


Chapter 2


Theoretical Foundations of Large Language Models

Yashodeep Bharat Deshmukh

 <https://orcid.org/0009-0007-4873-2269>

Defence Institute of Advanced Technology, India

Abhishek Mukhopadhyay

 <https://orcid.org/0000-0002-4341-0523>

Amity University, Kolkata, India

ABSTRACT

Large Language Models (LLMs) are powerful tools in natural language processing with complex principles. This chapter delves into the theoretical foundations of LLMs, covering essential concepts, algorithms, and intuition, aiming to offer readers a strong grasp of the framework that empowers LLMs to perform varied language tasks, thus laying groundwork for advanced applications. Beginning with NLP fundamentals and statistical techniques, the chapter traces the evolution of language models from n -gram models to dense word embeddings like Word2Vec and GloVe. It then examines neural network architectures, progressing Transformers, dissecting key components such as attention mechanisms and encoder-decoder structures, including models like BERT and GPT. The chapter emphasizes the mathematical and algorithmic principles enabling LLMs' capabilities. The chapter synthesizes key theoretical principles behind LLMs, highlights their strengths and limitations, and explores future research areas like multimodal and few-shot learning, offering a comprehensive grounding for readers.

DOI: 10.4018/979-8-3693-8387-2.ch002

Copyright © 2026, IGI Global Scientific Publishing. Copying or distributing in print or electronic forms without written permission of IGI Global Scientific Publishing is prohibited. Use of this chapter to train generative artificial intelligence (AI) technologies is expressly prohibited. The publisher reserves all rights to license its use for generative AI training and machine learning model development.

INTRODUCTION

Large Language Models (LLMs) have revolutionized the field of natural language processing (NLP) in recent years. These powerful neural networks, trained on vast amounts of text data, have demonstrated remarkable capabilities in language understanding, generation, and reasoning. From chatbots and machine translation to question answering and content creation, LLMs are driving breakthroughs across a wide range of NLP applications.

This chapter delves into the theoretical foundations underpinning LLMs, providing a comprehensive overview of the key concepts, techniques, and architectures that enable their impressive performance. We begin by exploring the fundamentals of language modeling, including text representation, tokenization, and various statistical and neural network-based approaches. Next, we examine the transformative role of the attention mechanism in LLMs, highlighting its ability to capture long-range dependencies and enable efficient parallelization.

Building upon these foundational concepts, we introduce the transformer architecture, which has become the backbone of state-of-the-art LLMs. We discuss the core components of transformers, such as the encoder-decoder structure, positional encoding, and layer normalization, and explore how these elements contribute to their scalability and effectiveness.

We then dive into the world of LLMs and their variants, including encoder-only models like BERT, (Vaswani et al., 2017), decoder-only models like the GPT family, and encoder-decoder models such as T5. We examine the scaling laws and model sizing considerations that have propelled the development of increasingly large and capable LLMs, as well as their impressive few-shot and zero-shot learning abilities.

Finally, we conclude by discussing the current challenges and future directions in LLM research and development. We explore ethical considerations, such as bias and fairness, computational efficiency and environmental impact, and emerging architectures and techniques that promise to push the boundaries of what LLMs can achieve.

By the end of this chapter, readers will have a solid understanding of the theoretical foundations that underpin the remarkable advances in LLMs. Armed with this knowledge, they will be well-equipped to explore the exciting possibilities and tackle the challenges that lie ahead in this rapidly evolving field.

40 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/chapter/theoretical-foundations-of-large-language-models/390561

Related Content

Batched Computing

(2023). *Developing Linear Algebra Codes on Modern Processors: Emerging Research and Opportunities* (pp. 138-160).

www.irma-international.org/chapter/batched-computing/313456

Building With LLMs: Strategies for Real-World Business Implementation

Charvi Sanjay Suri, Isha Sanjay Suri and Gaurav M. Divtelwar (2026). *Leveraging LLMs for Business Innovation: Practical Solutions and Future Trends* (pp. 41-68).

www.irma-international.org/chapter/building-with-llms/401811

Technology in Education for People with Disabilities: Coding for Inclusion and Accessible Learning Design

Abhinav Varshney, Anushka Samaiya, Shrinu Varshney, Kunj Bihari Meena and Vipin Tyagi (2026). *Effective Coding Skill Development in Education* (pp. 253-290).

www.irma-international.org/chapter/technology-in-education-for-people-with-disabilities/398441

LLMs in ERP: Enhancing Financial Reporting

Satwik Jambula (2026). *Leveraging LLMs for Business Innovation: Practical Solutions and Future Trends* (pp. 227-250).

www.irma-international.org/chapter/llms-in-erp/401817

Empowering Scientific Computing and Data Manipulation With Numerical Python (NumPy)

Tesfaye Fufa Gedefa, Galety Mohammed Gouse and Garamu Tilahun Iticha (2023). *Advanced Applications of Python Data Structures and Algorithms* (pp. 147-161).

www.irma-international.org/chapter/empowering-scientific-computing-and-data-manipulation-with-numerical-python-numpy/326082