# Advanced Deep Learning - Transformer Models

—

## Multi-Head Attention Mechanism

# About the Instructor

Charles E. Gormley

Areas of Study: Data Science & Economics



**Career**

-Economics Researcher

-Machine Learning Engineer - Fintech & Medtech Companies

-Software Engineer - Major Finance Company

-Machine Learning Engineer - Consulting

# Agenda

➢ PreRequisites
➢ Overview of Transformer Model
➢ Attention
➢ Attention Mecahnism
  ○ Linear Function
  ○ Query, Key & Value
➢ Multiple Heads - Why?
  ○ Attention Filter
➢ Recap
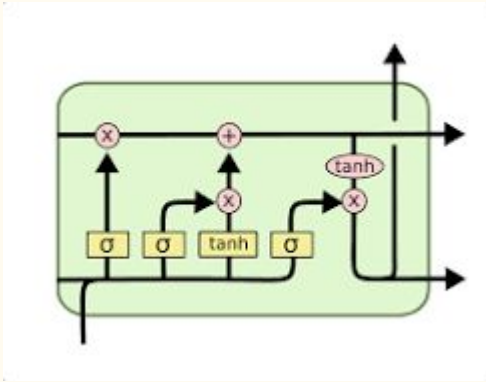➢ What's Next →

# PreRequisites

## Mathematics & Modeling

- Linear Algebra: Matrix Multiplication
- Linear Algebra: Scaling Matrices
- Multi-Layer Perceptrons
- Activation Functions
- AutoEncoders
- Positional Encoding
- Tokenization
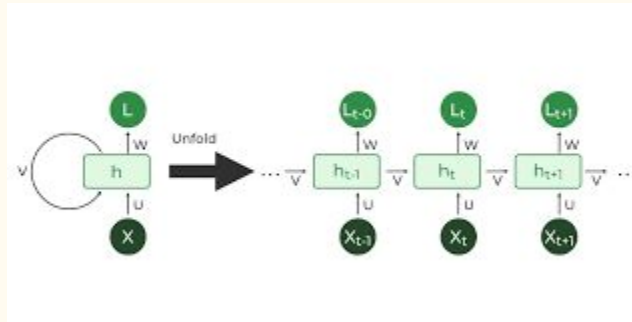- Positional Embeddings
- Back-Propagation

## Programming

- Basic Python 3.x
- Object Oriented Programming
- Numpy
- SkLearn
- PyTorch & Keras

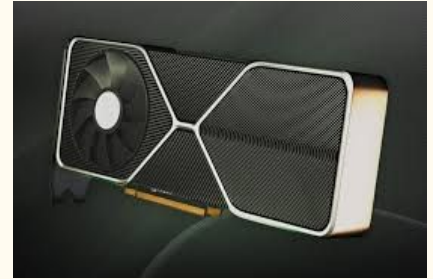# Overview of Transformer Models: Definition

What is a Transformer Model: Neural Network with extensive ability to handle long range dependencies in sequential data & their capacity for parallelization.



LSTM



RNN



Parallelization

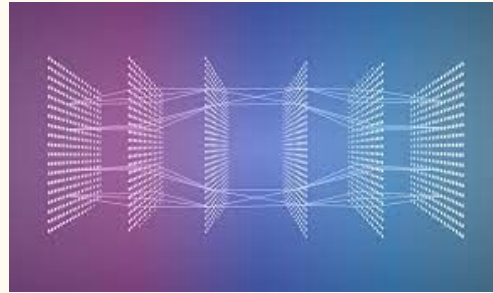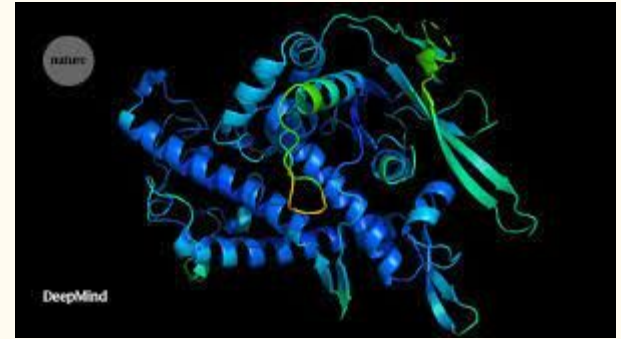# Overview of Transformer Models: Industry Use Cases
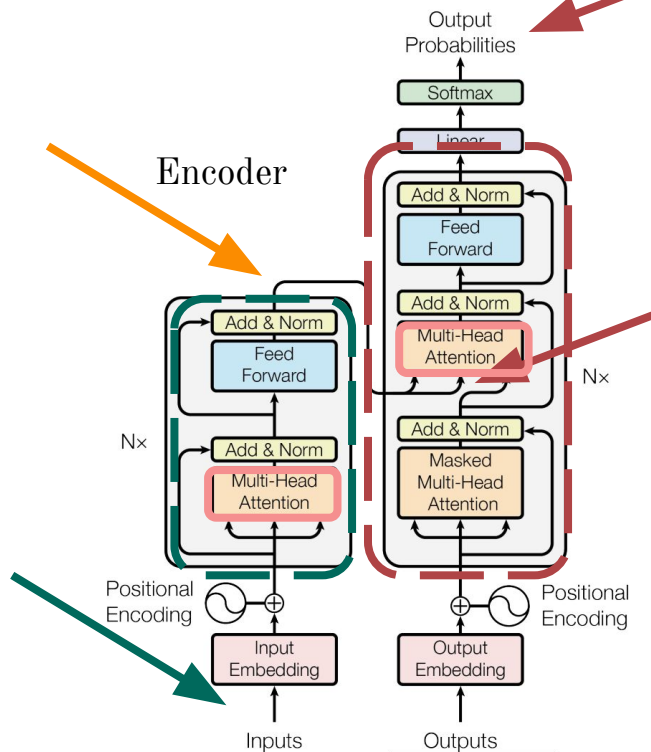


GPT Algorithms



Image Transformers (VIT)



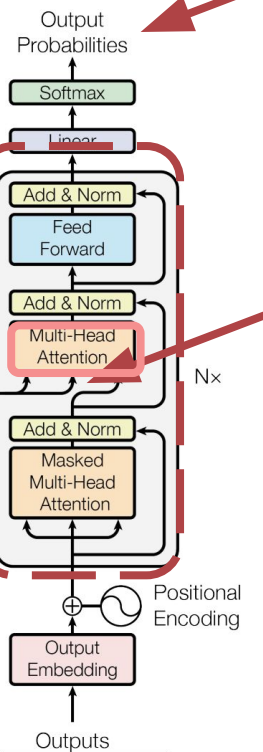Protein Folding

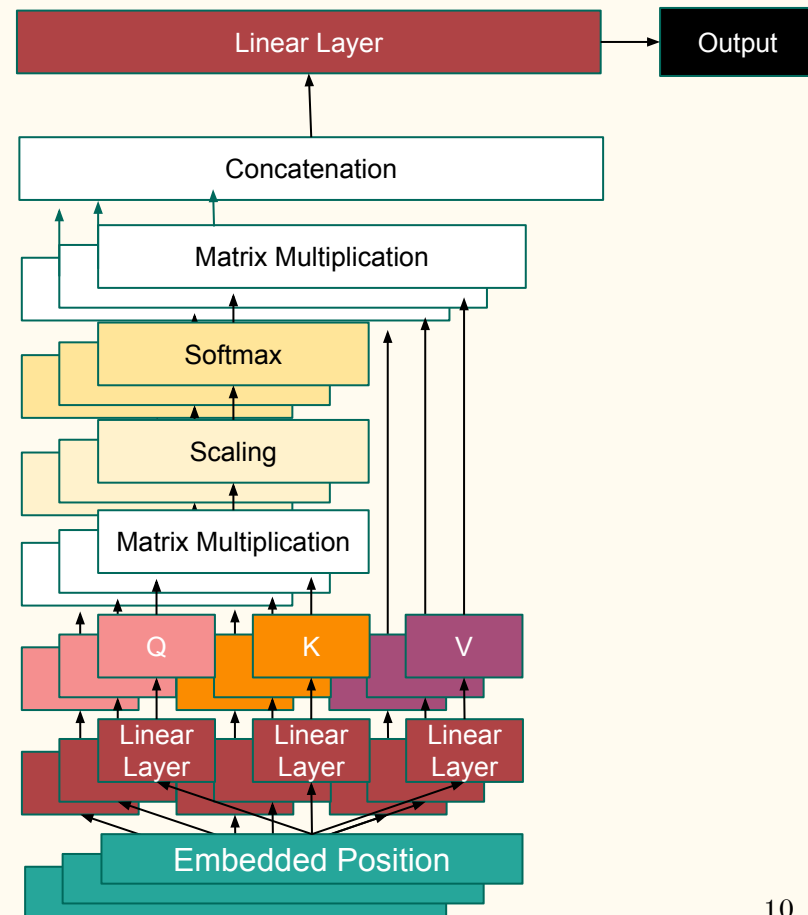# Overview of Transformer Models: Architecture

# What is Attention?

*"**Winter is coming,** we know what's coming with it. We can learn to live with the **wildlings** or we can add them to the **army of the dead**"*
*~Jon Snow*

# What is Self-Attention

What is the meaning of the following **sentence**:

*"The man who passes the* **sentence** *should swing the sword" ~ Ned Stark*
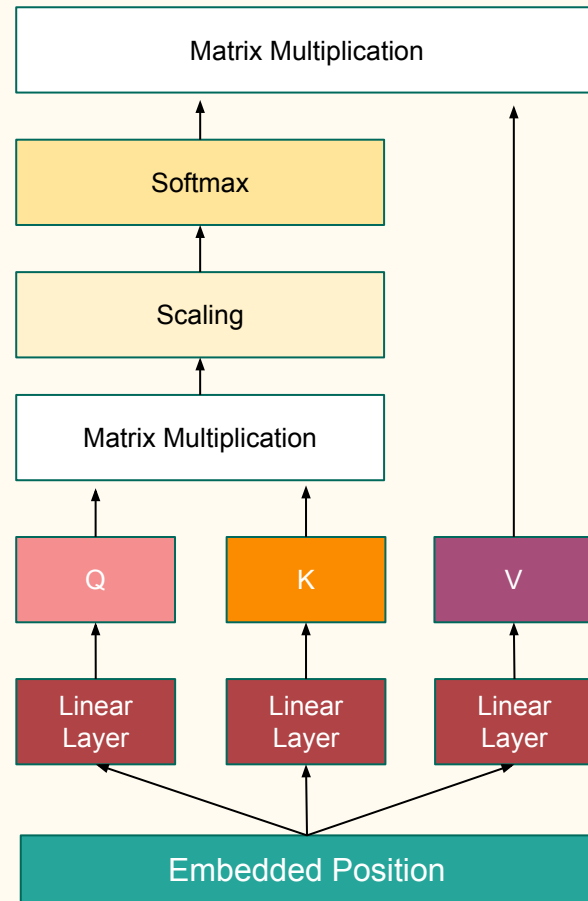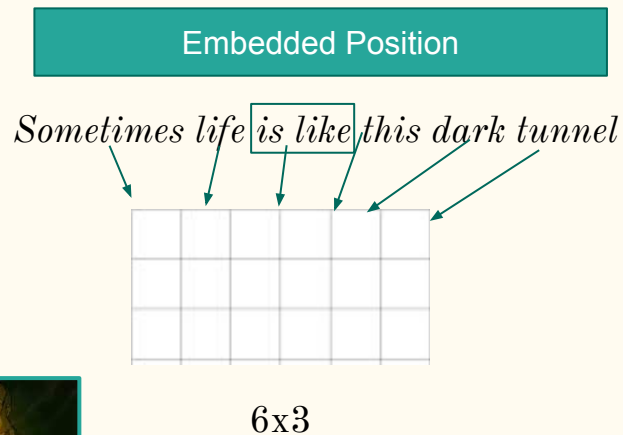
# Multi-Head Attention

# Single Attention Mechanism

Sections
1. Embedded Position Layer
2. Linear Layers
3. Query, Key, Value Matrices
4. Query x Key - MatMul
5. Scaler
6. Softmax - Normalization

# Single Attention Mechanism

**Embedded Position**

*Sometimes life is like this dark tunnel*

6x3

| Matrix Multiplication |
| Softmax |
| Scaling |
| Matrix Multiplication |

Q | K | V

Linear Layer | Linear Layer | Linear Layer

**Embedded Position**

# What is Query, Key, Value

# Single Attention Mechanism

Linear Layer    Linear Layer    Linear Layer

xW+b=y

6x6     6x3     +bias = Query

6x3

+bias = Key

+bias = Value

Matrix Multiplication

Softmax

Scaling

Matrix Multiplication

Q    K    V

Linear Layer    Linear Layer    Linear Layer

Embedded Position

14

# Single Attention Mechanism

# Single Attention Mechanism -Attention Filter



*Sometimes life is like this dark tunnel*

# Single Attention Mechanism

Softmax

Scaling

Scaling | Softmax

$1 \ / \ 6^{1/2}$

Helps Normalize
Values between $0 \rightarrow 1$

$6 =$ Key Vector Length

Matrix Multiplication

Softmax

Scaling

Matrix Multiplication

Q

K

V

Linear Layer

Linear Layer

Linear Layer

Embedded Position
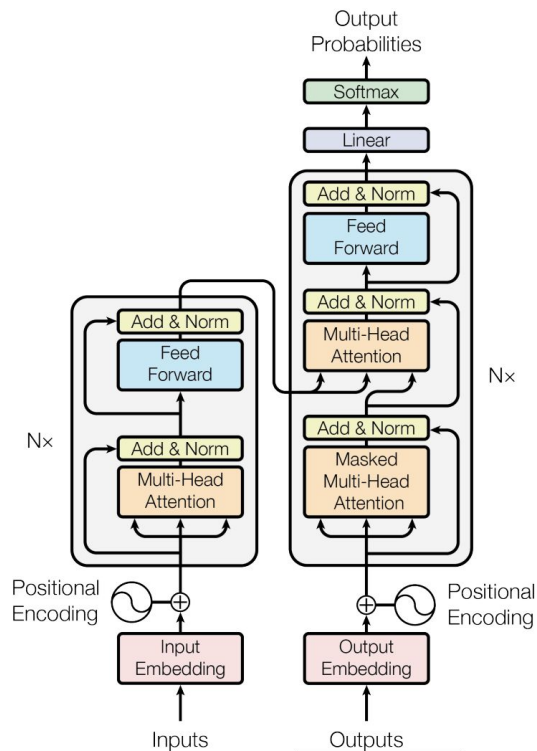
17

# Multi-Head - Self-Attention Filters

# Multi-Head Attention

# What's Next..



**BERT**

Encoder

**GPT**

Decoder

# What We Covered Today

➢ Overview of Transformer Model
➢ What is Attention
➢ Attention Mechanisms
➢ Multiple Head Attention