# TransforLearn: Interactive Visual Tutorial for the Transformer Model

Lin Gao[1], Zekai Shao[1], Ziqin Luo[1], Haibo Hu[2], Cagatay Turkay[3], Siming Chen[1,4]

[1] School of Data Science, Fudan University

[2] School of Big Data & Software Engineering, Chongqing University

[3] Centre for Interdisciplinary Methodologies, University of Warwick
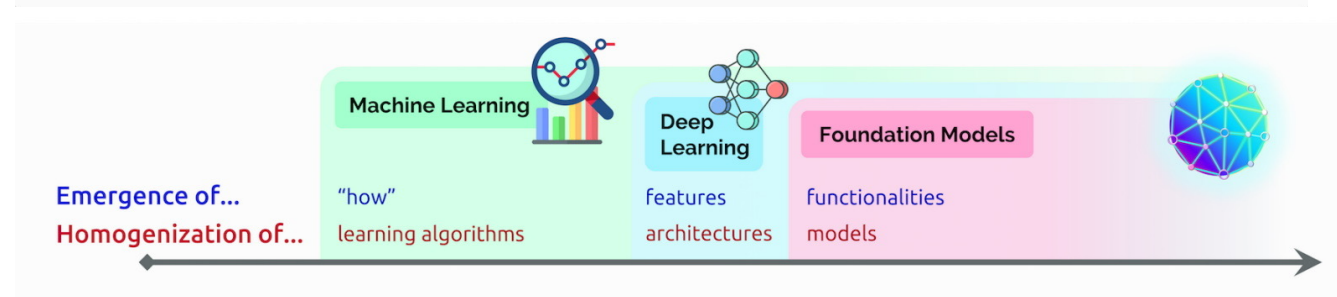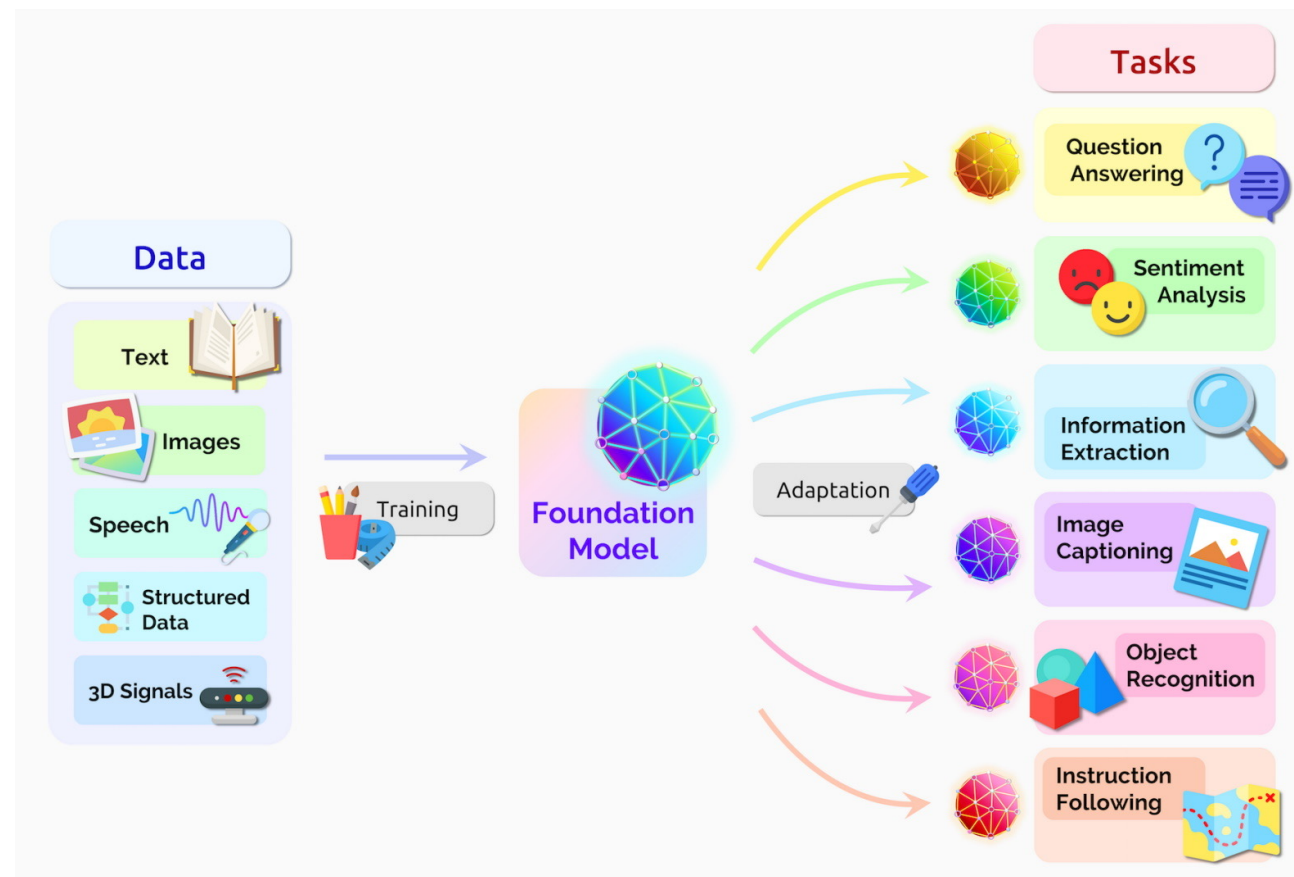
[4] Shanghai Key Laboratory of Data Science

# Background

Transformers are already used with many data sources for applications.

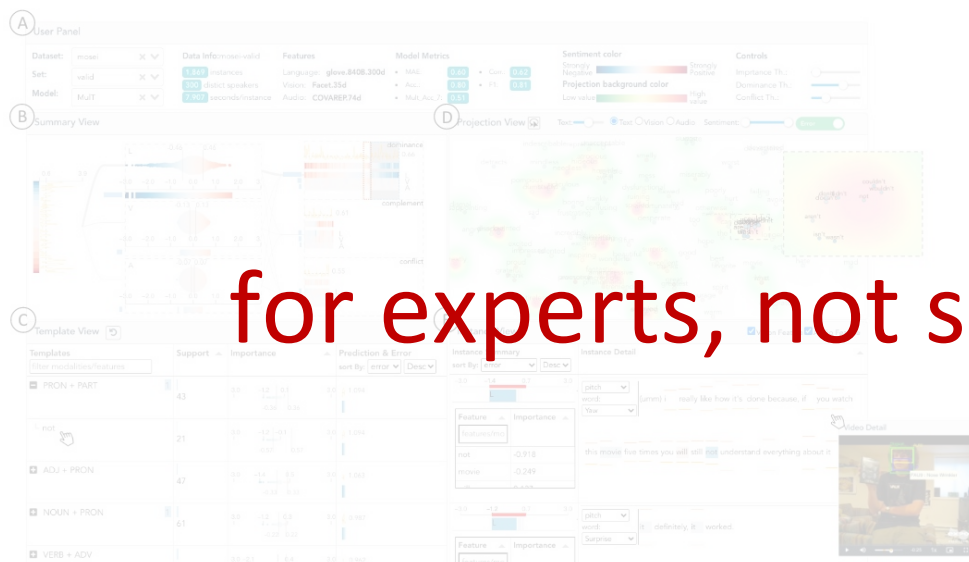Transformers mark the next stage of AI's development, what some call the era of transformer AI.

The popularity of Transformer has sparked significant interest in learning its working mechanisms.



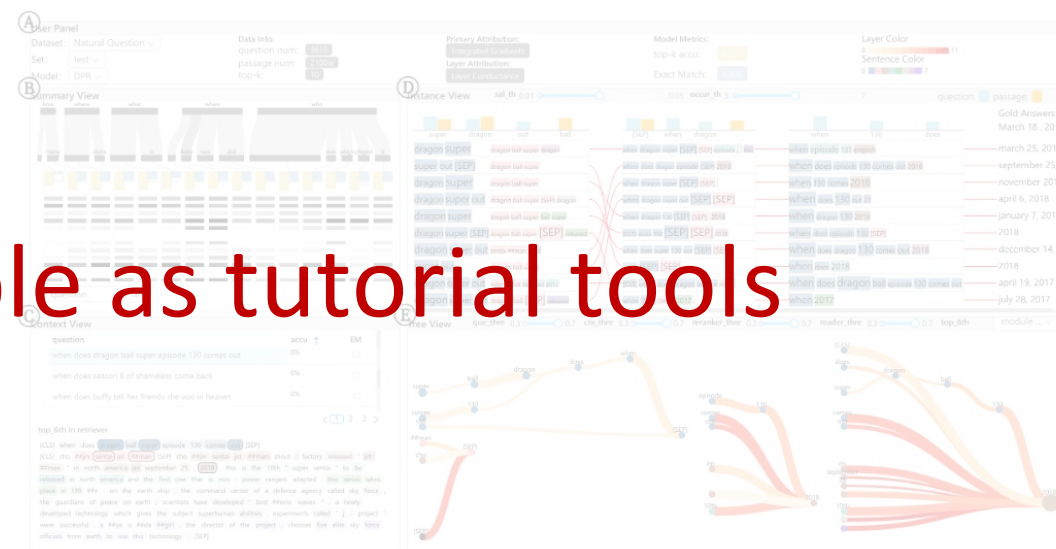[1] https://blogs.nvidia.com/blog/what-is-a-transformer-model/

# Background

## Visualization for understanding deep learning models

- how the models make decisions & what they learned

- model improvement & debugging



for experts, not suitable as tutorial tools
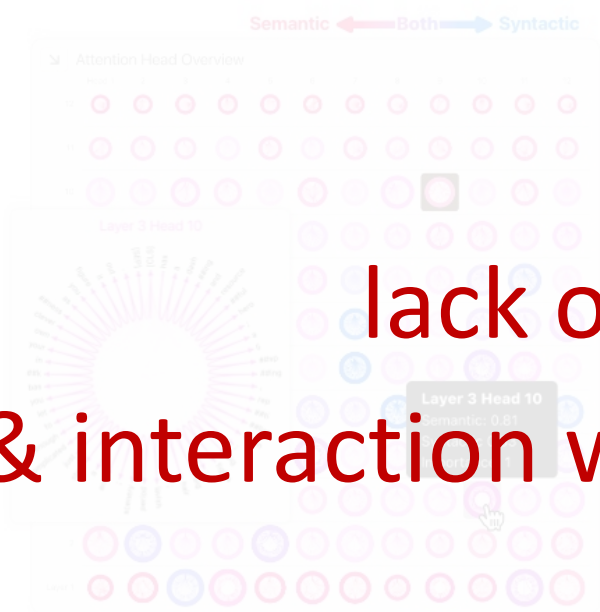
M2lens [1] (TVCG 2021)

VEQA[2] (TVCG 2023)

[1] Wang X, He J, Jin Z, et al. M2lens: Visualizing and explaining multimodal models for sentiment analysis[J]. IEEE Transactions on Visualization and Computer Graphics, 2021, 28(1): 802-812.
[2] Shao Z, Sun S, Zhao Y, et al. Visual Explanation for Open-domain Question Answering with BERT[J]. IEEE Transactions on Visualization and Computer Graphics, 2023.

Visual interpretation of Transformers

- interpretation of embedding and attention mechanisms

- blogs & videos for tutorial



lack of mathematical details
& interaction with the actual data flow or task
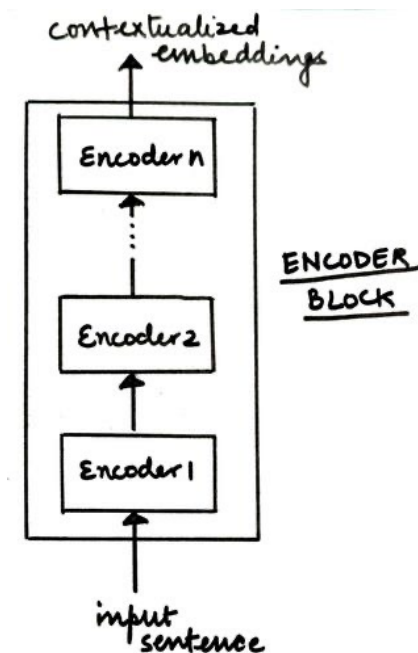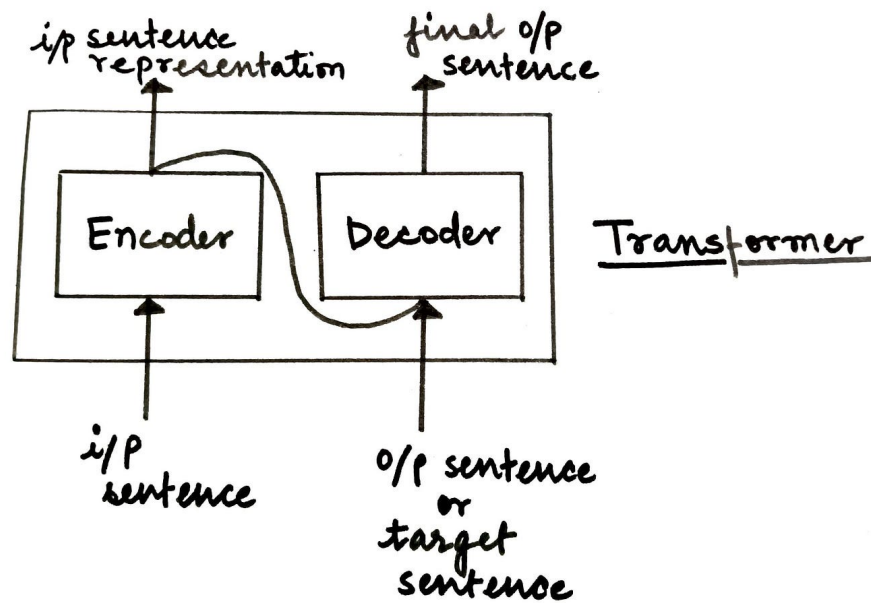
Dodrio[1] (ACL 2021)

Jalammar's blogs [2]

[1] Wang Z J, Turko R, Chau D H. Dodrio: Exploring transformer models with interactive visualization[J]. arXiv preprint arXiv:2103.14625, 2021.
[2] https://jalammar.github.io/a-visual-guide-to-using-bert-for-the-first-time/

# Preliminary Study

For lecturers, they need to manually break down Transformer into multiple steps and discuss them in a sequence of slides.
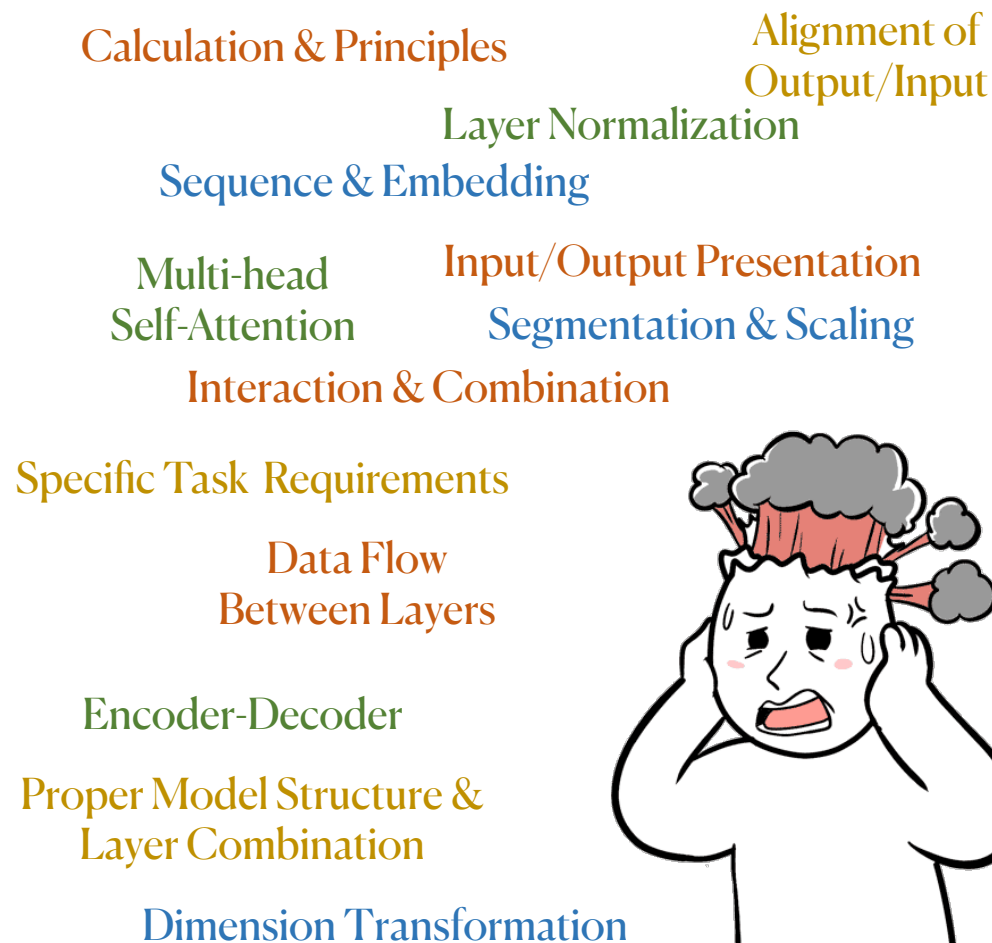
- Theoretical learning -> Dynamic thinking combined with practice

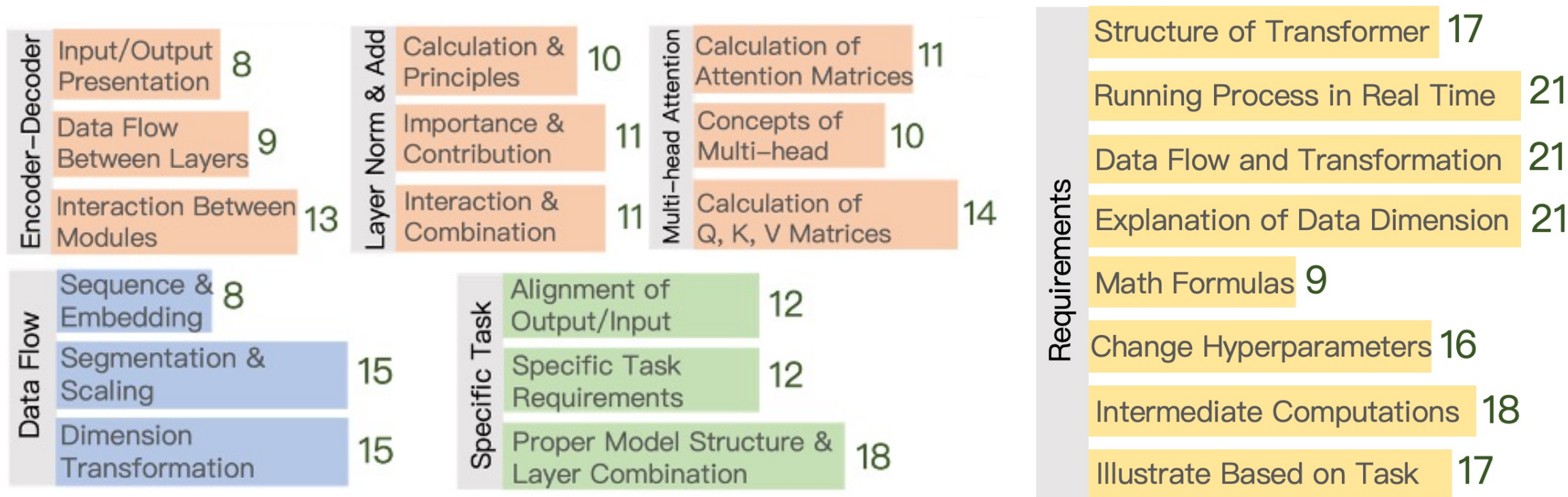- Contextual interactions

- Class engagement

# Preliminary Study

Beginners face difficulties in comprehending and learning Transformers due to its complex structure, data transformation and abstract downstream task.

- Encoder/Decoder, Attention ……
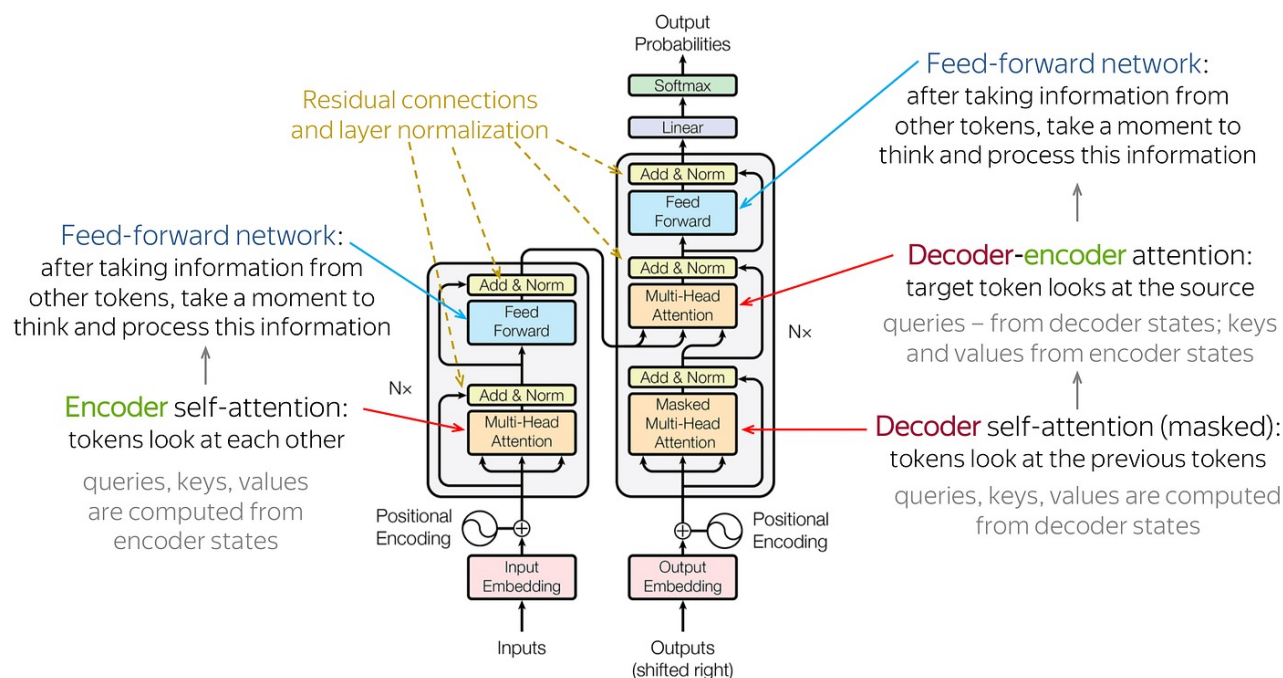
- Embedding, Dimension ……

- Alignment, Process ……

Calculation & Principles

Alignment of Output/Input

Layer Normalization

Sequence & Embedding

Multi-head Self-Attention

Input/Output Presentation

Segmentation & Scaling

Interaction & Combination

Specific Task Requirements

Data Flow Between Layers

Encoder-Decoder

Proper Model Structure & Layer Combination

Dimension Transformation

The survey asked about the key challenges in learning and applying Transformers from various aspects, and what features would be helpful in an interactive tool for beginners.

Consequently, an interactive visual tutorial is needed for deep learning beginners and non-experts to comprehensively learn about Transformers.
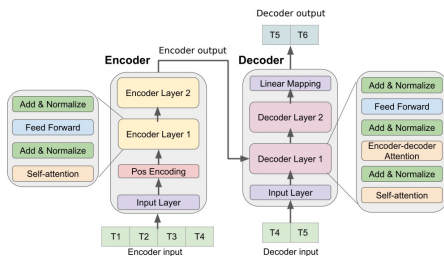
What can TransforLearn do?



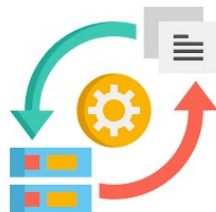Why does Transformer has such a complex architecture[1]

[1] https://stats.stackexchange.com/questions/512242/why-does-transformer-has-such-a-complex-architecture

# Tasks & Requirements

| Task-1 | Task-2 | Task-3 | Task-4 |
|---|---|---|---|



complex structure & layer operations

data flow & transformation

practical use in downstream tasks

guidance & feedback

| Requirement-1 | Requirement-2 | Requirement-3 | Requirement-4 |
|---|---|---|---|

A visual summary of the model architecture and data flow.

An interactive interface for layer operations and mathematical formulas.

Exploration mode between module levels based on downstream tasks.

Self-directed and immersive learning experiences.
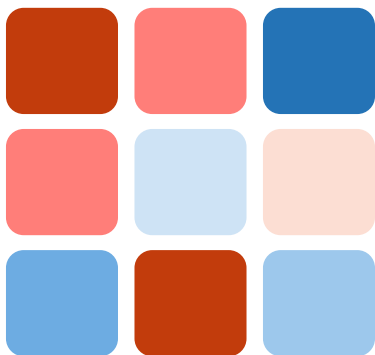
# Visual Design - Overview

**Architecture Overview** ⟹ **Module Detailed Views**

$$\begin{bmatrix} 11 & 20 & 109 \\ 21 & 54 & 37 \\ 74 & 11 & 60 \end{bmatrix}$$

Sequence Data
Parameter Data

⟹ 

Breaking text into individual word segmentations.
Word Segmentations

Mapping words to dense vector representations.
embeddings (4,512)

Index in word token dictionary.
Word Tokens

Add positional information to original embeddings.
after positional encoding (4,512)

Add the original *layer normalizations* into *attention scores*

Add the original *layer normalizations* into *feed-forward network results*
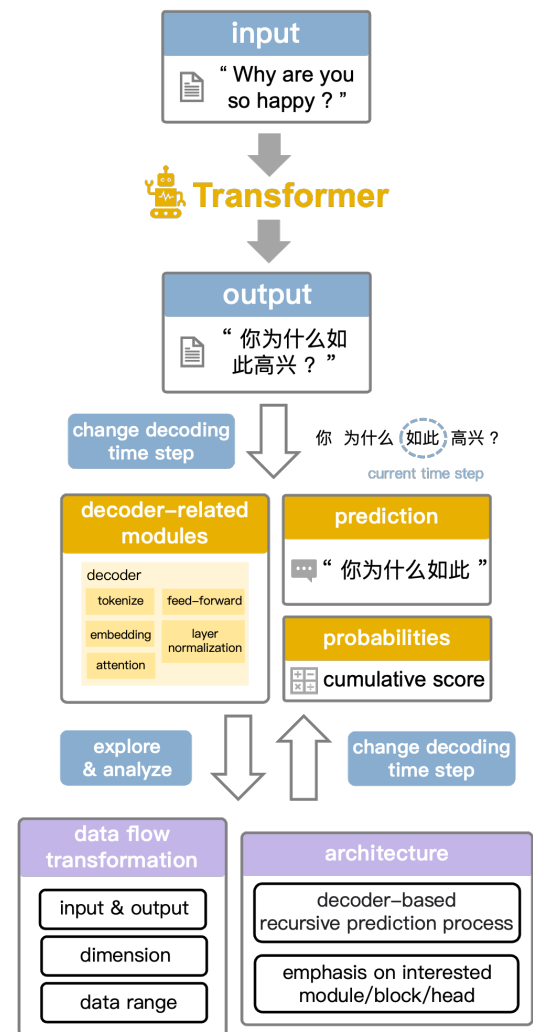
🖑 **Click** to explore each layer operation!
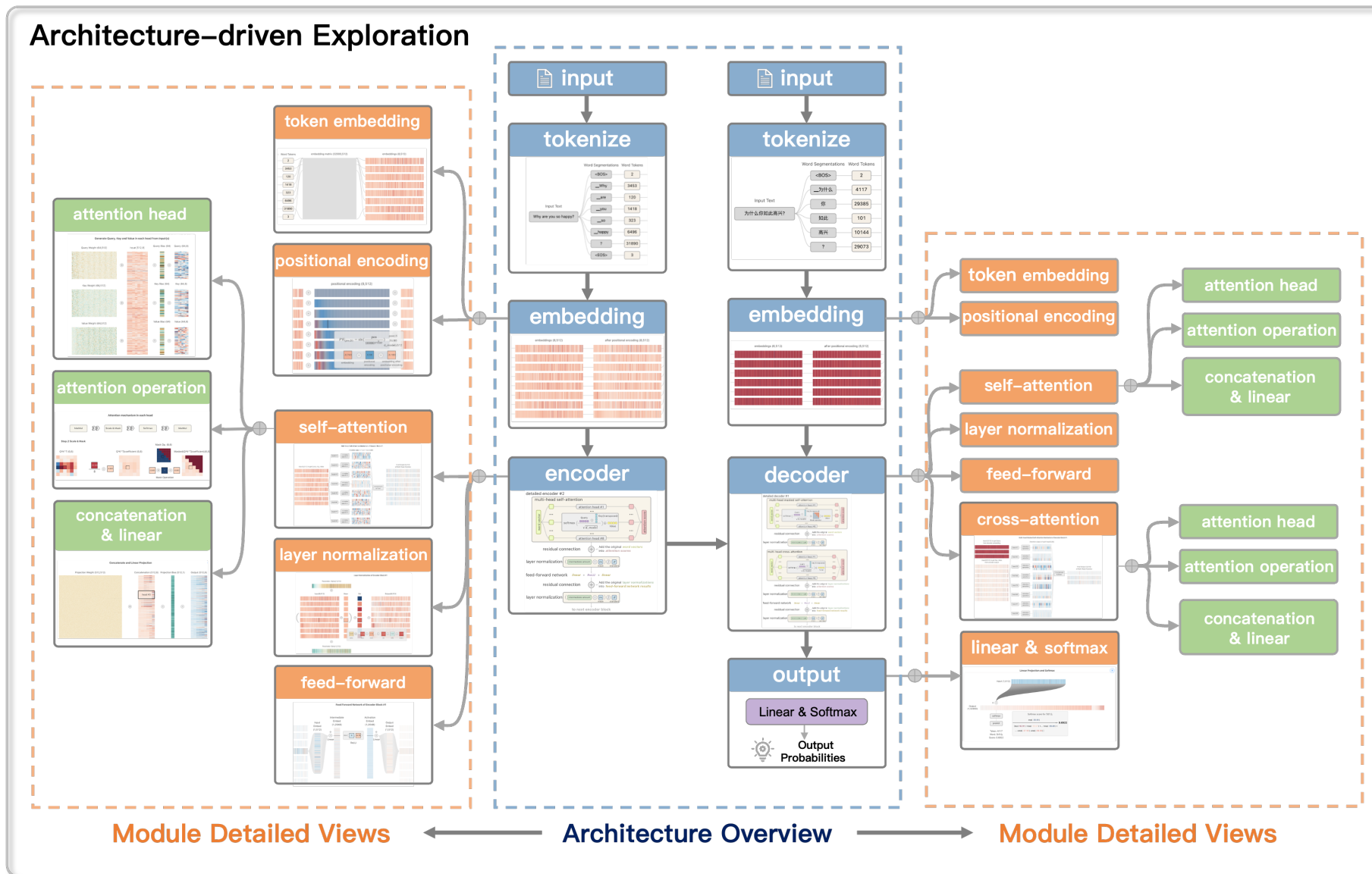
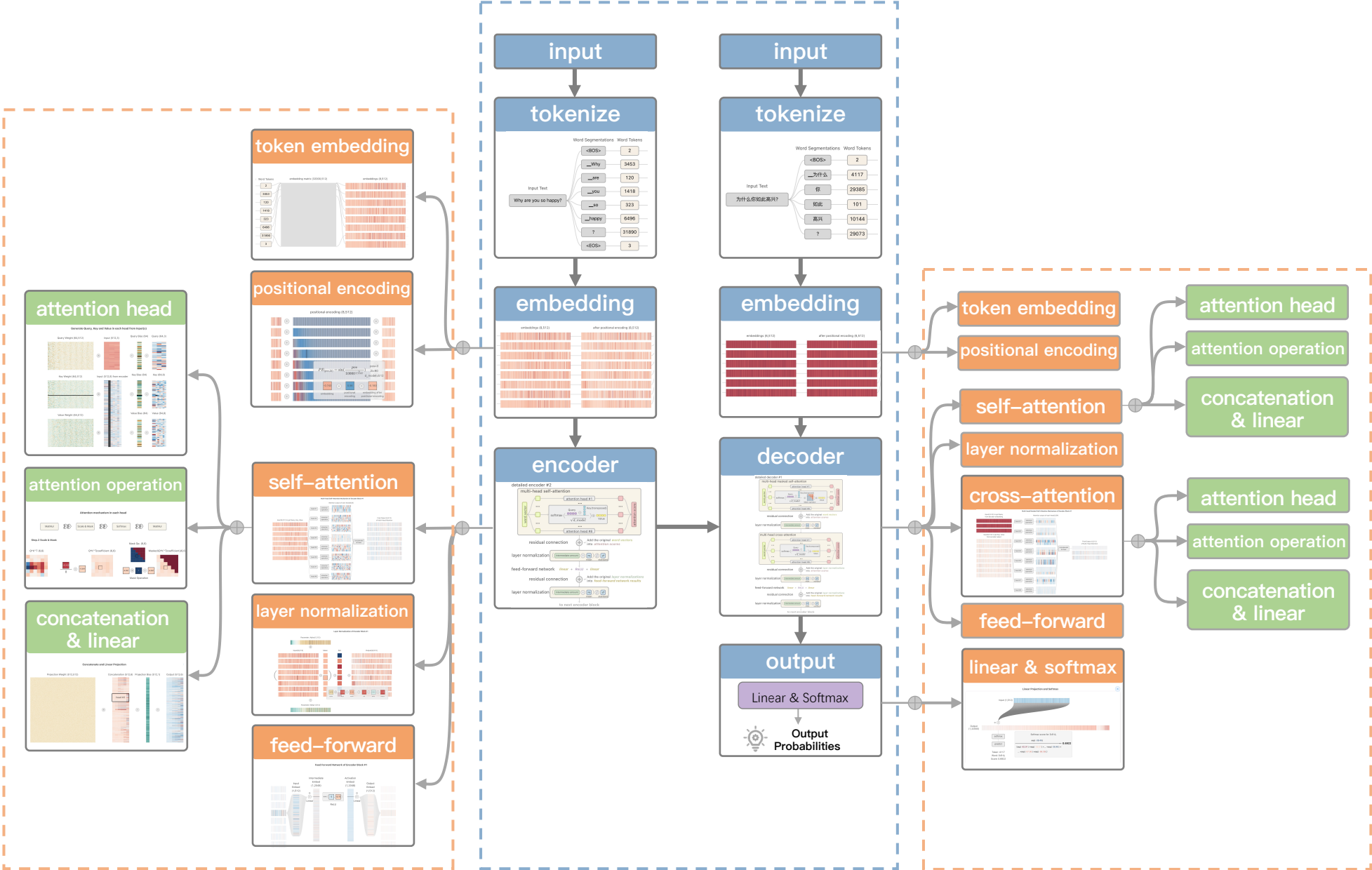🖑 **Click** to learn how to generate Q,K,V!

# Architecture-driven Exploration



Module Detailed Views ← Architecture Overview → Module Detailed Views

# Encoder

Residual connections and layer normalization

Feed-forward network: after taking information from other tokens, take a moment to think and process this information

Encoder self-attention: tokens look at each other

queries, keys, values are computed from encoder states



Encoder Block

# Encoder

inputs

"Why are you so happy?"

tokenize

embedding

encoder



TransforLearn: Interactive Visual Tutorial for the Transformer Model

**Input View**

Sentence to be translated:    Why are you so happy?

**Translation View**

Translation:    为什么你如此高兴?

Current translation:    为什么你如此高兴

Prediction in current iteration:    **?**
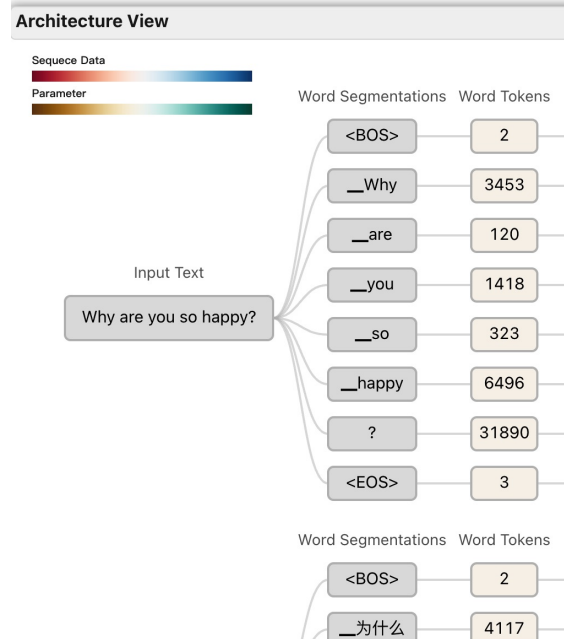
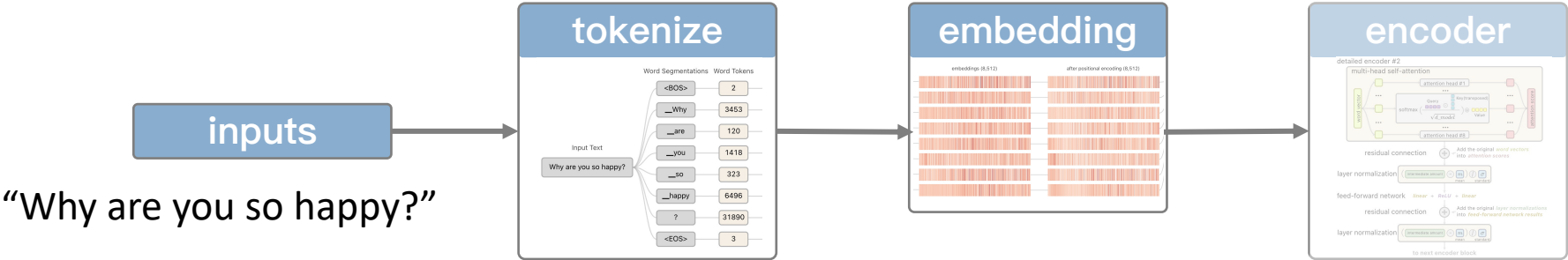Cumulative Score (Probabilities):    **0.039872859081467996**

**Architecture View**

# Encoder

inputs

"Why are you so happy?"

### tokenize

Word Segmentations | Word Tokens

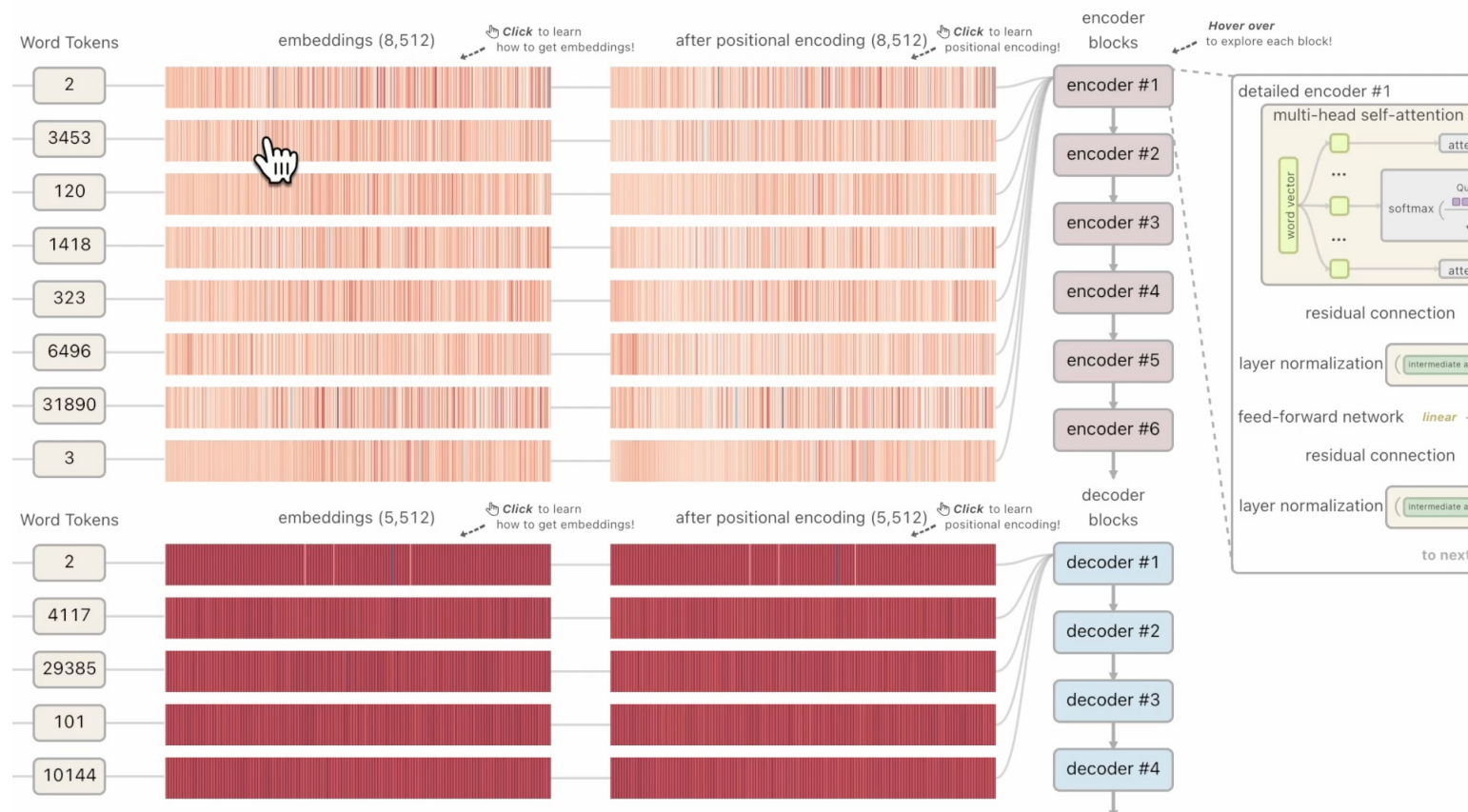| | |
|---|---|
| \<BOS\> | 2 |
| __Why | 3453 |
| __are | 120 |
| __you | 1418 |
| __so | 323 |
| __happy | 6496 |
| ? | 31890 |
| \<EOS\> | 3 |

Input Text
Why are you so happy?

### embedding

embeddings (8,512)    after positional encoding (8,512)

### encoder

detailed encoder #2
multi-head self-attention
attention head #1
softmax
attention head #8
residual connection    Add the original word vectors into attention scores
layer normalization
feed-forward network    linear + relu + linear
residual connection    Add the original layer normalization into feed-forward network results
layer normalization

to next encoder block

---

## TransforLearn: Interactive Visual Tutorial for the Transformer Model

**Input View**

Sentence to be translated:    Why are you so happy?

**Translation View**

Translation:    为什么你如此高兴?

Current translation:    为什么你如此高兴    ‹ ›

Prediction in current iteration:    **?**

Cumulative Score (Probabilities):    **0.039872859081467996**

**Architecture View**

Sequece Data

Parameter

Input Text

Why are you so happy?

# Encoder

inputs

"Why are you so happy?"



tokenize

| Word Segmentations | Word Tokens |
|---|---|
| <BOS> | 2 |
| __Why | 3453 |
| __are | 120 |
| __you | 1418 |
| __so | 323 |
| __happy | 6496 |
| ? | 31890 |
| <EOS> | 3 |

Input Text: Why are you so happy?

embedding

encoder

detailed encoder #2
multi-head self-attention

## TransforLearn: Interactive Visual Tutorial for the Transformer Model

**Input View**

Sentence to be translated: Why are you so happy?

**Translation View**

Translation: 为什么你如此高兴?

Current translation: 为什么你如此高兴

Prediction in current iteration: ?

Cumulative Score (Probabilities): 0.039872859081467996

**Architecture View**

Sequece Data

Parameter

| Word Segmentations | Word Tokens |
|---|---|
| <BOS> | 2 |
| __Why | 3453 |
| __are | 120 |
| __you | 1418 |
| __so | 323 |
| __happy | 6496 |
| ? | 31890 |
| <EOS> | 3 |

Input Text: Why are you so happy?

| Word Segmentations | Word Tokens |
|---|---|
| <BOS> | 2 |
| __为什么 | 4117 |

# Encoder

inputs

"Why are you so happy?"

tokenize

embedding

encoder



## TransforLearn: Interactive Visual Tutorial for the Transformer Model

**Input View**

Sentence to be translated:

🔍 Why are you so happy?

**Translation View**

Translation: 为什么你如此高兴?

Current translation: 为什么你如此高兴?  ‹  ›

Prediction in current iteration:  **?**

Cumulative Score (Probabilities):  **0.039872859081467996**

**Architecture View**

Sequece Data

Parameter

Word Segmentations  Word Tokens

embeddings (8,512)  🖱 *Click* to learn how to get embeddings!

after positional encoding (8,512)  🖱 *Click* to lea positional en

| <BOS> | 2 |
| __Why | 3453 |
| __are | 120 |
| __you | 1418 |
| __so | 323 |
| __happy | 6496 |
| ? | 31890 |
| <EOS> | 3 |

Input Text

Why are you so happy?

Word Segmentations  Word Tokens

embeddings (5,512)  🖱 *Click* to learn how to get embeddings!

after positional encoding (5,512)  🖱 *Click* to lea positional en

| <BOS> | 2 |
| __为什 | 4117 |

# Encoder

inputs

"Why are you so happy?"

tokenize

embedding

encoder

---

**TransforLearn: Interactive Visual Tutorial for the Transformer Model**

**Input View**

Sentence to be translated: Why are you so happy?

**Translation View**

Translation: 为什么你如此高兴？

Current translation: 为什么你如此高兴？

Prediction in current iteration: ?

Cumulative Score (Probabilities): 0.039872859081467996

**Architecture View**

Sequece Data

Parameter

| Word Segmentations | Word Tokens |
| --- | --- |
| <BOS> | 2 |
| __Why | 3453 |
| __are | 120 |
| __you | 1418 |
| __so | 323 |
| __happy | 6496 |
| ? | 31890 |
| <EOS> | 3 |

Input Text: Why are you so happy?

embeddings (8,512)   Click to learn how to get embeddings!

after positional encoding (8,512)   Click to lea positional en

| Word Segmentations | Word Tokens |
| --- | --- |
| <BOS> | 2 |
| __为什 | 4117 |

embeddings (5,512)   Click to learn how to get embeddings!

after positional encoding (5,512)   Click to lea positional en

# Encoder

inputs

"Why are you so happy?"



tokenize

embedding

encoder

---

## TransforLearn: Interactive Visual Tutorial for the Transformer Model

**Input View**

Sentence to be translated: Why are you so happy?

**Translation View**

Translation: 为什么你如此高兴?

Current translation: 为什么你如此高兴

Prediction in current iteration: ?

Cumulative Score (Probabilities): 0.039872859081467996

**Architecture View**

# Encoder



self-attention

layer normalization

feed-forward

encoder

TransforLearn: Interactive Visual Tutorial for the Transformer Model

# Encoder - Attention

# Encoder



self–attention

layer normalization

feed–forward

# Encoder



Layer Normalization of Encoder Block #1

# Encoder

# Encoder



Feed Forward Network of Encoder Block #1

Input (8,512)     Input Embed (1,512)     Intermediate Embed (1,2048)     Activation Embed (1,2048)     Output Embed (1,512)     Output (8,512)

Linear    $\max(0, -5.75)$    ReLU    Linear

feed–forward

# Visual Design - Overview



**Architecture-driven Exploration**

Module Detailed Views ← Architecture Overview → Module Detailed Views

**Task-driven Exploration**

# Task- driven Exploration



Task-driven Exploration

Explore data flow changes

- Input and output, data dimension, data range

Analyze structural features

- Decoding time step -> translation progress
- Focus on a specific module or head

# Task- driven Exploration

Explore data flow changes

- Input and output, data dimension, data range

Analyze structural features

- Decoding time step -> translation progress

- Focus on a specific module or head

# Usage Scenario

Self-study guidance for a beginner

- utilize Transformer to extract features from sequence data

- the concept and generation process of the Q, K, and V matrices

- the use of decoders for prediction

Teaching aid for lectures

- better summarize and present the teaching points

- increases the practicality and vividness of the entire teaching process

# Usage Scenario

Usage Scenario

Self-study guidance for beginners

# Evaluation

## User-controlled Experiment



User Portrait → System Introduction → System Tutorial → Observable Study → Exit Questionnaire

R-1  visual summary

R-2  interactive interface

R-3  exploration mode

R-4  self-directed & immersive

| Level | Goal | Question |
|---|---|---|
| easy | G1 | Q1: Components and data flow of feed-forward network. |
| easy | G3 | Q2: Identify key words from attention matrix. |
| easy | G3 | Q3: Final output in translation task and its derivation. |
| medium | G1 | Q4: Differences between cross- and self-attention. |
| medium | G2 | Q5: Add & LN significance and implementation. |
| medium | G1 | Q6: Parallelism in Transformer. |
| hard | G2 | Q7: Reasons for scaling before softmax. |
| hard | G2 | Q8: Process of calculating PE & variation with position. |

# Evaluation

- improve users' <span style="color:red">understanding</span> of structures and tasks
- bring more <span style="color:red">activity</span>, <span style="color:red">autonomy</span> and <span style="color:red">divergent thinking</span>
- enhancing users' <span style="color:red">efficiency</span> in learning through a broader coverage and enhanced interaction



(A) Comparison of Questions' Answers Among Groups

(B) Comparison of Answer Time Among Groups

(C) Comparison of Learning Efficiency Index Among Groups

| $E_{GroupX,i}$ | $i=1$ | $i=2$ | $i=3$ | $i=4$ | $i=5$ | $i=6$ | $i=7$ | $i=8$ | Mean | Std |
|---|---|---|---|---|---|---|---|---|---|---|
| $X=B$ | 0.737 | 1.005 | 0.680 | 0.839 | 0.981 | 0.824 | 0.965 | 0.851 | 0.851 | 0.121 |
| $X=T$ | 1.421 | 1.702 | 1.631 | 1.402 | 1.263 | 1.385 | 1.381 | 1.542 | **1.466** | 0.146 |

# Evaluation

User interviews

<span style="color:green">Implication</span>

- Usability and effectiveness.
- Validating the knowledge for experts.

<span style="color:red">Limitation</span>

- Different appropriate learning resources for different needs.
- Need for more instructions, animations, and comparisons.

# Thanks for your listening !

TransforLearn: https://trans-for-learn.github.io/

Welcome to our homepage: http://fduvis.net/

Email: leenagao0430@gmail.com

## TransforLearn

### Interactive Visual Tutorial for the Transformer Model

Lin Gao[1], Zekai Shao[1], Ziqing Luo[1], Haibo Hu[2], Cagatay Turkay[3], Siming Chen[1,4]

[1]School of Data Science, Fudan University

[2]School of Big Data & Software Engineering, Chongqing University

[3]Centre for Interdisciplinary Methodologies, University of Warwick

[4]Shanghai Key Laboratory of Data Science

📄 Paper    Code    🖥 Demo