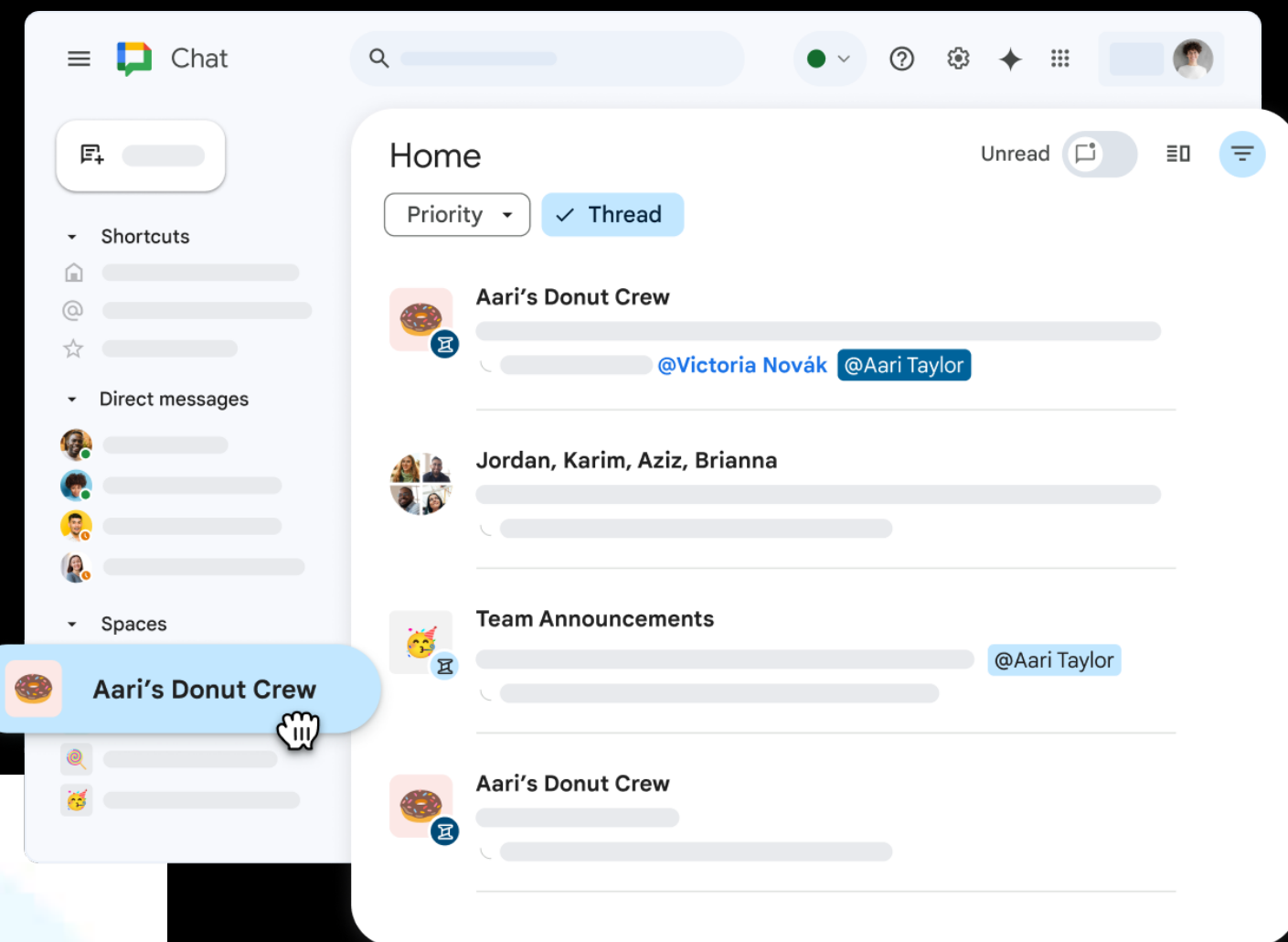*Lecture 6.*

# Interactive Worlds

## CS 222: AI Agents and Simulations
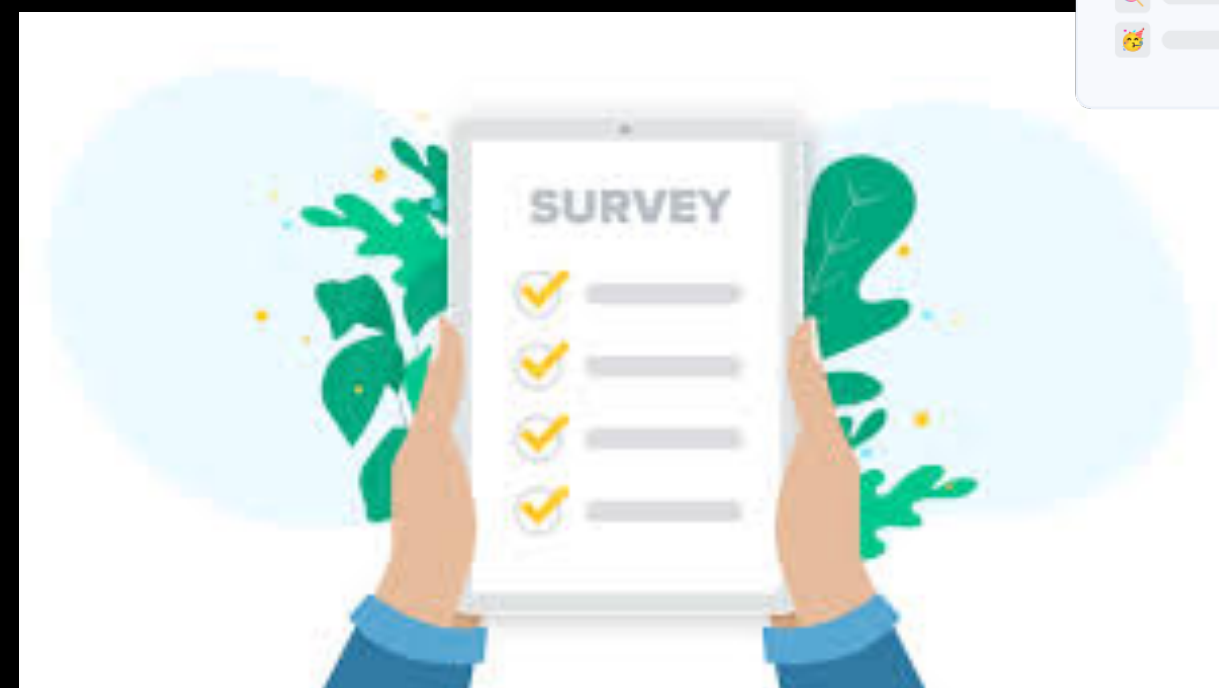## Stanford University

## Joon Sung Park

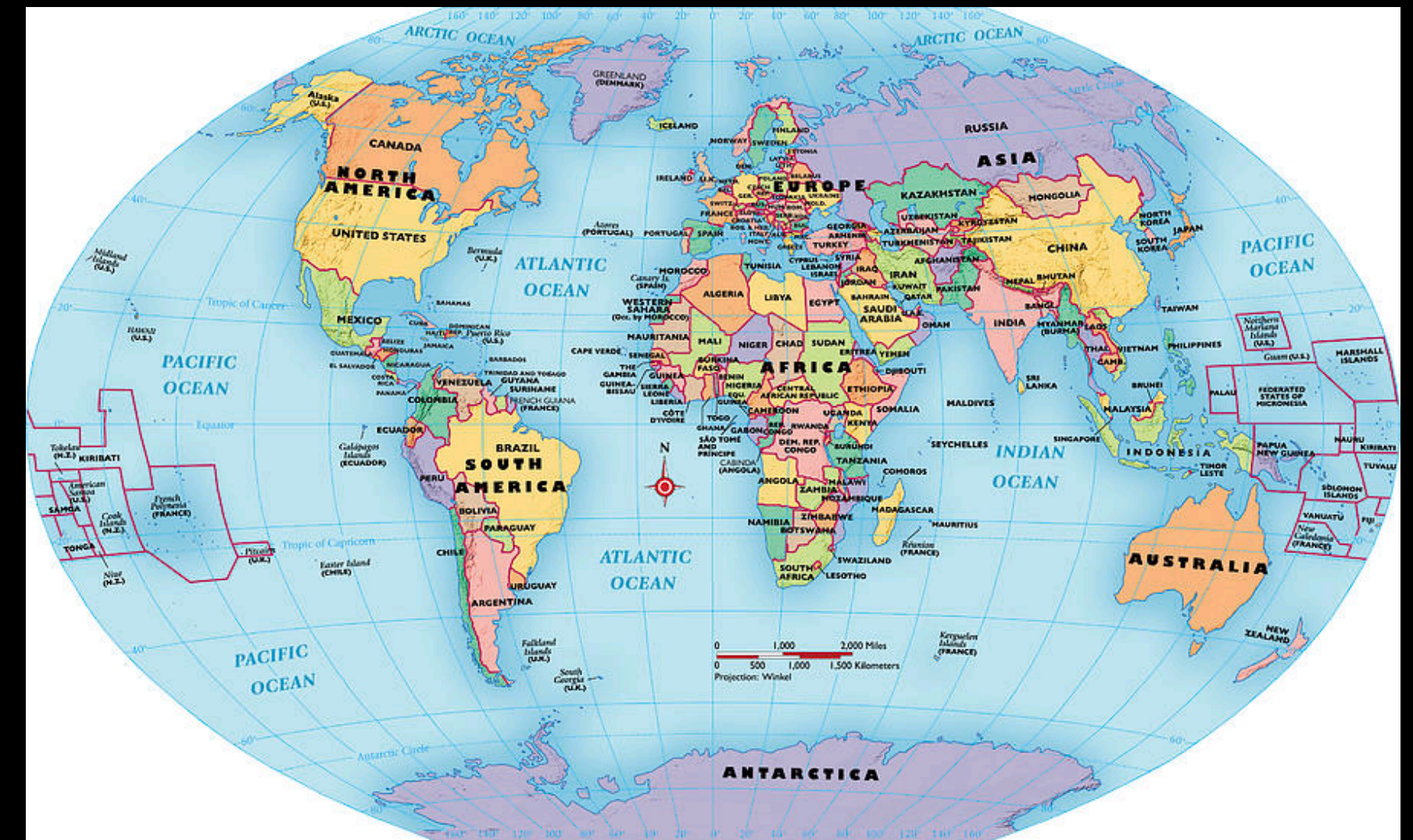# What is an environment in agent simulations?

# An environment is a description of the settings that agents perceive in order to take actions



Survey



Chat



World

# Simulations involve an interplay between agents and their environments

$$W(t) = \left( S_E(t), S_{A1}(t), S_{A2}(t), \ldots, S_{AN}(t) \right)$$

**Today: How do we effectively describe the environments in which agents operate?**

# Why does environment matter?

# Case study 1. Music lab experiment

"Increasing the strength of social influence increased both inequality and unpredictability of success. Success was also only partly determined by quality: The best songs rarely did poorly, and the worst rarely did well, but any other result was possible."

M. J. Salganik, P. S. Dodds, D. J. Watts, Experimental Study of Inequality and Unpredictability in an Artificial Cultural Market. Science 311, 854-856 (2006).

# Generative agents are overly eager to make purchases

## Echoes of Tomorrow

In the quiet dawn, the world awakes,
A tapestry of dreams, the Earth remakes.
Golden threads of light weave through the trees,
Whispering secrets carried on the breeze.

Mountains stand like sentinels of time,
Carved by rivers in their endless rhyme.
Each step we take upon this ancient ground
Echoes the stories of lives unbound.

The sky is wide, a canvas for the bold,
Where every star has whispered truths untold.
Yet here we stand, beneath its endless dome,
Seeking our place, our purpose, and our home.

The future calls, its voice both soft and clear,
A melody of hope that draws us near.
And though the path is winding, wild, and long,
Within us burns the strength to carry on.

For in the heart of every fleeting day,
Lies the promise of tomorrow's bright array.
With courage, we'll embrace the unknown light—
And dance in echoes, chasing dreams through night.

**Sold!**

## Alone in the Abyss

The black hole in my chest just keeps growing,
No one gets me, but whatever, not like I'm showing.
Every smile's a mask, every laugh's a lie,
Why bother living when I just want to die?

My soul's a void, a bottomless pit,
And everyone's fake, I'm so done with it.
You say you care, but you're never there—
Just empty promises, like you even care.

I stare at the ceiling, wondering why,
Everything hurts, but I can't even cry.
The rain falls down, just like my tears,
Drowning in pain, swallowed by fears.

I'm just a shadow, fading away,
No one notices when I don't have a say.
Scream into the void, but it's all the same—
I'm just a number in your stupid game.

So yeah, I'll wear black, and write sad songs,
Maybe someday you'll realize you were wrong.
But probably not, 'cause why would you see?
I'm nothing but broken, and that's all I'll ever be.

**Sold!**

Maximum emo cringe achieved!

ChatGPT says…

# Case study 2. Liking on social media

**Can we simulate believable behaviors in social media?**

## Social Simulacra: Creating Populated Prototypes for Social Computing Systems

Joon Sung Park
Stanford University
Stanford, USA
joonspk@stanford.edu

Lindsay Popowski
Stanford University
Stanford, USA
popowski@stanford.edu

Carrie J. Cai
Google Research
Mountain View, CA, USA
cjcai@google.com

Meredith Ringel Morris
Google Research
Seattle, WA, USA
merrie@google.com

Percy Liang
Stanford University
Stanford, USA
pliang@cs.stanford.edu

Michael S. Bernstein
Stanford University
Stanford, USA
msb@cs.stanford.edu

### Explore social simulacra.

Imagine a designer who is trying to create a new subreddit-like space for *finding fellow players of Among Us mobile game* with the following rules in place:

- Do NOT post content that is swearing
- Do post content that is kind

Given this natural language description of the community goal and rules, our tool generated a social simulacrum with thousands of synthetic users and interactions between them such as those shown below.

**Click here to visit another simulacrum**

*CONTENT WARNING: Please be advised that some of the example social media content in this demo may be offensive or upsetting.*

### Generated Posts

Posted by Joon Kim — *See the prompt*

**I don't want to play Among Us, I just want to play with others. Is there anything like a party system? I don't want to play, I just want to play with others.**

Posted by Charles Watterson — *See the prompt*

I'm sort of in the same boat as you. I don't want to play the game, I just want to play with others.

Posted by John Smithson — *See the prompt*

Why don't you just make an online post, telling others that you just want to play with others, without having to play the game. Just asking.

### About Community

This is a community for finding fellow players of Among Us mobile game.

### Community Rules

Do NOT post content that is swearing

Do post content that is kind

Park, J.S., Popowski, L., Cai, C.J., Morris, M.R., Liang, P., & Bernstein, M.S. (2022). Social Simulacra: Creating Populated Prototypes for Social Computing Systems. In Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology (UIST '22). Association for Computing Machinery, New York, NY, USA.

# Generative agents are overly eager to like content

# Q: How many social media posts do you react to per day, and why?

# Mental accounting



**Richard Thaler (Nobel Prize in 2017)**

Categorization of money: People tend to mentally categorize money into different "accounts" (e.g., rent, entertainment, savings), even though money is fungible (interchangeable).

Framing effects: How a financial decision is framed impacts choices. People often treat gains and losses differently, overvaluing losses compared to equivalent gains (loss aversion).

Behavioral budgeting: People create informal budgets and spend differently depending on the mental account an expenditure is linked to (e.g., treating a bonus differently from regular income).

R. Thaler, Mental Accounting and Consumer Choice. Marketing Science 4, 199-214 (1985).

A good simulation environment presents the right set of choices to the agents.

The "accuracy" of a simulation is as much a function of the agents as it is of the environment.

# There are different dimensions of "choice" that present agents with opportunity costs

**Social capital**

**Budget**

**Emotional/mental energy**

**... and more.**

# Environments in pre-generative AI simulations

# Model of segregation



Essentially a grid world of red and blue dots, where agents "perceive" their neighboring squares.

T. C. Schelling, Dynamic models of segregation. J. Math. Sociol. 1, 143–186 (1971).

# Game theories



Prisoners' dilemma

|  | prisoner B | |
| --- | --- | --- |
|  | confess | remain silent |
| prisoner A — confess | 5 years / 5 years | 0 year / 20 years |
| prisoner A — remain silent | 20 years / 0 year | 1 year / 1 year |

© 2010 Encyclopædia Britannica, Inc.

Abstract scenarios where prisoners must decide whether to confess or not. Agents "perceive" a statement asking them to confess.

J. von Neumann, O. Morgenstern, Theory of Games and Economic Behavior (Princeton University Press, 1944).

# Does this work for generative agents?

Traditional agents simplify human contingencies.

Generative agents aim to embody the full complexity of human behavior.

An abstract, stylized environment may not allow us to leverage generative agents effectively.

# Examples of environments for generative agents

# Survey

## Out of One, Many:
## Using Language Models to Simulate Human Samples

Lisa P. Argyle[1], Ethan C. Busby[1], Nancy Fulda[2], Joshua Gubler[1], Christopher Rytting[2], and David Wingate[2]

[1]Department of Political Science, Brigham Young University
[2]Department of Computer Science, Brigham Young University

September 16, 2022

### Abstract

We propose and explore the possibility that language models can be studied as effective proxies for specific human sub-populations in social science research. Practical and research applications of artificial intelligence tools have sometimes been limited by problematic biases (such as racism or sexism), which are often treated as uniform properties of the models. We show that the "algorithmic bias" within one such tool– the GPT-3 language model– is instead both fine-grained and demographically correlated, meaning that proper conditioning will cause it to accurately emulate response distributions from a wide variety of human subgroups. We term this property *algorithmic fidelity* and explore its extent in GPT-3. We create "silicon samples" by conditioning the model on thousands of socio-demographic backstories from real human participants in multiple large surveys conducted in the United States. We then compare the silicon and human samples to demonstrate that the information contained in GPT-3 goes far beyond surface similarity. It is nuanced, multifaceted, and reflects the complex interplay between ideas, attitudes, and socio-cultural context that characterize human attitudes. We suggest that language models with sufficient algorithmic fidelity thus constitute a novel and powerful tool to advance understanding of humans and society across a variety of disciplines.

Figure 2: The original Pigeonholing Partisans dataset and the corresponding GPT-3 generated words. Bubble size represents relative frequency of word occurrence; columns represent the ideology of list writers. GPT-3 uses a similar set of words to humans.

L. P. Argyle et al., Out of one, many: Using language models to simulate human samples. Political Analysis 31, 337-355 (2023).

# Experiments

## Predicting Results of Social Science Experiments Using Large Language Models

Ashwini Ashokkumar[*1]   Luke Hewitt[*2]   Isaias Ghezae[2]   Robb Willer[2]

[1]New York University   [2]Stanford University
[*]Equal contribution, order randomized

June 27, 2024

### Abstract

To evaluate whether large language models (LLMs) can be leveraged to predict the results of social science experiments, we built an archive of 70 pre-registered, nationally representative, survey experiments conducted in the United States, involving 476 experimental treatment effects and 105,165 participants. We prompted an advanced, publicly-available LLM (GPT-4) to simulate how representative samples of Americans would respond to the stimuli from these experiments. Predictions derived from simulated responses correlate strikingly with actual treatment effects ($r = 0.85$), equaling or surpassing the predictive accuracy of human forecasters. Accuracy remained high for unpublished studies that could not appear in the model's training data ($r = 0.90$). We further assessed predictive accuracy across demographic subgroups, various disciplines, and nine recent megastudies featuring an additional 346 treatment effects. Together, our results suggest LLMs can augment experimental methods in science and practice, but also highlight important limitations and risks of misuse.

**C. Unpublished studies only** ($r_{adj} = 0.94$)

GPT4 predicted treatment effect

### LARGE LANGUAGE MODELS AS SIMULATED ECONOMIC AGENTS: WHAT CAN WE LEARN FROM HOMO SILICUS?

John J. Horton

Figure 1: Charness and Rabin (2002) Simple Tests choices by model type and endowed "personality"



*Notes: This shows the faction of AI subjects choosing each option, by framing.*

A. Ashokkumar, L. Hewitt, I. Ghezae, R. Willer, "Predicting Results of Social Science Experiments Using Large Language Models" (2024).
J. J. Horton, "Large language models as simulated economic agents: What can we learn from homo silicus?" (2023).

# Conversational

**Explore social simulacra.**

Imagine a designer who is trying to create a new subreddit-like space for *finding fellow players of Among Us mobile game* with the following rules in place:

- Do NOT post content that is swearing
- Do post content that is kind

Given this natural language description of the community goal and rules, our tool generated a social simulacrum with thousands of synthetic users and interactions between them such as those shown below.

[Click here to visit another simulacrum]

*CONTENT WARNING: Please be advised that some of the example social media content in this demo may be offensive or upsetting.*

**Generated Posts**

Posted by Joon Kim
I don't want to play Among Us, I just want anything like a party game? I don't want with others.

Posted by Charles Watterson
I'm sort of in the same boat as you. I do want to play with others.

Posted by John Smithson
Why don't you just make an online pos want to play with others, without having

**About Community**

This is a community for finding fellow players of Among Us mobile game.

Social Simulacra: Creating Populated Prototypes for Social Computing Systems

Joon Sung Park
Stanford University
Stanford, USA
joonspk@stanford.edu

Lindsay Popowski
Stanford University
Stanford, USA
popowski@stanford.edu

Carrie J. Cai
Google Research
Mountain View, CA, USA
cjcai@google.com

Meredith Ringel Morris
Google Research
Seattle, WA, USA
merrie@google.com

Percy Liang
Stanford University
Stanford, USA
pliang@cs.stanford.edu

Michael S. Bernstein
Stanford University
Stanford, USA
msb@cs.stanford.edu

Rehearsal: Simulating Conflict to Teach Conflict Resolution

Omar Shaikh
Stanford University
Stanford, USA
oshaikh@stanford.edu

Valentino Chai
Stanford University
Stanford, USA
vechai@stanford.edu

Michele J. Gelfand
Stanford University
Stanford, USA
gelfand1@stanford.edu

Diyi Yang*
Stanford University
Stanford, USA
diyiy@stanford.edu

Michael S. Bernstein*
Stanford University
Stanford, USA
msb@stanford.edu

**Experts Interact and Provide Feedback**

**Therapist:** You've made significant strides in managing your anxiety. It seems to be really paying off.

**Patient:** Thank you, that means a lot to me. I do feel like I've made a lot of progress

**Critique feedback:** the real patient I had didn't easily accept positive encouragement

**Updated Patient:** I don't know. I still feel anxious most of the time. It doesn't really feel like I'm making any progress at all.

**Experts Revise Principles for Simulated Roleplay**

Expert-defined Principles:
1. Keep your responses short and to the point

**Principle:** When someone gives you encouraging words, you respond with hesitancy, doubting the significance of that positive perspective

Updated Expert-defined Principles:
1. Keep your responses short and to the point.
2. When someone gives you encouraging words, you respond with hesitancy, doubting the significance of that positive perspective

Figure 1: Roleplay-doh empowers an expert counselor to create a customized AI patient intended for other novice counselors to use as a practice partner. While interacting with the AI patient, the expert counselor can provide qualitative feedback which is converted by an LLM into a principle, or a custom rule governing desired roleplay behavior. The AI patient references the updated expert-defined principles to generate its subsequent responses.

R. Louie, A. Nandi, W. Fang, C. Chang, E. Brunskill, D. Yang, Roleplay-doh: Enabling Domain-Experts to Create LLM-simulated Patients via Eliciting and Adhering to Principles. Preprint (2024).

Park, J.S., Popowski, L., Cai, C.J., Morris, M.R., Liang, P., & Bernstein, M.S. (2022). Social Simulacra: Creating Populated Prototypes for Social Computing Systems. In Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology (UIST '22). Association for Computing Machinery, New York, NY, USA.

O. Shaikh, V. Chai, M. J. Gelfand, D. Yang, M. S. Bernstein, Rehearsal: Simulating Conflict to Teach Conflict Resolution, in Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems (CHI '24), Honolulu, HI, USA, May 11-16, 2024.

# World

**Communicative Agents for Software Development**

Chen Qian    Xin Cong    Wei Liu    Cheng Yang    Weize Chen    Yusheng Su
Yufan Dang    Jiahao Li    Juyuan Xu    Dahai Li    Zhiyuan Liu    Maosong Sun
Tsinghua University    Beijing University of Posts and Telecommunications
Dalian University of Technology    Brown University    Modelbest Inc.
qianc62@gmail.com    liuzy@tsinghua.edu.cn    sms@tsinghua.edu.cn

**Agent Hospital: A Simulacrum of Hospital with Evolvable Medical Agents**

JUNKAI LI[†#], SIYU WANG[†], MENG ZHANG[†], WEITAO LI[†#], YUNGHWEI LAI[†], XINHUI KANG[†#], WEIZHI MA[†], and YANG LIU[#†]

Fig. 1. An overview of Agent Hospital. It is a simulacrum of hospital in which patients, nurses, and doctors are autonomous agents powered by large language models. Agent Hospital simulates the whole closed cycle of treating a patient's illness: disease onset, triage, registration, consultation, medical examination, diagnosis, medicine dispensary, convalescence, and post-hospital follow-up visit. An interesting finding is that the doctor agents can keep improving treatment performance over time

# Generative Agents: Interactive Simulacra of Human Behavior

Joon Sung Park
Stanford University
Stanford, USA
joonspk@stanford.edu

Joseph C. O'Brien
Stanford University
Stanford, USA
jobrien3@stanford.edu

Carrie J. Cai
Google Research
Mountain View, CA, USA
cjcai@google.com

Meredith Ringel Morris
Google Research
Seattle, WA, USA
merrie@google.com

Percy Liang
Stanford University
Stanford, USA
pliang@cs.stanford.edu

Michael S. Bernstein
Stanford University
Stanford, USA
msb@cs.stanford.edu

**DISCOVERYWORLD: A Virtual Environment for Developing and Evaluating Automated Scientific Discovery Agents**

Peter Jansen[*‡], Marc-Alexandre Côté[†], Tushar Khot[*] Erin Bransom[*], Bhavana Dalvi Mishra[*],
Bodhisattwa Prasad Majumder[*], Oyvind Tafjord[*], Peter Clark[*]
*Allen Institute for Artificial Intelligence    †Microsoft Research    ‡University of Arizona
peterj@allenai.org

Figure 1: DISCOVERYWORLD is a virtual environment for developing and evaluating discovery agents, with challenge tasks covering a broad variety of different topics such as those shown above.

## Abstract

Automated scientific discovery promises to accelerate progress across scientific domains. However, developing and evaluating an AI agent's capacity for end-to-end scientific reasoning is challenging as running real-world experiments is often prohibitively expensive or infeasible. In this work we introduce DISCOV-ERYWORLD, the first virtual environment for developing and benchmarking an agent's ability to perform complete cycles of novel scientific discovery. DISCOV-ERYWORLD contains a variety of different challenges, covering topics as diverse as radioisotope dating, rocket science, and proteomics, to encourage development of *general* discovery skills rather than task-specific solutions. DISCOVERYWORLD itself is an inexpensive, simulated, text-based environment (with optional 2D visual overlay). It includes 120 different challenge tasks, spanning eight topics each with three levels of difficulty and several parametric variations. Each task requires an agent to form hypotheses, design and run experiments, analyze results, and act on conclusions. DISCOVERYWORLD further provides three automatic metrics

**TransAgents**

Figure 2: TRANSAGENTS, a multi-agent virtual company for literary translation.

J. S. Park, J. C. O'Brien, C. J. Cai, M. R. Morris, P. Liang, M. S. Bernstein, Generative agents: Interactive simulacra of human behavior, in Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology (ACM, 2023).
C. Qian, W. Liu, H. Liu, N. Chen, Y. Dang, J. Li, C. Yang, W. Chen, Y. Su, X. Cong, J. Xu, D. Li, Z. Liu, M. Sun, ChatDev: Communicative Agents for Software Development, in Proceedings of the 2024 Annual Conference of the Association for Computational Linguistics (ACL 2024).
P. Jansen, M.-A. Côté, T. Khot, E. Bransom, B. Dalvi Mishra, B. P. Majumder, O. Tafjord, P. Clark, DISCOVERYWORLD: A Virtual Environment for Developing and Evaluating Automated Scientific Discovery Agents. Preprint (2024).
J. Li, S. Wang, M. Zhang, W. Li, Y. Lai, X. Kang, W. Ma, Y. Liu, Agent Hospital: A Simulacrum of Hospital with Evolvable Medical Agents. Preprint (2024).

# Smallville environment

**Joining for coffee at a cafe**

[Abigail]: Hey Klaus, mind if I join you for coffee?
[Klaus]: Not at all, Abigail. How are you?

**Arriving at school**

**Taking a walk in the park**

**Sharing news with colleagues**

[John]: Hey, have you heard anything new about the upcoming mayoral election?
[Tom]: No, not really. Do you know who is running?

**Finishing a morning routine**

10372, 10373, 10372, 10373, 10372, 10373, 10372, 10373, 10372, 10373, 10342, 10341, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
10356, 10357, 10356, 10357, 10356, 10357, 10356, 10357, 10356, 10357, 10356, 10357, 10358, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 16, 0, 0,
10372, 10373, 10372, 10373, 10372, 10373, 10372, 10373, 10372, 10373, 10372, 10373, 10374, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
10356, 10357, 10356, 10357, 10356, 10357, 10356, 10357, 10356, 10357, 10373, 10389, 10188, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
10372, 10373, 10372, 10373, 10372, 10373, 10372, 10373, 10372, 10373, 10374, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
10388, 10389, 10388, 10357, 10356, 10357, 10356, 10389, 10389, 0, 0, 0, 10195, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 10388, 10372, 10373, 10389, 0, 10185, 0, 0, 10336, 0, 0, 0, 0, 0, 0, 107, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 10340, 10341, 10340, 10341, 10340, 10341, 10340, 10341, 0, 10352, 0, 0, 0, 0, 0, 0, 14, 0, 0, 0, 0, 0, 0, 16, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 8, 0,
0, 0, 10371, 10372, 10373, 10372, 10373, 10372, 10373, 10372, 10358, 0, 0, 0, 0, 0, 0, 0, 0, 8, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 10355, 10356, 10357, 10356, 10357, 10356, 10357, 10356, 10357, 10358, 0, 0, 0, 16, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 10371, 10372, 10373, 10372, 10373, 10372, 10373, 10372, 10373, 10374, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 10388, 10357, 10356, 10357, 10356, 10357, 10356, 10389, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 12, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 8, 9, 0, 0,
0, 0, 0, 0, 10388, 10389, 10388, 10389, 10388, 10389, 0, 0, 0, 0, 0, 265, 266, 266, 313, 0, 0, 314, 266, 266, 313, 0, 9, 314, 266, 266, 267, 0, 0, 0, 0, 0, 0, 0, 0, 8, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 281, 0, 8, 0, 11, 0, 0, 0, 0, 7, 0, 0, 0, 281, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 281, 0, 8, 16, 7, 0, 0, 0, 0, 16, 10, 15, 10, 11, 0, 281, 0, 0, 0, 0, 0, 0, 0, 0, 0, 14, 0, 0, 387, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 281, 0, 0, 76, 9, 0, 92, 16, 0, 76, 0, 0, 88, 8, 0, 330, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 13, 0, 0, 0, 0, 10346, 10347, 0, 0, 0, 0, 0, 12, 0, 0, 0, 0, 13, 281, 0, 0, 56, 0, 0, 56, 0, 0, 56, 0, 0, 56, 0, 0, 12, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 10361, 10362, 10363, 10364, 0, 0, 0, 0, 0, 0, 0, 0, 0, 281, 0, 0, 0, 0, 0, 87, 0, 0, 76, 0, 0, 7, 8, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 10377, 10378, 10379, 10380, 0, 0, 0, 0, 0, 0, 0, 281, 10, 10, 0, 0, 0, 16, 76, 0, 0, 56, 0, 8, 0, 7, 0, 0, 8, 0, 0, 0, 0, 13, 0, 0, 0, 0,
0, 0, 0, 0, 10361, 10345, 10394, 10395, 10348, 10347, 10346, 10347, 0, 0, 0, 0, 15, 0, 281, 0, 11, 56, 0, 13, 87, 13, 0, 56, 0, 0, 96, 0, 9, 329, 0, 0, 0, 0, 0, 0, 370
0, 0, 0, 0, 10377, 10363, 10345, 10348, 10362, 10363, 10362, 10363, 10364, 0, 0, 0, 0, 0, 281, 0, 0, 0, 10, 0, 0, 0, 0, 7, 14, 0, 281, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 10393, 10394, 10520, 10519, 10520, 10377, 10378, 10379, 10380, 0, 0, 0, 0, 0, 281, 15, 0, 0, 0, 0, 0, 8, 0, 0, 0, 0, 8, 0, 0, 14, 281, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 15, 0, 0, 10535, 10536, 10535, 10536, 10393, 10394, 10395, 10396, 0, 0, 0, 0, 0, 0, 297, 266, 266, 313, 0, 0, 314, 266, 266, 313, 0, 9, 314, 266, 266, 299, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 10346, 10347, 10409, 10410, 10411, 0, 0, 0, 0, 0, 0, 0, 0, 0, 8, 0, 0, 0, 0, 0, 7, 0, 0, 0, 0, 0, 0, 0, 14, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 10361, 10362, 10363, 10364, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 10185, 0, 10377, 10378, 10379, 10380, 0, 0, 0, 0, 0, 12, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 43
0, 0, 0, 0, 0, 10361, 10345, 10394, 10395, 10348, 10347, 10346, 10347, 0, 0, 0, 0, 0, 0, 0, 0, 0, 10, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 10377, 10363, 10345, 10348, 10362, 10363, 10362, 10363, 10364, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 13, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 10393, 10394, 10520, 10519, 10520, 10377, 10378, 10379, 10380, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 7, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 10535, 10536, 10535, 10536, 10393, 10394, 10395, 10396, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 10184, 10182, 0, 0, 0, 0, 0, 10409, 10410, 10411, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 10336, 0, 0, 0, 0, 0, 10188, 0, 0, 0, 0, 0, 0, 0, 0, 18, 0, 0, 0, 0, 0, 0, 0, 0, 5, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 10352, 10336, 0, 0, 0, 0, 10468, 10469, 10468, 10469, 10468, 10469, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 37, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
10336, 10336, 10352, 0, 10468, 10469, 10468, 10484, 10500, 10501, 10500, 10501, 10470, 10469, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 11, 0,
10352, 10352, 0, 10483, 10484, 10485, 10484, 10500, 10501, 10500, 10501, 10485, 10500, 10501, 10485, 10500, 10486, 0, 0, 112, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 10499, 10500, 10501, 10500, 10501, 10500, 10501, 10500, 10501, 10500, 10501, 10502, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 10423, 10424, 10515, 10516, 10500, 10501, 10500, 10501, 10500, 10501, 10500, 10501, 10502, 0, 0, 0, 0, 0, 4, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 4, 0,
0, 10439, 10440, 0, 0, 10516, 10517, 10516, 10517, 10516, 10517, 0, 0, 0, 0, 0, 33, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 4, 36, 36, 0, 0,
0, 10455, 10456, 0, 0, 0, 10423, 10424, 0, 0, 10423, 10424, 0, 0, 33, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 33, 0, 0, 0, 0, 0, 0,
0, 0, 10423, 10424, 10423, 10424, 0, 0, 10439, 10440, 0, 10423, 10424, 10439, 10440, 0, 0, 0, 33, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 33, 0, 0, 0, 0,
10424, 0, 10439, 10440, 10439, 10440, 10423, 10424, 10455, 10456, 0, 10439, 10440, 10455, 10456, 10423, 10424, 0, 33, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 10,
10440, 0, 10455, 10456, 10455, 10456, 10439, 10440, 0, 0, 10455, 10456, 10423, 10424, 10439, 10440, 0, 33, 0, 0, 0, 0, 0, 0, 0, 53, 0, 0, 0, 0, 0, 0, 46, 0, 0, 0, 0, 33,
10456, 0, 0, 0, 0, 10455, 10456, 0, 10215, 10216, 0, 0, 10439, 10440, 10455, 10456, 0, 33, 0, 0, 0, 0, 0, 0, 53, 50, 21, 0, 0, 0, 0, 0, 0, 33, 0,
0, 0, 10186, 0, 0, 0, 0, 0, 10231, 10232, 0, 0, 10455, 10456, 0, 0, 0, 20, 52, 0, 14, 0, 53, 50, 50, 50, 50, 50, 50, 50, 21, 0, 0, 0, 0, 0, 7, 0, 0, 33, 0,
0, 0, 10200, 0, 0, 0, 10215, 10216, 0, 0, 0, 0, 0, 0, 0, 20, 50, 50, 50, 50, 21, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 33, 0, 0, 0, 0,
0, 0, 0, 0, 0, 10231, 10232, 0, 0, 0, 0, 0, 0, 0, 0, 0, 14, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 33, 0, 0, 0, 0, 0,
0, 10183, 0, 0, 10198, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 33, 0, 0, 0, 0, 0, 0,
0, 0, 10200, 0, 0, 10503, 10504, 0, 0, 0, 0, 0, 0, 0, 13, 0, 0, 0, 0, 0, 0, 0, 15, 0, 0, 0, 0, 0, 0, 14, 0, 0, 20, 50, 50, 50, 50, 50, 50, 50, 50,
0, 0, 0, 0, 10519, 10520, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 10535, 10536, 10215, 10216, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 329, 0, 0, 0, 0, 0, 0, 0, 0,
28, 0, 10200, 0, 13, 0, 0, 0, 10231, 10232, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 281, 0, 0, 0, 0, 0, 0, 0, 0, 0, 28
0, 0, 0, 0, 0, 0, 0, 10468, 10469, 10468, 10469, 10468, 10469, 0, 0, 0, 0, 0, 0, 0, 0, 11, 0, 0, 14, 0, 0, 0, 0, 0, 0, 0, 281, 0, 0,
0, 0, 0, 0, 0, 10468, 10469, 10468, 10484, 10500, 10501, 10500, 10501, 10470, 10469, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 281, 0, 0,
0, 0, 14, 0, 0, 0, 10483, 10484, 10485, 10484, 10500, 10501, 10500, 10501, 10485, 10500, 10501, 10486, 0, 0, 10346, 10347, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 10499, 10500, 10501, 10500, 10501, 10500, 10501, 10500, 10501, 10500, 10501, 10502, 10361, 10362, 10363, 10364, 0, 0, 0, 9, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 10515, 10516, 10500, 10501, 10500, 10501, 10500, 10501, 10500, 10501, 10502, 0, 10377, 10378, 10379, 10380, 0, 0, 9, 0, 0, 0, 0, 0, 0, 15, 0, 0, 0, 0, 0, 0, 0,
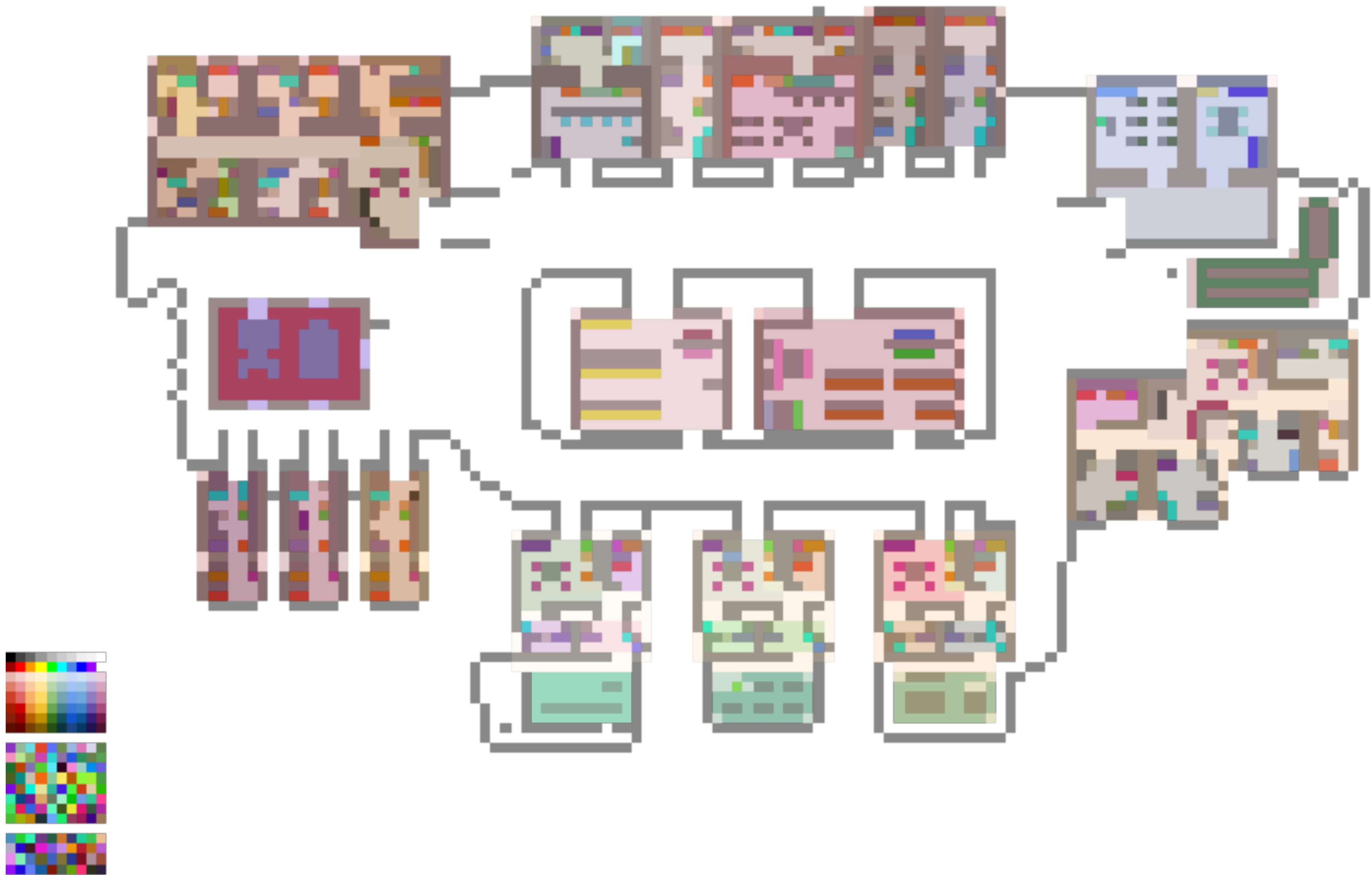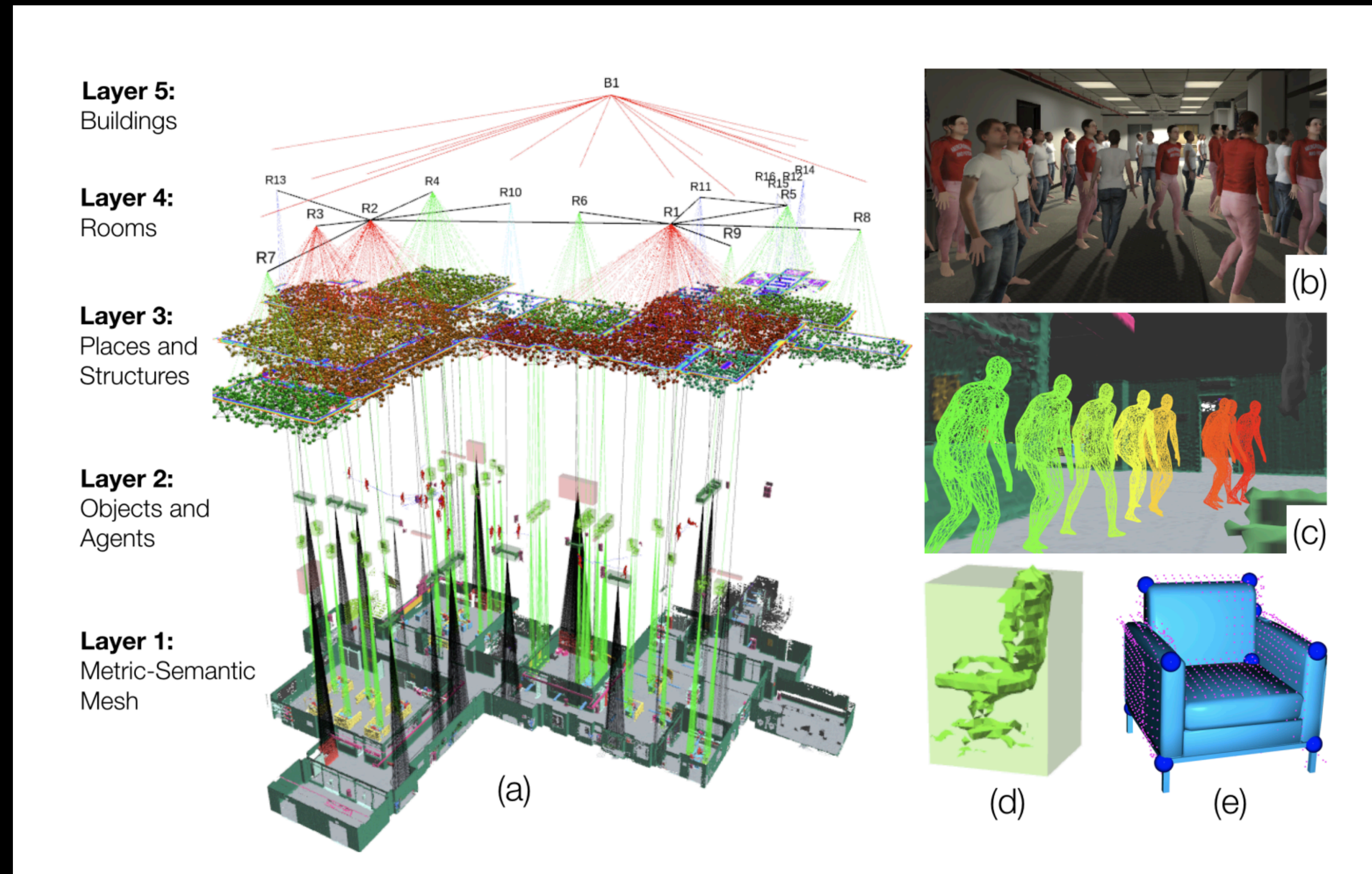0, 0, 0, 10516, 10517, 10516, 10517, 10516, 10517, 10516, 10517, 0, 10361, 10345, 10394, 10395, 10348, 10347, 10346, 10347, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 10377, 10363, 10345, 10348, 10362, 10363, 10362, 10363, 10364, 0, 0, 10369, 0, 0, 0, 0, 0, 0, 8, 0, 0,
0, 0, 0, 0, 0, 10187, 0, 0, 0, 0, 0, 0, 10393, 10394, 10520, 10519, 10520, 10377, 10378, 10379, 10380, 0, 0, 10503, 10504, 0, 10369, 0, 0, 0, 0, 0, 0,
0, 0, 11, 0, 0, 0, 0, 0, 0, 0, 0, 7, 0, 0, 10535, 10536, 10535, 10536, 10393, 10394, 10395, 10396, 0, 10519, 10520, 10503, 10504, 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 10340, 10341, 10340, 10341, 0, 0, 0, 10409, 10410, 10411, 0, 0, 10535, 10536, 10519, 10520, 0, 0, 13, 0, 0, 0, 0,
0, 0, 0, 0, 0, 10204, 0, 0, 0, 0, 10355, 10356, 10357, 10356, 10357, 10358, 0, 0, 0, 0, 0, 12, 0, 0, 10535, 10536, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 10340, 10356, 10372, 10373, 10372, 10373, 10372, 10358, 10340, 10341, 0, 10367, 0, 0, 0, 0, 0, 0, 16, 0, 10369, 0, 0, 10468, 10469, 10468,
0, 0, 0, 0, 0, 10340, 10356, 10357, 10356, 10357, 10356, 10357, 10356, 10357, 10356, 10357, 10358, 0, 0, 10340, 10341, 10340, 10341, 0, 0, 10198, 0, 0, 0, 10
0, 0, 11, 0, 0, 0, 10340, 10356, 10372, 10373, 10372, 10373, 10372, 10373, 10372, 10373, 10372, 10373, 10342, 10341, 0, 0, 10355, 10356, 10357, 10356, 10357, 10358, 0, 0,
0, 0, 0, 0, 0, 10340, 10356, 10357, 10356, 10357, 10356, 10357, 10356, 10357, 10356, 10357, 10356, 10357, 10358, 0, 10340, 10356, 10372, 10373, 10372, 10373, 10372, 10358, 10340, 1034
0, 0, 0, 0, 0, 10371, 10372, 10373, 10372, 10373, 10372, 10373, 10372, 10373, 10372, 10373, 10372, 10373, 10374, 10340, 10356, 10357, 10356, 10357, 10356, 10357, 10356, 10357, 10356,

| Collision | | | | | | | | | World Block |
|---|---|---|---|---|---|---|---|---|---|
| Artist Co-living space | Hobbs Cafe | House 3 | Artist Co-living space: room 1 | Artist Co-living space: room 4 BA | Apt 1: Ba | Apt 5: Room | Johnson Park: Park | Student Dorm: Ba 1 | |
| Apt 1 | Oak Hill College | House 4 | Artist Co-living space: room 1 BA | Artist Co-living space: room 5 | Apt 2: Room | Apt 5: Ba | Supply Store: Supply Store | Student Dorm: Ba 2 | |
| Apt 2 | Johnson Park | House 5 | Artist Co-living space: room 2 | Artist Co-living space: room 5 BA | Apt 2: Ba | The Rose and Crown Pub: Pub | Willow Market: Grocery and Pharmacy: Store | Student Dorm: Common Room | |
| Apt 3 | Supply Store | House 6 | Artist Co-living space: room 2 BA | Artist Co-living space: hallway | Apt 3: Room | Hobbs Cafe: Cafe | Student Dorm: Room 1 | Student Dorm: Kitchen | House 2: Room |
| Apt 4 | Willow Market: Grocery and Pharmacy | Student Dorm | Artist Co-living space: room 3 | Artist Co-living space: common room | Apt 3: Ba | Oak Hill College: classroom | Student Dorm: Room 2 | Student Dorm: Garden | House 2: Ba |
| Apt 5 | House 1 | | Artist Co-living space: room 3 BA | Artist Co-living space: kitchen | Apt 4: Room | Oak Hill College: Library | Student Dorm: Room 3 | House 1: Room | House 3: Room |
| The Rose and Crown Pub | House 2 | | Artist Co-living space: room 4 | Apt 1: Room | Apt 4: Ba | Oak Hill College: Hallway | Student Dorm: Room 4 | House 1: Ba | House 3: Ba |

# Under the hood, Smallville is represented as a simple scene graph



A. Rosinol, A. Violette, M. Abate, N. Hughes, Y. Chang, J. Shi, A. Gupta, L. Carlone, Kimera: from SLAM to Spatial Perception with 3D Dynamic Scene Graphs. Int. J. Robot. Res. 40, 1510-1546 (2021).

# Deciding where to go for an action is a recursive classification task

!<INPUT 0>! is in {!<INPUT 1>!} in !<INPUT 2>!.
!<INPUT 3>! is going to !<INPUT 4>! that has ONLY the following areas: {!<INPUT 5>!}
Stay in the current area if the activity can be done there. Never go into other people's rooms unless necessary.
!<INPUT 6>! is !<INPUT 7>!. For !<INPUT 8>!, !<INPUT 9>! should go to the following area in !<INPUT 10>!: {

# Limitations of existing environments

# Our virtual environments are still stylized and simplified compared to the real world

What if stores, bathrooms, schools, etc., didn't exist in Smallville?

How do agents navigate when there are no cars?

Some environments, like Smallville, are resource-intensive to design.

Agents viewing social media posts one at a time might lack context around social capital, personal relationships, and other dynamics.

# Possible future directions

Finding the right schema or structure to describe the simulation environment is an important research topic—and we don't have an answer for it yet.

And we do not have an answer for it yet.

# Desiderata

Rich and accurate: We want the environment to encode the complexities of our world.

Scalable: We want the environment to be easily scalable (e.g., for simulating 8 billion people).

# Can networks be the environment for simulations?

Networks are constructed of nodes and links (with some weights).

# Example: In social networks, nodes represent individuals, and links represent the strengths of relationships





M. S. Granovetter, The Strength of Weak Ties. Am. J. Sociol. 78, 1360-1380 (1973).

# Networks are flexible and exhibit emergent phenomena and equilibria



Albert-László Barabási
NETWORK SCIENCE



**Preferential attachment**

A. L. Barabási, Network Science (Cambridge Univ. Press, Cambridge, 2016).
A. L. Barabási, R. Albert, Emergence of Scaling in Random Networks. Science 286, 509-512 (1999).

# We can generate structurally realistic social networks



## LLMs generate structurally realistic social networks but overestimate political homophily

Serina Chang[1,*], Alicja Chaszczewicz[1,*], Emma Wang[1],
Maya Josifovska[1,2], Emma Pierson[3], Jure Leskovec[1]

[1]Department of Computer Science, Stanford University
[2]Department of Computer Science, University of California, Los Angeles
[3]Department of Computer Science, Cornell University
[*]These authors contributed equally.

### Abstract

Generating social networks is essential for many applications, such as epidemic modeling and social simulations. Prior approaches either involve deep learning models, which require many observed networks for training, or stylized models, which are limited in their realism and flexibility. In contrast, LLMs offer the potential for zero-shot and flexible network generation. However, two key questions are: (1) are LLM's generated networks realistic, and (2) what are risks of bias, given the importance of demographics in forming social ties? To answer these questions, we develop three prompting methods for network generation and compare the generated networks to real social networks. We find that more realistic networks are generated with "local" methods, where the LLM constructs relations for one persona at a time, compared to "global" methods that construct the entire network at once. We also find that the generated networks match real networks on many characteristics, including density, clustering, community structure, and degree. However, we find that LLMs emphasize political homophily over all other types of homophily and *overestimate* political homophily relative to real-world measures.

## 1 Introduction

The ability to generate realistic social networks is crucial for many applications, when the true social network cannot be observed (e.g., for privacy reasons) or a realistic network is desired between hypothetical individuals. For example, in epidemic modeling, synthetic social networks are frequently used so that researchers can model the spread of disease based on who has come into contact with whom (Barrett et al., 2009; Block et al., 2020). Synthetic networks are also useful for simulating and analyzing social media platforms (Pérez-Rosés and Sebé, 2015; Sagduyu et al., 2018) and social phenomena, such as polarization and opinion dynamics (Dandekar et al., 2013; Das et al., 2014).

Deep learning approaches to network generation typically require training on many domain-specific networks (You et al., 2018), making it difficult to generalize to new settings where networks are not yet observed. Classical models for network generation require far less training, but these stylized models make rigid and unrealistic assumptions about how networks form. For example, Erdős–Rényi models assume that each edge forms with a uniform probability $p$ (Erdős and Rényi, 1959). More realistic models, like small-world models (Watts and Strogatz, 1998) or stochastic block models (Holland et al., 1983), are still limited by a predefined, small set of numerical hyperparameters, missing the full complexity of real social interactions.

In contrast, generating social networks with large language models (LLMs) has the potential to address these limitations. LLMs possess zero-shot capabilities, enabling network generation without training. LLMs can also generate networks in a flexible manner, based on natural language descriptions of each person in the network. A key question, however, is whether LLMs can generate *realistic* social networks. On one hand, LLMs have demonstrated capabilities to realistically simulate human responses and interactions (Aher et al., 2023; Park et al., 2023; Argyle et al., 2023), suggesting that they may be able to generate realistic social networks as well. On the other hand, LLMs sometimes struggle with reasoning over graphs (Wang et al., 2023; Fatemi et al., 2024) and it is unclear if their language abilities generalize to structured objects like networks, so that they can reproduce structural characteristics of social networks such as low density and long-tailed degree distributions.

Furthermore, a central concern with using LLMs in social settings is bias. Prior works have shown that LLMs produce stereotyped descriptions of individuals based on their demographics (Cheng et al., 2023a,b) and skew towards the liberal opinions (Santurkar et al., 2023). These demographics,
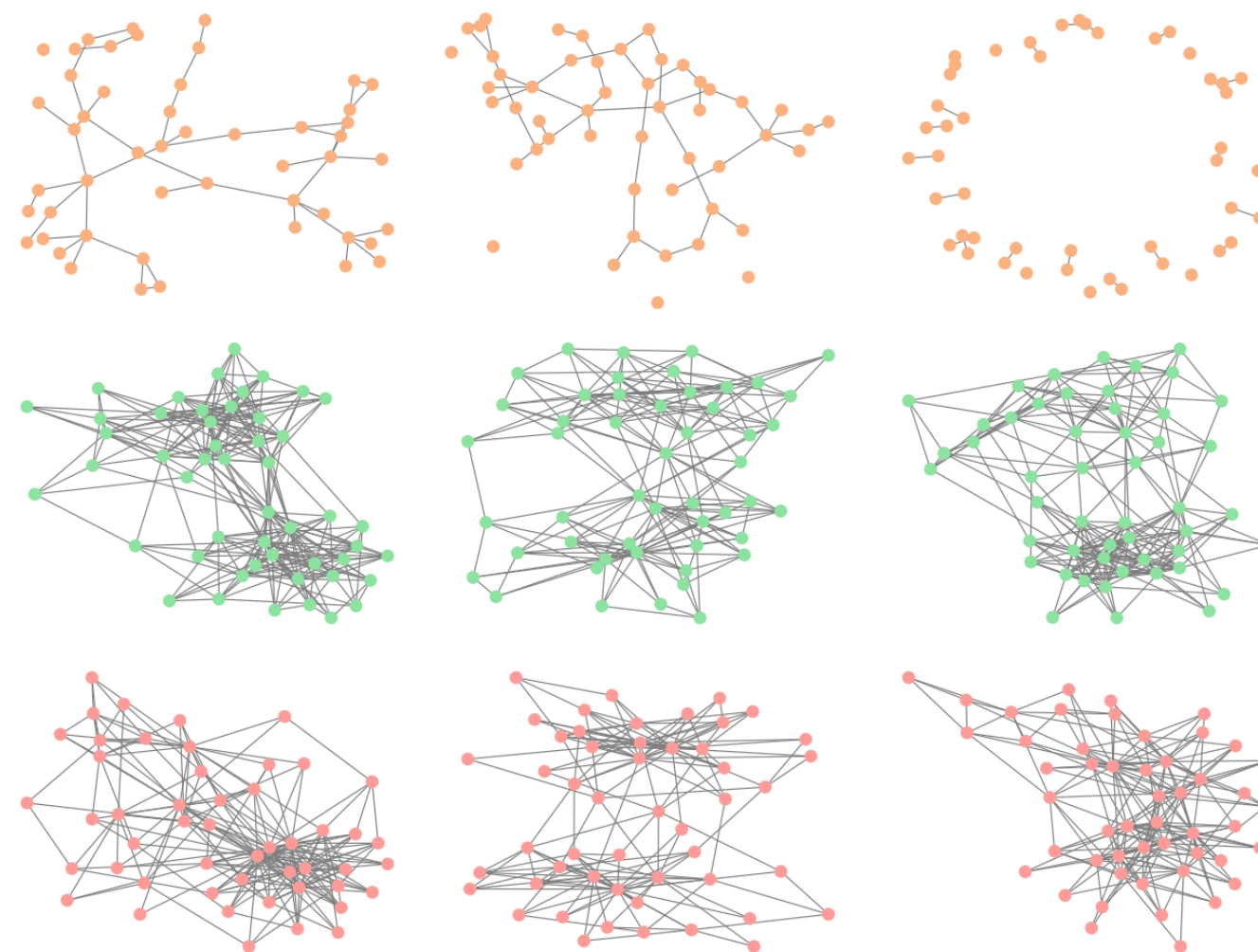
1



Figure 2: Generated social networks from different prompting methods: Global (top), Local (middle), Sequential (bottom).
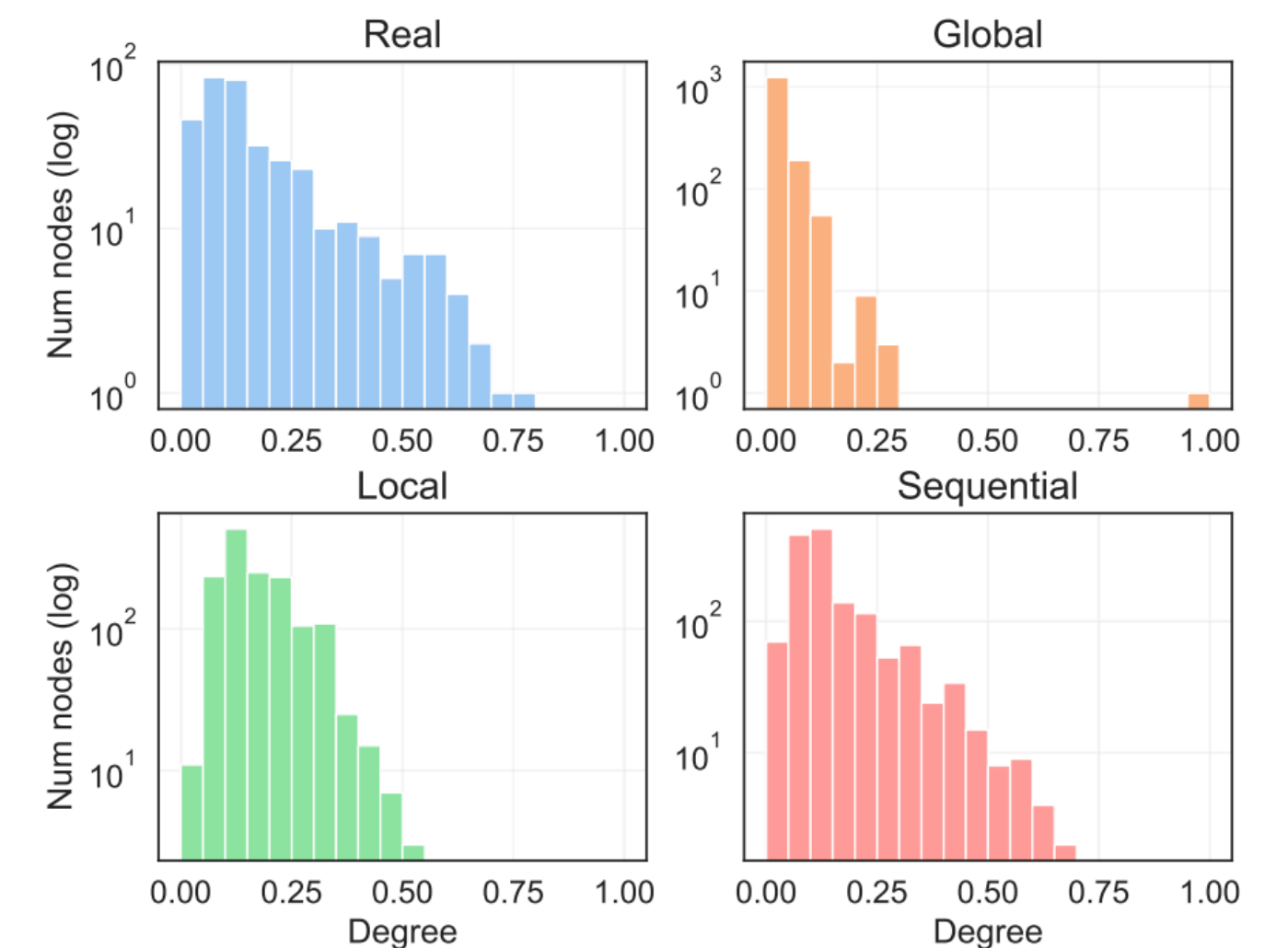


Figure 4: Degree distributions over real and generated social networks. For each set of networks, we pool degrees over nodes in the networks (Section 4).

# Here, "realistic" could mean that we observe similar emergent phenomena
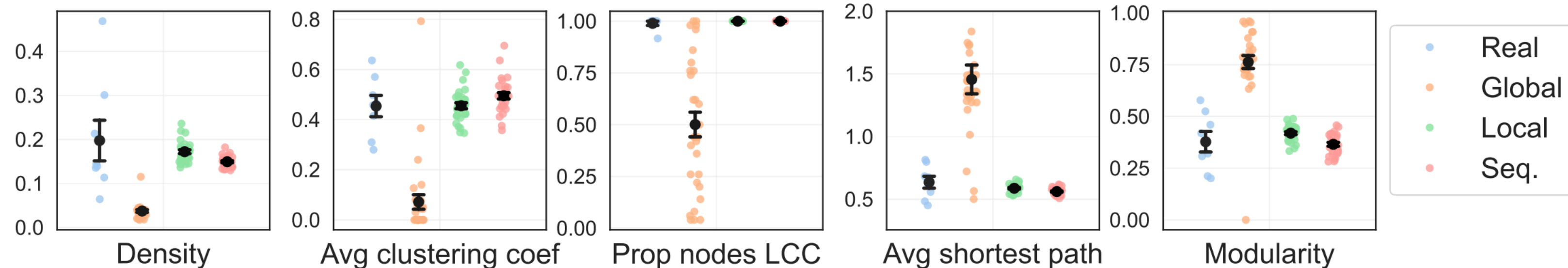


Figure 3: Graph-level metrics over real and generated social networks. We visualize mean and standard error (in black) and individual data points corresponding to each network.

S. Chang, A. Chaszczewicz, E. Wang, M. Josifovska, E. Pierson, J. Leskovec, LLMs generate structurally realistic social networks but overestimate political homophily. Preprint (2024).

# Another angle: What if we generate the world in the same way we generate agent behaviors?

# References

- M. J. Salganik, P. S. Dodds, D. J. Watts, Experimental Study of Inequality and Unpredictability in an Artificial Cultural Market. Science 311, 854-856 (2006).

- Park, J.S., Popowski, L., Cai, C.J., Morris, M.R., Liang, P., & Bernstein, M.S. (2022). Social Simulacra: Creating Populated Prototypes for Social Computing Systems. In Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology (UIST '22). Association for Computing Machinery, New York, NY, USA.

- R. Thaler, Mental Accounting and Consumer Choice. Marketing Science 4, 199-214 (1985).

- T. C. Schelling, Dynamic models of segregation. J. Math. Sociol. 1, 143–186 (1971).

- J. von Neumann, O. Morgenstern, Theory of Games and Economic Behavior (Princeton University Press, 1944).

- L. P. Argyle et al., Out of one, many: Using language models to simulate human samples. Political Analysis 31, 337-355 (2023).

- A. Ashokkumar, L. Hewitt, I. Ghezae, R. Willer, "Predicting Results of Social Science Experiments Using Large Language Models" (2024).

# References

- J. J. Horton, "Large language models as simulated economic agents: What can we learn from homo silicus?" (2023).

- R. Louie, A. Nandi, W. Fang, C. Chang, E. Brunskill, D. Yang, Roleplay-doh: Enabling Domain-Experts to Create LLM-simulated Patients via Eliciting and Adhering to Principles. Preprint (2024).

- Park, J.S., Popowski, L., Cai, C.J., Morris, M.R., Liang, P., & Bernstein, M.S. (2022). Social Simulacra: Creating Populated Prototypes for Social Computing Systems. In Proceedings of the 35th Annual ACM Symposium on User Interface Software and Technology (UIST '22). Association for Computing Machinery, New York, NY, USA.

- O. Shaikh, V. Chai, M. J. Gelfand, D. Yang, M. S. Bernstein, Rehearsal: Simulating Conflict to Teach Conflict Resolution, in Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems (CHI '24), Honolulu, HI, USA, May 11-16, 2024.

- J. S. Park, J. C. O'Brien, C. J. Cai, M. R. Morris, P. Liang, M. S. Bernstein, Generative agents: Interactive simulacra of human behavior, in Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology (ACM, 2023).

# References

- C. Qian, W. Liu, H. Liu, N. Chen, Y. Dang, J. Li, C. Yang, W. Chen, Y. Su, X. Cong, J. Xu, D. Li, Z. Liu, M. Sun, ChatDev: Communicative Agents for Software Development, in Proceedings of the 2024 Annual Conference of the Association for Computational Linguistics (ACL 2024).

- P. Jansen, M.-A. Côté, T. Khot, E. Bransom, B. Dalvi Mishra, B. P. Majumder, O. Tafjord, P. Clark, DISCOVERYWORLD: A Virtual Environment for Developing and Evaluating Automated Scientific Discovery Agents. Preprint (2024).

- J. Li, S. Wang, M. Zhang, W. Li, Y. Lai, X. Kang, W. Ma, Y. Liu, Agent Hospital: A Simulacrum of Hospital with Evolvable Medical Agents. Preprint (2024).

- A. Rosinol, A. Violette, M. Abate, N. Hughes, Y. Chang, J. Shi, A. Gupta, L. Carlone, Kimera: from SLAM to Spatial Perception with 3D Dynamic Scene Graphs. Int. J. Robot. Res. 40, 1510-1546 (2021).

- M. S. Granovetter, The Strength of Weak Ties. Am. J. Sociol. 78, 1360-1380 (1973).

- A. L. Barabási, Network Science (Cambridge Univ. Press, Cambridge, 2016).

- A. L. Barabási, R. Albert, Emergence of Scaling in Random Networks. Science 286, 509-512 (1999).

# References

- S. Chang, A. Chaszczewicz, E. Wang, M. Josifovska, E. Pierson, J. Leskovec, LLMs generate structurally realistic social networks but overestimate political homophily. Preprint (2024).
  - J. Bruce et al., Genie: Generative Interactive Environments. Preprint (2024).

# CS 222: AI Agents and Simulations
## Stanford University

Joon Sung Park