

# Analyse économique des coûts d'inférence des LLM (2023–2026)

---

## Introduction

---

Les modèles de langage de grande taille (LLM) tels que ceux développés par **OpenAI** et **Anthropic** connaissent une adoption massive. Pourtant, leur exploitation repose sur une infrastructure coûteuse : GPU haut de gamme, consommation électrique élevée, et optimisation énergétique limitée en mono agent.

L'objectif de ce projet est de :

- Évaluer l'évolution des **coûts d'inférence** dans le temps.
  - Identifier les **seuils de rentabilité (break-even)** pour différents profils d'abonnement.
  - Proposer une **projection à 5–10 ans** des coûts et des prix nécessaires à la rentabilité.
- 

## Méthodologie

---

### 1. Sources de données :

- **EIA** (Energy Information Administration, USA) pour les prix de l'électricité commerciale.
- Estimations de prix GPU (H100, L4) via **marchés publics** et overrides manuels.
- Calculs internes : consommation énergétique, PUE (Power Usage Effectiveness), throughput (tokens/sec).

### 2. Pipeline ETL :

- Extraction (Python, API EIA + scrapers GPU).
- Transformation (nettoyage, harmonisation temporelle, intégration coûts GPU/électricité).
- Chargement dans **PostgreSQL** via `\copy`.

### 3. Analyse :

- Calcul du **coût par million de tokens** (électricité + GPU).
  - Définition de profils d'abonnements :
    - **Lite** : 200k tokens/mois.
    - **Standard** : 1M tokens/mois.
    - **Pro** : 5M tokens/mois.
  - Application d'une marge cible de **70%** pour estimer les prix break-even.
-

# Analyse

## 1. Coûts observés (2023–2026)

- Le **coût moyen par million de tokens** varie entre **1.8 et 2.2 USD** selon la période.
- Les coûts fluctuent avec :
  - L'évolution du prix de l'électricité (source EIA).
  - Les ajustements sur le prix horaire des GPU (H100, L4).

## 2. Seuils de rentabilité (break-even)

- **Lite (200k tokens/mois)** : break-even  $\approx$  **1.2–1.4 USD/mois**.
- **Standard (1M tokens/mois)** : break-even  $\approx$  **6–7 USD/mois**.
- **Pro (5M tokens/mois)** : break-even  $\approx$  **31–34 USD/mois**.

Ces niveaux restent **largement en dessous des abonnements actuels** (ChatGPT Plus : 20 USD/mois, Claude Pro : 20 USD/mois).

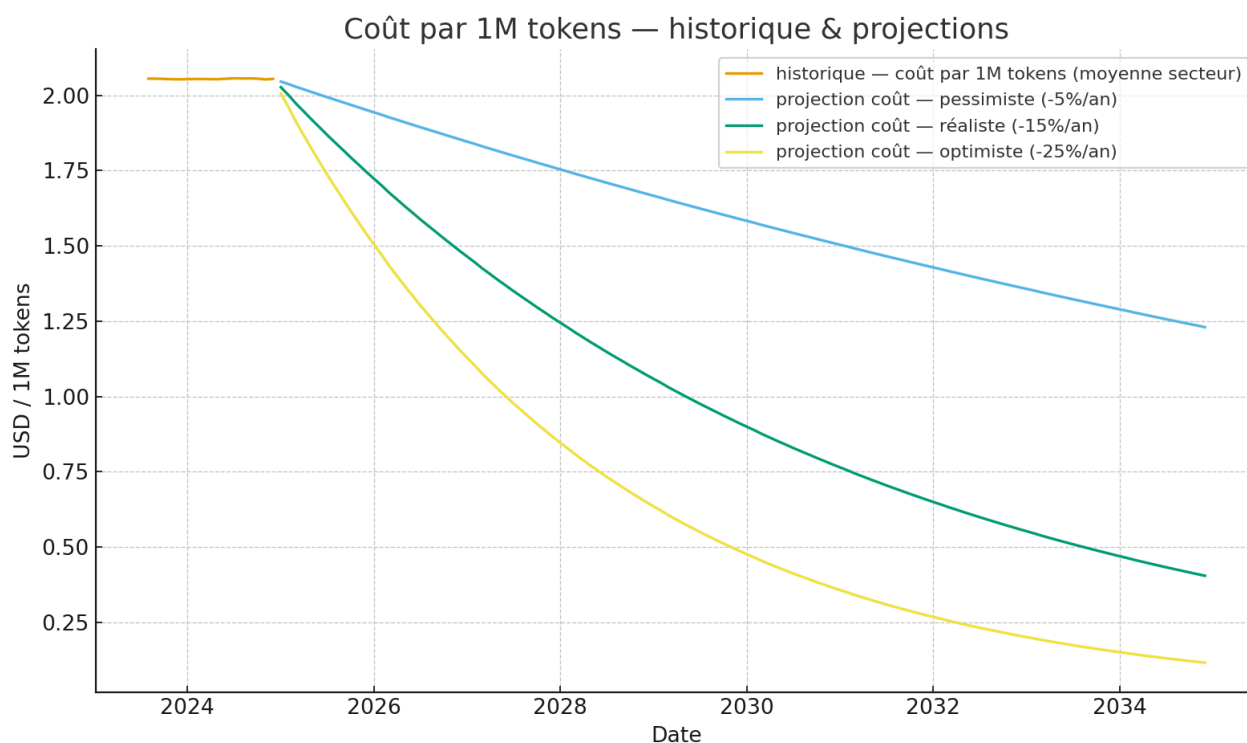
Cela indique que les entreprises subventionnent massivement l'accès, malgré des coûts importants (ici on ne parle de l'électricité mais il y a d'autres facteurs à prendre en compte.)

## 3. Déficit structurel

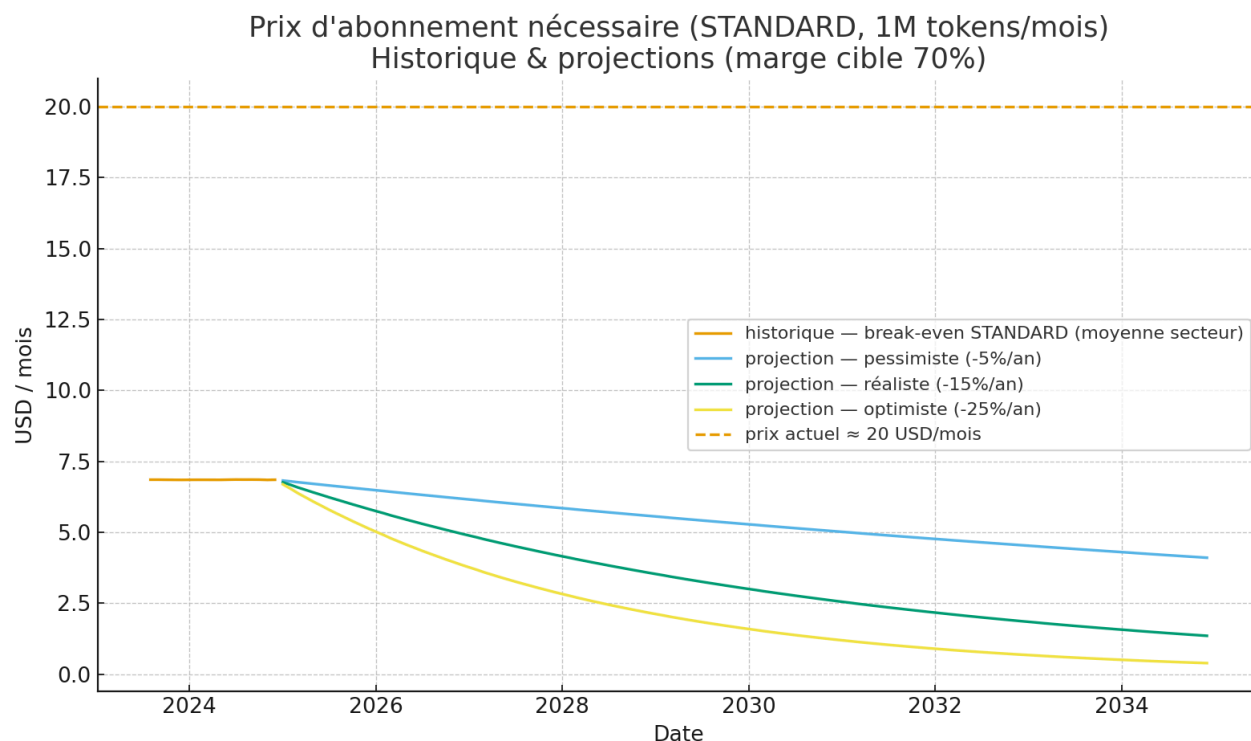
- Même sans inclure **R&D, salaires, serveurs et stockage**, l'activité reste **déficitaire**.
- Le prix actuel des abonnements ne couvre pas le coût marginal réel de l'inférence.

## Visualisations sur 5-10 ans

### 1. Courbe du coût par million de tokens (2023–2026).



## 2. Courbe fictive du prix d'abonnement nécessaire pour atteindre la rentabilité.



(Les graphiques sont produits via matplotlib, disponibles en notebook Python.)

## Discussion

Pourquoi OpenAI et Anthropic continuent malgré les pertes ?

- **Effet de réseau** : plus d'utilisateurs = plus de données = meilleurs modèles.
- **Course stratégique** : l'IA est un secteur winner-takes-all ; être leader prime sur la rentabilité immédiate.
- **Subventions massives** : Microsoft, Google et Amazon financent lourdement ces acteurs.
- **Pari sur le futur** : les coûts GPU/énergie pourraient baisser, ou les prix des abonnements être relevés.
- **Se rendre indispensable** sur le marché.

À terme, il est probable que :

- Les abonnements **augmentent progressivement** (25–40 USD/mois).
- Les offres se **segmentent davantage** (Lite, Pro, Entreprise).
- Les entreprises misent sur des **revenus annexes** (API, intégrations, produits SaaS).

## Conclusion & perspectives

- Les LLM sont exploités à perte, même en ne considérant que **GPU + électricité**.
- La viabilité économique à long terme dépendra de :
  - **Hausse des prix** côté abonnements.
  - **Optimisation énergétique** (PUE, chips spécialisés).
  - **Effet d'échelle** et amélioration des throughput.

- Nos projections indiquent une **hausse nécessaire des abonnements d'ici 5-10 ans**, sous peine de déficits insoutenables.

**Prochaines étapes :**

- Approfondir les projections sur 2030-2035 avec scénarios optimistes/pessimistes.
- Étendre l'analyse à d'autres acteurs (Mistral, Meta, Google DeepMind).
- Simuler l'impact d'une taxation carbone sur le coût final des LLM.