

# Inference Engine

## Introduction

An inference engine is a critical software component that serves as the decision-making brain of artificial intelligence systems. It processes facts and applies logical rules to deduce new information, make predictions, and solve problems.

Whether powering expert systems, chatbots, predictive analytics platforms, or autonomous vehicles, inference engines translate raw data and predefined rules into actionable insights. As AI adoption accelerates across industries, understanding inference engines becomes essential for organizations seeking to implement robust, decision-making systems.

## Basic Concept

An inference engine works with two main elements:

- A knowledge base that stores facts and rules.
- A reasoning mechanism that applies these rules to the facts to derive new information.

For example, if the system knows the facts “Patient has fever” and “Patient has cough,” and it has a rule “IF fever AND cough THEN possible flu,” the inference engine can conclude “possible flu” from the given facts.

## How it Works

Most inference engines follow a cycle:

- Match: Compare current facts to rules and find which rules can be applied.
- Execute: Apply the selected rule and add or update facts in working memory.

This cycle continues until no more rules can be applied or a goal condition is satisfied.

## Types of Inference

There are two common reasoning strategies:

- **Forward Chaining:**

It begins with known facts and applies inference rules to progressively generate new facts until a goal is reached. This approach is useful when exploring what conclusions can be derived from available data. For example, in medical diagnosis: if a patient has a fever AND a cough, the engine concludes the patient may have influenza.

- **Backward Chaining:**

Backward chaining starts with a hypothesis or goal and works backward through inference rules to find supporting facts needed to achieve it[1]. This method is effective when seeking specific answers. For example: to prove "patient has influenza," the engine works backward to determine what conditions (fever, cough, etc.) must be present.

## Components of an Inference Engine

Typical components include:

- **Knowledge base:** Stores domain facts and rules (for example, medical rules or home automation rules).
- **Working memory:** Holds current facts and intermediate conclusions during reasoning.
- **Inference mechanism** (rule interpreter): Matches facts with rules and decides which rule to fire next.
- **Control strategy:** Determines rule priority when multiple rules can be applied.

These components work together to allow the engine to reason systematically over the stored knowledge.

## Applications

Inference engines appear in many areas of AI and decision support:

- Expert systems (medical diagnosis, troubleshooting, technical support).
- Business rule engines (loan approval, fraud detection).
- Smart devices and IoT (smart homes, industrial monitoring).
- Modern machine learning serving systems, where the “inference engine” runs trained models to make predictions in real time.

In each case, the engine’s job is to turn available data plus rules or models into concrete decisions or recommendations.

# Modern LLM Inference Engines

While traditional inference engines work with rule-based systems, modern Large Language Model (LLM) inference engines are specialized frameworks designed to efficiently run neural network models for text generation and natural language processing. Three prominent examples are llama.cpp, vLLM, and SGLang, each optimized for different use cases and performance characteristics.

## LLama.cpp

llama.cpp is a lightweight, C/C++-based inference engine designed to run LLM models on CPUs and resource-constrained devices without requiring heavy dependencies like PyTorch or CUDA. It prioritizes efficiency and portability, enabling deployment on laptops, mobile devices, and even Raspberry Pi boards.

### Key Features

- Quantization support: Performs ahead-of-time model quantization and on-the-fly key-value cache quantization to reduce memory usage
- Efficient model format: Uses the GGUF file format for optimized model storage and loading
- Speculative decoding: Implements advanced optimization techniques that can double or triple inference speed without sacrificing output quality
- Flexible deployment: Supports partial offloading of model layers to system RAM, allowing devices to run models larger than available GPU VRAM
- API compatibility: Provides OpenAI-compatible endpoints and grammar-based output formatting

### Best Use Cases

llama.cpp excels in edge computing scenarios, embedded systems, and situations where running models locally without cloud infrastructure is essential. It's ideal for developers who need fine-grained control over inference settings and want to optimize for resource-constrained environments.

# vLLM

vLLM is a high-performance inference engine focused on maximizing throughput and efficiency for serving LLMs at scale, particularly on GPUs. It introduces innovative memory management techniques and has become a popular choice for production deployments.

## Key Features

- PagedAttention mechanism: Revolutionary memory management approach that delivers up to 24x higher throughput than traditional methods by efficiently managing key-value cache memory
- Separated architecture: Decouples the API server from the inference engine using separate processes connected via ZMQ, eliminating Python GIL contention
- Outstanding performance: Version 0.6.0 achieves 2.7x higher throughput and 5x faster time-per-output-token (TPOT) on Llama 8B models compared to earlier versions
- GPU optimization: Specifically tuned for modern GPUs like NVIDIA H100, where it achieves highest throughput when serving Llama models

## Best Use Cases

vLLM is best suited for production environments requiring high-throughput serving, real-time responses, and efficient memory utilization. It handles cloud deployments with multiple concurrent requests effectively, making it ideal for API services and enterprise applications.

# SGLang

SGLang (Structured Generation Language) is a domain-specific language embedded in Python that prioritizes low-latency inference and excels at handling complex, structured tasks.

## Key Features

- RadixAttention: Proprietary computation reuse algorithm that significantly accelerates inference speed through intelligent caching of intermediate results
- Compressed Finite State Machines: Advanced architecture designed for enhanced structured output generation, such as JSON formatting and form filling
- Complex task optimization: Achieves 3-5x faster performance than competitors on structured tasks requiring specific output formats
- Dynamic workload distribution: Efficiently manages varying workloads with sophisticated scheduling

## Performance Characteristics

SGLang handles approximately 120-160 requests per second with response times of 150-200 milliseconds. While raw throughput may be lower than vLLM in simple generation tasks, it significantly outperforms on complex structured generation workloads.

## Best Use Cases

SGLang is optimal for applications requiring structured responses, such as form filling, function calling, constrained generation, and scenarios where output must follow specific formats or schemas. It's particularly well-suited for cloud environments with complex task requirements.

# Comparison Summary

Feature	llama.cpp	vLLM	SGLang
<b>Primary focus</b>	CPU efficiency, edge deployment	GPU throughput, memory efficiency	Low latency, structured outputs
<b>Hardware target</b>	CPUs, resource-constrained devices	High-end GPUs (H100, A100)	Cloud GPUs
<b>Throughput</b>	Moderate	Highest (2.7x improvements)	120-160 req/sec
<b>Latency</b>	Variable	Very low TPOT	150-200ms
<b>Best for</b>	Local/edge inference, embedded systems	High-scale API serving, production	Complex structured tasks

## Model Used

The Qwen2-VL 2B-Instruct model is used because it offers a strong balance between multimodal capability and computational efficiency. This model is designed to handle both text and image inputs, which makes it suitable for applications that require understanding or reasoning over visual content as well as natural language instructions. Its relatively small size (2B parameters) means it can run more comfortably on limited hardware resources, especially when paired with an efficient inference engine like llama.cpp, while still providing high-quality, instruction-following behavior. This combination allows the system to deliver useful, context-aware responses in multimodal scenarios without demanding large-scale GPU infrastructure.

## Inference Engine

llama.cpp as an inference engine is selected due to its lightweight architecture and exceptional efficiency in running large language models on diverse hardware configurations. Built in C/C++, llama.cpp eliminates heavy dependencies like PyTorch or CUDA, making it ideal for deployments where resource optimization is critical. The engine's support for quantization through the GGUF format significantly reduces memory requirements, allowing the Qwen2-VL 2B model to run smoothly even on

systems with limited GPU memory. Additionally, llama.cpp's ability to offload model layers between GPU VRAM and system RAM provides flexibility in managing computational resources. Its speculative decoding feature can double or triple inference speed without sacrificing output quality, ensuring responsive performance for real-time applications. The OpenAI-compatible API endpoints further simplify integration into existing workflows, while enabling complete local inference maintains data privacy and eliminates cloud dependency costs.

## Model Performance

### English Printed Text Performance

The model demonstrated exceptional performance on English printed documents, achieving a perfect accuracy of 100% in 16.33 seconds. The extracted text precisely matched the expected output: "Today was a wonderful day. I went for a long walk in the park," showing the model's strong capability in recognizing clean, machine-printed English characters. This high accuracy validates the effectiveness of the Qwen2-VL 2B model when processing standard printed text in well-structured formats.

```
=====
🚀 STARTING MODEL EVALUATION
=====

[1/4] Processing...

📄 Testing: English Printed - test_images/english_printed.png
Expected: Today was a wonderful day. I went for a long walk in the par
Extracted: Today was a wonderful day. I went for a long walk in the par
Accuracy: 100.0%
Time: 16.33s
Status: ✅ Pass
```

### English Handwritten Text Performance

For English handwritten documents, the model achieved 87.92% accuracy with a processing time of 19.33 seconds. The expected text was "Once there was a dog. He was very hungry. He found a piece of," and the model successfully extracted most of this content, though with minor truncation at the end. This strong performance on handwritten English demonstrates the model's robustness in handling informal writing styles and variable character shapes, which are typically more challenging than printed text.

[2/4] Processing...

```
📄 Testing: English Handwritten - test_images/english_handwritten.png
Expected: Once there was a dog.He was very hungry.He found a piece of
Extracted: Once there was a dog.
He was very hungry.
He found a piece o
Accuracy: 87.92%
Time: 19.33s
Status: ✅ Pass
```

## Nepali Printed Text Performance

The model showed moderate performance on Nepali printed text, achieving 59.19% accuracy in 9.02 seconds. While the processing was relatively fast, the accuracy fell below acceptable thresholds, indicating challenges in recognizing Devanagari script characters in printed format. The extracted text showed partial recognition but contained significant deviations from the expected Nepali content, suggesting that the model requires additional fine-tuning or optimization for non-Latin scripts.

[3/4] Processing...

```
📄 Testing: Nepali Printed - test_images/nepali_printed.png
Expected: अृथ हुन खाल की सम्बन्ध हुन खाली स्वतन्त्र हुन खाली या प
Extracted: अृथ हुत खाल की सामन्य हुत खाल
स्वतन्त्र हुत खाल या दरव
Accuracy: 59.19%
Time: 9.02s
Status: ❌ Fail
```

## Nepali Handwritten Text Performance

Nepali handwritten documents presented the most significant challenge, with the model achieving only 28.92% accuracy despite taking 52.57 seconds to process. This represents the lowest performance across all test categories, with the model struggling to accurately recognize handwritten Devanagari characters. The extended processing time combined with low accuracy indicates that handwritten non-Latin scripts require substantial improvements, possibly through specialized training data or alternative model architectures better suited for complex script recognition.

[4/4] Processing...

```
📄 Testing: Nepali Handwritten - test_images/nepali_handwritten.png
Expected: होस्पियरको चेतावनी स्कैलाइ दिए नानीधुलाई घरमै ढोडे दिन
Extracted: होस्पियरको चेतावनी स्कैलाइ दिए नानीधुलाई घरमें ढोडे दिनभ
Accuracy: 28.92%
Time: 52.57s
Status: ✗ Fail
```

## Overall Performance Summary

Across all four test cases, the model achieved an average accuracy of 69.01% with an average processing time of 24.31 seconds. The results clearly show a performance gradient: the model excels with English printed text (100%), performs well on English handwritten text (87.92%), faces challenges with Nepali printed text (59.19%), and struggles significantly with Nepali handwritten text (28.92%). This pattern indicates that while the Qwen2-VL 2B model is highly capable for English-language document processing, it requires additional optimization or model selection adjustments for reliable multilingual support, particularly for languages using non-Latin scripts like Devanagari.

```
=====
📊 SUMMARY STATISTICS
=====
Total Tests: 4
Average Accuracy: 69.01%
Average Time: 24.31 seconds

=====
📋 DETAILED RESULTS
=====


| Document Type       | Accuracy | Time   | Status |
|---------------------|----------|--------|--------|
| English Printed     | 100.00%  | 16.33s | ✓ Pass |
| English Handwritten | 87.92%   | 19.33s | ✓ Pass |
| Nepali Printed      | 59.19%   | 9.02s  | ✗ Fail |
| Nepali Handwritten  | 28.92%   | 52.57s | ✗ Fail |


=====
```

