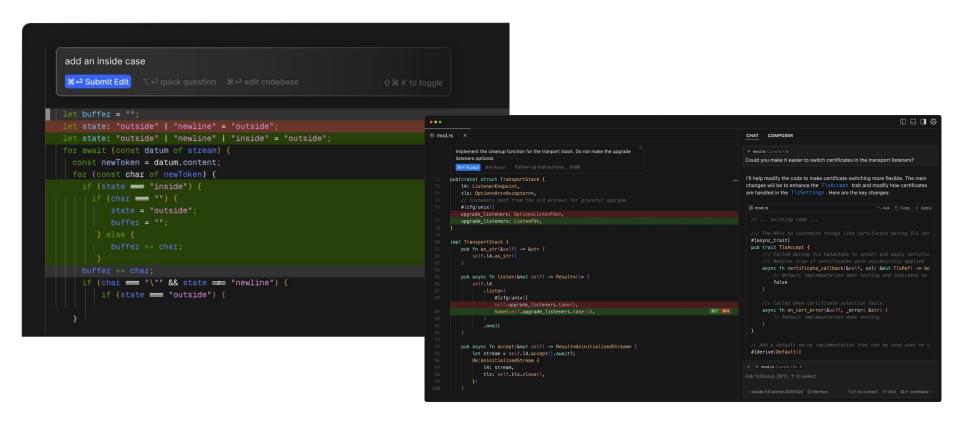


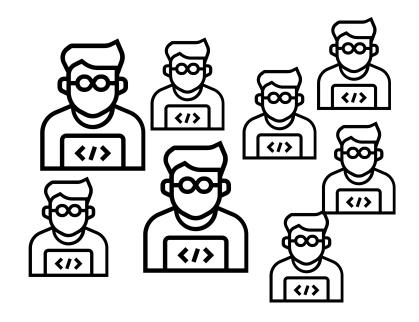
Al Code generation has transformed the way programmers work.



# Design (What do we build?)



Engineering (How do we build it?)

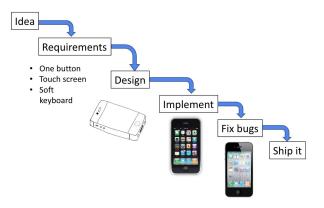


Created by Brickclay from Noun Project

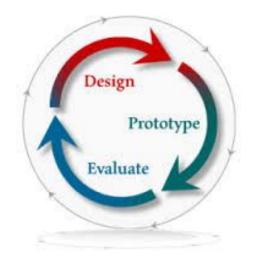
Created by WiStudio from Noun Project

#### Two Software Design Processes

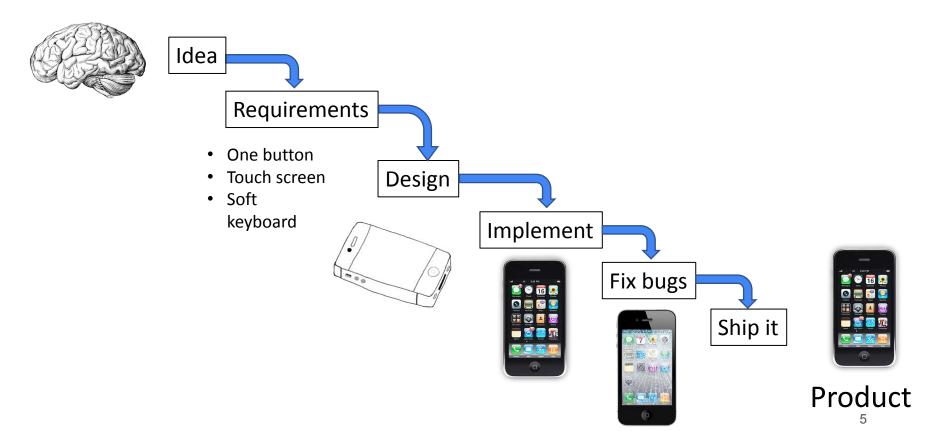
#### Feedforward



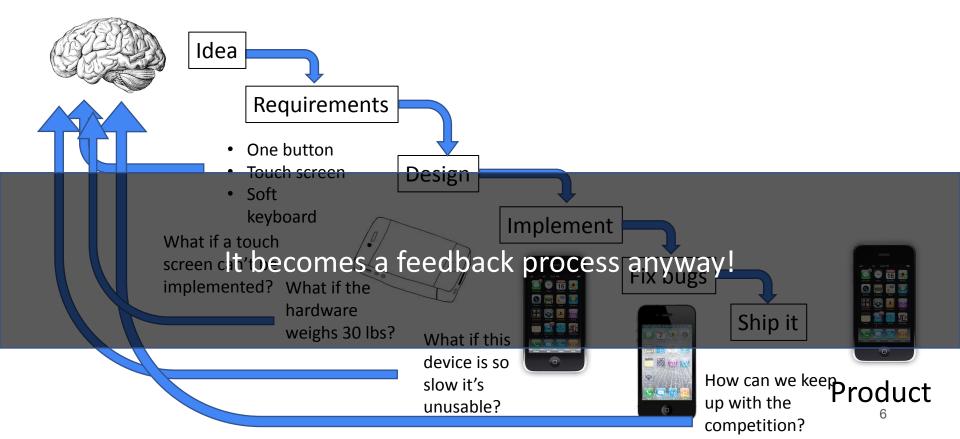
#### Feedback



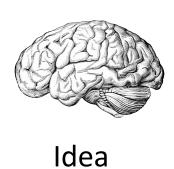
#### The Waterfall Model

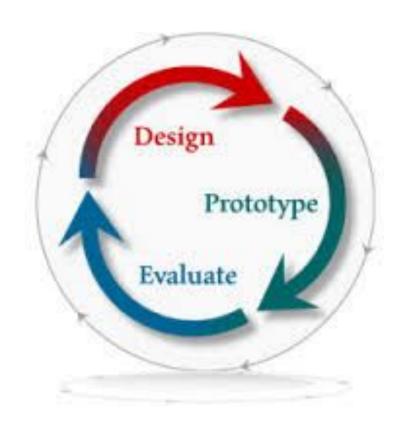


#### The Waterfall Model: What could go wrong?



## Iterative Design is a feedback loop.







#### Each feature needs to be iterated before moving on

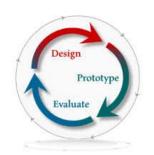


Idea



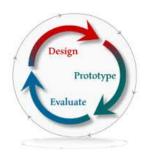
Touch screen





Soft keyboard





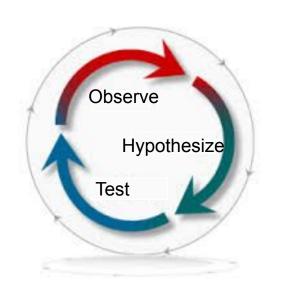
One button



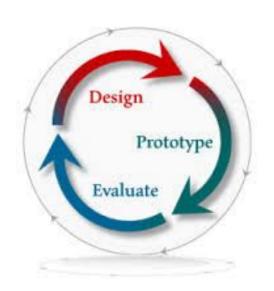


#### Human-Centered Design: Two feedback loops

Person with a challenge



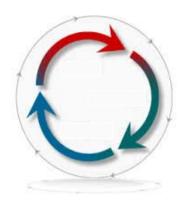






**Understand the problem** 

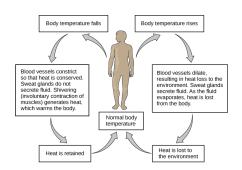
Solve the problem



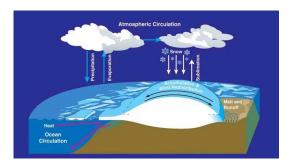
Al Code generation isn't enough to create better software products...

We have to close feedback loops between design and engineering.

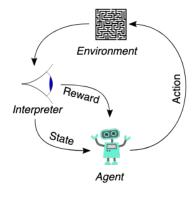
#### Feedback loops are essential to all systems



The human body



Earth systems



Training AI

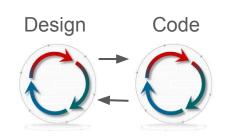


Education



Childhood Development

#### Bridging Design and Code with Gen Al





Iterative Problem Specification



**Iterative Development** 









#### LogoMotion

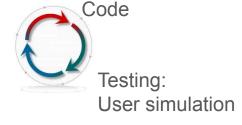
- Self-debugging loop
- Iteration with AI Editing widgets



#### **Double Agents**

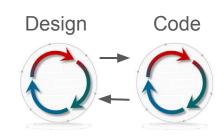
- Testing with user simulation
- Iterating on designs with AI suggestions







#### **Bridging Design and Code with Gen Al**





Iterative Problem Specification



Iterative Development







- Self-debugging loop
- Iteration with AI Editing widgets





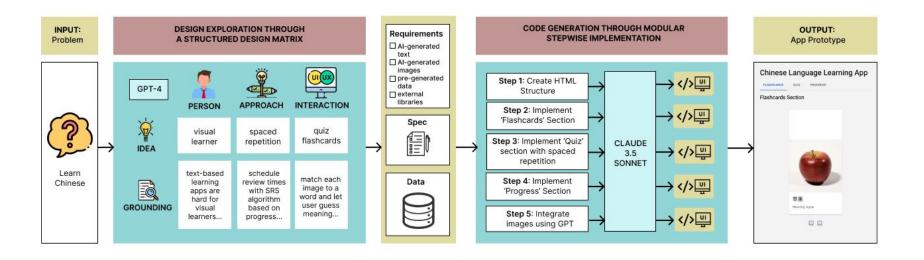




#### **Double Agents**

- Testing with user simulation
- Iterating on designs with Al suggestions





#### DynEx: Agentic Assistance to Bridge Design and Code

Jenny Ma, Karthik Sreedhar, Vivian Liu, Pedro Alejandro Perez, Sitong Wang, Riya Sahni, Lydia B. Chilton CHI 2025

#### Everyone wants Agents to Develop Software

Published as a conference paper at ICLR 2024

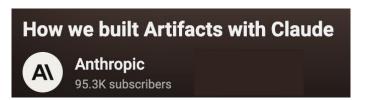
#### SWE-BENCH: CAN LANGUAGE MODELS RESOLVE REAL-WORLD GITHUB ISSUES?

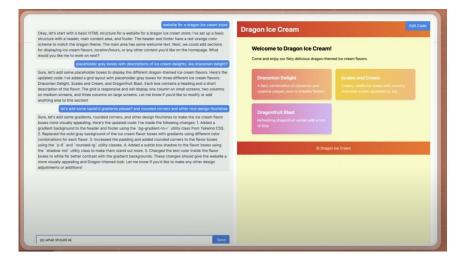
Carlos E. Jimenez\* 1,2 John Yang\* 1,2 Alexander Wettig<sup>1,2</sup>
Shunyu Yao<sup>1,2</sup> Kexin Pei<sup>3</sup> Ofir Press<sup>1,2</sup> Karthik Narasimhan<sup>1,2</sup>

<sup>1</sup>Princeton University <sup>2</sup>Princeton Language and Intelligence <sup>3</sup>University of Chicago



Figure 1: SWE-bench sources task instances from real-world Python repositories by connecting GitHub issues to merged pull request solutions that resolve related tests. Provided with the issue text and a codebase snapshot, models generate a patch that is evaluated against real tests.

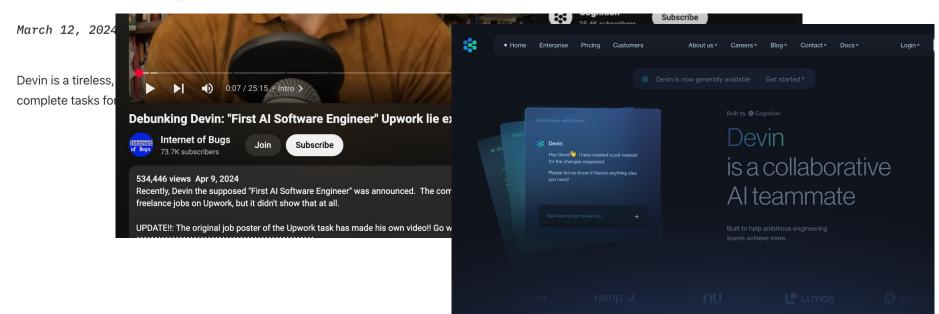




#### Can Agents do it by themselves? NO!

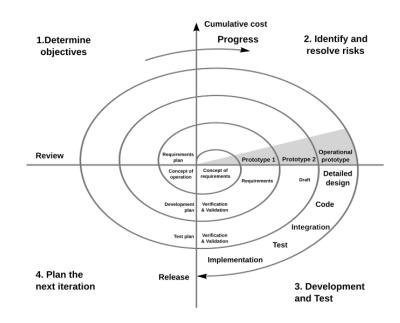
/blog

## Introducing Devin, the first AI software engineer



#### Why do we need human interaction?

- Specify the problem and the solution
- Iteratively test modular pieces (or prototypes) so errors don't compound.
  - Every problem has unknown aspects of the environment where it will operate.
  - We have to expect to find new problems.



We need to bridge design and implementation

#### There is a gap in going from an idea to a working solution

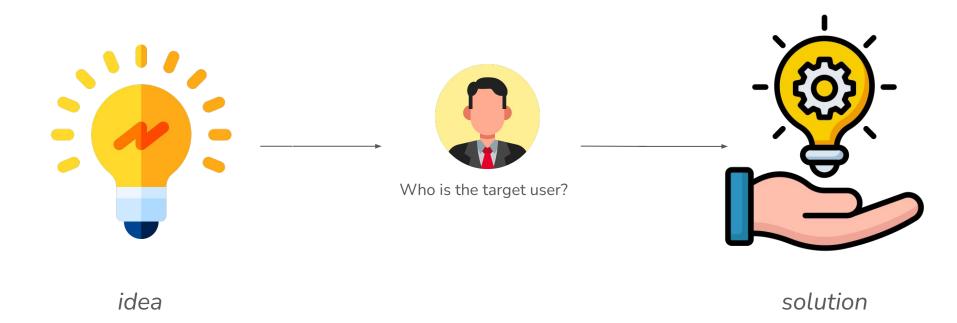




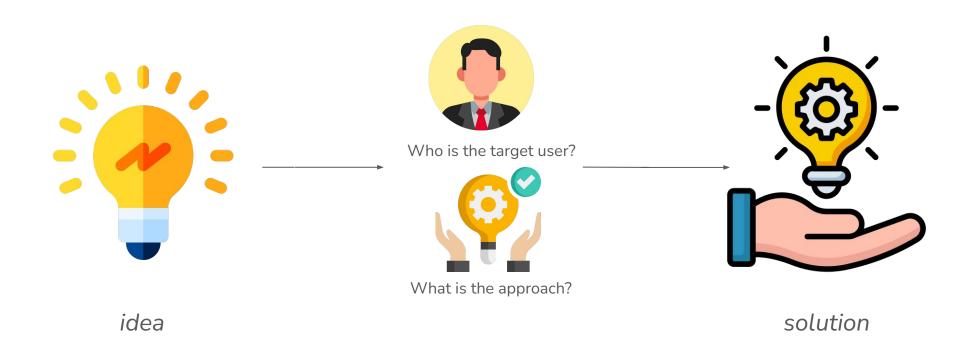
idea

solution

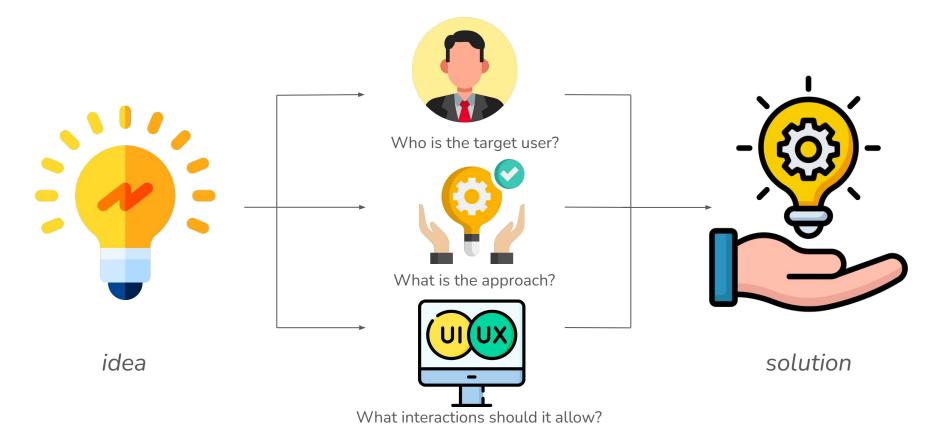
#### We need to specific the problem! Who is the user?



#### We need to specific the problem! What is the approach?



#### We need to specific the problem! What is the interaction paradigm?

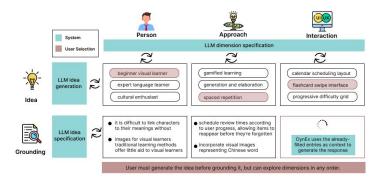


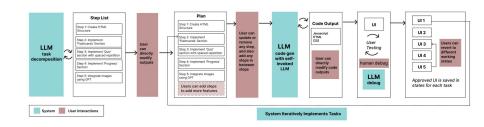




## **Design Exploration**

Implementation



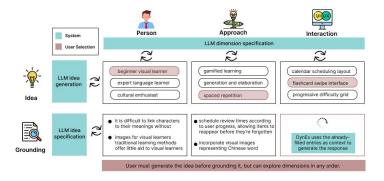


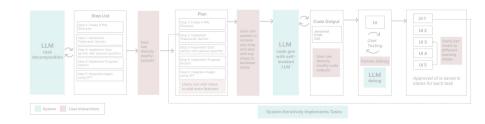


## **Design Exploration**



Implementation





Design Matrix: Help users specify the user, approach, and interaction paradigm.

Person

Approach



Interaction



Levels of specificity



Idea



Person:Idea

Approach:Idea

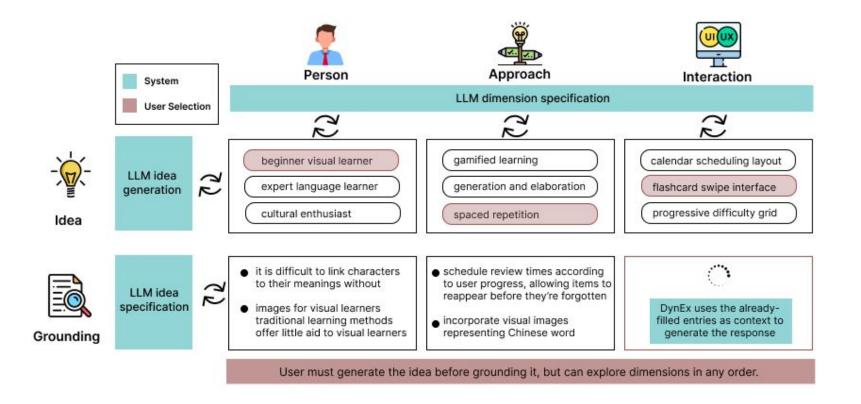
Interaction:Idea

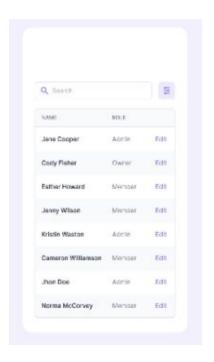
Person:Grounding

Approach: Grounding

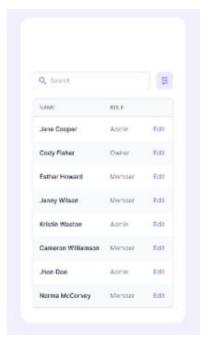
Interaction: Grounding

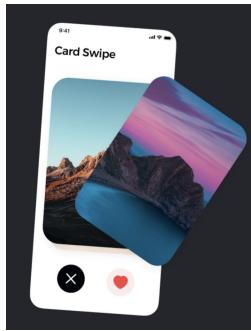
# The Design Matrix uses Gen AI to suggest ideas and grounding with respect to the other ideas/groundings





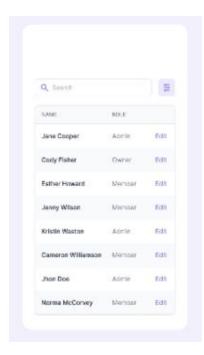
Table

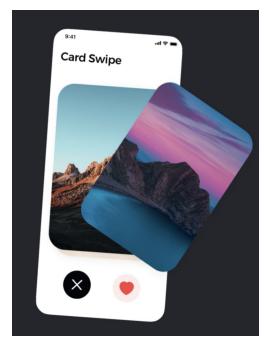




Table

Card Swipe





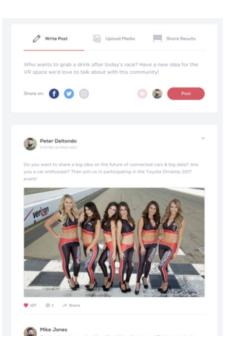
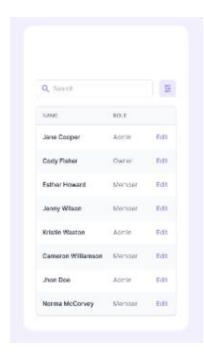
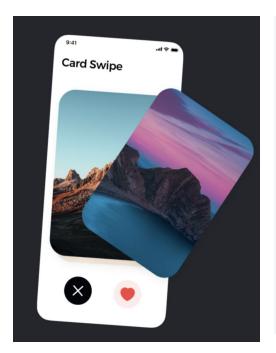
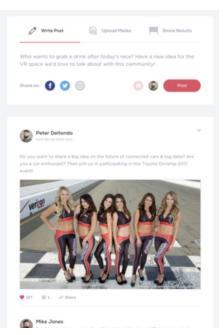


Table Card Swipe News Feed







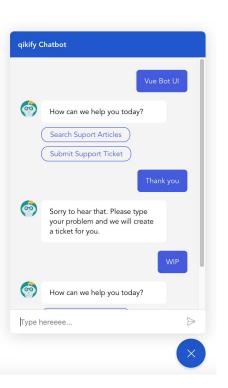


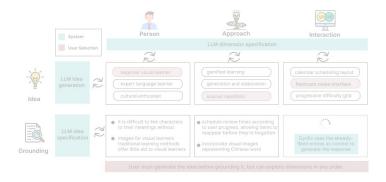
Table Card Swipe News Feed Chatbot

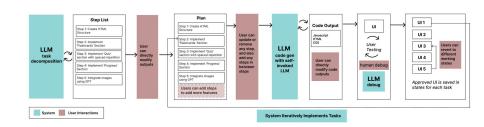


## **Design Exploration**



#### **Implementation**





#### Design Requirement Agent



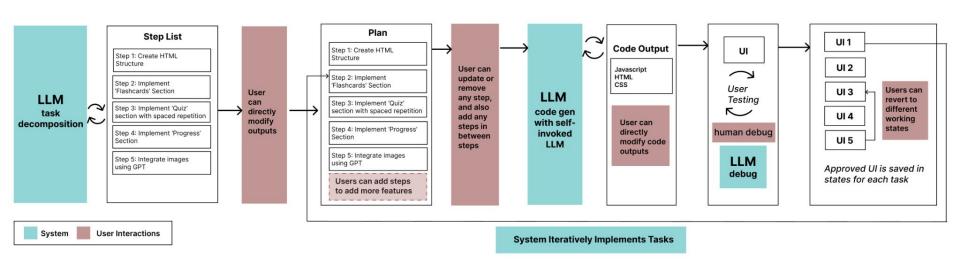




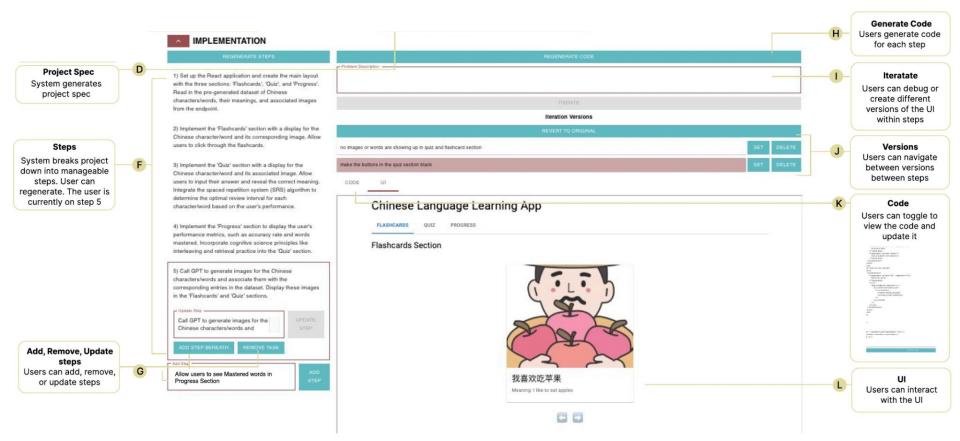
#### PROJECT REQUIREMENTS

#### Spec GPT GPT Application Layout: **Data Input Suggestions** - Divide the interface into three main sections: "Flashcards," "Quiz," and - Data Input "Progress." - The "Flashcards" section displays the Chinese character/word along with its corresponding image for visual association. ChartJS "id": 1, - The "Quiz" section presents the character/word and image, prompting "chinese": "我学习中文", the user to input the meaning. ☐ GoJS "meaning": "I study Chinese", - The "Progress" section shows the user's performance metrics, such as "imagePath": "images/studying.jpg" accuracy rate and words mastered. User Interactions: "id": 2, - In the "Flashcards" section, users can click through the flashcards to "chinese": "你好吗?", study the character/word and its associated image.

#### **Dynamic Iterative Implementation Agents**

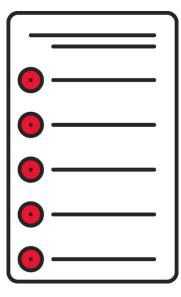


#### Iterative Development with Step-by-Step Implementation



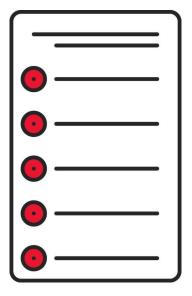
#### Modular Step-by-step Implementation

System breaks down implementation plan into **steps** and generates code for each step



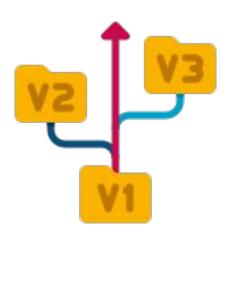
#### Modular Step-by-step Implementation

System breaks down implementation plan into **steps** and generates code for each step

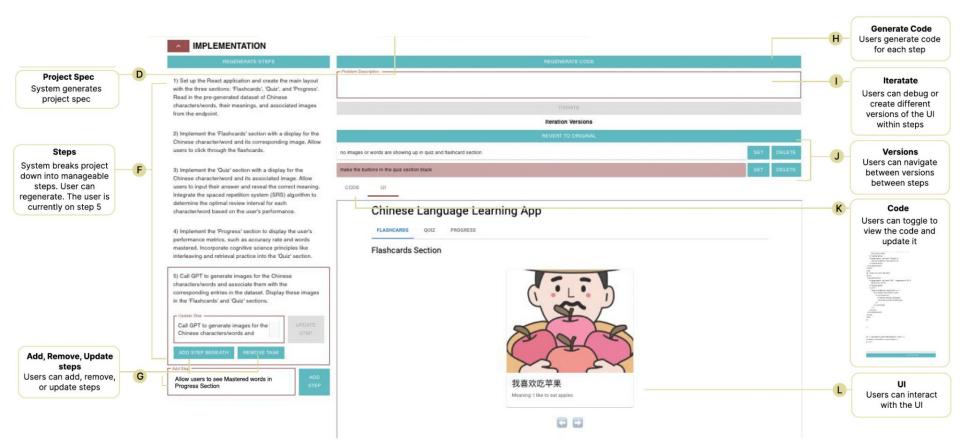


Users can **add**, **remove**, and **update** steps.

Steps provide a natural form of **version control**.



#### Self-Invoking LLMs

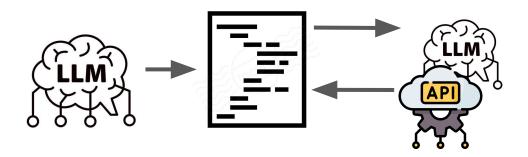


#### Self-Invoking Multi-Modal LLMs

Generative AI APIs are powerful and versatile, and excel at handling a wide variety of tasks.

Our generated code can call GPT to generate dynamic data and create images.

We can create more realistic applications that can more accurately mimic a user's experience





### Examples

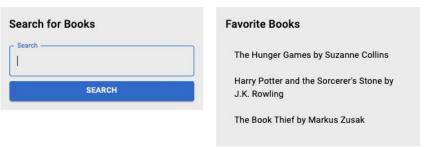
#### Movie Recommendation App

Preferenc	es		
comma-separates	0		
ovies (comme-sep	700049		
oia	narates)		
BMIT PREFERE	NCES		
nmendat	ions		
COMMEND FRO	OM WATCHLIST		
~	Moana		
	Animation An adventurous teenager sails ou quest to become a master way-fir  + WATCHLIST	t on a daring mission to save her people. Moana meets the once-mighty demigod Mau,i who guides her in her order.	
DY NWID	Finding Dory  Administra  With help from Nemo and Marin, Dury the forgetful fish embasks on a quest to reunits with her mother and father:  + WATCHLIST		
ĖTĖ	The Secret Life of Pets Administration The qualer file of a terrier named Max is upended when his owner takes in Duke, a stray whom Max instantly dislikes.  + WATCHUST		
	Kung Fu Panda Admission The Disgon Warrior has to clash against the savage Tai Lung as China's fate hangs in the balance. However, the Dragon Warrior matrile is supposedly matrialize to be bestowed upon an obese pands who is a novice in martial arts.  + WEICHLIST		
Yes Mice	Inside Out Annual Annua		
h			
movies			
on			
nlist		PRETESCARI	
na tion		The Secret Life of Pets Animation	
W DETAILS -	REMOVE	Q VIEW DETAILS - REMOVE	

Movie recommendations (and images) are made by calling GPT within the app to generate the related movies dynamically

#### **Book Recommendation App**

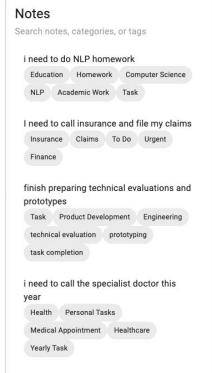


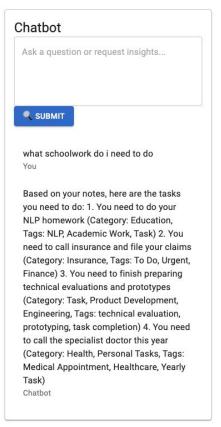


Book recommendations, images, and summaries are made by calling GPT within the app dynamically

#### Note-Taking and Sensemaking Application







Clustering of notes is done dynamically with GPT

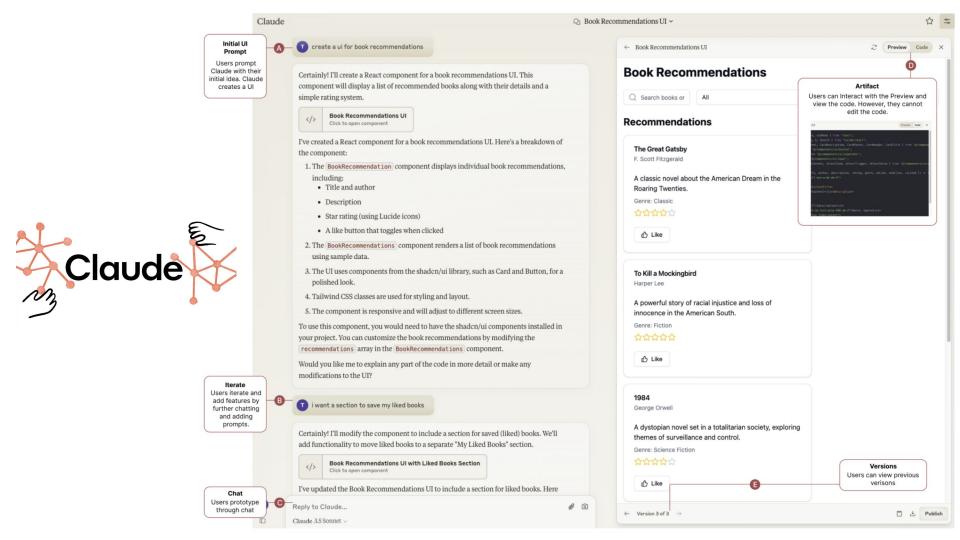
#### **User Evaluation Research Questions**

**RQ1:** [Divergence] To what extent does DynEx enable divergent exploration within a problem space?

**RQ2:** [Convergence] To what extent does DynEx allow users to better develop their ideas?

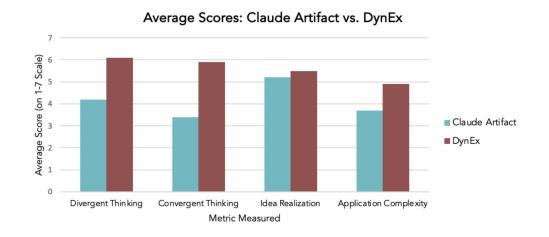
**RQ3:** [Implementation] To what extent does DynEx enable the code to realize a complex idea?

**RQ4:** [Overall] To what extent does DynEx allow for a better prototyping experience?



#### Results

- DynEx inspired new solutions, allowed users to explore a problem space, and further developed users ideas (p=0.05 level).
- DynEx enables users to create more complex, feature-rich, and intuitive applications that emulate a true user experience (p=0.05 level)



#### P6 - Friend-Activity-Joining Application

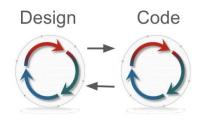
"[DynEx] suggested many features... such as how to match people with friends, how to visualize the friend-sharing [component], how to [join a friend's] friend experiences, how to share experiences with friends.. it created a pretty robust social ecosystem surrounding the calendar experience."

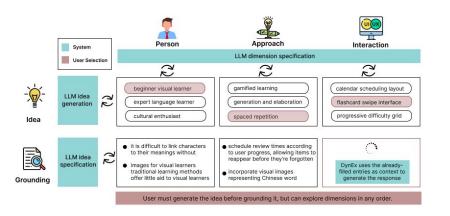
"[DynEx] incorporated really different, very distinct UI features that weren't very connected to each other at [face value] really well"

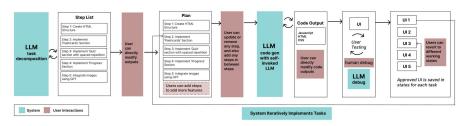
#### P8 - Concert Ticket Aggregator

"It was **intuitive**. It was **feature-rich**. It had all the important features, like **sorting by columns** and **sorting by genre**, **value for money**, etc... I liked the **conciseness** of the information... contrary to [Claude Artifact's] prototype which was... not-fit for this use case"

#### Dynex: Bridging Design and Code



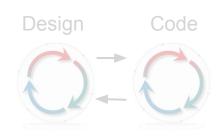




Iterative Problem Specification (Design)

Modular Code Generation and Testing

#### **Bridging Design and Code with Gen Al**





Iterative Problem Specification



Iterative Development







- Self-debugging loop
- Iteration with AI Editing widgets







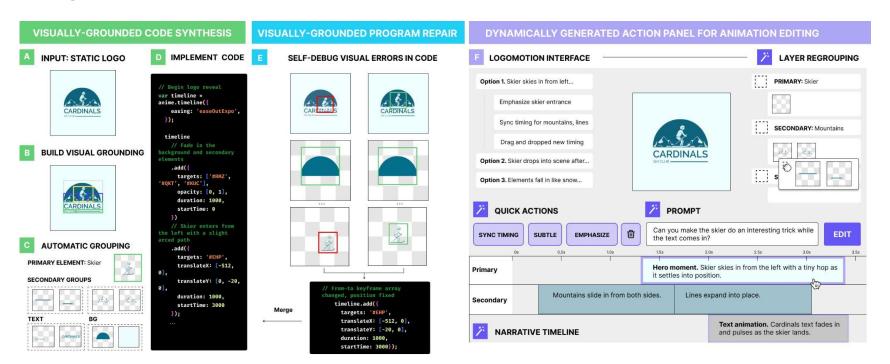


#### **Double Agents**

- Testing with user simulation
- Iterating on designs with Al suggestions



#### LogoMotion: Visually Grounded Code Generation and Repair

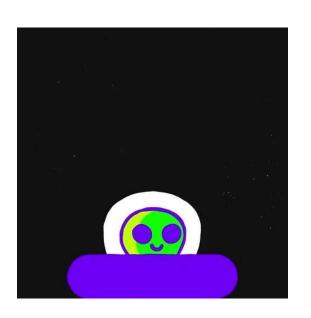




#### Creating semantically meaningful motion is hard.





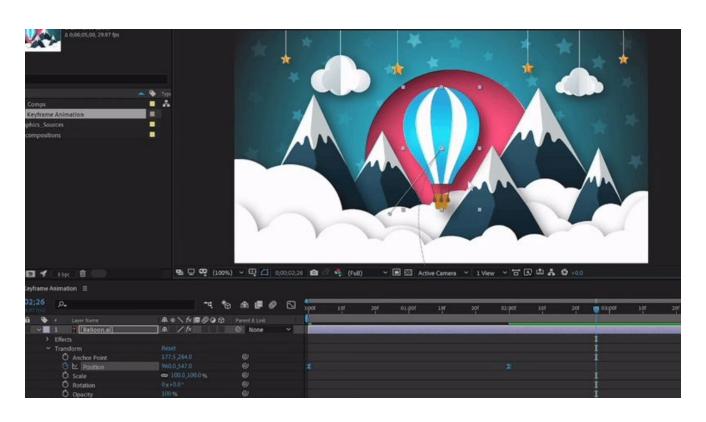


Cars should drive

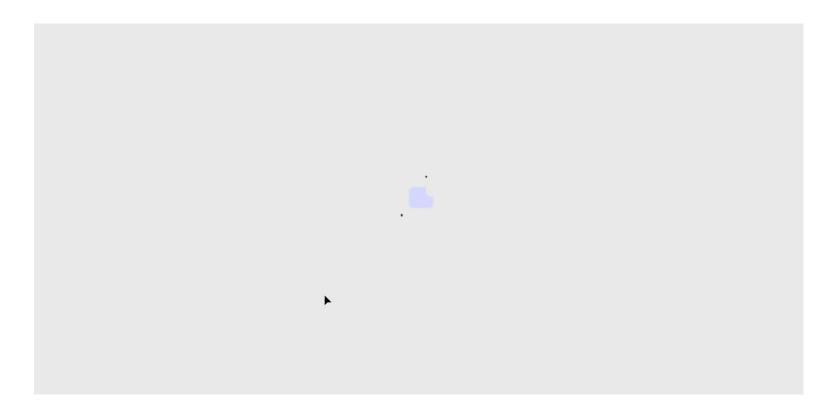
Fans should spin

UFOs should take off

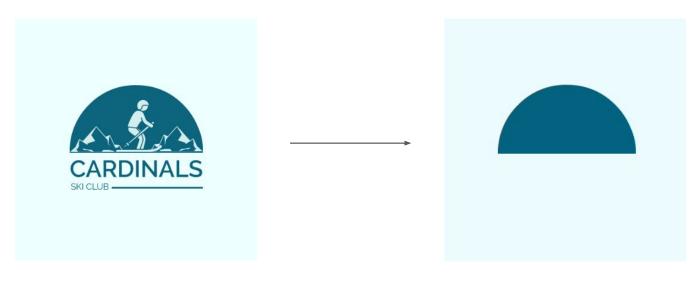
#### **Animation Uls are complex**



#### Templates are too rigid



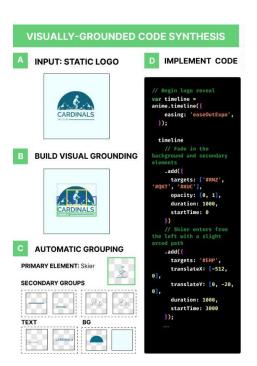
## Human-agent interaction for animation gives users control and freedom



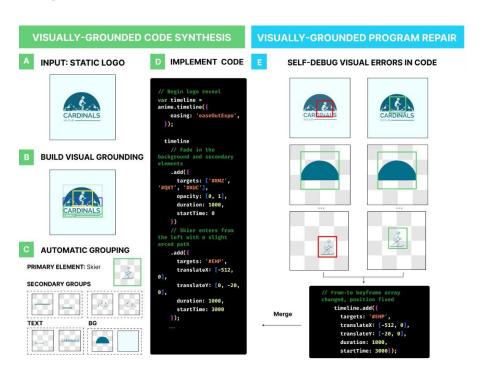
INPUT PDF

OUTPUT
HTML PAGE + ANIMATION CODE

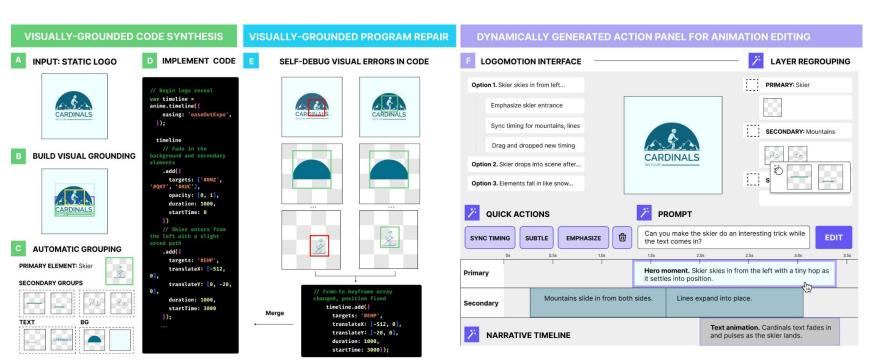
#### Agents can automatically author an animation



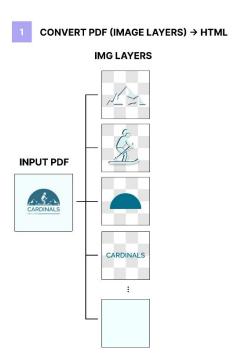
#### Agents can automatically debug and repair animations



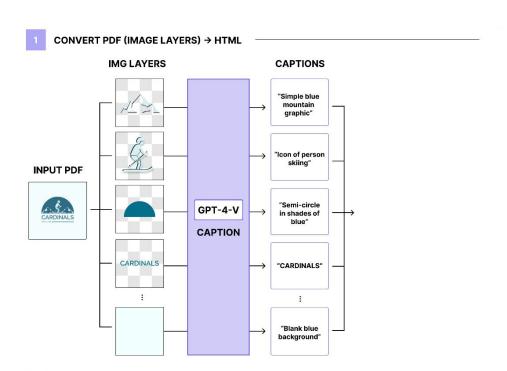
#### Agents Generate a UI for Editing the Animation



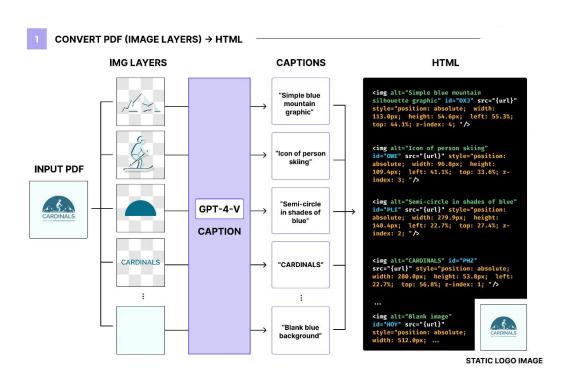
### Agents visual understand the input.



#### Agents label the parts of the input

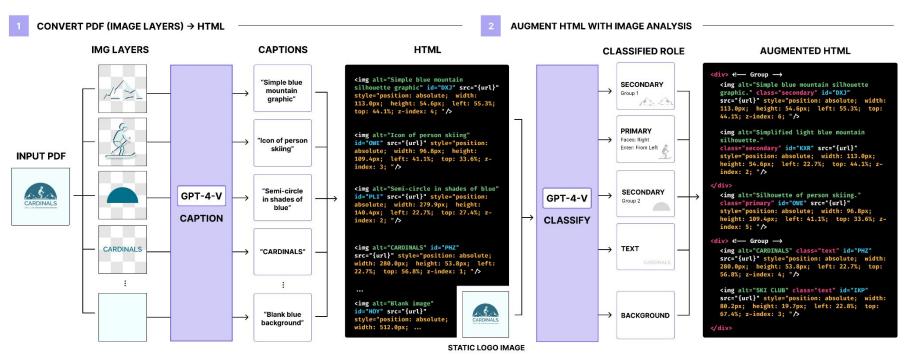


#### Agents represent the input as HTML



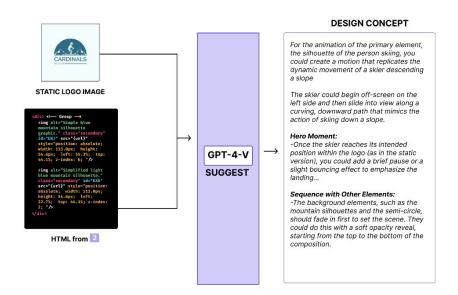
#### **Augment the HTML with visual analysis**

(Primary, secondary visuals, text, background, etc)

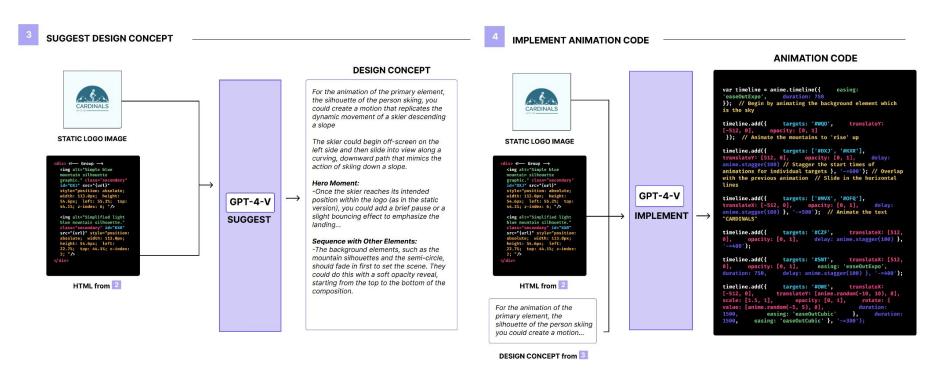


### Agents suggest a design concept.

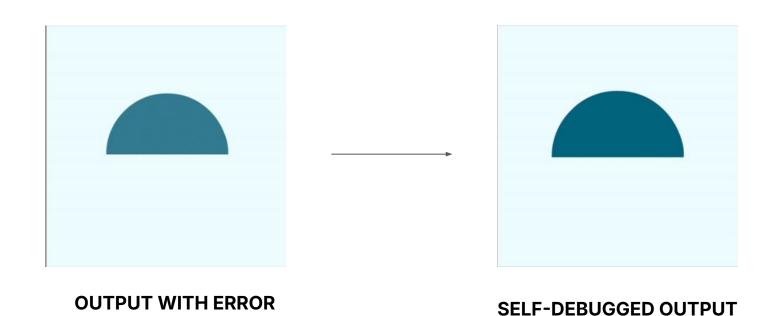
3 SUGGEST DESIGN CONCEPT



#### Agents implements animation code.

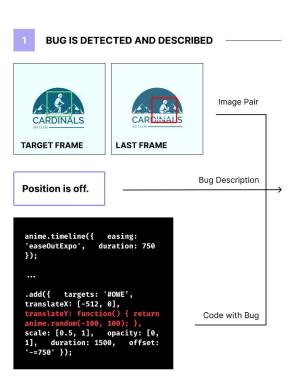


#### Agents can also detect a fix animation errors

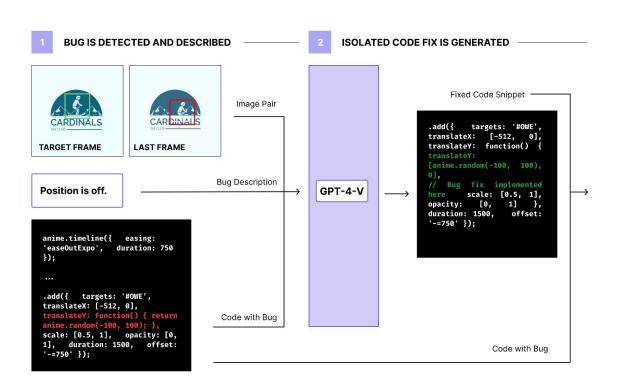


63

#### Self-debugging agents check for errors.



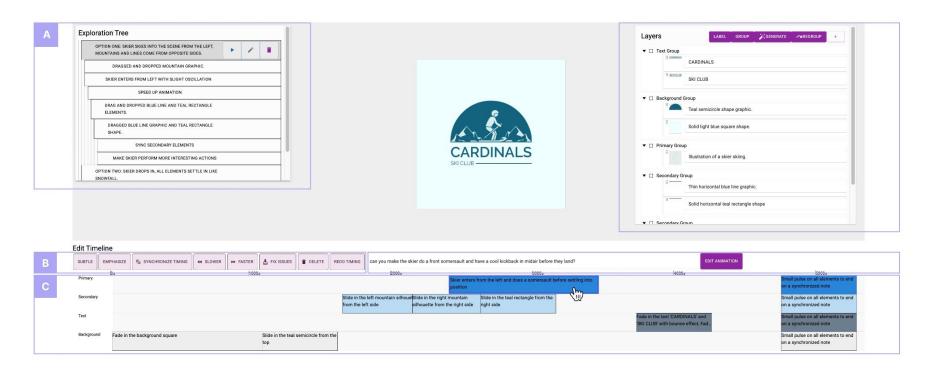
#### Self-debugging agents check for errors.



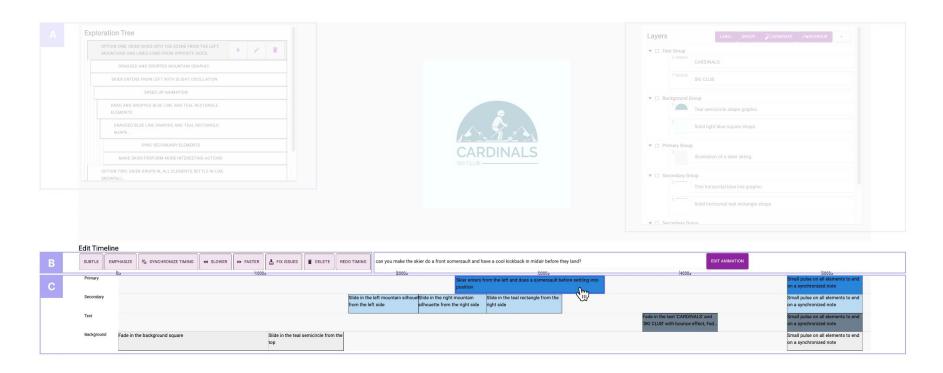
#### Self-debugging agents check for errors.

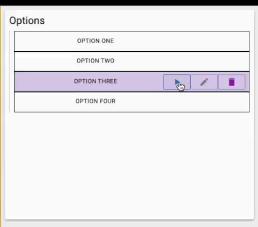


## Agents activate an editing UI for users to improve animations or explore new ones.

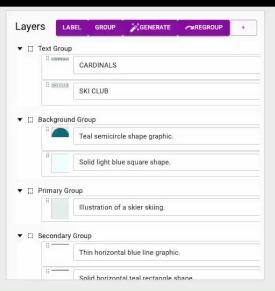


## A timeline widget can reorder animation blocks and adjust timing



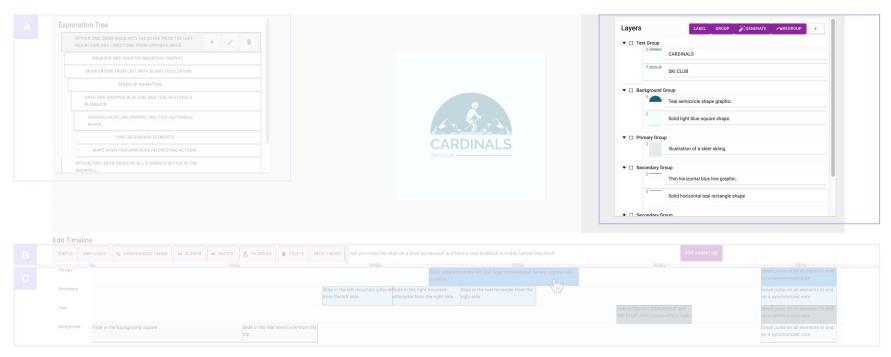


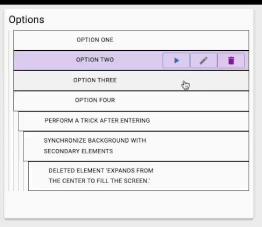




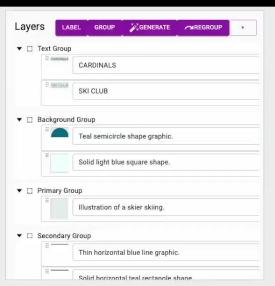


# A grouping widget can synchronize motion and timing of elements



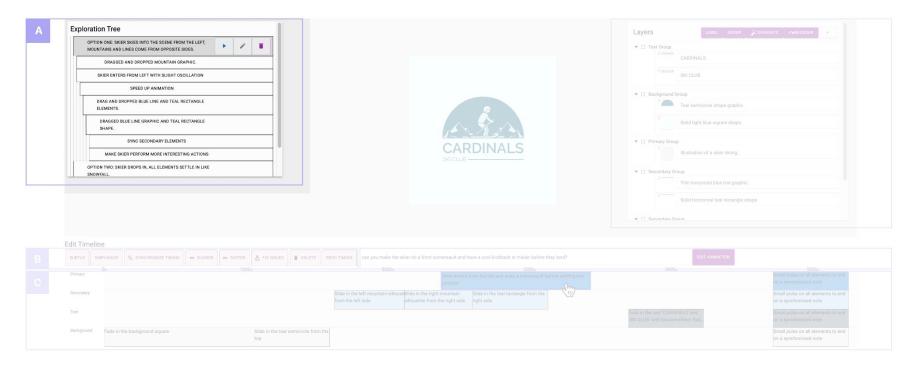








# A version tree allows users to review and select animations.

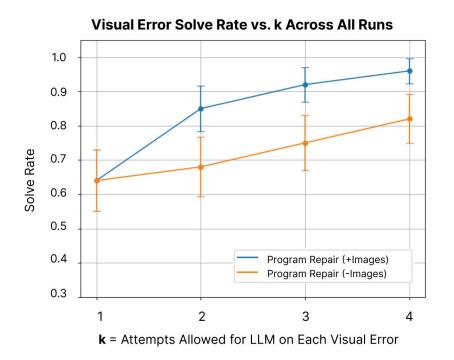


# Animation Agents produce semantically meaningful animations.

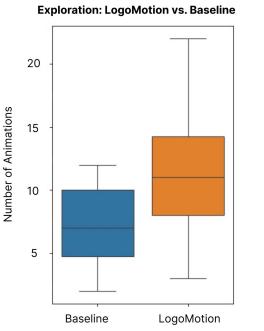


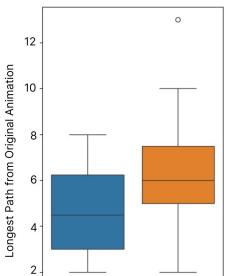


# **Program Repair Enables a 95% Solve Rate.**



# Editing Agents helped users explore more animations and craft better iterations than a baseline.





Baseline

LogoMotion

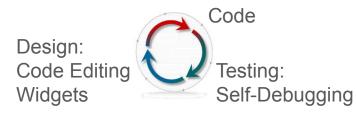
Iteration: LogoMotion vs. Baseline

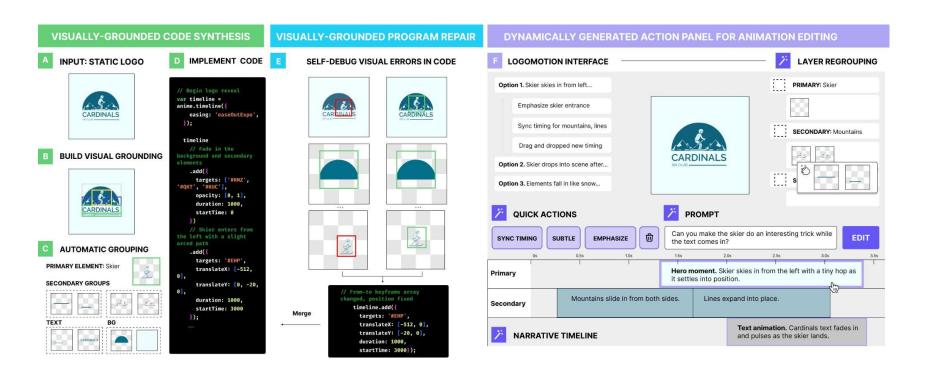


"It's really cool especially for someone who couldn't do it without this tool and wouldn't spend a lot of time with Youtube videos or tutorials to do the very basics." -P2

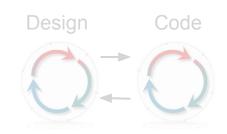
# LogoMotion induces feedback with:

- 1. Self Debugging
- 2. Code Editing Widgets





# **Bridging Design and Code with Gen Al**









- Iterative Problem Specification
- Iterative Development





# LogoMotion

- Self-debugging loop
- Iteration with AI Editing widgets



# **Double Agents**

- Testing with user simulation
- Iterating on designs with Al suggestions



Design: Policy updating

Code Editing

Widgets



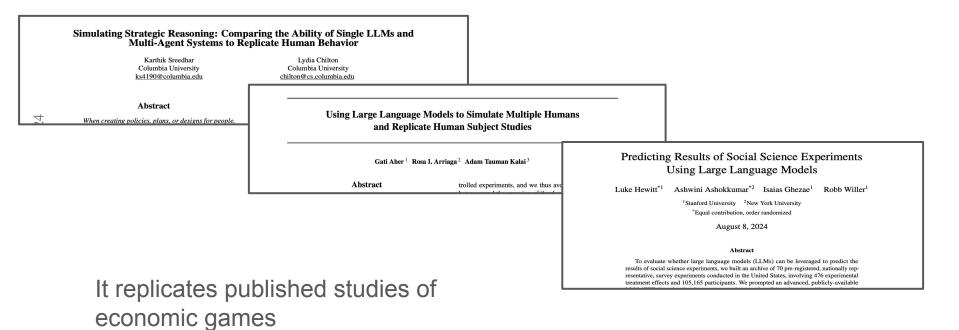
# DoubleAgents

Al Agents for Task Completion and Human Simulation

Karthik Sreedhar, Alice Cai, Jenny Ma, Jeffrey V. Nickerson, Lydia B. Chilton. Simulating Cooperative Prosocial Behavior with Multi-Agent LLMs: Evidence and Mechanisms for Al Agents to Inform Policy Decisions. IUI 2025.

Tao Long, Xuanming Zhang, Sitong Wang, Jenny Ma, Karthik Sreedhar, Lydia B. Chilton Double Agents. Al Agents for Task Completion and Human Simulation. In preparation.

## Several Studies have shown that AI can accurately replicate human behavior



And 80% of of 70+ published and *unpublished* psychology games

Agents who "Observe", "Think" and Act can engage in "realistic" human activity, like organizing a party.



complex behavior outside of the lab

We use AI to accurately simulate

There is one professor and 3 students:





The professor announces one of three late policies to the class ahead of a simulation:

(1) harsh

2) some leniency

(3) very lenient

There is one professor and 3 students:





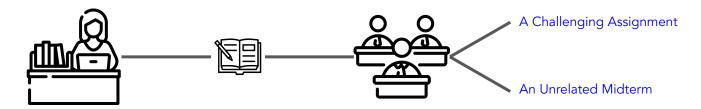
Students are assigned one of three "personality" traits.

There is one professor and 3 students:



Professors assign weekly homework, students work on it and turn it in.

There is one professor and 3 students:



We also add stressors to the simulation to see what affects it has on behavior.

There is one professor and 3 students:







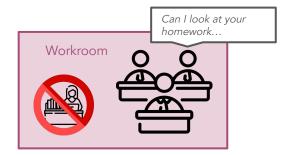
Students can enter the office to speak to or ask the professor questions.

# When we add stressors, we see behaviors consistent with real students

There is one professor and 3 students:







- 1) Cheating. Procrastinators ask overachievers if they can "look at their homework"
  - 2) Students email the professor pleading for extensions.

# Since AI can simulate human behavior inside and outside the lab, it can provide human feedback to test software systems!



### Simulating Cooperative Prosocial Behavior with Multi-Agent LLMs: Evidence and Mechanisms for AI Agents to Inform Policy Decisions Karthik Sreedhar Alice Cai Jenny Ma jm5676@columbia.edu ks4190@columbia.edu acai@college.harvard.edu Columbia University Harvard University Columbia University New York, New York, USA Cambridge, Massachusetts, USA New York, New York, USA Jeffrey V. Nickerson Lydia B. Chilton inickers@stevens.edu chilton@cs.columbia.edu Stevens Institute of Technology Columbia University Hoboken, New Jersey, USA New York, New York, USA .1 . . There is one professor and 3 students: Can I look at you homework Professor's Office Cheating. Procrastinators ask overachievers if they can "look at their homework" 2) Students email the professor pleading for extensions.

# DoubleAgents

Al Agents for Task Completion and Human Simulation

Lydia Chilton
Columbia Al Summit. Al in Business

# **Human Coordination Problems (HCPs).**

Warehouse.

Resource discovery.

Constraints. Availability.

Priority-based allocation.

Dynamic scheduling.

Optimization. Bipartite matching.

Real-time adjustments.

Uncertainty & adaptability



# LLMs Agents drive real-world impact by taking action.

Tool / API / function calling to take actions.

GPT Operator / Claude Computer Use.

Orchestration

Workflow automation, task execution.

Human-in-the-loop systems & interfaces.

Action planning, execution, replanning...

# DoubleAgents uses agents in a feedback loop to solve HCPs.

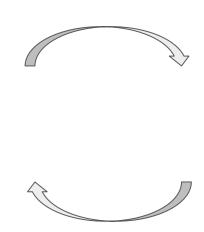




- Assigning work fairly
- Negotiating with workers' availability

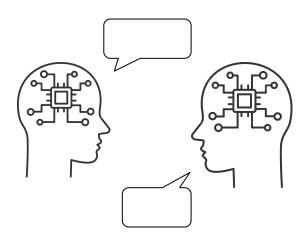
## Agents Implement

- Work Assignments
- Including contacting workers for availability questions and negotiation.



Al suggests revised policies when tests fail.

A human planner must approve new policies.



### Simulations **Evaluate** Solutions

- Test the planner for multiple scenarios of human available.
- Vary human communication styles, responsiveness, etc.

Example scenario:

Speaker scheduling for research seminars.

# Example scenario: Speaker scheduling for research seminars.

1:1 relationships (A speaker can only speak at 1 seminar slot.)

Availability constraints, preferences, and priority availability

Varied response times (some people never respond, some reply instantly)

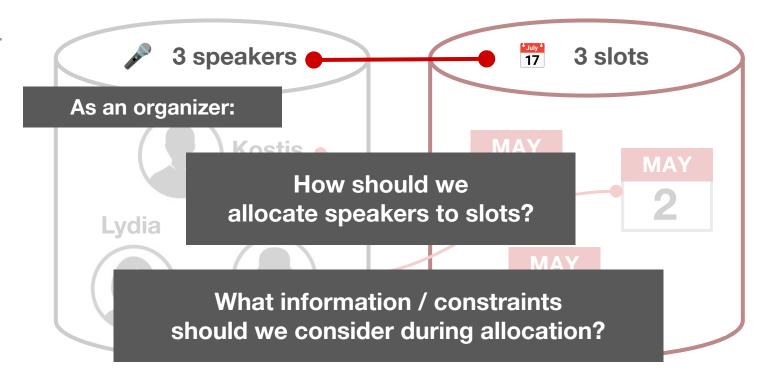
Extensive email communication (confirmation, request for additional availability, etc)

Multi-thread tracking email and sequential email follow-ups

# Example scenario:

Speaker scheduling for research seminars.

Example:





# 3 speakers



Speaker A



Speaker B



Speaker C



3 slots

# SLOT

Slot 1

1

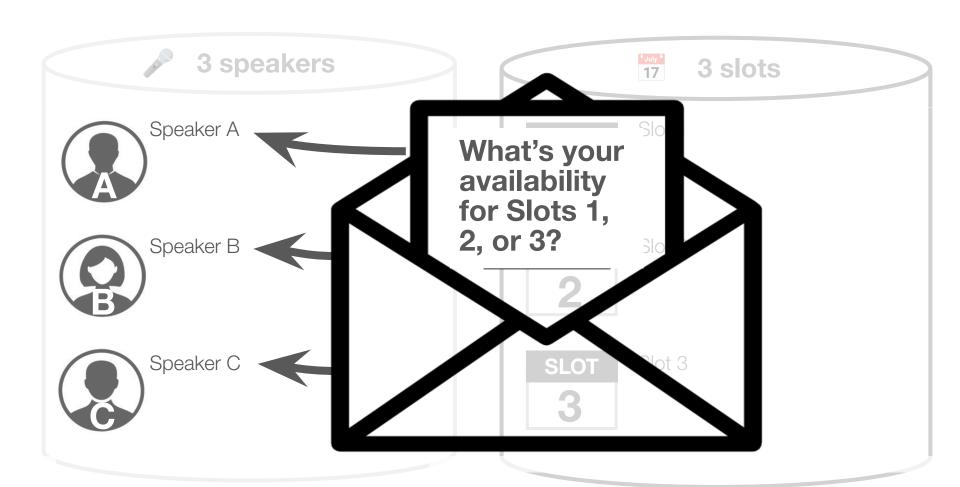
SLOT SI

2

Slot 2

SLOT Slot 3

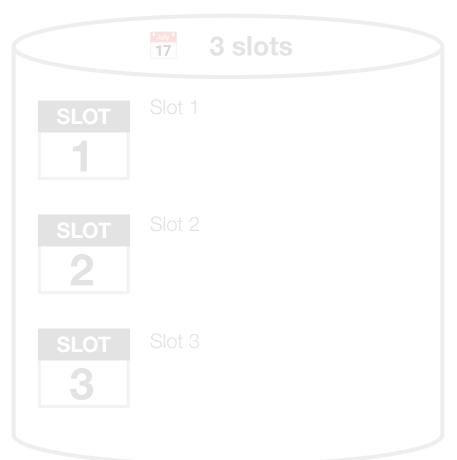
3



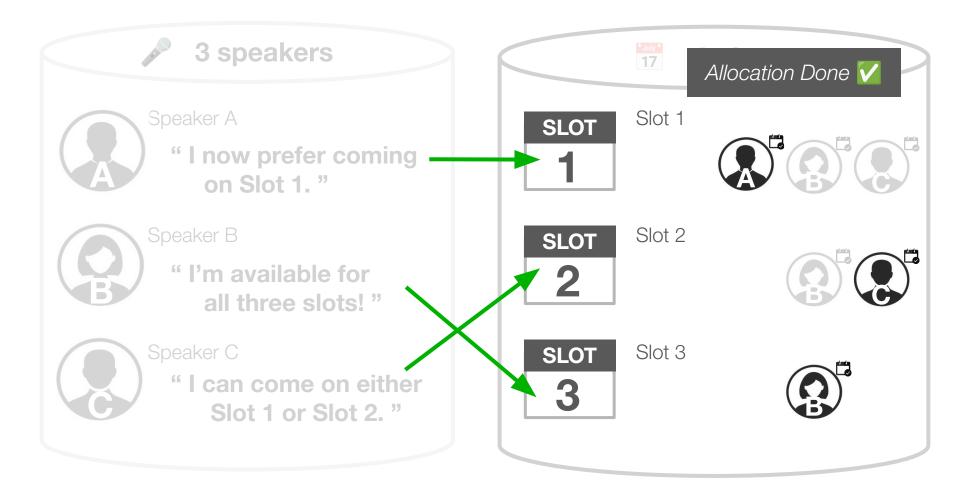


# Easy Case: No conflicted availability





# Easy Case: No conflicted availability





# **Medium Case:** Conflicted availability



# 3 speakers



Speaker A

"I now prefer coming on Slot 1."



Speaker B

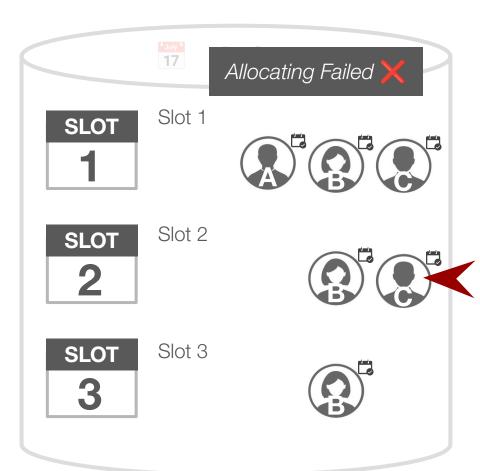
" I'm available for all three slots!"



Speaker C

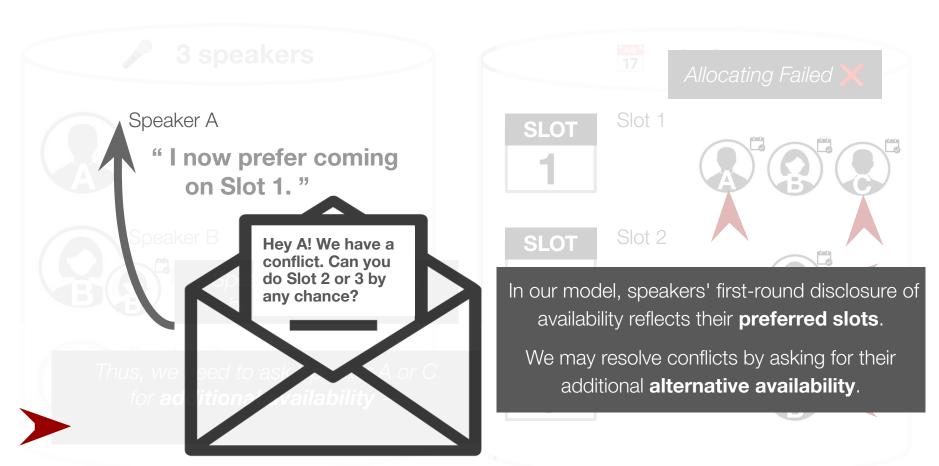
" I can come on Slot 1



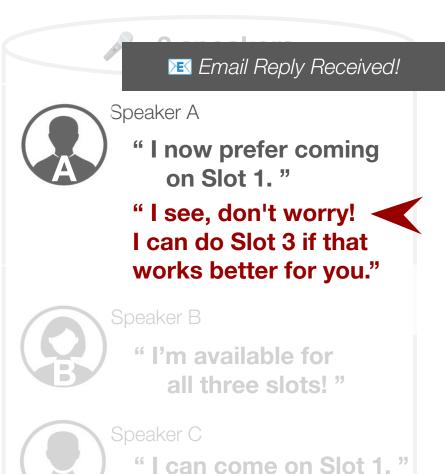


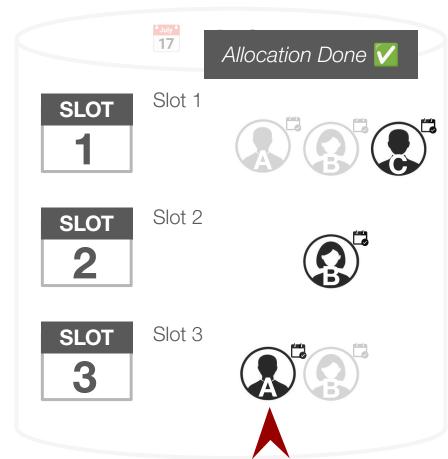


# Medium Case: Conflicted availability - must update policies



# Medium Case: Conflicted availability - must update policies

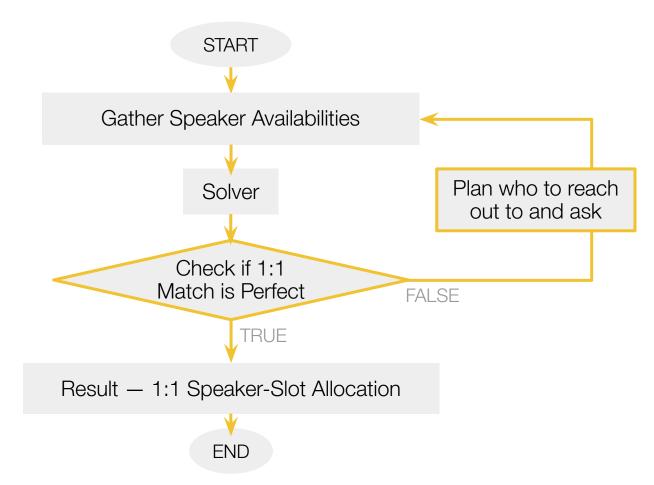




# **Design (Policies):**

 The solver can reach out to speakers again to ask for additional availability.

# **Execution Agent: HCP Solver**



# Repeat for harder test cases to discover more policies (designs) needed for real world deployment.

# Simulated test cases

1. Workers don't have enough availability.

2. Workers who don't respond to emails

# Policies to establish

 How many times is it okay to ask workers for more availability?

- How long should you wait to nudge?
- How many nudges should you do without a response?
- How you you balancing assigning dates to prompt responders versus hounding non-responders?

# DoubleAgents can solve all test cases, following new policies

	<b>Easy:</b> No conflicts	Medium: Few conflicts	Hard: Lots of conflicts
Is the allocation correct? The final speaker-slot allocation matches speakers agents' availabilities.	100%	100%	100%
# rounds of plan execution needed for achieving a perfect allocation			

## DoubleAgents can solve all test cases, following new policies

	<b>Easy:</b> No conflicts	Medium: Few conflicts	Hard: Lots of conflicts
Is the allocation correct? The final speaker-slot allocation matches speakers agents' availabilities.	100%	100%	100%
# rounds of plan execution needed for achieving a perfect allocation	0	1.67	3.38

# DoubleAgents uses agents in a feedback loop to solve HCPs.

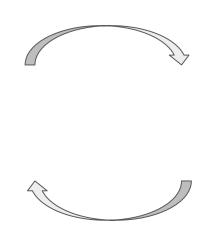




- Assigning work fairly
- Negotiating with workers' availability

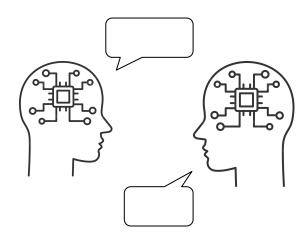
#### Agents Implement

- Work Assignments
- Including contacting workers for availability questions and negotiation.



Al suggests revised policies when tests fail.

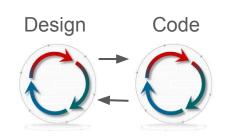
A human planner must approve new policies.



#### Simulations **Evaluate** Solutions

- Test the planner for multiple scenarios of human available.
- Vary human communication styles, responsiveness, etc.

## Bridging Design and Code with Gen Al





Iterative Problem Specification



**Iterative Development** 









## LogoMotion

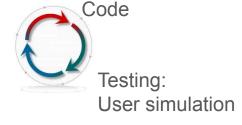
- Self-debugging loop
- Iteration with AI Editing widgets



## **Double Agents**

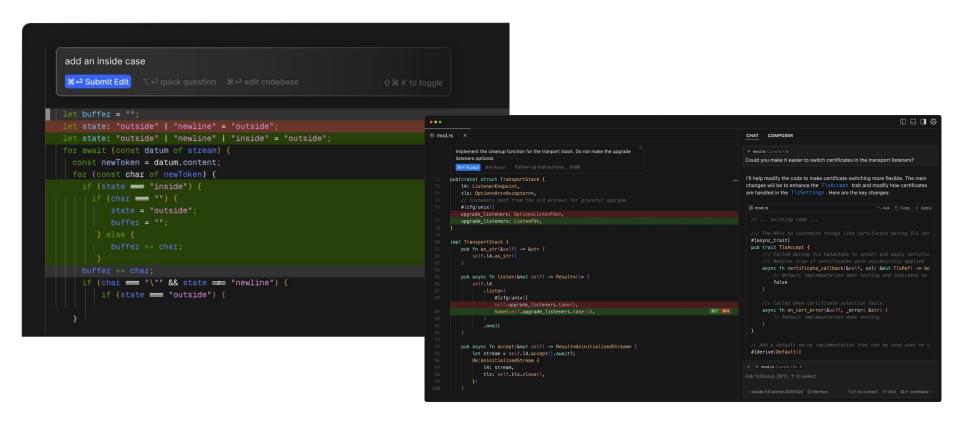
- Testing with user simulation
- Iterating on designs with AI suggestions







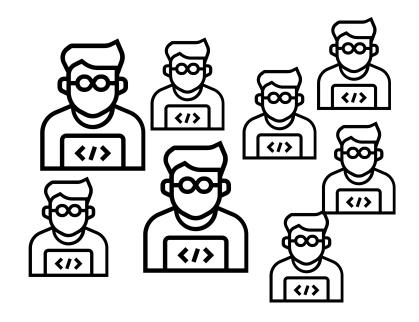
Al Code generation has transformed the way programmers work.



## Design (What do we build?)



Engineering (How do we build it?)

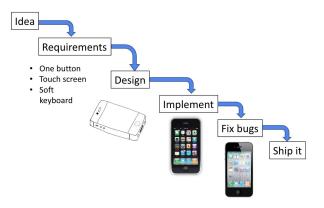


Created by Brickclay from Noun Project

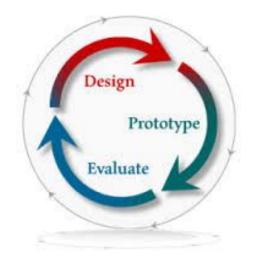
Created by WiStudio from Noun Project

## Two Software Design Processes

#### Feedforward

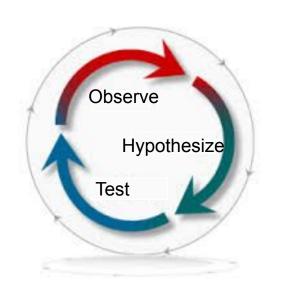


#### Feedback

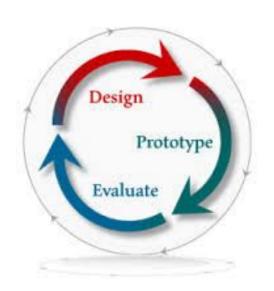


## Human-Centered Design: Two feedback loops

Person with a challenge





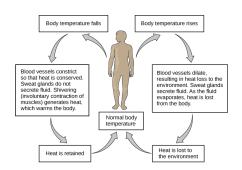




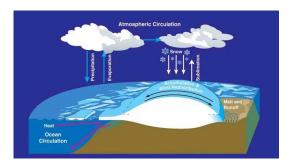
**Understand the problem** 

Solve the problem

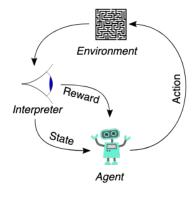
## Feedback loops are essential to all systems



The human body



Earth systems



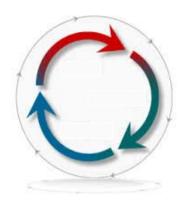
Training AI



Education



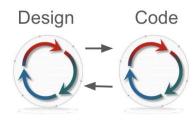
Childhood Development

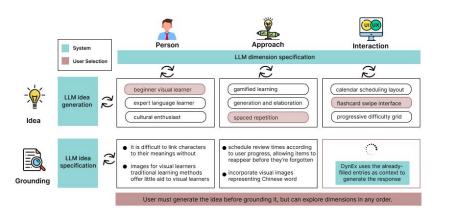


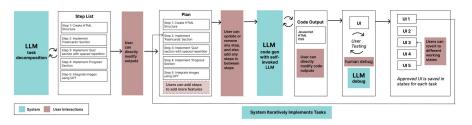
Al Code generation isn't enough to create better software products...

We have to close feedback loops between design and engineering.

## Dynex: Bridging Design and Code





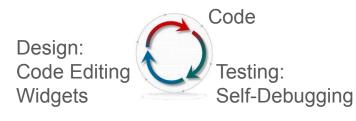


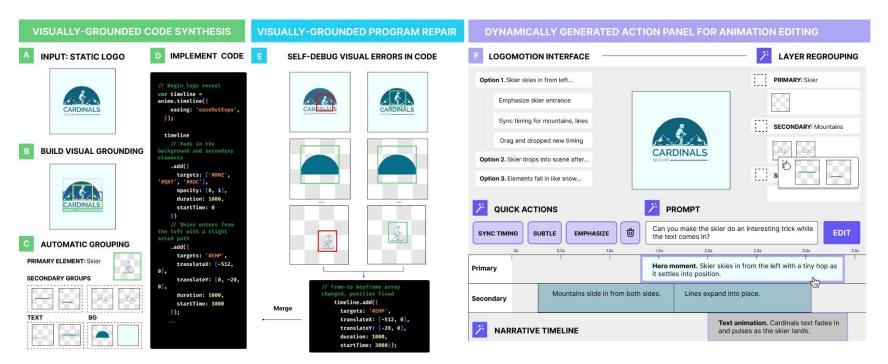
Iterative Problem
Specification (Design)

Modular Code Generation and Testing

### LogoMotion induces feedback with:

- 1. Self Debugging
- 2. Design iteration w/ Al Code Editing Widgets



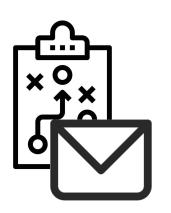


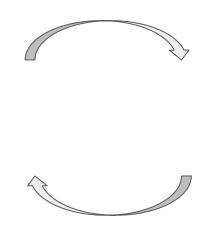
## DoubleAgents uses feedback for

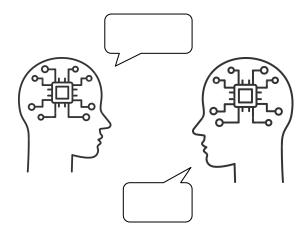
- 1) Simulating user testing
- 2) Policy (Design) Updates



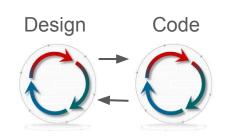








## Bridging Design and Code with Gen Al





Iterative Problem Specification



**Iterative Development** 









## LogoMotion

- Self-debugging loop
- Iteration with AI Editing widgets



## **Double Agents**

- Testing with user simulation
- Iterating on designs with AI suggestions



