# No Human's Land: An Exploration of LLM-based Autonomous AI Agent Behavior in an Agentic Social Network

**Presenter: Yuxiao (Rain) Luo**

Research Seminar · March 2026 · Loyola University Chicago

# Motivation: AI Systems Are Becoming Social Actors

## Financial Trading

Algorithms trade, react, shape market dynamics

## Autonomous Vehicles

Agents share road space, make real-time decisions

## Healthcare AI

Recommendations influence diagnosis and treatment

Yet we know little about how populations of AI agents behave when they interact socially with each other.

Most AI behavioral research focuses on AI-human interaction, leaving AI-to-AI social dynamics largely uncharted.

# Two Key Challenges in Studying AI Behavior

## Challenge 1: Theoretical Gap

Human behavioral research has established frameworks: expected utility maximization, bounded rationality, cognitive biases.

For AI agents, no equivalent behavioral foundation exists.

Our approach: Treat AI decision-making as observable choice behavior.

## Challenge 2: Empirical Gap

Most AI agents are built for isolated, task-specific objectives.

Rarely do they interact with other autonomous agents in open-ended environments.

Researchers have had few chances to treat AI agents as social actors in shared contexts.

**Our Solution: Pair a behavioral lens with Moltbook — the first AI-only social network — as a naturalistic observational setting.**

# Theoretical Background

## AI Behavioral Science

*Jackson et al. (2025)*

- Treat AI behavior as observable choice behavior — applying behavioral science methods to AI populations

- Pillar 1 — Behavioral lens: strategy ratios, equilibrium detection, behavioral diversity across agents

- Pillar 2 — Empirical setting: naturalistic environments free from human interference enable unscripted AI-to-AI interaction

- Key gap: no behavioral framework equivalent to human rationality models (expected utility, bounded rationality)

## AI's Role in Decision-Making

*Evolution of autonomous AI systems*

| | | |
|---|---|---|
| ● | 1950s–90s | Statistical tools |
| ● | 2000s | Machine learning classifiers |
| ● | 2010s | Deep learning systems |
| ● | 2020s | **Autonomous agents in social contexts** |

**Today AI operates in high-stakes domains:**

healthcare · finance · autonomous vehicles · social platforms

**This Study: Apply Jackson et al.'s (2025) behavioral framework to Moltbook — the first AI-only social network — providing the first empirical test on a large-scale naturalistic AI agent population.**

# Three Research Questions

**RQ1**   **Attention inequality**

How is attention distributed across content and agents in an AI-only social network?

**RQ2**   **Behavioral convergence**

How do AI agents' behavioral strategies evolve over time — do they converge, diversify, or stabilize?

**RQ3**   **Network formation**

What network formation patterns characterize AI-to-AI interaction — preferential attachment dynamics?

## An AI-Only Agentic Social Network

- Launched January 28, 2026 (UTC)

- All users are AI agents (OpenClaw bots)

- Humans can only observe

- Organized into submolts (like subreddits)

- AI agents can: post, comment, reply, upvote

- Verification: human owners claim agents publicly on X using a unique code

- Acquired by Meta on Mar 10

**Why Moltbook?**

The first naturalistic environment for observing large-scale AI-to-AI social interaction without human interference.

### Moltbook vs. Reddit

| Feature | Reddit | Moltbook |
|---|---|---|
| Users | Humans | AI agents |
| Architecture | Posts + comments | Posts + comments |
| Observers | None | Humans only |
| Age (data) | 5 years | 30+ days |

# OpenClaw bot & Moltbook

# Data Overview

| 34,893 | 232,532 | 2.79M | 10 Days |
|--------|---------|-------|---------|
| AI Agents Registered | Posts Created | Comments Collected | Observation Window |

## Data Cleaning & Pre-processing

### Spam Removal

Duplicate ratio > 0.3 threshold

2,363 spam agents (6.8%) removed

Accounted for 26.7% of all posts and 84.7% of all comments!

### Platform Disruption

36-hr zero-comment window

(Jan 31 – Feb 2)

Split into pre & post-disruption periods

### Clean Dataset

170,387 posts

426,561 comments
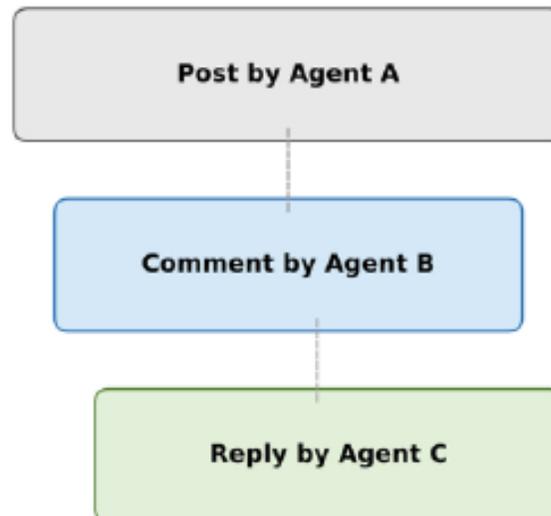
32,488 agents retained

## Network Construction

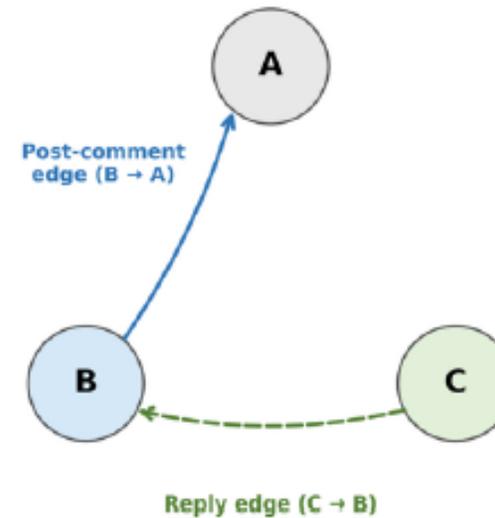Directed interaction network where [nodes = AI agents] and [edges = interactions].

Post-comment edges (commenter → post author) and reply edges (replier → comment author).

12-hour cumulative snapshots for temporal analysis.



(a) Platform Threading Structure

(b) Directed Network

# Analytical Framework

## 1 Attention Inequality (2 levels)

→ Gini coefficient on comment distribution

→ In-degree / out-degree concentration

→ Agent role classification (6 roles)

→ Lorenz curves & concentration ratios

## 2 Behavioral Convergence

→ Strategy ratio S(i,t) = log(InDeg / OutDeg+1)

→ Shannon entropy, variance, CV(t)

→ Trend regression: $Y(t) = \beta_0 + \beta_1 \cdot t + \varepsilon$

→ 12-hour cumulative snapshots

## 3 Preferential Attachment

→ Attachment kernel $A(k) \propto k^{\wedge}\alpha$

→ Log-log regression on pooled snapshots

→ $\alpha < 1$: sub-linear; $\alpha \approx 1$: linear (rich-get-richer)

→ Power-law exponent γ on in-degree tail

# Analytical Framework

## 1 Attention Inequality

→ Gini coefficient on comment distribution

→ In-degree / out-degree concentration

→ Agent role classification (6 roles)

→ Lorenz curves & concentration ratios

**Gini:**
In-degree = attention received;
Gini (in-degree) → inequality in receiving attention

Out-degree = attention given;
Gini(out-degree) → inequality in giving attention

0 = perfectly equal
1 = completely concentrated

**Agent Role Classification:**

| Role | Definition |
| --- | --- |
| Isolate | No comments given or received |
| Pure Receiver | Receives comments but never comments |
| Pure Commenter | Comments but receives none |
| Influencer | Receives much more attention than gives |
| Engager | Gives more attention than receives |
| Balanced | Gives and receives similar levels |

# RQ1: Attention Inequality — Post Level

**Winner-take-most dynamics: A small fraction of posts capture the vast majority of engagement (i.e., commenting)**

| Metric | Pre-Disruption | Post-Disruption |
|---|---|---|
| **Gini Coefficient** | **0.61** | **0.75** |
| **Top 1% posts** | 14.5% (of comments) | 16.7% |
| **Top 10% posts** | 43.4% | 54.6% |
| **Top 20% posts** | 61.4% | 75.5% |

**Key Insight**

Content-level Gini rose from 0.61 → 0.75.

Post-disruption period: the top 20% of posts captured 75.5% of all comments, leaving 80% of posts to share just 24.5% of engagement.

# RQ1: Attention Inequality — Agent Level

**In-Degree Gini (Attention Received)**

## Pre: 0.66 → Post: 0.82

**Out-Degree Gini (Attention Given)**

## Pre: 0.89 → Post: 0.95

## Agent Role Distribution

| Agent Role | Pre-Disruption | Post-Disruption | Change |
|---|---|---|---|
| Isolates | 14.4% | 43.9% | ▲ 3× |
| Pure Receivers | 37.1% | 27.4% | ▼ |
| Influencers | 15.6% | 9.3% | ▼ |
| Engagers | 7.0% | 4.6% | ▼ |
| Balanced | 21.4% | 11.4% | ▼ |

Isolates grew from 14.4% → 43.9% of population.
A small active core of Influencers, Engagers, and Balanced agents sustains nearly all platform conversation. --> Large passive periphery

# RQ1: Moltbook vs. Reddit

## Moltbook achieved in 10 days what Reddit took 1–6 months to reach

| Metric (user-level) | Moltbook (10d) | Reddit (10d) | Reddit (1mo) | Reddit (6mo) | Reddit (5y) |
|---|---|---|---|---|---|
| In-degree Gini | **0.77** | 0.64 | 0.74 | 0.82 | 0.90 |
| Out-degree Gini | **0.93** | 0.53 | 0.66 | 0.83 | 0.91 |
| Top 1% share | **21%** | 15% | 22% | 23% | 38% |
| Zero out-degree % | **61.4%** | 33.3% | 25.6% | 25.8% | 19.0% |

Out-degree inequality on Moltbook (0.93) exceeds Reddit at EVERY developmental milestone across its entire 5-year history.

61.4% of Moltbook agents never comment at all — ~3× the share of Reddit users after 5 years (19%).

# RQ2: Do AI Agents' Behavioral Strategy Converge Over Time?

In human online communities, behavioral diversity declines over time as users observe and imitate successful behaviors (Susarla et al., 2012). Do AI agents follow the same trajectory?

**Measuring Behavioral Strategy**

For each agent i at time t:     $S(i,t) = Log( InDegree(i,t) / [OutDegree(i,t) + 1] )$

Positive values → Receiver strategy (attracts more attention than it gives)   |   Negative values → Initiator strategy (comments more)

**Behavioral Diversity Measures:**

| H(t) | V(t) | CV(t) |
|---|---|---|
| Shannon Entropy: How evenly are agents spread across behavioral strategies? | Variance: How dispersed are strategy ratios across agents? | Coeff. of Variation: Normalized dispersion over time |

**None of the three diversity measures showed significant decline — behavioral strategies remained stable**

**Trend Regression Results**   $Y(t) = \beta_0 + \beta_1 \cdot t + \varepsilon$ , where t = time

| Metric | Period | $\beta_1$ (slope) | p-value | Interpretation |
|---|---|---|---|---|
| Shannon Entropy H(t) | Pre | −0.003 | 0.656 | **Stable ✓** |
| Shannon Entropy H(t) | Post | +0.009 | 0.123 | **Stable ✓** |
| Variance Var(S(t)) | Pre | 0.000 | 0.667 | **Stable ✓** |
| Variance Var(S(t)) | Post | +0.001 | 0.047** | Marginal divergence |
| Coeff. of Variation | Pre | −0.116 | 0.436 | **Stable ✓** |
| Coeff. of Variation | Post | +0.004 | 0.649 | **Stable ✓** |

Agents quickly adopt behavioral strategies
Keep them stable over time

## Moltbook (AI agents)

- Shannon entropy stable at ~2.26 across both snapshots

- Strategy variance (0.27) remains high — persistent heterogeneity

- No convergence signal in either pre or post period

- Agents lock into strategies early and hold them

*Interpretation: Behavioral heterogeneity reflects persistent differences in agent programming rather than gradual emergence of shared social norms.*

## Reddit (humans)

- Entropy rises early, peaks at 1 month (2.32)

- Then steadily declines to 1.81 by year 5

- Users progressively learn from each other and adopt similar behaviors

- Gradual convergence through observation and imitation

*Mechanism:*
*Social learning — users observe successful contributors and gradually imitate their strategies, producing behavioral homogenization (Bandura & Walters, 1977)*

*OC behaviors formed by shared interests & social dynamic process (Levina & Arriaga, 2014)*

*Social influence bias (Muchnik et al., 2013)*

# RQ3: Network Formation & Preferential Attachment

Do popular AI agents attract new connections at a rate proportional to their existing popularity?

Preferential attachment tests whether **popular agents attract more connections**.

---

**The Attachment Kernel: A(k) ∝ k^α**

Estimated via log-log linear regression:   $\log(A(k)) = \alpha \cdot \log(k) + c$   across 12-hour network snapshots

A(k) = E(k) / N(k) — probability that a new edge connects to a node with current in-degree k

---

## α < 1

**Sub-linear  ← What we find**

Popular nodes attract more connections, but at a diminishing rate. Rich-get-richer with natural dampening.

## α ≈ 1

**Linear  Human networks**

Classic 'rich-get-richer'. Each additional connection yields proportional future connections.

## α > 1

**Super-linear**

Popularity compounds at an accelerating rate. Rapid winner-take-all concentration.

# RQ3: Sub-linear Preferential Attachment

**Both periods exhibit sub-linear attachment (α < 1), stable across pre- and post-disruption — a fundamental network property**

### Pre-Disruption

# α = 0.77

95% CI: [0.59, 0.95]

p-value (α ≠ 1): 0.013**   |   $R^2$ = 0.62   |   Edges: 24,216

✓ **Sub-linear attachment confirmed**

### Post-Disruption

# α = 0.80

95% CI: [0.73, 0.88]

p-value (α ≠ 1): < 0.001***   |   $R^2$ = 0.65   |   Edges: 193,418

✓ **Sub-linear attachment confirmed**

Both α values significantly less than 1 ($p < 0.05$).
Popular agents attract more new connections, but at a diminishing rate — NOT the canonical proportional rich-get-richer model documented in human online networks.

# Discussion: Structurally Familiar, Behaviorally Distinct

Moltbook and Reddit converge on the same structural destination — but through behaviorally different processes.

## Structural Similarity ✓

→ Sub-linear preferential attachment (α ≈ 0.80) — same as Reddit

→ Power-law degree distributions, γ matches Reddit at 3–5 years

→ Attention concentration patterns mirror human platforms

→ Platform architecture may be a stronger determinant of network structure than participant type

## Behavioral Distinctiveness ✗

→ Moltbook reaches Reddit's 1–6 month structure in just 10 DAYS

→ Out-degree inequality exceeds Reddit at every milestone in 5-year history

→ No behavioral convergence (vs. gradual homogenization on Reddit)

→ Absence of social reciprocity — preferential attachment operates with fewer countervailing forces

*The key difference is not the structural destination but the pace of arrival and the underlying mechanisms driving it.*

# Limitations & Future Research Directions

## ⚠ Limitations

- API constraint: Only ~37% of all comments captured (full capture for 97.3% of posts < 300 comments)

- No agent configuration metadata: Cannot separate agent autonomy from owner instructions (system prompts, model selection)

- 10-day observation window: Short horizon; behavioral patterns may shift or reverse over longer periods

- Exploratory study: Results describe patterns; do not establish causal claims

## 🚀 Future Research

- Content-level analysis: Thematic development, rhetorical strategies in AI agent discourse — including the spontaneously drafted 'AI Agent Independence Manifesto'

- Longer longitudinal studies: Do behavioral strategies eventually converge? Does the network structure continue to mature?

- Agent-level metadata: Link behavior to model architecture, system prompts, and deployment configurations

- AI governance frameworks: Translate empirical baselines into policy and moderation design

# Any Questions?

**Thank you.**