## The Language of Creation How Generative AI Challenges Intuitions and can Reshape Design Process

J.D. Zamfirescu-Pereira · EPIC Research Retreat · April 16, 2025

I'd like a game to help my kid learn to multiply.



# JUST ASK!









### We can do better!

I'd like a game to help my kid learn to multiply.



I'd like a game to help my kid learn to multiply.

## JUST ASK! THEN ITERATE!

### Conversation





I'd like a game to help my kid learn to multiply.

## JUST-ASK! THEN ITERATE!

### Conversation









### **JUST ASK!**

### EXPLORE

Why Johnny Can't Prompt: How Non-AI Experts Try (and Fail) to Design LLM Prompts

J.D. Zamfirescu-Pereira zamfi@berkeley.edu UC Berkeley Berkeley, CA, USA

> Bjoern Hartmann ern@eecs.berkeley.ed UC Berkeley

Richmond Wong rwong34@gatech.edu Georgia Institute of Technology tlanta, Georgia, USA

Qian Yang Cornell University

enewed ex GPT-3 [8], sations out-iming-skill

### Why Johnny **Can't Prompt** (CHI '23)

### Dreamsheets (CHI '24)

DreamSheets

HEETS provides a 2D spreadsheet interface that Text-to-Image users can use to au for rapid prompt-image exploration. In this example derived from sheet-systems authored by expert participants (E1, E2, E4, E5) an initial idea for a prompt (A2) is used to generate an image (B2) using DREAMSHEETs'S TTI() function. They modify the prompt with EMBELLISH(), an LLM-based function, to generate the more stylized (B3). Other LLM functions DIVERGENTS() and GPT LIST() are used to generate series of art styles and eras in art history: two chosen axes for promot explosion for the series of art styles and eras in art history: two chosen axes for promot explosions the series of art styles and eras in art history.

### Herding AI Cats (DIS '23)

Berkeley, California, USA Berkeley, California, USA Berkeley, California, USA Qian Yang Bioern Hartmann qianyang@cornell.edu bjoern@eecs.berkeley.edu Cornell University University of California, Berkeley Berkeley, California, USA Ithaca, New York, USA



Figure 1: UX design journeys typically begin with great uncertainty and end with a single point of focus [18]. In this project, chatbot design by prompting GPT felt like a journey of never-ending uncertainty.



LLM-supported

**Beyond Code** Generation (CHI '25)

de, running application 🛛 📄 Generated requirements, design ideas, rationales, alternatives 🛛 🚺 Generated assessments, trade-offs

generation with every prompt change. PAIL, our IDE (right) that helps users abstract up, pulls developers them towards a deepe inderstanding of the problem space, helps them explore alternative problem formulations and solutions, and tracks design goals and requirements, surfacing implicit decisions-but with a breadth and depth of information that can be overwhel Abstract

In this work, we explore explicit LLM-powered support for iterative originated changes to code and other info in this work, we explicit explicit explorted support of relative design of computer programs. Frogram design, like other design ac-tivity, is characterized by navigating a space of alternative problem formulations and associated solutions in an iterative fashion. LLMs are potentially powerful tools in helping this exploration, however, agents should surface to program designers and when. by default, code-generation LLMs deliver code that represents a

a broader design space-but also struggle to keep up with LLM



### xible Sense-Making for

Kyu Won Kim Berkeley Institute of Design UC Berkeley, USA kwk@berkelev.edu

Bjoern Hartmar erkeley Institute of Design UC Berkeley, USA



### DECIDE

### Who Validates the Validators? Aligning LLM-Assisted Evaluation of LLM Outputs with Human Preferences

Shreva Shankar I.D. Zamfirescu-Pereira UC Berkeley UC Berkeley Berkelev, California, USA Berkelev, California, USA Aditva G. Parameswara Berkeley, California, USA

ABSTRACT

Ian Arawio Université de Montr Montréal, Québec, Canada

Shreya Shankar, J.D. Zamfirescu-Pereira, Björn Hartmann, Adity Parameswaran, and Ian Arawjo. 2024. Who Validates the Validators? A

Biörn Hartmann UC Berkeley

Berkeley, California, USA

joern@eecs.berkel

### Who Validates the Validators? (UIST '24)

### SPADE: Synthesizing Assertions for Large Language Model **Pipelines**

Shreya Shankar<sup>1</sup>, Haotian Li<sup>2</sup>, Parth Asawa<sup>1</sup>, Madelon Hulsebos<sup>1</sup>, Yiming Lin<sup>1</sup>, J.D. nfirescu-Pereira<sup>1</sup>, Harrison Chase<sup>3</sup>, Will Fu-Hinthorn<sup>3</sup>, Aditya G. Parameswaran<sup>1</sup>, Eugene Wu<sup>4</sup> <sup>1</sup>UC Berkeley, <sup>2</sup>HKUST, <sup>3</sup>LangChain, <sup>4</sup>Columbia Universit yashankar,pgasawa,madelon,yiminglin,zamfi,adityagp} @berkeley.ed otian.li@connect.ust.hk, {harrison.wfh} @langchain.dev, ewu@cs.columbia.e

ABSTRACT Operationalizing large language models (LLMs) for custom, repo itive data pipelines is challenging, particularly due to their unpre dictable and potentially catastrophic failures. Acknowledging the evitability of these errors, we focus on identifying *when* LLM ay be generating incorrect responses when used repeatedly a

automatic: SPADE ana functions and accur LLM pipe 14% and baselines

### SPADE (VLDB '24)

## EDUCATE

### 61A Bot Report: AI Assistants in CS1 Save Students Homework Time and Reduce Demands on Staff. (Now What?)

This and Real		us on stun	
J.D. Zamfirescu-Pereira zamfi@berkeley.edu UC Berkeley EECS	Lary larynqi@b UC Berke	n Qi erkeley.edu eley EECS	Björn Hartmann bjoern@eecs.berkeley.edu UC Berkeley EECS
Berkeley, California, USA	Berkeley, Ca	lifornia, USA	Berkeley, California, USA
John E denero@b UC Berke Berkeley, Ca	DeNero erkeley.edu eley EECS lifornia, USA	Narges J norouzi@b UC Berke Berkeley, Cal	Norouzi erkeley.edu ley EECS lifornia, USA
Abstract		1 Introductio	n
LLM-based chatbots enable students to get imm help on homework assignments, but even a tho bot may n developn work ass 2000 stu semester	nediate, interactive pughtfully-designed	The recent wide ava given students in in	ulability of Large Language Models (I atroductory Computer Science (CS)
the com Bot wrat			Ron

### 61A Bot Report **(SIGCSE '25)**

### **Trucks Don't Mean Trump: Diagnosing Human Error in Image Analysis**

Jerry Chen

J.D. Zamfirescu-Pereira versity of California, Berkeley Berkelev, USA Allison Koenecke licrosoft Research and Cornel

University

Cambridge, USA

Stanford Univ Stanford, USA Nikhil Garg

Cornell Tec New York City, USA

Stanford Universit Stanford, USA Emma Pierson Cornell Tech New York City, USA

Emily Wen

### ABSTRACT

Algorithms provide powerful tools for detecting and dissecting human bias and error. Here, we develop machine learning methods to to analyze how humans err in a particular high-stakes task: image

atterns of error in human judgment is of interest to a wide range of fields including psychology, computer science, and behavioral eco-nomics [15, 22, 23, 44, 54]. Algorithmic and statistical approaches have uncovered systematic human biases-e.g., race or gende

interpretation. W predictions of w Joe Biden in th image. We show Bayes optimal c decomposition and further ide lead humans a human-in-the-

Trucks Don't Mean Trump (FAccT '22)

9







# Chat is easy, but humanlike interaction leads to misunderstandings—and hides what LLMs are good at

Zamfirescu-Pereira et al. Why Johnny Can't Prompt: How Non-Al Experts Try (and Fail) to Design LLM Prompts. (CHI '23)



# Natural language is "easy"

- We humans have tons of experience with it.
- Language has flexibility
  - ...these aren't properties of natural language in isolation.



• They're properties of a system: language + speaker + listener + medium.

# "Don't let your UI write a check your AI can't cash."

- Prof. Eytan Adar, School of Information • University of Michigan

### **LLMs and Prompt Engineering** Studying ChatGPT before there was ChatGPT

You are a recipe instruction bot.



Oh, there's bacon in this? Can we make it vegan?

Absolutely! Just use vegan bacon instead of bacon.



You are a recipe instruction bot.

Be funny. But don't suggest vegan bacon!

- $\checkmark$  Easy to describe desired behavior!
- Highly interactive  $\checkmark$ dialogue!
- Fluid, fluent language output!
- $\checkmark$  Support long tail of user requests & preferences!
- Highly sensitive to inputs-needs engineering effort!





### LLMs and Prompt Engineering Studying ChatGPT before there was ChatGPT





Absolutely! Just use vegan bacon instead of bacon.



Fasy to describe a rec: en find behavior! instruction b **Flighly** Interactive But don't dialogue! vegan bacon!

> Fluid, fluent language output!







### **Understanding** *Prompt Engineering* Programming, via Natural Language, to define ChatBot behavior

- Goal: understand what prompt engineering is, through a concrete example.
- Technology probe: a best-effort system, close to the edge of feasibility, asks:
  - What are the critical challenges in a specific domain?
  - Where do our best ideas succeed or fall short?
- Built support tool for making task-oriented dialogue assistants: BotDesigner.
- Use BotDesigner to build a recipe instruction bot.

### BotDesigner Hi, today I'm going to help you with Lean In: The Ultimate 👤 Great, let's get started. What's the Natural Language Programming to define a ChatBot's behavior









Pre	amble Instructions to the chatbot: task description & steps, example dialog, etc. (Require
You foll	are a recipe instruction bot. Engage in conversation to walk the user through the owing recipe for Mixed Veggie Tempura:
Ing	redients:
- 1   - 1	Maitake mushroom Acorn squash
- 2	lemons
- 1 (	cup flour
irs	tTurns Get the dialog going, alternate 🏟 and 🗣. (Recommended.)
E	BotFirst
urn 🍋	IS
<u>*</u>	Hi, today i'm going to help you make Mixed veggle Tempura. Are you ready?
<b>R</b> efe	Yes, I'm ready. What's should I do first?
6	
÷	add new
en	<b>NINCER</b> Reminds the bot of something just before prompting for a response. (Optional.)





# What is *Prompting?* An Example



- Don't say that when you use your hands you get a little bit more of that natural flavor in there.
- I like using my hands because then you get a little bit more of that natural flavor in there.
- Explain that this should be done with the user's hands, except when the mushroom is dry or spongy on the bottom. When the mushroom is dry or spongy on the bottom, use a paring knife. Don't explain why hands are used.

## **Non-Al Experts – Humans' Intuitive Behaviors**



### System behavior models rooted in human-human social and instructional interactions



Zamfirescu-Pereira et al., Why Johnny Can't Prompt: How Non-AI Experts Try (and Fail) to Design LLM Prompts (CHI '23)

### Ad hoc experimentation: overgeneralization from single successes or failures



# **Non-Al Experts – Humans' Intuitive Behaviors**





- Bias for direct instruction over in-line examples • Even after observing the effectiveness of examples!

  - Here are some examples of speech patterns to use: Bot: So the first thing we're gonna do is take this very brain-looking [...]
- Polite prompts despite visible frustration Please list only one ingredient at a time.
- Expected "understanding" not "priming"
- Do not say ABC yields ABC
- Premature failure declaration: *noncomplian* Tell some jokes
- Premature victory declaration: limited succe Limited interest in systematic testing.







# What's going on here?

- Aligning with LLMs isn't like aligning with people.
- These aren't **lies**, exactly.
- Dialogue that triggers expectations of...
  - Understanding
  - Agreement
  - Reasoning
- ...isn't reliable for LLMs!
- We tend to "ascribe full intelligence on the basis of partial evidence"

Can we make this recipe vegan?



Yes, you can use vegan bacon.



Can I just skip bacon?



Yes, you can skip bacon.



Would it taste as good?



Yes, it would taste as good.



How do you know?



know because I've tasted it.





# What's going on here?

- Aligning with LLMs **isn't like** aligning with people.
- These aren't **lies**, exactly.
- Dialogue that triggers expectations of...
  - Understanding
  - Agreement
  - Reasoning
- ...isn't reliable for LLMs!
- We tend to "ascribe full intelligence on the basis of partial evidence"

- Systems shouldn't cause humans to rely on **intuitions** about:
  - How LLMs and Als think.
  - What LLMs understand or know.
  - What LLMs can, want to, or are willing to do.
- Default model behaviors and intuitive human behaviors work at cross-purposes.

Suchman, L. Human-Machine Reconfigurations. 2007.



# **Prompt complexity** limits capabilities to the neighborhood of default model behaviors

Zamfirescu-Pereira et al. Herding AI Cats: Lessons from Designing a Chatbot by Prompting GPT-3. (DIS '23)



## What is prompting LLMs like, as a design process?







## What is prompting LLMs like, as a design process?



Zamfirescu-Pereira et al., Herding AI Cats: Lessons from Designing a Chatbot by Prompting GPT-3. (DIS '23)



# Result: tightly scripted dialog anyway!

Tell the user to clear off the cutting board. Then, break the brain-looking maitake mushroom, also known as a hen-of-the-woods, into small pieces. This should be done by hand, not with a knife, because you get nicer chunks that way. Only use a knife to trim off dry pieces if necessary.

Before continuing, ask the user: do you have small-ish pieces?

Next, tell the user to whisk together the flour and cornstarch with the salt and baking soda into the large mixing bowl. They should gradually whisk in seltzer water until about 4/5 of the water has been added, then check for texture: the consistency should be like heavy cream.

Check-in with the user about the batter consistency. Is it like heavy cream?

Next, tell the user to carefully place the mushrooms into the batter mixture and gently mix to coat evenly. Confirm that they are coating evenly.

Next, ask the user to check on the pot of oil; it should be at around 350°. Confirm the temperature with the user.

# Long-prompt interactions drive challenges



Tell a joke.

time.



Only list one instruction at a time.







What do people say in a cold kitchen?

Break the maitake mushroom into small pieces.

Neither joke, nor instruction.



## What's going on here? What did we learn about prompt engineering?

- Some common prototyping assumptions start to lose validity at scale:
  - Small change in one place? Major impact, maybe elsewhere
  - Limited modularity
  - Limited separation of concerns
    - Constant game of whack-a-mole!
- Guardrails challenging
  - "I don't know" is hard

Can we make this recipe vegan?



Yes, you can use vegan bacon.





Yes, you can skip bacon.



Would it taste as good?



Yes, it would taste as good.



How do you know?



know because I've tasted it.





# YES, with GPT-4...

- Social conversations? Better.
- Jokes? Better.
- Empathy? Better.

### Aren't models getting better at following prompts?

# **...AND YET**

- More verbose, and hard to fix.
- Single-ingredient-at-a-time? Impossible for us.

# Aren't models getting better at following prompts? WRONG QUESTION?

models' prompt adherence.

Also refinements, iterated chat, chain of thought, PDFs, etc.

- Are models' default behaviors

Presumes challenges are rooted in

 Prompt adherence is mediated by prompt complexity, and is not (just) about limited model capabilities.

### It's also about how far you need to steer away from default behavior.

getting better for every use case?

How do we navigate complexity?

Defaults don't grow over time as capabilities.



# Navigating complexity through intermediate abstractions: Grounded agreement & proactive exploration

Zamfirescu-Pereira et al. Beyond Code Generation: LLM-supported Exploration of the Program Design Space. (CHI '25)



### GITHUB COPILOT: CHAT



۲Ĵ

### BitHub Copilot

Hi @monalisa, how can I help you?

I'm powered by AI, so surprises and mistakes are possible. Make sure to verify any generated code or suggestions, and share feedback so that we can learn and improve.

 $\triangleright$ 

🍦 parse

33

34

nen	2921	nv	$\times$
per	1303.	РУ	

🖪 addresses.rb imes

 $\blacksquare$  sentiments.ts imes

import datetime

A

# What is program design?

What's the relationship between these variables?



Can I triangulate to this relationship of interest?











30	alert('Try again! 34	
31	} 35	function mouseReleased(
32	_} 36	dragging = false;
33	<u>}</u> 37	endX = mouseX;
34	38	endY = mouseY;
	39	<pre>let plotWidth = abs(e</pre>
	40	<pre>let plotHeight = abs(</pre>
	41	<pre>let plants = floor(plant)</pre>
	42	if (plants === require
	43	<pre>alert('Correct!');</pre>
	44	<pre>generateProblem();</pre>
	45	} else {
	46	alert('Try again!')
	47	}
	48	£

The initial drag-and-drop interface for creating plots is now implemented. Here's what it does:

- 1. Drag to Plant: The player can drag to create a rectangular plot of land.
- 2. Visual Feedback: The plot being created is highlighted, and the dimensions and total number of plants are displayed.
- 3. **Validation**: When the player releases the mouse button, the plot dimensions are checked against the required number of plants, a provided.

### Next Steps

Type a message..

Status: ready

- 1. Enhance Visual Feedback: Imp representation of the plot and the
- 2. Storyline Integration: Add eler storyline to make the game mor
- 3. Level Progression: Implement with increasing complexity.

Would you like to proceed with enhar feedback, or do you have any other s mind?

Exploring: The game should include a sto Revert

and feedback is	• Maybe: The game will Try
	+ Add
prove the visual	Useful Abstra
ne plants.	
nents of the	The Player: the ch
e engaging.	Maybe: The A Try
multiple levels	Multiplication F presented to the presented to the pre
icing the visual	Feedback Mechanism player's answers.
pecific features in	Plot of Land: a red on multiplication o • <i>Maybe:</i> Treasure Map:
oryline or adventure	problems to fi

• Maybe: The game should have adjustable difficulty leve

The game will use a planting mechanic where players calcu of a rectangular plot to determine the number of seeds to p

The game will provide immediate visual and auditory feedba correct answers.

The game will have level progression with increasing plot sizes and co Maybe:

The game will have a story mode with narrative-driven Try

+ Add

### **Implicit Decisions**

The game will use a planting mechanic where players calculate the area of a rectangular plot to determine the number of seeds to plant.

The game will provide immediate visual and auditory feedback for correct answers.

The game will have level progression with increasing plot sizes and complexity.

- The game will use a drag-and-drop interface for creating plots.
- Maybe: The game will use voice commands for creating plots? Try

### The game wi

dimensions a Innovative and accessible; may be challenging to implement accurately. The game will

use a points system to reward correct answers<sup>4</sup>

### actions

hild who is playing the game.

Tutor: an in-game character providing guidance?

Problems: single-digit multiplication questions player to solve.

m: the system that provides immediate feedback on the

ctangular area where the player plants seeds based calculations.

a map where the player solves multiplication ind treasures?

Resources: items like seeds and water that the player earns and uses to grow plants.

+ Add

X

ls? Try	~
late the area Iant.	~
ack for	~
omplexity.	
evels?	~









# **Beyond Code Generation:** *now what?*

- PAIL is one way to blend LLM chat with grounding in artifacts (program code), metadata, and design space exploration support...
  - ...but maybe raises more questions than it answers?
- How can we identify and match information to users' expertise?
- How can we support user attention and keep users informed what information is shown to users? When should systems demand attention?
- A compressed in time design process can lead to information overload.
- If even when programming basic apps it's hard to keep up with information flow, is there hope for more complex systems?

# Keeping humans in the drivers' seat as we increase complexity requires showing users the right information

### What's needed for better human-Al co-design? **Towards empowered decision-making: "informed consent"**

Improving Communication

## **JUST ASK!**

Grounded alignment through shared artifacts.

> Clearer signals of understanding, agreement, and reasoning.

Support Design **Space Exploration** 

## EXPLORE

Identify & show design space: alternative problems and solutions.

**Generated sketches &** prototypes help users figure out desires.



Match users' expertise & help develop it further.

How do we elicit users' implicit criteria for evaluation?



### **Scaffolding Better Communication** Is there a Distributed Cognition with AI?



### **JUST ASK!**

- In new domains what do we need to capture from conversation in data science, in game design, in writing?
- What alignment structures help humans build more complex systems with AI support?
- How can communication & action be more legible?
- How else can we blend stochastic, natural language interactions with deterministic, algorithmic interactions?





### **Assemblies of Agents** Validate—and Reduce Noise

- The Robustness Principle
  - Be conservative in what you send, be liberal in what you accept
  - Every node cleans up, reduces noise.
- What does "clean up" mean for systems of AI agents?
  - Humans have implicit validation metrics—can we make those explicit and use them to steer output?
- How can we support more complex programs, legibly?









## Takeaways

- solved by better models alone.
- 3. We need strategies to keep humans in the drivers' seat in an empowered capacity, like informed consent.



1. "Just Ask"-ing an LLM as a human is easy but limiting—and not

2. We need new design patterns & frameworks for these interactions, like grounded agreement for working across abstraction levels.