Open in app ↗

# Medium

🔍 Search      ✍ Write    🔔    J✦

# A Simple Guide to DeepSeek R1: Architecture, Training, Local Deployment, and Hardware Requirements

Isaak Kamau · Follow

7 min read · Jan 23, 2025

👏 2.9K    💬 26         🔖 ▷ ⬆ •••

| Benchmark (Metric) | Claude-3.5-Sonnet-1022 | GPT-4o 0513 | DeepSeek V3 | OpenAI o1-mini | OpenAI o1-1217 | DeepSeek R1 |
|---|---|---|---|---|---|---|
| Architecture | - | - | MoE | - | - | MoE |
| # Activated Params | - | - | 37B | - | - | 37B |
| # Total Params | - | - | 671B | - | - | 671B |
| **English** MMLU (Pass@1) | 88.3 | 87.2 | 88.5 | 85.2 | **91.8** | 90.8 |
| MMLU-Redux (EM) | 88.9 | 88.0 | 89.1 | 86.7 | - | **92.9** |
| MMLU-Pro (EM) | 78.0 | 72.6 | 75.9 | 80.3 | - | **84.0** |
| DROP (3-shot F1) | 88.3 | 83.7 | 91.6 | 83.9 | 90.2 | **92.2** |
| IF-Eval (Prompt Strict) | **86.5** | 84.3 | 86.1 | 84.8 | - | 83.3 |
| GPQA Diamond (Pass@1) | 65.0 | 49.9 | 59.1 | 60.0 | **75.7** | 71.5 |
| SimpleQA (Correct) | 28.4 | 38.2 | 24.9 | 7.0 | **47.0** | 30.1 |
| FRAMES (Acc.) | 72.5 | 80.5 | 73.3 | 76.9 | - | **82.5** |
| AlpacaEval2.0 (LC-winrate) | 52.0 | 51.1 | 70.0 | 57.8 | - | **87.6** |
| ArenaHard (GPT-4-1106) | 85.2 | 80.4 | 85.5 | 92.0 | - | **92.3** |
| **Code** LiveCodeBench (Pass@1-COT) | 38.9 | 32.9 | 36.2 | 53.8 | 63.4 | **65.9** |
| Codeforces (Percentile) | 20.3 | 23.6 | 58.7 | 93.4 | **96.6** | 96.3 |
| Codeforces (Rating) | 717 | 759 | 1134 | 1820 | **2061** | 2029 |
| SWE Verified (Resolved) | **50.8** | 38.8 | 42.0 | 41.6 | 48.9 | 49.2 |
| Aider-Polyglot (Acc.) | 45.3 | 16.0 | 49.6 | 32.9 | **61.7** | 53.3 |
| **Math** AIME 2024 (Pass@1) | 16.0 | 9.3 | 39.2 | 63.6 | 79.2 | **79.8** |
| MATH-500 (Pass@1) | 78.3 | 74.6 | 90.2 | 90.0 | 96.4 | **97.3** |
| CNMO 2024 (Pass@1) | 13.1 | 10.8 | 43.2 | 67.6 | - | **78.8** |
| **Chinese** CLUEWSC (EM) | 85.4 | 87.9 | 90.9 | 89.9 | - | **92.8** |
| C-Eval (EM) | 76.7 | 76.0 | 86.5 | 68.9 | - | **91.8** |
| C-SimpleQA (Correct) | 55.4 | 58.7 | **68.0** | 40.3 | - | 63.7 |

*Poonam Soni*

Table 4 | Comparison between DeepSeek-R1 and other representative models.

## DeepSeek's Novel Approach to LLM Reasoning

DeepSeek has introduced an innovative approach to improving the reasoning capabilities of large language models (LLMs) through reinforcement learning (RL), detailed in their recent paper on DeepSeek-R1. This research represents a significant advancement in how we can enhance LLMs' ability to solve complex problems through pure reinforcement learning, without relying heavily on supervised fine-tuning.

*Before we proceed if you like this topic and you want to support me:*

1. ***Clap*** *my article 10 times; that will help me out.* 👏

2. ***Follow*** *me on Medium to get my latest articles* 🤝

**Technical Overview of DeepSeek-R1**

**Model Architecture:**

DeepSeek-R1 is not a singular model but a family of models, encompassing: **DeepSeek-R1-Zero** and **DeepSeek-R1**

*Let me clarify the key differences between DeepSeek-R1 and DeepSeek-R1-Zero:*

The Primary Distinction

**DeepSeek-R1-Zero** represents the team's initial experiment using pure reinforcement learning without any supervised fine-tuning. They started with their base model and applied reinforcement learning directly, letting the model develop reasoning capabilities through trial and error. While this approach achieved impressive results (71% accuracy on AIME 2024), it had some significant limitations, particularly in readability and language consistency. It features 671 billion parameters, utilizing a mixture-of-experts (MoE) architecture where each token activates parameters equivalent to 37 billion. This model showcases emergent reasoning behaviors, such as self-verification, reflection, and long chain-of-thought (CoT) reasoning.

**DeepSeek-R1,** in contrast, uses a more sophisticated multi-stage training approach. Instead of pure reinforcement learning, it begins with supervised

fine-tuning on a small set of carefully curated examples (called "cold-start data") before applying reinforcement learning. This approach addresses the limitations of DeepSeek-R1-Zero while achieving even better performance. This model also maintains the 671 billion parameter count but achieves better readability and coherence in responses.

## The Training Process Comparison

**Training Methodology:**

- **Reinforcement Learning**: Unlike traditional models that predominantly rely on supervised learning, DeepSeek-R1 uses RL extensively. The training leverages group relative policy optimization (GRPO), focusing on accuracy and format rewards to enhance reasoning capabilities without the need for extensive labeled data.

- **Distillation Techniques**: To democratize access to high-performing models, DeepSeek has also released distilled versions of R1, ranging from 1.5 billion to 70 billion parameters. These models are based on architectures like Qwen and Llama, showing that complex reasoning can be encapsulated in smaller, more efficient models. The distillation process involves fine-tuning these smaller models with synthetic reasoning data generated by the full DeepSeek-R1, thus preserving high performance at reduced computational cost.

**DeepSeek-R1-Zero's training process is straightforward:**

- Start with base model

- Apply reinforcement learning directly

- Use simple rewards based on accuracy and format

## DeepSeek-R1's training process has four distinct stages:

1. Initial supervised fine-tuning with thousands of high-quality examples

2. Reinforcement learning focused on reasoning tasks

3. Collection of new training data through rejection sampling

4. Final reinforcement learning across all types of tasks

## Performance Metrics:

- **Reasoning Benchmarks**: DeepSeek-R1 has shown impressive results on various benchmarks:

- **AIME 2024:** Achieved a 79.8% pass rate, compared to 79.2% by OpenAI's o1–1217.

- **MATH-500:** Scored an impressive 97.3%, slightly ahead of o1–1217's 96.4%.

- **SWE-bench Verified:** Outperformed in programming tasks, showcasing its coding proficiency.

- **Cost Efficiency:** The API for DeepSeek-R1 is priced at $0.14 per million input tokens for cache hits, making it significantly cheaper than comparable models like OpenAI's o1.

## Limitations and Future Work

The paper acknowledges several areas for improvement:

- The model sometimes struggles with tasks requiring specific output formats

- Performance on software engineering tasks could be enhanced

- There are challenges with language mixing in multilingual contexts

- Few-shot prompting consistently degrades performance

Future work will focus on addressing these limitations and expanding the model's capabilities in areas like function calling, multi-turn interactions, and complex role-playing scenarios.

## Deployment and Accessibility

- **Open Source and Licensing**: DeepSeek-R1 and its variants are released under the MIT License, promoting open-source collaboration and commercial use, including model distillation. This move is pivotal for fostering innovation and reducing the entry barriers in AI model development.

- **Model Formats:**

- Both models and their distilled versions are available in formats like GGML, GGUF, GPTQ, and HF, allowing flexibility in how they are deployed locally.

## 1. Web Access via DeepSeek Chat Platform:

The DeepSeek Chat platform provides a user-friendly interface to interact with DeepSeek-R1 without any setup requirements.

- **Steps to Access:**

- Navigate to the DeepSeek Chat platform

- Register for an account or log in if you already have one.

- After logging in, select the "Deep Think" mode to experience DeepSeek-R1's step-by-step reasoning capabilities.



DeepSeek Chat Platform

## 2. Access via DeepSeek API:

For programmatic access, DeepSeek offers an API compatible with OpenAI's format, allowing integration into various applications.

**Steps to Use the API:**

**a. Obtain an API Key:**

- Visit the <u>DeepSeek API platform</u> to create an account and generate your unique API key.

**b. Configure Your Environment:**

- Set the `base_url` to https://api.deepseek.com/v1.

- Use your API key for authentication, typically via Bearer Token in the HTTP header.

## c. Make API Calls:

- Utilize the API to send prompts and receive responses from DeepSeek-R1.

- Detailed documentation and examples are available in the DeepSeek API Docs.

```python
# Please install OpenAI SDK first: `pip3 install openai`

from openai import OpenAI

client = OpenAI(api_key="<DeepSeek API Key>", base_url="https://api.deepseek.com")

response = client.chat.completions.create(
    model="deepseek-chat",
    messages=[
        {"role": "system", "content": "You are a helpful assistant"},
        {"role": "user", "content": "Hello"},
    ],
    stream=False
)

print(response.choices[0].message.content)
```

DeepSeek API call example

## 3. Running DeepSeek-R1 Locally:

**Both Models (R1 and R1-Zero):**

- **Hardware Requirements:** The full models require significant hardware due to their size. A GPU with substantial VRAM (like Nvidia RTX 3090 or

higher) is recommended. For CPU use, you'd need at least 48GB of RAM and 250GB of disk space, although performance would be slow without GPU acceleration.

- **Distilled Models:** For local deployment with less resource-intensive hardware, DeepSeek provides distilled versions. These range from 1.5B to 70B parameters, making them suitable for systems with more modest hardware. For instance, the 7B model can run on a GPU with at least 6GB VRAM or on a CPU with about 4GB RAM for the GGML/GGUF format.

**Software Tools for Local Running:**

1. **Ollama:**

You can use <u>Ollama</u> to serve the models locally: (Ollama Is a tool for running open-source AI models locally on your machine. Grab it here: <u>https://ollama.com/download</u>)

# Get up and running with large language models.

Run Llama 3.3, Phi 4, Mistral, Gemma 2, and other models. Customize and create your own.

Download ↓

Available for macOS, Linux, and Windows

**Next, you'll need to pull and run the DeepSeek R1 model locally.**

Ollama offers different model sizes — basically, bigger models equal to smarter AI, but need better GPU. Here's the lineup:

```
1.5B version (smallest):
ollama run deepseek-r1:1.5b

8B version:
ollama run deepseek-r1:8b

14B version:
ollama run deepseek-r1:14b

32B version:
ollama run deepseek-r1:32b

70B version (biggest/smartest):
ollama run deepseek-r1:70b
```

To begin experimenting with DeepSeek-R1, it is advisable to start with a smaller model to familiarize yourself with the setup and ensure compatibility with your hardware. You can initiate this process by opening your terminal and executing the following command:

```
ollama run deepseek-r1:8b
```

Image courtesy from Reddit, via r/macapps

## Sending Requests to locally downloaded DeepSeek-R1 via Ollama:

Ollama provides an API endpoint to interact with DeepSeek-R1 programmatically. Ensure that the Ollama server is running locally before making API requests. You can start the server by running:

```
ollama serve
```

Once the server is active, you can send a request using `curl` as follows:

```
curl -X POST http://localhost:11434/api/generate -d '{
  "model": "deepseek-r1",
  "prompt": "Your question or prompt here"
}'
```

Replace `"Your question or prompt here"` with the actual input you wish to provide to the model. This command sends a POST request to the local Ollama server, which processes the prompt using the specified DeepSeek-R1 model and returns the generated response.

## Other methods to run/Access the models locally are:

**vLLM/SGLang:** Used for serving the models locally. Commands like vllm serve deepseek-ai/DeepSeek-R1-Distill-Qwen-32B — tensor-parallel-size 2 — max-model-len 32768 — enforce-eager can be used for the distilled versions.



Courtesy: HuggingFace

**llama.cpp:** You can also use llama.cpp to run the models locally.

## See What Others Are Building with DeepSeek-R1:

1. Running DeepSeek R1 across my 7 M4 Pro Mac Minis and 1 M4 Max MacBook Pro:

DeepSeek R1 1.5B running fully locally in your browser at 60 tok/ sec powered by WebGPU:

2/1/25, 7:43 PM

A Simple Guide to DeepSeek R1: Architecture, Training, Local Deployment, and Hardware Requirements | by Isaak Kamau | Jan, 2025 | Medium

RAG app to chat with your PDF files using the DeepSeek R1 model, running locally on your computer.

Running DeepSeek R1 version 1.5B perfectly locally on phone:

Cracking complex math problems with ease! (Thought for ~3200 tokens in about 35 seconds on M4 Max with mlx-lm.):

## Conclusions:

This progression from DeepSeek-R1-Zero to DeepSeek-R1 represents an important learning journey in the research. While DeepSeek-R1-Zero proved that pure reinforcement learning could work, DeepSeek-R1 showed how combining supervised learning with reinforcement learning could create an even more capable and practical model.

**Collaborations** 🤝 : Have an interesting AI project in mind? Let's team up! I'm available for collaboration on AI and machine learning initiatives, and keen to connect with other professionals in the field.

**Written by Isaak Kamau**

766 Followers  ·  25 Following

Follow

AI Developer: https://www.linkedin.com/in/isaak-mwangi/

# Responses (26)

What are your thoughts?

Respond

**Dr. Derek Austin** 🥳 he/him
5 days ago

Useful, thanks Isaak!

👏 38          💬 1 reply          Reply

**mou nandi**
4 days ago

Thank you for the quick write up.

👏 14          💬 1 reply          Reply

**Asrofi Al Kindi**
5 days ago

great article, thanks

👏 10          💬 1 reply          Reply

See all responses

# More from Isaak Kamau





👤 Isaak Kamau

👤 Isaak Kamau

## Microsoft Open-Sources 1-bit LLMs: Run 100B Parameter Model...

## Test-Time Compute Scaling: How to make an LLM "think longer" on...

https://github.com/microsoft/BitNet

Scaling Test-Time Compute: Unlocking Efficient Performance in Language Models

⭐ Oct 22, 2024    👏 970    💬 9              🔖⁺  •••

⭐ Dec 23, 2024    👏 113              🔖⁺  •••

👤 Isaak Kamau                      👤 Isaak Kamau

## When to Use XGBoost over LLM: A Case Study in Stock Market...

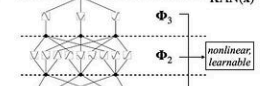Choosing between XGBoost (eXtreme Gradient Boosting) and a Large Language...

## A Simplified Explanation Of The New Kolmogorov-Arnold Network...

A screenshot from: https://arxiv.org/abs/2404.19756

Aug 4, 2024    ✋ 146          🔖 ⋯      ⭐ May 1, 2024    ✋ 3.3K    💬 8      🔖 ⋯

---

( See all from Isaak Kamau )

---

## Recommended from Medium

The Nam

In AI Advances by Kenny Vaneetvelde

## DeepSeek-R1: Architecture and training explain

Can we train a powerful reasoning LLM without Supervised Fine-Tuning?

## Want to Build AI Agents? Tired of LangChain, CrewAI, AutoGen &...

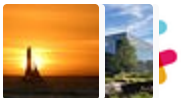Frameworks like LangChain, CrewAI, and AutoGen have gained popularity by promisin...

Jan 25   207   3

Jan 19   1K   22

## Lists

| | Staff picks |
|---|---|
| | 806 stories · 1603 saves |

| | Stories to Help You Level-Up at Work |
|---|---|
| | 19 stories · 930 saves |

| | Self-Improvement 101 |
|---|---|
| | 20 stories · 3258 saves |

| | Productivity 101 |
|---|---|
| | 20 stories · 2754 saves |

In Level Up Coding by Dr. Ashish Bamania

Abhishek Maheshwarappa

## DeepSeek-R1 Beats OpenAI's o1, Revealing All Its Training Secrets...

A deep dive into how DeepSeek-R1 was trained from scratch and how this open-...

## Fine-Tuning DeepSeek LLM: Adapting Open-Source AI for You...

Introduction

5d ago   1.1K   28

Jan 21   110   4

Mohit Vaswani

AI Papers Academy

## 6 AI Agents That Are So Good, They Feel Illegal

## DeepSeek-R1 Paper Explained — A New RL LLMs Era in AI?

AI agents are the future because they can replace all the manual work with automation...

Dive into the groundbreaking DeepSeek-R1 research paper, introduces open-source...

Jan 11    2.1K    81

Jan 25    228    3

See more recommendations