

# A report on LLMs

## Introduction

Built on the transformer architecture, Large Language Models (LLMs) are advanced AI systems trained on vast amounts of text data to understand, generate, and manipulate human language, in order to perform tasks like text generation, translation, summarization, and question answering. Their ability to grasp context and semantics has turned them into a main asset to make AI more accessible. Thus, models like GPT-4, LLaMA, and PaLM have become central to AI applications across industries.

This report aims at analyzing and synthesizing the contributions, differences, similarities, and limitations of 4 research papers related to Large Language Models (LLMs)..

## Themes /topics dealt with in the four selected papers are:

- An introduction to Large language Model
- The building of an AI-powered newsletter
- Language Models, are they reasonable?
- The Behind ChatGPT and Large Language Models,

## Titles and publication sources of the papers

1. *Introduction to Large Language Models (LLMs)*, Medium (published in *Agentman* blog series)
2. *What I Learned by Building an AI-Driven Newsletter*: Toloka Blog / Toloka AI
3. *Unreasonable Language Models?* Medium.com
4. *The Story Behind ChatGPT and Large Language Models*, Medium.com

## Paper Summaries

For each paper (3 to 5 total):

- Provide the full citation (author, year, title, venue).
- Summarize the research problem, proposed solution, and main results.
- Mention datasets used, model architecture, and evaluation metrics.

## Paper summaries

### ***Introduction to Large Language Models (LLMs)***

This article provides a beginner-friendly introduction to large language models (LLMs), aimed at full-stack developers. It explains what LLMs are, how they evolved from GPT to GPT-4, and their core underlying technologies like transformers, attention mechanisms, and tokenization. The author emphasizes how LLMs have democratized AI development, making it more accessible to non-experts. It outlines common applications such as Q&A systems, summarization, sentiment analysis, and translation. While the article is educational rather than research-focused, it lays the groundwork for building practical LLM-based systems. It highlights the transformative role of LLMs in modern AI and sets the stage for future articles on implementation and fine-tuning.

Prasad Thammineni, 2023, *Introduction to Large Language Models (LLMs)*: Medium / Agentman

## ***What I Learned by Building an AI-Driven Newsletter***

The article discusses Dr. Jack Saunders' experience building an AI-powered newsletter that summarizes new research papers. The challenge was building a system that reliably automates content curation and summarization using LLMs. He experimented with agent-based approaches but found them error-prone and unnecessarily complex. A fixed workflow, using APIs and prompt-based LLMs, proved more effective. Key insights included the non-determinism of LLMs, hallucinations, and the limits of even well-crafted prompts. Tools like ChatGPT, Brevo, and the arXiv API were leveraged. Human oversight remained essential for catching mistakes. The project highlighted LLMs' utility as assistants rather than autonomous agents. Saunders emphasizes practical use over hype and calls for simplicity and error resilience.

Dr. Jack Saunders, 2023, *What I Learned by Building an AI-Driven Newsletter*; Toloka AI Blog, July 26, 2023

## ***Unreasonable Language Models?***

The article discusses the limitations of current large language models (LLMs) in reasoning tasks, referencing Apple's paper "*The Illusion of Thinking*." It critiques the assumption that increasing model size and inference time leads to better reasoning. Through puzzles like Tower of Hanoi and River Crossing, the Apple study shows that LLMs perform well on medium-complexity tasks but fail on high-complexity ones. Even with ample tokens, models often reduce output rather than increase reasoning depth. It questions current evaluation metrics that prioritize final answers over reasoning steps. Models often generate correct algorithms but fail to apply them to large-scale inputs. Tool limitations (e.g., no calculator) hinder model reasoning. Additionally, lack of web content might influence performance. Overall, it challenges whether LLMs truly "reason" or just excel at pattern recognition.

Apoorv Jain, 2025, *Unreasonable Language Models?* Medium (Online publication platform)

## ***The Story Behind ChatGPT and Large Language Models***

This article traces the evolution of large language models (LLMs) from early probabilistic models to today's advanced systems like ChatGPT. It highlights the limitations of early models in handling long-range dependencies and explains the transformative impact of the Transformer architecture. The LLM development pipeline includes three main phases: **Pretraining**, **Supervised Fine-Tuning (SFT)**, and **Reinforcement Learning (RL)**. Pretraining captures general language understanding, SFT teaches task-specific behavior, and RL aligns responses with human preferences. The article emphasizes the importance of reward models in RL and the risk of "reward hacking", where models exploit flaws in the system. It concludes by underlining the complexity of aligning LLMs with human values, not just scaling their size.

Abderrahouf Lahmar, 2025, *The Story Behind ChatGPT and Large Language Models*, Medium

## **Comparative Analysis**

Compare the papers across key aspects such as:

- Objectives and problem domains
- Model architectures and innovations
- Training or fine-tuning strategies
- Benchmarks and evaluation
- Strengths, limitations, and reproducibility

Use tables or charts if helpful for comparison.

## Strengths, Limitations, and Reproducibility

- **Strengths:** Structured training pipeline; high language fluency
- **Limitations:** Reward hacking; sparse or flawed reward signals
- **Reproducibility:** Not detailed; high-resource requirement for replication

## Insights and Reflection

- What trends or patterns emerge across the papers?
- Which methods or approaches seem most promising or innovative?
- What limitations or challenges are commonly acknowledged?
- What are potential future directions in this research area?

## Trends or Patterns

- Clear pipeline evolution: Pretraining → SFT → RLHF
- Increasing focus on human alignment, not just model performance

## Most Promising Methods

- Reinforcement Learning from Human Feedback (RLHF)
- Transformer-based architectures for scalability and efficiency

## Common Limitations

- Misalignment due to flawed reward signals
- Reward hacking, where models game the system
- Lack of true understanding or reasoning ability in LLMs

## Potential Future Directions

- Better reward modeling for nuanced feedback
- Hybrid models combining symbolic reasoning and LLMs
- More interpretable and controllable LLM behaviour

A Comparison of 4 papers dealing with Large Language Models (LLM)				
Title & Author	<i>Introduction to Large Language Models (LLMs)</i> , Prasad Thammineni	<i>What I Learned by Building an AI-Driven Newsletter</i> , Dr. Jack Saunders	<i>Unreasonable Language Models?</i> Apoorv Jain	<i>The Story Behind ChatGPT and Large Language Models</i> , Abderraouf Lahmar
Key aspects				
Dataset	Not explicitly mentioned; generally refers to large text corpora	arXiv papers, news articles via TheNewsAPI, email content via Gmail API	Synthetic reasoning problems (e.g., Tower of Hanoi, River Crossing) with adjustable complexity.	Pretraining: Large-scale text corpora (books, internet, articles) SFT: Task-specific datasets (e.g., math problems with solutions) RL: Human preference-based feedback data

<b>Model Architecture</b>	Transformer (based on attention mechanism)	LLMs (likely GPT-based) accessed via prompting	Not explicitly stated, but implies standard LLMs and reasoning-augmented LLMs.	Transformer (based on <i>Attention is All You Need</i> , 2017)
<b>Evaluation Metrics</b>	Not discussed in detail; to be covered in future articles	Not formalized — manual review of newsletter quality and JSON formatting accuracy Note: This was more of a system-building case study than an empirical ML research paper	Success rate on problem solutions, token output analysis, step-by-step reasoning accuracy.	Coherence, helpfulness, harmlessness (qualitative) Reward signal performance during RL Human feedback alignment
<b>Objectives and Problem Domains</b>	Educate developers about LLMs and their potential applications  Bridge the gap between full-stack development and AI	Automate paper/news summarization for a daily newsletter  Build a reliable, semi-autonomous pipeline using LLMs	Understand limitations of LLMs in structured reasoning tasks. Evaluate generalization and stepwise problem-solving under increasing complexity.	Explain how LLMs like ChatGPT are trained Break down the three-stage process of language model development Address the challenges of aligning models with human values
<b>Model Architectures &amp; Innovations</b>	Focus on Transformers, attention mechanisms, and tokenization Overview of OpenAI's GPT series and mention of Meta's LLaMA, Google's PaLM	Prompt-driven LLMs Avoided autonomous agents in favor of structured workflows	Comparison between standard LLMs and reasoning-optimized models. Focus on models' ability to simulate recursive and algorithmic thought.	Model Architectures and Innovations Transformer architecture enables better long-range dependency handling RLHF (Reinforcement Learning from Human Feedback) improves alignment
<b>Training or Fine-Tuning Strategies</b>	Describes fine-tuning as an accessible method for adapting LLMs to specific tasks More in-depth strategies are promised in upcoming articles	No custom training; relied on pre-trained LLMs and API calls	Not explicitly discussed; assumes standard pre-trained LLMs.	Self-supervised pretraining for language structure Supervised fine-tuning on domain-specific tasks Reinforcement learning to refine model outputs based on human preferences

<b>Benchmarks &amp; Evaluation</b>	Not detailed in this article — intended for future parts of the series	Manual quality control and prompt structure validation Focus on robustness over formal metrics	Reasoning tasks split by difficulty (Low, Medium, High). Performance compared on stepwise output and final results	Human-centric evaluation of helpfulness, clarity, and safety Sensitivity to user intent and format adherence Behaviour under ethical constraints
<b>Strengths, Limitations, and Reproducibility</b>				
<b>Strengths</b>	Clear and accessible explanation of LLMs and their uses	Real-world insights into LLM deployment Emphasis on human-in-the-loop validation Use of existing APIs saved time	Highlights nuanced failures; real task modeling.	Structured training pipeline; high language fluency
<b>Limitations</b>	Lacks technical depth and empirical data	Reliance on external tools limits generalizability Prompt instability and JSON formatting issues Lack of formal evaluation metrics	Lacks tool access (e.g., calculator); strict output formatting affects fairness	Reward hacking; sparse or flawed reward signals
<b>Reproducibility</b>	Not applicable — this is a conceptual overview, not an experiment	High, as it uses publicly available APIs and tools	Dependent on synthetic problem generators and clear complexity tuning	Not detailed; high-resource requirement for replication
<b>Insights and Reflection</b>				
<b>Trends or Patterns</b>	Focus on making LLMs more accessible to broader developer communities Emphasis on real-world applications like Q&A, summarization, and translation	LLMs are increasingly used in workflows with human oversight Shift from large, monolithic agents to modular pipelines	LLMs excel at moderate reasoning but fail at higher complexities. Reasoning capability does not scale linearly with tokens or model size.	Clear pipeline evolution: Pretraining → SFT → RLHF Increasing focus on human alignment, not just model performance
<b>Most Promising Methods</b>	Transformer architecture continues to dominate LLM development Fine-tuning pretrained models for task-specific applications	RLHF for safety and customization Prompt engineering with robust error handling Lightweight workflows over complex agent-based systems	Focusing on intermediate reasoning steps and algorithm simulation. Building models that integrate tool use (e.g., calculators or code interpreters).	Reinforcement Learning from Human Feedback (RLHF) Transformer-based architectures for scalability and efficiency

<b>Common Limitations or Challenges</b>	The article doesn't explore limitations directly but implies a gap in developer understanding Future challenges may involve model selection, deployment, and evaluation	Hallucinations, bias, formatting errors Fragile prompts and overreliance on LLM behavior consistency	Models failing to "dry run" logical steps for large inputs. Evaluation overly focused on final answers. Lack of interpretability and consistency in reasoning depth.	Misalignment due to flawed reward signals Reward hacking, where models game the system Lack of true understanding or reasoning ability in LLMs
<b>Potential Future Directions</b>	Building full-stack LLM applications (starting with Q&A systems) Exploring system architecture, best practices, and evaluation in depth	Smaller, efficient models (e.g., Chinchilla, RETRO) Self-verifying and steerable LLMs Integration with retrieval and API systems Emphasis on human values and alignment (via RLHF)	Incorporate external tools to simulate real-world problem-solving. Shift evaluation to multi-step reasoning traceability. Focus on curriculum learning or hybrid neuro-symbolic models for deep reasoning.	Future progress hinges on improving training strategies and evaluation methods to ensure models behave ethically, safely, and usefully.

## Conclusion

### Meta-Analysis of LLM Articles

These four articles collectively highlight the rapid evolution and growing capabilities of Large Language Models (LLMs). Key findings emphasize the foundational role of the transformer architecture, the importance of pretraining followed by supervised fine-tuning and reinforcement learning, and the ongoing challenge of aligning LLM outputs with human intent. While LLMs excel at language understanding and generation, their reasoning abilities remain limited, especially in complex tasks. Reward hacking, lack of interpretability, and reliance on large datasets are recurring challenges. However, innovations like hierarchical reasoning, fine-tuning strategies, and tool integration (e.g., calculators, code execution) show promise. The field is moving toward more robust, aligned, and generalizable models. As LLMs continue to integrate into real-world systems, their evolution is increasingly shaped by user interaction, safety, and ethical deployment.