

# DeepSeek V4

## Technical Documentation

Publication date: April 27, 2026

Updated date: April 24, 2026

---

---

## DeepSeek V4 - Model Card

---

### General Information

<b>Model Provider</b>	DeepSeek AI
<b>Model name</b>	DeepSeek V4, including: <ul style="list-style-type: none"><li>• DeepSeek-V4-Pro</li><li>• DeepSeek-V4-Flash</li></ul>
<b>Release date</b>	April 24, 2026
<b>Model dependencies</b>	N/A

---

## Model properties

<b>Model architecture</b>	<p>DeepSeek V4 is a Mixture-of-Experts (MoE) language model series that retains the DeepSeekMoE framework and Multi-Token Prediction (MTP) strategy from DeepSeek-V3, while introducing several key architectural innovations:</p> <p>(1) Hybrid Attention Architecture: combines Compressed Sparse Attention (CSA) and Heavily Compressed Attention (HCA) to dramatically improve long-context efficiency. CSA compresses KV caches along the sequence dimension and applies DeepSeek Sparse Attention (DSA); HCA applies heavier compression with dense attention.</p> <p>(2) Manifold-Constrained Hyper-Connections (mHC): constrains residual mapping onto the manifold of doubly stochastic matrices (Birkhoff polytope), enhancing signal propagation stability while preserving model expressivity.</p> <p>(3) Muon Optimizer: employs the Muon optimizer for faster convergence and improved training stability.</p>
<b>Modalities</b>	<p>Text, three reasoning modes available:</p> <ul style="list-style-type: none"> <li>• Non-think: Fast, intuitive responses</li> <li>• Think High: Deliberate logical analysis</li> <li>• Think Max: Extended reasoning at maximum capacity</li> </ul>
<b>Context length</b>	<p>1M</p>
<b>Total model size</b>	<p>Pro: 1.6T parameters, of which 49B are activated for each token</p> <p>Flash: 285B parameters, of which 13B are activated for each token</p>

---

## Methods of distribution and licenses

<b>Distribution channels</b>	DeepSeek V4 is distributed primarily through two channels to accommodate different deployment needs: Open-source Repositories and application program interface(API).
<b>License</b>	MIT License: The assets distributed via open-source repositories (including model weights and code) are licensed under the MIT License as a free and open-source license.

---

## Use

	API deployment	Open-source deployment
<b>Acceptable Use Policy</b>	Access to the model via the API is governed by the <a href="#">DeepSeek Open Platform Terms of Service</a> .	This repository and the model weights are licensed under the <a href="#">MIT License</a> .
<b>Intended uses</b>	DeepSeek-V4, as a general-purpose AI model, aims to balance reasoning capabilities with output length, making it suitable for everyday use, such as question-answering scenarios and general agent tasks.	
<b>Technical means for model integration</b>	The technical means required for the model to be integrated into AI Systems or with models are described in the <a href="#">DeepSeek API Docs</a> .	Open-source deployment: The technical means required for the model to be integrated into AI Systems or with models are described in the model card available on the relevant open-source repository such as <a href="#">Hugging Face</a> .
<b>Required hardware</b>	N/A	See the model card available on Hugging Face.
<b>Required software</b>	N/A	See the model card available on Hugging Face.

## Model data

The capabilities of DeepSeek models are built on high-quality, large-scale, and diverse data sources. We place great emphasis on and strictly comply with laws and regulations related to intellectual property, trade secrets, and personal privacy, ensuring that all data acquisition and usage occur within a legal and compliant framework.

### 1. Pre-training Phase

During the pre-training phase, corpus data is required for training. This phase primarily uses the following two categories of data:

- **Public Data:** We use publicly available information on the internet to build the model's broad understanding of world knowledge. We employ technical methods to acquire and filter these freely accessible data to enrich the model's knowledge base.
- **Licensed Data:** We collaborate with third-party data providers to obtain proprietary datasets through legally signed agreements. We ensure all collaborations are based on lawful authorization.

The pre-training phase does not require personal information for training. Therefore, we do not intentionally collect personal information to associate with any specific account or individual, nor do we proactively use it to train our models. We exclude sensitive information, credit card numbers, or unique identification information from our training data sources to minimize the risk of collecting any personal information. However, due to the vast scale of pre-training data, some publicly available online content or licensed data from other providers may incidentally contain personal information. We employ technical measures to screen and remove such information from the training data as much as possible and conduct tests before using the data for training.

### 2. Optimization Training Phase

During the optimization training phase, we typically need to construct or annotate a set of question-answer pair data manually or automatically to train the model. These question-answer pairs are produced by our research team, with a small portion potentially based on

user input. If user input is used to construct training data, we apply secure encryption, strict de-identification, and anonymization to make it cannot be linked to any specific individual.

---