

Rate limits

The token limits include the input + output tokens. See the pricing [here](#) .

Model Tier RPM ¹ RPD ² TPM ³ TPD ⁴ Verified Inference "Free" 60 33,000 40k 4M (1) RPM: Request Per Minute (2) RPD: Request Per Day (3) TPM: Tokens Per Minute (4) TPD: Tokens Per Day

To access "Infinite" tier, please [apply here](#)

Rate limit headers

We set the following x-ratelimit headers to inform you on current rate limits applicable to you.

The following headers are set (values are illustrative):

Header Value Notes
retry-after 2 Seconds to wait until retrying*
x-ratelimit-limit-requests 28800 Requests per day allowed
x-ratelimit-limit-tokens 40000 Tokens per minute allowed
x-ratelimit-remaining-requests 123 Requests remaining for the day
x-ratelimit-remaining-tokens 1337 Tokens remaining for this minute
x-ratelimit-reset-requests 1337s Seconds until the daily rate limit resets
x-ratelimit-reset-tokens 1s Seconds until the minute based token limit resets
* The retry-after header is only returned if the response status code is 429 and the request was rate limited

[Models](#) [Eliza](#) [twitter](#) [github](#) [discord](#) [Powered by Mintlify](#)