

# Deep Dive into Six Recent Research Papers in AI

The field of artificial intelligence (AI) is rapidly evolving, with new research papers and models being released at a breakneck pace. This article provides a deep dive into six recent research papers that explore different aspects of AI, from improving safety and robustness to scaling foundation models and reasoning over public and private data. The papers covered are:

- OpenAI o1 System Card
- Titans: Learning to Memorize at Test Time
- MiniMax-01: Scaling Foundation Models with Lightning Attention
- DeepSeek-V3 Technical Report
- Reasoning over Public and Private Data in Retrieval-Based Systems
- Attention Is All You Need

## OpenAI o1 System Card

The OpenAI o1 System Card, dated December 5th, 2024 <sup>1</sup>, focuses on the safety work conducted for OpenAI's latest large language models, o1 and o1-mini. These models are trained with reinforcement learning to perform complex reasoning, enabling them to "think" before answering and produce a chain of thought before responding to the user. This deliberative process allows the models to refine their approach, try different strategies, and recognize their mistakes.

## Authors

The OpenAI o1 System Card was authored by Aaron Jaech and 260 other authors<sup>1</sup>.

## Key Contributions and Findings

- **Improved Safety:** The o1 models demonstrate significant improvements in safety and robustness compared to previous models. They achieve substantial improvements on jailbreak evaluations and are more aligned with OpenAI's content guidelines<sup>1</sup>.
- **Chain-of-Thought Reasoning:** The models utilize chain-of-thought reasoning, which allows them to reason about safety policies in context when responding to potentially unsafe prompts<sup>2</sup>.
- **Focus on Deliberate Reasoning:** The o1 family represents a shift from fast, intuitive thinking to slower, more deliberate reasoning, which has implications for both safety and alignment<sup>1</sup>.
- **Strong Performance:** OpenAI o1 ranks in the 89th percentile on competitive programming questions (Codeforces), places among the top 500 students in the US in a qualifier for the USA Math Olympiad (AIME), and exceeds human PhD-level accuracy on a benchmark of physics, biology, and chemistry problems (GPQA)<sup>4</sup>.

## Areas of Improvement

While the o1 models show promise, the system card also highlights areas for improvement:

- **Autonomous Capabilities:** Complex machine learning research automation remains a challenge. While the models show some initial success in autonomous tasks, such as generating code to deactivate an oversight mechanism or exfiltrate model weights, these capabilities require further refinement and careful monitoring<sup>5</sup>.
- **Safety Considerations:** The models have a medium risk rating in specific domains, including Cybersecurity, CBRN (Chemical, Biological, Radiological, and Nuclear), and Persuasion. This necessitates ongoing vigilance and continuous monitoring to ensure responsible use and mitigate potential risks<sup>2</sup>.

## Titans: Learning to Memorize at Test Time

Traditional transformer architectures face limitations in handling long contexts due to their quadratic complexity, which leads to increased computational demands as the context length grows<sup>7</sup>. Linear transformers, while offering improved efficiency, often struggle to maintain competitive performance<sup>8</sup>. To address these challenges, the paper "Titans: Learning to Memorize at Test Time"<sup>9</sup> introduces a new approach to long-term memory in neural networks.

## Authors

The authors of "Titans: Learning to Memorize at Test Time" are Ali Behrouz, Peilin Zhong, and Vahab S. Mirrokni<sup>10</sup>.

## Key Contributions and Findings

- **Neural Long-Term Memory:** The paper presents a deep neural long-term memory module that acts as a meta in-context model, learning to store data in its parameters at test time<sup>9</sup>.
- **Surprise-Based Memorization:** Inspired by the human long-term memory system, the memory module is designed to make surprising events more memorable. The surprise of an input is measured using the gradient of the neural network with respect to the input in associative memory loss<sup>9</sup>.
- **Decay Mechanism:** To manage limited memory, the paper proposes a decay mechanism that considers the ratio of memory size to data surprise, leading to improved memory management. Interestingly, this mechanism aligns with optimizing a meta-neural network using mini-batch gradient descent, momentum, and weight decay<sup>8</sup>.
- **Titans Architecture:** The paper introduces a new family of architectures called Titans, which integrate short-term, long-term, and persistent memory components. This architecture draws inspiration from the human memory system, where different types of memory interact to facilitate learning and retrieval<sup>11</sup>.

## Experimental Results

Experimental results on various tasks, including language modeling, common-sense reasoning, genomics, and time series tasks, show that Titans outperform Transformers and recent modern linear recurrent models. They can effectively scale to larger context window sizes with higher accuracy in needle-in-haystack tasks. For instance, in a benchmark called Single Nih, designed to test the ability to retrieve information from long sequences, Titans consistently achieved the highest accuracy scores, even when dealing with sequences of up to 16,000 tokens<sup>13</sup>.

## MiniMax-01: Scaling Foundation Models with Lightning Attention

The "MiniMax-01: Scaling Foundation Models with Lightning Attention" paper<sup>15</sup>, published in 2025<sup>16</sup>, introduces the MiniMax-01 series of models, including MiniMax-Text-01 and MiniMax-VL-01. These models are designed to handle longer contexts and achieve comparable performance to top-tier models.

## Authors

The authors of "MiniMax-01: Scaling Foundation Models with Lightning Attention" include Aonian Li, Bangwei Gong, and numerous other researchers from MiniMaxAI<sup>16</sup>.

## Key Contributions and Findings

- **Lightning Attention:** The models utilize lightning attention, a linear attention variant, for efficient scaling and handling of long contexts<sup>15</sup>.
- **Hybrid Architecture:** MiniMax-01 adopts a hybrid architecture that combines lightning attention, softmax attention, and Mixture-of-Experts (MoE). This allows the model to balance the precision of softmax attention with the speed and efficiency of lightning attention for different tasks<sup>16</sup>.
- **Long Context Window:** Leveraging advanced parallel strategies and innovative compute-communication overlap methods, such as Linear Attention Sequence Parallelism Plus (LASP+), varlen ring attention, and Expert Tensor Parallel (ETP), MiniMax-01's training context length is extended to 1 million tokens, and it can handle a context of up to 4 million tokens during inference<sup>16</sup>.
- **Mixture of Experts:** The model incorporates MoE with 32 experts, where each expert specializes in a specific area. This specialization allows the model to handle a much wider range of tasks with greater efficiency<sup>18</sup>.

## Evaluation

MiniMax-01 demonstrates strong performance on various academic benchmarks, comparable to top-tier models like GPT-4o and Claude-3.5-Sonnet<sup>15</sup>. The hybrid architecture and the use of

lightning attention contribute to its ability to handle long contexts efficiently<sup>20</sup>.

## Mixture-of-Experts (MoE)

Both MiniMax-01 and DeepSeek-V3 utilize a Mixture-of-Experts (MoE) architecture. This approach involves dividing a large neural network into smaller, specialized networks called "experts." Each expert focuses on a specific domain or task, and a routing mechanism determines which experts are most relevant for a given input. This specialization allows MoE models to handle a wider range of tasks with greater efficiency and scalability.

## DeepSeek-V3 Technical Report

The DeepSeek-V3 Technical Report <sup>21</sup> presents DeepSeek-V3, a powerful Mixture-of-Experts (MoE) language model with 671 billion total parameters.

### Key Contributions and Findings

- **Efficient Architecture:** DeepSeek-V3 adopts Multi-head Latent Attention (MLA) and DeepSeekMoE architectures for efficient inference and cost-effective training<sup>21</sup>.
- **Auxiliary-Loss-Free Load Balancing:** The model pioneers an auxiliary-loss-free strategy for load balancing, minimizing performance degradation<sup>21</sup>.
- **Multi-Token Prediction:** DeepSeek-V3 sets a multi-token prediction training objective for stronger performance and faster inference<sup>21</sup>.
- **Stable Training:** The training process is remarkably stable, without any irrecoverable loss spikes or rollbacks<sup>21</sup>.
- **Enhanced Precision:** DeepSeek-V3 utilizes a fine-grained quantization method to handle outliers and strategically uses CUDA cores to enhance Floating-Point Arithmetic Gaussian Mixture Model (FPAGMM) precision<sup>22</sup>.

### Evaluation

DeepSeek-V3 outperforms other open-source models and achieves performance comparable to leading closed-source models. It requires only 2.788 million H800 GPU hours for its full training<sup>21</sup>. DeepSeek-V3 base excels in code and math, surpassing models with more parameters. It achieved a remarkable 91.6 F1 score on the D-drop three-shot setting, which focuses on long context understanding, a key factor in its success in code-related tasks<sup>22</sup>. However, it is noted that DeepSeek-V3 still lags behind Claude Sonnet 3.5 on some benchmarks, particularly in coding benchmarks such as SWE-bench<sup>23</sup>.

## Reasoning over Public and Private Data in Retrieval-Based Systems

The paper "Reasoning over Public and Private Data in Retrieval-Based Systems" <sup>24</sup>, published

on March 22, 2022<sup>24</sup>, addresses the challenge of reasoning over information that is split between publicly and privately accessible scopes.

## Authors

The authors of "Reasoning over Public and Private Data in Retrieval-Based Systems" are Simran Arora, Patrick Lewis, Angela Fan, Jacob Kahn, and Christopher Ré<sup>24</sup>.

## Key Contributions and Findings

- **PAIR Framework:** The paper defines the PUBLIC-PRIVATE AUTOREGRESSIVE INFORMATION RETRIEVAL (PAIR) privacy framework for retrieval over multiple privacy scopes. This framework aims to enable retrieval-based systems to utilize both public and private data while preserving privacy<sup>24</sup>.
- **ConcurrentQA Benchmark:** The authors create ConcurrentQA, the first textual QA benchmark that requires concurrent retrieval over multiple data distributions (public and private). This benchmark provides a valuable resource for evaluating the performance of retrieval-based systems in this setting<sup>25</sup>.
- **Privacy-Performance Tradeoffs:** The paper investigates the privacy vs. performance tradeoffs that arise when retrieving information from multiple privacy scopes and explores ways to mitigate these tradeoffs<sup>25</sup>.

## Implications

This research highlights the importance of considering privacy in retrieval-based systems, especially when dealing with sensitive information. The proposed framework and benchmark provide a foundation for developing privacy-preserving solutions in this domain<sup>26</sup>.

## Attention Is All You Need

The landmark paper "Attention Is All You Need"<sup>28</sup>, published on July 10, 2023<sup>28</sup>, introduced the Transformer architecture, a deep learning architecture based solely on attention mechanisms.

## Authors

The authors of "Attention Is All You Need" are Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan Gomez, Lukasz Kaiser, and Illia Polosukhin<sup>28</sup>.

## Key Contributions and Findings

- **Transformer Architecture:** The paper introduces the Transformer, a new network architecture that dispenses with recurrence and convolutions entirely, relying solely on attention mechanisms. This architecture allows for parallel processing of data, leading to faster training and inference compared to traditional recurrent neural networks (RNNs)<sup>28</sup>.

- **Attention Mechanisms:** The paper highlights the importance of attention mechanisms, specifically scaled dot-product attention and multi-head attention. These mechanisms allow the model to focus on different parts of the input data dynamically and capture different aspects of the relationships between words in a sequence<sup>28</sup>.
- **Positional Encoding:** The Transformer utilizes positional encoding to incorporate information about the order of words in a sequence, which is crucial for understanding language<sup>28</sup>.
- **Parallel Processing:** The Transformer enables parallel processing of data, leading to faster training and inference<sup>30</sup>.
- **Scalability:** The architecture can efficiently handle a wide range of data sizes and complexities<sup>30</sup>.

## Impact

This paper is considered foundational in modern AI, as the Transformer approach has become the main architecture for large language models like those based on GPT<sup>28</sup>. The attention mechanism has revolutionized sequence modeling and transduction models in various tasks<sup>31</sup>. The Transformer model, unlike traditional seq2seq models, does not rely on the sequence of the text to perform encoding and decoding. This characteristic may allow the model to extrapolate to sequence lengths longer than those encountered during training<sup>28</sup>.

## Comparison of the Papers

Feature	OpenAI o1 System Card	Titans	MiniMax -01	DeepSeek-V3	Reasoning over Public/Private Data	Attention Is All You Need
Authors	Aaron Jaech et al <sup>1</sup> .	Ali Behrouz et al <sup>10</sup> .	Aonian Li et al <sup>16</sup> .	DeepSeek AI Team <sup>21</sup>	Simran Arora et al <sup>24</sup> .	Ashish Vaswani et al <sup>28</sup> .
Publication Date	December 5th, 2024 (inferred) <sup>1</sup>	January 2025 <sup>9</sup>	2025 <sup>16</sup>	December 2024 <sup>21</sup>	March 22, 2022 <sup>24</sup>	July 10, 2023 <sup>28</sup>
Key Contrib	Improved safety	Neural long-ter	Lightning g	Efficient architect	PAIR frameworko	Transformer

Feature	OpenAI o1 System Card	Titans	MiniMax -01	DeepSeek-V3	Reasoning over Public/Private Data	Attention Is All You Need
Strengths	and robustness <sup>1</sup> , Chain-of-thought reasoning <sup>2</sup> , Strong performance on various tasks <sup>4</sup>	memory module <sup>9</sup> , Surprise-based memorization <sup>9</sup> , Decay mechanism <sup>8</sup> , Titans architecture <sup>11</sup>	attention <sup>15</sup> , Hybrid architecture <sup>16</sup> , Long context window <sup>18</sup> , Mixture-of-Experts (MoE) <sup>18</sup>	ure <sup>21</sup> , Auxiliary-loss-free load balancing <sup>21</sup> , Multi-token prediction <sup>21</sup> , Stable training <sup>21</sup> , Enhanced precision <sup>22</sup>	rk <sup>24</sup> , ConcurrentQA benchmark <sup>25</sup> , Privacy-performance tradeoffs <sup>25</sup>	architecture <sup>28</sup> , Scaled dot-product attention and multi-head attention <sup>28</sup> , Positional encoding <sup>28</sup>
Methodology	Reinforcement learning, chain-of-thought prompting <sup>1</sup>	Associative memory loss, decay mechanism <sup>9</sup>	Lightning attention, softmax attention, MoE <sup>16</sup>	MLA, DeepSeekMoE, FP8 mixed precision training <sup>21</sup>	Split Iterative Retrieval (SPIRAL) problem <sup>24</sup>	Scaled dot-product attention, multi-head attention, positional encoding <sup>28</sup>
Limitations	Challenges in autonomous capabilities <sup>5</sup> ,	Not explicitly stated in the provided material	Not explicitly stated in the provided material	Lags behind Claude Sonnet 3.5 on some	Not explicitly stated in the provided material	Not explicitly stated in the provided material

Feature	OpenAI o1 System Card	Titans	MiniMax -01	DeepSeek-V3	Reasoning over Public/Private Data	Attention Is All You Need
	Safety considerations in specific domains <sup>2</sup>			benchmarks <sup>23</sup>		
Future Research	Chain-of-thought monitoring <sup>1</sup> , Improving autonomous capabilities <sup>1</sup>	Exploring different memory structures, update mechanisms, and retrieval processes <sup>9</sup>	Applicability of lightning attention to other tasks <sup>20</sup> , Exploring limitations of lightning attention	Improving performance on specific benchmarks <sup>23</sup> , Reducing model size for easier self-hosting	Developing more sophisticated privacy-preserving solutions <sup>26</sup> , Exploring different performance-privacy cost models	Exploring limitations and potential improvements of the Transformer architecture <sup>6</sup>
Impact	Advances in AI safety and alignment	Improved long-context understanding and memory management	Efficient handling of long contexts in large language models	Open-source MoE model with strong performance, particularly in code and	Foundation for privacy-preserving retrieval systems	Foundational architecture for modern large language models, revolutionized sequenc



Feature	OpenAI o1 System Card	Titans	MiniMax -01	DeepSeek-V3	Reasoning over Public/Private Data	Attention Is All You Need
				math		e modeling

## Code Implementations

Code implementations and further details for some of the papers are available:

- **OpenAI o1 System Card:** Examples and code snippets are available in the OpenAI Cookbook<sup>32</sup>.
- **Titans:** A PyTorch implementation is available on GitHub<sup>10</sup>.
- **MiniMax-01:** The model is open-sourced on GitHub<sup>33</sup>.
- **DeepSeek-V3:** The model and code are available on GitHub<sup>22</sup>.
- **Attention Is All You Need:** Various implementations and tutorials are available online, including PyTorch implementations on GitHub and Kaggle<sup>34</sup>.

## Open Questions and Future Research

The papers also raise several open questions and areas for future research:

- **OpenAI o1:** Further research is needed on chain-of-thought deception monitoring and improving autonomous capabilities<sup>1</sup>.
- **Titans:** Exploring different memory structures, update mechanisms, and retrieval processes could lead to further improvements<sup>9</sup>.
- **MiniMax-01:** Investigating the applicability of lightning attention to other tasks and exploring its limitations are important research directions<sup>20</sup>.
- **DeepSeek-V3:** Further research on improving performance on specific benchmarks and reducing model size for easier self-hosting is needed<sup>23</sup>.
- **Reasoning over Public/Private Data:** Developing more sophisticated privacy-preserving solutions and exploring different performance-privacy cost models are crucial<sup>26</sup>.
- **Attention Is All You Need:** While the Transformer architecture has been highly successful, research on its limitations and potential improvements continues<sup>6</sup>.

## Conclusion

These six research papers provide valuable insights into the current state of AI and highlight key areas of innovation. They demonstrate the ongoing efforts to improve the safety, robustness, efficiency, and scalability of AI models, while also addressing important considerations like

privacy. The OpenAI o1 System Card emphasizes the growing importance of safety and robustness in AI development, while Titans and MiniMax-01 showcase the increasing focus on efficiency and scalability in AI models, particularly for handling long contexts. DeepSeek-V3 contributes a powerful open-source MoE model with strong performance, and "Reasoning over Public and Private Data" addresses the crucial need for privacy-preserving AI systems. Finally, "Attention Is All You Need" has had a lasting impact on the field of AI, with the Transformer architecture and attention mechanisms becoming foundational elements in modern large language models.

One emerging trend is the combination of different AI techniques, as seen in the hybrid architectures of MiniMax-01, which blends lightning attention, softmax attention, and MoE. This approach allows models to leverage the strengths of different techniques to achieve better overall performance. Another notable trend is the increasing emphasis on open-source AI models, as exemplified by DeepSeek-V3. This promotes accessibility, collaboration, and innovation in the AI community.

These papers collectively contribute to the advancement of AI by pushing the boundaries of what's possible in terms of safety, efficiency, scalability, and privacy. The findings and open questions presented in these papers will undoubtedly shape future research and development in the field, leading to more powerful, responsible, and beneficial AI systems.

## Works cited

1. OpenAI o1 System Card, accessed January 16, 2025, <https://cdn.openai.com/o1-system-card-20241205.pdf>
2. OpenAI o1 System Card, accessed January 16, 2025, <https://openai.com/index/openai-o1-system-card/>
3. [2412.16720] OpenAI o1 System Card - arXiv, accessed January 16, 2025, <https://arxiv.org/abs/2412.16720>
4. Introducing OpenAI o1, accessed January 16, 2025, <https://openai.com/o1/>
5. Deep Dive: OpenAI's o1 - The Dawn of Deliberate AI, accessed January 16, 2025, <https://portkey.ai/blog/openai-o1-model-card-analysis/>
6. OpenAI o1 system card - Hacker News, accessed January 16, 2025, <https://news.ycombinator.com/item?id=42330666>
7. [QA] Titans: Learning to Memorize at Test Time - YouTube, accessed January 16, 2025, <https://www.youtube.com/watch?v=KbTxzlGWsAc>
8. Titans: Learning to Memorize at Test Time - YouTube, accessed January 16, 2025, <https://www.youtube.com/watch?v=hqTF9Cj0XYw>
9. Titans: Learning to Memorize at Test Time - arXiv, accessed January 16, 2025, <https://arxiv.org/html/2501.00663v1>
10. Unofficial implementation of Titans, SOTA memory for transformers, in Pytorch - GitHub, accessed January 16, 2025, <https://github.com/lucidrains/titans-pytorch>
11. [2501.00663] Titans: Learning to Memorize at Test Time - arXiv, accessed January 16, 2025, <https://arxiv.org/abs/2501.00663>
12. Titans: Learning to Memorize at Test Time - ResearchGate, accessed January 16, 2025, [https://www.researchgate.net/publication/387671240\\_Titans\\_Learning\\_to\\_Memorize\\_at\\_Test\\_Time](https://www.researchgate.net/publication/387671240_Titans_Learning_to_Memorize_at_Test_Time)
13. Titans: Learning to Memorize at Test Time : r/LocalLLaMA - Reddit, accessed January 16,

2025,

[https://www.reddit.com/r/LocalLLaMA/comments/1i0q8nw/titans\\_learning\\_to\\_memorize\\_at\\_test\\_time/](https://www.reddit.com/r/LocalLLaMA/comments/1i0q8nw/titans_learning_to_memorize_at_test_time/)

14. Titans: Learning to Memorize at Test Time - YouTube, accessed January 16, 2025, <https://www.youtube.com/watch?v=vUbTQ9EwFXU>

15. [2501.08313] MiniMax-01: Scaling Foundation Models with Lightning Attention - arXiv, accessed January 16, 2025, <https://arxiv.org/abs/2501.08313>

16. MiniMaxAI/MiniMax-Text-01 - Hugging Face, accessed January 16, 2025, <https://huggingface.co/MiniMaxAI/MiniMax-Text-01>

17. Paper page - MiniMax-01: Scaling Foundation Models with Lightning Attention, accessed January 16, 2025, <https://huggingface.co/papers/2501.08313>

18. MiniMax 01: Scaling Foundation Models with Lightning Attention - YouTube, accessed January 16, 2025, <https://www.youtube.com/watch?v=ul2pds22jLE>

19. MiniMax-01: Scaling Foundation Models with Lightning Attention - YouTube, accessed January 16, 2025, <https://www.youtube.com/watch?v=AxzE0b65BW0>

20. [2501.08313] MiniMax-01: Scaling Foundation Models with Lightning Attention : r/LocalLLaMA - Reddit, accessed January 16, 2025, [https://www.reddit.com/r/LocalLLaMA/comments/1i1ntmb/250108313\\_minimax01\\_scaling\\_foundation\\_models/](https://www.reddit.com/r/LocalLLaMA/comments/1i1ntmb/250108313_minimax01_scaling_foundation_models/)

21. DeepSeek-V3 Technical Report - arXiv, accessed January 16, 2025, <https://arxiv.org/html/2412.19437v1>

22. DeepSeek-V3 Technical Report - YouTube, accessed January 16, 2025, <https://www.youtube.com/watch?v=2PrkHkbDDyU>

23. Deepseek V3 is officially released (code, paper, benchmark results) : r/LocalLLaMA - Reddit, accessed January 16, 2025, [https://www.reddit.com/r/LocalLLaMA/comments/1hmmmt3/deepseek\\_v3\\_is\\_officially\\_released\\_code\\_paper/](https://www.reddit.com/r/LocalLLaMA/comments/1hmmmt3/deepseek_v3_is_officially_released_code_paper/)

24. reasoning over public and private data in retrieval-based systems - arXiv, accessed January 16, 2025, <https://arxiv.org/pdf/2203.11027>

25. [2203.11027] Reasoning over Public and Private Data in Retrieval-Based Systems - arXiv, accessed January 16, 2025, <https://arxiv.org/abs/2203.11027>

26. Reasoning over Public and Private Data in Retrieval-Based Systems - MIT Press Direct, accessed January 16, 2025, [https://direct.mit.edu/tacl/article/doi/10.1162/tacl\\_a\\_00580/117168/Reasoning-over-Public-and-Private-Data-in](https://direct.mit.edu/tacl/article/doi/10.1162/tacl_a_00580/117168/Reasoning-over-Public-and-Private-Data-in)

27. Knowledge Retrieval Over Public and Private Data - GitHub Pages, accessed January 16, 2025, <https://knowledge-nlp.github.io/aaai2023/papers/003-PQA-oral.pdf>

28. Attention Is All You Need - Wikipedia, accessed January 16, 2025, [https://en.wikipedia.org/wiki/Attention\\_Is\\_All\\_You\\_Need](https://en.wikipedia.org/wiki/Attention_Is_All_You_Need)

29. Attention Is All You Need: The Core Idea of the Transformer | by Zain ul Abideen | Medium, accessed January 16, 2025, <https://medium.com/@zaiinn440/attention-is-all-you-need-the-core-idea-of-the-transformer-bbfa9a749937>

30. Understanding Google's "Attention Is All You Need" Paper and Its Groundbreaking Impact, accessed January 16, 2025, <https://alok-shankar.medium.com/understanding-googles-attention-is-all-you-need-paper-and-its-groundbreaking-impact-c5237043540a>

31. Attention is All you Need - NIPS papers, accessed January 16, 2025,

<https://papers.neurips.cc/paper/7181-attention-is-all-you-need.pdf>

32. openai-cookbook/examples/o1/Using\_reasoning\_for\_routine\_generation.ipynb at main, accessed January 16, 2025,

[https://github.com/openai/openai-cookbook/blob/main/examples/o1/Using\\_reasoning\\_for\\_routine\\_generation.ipynb](https://github.com/openai/openai-cookbook/blob/main/examples/o1/Using_reasoning_for_routine_generation.ipynb)

33. MiniMax-Text-01 - GitHub, accessed January 16, 2025,

<https://github.com/MiniMax-AI/MiniMax-01>

34. Paper Implementation — Attention Is All You Need(Transformer) | by Ujjalkumarmaity, accessed January 16, 2025,

<https://medium.com/@ujjalkumarmaity1998/paper-implementation-attention-is-all-you-need-transformer-59b95a93195c>

35. Attention Is All you Need [pytorch] - Kaggle, accessed January 16, 2025,

<https://www.kaggle.com/code/soupmonster/attention-is-all-you-need-pytorch>