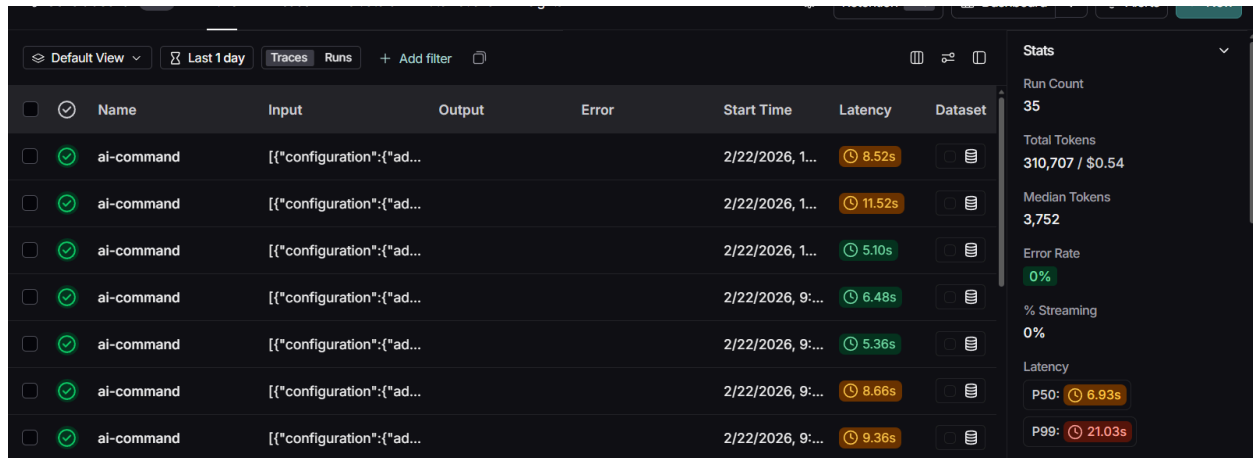When using Claude Haiku for AI Chat commands/tool/calling/etc:
The rough average of a command issued by a typical end user to my AI chat is 1.25 cents.



Assumptions:
User submits 20 commands an hour, so 25 cents an hour per user.
$6 a day per user
100 users:$18,000
1000 users:$180,000
10000 users:$1,800,000
100000 users:$18,000,000
A lot of money. It will be important to optimize everything as much as possible.

I am using the $200/month claude clode subscription and I used about 20% of my limit this week while coding, although I coded serially with one claude code session 99% of the time.