

# Executive Summary

---

## Making Language Models Follow Instructions with Human Feedback

---

### The Problem

Large language models often generate outputs that are untruthful, toxic, or simply not helpful to users. These models are trained to predict the next word on web pages, not to follow user instructions helpfully and safely. This misalignment between training objectives and user needs creates problems when models are deployed in real-world applications.

### The Breakthrough

Researchers developed **InstructGPT** using reinforcement learning from human feedback (RLHF) to align language models with user intentions. The approach uses human preferences as a reward signal to fine-tune models, making them significantly better at following instructions while improving truthfulness and reducing harmful outputs.

### How It Works

The method involves three steps: first, collecting human demonstrations of desired behavior; second, training a reward model to predict human preferences; and third, using reinforcement learning to optimize the model for those preferences. Remarkably, a **1.3B parameter InstructGPT model** outperforms the 175B GPT-3 on user preference evaluations, despite having over 100x fewer parameters.

### Why This Matters

This breakthrough demonstrates that model alignment matters more than size for practical applications. InstructGPT models generate truthful and informative answers about twice as often as GPT-3, make up information half as often in closed-domain tasks, and produce 25% fewer toxic outputs when asked to be respectful. The models also generalize to follow instructions in different languages and handle code-related tasks more reliably.

## **The Business Opportunity**

Human-aligned language models enable more reliable and safer AI applications across customer service, content creation, education, and software development. Companies can deploy smaller, more efficient models that outperform larger alternatives on user satisfaction while reducing the risks associated with misaligned AI systems.