

REFERENCE DOCUMENT

Estimating the Inference Carbon Intensity of Mistral AI Models: Methodology, Estimates, and Justification

Prepared by: **InferenceCarbon.ai Reference Team**

Date: 13 March 2026 | Version: 1.4 | Classification: For publication — subject to stated caveats

Important note: *Mistral AI is the only major AI provider to have published a peer-reviewed lifecycle assessment (LCA) of one of its models.¹ This gives Mistral the strongest anchor in the InferenceCarbon series: a provider-published, third-party-audited figure of 1.14 gCO_{2e} per 400-token response for Mistral Large 2 (123B dense).² This figure is cross-validated by Jegham et al. (2025), whose independent estimate aligned within one standard deviation.³ However, only the Large 2 anchor is directly provider-disclosed and independently cross-checked. All estimates for current models (Large 3, Medium 3, Small 3, Magistral) are scaled scenario estimates derived from that anchor via throughput heuristics and are materially less certain. The Large 2-to-Large 3 transition is a dense-to-MoE architecture change that particularly weakens the scaling proxy. Confidence is MEDIUM for the Large 2 anchor, LOW to MEDIUM-LOW for current-model estimates. We invite Mistral to extend its LCA to current models.*

Scope notice: *This document covers operational inference emissions only. It excludes training, embodied hardware emissions, water consumption, and end-of-life disposal. Estimates represent grams of CO₂ equivalent (gCO_{2e}) per 1,000 tokens of model output, using location-based Scope 2 accounting as the primary metric.⁴ Mistral's own LCA includes upstream emissions (hardware manufacturing, transport) which our series excludes; our estimates are therefore not directly comparable to Mistral's published 1.14 gCO_{2e} figure, which covers a wider scope.*

1. Executive Summary

Mistral AI presents the strongest disclosure case in the InferenceCarbon reference series. The company published the first peer-reviewed LCA of any major LLM in July 2025,⁵ conducted by Carbone 4 with ADEME support and audited by Resilio and Hubblo.⁶ The LCA reported 1.14 gCO_{2e} per 400-token Le Chat response for Mistral Large 2. (NB: this includes operational energy, plus amortised hardware manufacturing, transport, and upstream supply chain, not just the inference energy that we calculate in this paper.)⁷ Jegham et al. independently measured ~1.09 gCO_{2e} — aligning within one standard deviation.⁸ This cross-validation is the strongest in the series.

A critical advantage is France's electricity grid: 21.7 gCO_{2e}/kWh in 2024⁹ — roughly 17× lower than the US average (~370)¹⁰ and 27× lower than China (~580). Mistral is building sovereign compute in Essonne,

¹Mistral AI (July 2025). "Our Contribution to a Global Environmental Standard for AI." Peer-reviewed LCA with Carbone 4 and ADEME. <https://mistral.ai/news/our-contribution-to-a-global-environmental-standard-for-ai>

²Mistral LCA: 1.14 gCO_{2e} per 400-token Le Chat response, 45 mL water, 0.16 mg Sb eq. Includes upstream emissions. Mistral AI (2025), *ibid*.

³Jegham et al. Section 5.3: Mistral Large 2 (~1.09 gCO_{2e} for 400 tokens) aligned within 1 SD of Mistral LCA (~1.14 gCO_{2e}). <https://arxiv.org/abs/2505.09598>

⁴GHG Protocol (2025). Scope 2 Guidance.

<https://ghgprotocol.org/blog/upcoming-scope-2-public-consultation-overview-revisions>

⁶Carbone 4 (2025). "A New Milestone for Environmental Transparency in Generative AI."

<https://www.carbone4.com/en/ia-generative-mission-mistral-ai>

⁹RTE France (2025). "French Annual Electricity Review 2024." CIF 21.7 gCO_{2e}/kWh.

<https://analyseetdonnees.rte-france.com/en/annual-review-2024/keyfindings>

¹⁰US EPA eGRID (2024). US avg grid ~370 gCO_{2e}/kWh. <https://www.epa.gov/egrid>

France (18,000 Grace Blackwell GPUs, 40MW)¹¹² and Sweden (\$1.2B EcoDataCenter, operations expected from 2027).¹³ However, current inference runs partly on Azure, AWS, and other providers.¹⁴

External adversarial review characterised this paper as “the strongest paper in the series because it starts from a real disclosed anchor, but it still lets that unusually good anchor confer more certainty on the current Mistral lineup than the evidence really warrants.” We accept this characterisation.

Model	Location-Based Range(gCO ₂ e/1k tokens)	Central Scenario	CEAF %	Adjusted Central	Confidence
Mistral Large 2	0.60 – 1.20	0.85	0.50	0.43	Medium (disclosed)
Mistral Large 3	0.40 – 0.50	0.45	0.50	0.23	Medium-Low
Mistral Medium 3	0.25 – 0.35	0.30	0.50	0.15	Low
Mistral Small 3	0.04 – 0.07	0.05	0.50	0.03	Low

Table 1: Summary estimates for current hosted inference. Location-based at Azure CIF 0.35 kgCO₂e/kWh. CEAF = 0.50 for Microsoft Azure, based on verified 100% annual renewable matching (2025), discounted for lack of published hourly CFE data. All current-model estimates are scaled scenarios derived from the Large 2 anchor via throughput heuristics. The Large 3 estimate crosses a dense-to-MoE architecture boundary. Magistral Medium (reasoning) excluded — see Section 7.

Prospective Sovereign Compute — France/Sweden Infrastructure (Not Yet Operational)

Large 3 (France, CIF 0.022): 0.02–0.04 gCO₂e/1k tokens (central: 0.03). Small 3 (France, CIF 0.022): 0.003–0.005 (central: 0.003). CEAF not applied — French grid (~22 gCO₂e/kWh, 95% low-carbon) is already so clean that procurement adjustments are immaterial. Mistral Compute under construction; EcoDataCenter Sweden expected 2027. Not current product footprint.

Magistral Medium (Reasoning) — Illustrative Only

No direct measurement. Illustrative range (Azure, location-based): 0.50–2.10 gCO₂e/1k visible tokens. CEAF-adjusted (50%): 0.25–1.05. Confidence: Very Low.

2. What Has Been Published

Published LCA: In July 2025, Mistral published the first comprehensive LCA of a major LLM, covering Mistral Large 2.¹⁵ Conducted by Carbone 4, supported by ADEME, peer-reviewed by Resilio and Hubblo.¹⁶ Compliant with ISO 14040/44, GHG Protocol Product Standard, AFNOR Frugal AI methodology.

Published per-query LCA figure: 1.14 gCO₂e per 400-token Le Chat response, 45 mL water, 0.16 mg Sb eq.¹⁷ The only provider-published, audited per-inference carbon figure in the industry.

Published training figure: 20.4 ktCO₂e, 281,000 m³ water, 660 kg Sb eq over 18 months.¹⁸

¹¹DCD (2025). "France AI data center build-out." Mistral Compute: 18,000 GB200 GPUs, 40MW Essonne. <https://www.datacenterdynamics.com/en/analysis/france-ai-data-center-build-out-emmanuel-macron/>

¹²NVIDIA Blog (2025). "France Bolsters National AI Strategy." Mistral Compute with Grace Blackwell. <https://blogs.nvidia.com/blog/france-sovereign-ai-infrastructure/>

¹³Reuters (Feb 2026). "France AI company Mistral invests \$1.4 billion in data centres in Sweden." Operations expected from 2027.

<https://www.reuters.com/sustainability/boards-policy-regulation/france-ai-company-mistral-invests-14-billion-data-centres-sweden-2026-02-11/>

¹⁴Mistral AI (Dec 2025). "Introducing Mistral 3." Available on Azure, AWS, Scaleway. <https://mistral.ai/news/mistral-3>

¹⁸Mistral LCA: Training Large 2: 20.4 ktCO₂e, 281,000 m³ water, 660 kg Sb eq. Mistral AI (2025), *ibid*.

Infrastructure transparency: Mistral has disclosed transition from Azure and Google Cloud¹⁹ to Mistral Compute (18,000 NVIDIA Grace Blackwell GPUs, 40MW, Essonne).^{20,21} Also: \$1.2B EcoDataCenter in Sweden (operations expected 2027).^{22,23}

Open model weights: Large 3 is Apache 2.0.²⁴

No per-query data for current models: The LCA covers Large 2 only. No updated figures for Large 3, Medium 3, Small 3, or Magistral.

3. Anchor Derivation

3.1 Dual Anchors: Provider LCA and Jegham et al.

Anchor A (Mistral LCA): 1.14 gCO_{2e} per 400-token response.²⁵ Includes upstream emissions not in our operational-only scope.²⁶

Anchor B (Jegham): ~1.09 gCO_{2e} per 400-token query on Azure (H200/H100, PUE(Power Usage Effectiveness) 1.12, CIF 0.35).^{27,28} Operational-only, consistent with our series methodology.

Cross-validation: 1.14 vs ~1.09 gCO_{2e} — within one standard deviation.²⁹ The strongest anchor cross-validation in the entire InferenceCarbon series.

Adversarial review caveat: This anchor quality applies to Large 2 only. It does not transfer to current models. The prestige of the anchor should not spill over onto models that do not deserve the same confidence.

3.2 Deriving gCO_{2e} per 1,000 Tokens

Anchor A (Mistral LCA, 1.14 gCO_{2e} per 400 tokens) is the authoritative headline figure, but it includes upstream emissions outside our operational-only scope. For operational scaling we therefore use Jegham Anchor B: ~1.09 gCO_{2e} for 400 tokens on Azure (operational only, CIF 0.35). Per 1,000 tokens: ~2.73 gCO_{2e}/1k tokens. This is the basis for all downstream estimates.

4. Scaling to Current Models

4.1 Throughput Data

Model	Parameters	Output TPS	Architecture
Large 2 July '24 (anchor)	123B dense	31.4	Dense
Large 3	675B (41B active)	49.6	MoE
Medium 3	Proprietary	54.3	Proprietary
Small 3	24B dense	161.0	Dense

¹⁹Introl (2025). "France AI Sovereignty Push." Mistral previously on Azure/Google Cloud.

<https://introl.com/blog/france-ai-sovereignty-mistral-sovereign-cloud-2025>

²²BeBeez (Feb 2026). "Mistral signs \$1.2bn Sweden EcoDataCenter deal."

<https://bebeez.eu/2026/02/11/french-ai-firm-mistral-signs-1-2bn-deal-to-build-ecodatacenter-facility-in-sweden/>

²⁷Jegham, N. et al. (2025). "How Hungry is AI?" arXiv:2505.09598v6. <https://arxiv.org/abs/2505.09598>

²⁸Jegham et al. Table 1: Mistral on Azure. DGX H200/H100, PUE 1.12, CIF 0.35. <https://arxiv.org/abs/2505.09598>

Sources: Large 3³⁰; Medium 3³¹; Small 3³²; overview.³³

4.2 The Dense-to-MoE Fault Line

Adversarial review caveat: The Large 2 (123B dense) to Large 3 (675B MoE, 41B active) transition is the paper’s central technical fault line. A 675B model with 41B active parameters is not just a bigger Large 2; it is a qualitatively different serving regime. MoE changes the relationship between parameter count, active compute per token, memory movement, batching behaviour, and realised energy per visible token. Even though Large 3 may be faster per visible token, it does not follow that its joules per token can be inferred cleanly by inverse TPS (Tokens per Second) with a GEC (Generational Efficiency Correction, see below) overlay. The Large 3 estimate should be read as a dense-to-MoE scenario estimate, not a scaled continuation of Large 2.

4.3 Generational Efficiency Correction (GEC)

GEC ranges: Large 3 vs Large 2: 0.70–0.90 (midpoint 0.80). Medium 3: 0.50–0.70 (midpoint 0.60). Small 3: 0.25–0.40 (midpoint 0.30). All GEC values are under-justified placeholders. The entire estimation chain beyond the anchor — throughput rescaling, GEC correction, architecture proxy — is where most uncertainty lives.

5. Infrastructure and Carbon Intensity

5.1 Current Hosted Inference

Mistral has historically relied on Microsoft Azure and Google Cloud for training and inference.³⁴ Current models are available via Mistral’s own API, Azure, and AWS Bedrock.³⁵

Provider	CIF (kgCO ₂ e/kWh)	CEAF Tier	CEAF
Microsoft Azure	0.35	Tier 1.5	0.50
AWS Bedrock	0.287	Tier 2	0.50
Google Cloud	~0.25	Tier 1	0.66

5.2 Prospective Sovereign Infrastructure (Not Yet Operational)

Mistral is building sovereign European infrastructure, but it is not yet serving production inference:

Facility	Status	CIF (kgCO ₂ e/kWh)	Notes
Mistral Compute (Essonne, France)	Under construction	~0.022	40MW, 18,000 GB200 GPUs
EcoDataCenter (Borlänge, Sweden)	Announced Feb 2026; ops ~2027	~0.015	\$1.2B, hydro/nuclear grid

³⁰Artificial Analysis: Mistral Large 3. 675B/41B active MoE. 49.6 TPS. Intelligence 23.

<https://artificialanalysis.ai/models/mistral-large-3>

³¹Artificial Analysis: Mistral Medium 3. 54.3 TPS. Intelligence 19. <https://artificialanalysis.ai/models/mistral-medium-3>

³²Artificial Analysis: Mistral Small 3. 24B dense. 161 TPS. Intelligence 13. <https://artificialanalysis.ai/models/mistral-small-3>

³³Artificial Analysis: Mistral provider overview. 21 models tracked. <https://artificialanalysis.ai/providers/mistral>

Sources: France CIF^{36,37}; Mistral Compute^{38,39,40}; Sweden^{41,42}

Adversarial review caveat: Mistral’s carbon advantage is not fully “Mistral” yet in a deployment sense. The France and Sweden facilities are infrastructure announcements and transition plans, not yet a fully evidenced current routing profile for production inference. The paper must not conflate current hosted inference, future sovereign compute, and provider identity. Disclosure leadership and future low-carbon infrastructure deserve credit, but customer-facing carbon values must reflect current operational evidence, not strategic direction. (That said, we strongly support Mistral’s direction on this, and we will update our estimates as soon as it is correct to do so.)

The France advantage: France’s nuclear-dominated grid (~70% nuclear, 95% low-carbon)⁴³ produces electricity at ~22 gCO₂e/kWh⁴⁴ — 17× lower than US average.⁴⁵ Once Mistral Compute is operational, inference on French infrastructure could produce carbon intensity an order of magnitude lower than any US-hosted provider. This is a directionally plausible order-of-magnitude statement, but it assumes energy per query stays comparable across infrastructures. A Blackwell-based French stack could change server efficiency, cooling overhead, scheduler utilisation, and model-variant performance in ways we cannot quantify from CIF conversion alone.

6. Per-Model Estimates (Current Hosted Inference)

All estimates in this section use Azure CIF 0.35, reflecting current inference routing. These are scaled scenario estimates derived from the Large 2 anchor, not provider-disclosed figures. For prospective France-hosted estimates, see the Executive Summary table.

6.1 Mistral Large 3 (Dense-to-MoE Scenario Estimate)

675B/41B active MoE. Intelligence 23. 49.6 TPS.^{46,47}

Formula:

$$\text{Anchor} \times (\text{Anchor TPS} \div \text{Model TPS}) \times \text{GEC} = 0.85 \times (31.4 \div 49.6) \times \text{GEC} = 0.85 \times 0.63 \times \text{GEC}.$$

High (GEC 0.90): 0.48 → 0.50.

Central (GEC 0.80): $0.85 \times 0.63 \times 0.80 = 0.43 \rightarrow \mathbf{0.45}$.

Low (GEC 0.70): 0.38 → 0.40.

Note: the dense-to-MoE transition means TPS may understate per-token energy — MoE introduces memory movement, expert routing, and batch scheduling overheads not captured by TPS scaling.

Range (Azure): 0.40 – 0.50 gCO₂e/1k tokens (central: 0.45).

Confidence: Medium-Low — dense-to-MoE change weakens proxy.

Status: Scaled scenario estimate from Large 2 anchor.

6.2 Mistral Medium 3

Proprietary architecture. Intelligence 19. 54.3 TPS.⁴⁸

³⁷Enerdata (Jan 2025). France 2024: record low 21.3 gCO₂e/kWh, 95% low-carbon.

<https://www.enerdata.net/publications/daily-energy-news/nuclear-and-renewables-raised-frances-2024-power-generation-5-year-high.html>

⁴⁰NVIDIA Newsroom (June 2025). "Europe Builds AI Infrastructure With NVIDIA." Mistral French compute platform announced. <https://nvidianews.nvidia.com/news/europe-ai-infrastructure>

⁴³Ember (2024). France: 5.2% fossil in 2025. <https://ember-energy.org/countries-and-regions/france/>

Formula:

$$\text{Anchor} \times (\text{Anchor TPS} \div \text{Model TPS}) \times \text{GEC} = 0.85 \times (31.4 \div 54.3) \times \text{GEC} = 0.85 \times 0.58 \times \text{GEC}.$$

High (GEC 0.70): $0.34 \rightarrow 0.35$.

Central (GEC 0.60): $0.85 \times 0.58 \times 0.60 = \mathbf{0.30}$.

Low (GEC 0.50): 0.25.

Range (Azure): 0.25 – 0.35 gCO₂e/1k tokens (central: 0.30).

Confidence: Low. Status: Scaled scenario estimate.

6.3 Mistral Small 3

24B dense. Intelligence 13. 161 TPS.⁴⁹

Formula:

$$\text{Anchor} \times (\text{Anchor TPS} \div \text{Model TPS}) \times \text{GEC} = 0.85 \times (31.4 \div 161.0) \times \text{GEC} = 0.85 \times 0.20 \times \text{GEC}.$$

High (GEC 0.40): 0.07.

Central (GEC 0.30): $0.85 \times 0.20 \times 0.30 = \mathbf{0.05}$.

Low (GEC 0.25): 0.04.

Range (Azure): 0.04 – 0.07 gCO₂e/1k tokens (central: 0.05).

Confidence: Low. Status: Scaled scenario estimate.

7. Reasoning Models — Illustrative Only

Following adversarial review, Magistral Medium is demoted from the main tables. No direct measurement exists. The thinking-token overhead is unobserved and task-dependent. Once hidden reasoning tokens are unobserved, “per 1k tokens” becomes slippery unless the paper specifies whether that means visible output only or some estimate of total internal compute-related tokens.

Illustrative range (Azure): 0.50 – 2.10 gCO₂e/1k visible tokens. No central estimate is provided. This figure should not be used for comparison, display, or planning without explicit qualification.

8. Three-Layer Disclosure Structure

Layer	Metric	Large 2 (Anchor)	Large 3 (Azure)	Small 3 (Azure)
Layer 1: Energy	Wh/short query	~1.10	~0.55	~0.07
Layer 2: Location-based	gCO ₂ e/1k range	0.60–1.20	0.40–0.50	0.04–0.07
Layer 3: CEAF-adjusted	gCO ₂ e/1k	0.30–0.60	0.20–0.25	0.02–0.04

9. Clean Energy Adjustment Factor (CEAF)

For current hosted inference on Azure/AWS, standard CEAF values apply (Tier 1.5/2, CEAF 0.50). For prospective Mistral Compute in France, the grid is already so low-carbon (~22 gCO₂e/kWh) that CEAF adjustment is largely immaterial. The cleanest sustainability communication for French-hosted inference

is: physical location-based figure first, no net-like embellishment, and a clear note that procurement attributes are secondary because the grid is already unusually low-carbon.

10. Cross-Validation

Against Mistral LCA: Jegham ~ 1.09 vs Mistral 1.14 gCO₂e for 400 tokens — within 1 SD.⁵⁰ The only case where a provider's own disclosure validates an independent measurement.

Against Jegham/Altman: Jegham GPT-4o (0.423 Wh) aligns within 19% of Altman 0.34 Wh.⁵¹

Against Google: Google measured 0.24 Wh median Gemini prompt⁵² — lower-bound reference for efficient inference on custom hardware.

11. Confidence Assessment

Estimate	Confidence	Dominant uncertainty
Anchor (Large 2, LCA)	Medium-High	Provider-published, audited; scope wider than series
Anchor (Large 2, Jegham)	Medium	Single source; Azure infrastructure
Large 3 (MoE scenario)	Medium-Low	Dense-to-MoE; throughput proxy weak
Medium 3	Low	Proprietary architecture
Small 3	Low	Two-step scaling
Magistral (reasoning)	Very Low	No measurement; illustrative only
France CIF	High	RTE official; stable nuclear base
Current routing split	Unknown	Azure/AWS/Mistral Compute shares undisclosed

12. Limitations

1. Scope difference. Mistral LCA includes upstream; our series uses operational-only. Not directly comparable.
2. LCA covers Large 2 only. No provider data for Large 3, Medium 3, Small 3, or Magistral.
3. Dense-to-MoE architecture change. Large 2 to Large 3 is a qualitatively different serving regime that weakens throughput scaling.
4. Infrastructure in transition. Routing split between Azure/AWS/Mistral Compute unknown. Sovereign compute not yet operational.
5. Sweden facility not yet built. EcoDataCenter operations expected from 2027.⁵³
6. Reasoning model unquantified. No Magistral thinking-token measurement.
7. France CIF advantage not yet fully realised. Mistral Compute under construction.
8. Operational inference only. Training (20.4 ktCO₂e), water, Scope 3 excluded.
9. Third-party hosting. Mistral models widely self-hosted (Apache 2.0).

⁵¹Altman, S. (June 2025). "The Gentle Singularity." 0.34 Wh avg ChatGPT query.
<https://blog.samaltman.com/the-gentle-singularity>

⁵²Google (Aug 2025). "Measuring Environmental Impact of AI Inference." 0.24 Wh. arXiv:2508.15734.
<https://arxiv.org/abs/2508.15734>

10. Disclosure leadership is not current product footprint. Mistral’s transparency advantage does not automatically mean the current operational footprint of every Mistral-served token is already the best in class.

13. Response to Adversarial Review

An external adversarial review was conducted by GPT-5.4 Thinking (OpenAI) on 13 March 2026, using dual personas: a PhD technologist specialising in data-centre technologies, and a senior sustainability professional. The review was conducted on v1.0 of this paper.

Summary verdict: “The strongest paper in the series because it starts from a real disclosed anchor, but it still lets that unusually good anchor confer more certainty on the current Mistral lineup than the evidence really warrants.” The technologist called it “the best-anchored paper in the series, but still not a measurement-grade framework for current Mistral models.” The sustainability reviewer called it “the most credible and publication-ready paper in the set, but still not fully safe for crisp model-by-model external claims beyond the anchor model.”

We accept these characterisations.

Three changes insisted upon and accepted:

- 1. Separate anchor quality from estimate quality.** Labels restructured: only Large 2 is “Provider-disclosed, independently cross-validated.” All current models are labelled “Scaled scenario estimates derived from Large 2 anchor.” Executive summary split into three separate tables: anchor, current hosted inference, and prospective sovereign compute. Implemented throughout v1.2.
- 2. Split current from future infrastructure.** Tables now separate current hosted inference (Azure CIF 0.35) from prospective sovereign compute (France CIF 0.022, Sweden CIF 0.015). France/Sweden figures explicitly labelled “not yet operational” and “prospective.” Caveat added that energy-per-query may differ on Blackwell-based French stack. Implemented throughout v1.2.
- 3. Demote Magistral from main tables.** Magistral removed from all summary and recommendation tables. Presented as range-only with no central estimate, explicitly labelled “illustrative only.” Moved to Section 7. Implemented in v1.2.

Additional changes: Dense-to-MoE transition elevated to named section (4.2) as “central technical fault line.” Anchor caveat added to Section 3 preventing prestige spillover. Reuters Sweden timeline added (operations 2027). NVIDIA Newsroom source added for French compute announcement. Limitation 10 added: disclosure leadership \neq current product footprint. LCA-adjacent terminology clarified throughout.⁵⁴⁵⁵

14. Recommendation

For InferenceCarbon.ai’s carbon transparency tool, we recommend three distinct display tiers for Mistral:

Tier 1: Provider-disclosed (Large 2 only): Label “Provider-disclosed, LCA-anchored.” Range 0.60–1.20 gCO₂e/1k tokens. This label should not be applied to any other Mistral model.

Tier 2: Scaled scenario estimates (Large 3, Medium 3, Small 3): Label “Scaled estimate from LCA anchor.” Display current hosted-inference ranges only. Do not display France-hosted figures as current product footprint.

Model	Current hosted range (Azure)	Label
Large 3	0.40 – 0.50 (c: 0.45)	Scaled estimate from LCA anchor

Medium 3	0.25 – 0.35 (c: 0.30)	Scaled estimate — no current data
Small 3	0.04 – 0.07 (c: 0.05)	Scaled estimate — no current data

Tier 3: Prospective sovereign compute: Display France/Sweden figures only in a secondary panel, clearly marked “Prospective — not yet operational.” Do not present as current product footprint.

Magistral: Do not display in main tables. If shown at all, range-only (0.50–2.10) in a reasoning-model secondary panel with “illustrative only” label.

We commit to updating when: (a) Mistral publishes LCA data for current models; (b) Mistral Compute routing becomes known; (c) EcoDataCenter Sweden becomes operational; or (d) new independent measurements become available.

15. Key Sources

Source	Used for
Mistral AI LCA (July 2025)	Provider-published anchor: 1.14 gCO ₂ e/400 tokens
Carbone 4 (2025)	LCA methodology, peer review
Jegham et al. arXiv:2505.09598v6	Independent anchor: ~1.09 gCO ₂ e; cross-validation
Artificial Analysis (2026)	Throughput for current models
RTE France (2025)	French grid CIF: 21.7 gCO ₂ e/kWh
DCD / NVIDIA (2025)	Mistral Compute infrastructure
Reuters / BeBeez (2026)	Sweden EcoDataCenter timeline
GHG Protocol Scope 2 (2025)	Location-based methodology

End of document. Version 1.4, 19 March 2026.

Prepared by InferenceCarbon.ai Reference Team