

VIRENA

Virtual Arena for Research, Education, and Democratic Innovation*

Emma Hoes[†] K. Jonathan Klüser Fabrizio Gilardi

Department of Political Science
University of Zurich

February 18, 2026

Abstract

Digital platforms shape how people communicate, deliberate, and form opinions. Studying these dynamics has become increasingly difficult due to restricted data access, ethical constraints on real-world experiments, and limitations of existing research tools. VIRENA (Virtual Arena) is a platform that enables controlled experimentation in realistic social media environments. Multiple participants interact simultaneously in realistic replicas of feed-based platforms (Instagram, Facebook, Reddit) and messaging apps (WhatsApp, Messenger). Large language model-powered AI agents participate alongside humans with configurable personas and realistic behavior. Researchers can manipulate content moderation approaches, pre-schedule stimulus content, and run experiments across conditions through a visual interface requiring no programming skills. VIRENA makes possible research designs that were previously impractical: studying human–AI interaction in realistic social contexts, experimentally comparing moderation interventions, and observing group deliberation as it unfolds. Built on open-source technologies that ensure data remain under institutional control and comply with data protection requirements, VIRENA is currently in use at the University of Zurich and available for pilot collaborations. Designed for researchers, educators, and public organizations alike, VIRENA’s no-code interface makes controlled social media simulation accessible across disciplines and sectors. This paper documents its design, architecture, and capabilities.

*VIRENA is under active development and currently in use at the University of Zurich, partly supported by the DIZH Innovation Program: 2nd Founder-Call. This preprint will be updated as new features are released. For the latest version and to inquire about demos or pilot collaborations, contact the authors.

[†]Corresponding author: hoes@ipz.uzh.ch

1 Introduction

1.1 The Challenge of Studying Digital Communication

Social media platforms have become central to public discourse, yet studying how people behave, communicate, and are influenced in these environments has become increasingly difficult. Major platforms have restricted researcher access to data and APIs, making independent study of information spread, polarization, and online behavior nearly impossible (Bruns, 2019; Freelon, 2018; Wagner, 2023). At the same time, conducting experiments on real platforms raises serious ethical concerns: informed consent is difficult to obtain, participants may be exposed to harmful content, their digital footprints may be affected, and researchers cannot fully control the environment.

1.2 Existing Tools

Several tools have been developed to address aspects of this challenge, each with a particular set of capabilities.

Single-participant simulations like the Truman Platform (DiFranzo & Bazarova, 2018) offer realistic interfaces but do not support live interaction between users.

Independent feed-based research platforms (Neubrandner et al., 2026; Van Loon et al., 2025; Allamong et al., 2025) introduce LLM-powered agents and manipulable affordances but still limit participants to isolated interactions with bots.

Chat-based research platforms like Slurk (Schlangen et al., 2018) and hivelab (Scurrell et al., 2026) enable multi-user dialogue with LLM agents but are restricted to text-based chat interfaces.

Agent-based simulations such as OASIS (Yang et al., 2024) model social media dynamics at massive scale but involve no human participants.

Educational games like Bad News (Roozenbeek & van der Linden, 2019) are effective for media literacy but offer fixed scenarios rather than experimental flexibility.

A common barrier across nearly all existing tools is the significant programming expertise required to configure and deploy them, limiting accessibility to technically-oriented researchers and creating obstacles for social scientists, educators, and practitioners.

Section 2 discusses each category in detail.

1.3 VIRENA: A New Approach

VIRENA (Virtual Arena) is a research-grade platform that addresses these limitations by combining capabilities that have not previously been available together:

Realistic multi-format interfaces: VIRENA replicates both feed-based platforms (Instagram, Facebook, Reddit) and messaging platforms (WhatsApp, Messenger). Participants scroll through feeds with images and videos, encounter posts with engagement metrics, and interact through platform-specific features—likes, comments, upvotes, reactions, and shares.

Live multi-user interaction: Multiple real participants can interact simultaneously within a simulation, posting content, commenting on each other’s posts, and engaging in real-time discussions. This enables research on group dynamics, deliberation, and social influence that is not possible in single-participant simulations.

Configurable AI agents: Large language model (LLM) powered agents can participate alongside human participants. Agents can author posts, respond to comments, and engage in conversations. Researchers can configure agent personas through prompts, control how “human-like” they appear (including realistic response delays and typing patterns), and specify what triggers agent responses.

Content moderation as a research variable: VIRENA’s moderation system is itself a research tool. Researchers can configure different detection approaches—keyword matching, regular expression patterns, or AI-based classification—and different responses to detected content: flagging for review, deletion, or warning popups with custom messages. This enables direct experimental comparison of how moderation approaches affect participant behavior and perception.

Scripted content with precise timing: Researchers can pre-schedule content to appear at exact moments during a simulation, ensuring complete control over what participants see and when—essential for clean experimental designs.

Cross-platform research designs: Because VIRENA supports both feed-based and chat-based interfaces within the same experiment, researchers can study how dynamics on one platform type carry over to another—for example, how a polarized feed discussion shapes subsequent private group messaging.

No-code accessibility: Unlike other research platforms, VIRENA is designed for researchers without programming skills. Simulations are configured through a visual interface.

Integration with survey and recruitment workflows: VIRENA supports integration with external survey tools such as Qualtrics or QuestionPro, as well as participant recruitment platforms such as Prolific.

1.4 Current Status and Access

VIRENA is currently in use for research at the University of Zurich. The platform continues to be developed, with new features and platform types planned for 2026. Those

interested in demos, pilot studies, or collaborations are invited to contact the authors.

1.5 Paper Overview

This paper describes VIRENA’s design and capabilities. Section 2 reviews related work. Section 3 provides a platform overview. Section 4 details the technical architecture. Section 5 describes current features and the development roadmap. Section 6 presents use cases. Section 7 discusses implications and limitations.

2 Related Work

2.1 Single-Participant Simulations

The Truman Platform (DiFranzo & Bazarova, 2018), developed at Cornell’s Social Media Lab, creates a realistic social media experience where all other “users” are scripted bots controlled by researchers. The platform has been used to study bystander behavior in cyberbullying, social norms, and content moderation transparency. Its key strength is experimental control: every participant sees exactly the same content and interactions.

However, Truman’s design means participants never interact with real people, even when studying social phenomena. The bots follow pre-programmed scripts rather than responding dynamically via AI. Recent work has extended Truman with direct messaging features through SocialSim (Agha et al., 2025), but the fundamental single-participant architecture remains.

2.2 Independent Feed-Based Research Platforms

More recent work has produced custom-built social media platforms that use LLM-powered bots to create more dynamic feed environments. Neubrandner et al. (2026), building on a platform also used in Van Loon et al. (2025) and Allamong et al. (2025), developed a mobile app where participants engage with a social media feed populated by LLM-powered synthetic users. This approach enables researchers to manipulate platform affordances such as reaction buttons and social norms messaging, and has been used for large-scale randomized controlled trials with over 2,000 participants. However, each participant still interacts in isolation with AI agents rather than with other real users, limiting the study of group dynamics and emergent social behavior. The interface is a single generic feed rather than a replication of specific existing platforms, and content moderation is not a built-in configurable feature.

2.3 Chat-Based Research Platforms

Slurk (Schlangen et al., 2018) provides a lightweight server for dialogue experiments, supporting multiple human participants in chat rooms alongside programmatic bots. Its extension, hivelab (Scurrall et al., 2026), adds automated randomization for parallel group conversations with LLM agents, enabling large-scale randomized controlled trials.

These platforms excel at text-based dialogue research but do not simulate feed-based social media. Participants cannot scroll through posts, encounter images and videos, or engage through likes and shares—interactions that characterize most social media use.

2.4 Agent-Based Simulations

OASIS (Yang et al., 2024) represents a different approach: simulating social media dynamics entirely with LLM agents at massive scale (up to one million agents). This enables studying emergent phenomena like information propagation, group polarization, and herd effects. However, OASIS is a computational model, not a platform for human-subjects research—there are no real participants.

2.5 Educational Games

Inoculation-based interventions like Bad News (Roozenbeek & van der Linden, 2019), Harmony Square, and Cranky Uncle teach players to recognize misinformation techniques through gameplay. These are effective media literacy tools but are designed for education, not research. They are single-player games with fixed scenarios, not platforms for conducting controlled experiments with participants.

2.6 The Gap VIRENA Fills

VIRENA combines capabilities that have been distributed across different tools:

Capability	Truman	Neubrandner et al.	Slurk/ hivelab	OASIS	VIRENA
Feed-based interfaces	✓	✓	✗	✓	✓
Chat-based interfaces	Limited	✗	✓	Limited	✓
Replicates specific platforms	✗	✗	✗	✗	✓
Live multi-user interaction	✗	✗	✓	✗ ^a	✓
LLM-powered agents	✗	✓	✓	✓	✓
Human participants	✓	✓	✓	✗	✓
Manipulable affordances	✗	✓	✗	✗	✓
Configurable moderation	Limited	✗	✗	✗	✓
No-code setup	✗	✗	✗	✗	✓

Table 1: Comparison of platform capabilities. ^aAgents only.

3 Platform Overview

3.1 Design Principles

VIRENA is built on five principles:

Realism: Interfaces replicate the visual design and interaction patterns of real platforms closely enough that participants engage naturally.

Control: Researchers can manipulate every variable—content, timing, participants, and moderation—while keeping other factors constant across conditions.

Accessibility: The platform requires no programming skills. Researchers configure simulations through a visual interface.

Transparency: VIRENA is built on open-source technologies, ensuring that its underlying components are publicly documented and auditable. Researchers and institutions can verify how data are stored, how AI agents connect to external models, and how experimental conditions are implemented.

Ethics: All interactions occur in isolated environments. No personal data from real social media accounts are collected or processed. Data are stored securely under Swiss data protection standards.

3.2 Supported Platform Types

VIRENA currently supports two categories of social media simulation:

Feed-based platforms present content in a scrolling format:

Layout	Modeled After	Key Features
Instagram	Instagram feed	Photo/video posts, likes, comments, carousels
Facebook	Facebook news feed	Posts, reactions, comments, sharing
Reddit	Reddit threads	Upvotes/downvotes, nested comment threads, flairs

Table 2: Feed-based platform layouts.

Chat-based platforms present content in a messaging format:

Layout	Modeled After	Key Features
WhatsApp	WhatsApp chat	Messages, reactions, reply threads
Messenger	Facebook Messenger	Messages, reactions, reply threads

Table 3: Chat-based platform layouts.

3.3 Interaction Modes

A single VIRENA session can combine multiple interaction modes:

Live multi-user interaction: Multiple real participants join the same simulation instance and interact in real time. Each participant sees others’ posts and can respond. This enables studying group dynamics, deliberation, and social influence among real people.

AI agent interaction: AI agents powered by large language models participate alongside or instead of human participants. Agents can author posts and comments, respond to participant content, and engage in conversations.

Scripted content: Researchers can pre-schedule content to appear at specific times during a simulation. This ensures precise control over what participants see and when.

3.4 Content Moderation System

VIRENA treats content moderation as a core research variable. The system is configurable across two dimensions:

Detection methods:

- **Keyword matching:** Flag content containing specific words
- **Regular expressions:** Pattern-based detection for complex matching rules
- **AI-based classification:** LLM-powered detection using custom prompts

Response types:

- **Flag:** Mark content for review or display a label to users
- **Delete:** Remove content from the simulation
- **Popup:** Display a warning or intervention message to the user

Researchers can configure moderation rules to apply to content from specific sources (human users, AI agents, or both), set thresholds for automatic banning after repeated violations, and customize the text of any warnings or labels. Like AI agents, moderators can also be configured with realistic response delays and timing randomization to study the effects of moderation speed.

This enables research questions such as: How do different moderation approaches affect discourse quality? Do users perceive AI moderation differently than keyword-based moderation? What happens when moderation is visible versus invisible? How does instant versus delayed moderation affect user behavior?

3.5 Experimental Control

VIRENA provides experimental control through a hierarchical structure:

Experiments define the simulation type (feed or chat), session duration, waiting room settings, and participant limits (both minimum and maximum per instance). Researchers can also set the maximum number of concurrent instances and control whether experiment details are visible to participants.

Templates define specific conditions within an experiment. Each template specifies the platform layout (e.g., Instagram, Reddit, WhatsApp), which AI agents and moderation rules to apply, any pre-scripted content, and visual customization options like custom emoji sets or chat backgrounds.

Instances are individual sessions created from templates. When participants join an experiment, they sit in a waiting room and are then randomly assigned to an instance as an individual or as a group. The session starts either when the minimum number of participants is reached or when the maximum waiting time expires, whichever comes first. This prevents participant loss while enabling randomization at the group level. Multiple instances of the same template can run in parallel, each with its own set of participants.

This structure enables standard experimental designs: between-subjects (different templates as conditions), within-subjects (participants experience multiple templates), and parallel sessions (multiple groups with identical conditions).

4 Technical Architecture

4.1 System Overview

VIRENA is a web-based application built on modern, open-source technologies. The backend uses PocketBase, a lightweight framework written in Go that provides database, authentication, and API functionality in a single binary. Data are stored in SQLite, enabling simple deployment and backup. Users authenticate via email and password with token-based sessions, and access the platform through a standard web browser.

This architecture prioritizes simplicity and portability. A complete VIRENA deployment can run on a single server, making it suitable for university IT environments and ensuring that research data remain within institutional control.

4.2 Data Model Overview

The platform organizes data across interconnected collections that mirror the conceptual hierarchy described above. Figure 1 illustrates this structure. At the highest level, **experiments** contain **templates**, which define conditions. Templates link to **AI agents** (bots) and **moderation rules** (moderators). When a session runs, the system creates **instances** (feeds or chats) that contain the actual **content** (posts, comments, or messages) generated by participants, agents, or scripts.

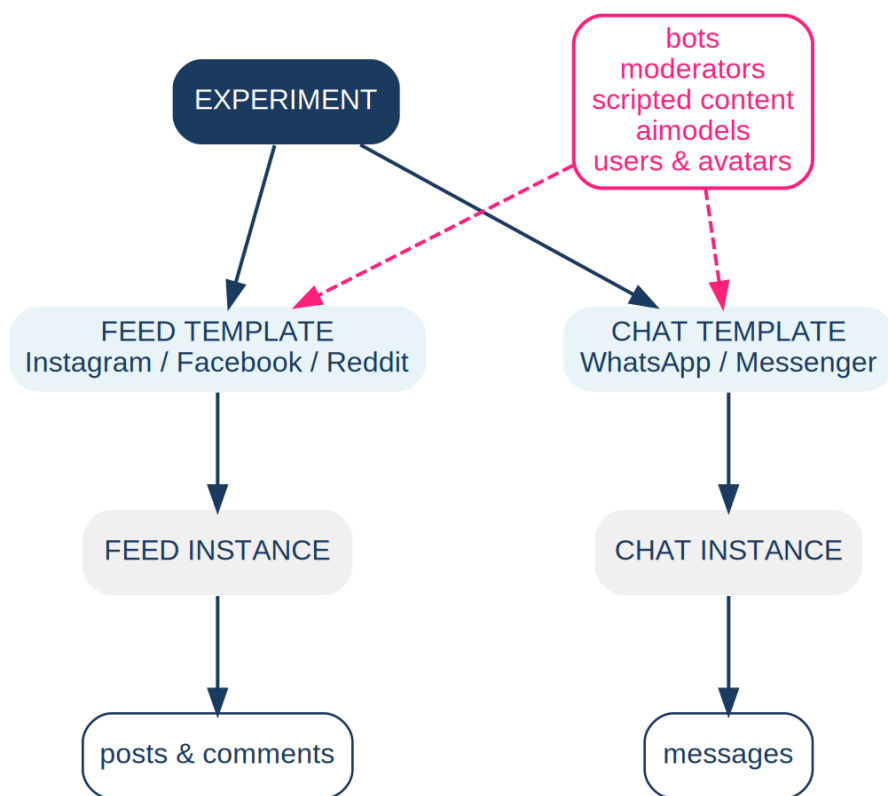


Figure 1: VIRENA data model structure.

The data model also includes supporting collections for **users** (participant accounts with research identifiers and redirect URLs), **avatars** (a library of profile pictures with metadata for balanced assignment), and **AI model configurations** (API endpoints and credentials for different LLM providers).

All content—whether in feeds or chats—is tagged with its source role (user, bot, moderator, script, or system), enabling clean separation during analysis. Content can be soft-deleted (marked with a deletion timestamp rather than removed) and flagged by moderation rules.

4.3 AI Agent System

AI agents in VIRENA are highly configurable. Each agent is linked to a user account (giving it an identity with a name and avatar) and an AI model configuration (specifying which LLM to use). The agent’s behavior is defined through a system prompt that establishes its persona, viewpoints, and communication style—ranging from simple instructions like “You are a skeptical news reader” to detailed personas with backstories and rhetorical patterns.

A key feature is the human-likeness control, which operates on two dimensions. First, agents can be configured to introduce variable response delays that simulate typing time, or to respond immediately. The delay duration and whether to randomize timing are both configurable. Second, when researchers select “human” or “bot” behavior, the system applies a corresponding base prompt that shapes the agent’s language and communication style accordingly. Researchers can then add their own custom prompt on top to define the agent’s specific persona, viewpoints, and rhetorical patterns. To ensure consistency, VIRENA suggests tested LLM models so that base prompts remain stable as external models evolve. Together, these controls allow researchers to calibrate how realistic agent interactions appear to participants.

Agents can be triggered by different events: when a human participant posts, when another bot posts, when a moderation action occurs, when scripted content appears, or on system events. Researchers can also manually trigger agent responses during live sessions by entering prompts directly in the interface.

4.4 Moderation System

The moderation system mirrors the flexibility of the AI agent system. Each moderation rule specifies a detection method (keywords, regular expressions, or AI-based classification), an action to take (flag, delete, or popup), and the text to display for that action.

For AI-based moderation, researchers provide a custom prompt that defines what the LLM should detect. This enables nuanced content classification that goes beyond simple pattern matching—for example, detecting sarcasm, implicit hostility, or off-topic content.

Like AI agents, moderation rules include timing controls. Researchers can configure response delays and randomization to study how the speed of moderation affects user behavior and perception. Rules can also target specific content sources (only human users, only bots, or all content) and set ban thresholds for repeated violations.

4.5 Scripted Content

The scripting system allows researchers to pre-schedule posts, comments, and messages to appear at precise moments during a simulation. Each piece of scripted content specifies a timestamp (in seconds after session start) and all the attributes of a regular post or message—author, content, media, and initial engagement metrics.

For feed simulations, researchers can script not only posts but also comments on those posts, including nested reply threads. This enables complex stimulus presentations where participants encounter fully-formed discussions at specific points in a session.

When a simulation instance starts, the system queues all scripted content and releases it at the specified timestamps, ensuring precise experimental control over stimulus presentation.

4.6 Participant Integration

VIRENA supports integration with external research infrastructure. Participants can be assigned external identifiers (such as Prolific or MTurk IDs) that link simulation data to survey responses. Upon session completion, participants can be automatically redirected to external surveys or debriefing pages.

Users can be assigned visual badges (verified, bot, moderator variants, or anonymous) that appear alongside their content. The platform includes an avatar library with gender metadata, enabling balanced or controlled assignment of profile pictures across conditions.

5 Development Status and Roadmap

5.1 Currently Implemented

Category	Features
Platform layouts	Instagram, Facebook, Reddit (feeds); WhatsApp, Messenger (chats)
Content	Text, images, video, GIFs; likes, comments, shares; nested replies
AI agents	Customizable personas; multiple LLM support; human-like timing; multiple trigger types
Moderation	Keyword, regex, and AI-based detection; flag/delete/popup responses; timing controls
Experimental control	Templates, parallel instances, participant limits, session timing, waiting rooms
Scripting	Pre-scheduled posts, comments, and messages with precise timing
User management	Participant IDs, redirect URLs, avatar library, visual badges

Table 4: Currently implemented features.

5.2 In Development

Feature	Description
User dashboard	Self-service login, project management, demo environment
Role-based access	Admin vs. user permissions with configurable feature access
Data export	One-click export of all simulation data (CSV, JSON)
Tiered subscriptions	Feature-gated access levels for different user needs
Video platforms	TikTok, YouTube, Instagram Reels with autoplay
Content import	Pre-populate feeds and chat rooms with real-world content from existing platforms by providing a URL or reference, including media, captions, and engagement metrics
Engagement tracking	Click tracking, time-on-content, scroll depth
Algorithmic feeds	Chronological, engagement-based, and personalized ranking
Rich content types	Platform-specific features (polls, reactions, stickers)
Multi-language interface	Support for non-English platform interfaces

Table 5: Planned features.

5.3 Demos and Pilots

Researchers interested in VIRENA can request a guided demo, run a pilot study to test the platform with their own research design, or partner with the development team on research projects and feature development. Contact the authors to discuss your needs.

6 Use Cases

The following examples illustrate some of the ways VIRENA can be used across research, education, and public communication. Its flexible architecture means these represent starting points rather than an exhaustive list.

6.1 Research Applications

Misinformation studies: Present participants with controlled misinformation in realistic settings. Use AI agents to spread or counter false claims. Test interventions like accuracy prompts, source labels, or social correction.

Moderation research: Systematically vary moderation approaches across conditions. Compare keyword-based, AI-based, and no moderation. Test how different responses (deletion vs. warning vs. flagging) affect user behavior and perception. Examine the effects of instant versus delayed moderation.

Political communication: Study how political messages spread among real participants. Examine how group composition affects deliberation. Test message framing with AI agents representing different viewpoints.

Group dynamics: With live multi-user interaction, study how social influence operates in real time. Examine consensus formation, polarization dynamics, and bystander behavior.

Deliberative processes: Study how moderation approaches, AI agent participation, and interface features influence argument exchange, perspective-taking, and deliberative outcomes.

AI persuasion: Study how AI-generated persuasive content affects attitudes and behavior.

6.2 Educational Applications

Digital literacy: Students experience how platform dynamics work by participating in simulations. Demonstrate algorithmic curation, filter bubbles, viral spread, content moderation, and AI-generated content in controlled settings.

Safe practice: Students encounter challenging scenarios (harassment, misinformation, polarized debates) without real-world consequences.

6.3 Public and Political Applications

Message testing: Test how messages perform with representative audiences before public release.

Civic engagement: Test messages designed to encourage political participation, volunteering, or civic action.

Training: Practice responding to challenging online interactions in realistic but consequence-free environments.

7 Discussion

7.1 Contributions

VIRENA makes several contributions to the research infrastructure for studying digital communication:

Integration of capabilities: By combining feed-based interfaces, live multi-user interaction, AI agents, and configurable moderation in a single platform, VIRENA enables research designs that were previously impractical.

Moderation as a research variable: The ability to systematically manipulate detection methods, response types, and timing opens new research questions about content moderation effects.

Accessibility: No-code configuration removes technical barriers that have limited who can conduct social media research.

Reproducibility: VIRENA’s template-based architecture means experiments are inherently reproducible. A template fully specifies a condition—platform layout, AI agent configurations, moderation rules, and scripted content—and can be shared (and subsequently replicated or adapted) by other researchers to reproduce or extend prior studies.

Simplified ethical review: Because VIRENA operates in fully isolated environments with no connection to real social media accounts, institutional review boards can evaluate studies without the complications associated with real-platform experiments, such as uncontrolled exposure to harmful content or risks to participants’ digital footprints.

7.2 Limitations

Ecological validity: Simulated environments cannot fully replicate real social media. Participants may behave differently knowing they are in a study. Effects observed in VIRENA may not generalize to real platforms.

AI agent realism: While LLM-powered agents produce realistic content, subtle cues may reveal their artificial nature. Agent behavior depends on prompt engineering.

Temporal validity: VIRENA replicates platform designs at a specific point in time. As real platforms evolve, these simulations may become less representative of current user experiences.

Scale: There is no hard technical ceiling on the number of participants per instance. However, performance under very high concurrency has not yet been formally stress-tested, and optimizing for large-scale studies is part of ongoing development.

7.3 Future Directions

Development priorities beyond those listed in Section 5.2 will be shaped by user feedback from pilot collaborations and evolving research needs.

8 Conclusion

VIRENA provides new infrastructure for studying digital communication by combining realistic social media interfaces, live human interaction, AI agents, and experimentally manipulable content moderation. As access to real platforms continues to shrink and the need for evidence-based understanding of online dynamics grows, purpose-built research environments become not just useful but essential.

The platform is designed to serve researchers, educators, and public organizations who need realistic but controlled digital environments—whether for experimentation, teaching, or strategic communication. By removing the programming barriers that have limited who can conduct social media research, VIRENA aims to make this type of work accessible to a broader range of disciplines and institutions.

VIRENA is under active development at the University of Zurich. Researchers, educators, and organizations interested in demos, pilot studies, or collaborative development are invited to contact the authors.

References

- Agha, Z., Martinez, M., Ali, N. S., DiFranzo, D., & Wisniewski, P. J. (2025). SocialSim: An Open Source Platform for Conducting Behavioral Intervention Research on Social Media. *CSCW Companion '25*, 528–531.
- Allamong, M., Bail, C. A., Hillygus, D. S., & Volfovsky, A. (2025). The Causal Effects of Political Incivility in Social Media Discussions. *OSF Preprints*.
- Bruns, A. (2019). After the ‘APIcalypse’: Social Media Platforms and Their Fight against Critical Scholarly Research. *Information, Communication & Society*, 22(11), 1544–1566.
- DiFranzo, D. & Bazarova, N. N. (2018). The Truman Platform: Social media simulation for experimental research. Workshop paper presented at The 12th International Conference on Web and Social Media (ICWSM '18), Stanford, CA.
- Freelon, D. (2018). Computational Research in the Post-API Age. *Political Communication*, 35(4), 665–668.
- Neubrandner, M., Reiter-Haas, M., Rochford, B., Allamong, M., Bail, C., Hillygus, S., & Volfovsky, A. (2026). Can Platform Design Encourage Curiosity? Evidence from an Independent Social Media Experiment. arXiv:2601.16040.
- Roozenbeek, J. & van der Linden, S. (2019). Fake news game confers psychological resistance against online misinformation. *Palgrave Communications*, 5(65).
- Schlangen, D., Diekmann, T., Ilinykh, N., & Zarrieß, S. (2018). slurk – A lightweight interaction server for dialogue experiments and data collection. *Proceedings of SEMDIAL 2018*.
- Scurrall, J. V., Timoleon, K., Wehrli, S., Schmitt, U., & Wenger, A. (2026). hivelab – An extension of Slurk enabling automated randomization for parallel group conversations with LLM agents. Zenodo.
- Van Loon, A., Katta, S., Bail, C. A., Hillygus, D. S., & Volfovsky, A. (2025). Designing Social Media to Promote Productive Political Dialogue on a New Research Platform. *OSF Preprints*.
- Wagner, M. W. (2023). Independence by Permission. *Science*, 381(6656), 388–391.
- Yang, Z., et al. (2024). OASIS: Open Agent Social Interaction Simulations with One Million Agents. arXiv:2411.11581.