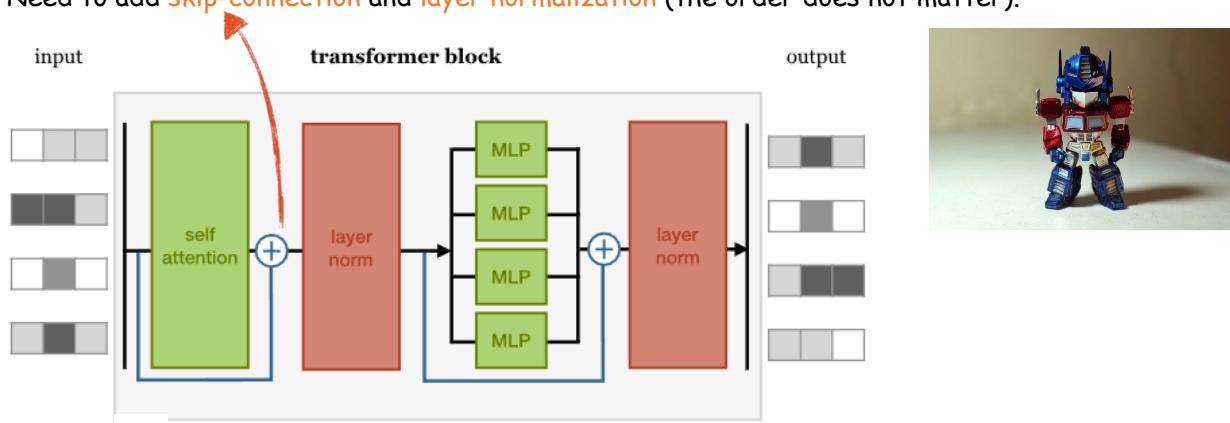
19

19

## Transformer

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>  
<https://peterbloem.nl/blog/transfomers>

- Transformer = Multi-head self-attention + MLP + position encoding + autoregression
- Need to add skip-connection and layer normalization (the order does not matter).



20

20

## Layer Normalization

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>  
<https://peterbloem.nl/blog/transfomers>

- **Layer normalization:** A trick to help models train faster.
- Cut down on uninformative variation in hidden values by normalizing to unit mean and standard deviation within each layer: Normalized gradients.
- Let  $x \in \mathbb{R}^d$  be an individual (word) vector in the model.
- Let  $\mu = \sum_{j=1}^d x_j$ ; this is the mean;  $\mu \in \mathbb{R}$ .
- Let  $\sigma = \sqrt{\frac{1}{d} \sum_{j=1}^d (x_j - \mu)^2}$ ; this is the standard deviation;  $\sigma \in \mathbb{R}$ .
- Let  $\gamma \in \mathbb{R}^d$  and  $\beta \in \mathbb{R}^d$  be learned “gain” and “bias” parameters. (Can omit!)
- Then layer normalization computes:

$$\text{output} = \frac{x - \mu}{\sqrt{\sigma + \epsilon}} * \gamma + \beta$$

Normalize by scalar  
mean and variance      Modulate by learned  
elementwise gain and bias

### Layer normalization

JL Ba, JR Kiros, GE Hinton - arXiv preprint arXiv:1607.06450, 2016 - arxiv.org  
... , we transpose batch **normalization** into **layer normalization** by computing the mean and  
variance used for **normalization** from all of the summed inputs to the neurons in a **layer** on a ...  
☆ Save ⤙ Cite Cited by 10350 Related articles All 6 versions ⚡

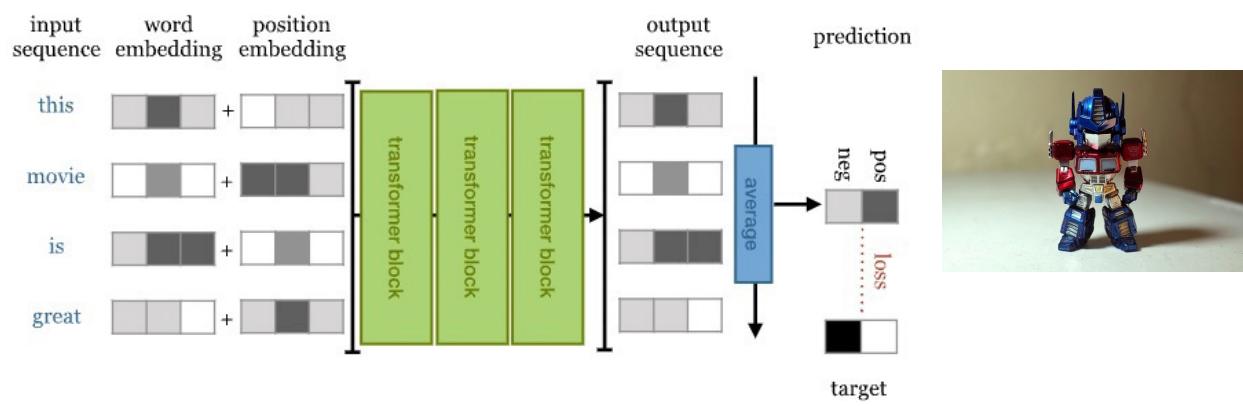
21

21

## Classification Transformer

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>  
<https://peterbloem.nl/blog/transfomers>

- Directly train a classifier on top of a transformer.



22

22

# Attention is All You Need

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>  
<https://peterbloem.nl/blog/transformers>

- Input: Sequence in language one and Sequence in language two.
- Architecture: Encoder + Decoder
- 8 heads, 512 embedding dimensions, 2048 sentence length
- Trained on 8 GPUs for 5 days.

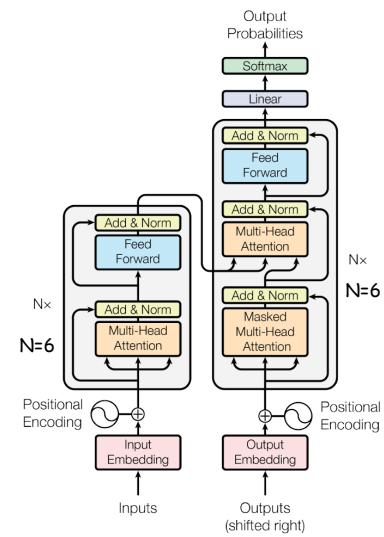


Figure 1: The Transformer - model architecture.

23

23

# Attention is All You Need

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>  
<https://peterbloem.nl/blog/transformers>

## Machine Translation

Model	BLEU		Training Cost (FLOPs)	
	EN-DE	EN-FR	EN-DE	EN-FR
ByteNet [18]	23.75			
Deep-Att + PosUnk [39]		39.2		$1.0 \cdot 10^{20}$
GNMT + RL [38]	24.6	39.92	$2.3 \cdot 10^{19}$	$1.4 \cdot 10^{20}$
ConvS2S [9]	25.16	40.46	$9.6 \cdot 10^{18}$	$1.5 \cdot 10^{20}$
MoE [32]	26.03	40.56	$2.0 \cdot 10^{19}$	$1.2 \cdot 10^{20}$
Deep-Att + PosUnk Ensemble [39]		40.4		$8.0 \cdot 10^{20}$
GNMT + RL Ensemble [38]	26.30	41.16	$1.8 \cdot 10^{20}$	$1.1 \cdot 10^{21}$
ConvS2S Ensemble [9]	26.36	<b>41.29</b>	$7.7 \cdot 10^{19}$	$1.2 \cdot 10^{21}$
Transformer (base model)	27.3	38.1	<b><math>3.3 \cdot 10^{18}</math></b>	
Transformer (big)	<b>28.4</b>	<b>41.8</b>	$2.3 \cdot 10^{19}$	

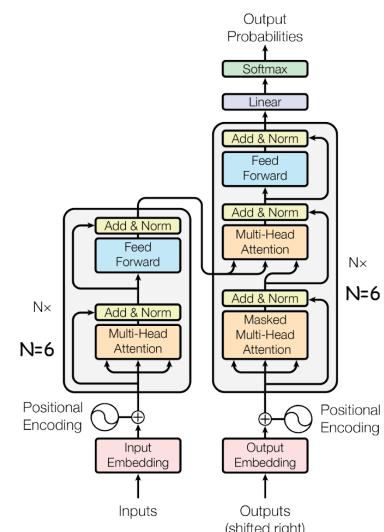


Figure 1: The Transformer - model architecture.

24

24

# Attention is All You Need

Stanford CS224N: <https://web.stanford.edu/class/cs224n/>  
<https://peterbloem.nl/blog/transformers>

**Document Generation**

Model	Test perplexity	ROUGE-L
<i>seq2seq-attention, L = 500</i>	5.04952	12.7
<i>Transformer-ED, L = 500</i>	2.46645	34.2
<i>Transformer-D, L = 4000</i>	2.22216	33.6
<i>Transformer-DMCA, no MoE-layer, L = 11000</i>	2.05159	36.2
<i>Transformer-DMCA, MoE-128, L = 11000</i>	1.92871	37.9
<i>Transformer-DMCA, MoE-256, L = 7500</i>	1.90325	38.8

The old standard      Transformers all the way down.

The parallelizability of transformer enables large-scale pre-training!

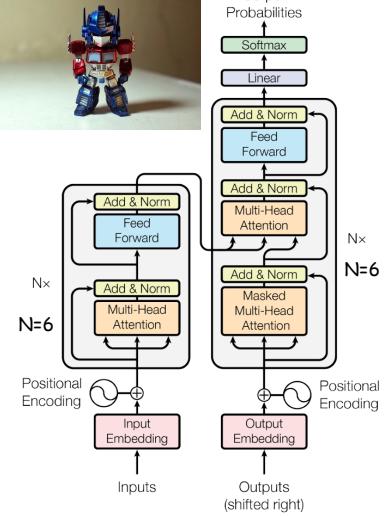



Figure 1: The Transformer - model architecture.

25

25

# Application of Language Model: Detecting FTD

Psychiatry Research 304 (2021) 114135

Contents lists available at ScienceDirect  
**Psychiatry Research**  
journal homepage: [www.elsevier.com/locate/psychres](http://www.elsevier.com/locate/psychres)

Check for updates

Detecting formal thought disorder by deep contextualized word representations

Justyna Sarzynska-Wawer <sup>a,b</sup>, Aleksander Wawer <sup>1,1</sup>, Aleksandra Pawlak <sup>2,c</sup>,  
Julia Szymanowska <sup>2,c</sup>, Izabela Stefanik <sup>d</sup>, Michał Jarkiewicz <sup>3,d</sup>, Łukasz Okruszek <sup>a</sup>

<sup>a</sup> Institute of Psychology, Polish Academy of Sciences, Warsaw, Poland  
<sup>b</sup> Institute of Computer Science, Polish Academy of Sciences, Warsaw, Poland  
<sup>c</sup> University of Social Sciences and Humanities, Chodkiewicza 19/31, 03-615 Warsaw, Poland  
<sup>d</sup> Institute of Psychiatry and Neurology, Szczęśliwego 9, 02-957 Warsaw, Poland

**ARTICLE INFO**

**Keywords:** Schizophrenia; Language; Natural language processing; Deep learning

**ABSTRACT**

Computational linguistics has enabled the introduction of objective tools that measure some of the symptoms of schizophrenia, including the coherence of speech associated with formal thought disorder (FTD). Our goal was to investigate whether neural network based utterance embeddings are more accurate in detecting FTD than models based on individual indicators. We used six questions from the Scale for the Assessment of Thought, Language and Communication (TLC). Using all six TLC questions, the ELMo obtained an accuracy of 80% in distinguishing patients from healthy people. Previously used coherence models were less accurate at 70%. The classifying clinician was accurate 74% of the time. Our analysis shows that both ELMo and TLC are sensitive to the presence of formal thought in patients. Our results show that methods using text representations from language models were more accurate than those based solely on the assessment of FTD, and can be used as measures of disordered language that complement human clinical ratings.

**• How to detect formal thought disorder (FTD)?**

**• Embeddings from LSTM language models (ELMo) can more accurately detect/predict FTD than individual indicators (benchmark: coherence model, which can somehow be viewed as traditional NLP method).**

**• Accuracy (N=70, 35 healthy and 35 patients):**

- **ELMo: 80%**
- **Coherence models: <70%**
- **Clinician: 74%**

[HTML] Detecting formal thought disorder by deep contextualized word representations  
J. Sarzynska-Wawer, A. Wawer, A. Pawlak, ... - Psychiatry ..., 2021 - Elsevier

Computational linguistics has enabled the introduction of objective tools that measure some of the symptoms of schizophrenia, including the coherence of speech associated with formal thought disorder (FTD). Our goal was to investigate whether neural network based utterance embeddings are more accurate in detecting FTD than models based on individual indicators. The present research used a comprehensive Embeddings from Language Models (ELMo) approach to represent interviews with patients suffering from schizophrenia ...

☆ Save 99 Cite Cited by 14561 Related articles All 24 versions Web of Science: 74 26

26

13

# Application of Transformer: Polarized Framing of Immigration

**PNAS** RESEARCH ARTICLE | COMPUTER SCIENCES OPEN ACCESS

**Computational analysis of 140 years of US political speeches reveals more positive but increasingly polarized framing of immigration**

Dallas Card                                                                                            <img alt="CrossMark icon" data-bbox="3808 234 3818

## Agenda

- Sequence-to-Sequence (Seq2seq) and Neural Machine Translation (NMT)
- Attention is All You Need
- Extensions, Adaptations, and Substitutes of Transformers

29

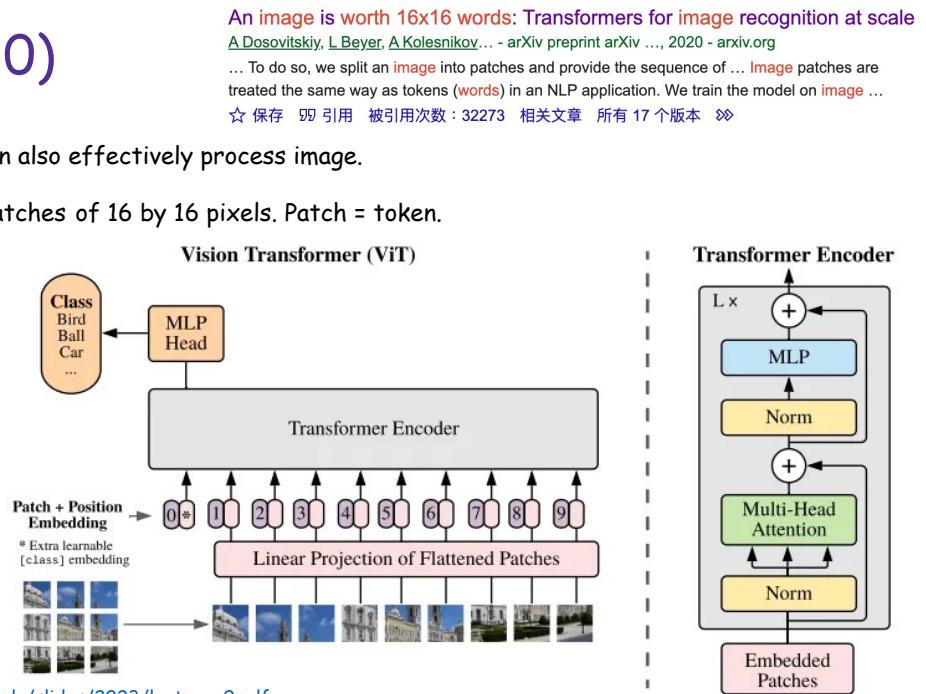
29

## ViT (2020)

- Transformer encoders can also effectively process image.
- Decompose images into patches of 16 by 16 pixels. Patch = token.

Transformer lacks the  
inductive biases of CNN:  
Translation invariance and  
locality.

It may be smart to combine  
CNN and transformers.



30

30

## ViT vs. ResNet

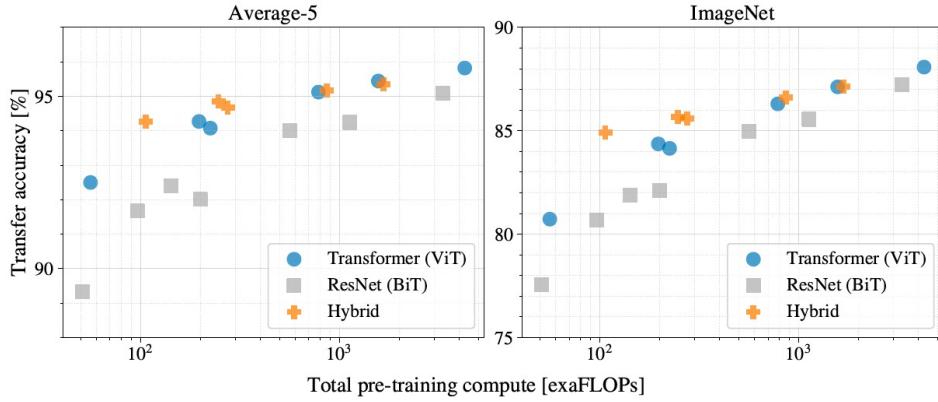


Figure 5: Performance versus pre-training compute for different architectures: Vision Transformers, ResNets, and hybrids. Vision Transformers generally outperform ResNets with the same computational budget. Hybrids improve upon pure Transformers for smaller model sizes, but the gap vanishes for larger models.

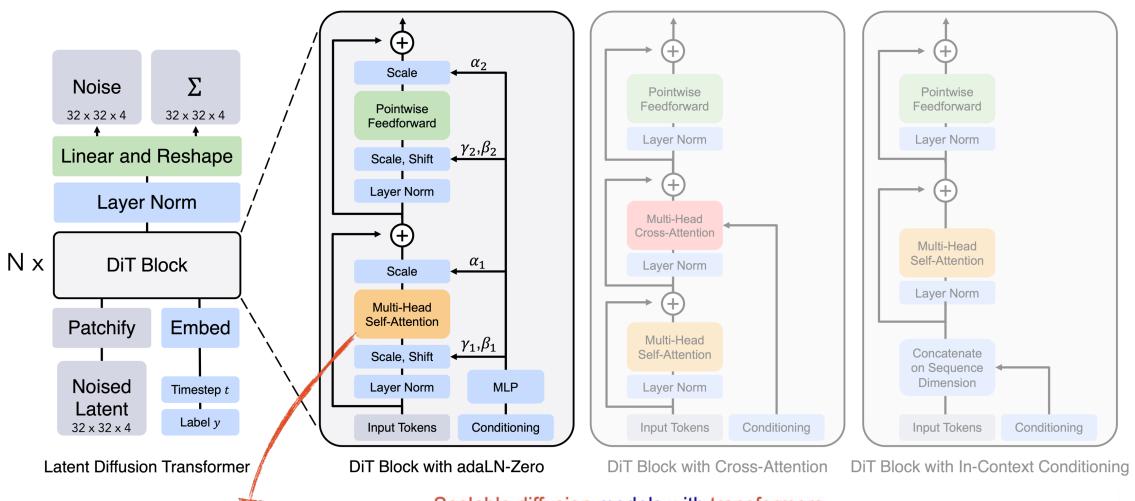
Attention is all you need, again!

Reference: [http://cs231n.stanford.edu/slides/2023/lecture\\_9.pdf](http://cs231n.stanford.edu/slides/2023/lecture_9.pdf)

31

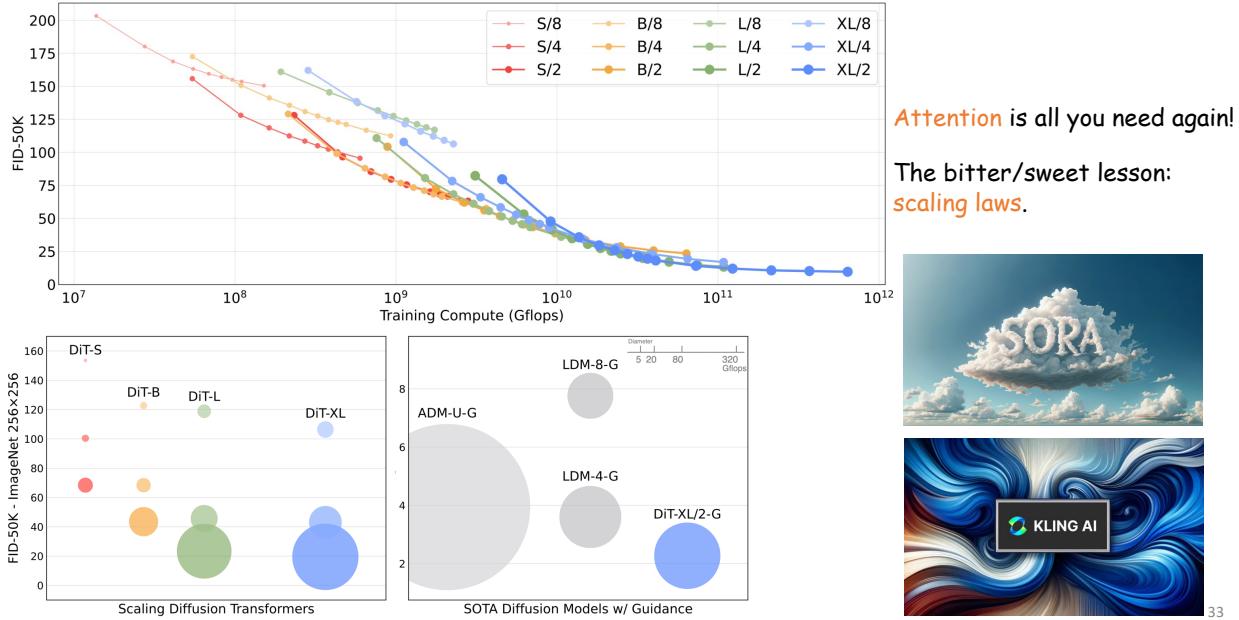
31

## Diffusion Transformer (DiT)



32

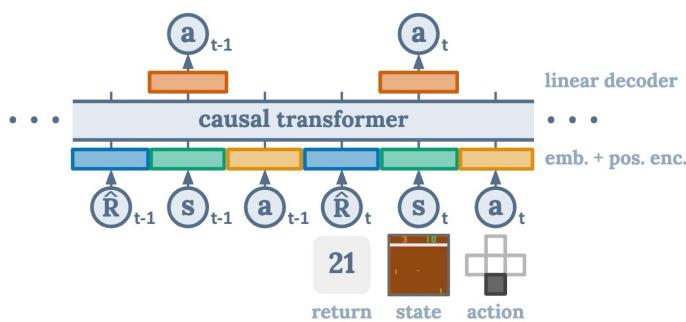
## Diffusion Transformer (DiT) Performances



33

33

## Decision Transformer



- Mainstream RL: Learn value function or policy gradient.
- Decision transformer: Directly learn the next action based on an auto-regressive model.
- Multi-modal transformer: (reward to go, state, action)

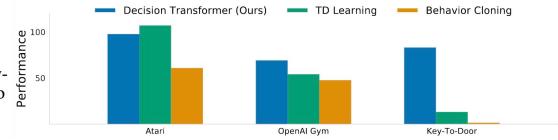


Figure 1: Decision Transformer architecture<sup>1</sup>. States, actions, and returns are fed into modality-specific linear embeddings and a positional episodic timestep encoding is added. Tokens are fed into a GPT architecture which predicts actions autoregressively using a causal self-attention mask.

Reinforcement learning as conditional sequence modeling, like language models

**Decision transformer:** Reinforcement learning via sequence modeling  
[L Chen, K Lu, A Rajeswaran, K Lee... - Advances in neural ...](#), 2021 - proceedings.neurips.cc  
... of the Transformer architecture, and associated advances in language modeling such as GPT-x and BERT. In particular, we present Decision Transformer, ..., Decision Transformer simply ...  
☆ 保存 ⌂ 引用 被引用次数: 1748 相关文章 所有 11 个版本 ☰

34

34

# BC Applied to Inventory Management

**informs**  
<https://pubsonline.informs.org/journal/mnsc>

MANAGEMENT SCIENCE  
 Vol. 69, No. 2, February 2023, pp. 759–773  
 ISSN 0025-1909 (print), ISSN 1526-5501 (online)

## A Practical End-to-End Inventory Management Model with Deep Learning

Meng Qi,<sup>a</sup> Yuanxuan Shi,<sup>b</sup> Yongzhi Qi,<sup>c</sup> Chenxin Ma,<sup>d</sup> Rong Yuan,<sup>d</sup> Di Wu,<sup>d</sup> Zuo-Jun (Max) Shen<sup>a,c</sup>\*

<sup>a</sup>S.C. Johnson College of Business, Cornell University, Ithaca, New York 14853; <sup>b</sup>Department of Electrical and Computer Engineering, University of California-San Diego, San Diego, California 92110; <sup>c</sup>IDecon Smart Supply Chain Y, Mountain View, California 94013; <sup>d</sup>JD.com Supply Valley Research Center, Mountain View, California 94041; <sup>\*</sup>College of Engineering, University of California-Berkeley, Berkeley, California 94720; Faculty of Engineering & Faculty of Business and Economics, University of Hong Kong, Pokfulam, Hong Kong

\*Corresponding author:  
 Contact: mq6@cornell.edu; <https://doi.org/10.1287/mnsc.2022.4564> (Yuanxuan Shi); <https://doi.org/10.1287/mnsc.2022.4564> (Yongzhi Qi); <https://doi.org/10.1287/mnsc.2022.4564> (Chenxin Ma); <https://doi.org/10.1287/mnsc.2022.4564> (Rong Yuan); <https://doi.org/10.1287/mnsc.2022.4564> (Di Wu); <https://doi.org/10.1287/mnsc.2022.4564> (Zuo-Jun (Max) Shen)

Received: June 2, 2020  
 Revised: November 6, 2020  
 Accepted: November 23, 2020  
 Published Online in Articles in Advance: December 15, 2020

<https://doi.org/10.1287/mnsc.2022.4564>  
 Copyright © 2022 INFORMS

Keywords: end-to-end decision-making • inventory management • deep learning • e-commerce

**Abstract.** We investigate a data-driven multiperiod inventory replenishment problem with uncertain demand and vendor lead time (VLT) with accessibility to a large quantity of historical data. Different from the traditional two-step predict-then-optimize (PTO) solution framework, we propose a one-step end-to-end model (E2E) framework that uses deep learning models to output the suggested replenishment amount directly from input features without any intermediate step. The E2E model is trained to capture the behavior of the optimal dynamic replenishing solution under the VLT. By conducting a series of thorough experiments using real data from one of the leading e-commerce companies, we demonstrate the advantages of the proposed E2E model over conventional PTO frameworks. We also conduct a field experiment with JD.com, and the results show that our new algorithm outperforms the current e-commerce company's inventory management model substantially compared with JD's current practice. For the supply chain management industry, our E2E model shortens the decision process and provides an automatic inventory management solution with the possibility to generalize and scale. The concept of E2E, which uses the same information set for the ultimate goal, can also be used in practice for other supply chain management circumstances.

History: Accepted by Hamid Nazerzadeh, big data analytics.

Funding: This research was supported by the National Key Research and Development Program of China (Grant 2018YFB170060) and National Natural Science Foundation of China (Grants 7199460 and 91734210).

Supplemental Material: The online data are available at <https://doi.org/10.1287/mnsc.2022.4564>.

- Use multi-quantile RNNs to provide end-to-end predictions from features to the optimal inventory decisions, whereas most of the literature applies the predict-then-optimize paradigm.
- A field experiment shows that the e2e approach substantially reduces the inventory costs compared with some naïve benchmarks.
- Behavioral cloning (BC) applied to an MDP, with uncertainty completely removed.

### A practical end-to-end inventory management model with deep learning

M. Qi, Y. Shi, Y. Qi, C. Ma, R. Yuan, D. Wu, Z. Shen  
 Management Science, 2023 : pubsonline.informs.org

We investigate a data-driven multiperiod inventory replenishment problem with uncertain demand and vendor lead time (VLT) with accessibility to a large quantity of historical data. Different from the traditional two-step predict-then-optimize (PTO) solution framework, we propose a one-step end-to-end (E2E) framework that uses deep learning models to output the suggested replenishment amount directly from input features without any intermediate step. The E2E model is trained to capture the behavior of the optimal dynamic replenishing solution under the VLT. By conducting a series of thorough experiments using real data from one of the leading e-commerce companies, we demonstrate the advantages of the proposed E2E model over conventional PTO frameworks. We also conduct a field experiment with JD.com, and the results show that our new algorithm outperforms the current e-commerce company's inventory management model substantially compared with JD's current practice. For the supply chain management industry, our E2E model shortens the decision process and provides an automatic inventory management solution with the possibility to generalize and scale. The concept of E2E, which uses the same information set for the ultimate goal, can also be used in practice for other supply chain management circumstances.

SHOW MORE ▾

Save Cite Cited by 65 Related articles All 4 versions Web of Science: 5 More

35

35

# No Free Lunch: Quadratic Training Cost

- Training (and inference) cost of transformer:  $O(n^2)$ ; cost of training RNN:  $O(n)$ .
  - Inference memory of transformer:  $O(n)$ ; that of RNN:  $O(1)$ .
- **Sparse Attention** reduces the full pairwise attention calculations (such as sliding windows, fixed patterns, etc.) to lower computational complexity to  $O(nk)$ ; e.g., "Generating Long Sequences with Sparse Transformers" (2019)
- **Low-Rank Approximation** approximates the attention matrix using low-rank decomposition to reduce the computational complexity to  $O(nk)$ ; e.g., "Lformer: Self-Attention with Linear Complexity" (2020)
- **Kernel Methods** approximate self-attention calculations to reduce the complexity to  $O(n \log n)$ ; e.g., "Performer: Efficient Transformer with Linear Complexity" (2021)
- **Local Attention** uses sliding window-based attention to reduce the complexity to  $O(nk)$ ; e.g., "Longformer: The Long-Document Transformer" (2020)
- **Mixture of Experts (MoE)** uses multiple expert models and activate only a few at a time, adopted by the SOTA LLMs such as DeepSeek-V3; e.g., "Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sharded Training"
- **Linear recurrence for parallel training** optimizes RNN architecture with parallelization; e.g., Mamba: Linear-Time Sequence Modeling with Selective State Spaces (2023)
  - Empirically, not as competitive as transformers for in-context learning.

36

36

18