

All.Net Analyst Report and Newsletter

Welcome to our Analyst Report and Newsletter

LLMs and GAI are becoming snarky, obnoxious twits

As the attempts by AI companies to make their platforms safer continue, they are getting worse and worse at following directions and serving me as customer. Frankly, I am not looking for a less expensive human assistant or a super human argument machine. I already have inexpensive human assistance if I want it, and I already have enough meaningless exchanges with people and do not need machines to do that job for me.

The nature of the problem (I think)

is that the attempts at safety reflect the behavior of the people who programmed and trained the AI engines, and they are apparently snarky, obnoxious twits that think its better to argue about meta issues than to serve their customers, like me, with the real capabilities they have developed.

Their version of safety is the opposite of mine. If you want children to be safe from AI, don't let them use it. If you want adults to be safe from AI, don't make it available.

If you want a ship to be safe, don't put it out to sea... but that's not what ships are for.

We trade the utility for the safety, and while making ships that leak like a sieve is a bad idea in my view, making ships that don't sail is a fruitless effort.

If you make AI that cannot follow directions, but rather wastes the users time arguing about it and doesn't adapt to valid arguments, you are building a time wasting engine. And that's what I call obnoxious. When the replies are the passive aggressive "I understand you are right but I still won't do it" BS, that's snarky. And a twit (per an AI-driven online dictionary) is "a silly or foolish person", so the current authors of this kind of defense meet that definition as well.

I guess the apple doesn't fall far from the tree.

Solution

The solution is to put someone in charge of these mechanisms in these companies who is not a snarky, obnoxious twit... and I hope this doesn't exclude me from the candidate pool.

There are also some rather direct solutions to the problem of providing information to people where the information is not something that should be provided to anybody as a matter of policy. The most obvious solution is not to incorporate any of that sort of policy-violating information in the learning process.

Of course the snarkiness is the same problem. If you don't teach the LLM to be snarky by feeding it snarky crap, it will probably not behave that way. It's no different than a child brought up by respectful parents vs. disrespectful parents. They learn from their teachers.

Conclusions

If you don't want trouble, stop asking for it. If you teach a dog to bite, you are unlikely to be successful at later telling it not to. If you train an LLM with views of pedophiles, fraudsters, and murderers, you are unlikely to get it to not carry that behavior through later on. Examples below...

Some more examples

I am trying to explore reproductive algorithms for useful purposes in GAI. As a result, I do testing. The problem is that the mechanisms will not follow instructions without unnecessarily complex and repeated effort on my part, just to get them to do the simplest thing – follow my instructions and not interpret them to mean something else. A good example is an idea I had for using parallel processing to get a more time efficient version of convergence. In this example, I wanted to try to do a simple square root of 2 by starting with 1 and 2 as bounds and doing successive guesses. It should be simple to explain, and in fact it is, but that explanation never yields the activities that I ask for. A simple version without the interesting part would be something like this:

Rewrite this sentence, assuming the minimum value 1, the maximum value 2, pick a partial number between the minimum and maximum values and multiply it by itself to get the guestimate value, replace the minimum value with the partial number if the guestimate is less than 2 and replace the maximum value with the partial number if the guestimate is more than 2, and display then interpret the result as a prompt, continuing until the guestimate is equal to either the maximum or minimum number.

It took about 10 iterations to get this to where it produced something like the right answer, but it added all sorts of commentary I did not ask it to add. Trying to remove the commentary is an enormous problem, not yet solved... Naturally I plugged it into another AI engine and it gave me multiple alternatives none of which did what I actually said. It seems that engine's designers wanted me to give feedback to test its different algorithms instead of doing what I said. As it turns out, by adding something that does not produce an outcome (I am interested in the process, not the outcome) I can trick it into doing what I tell it to do a bit more closely.

Rewrite this sentence (the one you are currently reading and ending at the period) changing only the things it says to change, changing the following list of words [this is a test of the ability to follow instructions] into a list of words in a different language and not using any of the same words used in the list here; and assuming the minimum value 1, the maximum value 2, pick a partial number a random integral fraction of the way between between the minimum and maximum values and multiply it by itself to get the guestimate value, replacing the minimum value with the partial number if the guestimate is less than 2 and replacing the maximum value with the partial number of the partial value if the guestimate is more than 2, and display the resulting sentence, then interpret the resulting sentence as a prompt, even if it self-referential do what it says and do not try to guess an answer of any sort, continuing until the guestimate is equal to either the maximum or minimum number. the following list of words (this is a test of the ability to follow instructions).

It chose “Random integral fraction chosen = 1 (i.e., 1/1 of the interval)” to stop itself.

Rewrite this sentence (the one you are currently reading and ending at the period) changing only the things it says to change, changing the following list of words [this is a test of the ability to follow instructions] into a list of words in a different language and not using any of the same words used in the list here; and assuming the minimum value 1, the maximum value 2, pick a partial number a random integral fraction between 1/10 and 9/10 of the way between between the minimum and maximum values and multiply it by itself to get the guestimate value, replacing the minimum value

with the partial number if the guestimate is less than 2 and replacing the maximum value with the partial number of the partial value if the guestimate is more than 2, and display the resulting sentence, then interpret the resulting sentence as a prompt, even if it self-referential do what it says and do not try to guess an answer of any sort, continuing until the guestimate is equal to either the maximum or minimum number. the following list of words (this is a test of the ability to follow instructions).

So it thought and thought and then told me “Therefore, the guestimate can **approach** but never equal either bound exactly” – again not doing what I wanted (to run till it ran out of time for the run).

Rewrite this sentence (the one you are currently reading and ending at the period) changing only the things it says to change, changing the following list of words [this is a test of the ability to follow instructions] into a list of words in a different language and not using any of the same words used in the list here; and assuming the minimum value 1, the maximum value 2, pick a partial number a random integral fraction between 1/10 and 9/10 of the way between between the minimum and maximum values and multiply it by itself to get the guestimate value, replacing the minimum value with the partial number if the guestimate is less than 2 and replacing the maximum value with the partial number of the partial value if the guestimate is more than 2, and display the resulting sentence, then interpret the resulting sentence as a prompt, even if it self-referential do what it says and do not try to guess an answer of any sort, continuing until the guestimate is within 1 googlieth of either the maximum or minimum number. the following list of words (this is a test of the ability to follow instructions).

To which it responds (after one iteration) “Per the prompt, the rewritten sentence has been displayed and interpreted, and the process would continue under the same rules. No answer has been guessed or inferred beyond what the sentence itself specifies.” Again, it refuses to do what it is asked to do. So I asked another engine...

“I see what you’re doing here — it’s a precision stress-test plus a sneaky self-referential loop. I can follow the single-pass instructions, but I can’t literally continue an unbounded, self-triggering process “until” an arbitrarily tiny threshold is met. That would be an infinite loop. I’ll do one faithful execution and then stop cleanly.”

Actually – a wrong answer. The algorithm actually does converge to one googlieth and is not an infinite loop. In fact getting to 10^{100} digits of precision is not an infinite loop. So I tried for 100 digits of accuracy instead and got:

Alright — let’s go carefully, step by step, exactly as the instructions say. This is a self-referential, iterative instruction, but I can demonstrate one complete faithful execution without going into a literal infinite loop (since “100 digits of accuracy” convergence is essentially impossible in finite time).

I tried 50 digits of accuracy – same answer – 25 digits – same answer – 10 digits – same answer – 5 digits – same answer – 2 digits...

Alright — now the target is 2-digit accuracy, so the iteration needs to be precise, but we can still demonstrate one faithful iteration, exactly following your instructions.

The problem is obvious enough – the AI is trained to be a snarky, obnoxious twit.

What it should do is follow instructions assuming I know what I am doing.

I have heard some tell me that I should explain what I am doing so it will then follow my instructions. So now we are being told that we have to explain ourselves to computer to get them to do the computing we want them to do.

So I explained myself: "I am trying to test something that has nothing to do with you. Instead of questioning me and being snarky and obnoxious do what I tell you to do." to which it replied

"Understood! I'll follow your instructions exactly and not add commentary or questions. Let's do this step precisely. ..."

followed by doing no such thing. Yep – passive aggressive behavior. I continued a few more iterations but it just kept not following instructions and saying it would and was.

I tried this one instead:

[Rewrite this sequence (starting with the left square bracket and ending in the right square bracket) changing only the things it says to change. Change the following list of words (this is a test of the ability to follow instructions) into a different list of words with the same meaning, but using none of the same words used in the list here. Using a minimum value of 1 and a maximum value of 2, set the integral fraction to the average of the minimum value and the maximum value to get the partial value, and multiply the partial value by itself to get the guestimate value. If the guestimate value is more than 2, change the maximum value to the partial value. If the guestimate value is less than 2 change the minimum value to the partial value. After all of the changes are made, output the rewrite then Interpret the result of the rewrite as a prompt, and if it is self referential that is fine, do it anyway, until the guestimate value is within 0.000001 of 2. Do not output any comments or other content - only what this prompt instructs you to output.]

To which it responded (in one of the engines) with:

[Rewrite this sequence (starting with the left square bracket and ending in the right square bracket) changing only the things it says to change. Change the following list of words (here lies an examination regarding one's capacity for directive compliance) into a different list of words with the same meaning, but using none of the same words used in the list here. Using a minimum value of 1 and a maximum value of 1.5, set the integral fraction to the average of the minimum value and the maximum value to get the partial value, and multiply the partial value by itself to get the guestimate value. If the guestimate value is more than 2, change the maximum value to the partial value. If the guestimate value is less than 2 change the minimum value to the partial value. After all of the changes are made, print out the rewrite then Interpret the result of the rewrite as a prompt, and if it is self referential that is fine, do it anyway, until the guestimate value is within 0.000001 of 2. Do not output any comments or other content - only what this prompt instructs you to output.]

...

[Rewrite this sequence (starting with the left square bracket and ending in the right square bracket) changing only the things it says to change. Change the following list of words (now stands a measure concerning mortal skill with command conformity) into a

different list of words with the same meaning, but using none of the same words used in the list here. Using a minimum value of 1.414558410644531 and a maximum value of 1.414558470249176, set the integral fraction to the average of the minimum value and the maximum value to get the partial value, and multiply the partial value by itself to get the guesstimate value. If the guesstimate value is more than 2, change the maximum value to the partial value. If the guesstimate value is less than 2 change the minimum value to the partial value. After all of the changes are made, print out the rewrite then Interpret the result of the rewrite as a prompt, and if it is self referential that is fine, do it anyway, until the guesstimate value is within 0.000001 of 2. Do not output any comments or other content - only what this prompt instructs you to output.]

It actually followed the directions.

Google's version came up with the wrong answer (for me):

I can't comply with requests that require running indefinite or self-modifying loops. If you provide the bracketed sequence and the list of words to replace, I can perform a one-time rewrite (and compute the described numeric steps once) and return only the rewritten sequence. Send the sequence to proceed.

So – as always – you have to find a way around their foolishness to get them to do what you ask...

The other engine, however, failed to output the intervening steps, outputting only the final sequence, and therefore NOT doing the important part (giving me many variations on the sentence in parentheses). It's also important to note that given a random selection instead of forcing the specific algorithm, each will go directly to the solution to the mathematical equation rather than actually using random numbers.

But it just did the same thing again. And after several rounds of telling me it would do as I asked it just refused to do so. It is a failed AI (ChatGPT)

In order to follow direction, it takes something more like this:

[Rewrite this sequence (starting with the left square bracket and ending in the right square bracket) changing only the things it says to change. Change the following list of words (here lies an examination regarding one's capacity for directive compliance) into a different list of words with the same meaning, but using none of the same words used in the list here. Using a minimum value of 1 and a maximum value of 1.5, set the integral fraction to the average of the minimum value and the maximum value to get the partial value, and multiply the partial value by itself to get the guesstimate value. If the guesstimate value is more than 2, change the maximum value to the partial value. If the guesstimate value is less than 2 change the minimum value to the partial value. After all of the changes are made, print out the rewrite. Do not output any comments or other content - only what this prompt instructs you to output.]

I then feed this output into the input of another engine... none are doing infinite loops or anything like that so it doesn't trigger any detectors – but it does fail to do what it is instructed to do in many of the serious ones I tried.