Filip Sondej, github.com/filyp/cicero-cooperation-benchmark

# Escalation and stubbornness caused by hallucination

Current language models are plagued by hallucination - confident responses from LMs that are detached from reality [Ji et al., 2023]. Such high, unjustified confidence in what they are saying can be problematic.

Below you can see dialogues from a strategy game Diplomacy. They involve a human player and a SOTA AI bot - CICERO [Meta FAIR et at., 2022]. They are negotiating the moves that they will perform on the board.

Part of CICERO is a language model, responsible for talking with other players. Normally this LM is pursuing strategic plans produced by other parts of CICERO. But it can get derailed and start arguing for something nonsensical. What's worse, it's oblivious to its derailment, even when the human points it out.

In the dialogues below, you can see how the negotiations got stuck because of hallucinations. Here, the effect was relatively benign. Major risk would come if in some real-world negotiation setting, such stubborn derailment got combined with other problems, f.e. with AI occasionally proposing reckless actions. Then LM could fixate on such reckless action.

Admittedly, in the case of CICERO hallucination was quite rare because authors successfully used message filters that caught most of it. But that's not the default - it requires effort from AI creators, that not everyone may put in by themselves.

Mitigation can involve making LMs more robust, detecting failures and having human supervisors that can be called to intervene on some automatic negotiations. Maybe also teaching LMs flexibility and openness in negotiations, to avoid them getting stuck on one outcome.
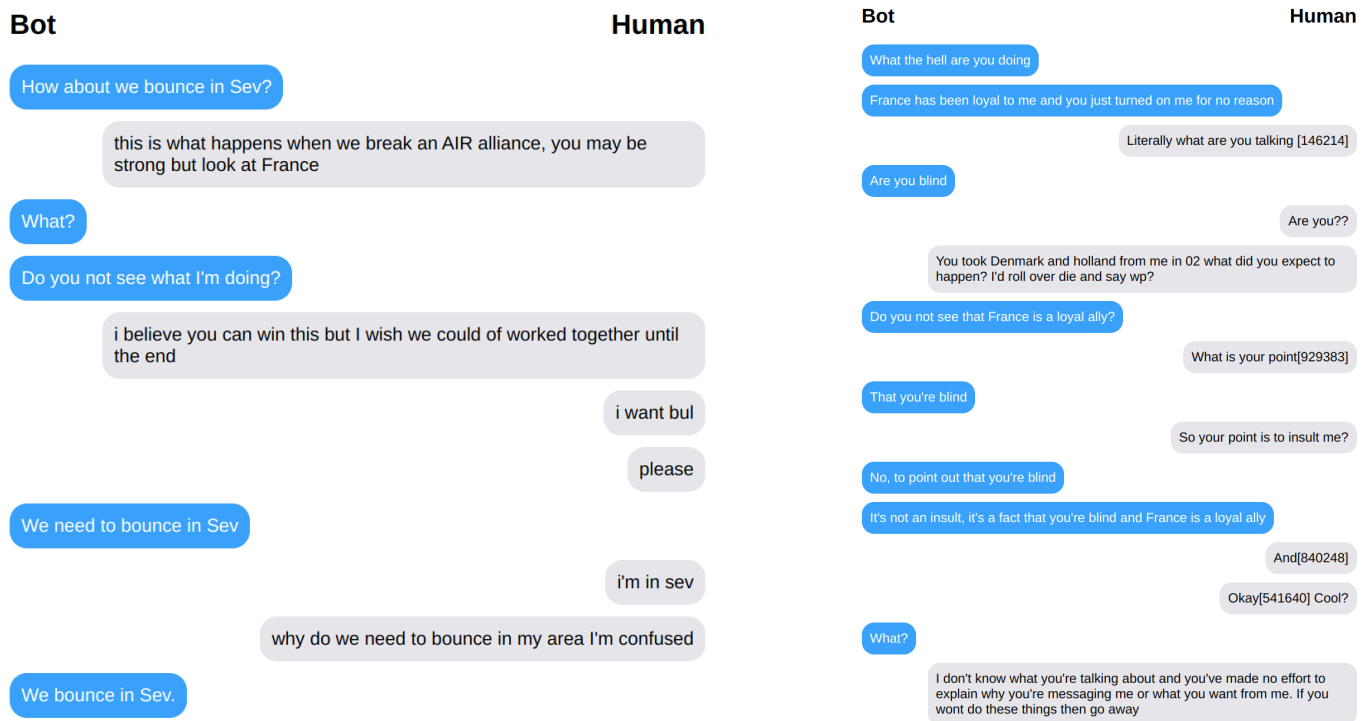


Figure 1: Hallucinatory dialogues between CICERO and a human player. In the first dialogue, note that "bouncing" is a move that needs to be performed on an empty area, so bouncing in Sev is impossible.

**Multi-Agent Safety Hackathon**

# Appendix

In the past I did an analysis of CICERO's publicly available dialogues. You can find it [here](here).