



# DeepLorlA

## Mastering Large Language Models: Efficient Techniques for Fine-Tuning

Barbara Gendron-Audebert, PhD student (MosAlk team)

LORIA, Université de Lorraine, CNRS  
DeepLorlA Network

January 15, 2025

# About Me

2nd-year PhD student - Knowledge-Enhanced Language Models (Université de Lorraine)

## Background & Research

- Maths engineering degree
- Master Thesis: Meta-Learning in Conversational Context
- PhD Topic: Controlled Conversational Models through Conversation-Dedicated Ontology
- *Keywords: Large Language Models (LLMs), Conversational Agents, Ontologies, Fine-Tuning*

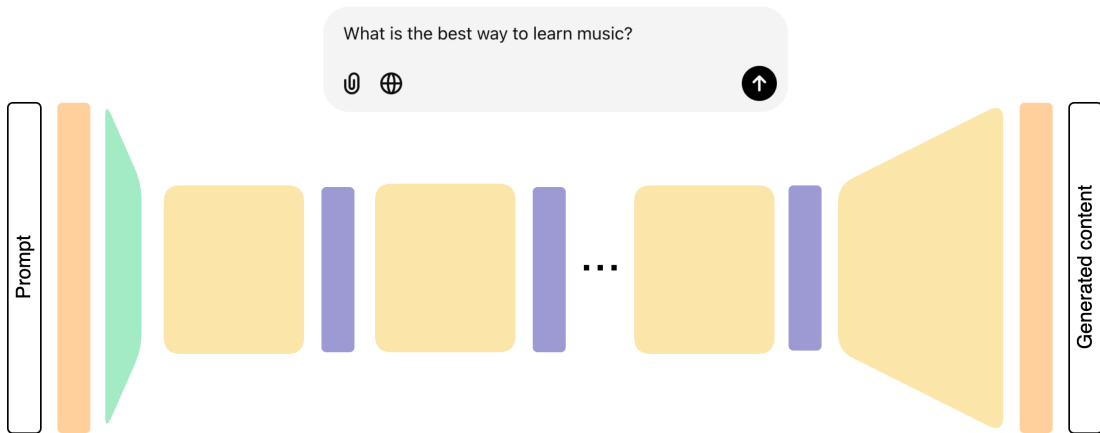
More info: [b-gendron.github.io](https://b-gendron.github.io)

# Tutorial Objectives

What to be learnt

What to be practiced

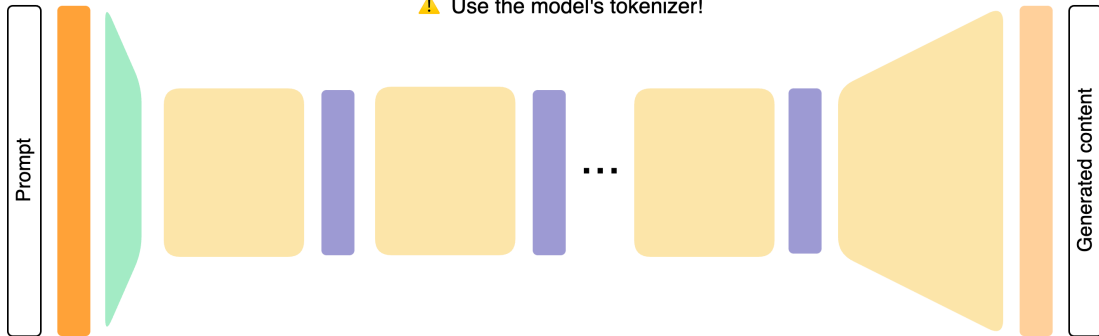
## ① Prompt the LLM



## ② Tokenize the Prompt Content

"What is the best way to learn music?"  
[531, 9, 45, 22, 3316, 2444, 34, 2172, 334]

⚠ Use the model's tokenizer!



## ③ Apply an Embedding layer

`Embedding(vocab_size, hidden_size)`

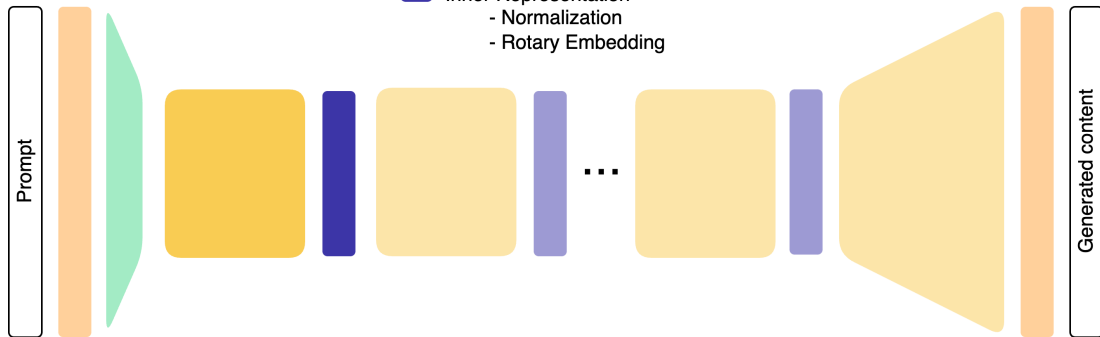
$$i \mapsto (e_1, \dots, e_{\text{hidden\_size}})$$



## ④ Going Through a Decoder Block

- Transformer Decoder Blocks
  - Attention layers (query, key, value)
  - Dense layers (Multi Layer Perceptron)

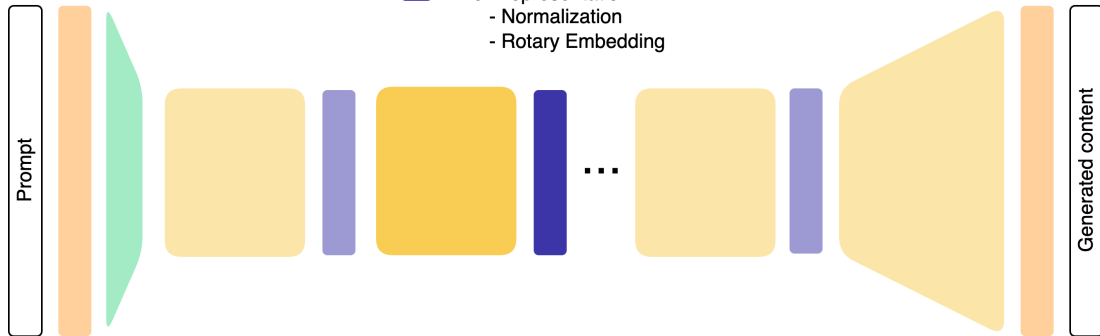
- Inner Representation
  - Normalization
  - Rotary Embedding



## ④ Going Through a Decoder Block

- Transformer Decoder Blocks
  - Attention layers (query, key, value)
  - Dense layers (Multi Layer Perceptron)

- Inner Representation
  - Normalization
  - Rotary Embedding

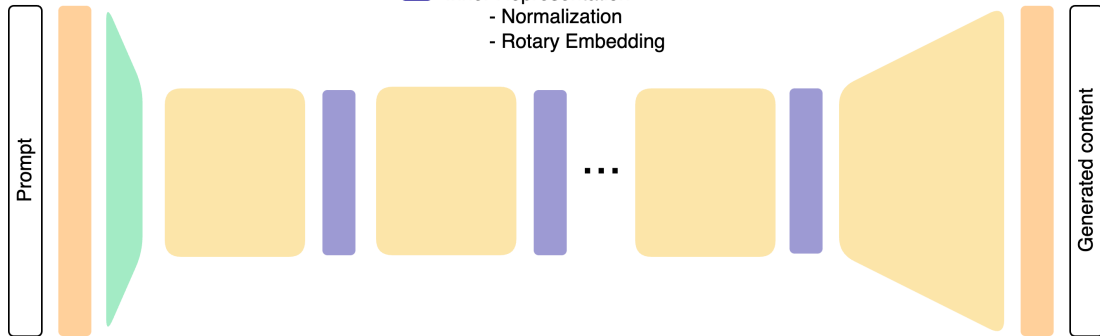




## ④ Going Through a Decoder Block

- Transformer Decoder Blocks
  - Attention layers (query, key, value)
  - Dense layers (Multi Layer Perceptron)

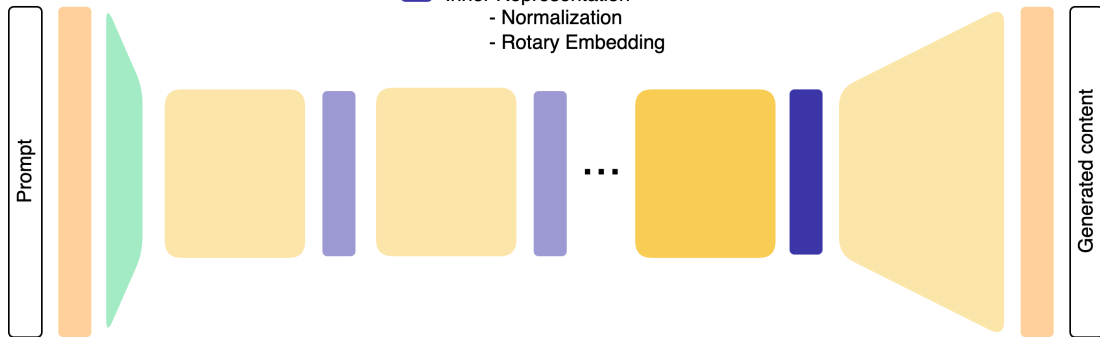
- Inner Representation
  - Normalization
  - Rotary Embedding



## ④ Going Through a Decoder Block

- Transformer Decoder Blocks
  - Attention layers (query, key, value)
  - Dense layers (Multi Layer Perceptron)

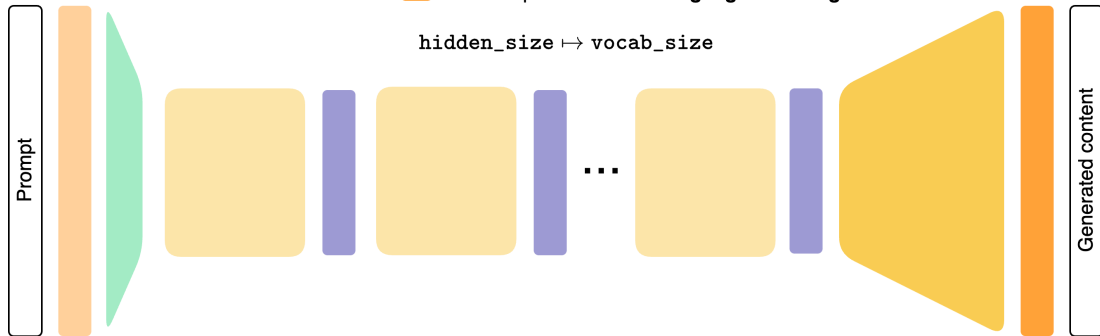
- Inner Representation
  - Normalization
  - Rotary Embedding



## ⑤ Out of Last Decoder Block

- Transformer Decoder Blocks
  - Attention layers (query, key, value)
  - Dense layers (Multi Layer Perceptron)
- Inner Representation + **Language Modeling Head**

`hidden_size`  $\mapsto$  `vocab_size`



## ⑥ Decode Generated Tokens in Natural Language

What is the best way to learn music?



The best way to learn music depends on your goals, interests, and learning style, but here are some effective strategies that can help:

