# Perplexity Model Council: the rise of Ensemble UI

For the last few years, AI adoption has felt oddly tribal. You were either on "Team GPT-4" or "Team Claude." You picked your ecosystem, paid your subscription, and learned to live with that specific model's quirks. If you wanted a second opinion, you had to copy-paste your prompt into a different tab.

Perplexity's new **Model Council** feature changes the topology of how we interact with these systems. By allowing users to query multiple frontier models (like GPT-4o, Claude 3.5 Sonnet, and Llama 3) simultaneously and view their outputs side-by-side, they are effectively bringing *ensemble methods* out of the machine learning backend and into the user interface.

## The shift to "Ensemble UI"

In machine learning, an "ensemble" combines the predictions of several base estimators to improve robustness. Usually, this happens invisibly: the system aggregates the results and serves you the winner. Model Council does something different: it leaves the aggregation to the user.

This is significant because LLMs are not search engines; they are reasoning engines. And reasoning often benefits from debate. When you see three distinct models tackle the same complex prompt, you aren't just looking for the "correct" answer. You are looking for consensus.

## Triangulation as a defense against hallucination

The most immediate implication is error correction. If three models agree on a fact and the fourth one diverges wildy, the outlier becomes obvious. This "visual triangulation" empowers the user to spot hallucinations that might have otherwise slipped through in a single-model interface.

This creates a new layer of trust. Instead of trusting a specific brand ("I trust OpenAI"), you trust the *intersection* of multiple independent systems. It acknowledges that no single model is the "God Model" anymore. They are just distinct perspectives on the same latent space.

## Commoditizing the intelligence layer

From a market perspective, this is aggressive commoditization. When models are displayed side-by-side in a uniform grid, the brand halo fades. You stop caring about the logo and start caring exclusively about the output quality. If the open-source Llama model consistently matches the performance of the proprietary models in your specific domain, the justification for expensive, closed-source moats begins to erode.

It forces a brutal meritocracy. The models are no longer protected by their walled gardens; they are fighting for your attention in real-time, on every single query.

## The user as the meta-learner

Finally, this feature forces the user to become a more active participant. We are shifting from "answer retrieval" to "answer synthesis." The cognitive load is slightly higher—you have to read and compare—but the potential for insight is deeper. You start to learn *how*

the models think: which one is verbose, which one is terse, which one is better at code, and which one is better at nuance.

We are moving away from the idea of an AI Oracle. We are moving toward an AI Committee, and for now, you are the chairperson.

**Contributor:** Alessandro Linzi