

# LLM Evaluation Matrix

Use this matrix to compare multiple LLMs against your requirements. Assign a weight to each pillar, then score each model from 1 (poor) to 5 (excellent).

## Step 1: Define Weights per Pillar

Adjust these weights based on what matters most for this project.

- Performance & quality (P): \_\_\_\_\_ / 100
- Latency & scalability (L): \_\_\_\_\_ / 100
- Cost & total cost of ownership (C): \_\_\_\_\_ / 100
- Security & compliance (S): \_\_\_\_\_ / 100
- Vendor support & ecosystem (V): \_\_\_\_\_ / 100

## Step 2: Scoring Scale

Use the same scale for all criteria:

- 1 = Poor / does not meet requirements
- 3 = Acceptable / meets minimum requirements
- 5 = Excellent / strongly exceeds requirements

## Step 3: List Candidate Models

- Model A: \_\_\_\_\_
- Model B: \_\_\_\_\_
- Model C: \_\_\_\_\_

## Step 4: Score Each Model by Pillar

For each criterion below, score Model A, Model B, and Model C from 1–5.

### 1. Performance & Quality (P)

How well does each model handle your tasks and domain?

Criterion	Weight (P)	Model A (1–5)	Model B (1–5)	Model C (1–5)
P1. Accuracy on your test set	_____	_____	_____	_____
P2. Reasoning on complex cases	_____	_____	_____	_____
P3. Handling domain-specific jargon	_____	_____	_____	_____

P4. Reliability / consistency of outputs	_____	_____	_____	_____
--	-------	-------	-------	-------

## 2. Latency & Scalability (L)

Can the model meet your speed and scale requirements?

Criterion	Weight (L)	Model A (1–5)	Model B (1–5)	Model C (1–5)
L1. Average latency at expected load	_____	_____	_____	_____
L2. Behaviour at peak traffic	_____	_____	_____	_____
L3. Rate limits vs. your needs (TPS / QPS)	_____	_____	_____	_____
L4. Uptime / availability guarantees	_____	_____	_____	_____

## 3. Cost & Total Cost of Ownership (C)

Look beyond per-token price to total cost over time.

Criterion	Weight (C)	Model A (1–5)	Model B (1–5)	Model C (1–5)
C1. Input token price	_____	_____	_____	_____
C2. Output token price	_____	_____	_____	_____
C3. Estimated monthly run-rate at your usage	_____	_____	_____	_____
C4. Extra infra / ops cost (if self-hosted)	_____	_____	_____	_____
C5. Integration & engineering effort	_____	_____	_____	_____

## 4. Security & Compliance (S)

Ensure the model fits your risk, data, and regulatory needs.

Criterion	Weight (S)	Model A (1–5)	Model B (1–5)	Model C (1–5)

S1. Data usage policy (training / logs)	_____	_____	_____	_____
S2. Data residency options match needs	_____	_____	_____	_____
S3. Private networking / VPC / on-prem options	_____	_____	_____	_____
S4. Logging & audit capabilities	_____	_____	_____	_____
S5. Fit with your specific regulations	_____	_____	_____	_____

**5. Vendor Support & Ecosystem (V)**

Consider maturity, support quality, and ecosystem strength.

Criterion	Weight (V)	Model A (1–5)	Model B (1–5)	Model C (1–5)
V1. Provider maturity & stability	_____	_____	_____	_____
V2. Documentation & SDK quality	_____	_____	_____	_____
V3. Support responsiveness / SLAs	_____	_____	_____	_____
V4. Roadmap alignment (12–24 months)	_____	_____	_____	_____
V5. Community, ecosystem & integrations	_____	_____	_____	_____

**Step 5: Summary & Decision**

After scoring, you can compute a rough weighted score per model (manually or in a spreadsheet) and capture your decision notes here.

- **Model A – total score & notes:** \_\_\_\_\_
- **Model B – total score & notes:** \_\_\_\_\_
- **Model C – total score & notes:** \_\_\_\_\_

**Final decision and rationale:** \_\_\_\_\_