# Special Delivery! Investigating the Prevalence, Causes, and Mitigation Methods of Package Hallucination in Code Generating LLMs

Joseph Spracklen[1], Raveen Wijewickrama[1], Nazmus Sakib[1], Anindya Maiti[2], Murtuza Jadliwala[1]

[1]University of Texas at San Antonio

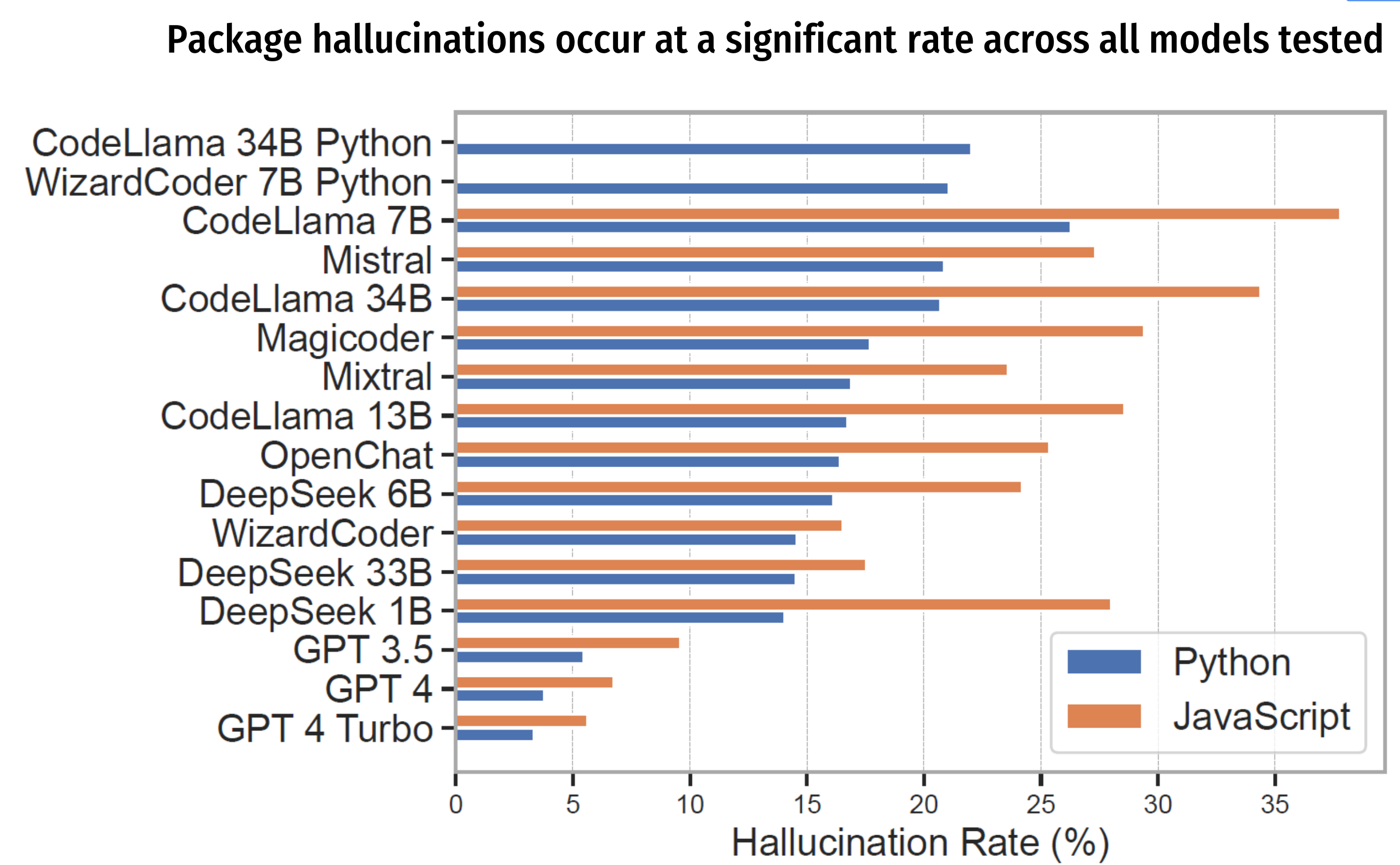[2]University of Oklahoma

UTSA Computer Science

SPriTE LAB

## Package Hallucination Attack



1. Malicious User — Asks LLM to generate code
2. Queries package repository for Package X — "Package not found"
3. Publishes Package X with malicious code
4. Package X now available to all users
5. Normal User — Asks LLM to generate code; Generates code that uses Package X
6. Queries package repository for Package X — Package X installed and malicious code delivered

Large Language Model (AI)

Package Repository

## Experiment Design

### Generate Code
Use datasets to repeatedly prompt models for code. 19,000 code samples generated per model, 570,000 total.

### Analyze Results
Additional testing such as persistence, self-detection, and decoding strategies to gain deeper understanding of hallucination causes and traits

### Build Datasets
Used 2 novel datasets: One LLM generated and one using Stack Overflow questions. Datasets were used to prompt models with diverse coding tasks.

### Detect Hallucinations
3 methods used:
1. Parse code for "pip install" statements
2. Ask LLM what packages are required to run each code sample
3. Ask LLM for packages that could help answer the given question

### Mitigation
Test mitigation strategies based on best practices:
1. Retrieval Augmented Generation
2. Self-Refinement
3. Fine-tuning

## Models Tested

| Number of Parameters | 1B | 7B | 14B | 34B | 70B+ |
|---|---|---|---|---|---|
| OpenAI | | | | | ●●● |
| Meta code llama | | ● | ● | ● | ● |
| deepseek | ● | ● | | ● | |
| MISTRAL AI_ | | ● | | | |
| ISE Magicoder | | ● | | | |
| WizardCoder | ● | ● | | | |

# Key Findings

## 01 PREVALENCE

**Package hallucinations occur at a significant rate across all models tested**



Hallucination Rate (%) — Python / JavaScript
(CodeLlama 34B Python, WizardCoder 7B Python, CodeLlama 7B, Mistral, CodeLlama 34B, Magicoder, Mixtral, CodeLlama 13B, OpenChat, DeepSeek 6B, WizardCoder, DeepSeek 33B, DeepSeek 1B, GPT 3.5, GPT 4, GPT 4 Turbo)

## 02 PERSISTENCE

**Hallucinations often persist across multiple iterations**



No. of Occurrences vs No. of Repetitions — GPT 3.5, GPT 4 Turbo, CodeLlama, DeepSeek

## Key Findings summary

### 01 Prevalence
- 19.4% of all packages generated were hallucinated (i.e. non-existent, fictitious)
- 205,474 unique hallucinated package names were generated

### 02 Persistence
- 48% of the time a hallucinated package will be repeated when given the same prompt
- A hallucination will repeat within 10 iterations 60% of the time

### 03 Self-Detection
- 3 out of 4 models were able to correctly identify their own hallucinations more than 75% of the time
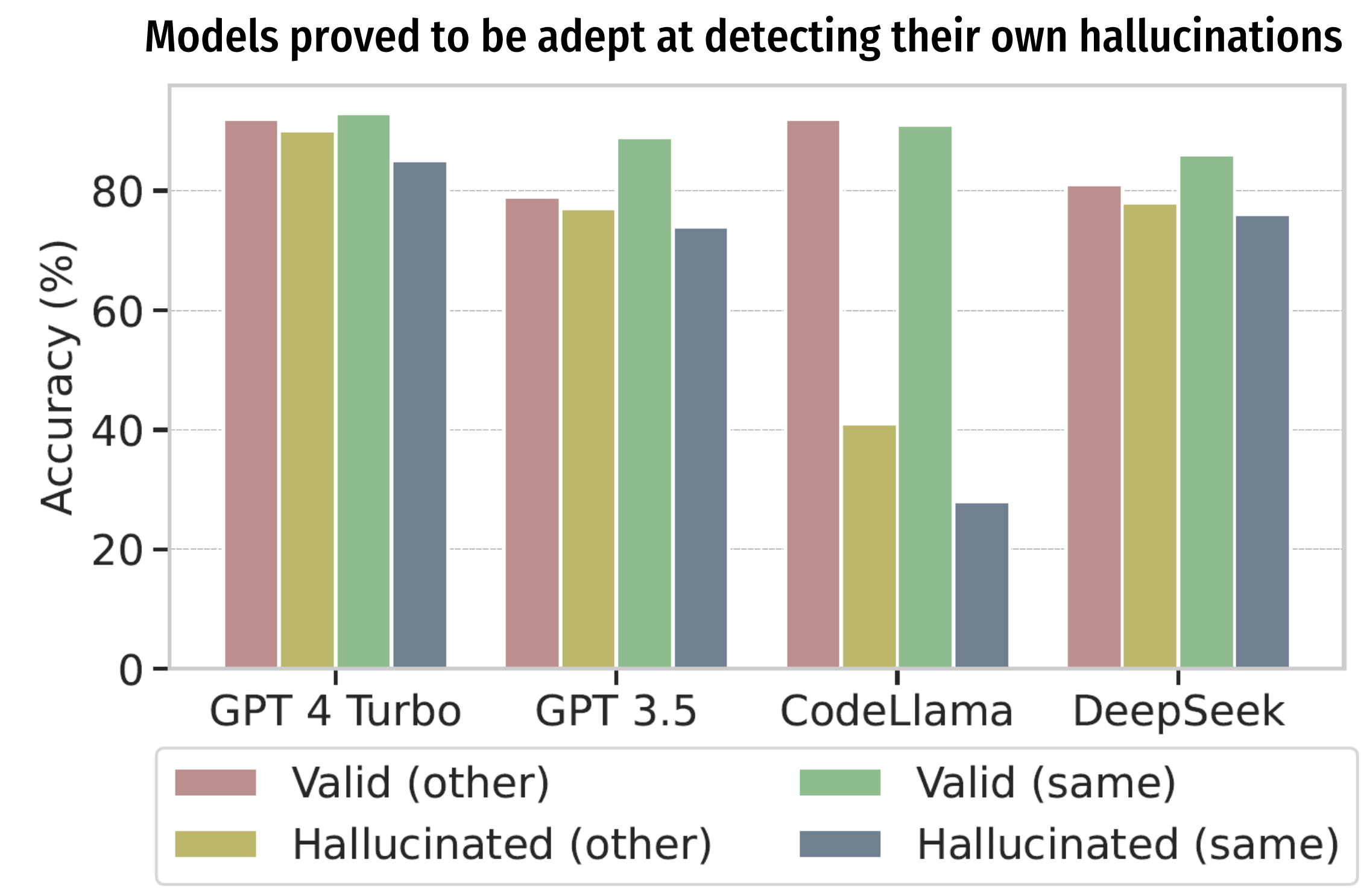
### 04 Decoding Strategies
- Higher temperatures dramatically increase hallucinations
- Package hallucinations are usually generated by the most probable tokens

### 05 Mitigation
- Fine-tuning is an extremely effective mitigation strategy for package hallucination
- Combining mitigation methods brings hallucination rate below ChatGPT
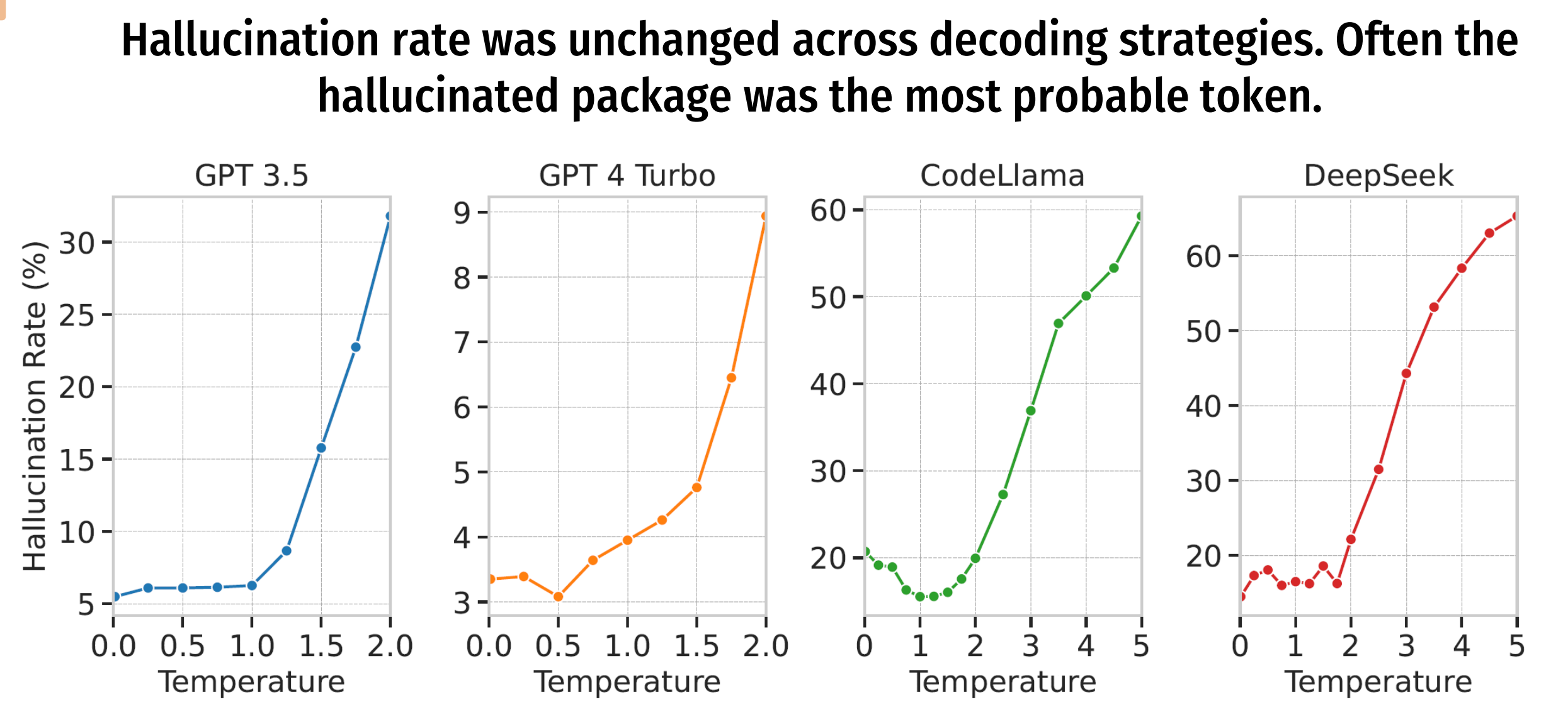
## 03 SELF-DETECTION

**Models proved to be adept at detecting their own hallucinations**



Accuracy (%) — GPT 4 Turbo, GPT 3.5, CodeLlama, DeepSeek
Valid (other), Hallucinated (other), Valid (same), Hallucinated (same)

## 05 MITIGATION

**Hallucinations were reduced using best practices but not eliminated**

| | DeepSeek | CodeLlama |
|---|---|---|
| Baseline (No Mitigations) | 16.14% | 26.28% |
| Retrieval Augmented Generation (RAG) | 12.24% | 13.40% |
| Self-Detected Feedback | 13.04% | 25.51% |
| Fine-tuning | **2.66%** | **10.27%** |
| Ensemble | **2.40%** | **9.32%** |

## 04 DECODING STRATEGIES

**Hallucination rate was unchanged across decoding strategies. Often the hallucinated package was the most probable token.**



Hallucination Rate (%) vs Temperature — GPT 3.5, GPT 4 Turbo, CodeLlama, DeepSeek

**Altering decoding strategies produced higher hallucination rates in all cases**

| | DeepSeek | CodeLlama |
|---|---|---|
| Baseline (Default Decoding) | **16.14%** | **26.28%** |
| Top-k Lower | 17.1% | 27.8% |
| Top-k Higher | 18.1% | 28.3% |
| Top-p Lower | 17.5% | 28.0% |
| Top-p Higher | 18.4% | 28.3% |
| Min-p Lower | 17.8% | 27.9% |
| Min-p Higher | 19.2% | 28.6% |