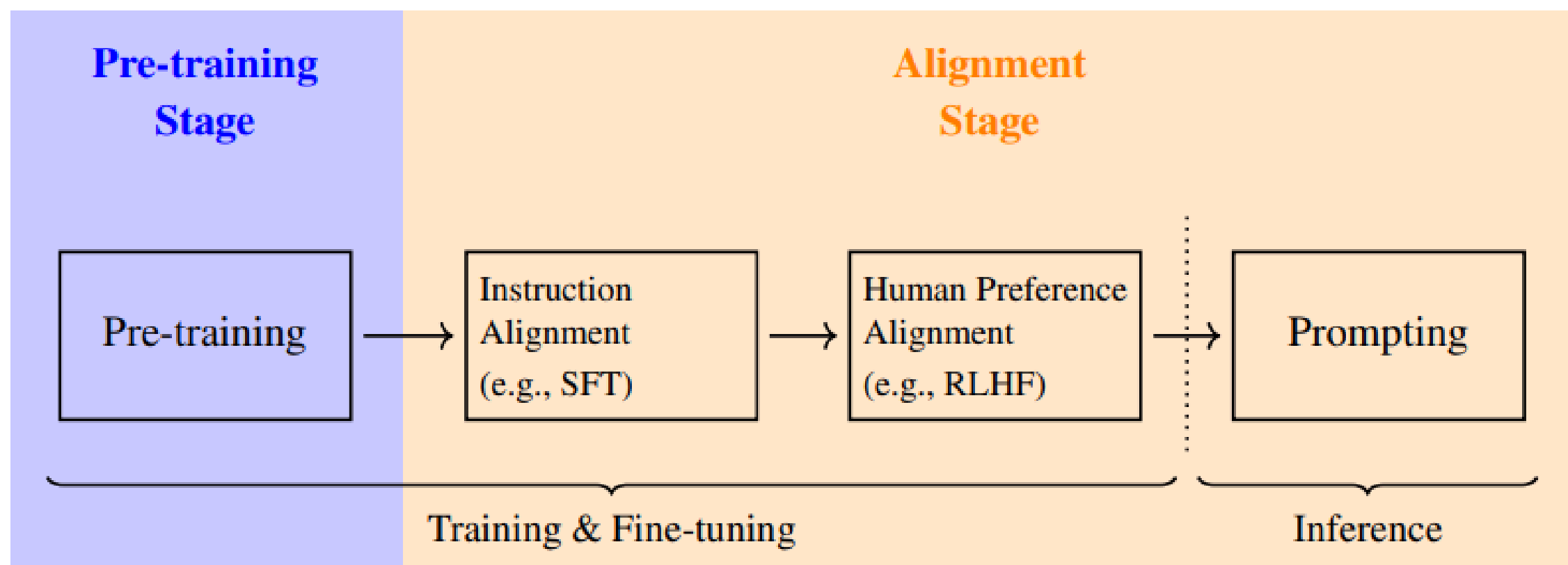


# LLM Foundations



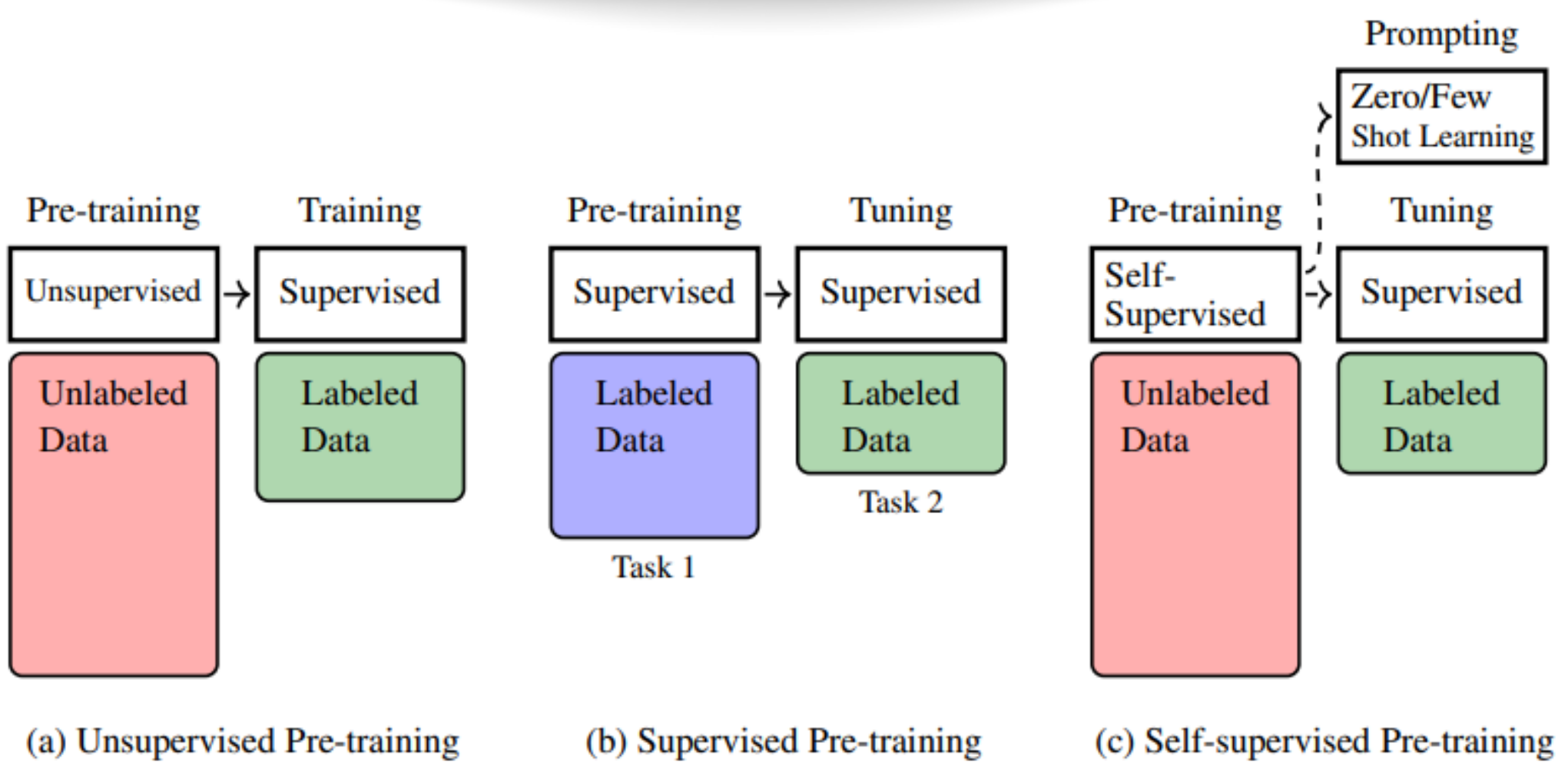
Curated by:  
Dr. Maryam Miradi

## Overview



# LLM Foundations

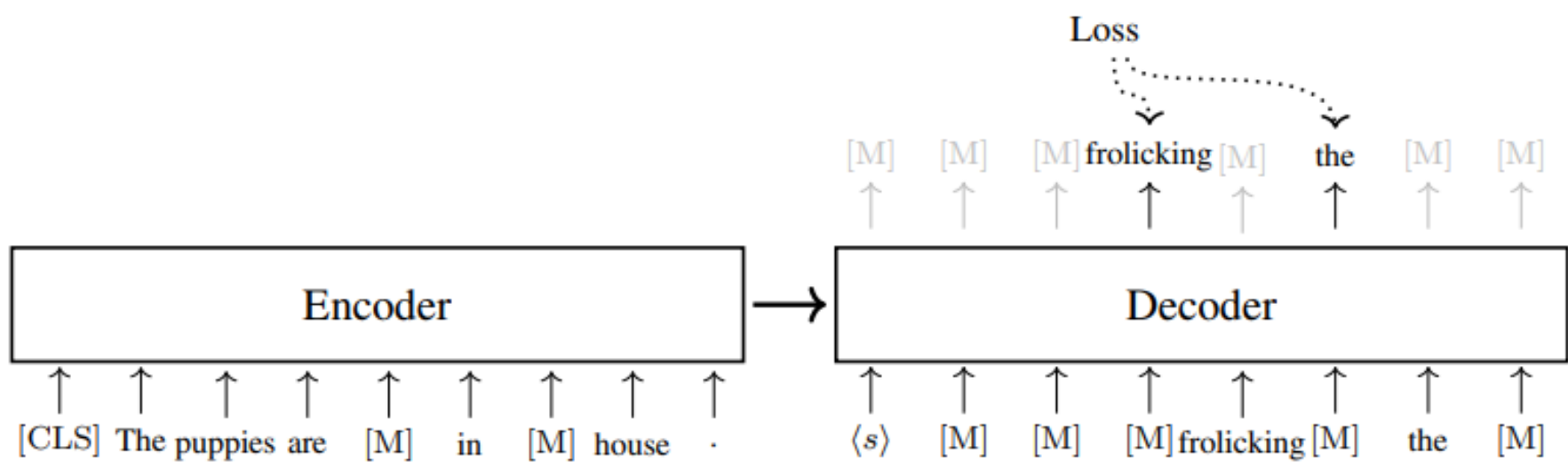
## Pre-training types



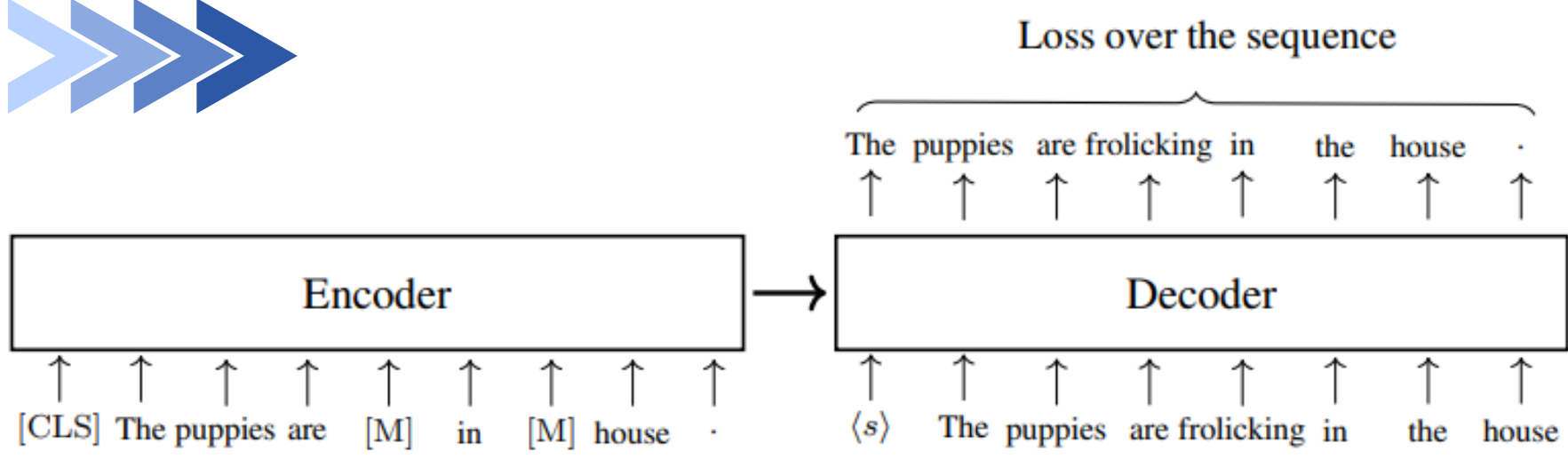
# LLM

# Foundations

## Encoders & Decoders



(a) Training an encoder-decoder model with BERT-style masked language modeling

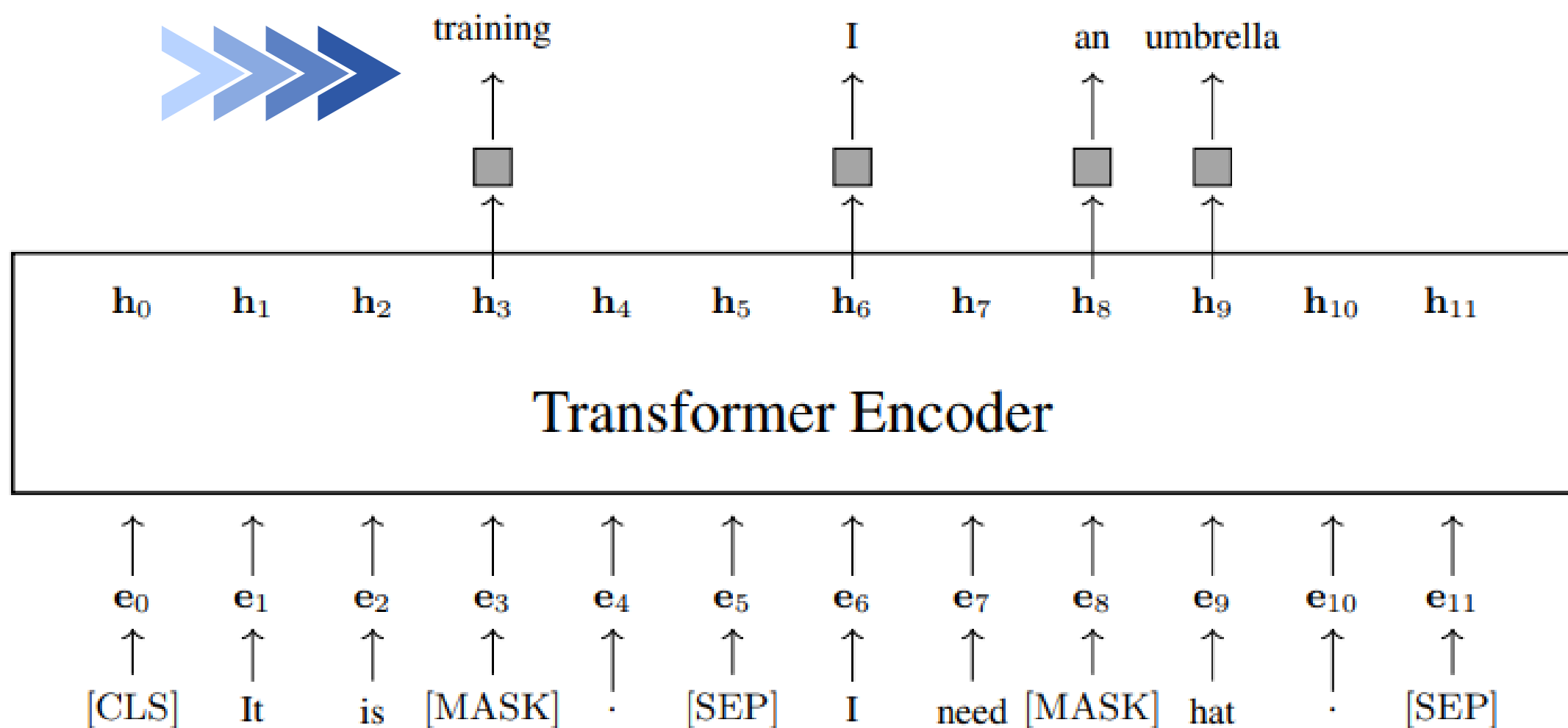


(b) Training an encoder-decoder model with denoising autoencoding



# LLM Foundations

## Transformers Pre-training

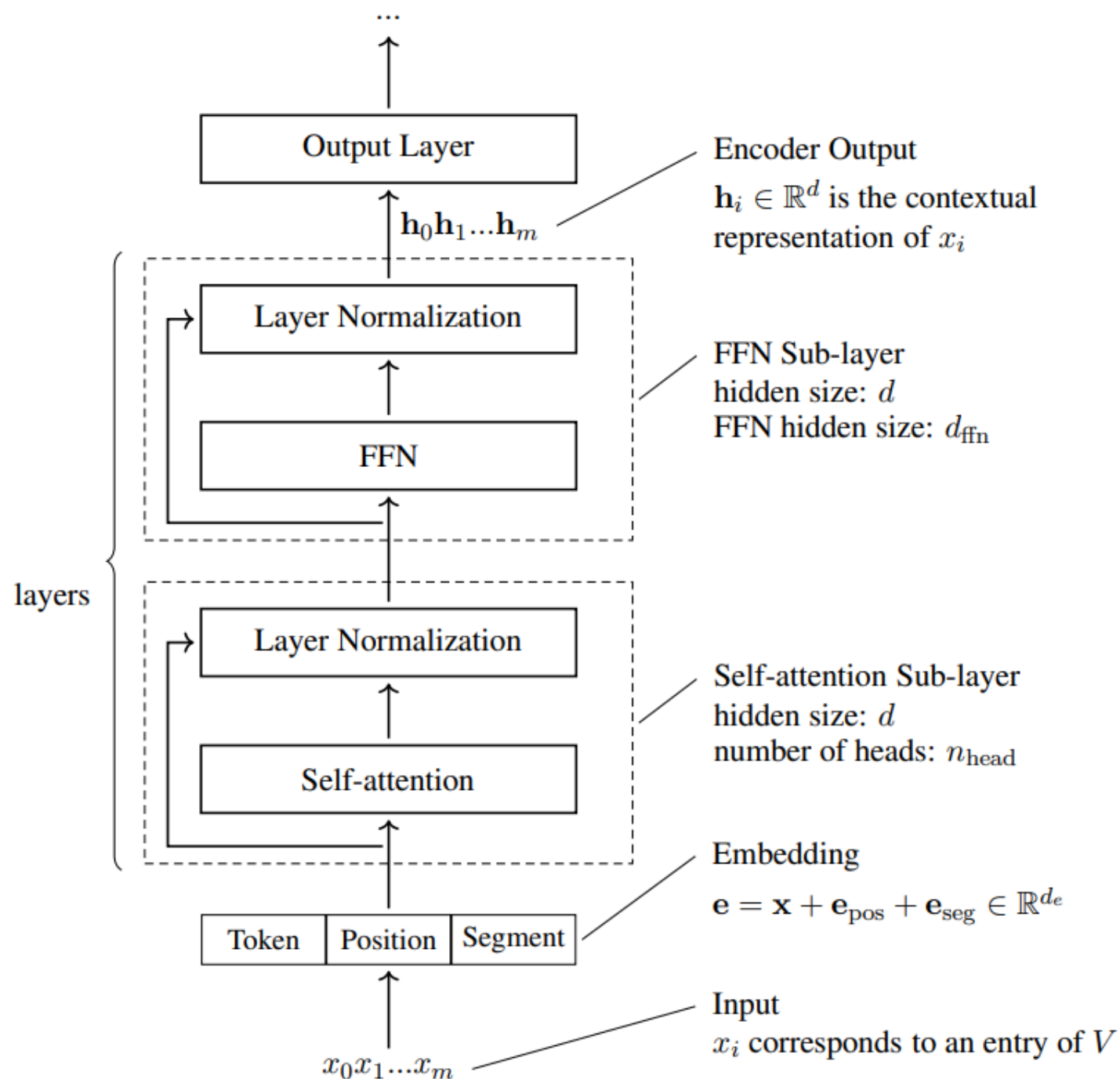


# LLM Foundations

## BERT



Curated by:  
Dr. Maryam Miradi

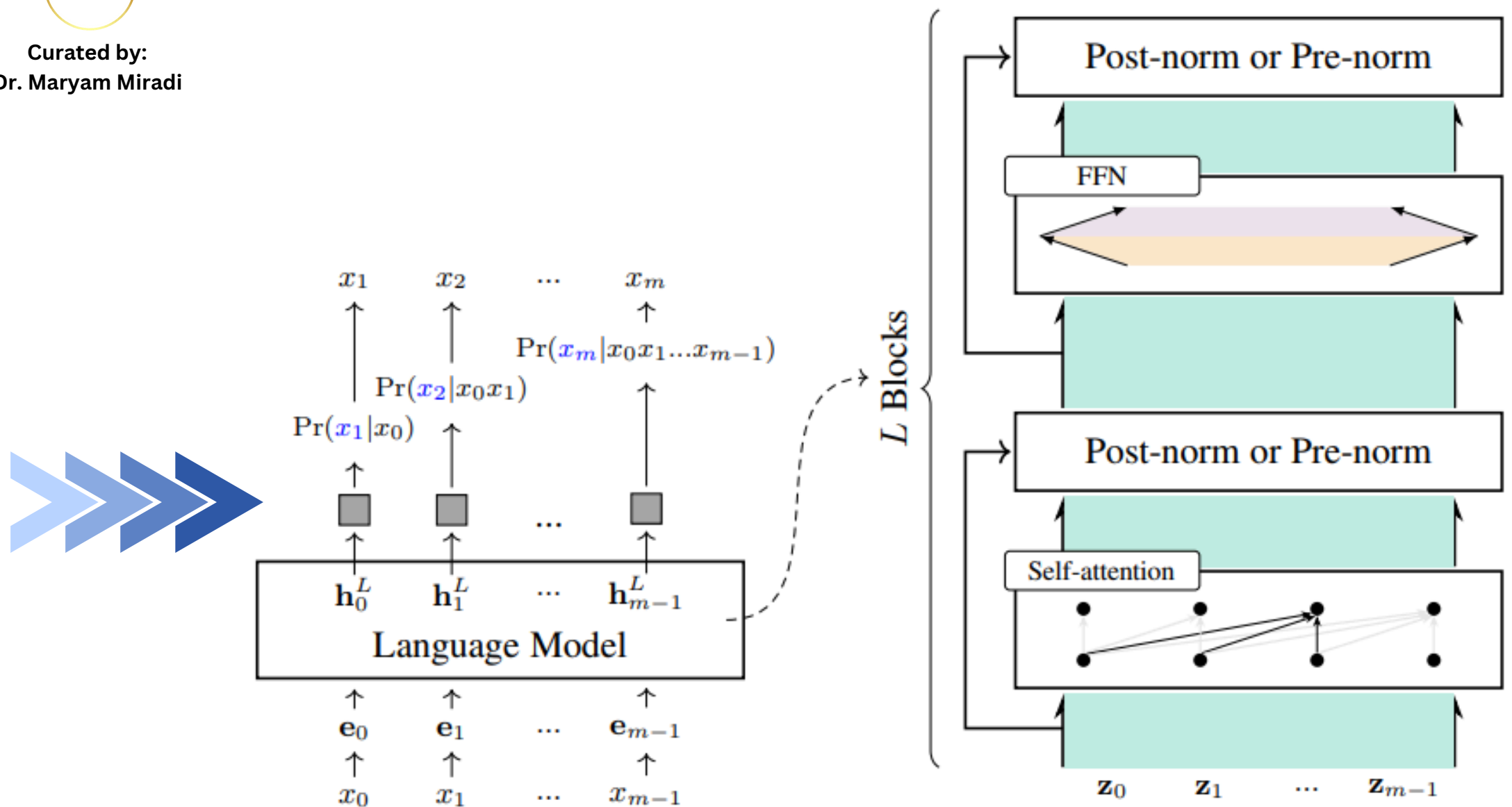


# LLM Foundations

## Transformer-decoder architecture



Curated by:  
Dr. Maryam Miradi

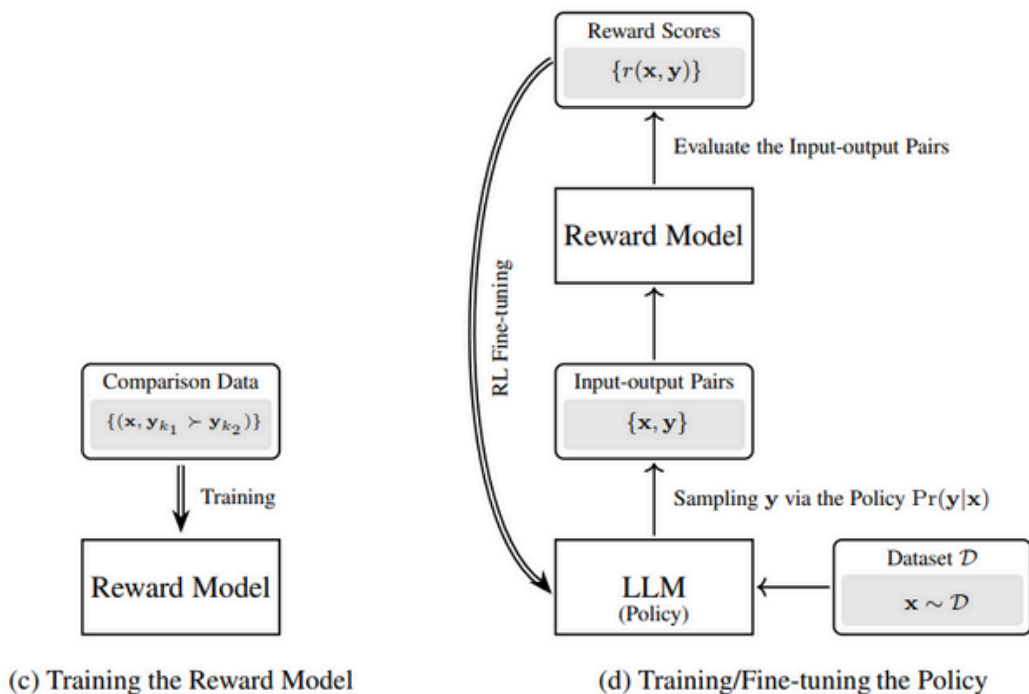
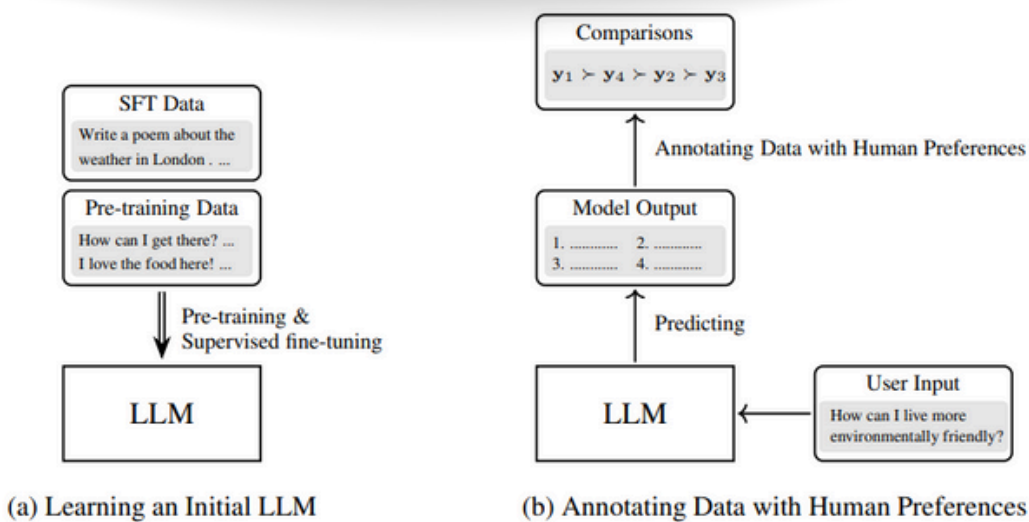


# LLM

# Foundations

7/16

## RLHF

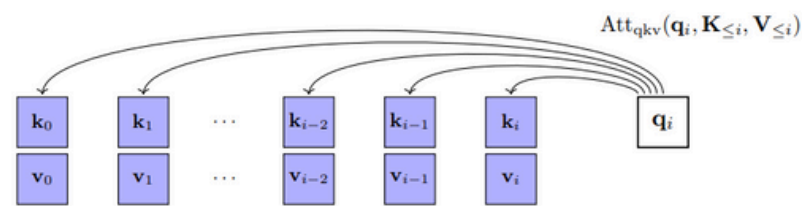


Curated by:  
Dr. Maryam Miradi

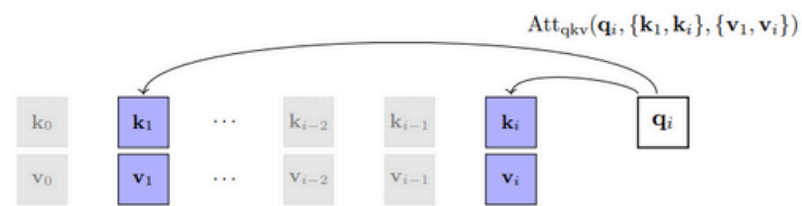
8/16

# LLM Foundations

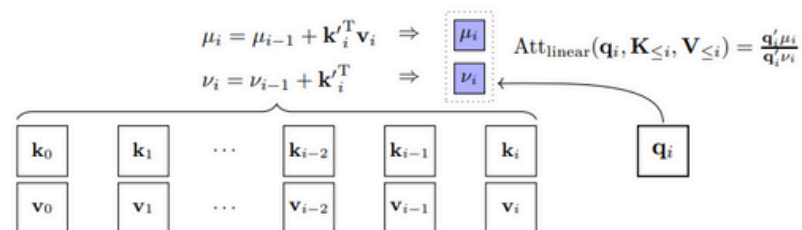
## Attention Types



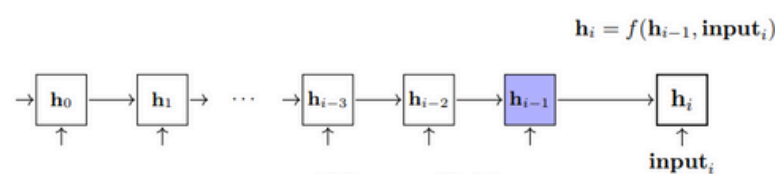
(a) Standard Self-attention



(b) Sparse Attention



(c) Linear Attention



(d) Recurrent Models



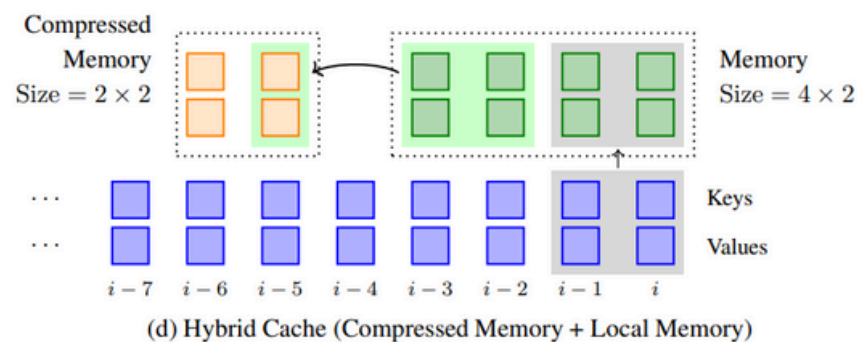
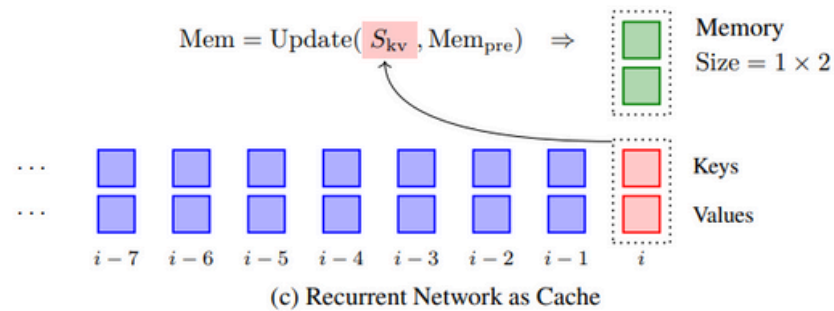
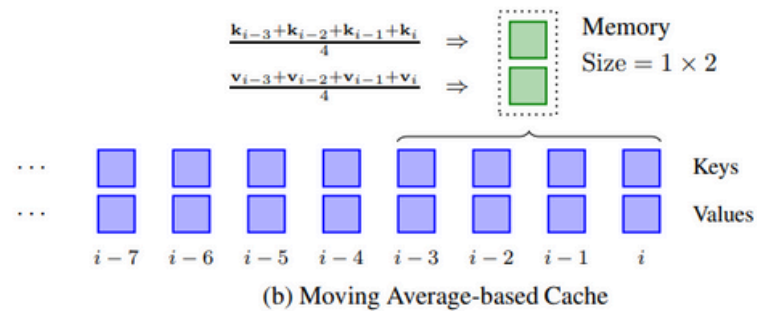
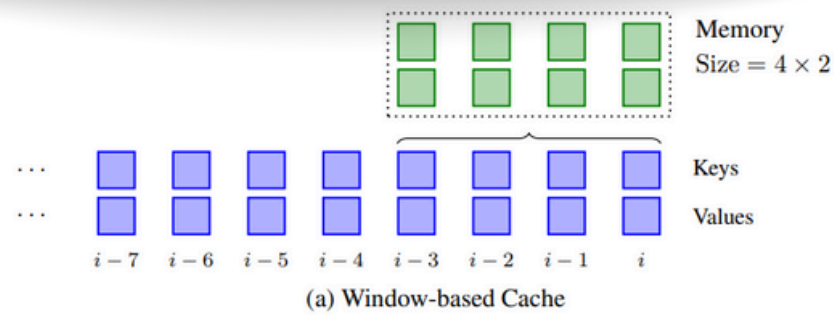
Curated by:  
Dr. Maryam Miradi



9/16

# LLM Foundations

# Memory



Curated by:  
Dr. Maryam Miradi

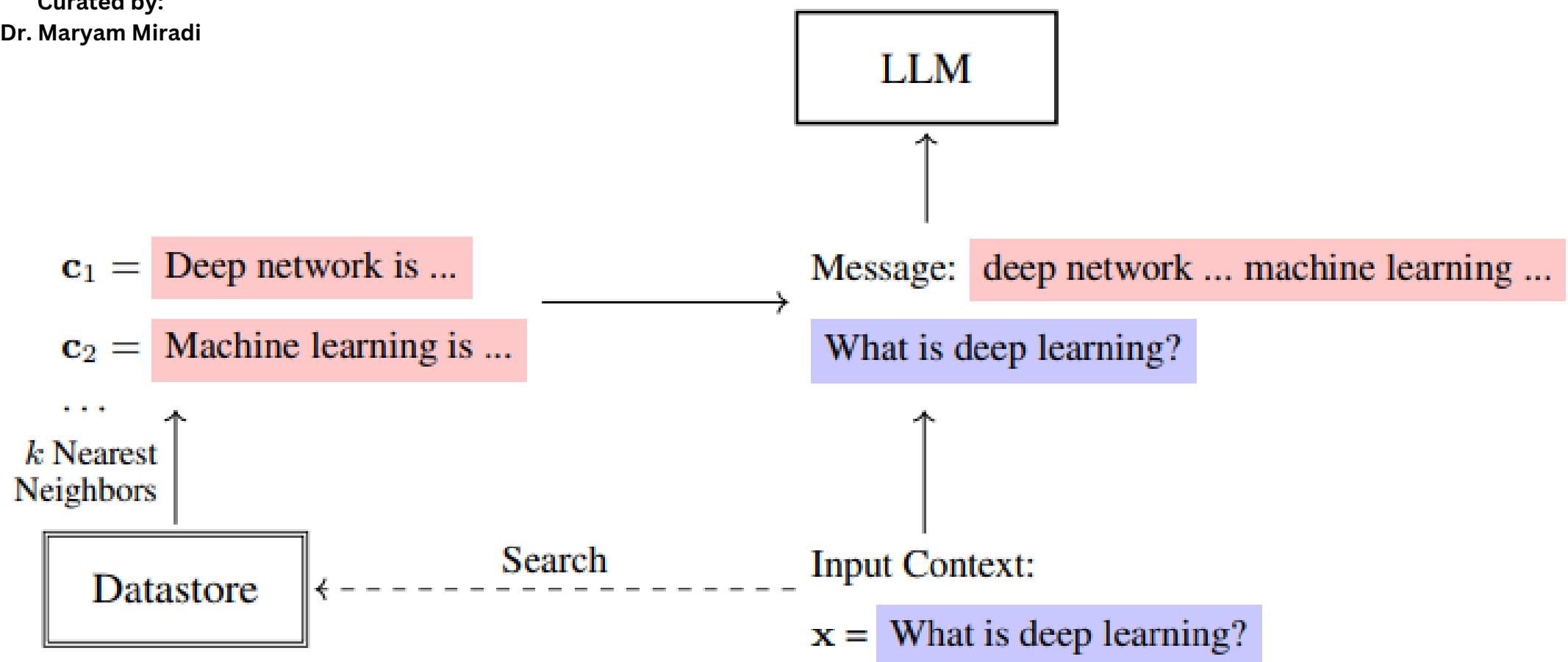
# LLM Foundations

## RAG

10/16



Curated by:  
Dr. Maryam Miradi



(c) Retrieval-augmented Generation

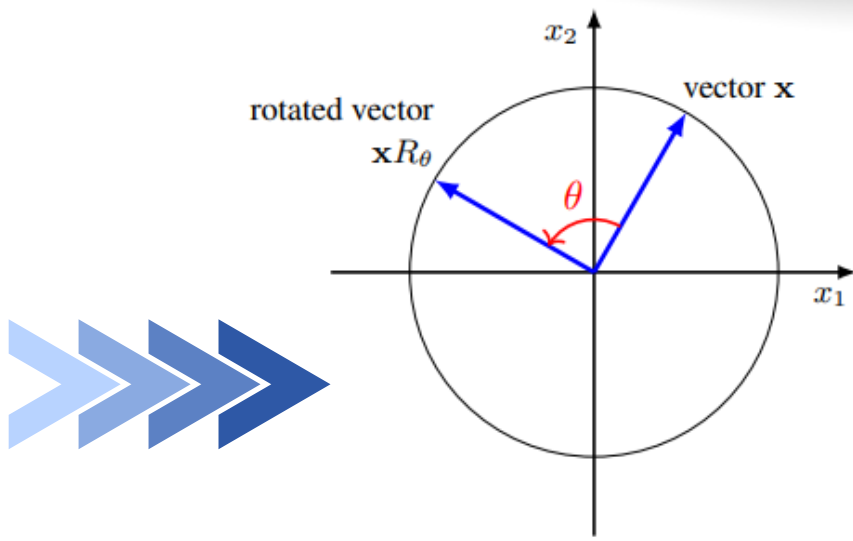
# LLM Foundations

11/16

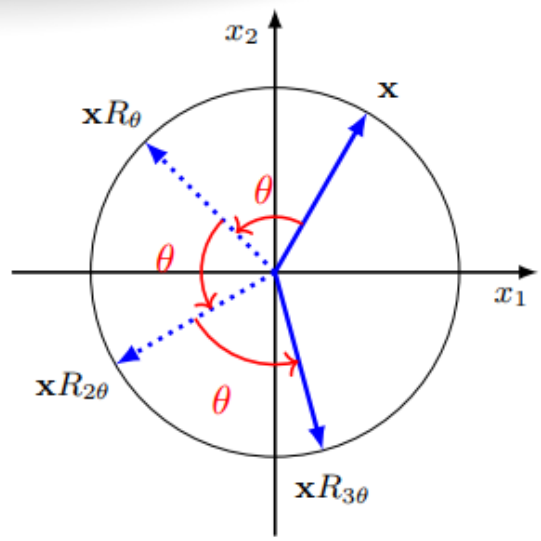
## Embedding



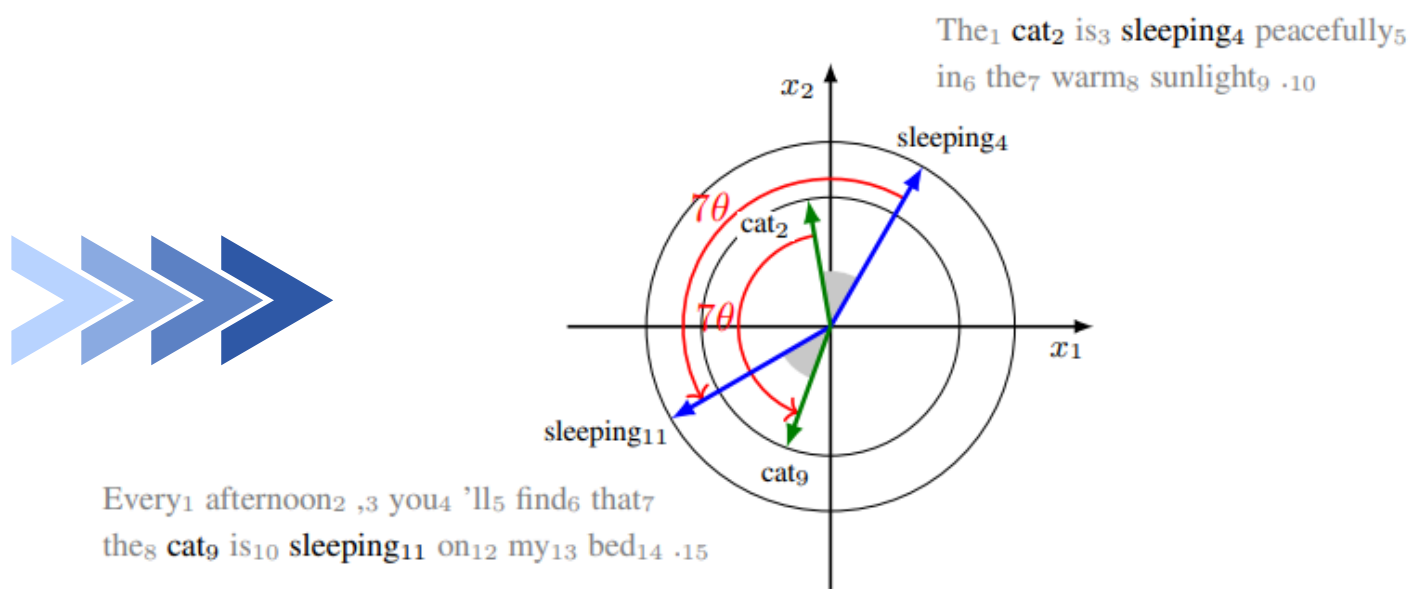
Curated by:  
Dr. Maryam Miradi



(a) Single-step Rotation



(b) Multi-step Rotation



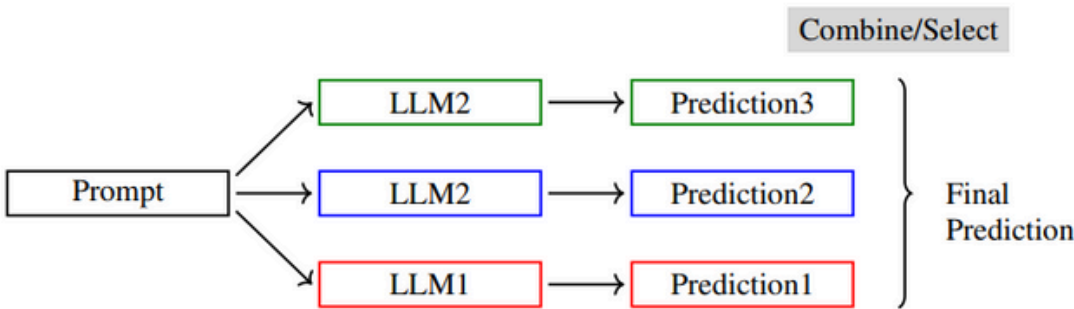
(c) Angles between embeddings of two tokens at different positions

# LLM

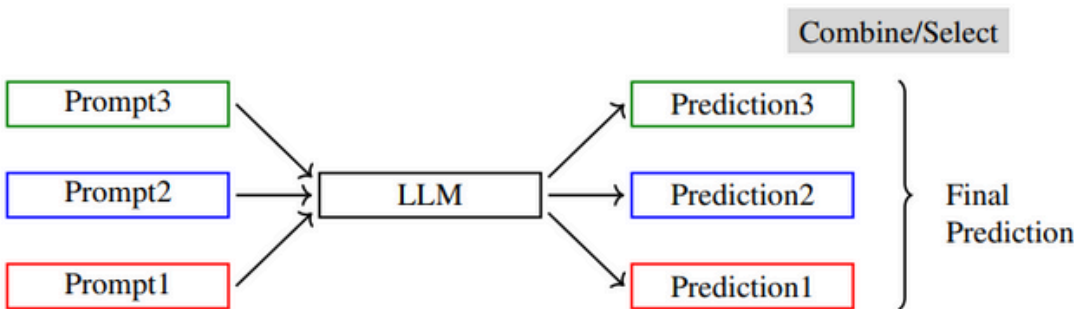
# Foundations

12/16

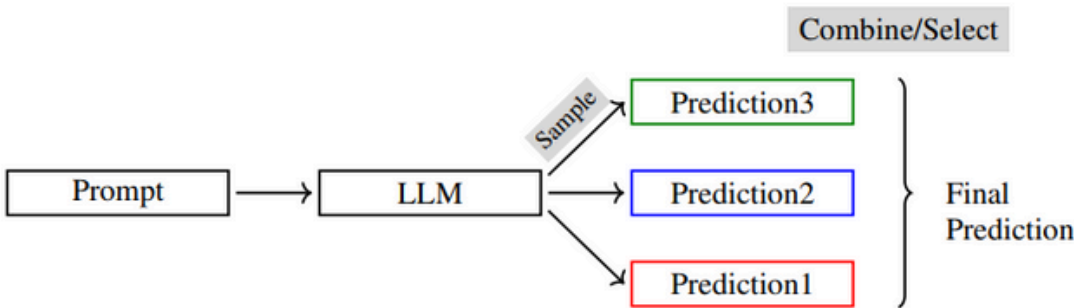
# Ensembling



(a) Model Ensembling



(b) Prompt Ensembling

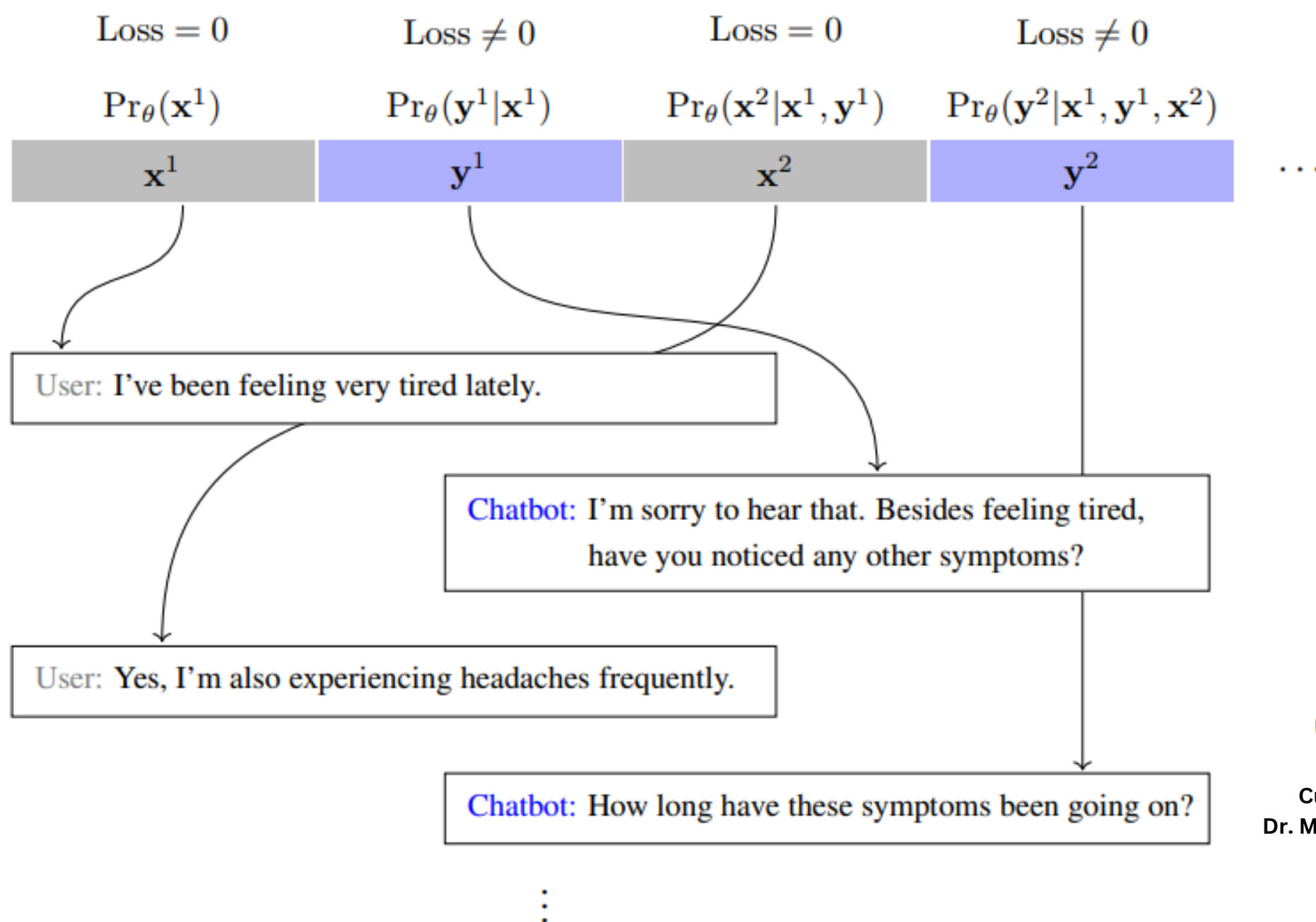


(c) Output Ensembling

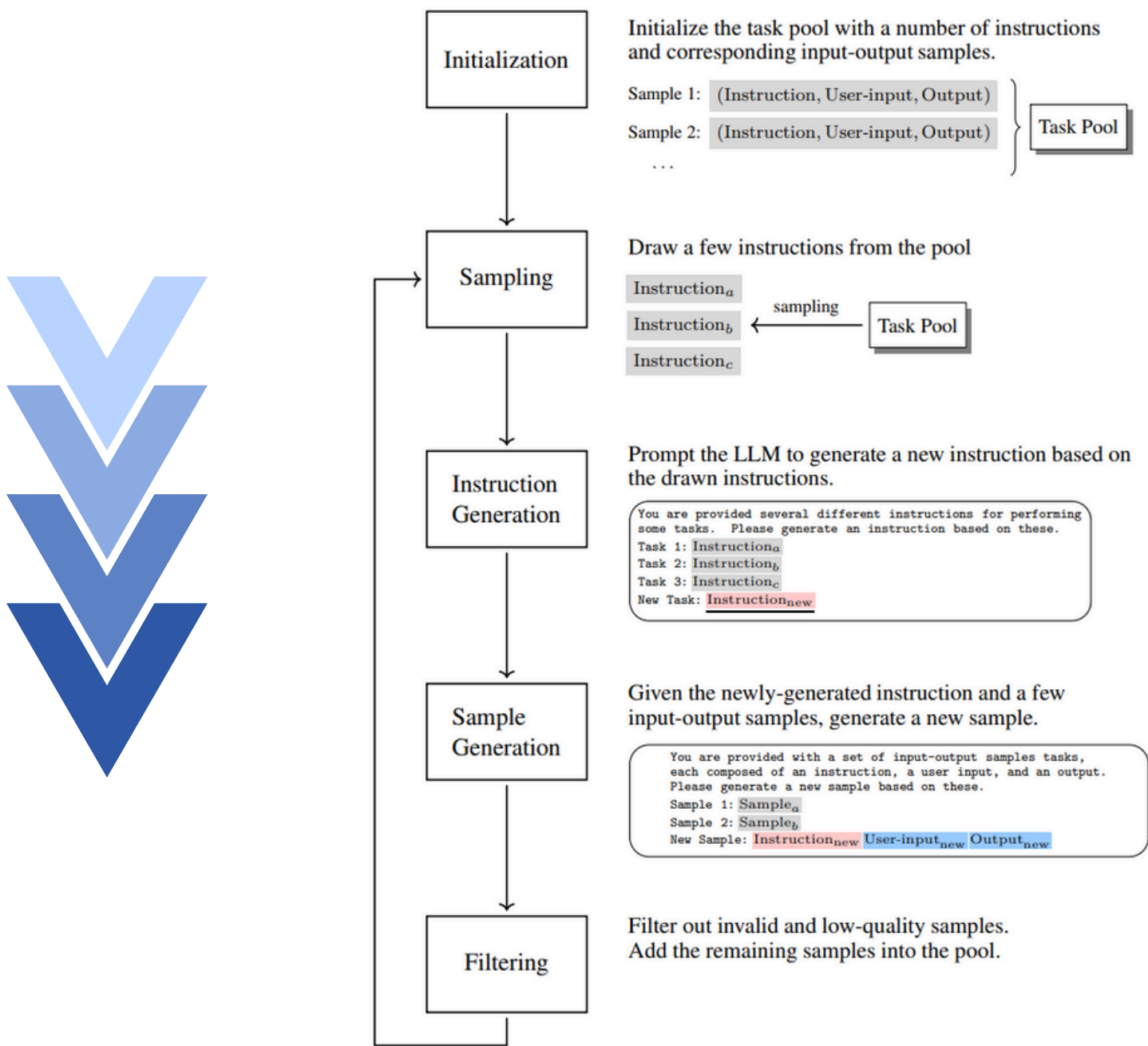


Curated by:  
Dr. Maryam Miradi

# Fine-Tuning



# Self-Instruct



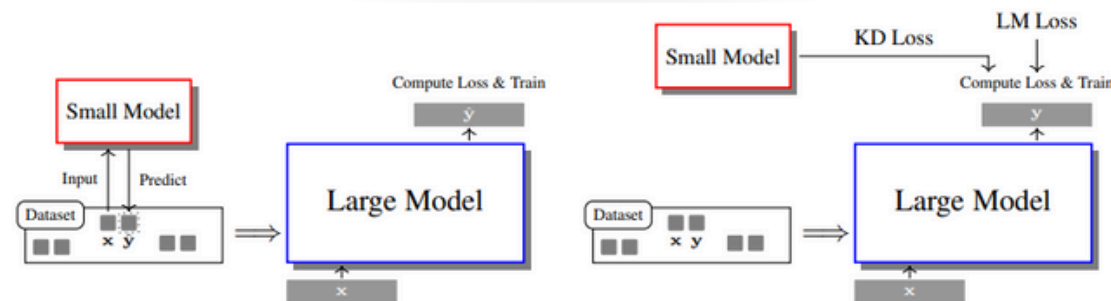
Curated by:  
Dr. Maryam Miradi

# LLM

# Foundations

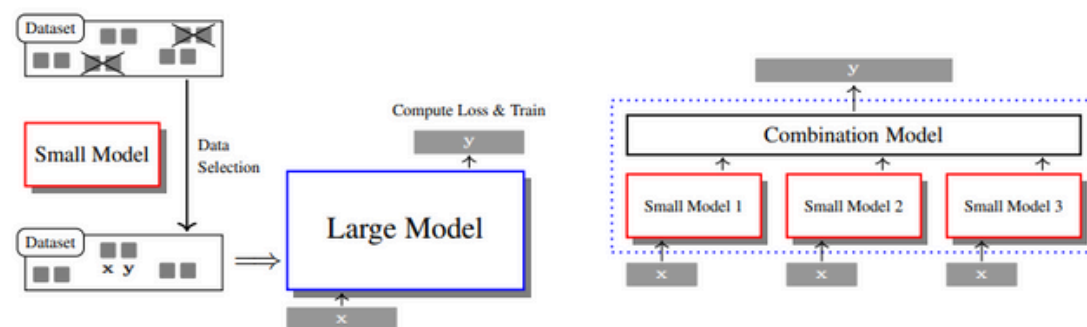
15/16

## Small-to-Large



(a) Fine-tuning on data generated by a small model (weak-to-strong generalization)

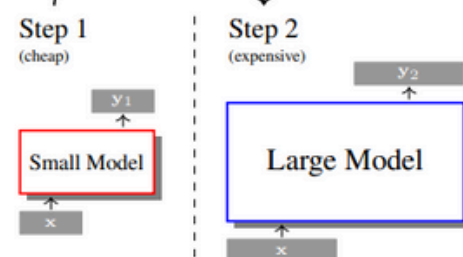
(b) Fine-tuning with KD Loss from a small model (weak-to-strong generalization)



(c) Data selection with a small model

(d) Ensemble of multiple small models

If Step 1 is not satisfactory, go to Step 2



(e) Cascading (at inference time)

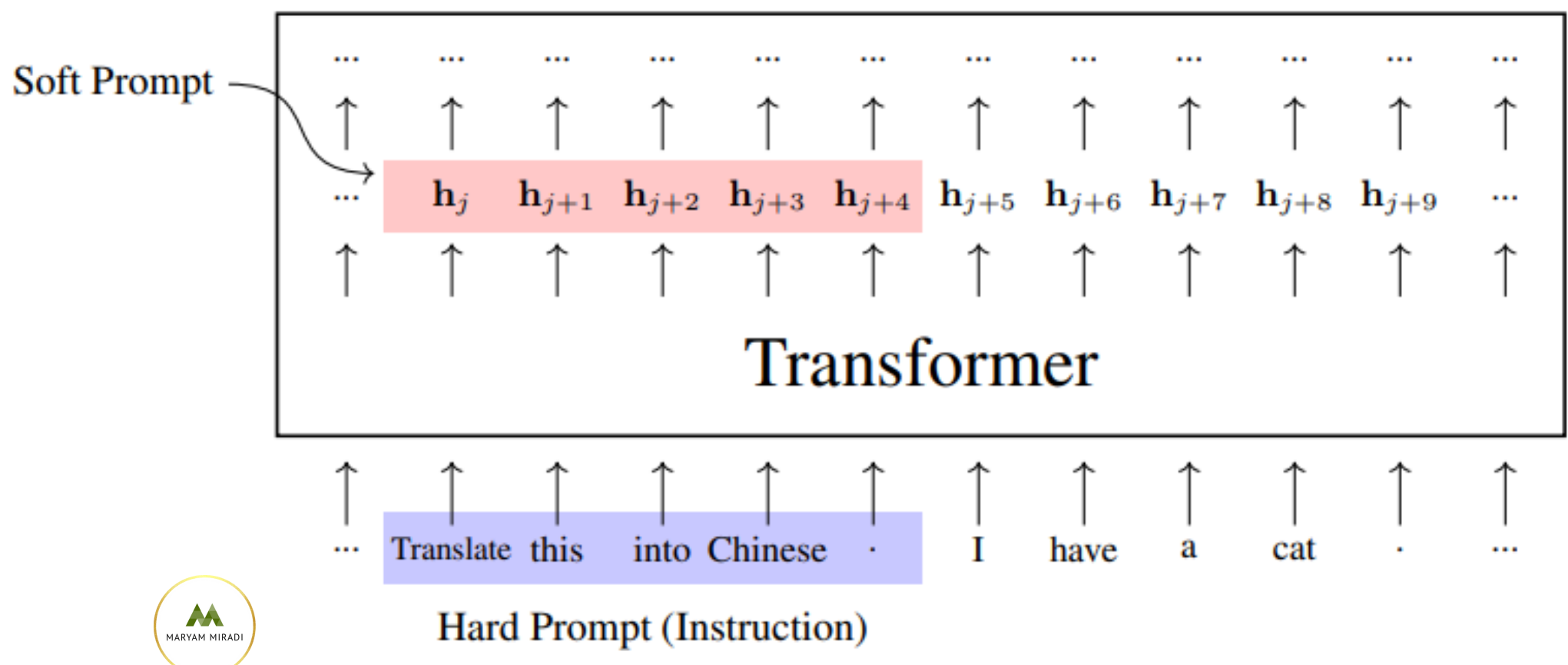


Curated by:  
Dr. Maryam Miradi

# LLM Foundations

16/16

## Soft Prompts



Curated by:  
Dr. Maryam Miradi

