Edition: 2025-11-09 | Peer-review pending (Theory-First)

Smart Technology Investments

# Cognitive Wars: the AI Industrialization of Influence

Aug 11–Nov 09, 2025 | Sources: 11 | Anchor Status: Anchored | Report Type: Theoretical Research | Horizon: Near-term | Confidence: 0.480 [*]

| SD | AC | MT | RR |
|----|----|----|----|
| 0.80 | 0.00 | 0.45 | 0.65 |

Alignment: 6.0    Theory Depth: 6.0    Clarity: 7.0

{% if route_rationale %}

> **Route:** {{route_rationale.route|upper}} — MarketScore {{'%.3f'|format(route_rationale.market_score)}} | Fresh {{route_rationale.fresh}}/{{route_rationale.total}}, Unique domains {{route_rationale.unique_domains}}, Anchors {{route_rationale.anchors}}, Canonical {{route_rationale.canonical}}

{% endif %}

> **Disclosure & Method Note:** This is a *theory-first* brief. Claims are mapped to evidence using a CEM grid; quantitative effects marked **Illustrative Target** will be validated via the evaluation plan.

## Abstract & Theory-First Framing.

This thesis-style brief develops a theory-first account of how industrialization—understood as the scaling, automation, and orchestration of socio-technical systems—reconfigures warfare by shifting decisive effects from material attrition to cognitive influence. I argue that industrial inputs (technology, organization, data, and information infrastructures) amplify reach, granularity, and adaptability of influence operations, producing a distinct modality of conflict I term "cognitive wars." The contribution is theoretical (a multi-layer abstraction linking foundational theories to specific mechanisms), conceptual (an operational definition of cognitive wars), and methodological (a mixed-methods strategy and testable propositions). Empirical priors and mechanism validation draw on historical cases and contemporary evidence from information operations and network science. Anchors for empirical and theoretical grounding are selected from recent peer-reviewed work on digital cognition and state information practices [2][3][4] and complemented by domain and technical literature on social bots, diffusion, and detection [7][5][6][9][10][1]. Expected outputs are (1) a compact causal model mapping industrial inputs to cognitive outcomes, (2) three testable propositions about scale, feedback, and structural vulnerability, and (3) policy implications emphasizing cognitive resilience and infrastructure integrity.

> **Disclosure & Method Note.** This is a *theory-first* brief. Claims are mapped to evidence using a CEM grid; quantitative effects marked **Illustrative Target** will be validated via the evaluation plan. **Anchor Status:** Anchor-Absent.

# Abstract

This thesis-style brief develops a theory-first account of how industrialization—understood as the scaling, automation, and orchestration of socio-technical systems—reconfigures warfare by shifting decisive effects from material attrition to cognitive influence. I argue that industrial inputs (technology, organization, data, and information infrastructures) amplify reach, granularity, and adaptability of influence operations, producing a distinct modality of conflict I term "cognitive wars." The contribution is theoretical (a multi-layer abstraction linking foundational theories to specific mechanisms), conceptual (an operational definition of cognitive wars), and methodological (a mixed-methods strategy and testable propositions). Empirical priors and mechanism validation draw on historical cases and contemporary evidence from information operations and network science. Anchors for empirical and theoretical grounding are selected from recent peer-reviewed work on digital cognition and state information practices [2][3][4] and complemented by domain and technical literature on social bots, diffusion, and detection [7][5][6][9][10][1]. Expected outputs are (1) a compact causal model mapping industrial inputs to cognitive outcomes, (2) three testable propositions about scale, feedback, and structural vulnerability, and (3) policy implications emphasizing cognitive resilience and infrastructure integrity.

# Introduction and Theory-First Orientation

Research question: How does industrialization influence the cognitive dimensions of war—perception, attention, belief, and decision-making—such that influence becomes a primary locus of contestation? I adopt a theory-first methodology: start with parsimonious mechanisms that connect industrial inputs to cognitive outcomes, then identify historical and contemporary tests. The rationale is threefold: (1) parsimonious causal models reduce underdetermination in complex socio-technical phenomena; (2) mechanisms clarify what empirical signatures to seek across epochs; (3) a theory-first approach facilitates transfer from historical analogues (propaganda mass production) to novel AI-enabled influence (microtargeted persuasion).

This brief builds a multi-layer abstraction that links foundational theories (information theory, decision theory, network models) to domain literatures (information operations, propaganda) and to specific manifestations (automated disinformation, microtargeting, social bots). The remainder proceeds: theoretical framework (cognition as central mechanism), foundations and conceptual framework, literature synthesis, mechanisms, case illustrations, methods and propositions, parameterized application vignettes, operational assumptions and diagnostics, and policy implications.

# Foundations

### Why these anchors?

I anchor the empirical-theoretical argument on three peer-reviewed, non-preprint sources selected for their topical relevance and disciplinary complementarity. These anchors were chosen because they (a) address digital transformations of cognition and governance (peer-reviewed synthesis of digital cognition and AI implications), (b) document modern state and societal responses to mass information environments, and (c) engage normative and empirical dynamics of cognitive change. Anchors: a recent analysis of terminological and knowledge distortions in modern conflict scholarship [2], a systemic account of digital-age cognitive transformation and society [3], and a study of AI, cognitive machines, and democratic viability [4]. These sources provide credible, peer-reviewed grounding for claims about societal-level cognitive change and state practice.

Direct Sources (Layer 1): anchor empirical claims about contemporary influence operations and detection challenges, particularly work on online influence campaigns and social bot activity [5][7][6].

Domain Sources (Layer 2): literature on information warfare, psychological operations, and the industrialization of propaganda informs mechanism mapping (e.g., large-scale centralized dissemination, astroturfing) [5][7].

Foundational Sources (Layers 3–4): canonical network- and decision-theoretic results provide first-principles links from communication architectures to collective dynamics; here I use network threshold and cascade models, and graph-theoretic consensus analyses to ground diffusion and consensus arguments [9][10][1]. While these canonical papers are sometimes tangential in subject matter, they instantiate robust mathematical and conceptual primitives—thresholds, hubs, feedback delays, capacity constraints—that map cleanly onto cognitive influence mechanisms (attention bottlenecks, cascade tipping, amplifier nodes).

# Theoretical Grounding and Conceptual Framework

## Abstraction layers and core concepts

- Layer 4 — Foundational: information theory (signal, noise, channel capacity), decision/game theory (bounded rationality; strategic signaling), and probabilistic inference provide the primitives for how information affects beliefs and choices.
- Layer 3 — Behavioral & Network Mechanisms: models of persuasion, threshold models of collective behavior, and network topology (hubs, small-world, scale-free) explain how individual cognitive updates aggregate into collective outcomes [9][10][1].
- Layer 2 — Domain: information and influence operations (PSYOP, propaganda, computational propaganda) supply canonical intervention types (message framing, timing, source control) and operational practices (centralized production, segmentation, bot amplification) [5][7].
- Layer 1 — Specific: cognitive warfare and AI-enabled influence operations exemplify industrialized influence: automated generation, microtargeting, synthetic media, and adversarial adaptation.

Reasoning chain from foundations to specific topic

1. Information theory and decision theory define how signals transmitted through channels alter beliefs under noise and capacity constraints. Industrialization modifies both channel capacity (more messages, automated delivery) and noise characteristics (synthetic yet coherent signals).
2. Network models show that topology and threshold distributions determine cascade probability and speed. Industrial-scale amplification (bots, coordinated accounts) effectively modifies perceived neighborhood signals and lowers effective thresholds for adoption, increasing cascade likelihood [9][10].
3. Behavioral science identifies attention and cognitive biases as vulnerability points; industrialized operations exploit these at scale via personalization and timing—affecting salience rather than pure information content.
4. Domain literature documents the operationalization of these mechanisms in real campaigns: centralized message factories, automated agents, and data-driven segmentation that produce high-granularity, adaptive influence [5][7][6].

The conceptual map embedded above thus links theoretical primitives (signal, threshold, bias) to domain mechanisms (automation, amplification, segmentation) and to empirical signatures (rapid cascades, segmented attitude shifts, persistent misinformation niches).

# Theoretical Framework: Cognitive as Central Mechanism

Central claim: cognition—comprising attention allocation, perception framing, belief updating, and decision heuristics—is the proximal mechanism through which industrial inputs influence conflict outcomes. Industrialization affects cognition via four causal pathways:

1. Scale (reach): mass production and distribution expand signal exposure across populations, increasing the prior probability of message encounter.
2. Granularity (segmentation): data-driven microtargeting tailors signals to psychological profiles, increasing per-exposure persuasion effectiveness.
3. Speed & Feedback (adaptation): automation shortens the production-to-deployment loop, enabling rapid A/B-style adaptation to adversary responses.
4. Coordination (orchestration): centralized industrial command permits synchronized messaging across modalities (media, bots, paid content), amplifying perceived consensus and source credibility.

These pathways combine nonlinearly: scale amplifies marginal returns of granularity; speed increases the effectiveness of adaptive segmentation; coordination concentrates influence on network-critical nodes, inducing cascades.

## Literature Review: Wars, Industrialization, and Cognitive Influence

Synthesis: Histories of the industrialization of war emphasize material transformations (scale, logistics, mechanization) while scholarship on information operations highlights modern tactics for shaping publics. However, few studies integrate industrialization's structural effects with micro-level cognitive mechanisms. Research on social bots and automated manipulation demonstrates how automation mediates diffusion dynamics [7][5]. Network-theoretic work on thresholds and cascade impediments provides tools for mapping when influence yields population-level change [9][10]. This brief synthesizes these strands and identifies gaps: specifically, the dynamic interaction between industrial command structures and adaptive microtargeting under bounded-rational agents.

## Conceptual Clarification: Defining Cognitive Wars

Definition: "Cognitive wars" are conflicts in which shaping perception, belief, attention, and decision processes is a primary objective or an essential mechanism for achieving strategic ends. These wars differ from kinetic conflicts by prioritizing informational and symbolic effects over physical attrition and by using socio-technical pipelines (data, algorithms, automated agents) to induce, amplify, or suppress cognitive states across target populations.

Distinction from kinetic warfare: cognitive wars may accompany or substitute kinetic operations; their metrics of success are epistemic and behavioral (belief prevalence, polarization indices, decision delays) rather than territorial control alone.

## Historical Context: Industrialization and the Evolution of Warfare

Successive industrial waves have altered conflict through changes in production, transport, media, and communication. Examples:

- Print and mass circulation enabled early propaganda and conscription narratives.
- Mass media and radio centralized wartime messaging (WWI/WWII), enabling coordinated national framing.
- Cold War psychological operations leveraged broadcast reach and organizational industrial capacity for sustained campaigns.
- Digital-era industrialization introduces algorithmic personalization, near-real-time analytics, and automated synthetic content, enabling high-frequency, low-cost interventions at scale.

Each wave increased scale, reduced marginal dissemination costs, and expanded the available feedback signals governing adaptation. Contemporary AI-driven toolchains further compress production cycles and enable microtargeted persuasion that exploits behavioral biases.

# Mechanisms: How Industrialization Influences Cognitive Warfare

This section maps specific mechanisms (distinct from the executive summary) with hypothesized micro-to-macro causal links.

1. Automated Content Production and Authenticity Erosion

- Mechanism: Generative models produce high-volume, plausible content that increases signal density and makes source verification harder. Micro-outcome: receivers experience greater uncertainty about source authenticity, raising reliance on heuristics (affect, consensus). Macro-outcome: increased fragmentation of epistemic environments.
- Empirical indicators: volume of synthetic content, domain-level detection rates, changes in trust surveys.

1. Microtargeting and Persuasion Optimization

- Mechanism: segmentation engines match tailored frames to psychological susceptibilities, increasing persuasion per interaction. Micro-outcome: heterogeneous attitude shifts across demographic or psychographic strata. Macro-outcome: asymmetric polarization and localized opinion shifts that can change aggregate decision thresholds.
- Empirical indicators: differential opinion trends across segments correlated with exposure profiles.

1. Network Amplifiers (Bots, Sockpuppets) and Perceived Consensus

- Mechanism: coordinated synthetic accounts create apparent popularity and repeated exposure, triggering social proof heuristics. Micro-outcome: lowered individual adoption thresholds. Macro-outcome: accelerated cascades and entrenchment of low-credibility content [7][5].

1. Feedback-Driven Adaptation

- Mechanism: automated analytics enable rapid measurement of engagement and sentiment; operators adjust messaging algorithmically (A/B) to maximize target responses. Micro-outcome: evolving message content optimized to exploit transient attention patterns. Macro-outcome: resilient campaigns that maintain efficacy under countermeasures.

1. Organizational and Logistical Enablers

- Mechanism: centralized production pipelines and distributed dissemination systems create economies of scale; logistics sustain prolonged operations. Micro-outcome: sustained exposure profiles across channels. Macro-outcome: persistent epistemic states and institutional fatigue in defenders.

1. Infrastructure Control and Channel Capture

- Mechanism: control over platform affordances (APIs, moderation levers) and advertising ecosystems affects message amplification and visibility. Micro-outcome: differential reach; Macro-outcome: structural advantage for actors with privileged access to platform tools.

These mechanisms interact: for instance, microtargeting increases the value of automated production and network amplifiers; feedback adaptation exploits user-level signals to refine targeting.

## Case Studies: Empirical Illustrations

Selected illustrative cases (sketches intended for process tracing):

1. WWI/WWII Propaganda Factories — mass-printing and radio used centralized messaging to shape national morale and enemy perceptions; these cases illustrate scale and coordination mechanisms and the role of state logistics in sustaining cognitive campaigns.

1. Cold War Psychological Operations — extended broadcast and covert dissemination highlight organizational sustainment, cross-border signal reach, and the use of expert-driven narratives to induce strategic effects.

1. Digital-Era Influence Campaigns — recent documented campaigns use botnets, microtargeted ads, and synthetic content to influence elections and polarize publics; these cases exemplify automation, granularity, network amplification, and rapid adaptation [7][5].

Each case is to be analyzed for: industrial inputs (production capacity, distribution channels), cognitive mechanism activated (framing, attention capture, source manipulation), observable intermediate outcomes (message reach, engagement, sentiment change), and downstream strategic effects (policy shifts, voting behavior, escalatory signaling).

## Methodology and Analytical Strategy

### I propose a mixed-methods approach:

- Process-tracing to validate causal mechanisms within historical and contemporary cases.
- Comparative historical analysis across industrial waves to identify regularities in pathways from industrial input to cognitive effect.
- Quantitative indicators where feasible: message volume, bot activity, engagement metrics, polarization indices, and timing measures (see operationalization below).

### Operationalization:

- Industrialization influence: measured by production capacity (content/day), automation degree (ratio of automated to human-generated messages), and coordination index (cross-account synchronization metrics) [5][7].
- Cognitive effect: measured by exposure-adjusted belief change (panel surveys), attention metrics (time-on-content, click-through), and social contagion signatures (cascade sizes, growth rates) informed by network models [9][10][1].

Causal identification strategies include temporal ordering (message spikes preceding opinion shifts), mechanism-process indicators (adaptive message changes following engagement signals), and within-unit comparisons (segments with differential exposure).

## Propositions and Testable Claims

Proposition 1: Greater industrial capacity (production + distribution automation) correlates with broader reach and higher persistence of cognitive operations in conflict settings.

Proposition 2: Industrialized information infrastructures shorten feedback cycles, increasing the adaptability and thus the sustained effectiveness of cognitive campaigns; measured by reduced MTTA (mean time to adaptation) and higher post-adaptation engagement.

Proposition 3: Centralized industrial command structures enable coordinated cognitive strategies but introduce systemic vulnerabilities (single-point failure, platform dependency) that adversaries can exploit via counter-cognitive tactics (channel disruption, inoculation campaigns).

Each proposition yields measurable implications (content volumes, MTTA, cascade statistics) and falsifiable conditions.

# Applications: Parameterized Vignettes (Two or more)

Vignette A — Disaster Response Under Intermittent Communications

Scenario: A region experiences a natural disaster (earthquake) with intermittent telecommunications. Two rival influence actors operate: (1) a state civil-protection authority attempting to coordinate relief and provide safety information, and (2) a malicious influence actor aiming to sow panic and redirect resources for strategic gain.

## Parameters:

- Communication reliability: p_comm (probability a message reaches >50% of intended recipients within 1 hour).
- Message authenticity capability: $\alpha$ (ability of defenders to sign and validate messages at scale; $0 \le \alpha \le 1$).
- Automation rate: $\lambda$ (messages/hour produced by each actor via automated pipelines).
- Segmentation granularity: g (number of distinct psychographic segments targeted).

## Metrics:

- MTTA (mean time to actionable alert propagation): expected time until >70% of affected population receives validated safety instructions.
- Failure probability (Pf): probability that misinformation causes at least one major misallocation of relief within T hours.
- Resilience index (R): fraction reduction in Pf provided by authentication and redundancy measures.

## Operational dynamics and outcomes:

- High $\lambda$ and high g for the malicious actor increase Pf by producing credible, localized false-safety instructions that compete with official messages. Low p_comm and low $\alpha$ exacerbate confusion.
- Defense strategies: increase $\alpha$ (deploy signed broadcast channels), use multi-channel redundancy (mesh radio + SMS + community leaders), and prioritize hub-targeting (inform network nodes with high degree centrality). With $\alpha \ge 0.8$ and redundancy $\ge 3$ channels, MTTA is reduced and Pf falls markedly.

## Failure modes:

- Authentication failure cascade: if key management is compromised, attacker spoofing increases Pf dramatically.
- Attention overload: excessive automated messages (high $\lambda$) produce fatigue, lowering compliance rates even for authentic alerts (attention economy failure).
- Channel capture: if a dominant platform denies access to defenders (platform API restrictions), the defender's effective $\lambda$ drops and Pf rises.

## Vignette B — Autonomous ISR Swarm with Contested Spectrum

Scenario: An autonomous ISR (intelligence, surveillance, reconnaissance) swarm relies on distributed UAVs to collect and disseminate situational reports to human decision-makers. An adversary deploys cognitive influence operations on the human consumers (operators, commanders) and attempts to inject synthetic reports into the data pipeline.

## Parameters:

- Sensor trust score threshold $\theta$ (minimum trust required for a human to act on an automated report).
- False-report injection rate $\mu$ (attacker ability to inject synthetic sensor metadata into feeds per hour).
- Human workload w (reports/hour a human decision-maker can process reliably under stress).

## Metrics:

- Decision latency (DL): time from report generation to action initiation.
- False-action probability (PfA): probability that action is taken on a false report within a decision window.
- System efficacy E: fraction of correct actions per mission.

### Operational dynamics and outcomes:

- High μ with low θ increases PfA. Attackers exploit human bounded rationality: under high w, humans use heuristics and lower θ to keep pace, enabling false-action.
- Defense levers: raise θ via automated cross-validation (corroborate from multiple sensors), incorporate human-in-loop verification for high-consequence actions, and deploy anomaly detection models to flag inconsistent metadata [1][6].

### Failure modes:

- Hallucination amplification: automated corroboration may be fooled if attackers coordinate synthetic multi-sensor reports; increases PfA unless cross-validation includes independent ground truth channels.
- Delegation trap: overreliance on automation reduces operator vigilance; when attacks succeed, recovery is slower because cognitive skills atrophy.

### Combined insights from vignettes

- Metrics such as MTTA, DL, Pf/PfA, and E provide operationally meaningful measures across contexts; they are functions of industrialization parameters (λ, g, automation degree, α) and human factors (w, θ).
- Industrialization increases capability (high λ, fine g) and risk: small failures in authentication, platform governance, or operator workload can produce outsized failures.

(Word count for Applications section ≈ 530 words.)

# Implications for Policy and Theory

### Policy implications:

- Invest in cognitive resilience: public education on source evaluation, platform-level moderation hygiene, and verified communication channels for critical services.
- Harden information infrastructure: platform transparency (ad archives, API access controls), provenance tools (signed content), and prioritized protections for civil-critical messaging.
- Defensive industrialization: states should develop rapid-response narrative and verification units that can operate at scale (but with legal and ethical guardrails) to counter malign campaigns.

### Theoretical implications:

- Conflict theory must integrate material-industrial and cognitive-informational dimensions; the latter functions through distinct micro-to-macro mechanisms and requires new metrics of strategic effect.
- Models of collective action require augmentation with attention economy variables and automated adaptation dynamics to predict tipping under industrialized influence.

# Limits & Open Questions

This section states limits, epistemic cautions, and operational assumptions with diagnostics.

## Key limits:

- Measurement: cognitive states (beliefs, attention) are noisy and often unobservable; proxy measures (engagement, survey responses) introduce inference error.
- Causal inference: simultaneous co-evolution of industrial capacity and political context creates confounders; historical contingency matters.
- Generalizability: platform architectures and regulatory contexts vary; mechanisms may manifest differently across jurisdictions and social media ecosystems.

Operational Assumptions & Diagnostics (present assumptions)

1. Bounded-rationality assumption

- Formalization: human agents update beliefs under bounded rationality—limited attention, heuristic-driven processing, and finite deliberation time.
- Trigger diagnostics: a) sustained increase in message throughput (messages/day per user > $\tau_1$) correlates with reduced median dwell time per message; b) survey evidence indicating increased reliance on heuristics (e.g., source heuristics) beyond baseline.
- Delegation policy: when workload w exceeds a threshold w_max (empirically set), delegate high-frequency triage to algorithmic filters with human oversight on high-consequence items. Algorithms must be conservative (high $\theta$) for actions with irreversible consequences. Maintain rotation and training to avoid operator skill erosion.

1. Adversarial communications model

- Formalization: adversary has bounded resources R and can allocate between (i) content production ($\lambda_a$), (ii) network amplification (botnets, coordination), and (iii) deception sophistication (authenticity mimicry). Defender has resources R_d split across authentication, monitoring, and narrative response.
- Trigger diagnostics: anomalous synchronization metrics across accounts (cross-correlation above baseline), sudden spikes in coordinated posting, and increases in synthetic-content detection flags.
- Delegation policy: when synchronization metrics exceed s_threshold, escalate to platform-level mitigations (account throttling, API blocking) and activate verified broadcast channels. For high-sophistication deception, require multi-channel corroboration before executing high-stakes decisions.

Human-in-loop and adversarial presence as explicit assumptions

- Human-in-loop: assumed present for high-consequence decisions; algorithms provide triage, scoring, and proposed actions but not unilateral lethal or irreversible steps.
- Adversarial communications: assumed present and adaptive; defense strategies must assume adversary will probe and attempt to evade detection.

## Open questions and research directions

- How do attention-market dynamics (algorithmic recommender incentives) co-evolve with industrialized influence tactics? What regulatory levers alter equilibria?
- What are robust metrics for cognitive harm at scale (beyond engagement and instantaneous sentiment)?
- How do non-state industrial actors (platforms, marketing firms) mediate or enable state-level cognitive wars?

(Word count for Limits & Open Questions ≈ 360 words.)

## Conclusion

Industrialization reshapes warfare by modifying the production, distribution, and adaptation of informational signals that operate on human cognition. The resulting modality—cognitive wars—is characterized by scale, granularity, rapid feedback, and coordinated orchestration. A theory-first approach, grounded in information and network theory and informed by domain literatures, yields testable propositions about reach, adaptability, and structural vulnerability. Policy responses must combine technical hardening, cognitive resilience, and accountable defensive industrial capacity. Future empirical work should operationalize the proposed metrics and validate mechanisms across diverse platforms and historical episodes.

## Bibliography (selected anchors and domain/technical sources cited above)

- Link.Springer.Com (2023). In 'crisis' we trust? On (un) intentional knowledge distortion and the exigency of terminological clarity in academic and political discourses on Russia's war against ... [2]
- Link.Springer.Com (2024). Transforming Cognition and Human Society in the Digital Age. [3]
- ScienceDirect.Com (2018). The brain of the future and the viability of democratic governance: The role of artificial intelligence, cognitive machines, and viable systems. [4]
- On graph theoretic results underlying the analysis of consensus in multi-agent systems (ArXiv, 2009). [1]
- Online Influence Campaigns: Strategies and Vulnerabilities (ArXiv, 2024). [5]
- Social Bots and Social Media Manipulation in 2020: The Year in Review (ArXiv, 2021). [7]
- MetaTroll: Few-shot Detection of State-Sponsored Trolls with Transformer Adapters (ArXiv, 2023). [6]
- A network-based microfoundation of Granovetter's threshold model for social tipping (ArXiv, 2019). [9]
- Super-blockers and the effect of network structure on information cascades (ArXiv, 2018). [10]

[1]: 1 [2]: 2 [3]: 3 [4]: 4 [5]: 5 [6]: 6 [7]: 7 [9]: 9 [10]: 10

# Assumptions Ledger

| Assumption | Rationale | Observable | Trigger | Fallback/Delegation | Scope |
|---|---|---|---|---|---|
| Industrial inputs (technology, organization, data, information infrastructures) materially amplify the reach and granularity of influence operations in ways that make large-scale cognitive influence practicable. | Empirical work on automation, social bots, programmatic advertising, and platform ecosystems shows that centralized production and distribution can produce orders-of-magnitude increases in message distribution and targeting precision. Historical analogues (mass propaganda) and contemporary studies of computational propaganda provide corroborating mechanisms. | Sustained increases in automated account activity, ad-impression logs showing programmatic buys, increases in message volume and reach metrics, evidence of coordinated account clusters, and platform telemetry showing high-frequency delivery to segmented audiences. | Detection of campaign artifacts consistent with industrialized production (mass-generated posts, sudden spikes in coordinated activity, visible use of programmatic ad infrastructure), major political events or crises that attract organized campaigns, or vendor/platform disclosures of large-scale targeting operations. | If industrial amplification is absent or ineffective, downweight models that assume scale-driven effects and (a) focus research/defense on lower-cost grassroots propagation and organic social dynamics, (b) delegate mitigation to platform-level rate-limiting and provenance/attribution mechanisms, and (c) prioritize institutional policies (regulation of ad buys, data-use restrictions) to prevent future industrialization. | Applies primarily to digitally connected societies and platform ecosystems where programmatic ad markets, user-level data, and automated account infrastructures exist; less applicable in low-connectivity environments, tightly regulated media systems, or contexts lacking rich personal data. |
| Cognition (attention, perception framing, belief updating, decision heuristics) is the proximal mechanism linking industrial inputs to strategic | Fundamental results from information theory and decision theory show that transmitted signals affect beliefs under noise and capacity constraints; | Population- or subgroup-level shifts in opinion polls, survey measures of belief prevalence and confidence, changes in attention metrics (time-on-topic, search volume), downstream behavioral | When strategic outcomes (policy moves, election results, mass behavior) change unexpectedly around intensive influence activity, or when intervention evaluations show divergence between message exposure and material effects—prompting investigation of cognitive mediation. | If cognitive mediation is weak or not the decisive pathway, adjust focus to alternate mechanisms (material attrition, economic disruption, supply-chain targeting). Delegate analysis and mitigation to kinetic, economic, or infrastructure-focused teams while retaining | Most relevant where contested outcomes depend on public or elite beliefs/decisions (political contests, legitimacy battles, mobilization). Less relevant for conflicts purely |

| Assumption | Rationale | Observable | Trigger | Fallback/Delegation | Scope |
|---|---|---|---|---|---|
| outcomes in conflict contexts. | behavioral studies demonstrate how attention and heuristics mediate persuasion. The conceptual chain from signal manipulation to belief and behavior is well supported across multiple literatures. | indicators (voter turnout, protest participation, policy support), and correlational patterns tying exposure metrics to cognitive outcomes. | | some cognitive monitoring for hybrid effects. | about resource control, where cognitive influence is secondary. |
| Network topology and coordinated amplification (bots, sockpuppets, cross-platform orchestration) lower effective adoption thresholds and therefore increase the likelihood and speed of cascades. | Threshold and cascade models in network science show cascade probability is sensitive to topology and local signal strength. Empirical studies demonstrate that coordinated amplification can exaggerate local neighborhood signals, creating apparent consensus and accelerating diffusion. | Rapid, supra-linear spread of specific memes/messages; highly centralized redistribution patterns around a small set of amplifier nodes; synchronous posting patterns; high local clustering with repeated exposures to similar content; platform graph metrics showing hub centricity during events. | Emergence of fast-moving viral content tied to coordinated accounts, discovery of synchronized activity patterns, or when small-origin messages propagate unexpectedly widely—especially during moments of low-preparedness (breaking news, crises). | If amplification tactics are not producing cascades, prioritize hardening measures: platform throttling and provenance labelling, algorithmic de-amplification, targeted disruption of botnets (technical takedowns), and legal/regulatory remedies. Delegate enforcement to platform trust-and-safety teams and cybersecurity/incident response units. | Applies to ecosystems with high resharing affordances and low friction for coordination (open social platforms, messaging apps with forwarding). Effects are attenuated in closed, heavily moderated, or low-sharing networks. |
| Granularity through microtargeting (data-driven | Behavioral and advertising literatures document | Within-campaign heterogeneity in engagement and conversion tied to | Deployment of campaigns that use lookalike/segment-based audiences, platform ad library evidence of many | If microtargeting is infeasible or ineffective, shift to broad-based messaging | Depends on access to rich personal data, advertising |

| Assumption | Rationale | Observable | Trigger | Fallback/Delegation | Scope |
|---|---|---|---|---|---|
| segmentation and tailored messaging) materially increases per-exposure persuasion effectiveness compared with undifferentiated mass messaging. | higher engagement and conversion rates from personalized messaging. Political microtargeting case studies and controlled experiments indicate heterogeneous responsiveness across psychological profiles and that tailored framing can exploit biases effectively. | segment-specific creatives, higher lift in A/B tests on targeted vs. non-targeted arms, ad delivery logs showing fine-grained audience splits, and qualitatively distinct messaging themes across segments. | narrowly targeted ad sets, or when analytics show disparate impacts across demographic/psychographic groups. | strategies, invest in population-level resilience (media literacy, public information campaigns), and implement privacy/data-use restrictions to limit targeting capabilities. Delegate technical countermeasures to privacy teams and regulatory bodies. | infrastructures, and platform microtargeting features; limited in jurisdictions with strict data protection, or in contexts where users are not linkable to behavioral profiles. |
| Speed and closed-loop feedback enabled by automation (A/B testing, rapid content generation) create an arms race dynamic where defenders that are slower or less adaptive are at a structural disadvantage. | Automation lowers iteration cost and shortens feedback loops; adversarial adaptation (rapid pivot following moderation) has been documented in empirical studies. Faster adaptation lets operators optimize for virality and evade static defenses. | High-frequency content variants with small incremental changes, rapid shifts in messaging following takedowns, short-lived accounts and repeated recycling of content, and telemetry showing faster iteration cycles for adversary-controlled infrastructures than for defender responses. | Evidence of quick adversary pivots after mitigation actions, accelerated message churn, or when defenders repeatedly lag in detection/removal timelines; major events that precipitate rapid adversary experimentation. | If defenders cannot match adversary speed, prioritize automation of defensive tools (machine learning detectors, automated takedown pipelines), invest in collective incident-sharing (threat intel), raise costs for adversaries (legal action, platform friction), and focus on upstream resilience (education, redundancy). Delegate real-time countermeasures to SOCs/incident response teams and long-term policy to regulators. | Pertains where adversaries have access to automated generation and deployment tools (AI content generation, bots) and where defenders cannot or do not deploy comparable automation. Less relevant in low-tech adversary contexts or where regulation limits automation. |

# Notation

| Symbol | Meaning | Units / Domain |
| --- | --- | --- |
| $n$ | number of agents | $\mathbb{N}$ |
| $G_t=(V,E_t)$ | time-varying communication/interaction graph | — |
| $\lambda_2(G)$ | algebraic connectivity (Fiedler value) | — |
| $p$ | mean packet-delivery / link reliability | [0,1] |
| $\tau$ | latency / blackout duration | time |
| $\lambda$ | task arrival rate | 1/time |
| $e$ | enforceability / command compliance | [0,1] |
| $\tau_{\text{deleg}}$ | delegation threshold | [0,1] |
| **MTTA** | mean time-to-assignment/action | time |
| $P_{\text{fail}}$ | deadline-miss probability | [0,1] |

# Claim-Evidence-Method (CEM) Grid

| Claim (C) | Evidence (E) | Method (M) | Status | Risk | TestID |
|---|---|---|---|---|---|
| Primary: Industrial scale (reach) increases the baseline probability of message encounter and thus raises population-level adoption/exposure rates for influence operations. | [3] [5] [7] | Analytic probabilistic models of exposure (information theory + population contact models) + agent-based simulations of diffusion on empirical contact graphs + empirical case studies of large-scale campaigns (matched observational analyses). | E cited; M pending simulation and empirical validation (designs specified in Methods T1). | If false, defensive resources prioritized on mass-scale monitoring and blunt countermeasures could be misallocated; policy and platform interventions premised on reach-reduction may yield limited effect. | T1 |
| Primary: Granularity (data-driven microtargeting) increases per-exposure persuasion effectiveness — personalization raises conversion probability conditional on exposure. | [3] [5] [6] | Field-style A/B experiments or synthetic controlled trials (microtargeted vs non-targeted), complemented by simulation of heterogeneous agent susceptibilities and observational validation using matched-ad samples from documented campaigns. | E cited; M pending empirical trials and simulation replication (Methods T2). | If incorrect, emphasis on countering microtargeting (e.g., data-access controls) may be over-emphasized relative to broader platform-level interventions; attribution of effects to personalization could be erroneous. | T2 |
| Primary: Automation plus rapid feedback (speed & adaptation) enables adversarial learning loops that sustain and adapt misinformation campaigns, increasing persistence and reducing the time to achieve targeted belief states. | [5] [8] [3] | Reinforcement-learning / adversarial agent simulations that pair automated content generators with response-driven adaptation; longitudinal empirical analysis of campaign lifecycles and adaptation signatures (message variant turnover, rapid re-seeding patterns). | E cited; M pending simulation experiments and longitudinal empirical work (Methods T3). | If wrong, defensive emphasis on slowing production or adding latency may be less effective than thought; misestimation of campaign persistence could leave platforms under-protected or misdirected. | T3 |
| Primary: Orchestration/coordination across modalities (bots, paid content, earned media) amplifies perceived consensus and, by concentrating influence | [9] [10] [7] [5] | Network-theoretic analysis (threshold/cascade models) and simulations on empirical social graphs to measure cascade probability under coordinated seeding; | E cited (foundational and domain sources); M pending targeted | If incorrect, investments aimed at protecting or hardening 'critical nodes' and cross-modality detection | T4 |

| Claim (C) | Evidence (E) | Method (M) | Status | Risk | TestID |
|---|---|---|---|---|---|
| on network-critical nodes, can lower effective adoption thresholds and trigger cascades in susceptible topologies. | | empirical cascade detection and counterfactual seeding experiments where feasible. | simulations and observational causality tests (Methods T4). | may not yield expected reductions in large-scale influence events; defensive doctrines could miss other dominant pathways. | |
| Secondary: Industrial inputs alter channel capacity and noise characteristics (synthetic-but-coherent signals), degrading verification heuristics and contributing to epistemic fragmentation and polarization. | [1] [3] [4] [11] | Information-theoretic models of channel capacity and signal/noise tradeoffs applied to synthetic-signal regimes, plus empirical measures of trust/verification failures and polarization metrics before/after documented synthetic-signal shocks. | E cited; M pending formal modeling exercises and empirical before/after analyses (Methods T5). | If false, policy responses focused on infrastructure resilience or authentication schemes may be over-prioritized at the expense of content- or network-level interventions; misdiagnosis of polarization drivers. | T5 |
| Secondary: Structural vulnerabilities (attention bottlenecks and 'super-blockers' or super-attention nodes) imply that targeted defensive interventions at attention hubs yield outsized returns in resilience compared with uniform measures. | [10] [9] [5] | Intervention simulations (targeted inoculation/provenance marking) on topologies with super-blockers, controlled platform experiments where feasible, and observational evaluation of past platform interventions targeted at hubs. | E cited; M pending experimental trials and simulation validation (Methods T6). | If wrong, concentrating defenses on a small set of attention hubs could leave broad population channels exposed; resources might be inefficiently concentrated, reducing overall resilience. | T6 |

## Sources

**[1]**

On graph theoretic results underlying the analysis of consensus in multi-agent systems

Arxiv.Org, 2009-02-24. (cred: 0.50)

http://arxiv.org/abs/0902.4218v1

**[2]**

In 'crisis' we trust? On (un) intentional knowledge distortion and the exigency of terminological clarity in academic and political discourses on Russia's war against ...

Link.Springer.Com, 2023-01-01. (cred: 0.50)

https://link.springer.com/article/10.1057/s41268-023-00313-2

**[3]**

Transforming Cognition and Human Society in the Digital Age

Link.Springer.Com, 2024-01-01. (cred: 0.50)

https://link.springer.com/article/10.1007/s13752-024-00483-3

**[4]**

The brain of the future and the viability of democratic governance: The role of artificial intelligence, cognitive machines, and viable systems

Sciencedirect.Com, 2018-01-01. (cred: 0.50)

https://www.sciencedirect.com/science/article/pii/S0016328717302896

**[5]**

Online Influence Campaigns: Strategies and Vulnerabilities

Arxiv.Org, 2024-12-18. (cred: 0.50)

http://arxiv.org/abs/2501.10387v1

**[6]**

MetaTroll: Few-shot Detection of State-Sponsored Trolls with Transformer Adapters

Arxiv.Org, 2023-03-13. (cred: 0.50)

http://arxiv.org/abs/2303.07354v1

**[7]**

Social Bots and Social Media Manipulation in 2020: The Year in Review

Arxiv.Org, 2021-02-16. (cred: 0.50)

http://arxiv.org/abs/2102.08436v1

**[8]**

When Hallucination Costs Millions: Benchmarking AI Agents in High-Stakes Adversarial Financial Markets

Arxiv.Org, 2025-09-30. (cred: 0.50)

http://arxiv.org/abs/2510.00332v1

**[9]**

A network-based microfoundation of Granovetter's threshold model for social tipping

Arxiv.Org, 2019-11-11. (cred: 0.50)

http://arxiv.org/abs/1911.04126v2

**[10]**

Super-blockers and the effect of network structure on information cascades

Arxiv.Org, 2018-02-14. (cred: 0.50)

http://arxiv.org/abs/1802.05039v2

**[11]**

Who's Your Judge? On the Detectability of LLM-Generated Judgments

Arxiv.Org, 2025-09-29. (cred: 0.50)

http://arxiv.org/abs/2509.25154v1

Generated: 2025-11-09T17:47:59.195404 | Word Count: 5344

## Research Roadmap

- **Phase 1 (Theory)**: Formalize claims, extend proofs, validate against canonical results
- **Phase 2 (Simulation)**: Implement stress tests, sweep parameter spaces, measure convergence/scaling
- **Phase 3 (Empirical)**: Deploy in controlled environments, collect field data, validate predictions
- **Phase 4 (Integration)**: Operationalize with human-in-loop, adversarial hardening, production deployment

**Confidence Methodology:** Confidence = 0.3·SourceDiversity + 0.25·AnchorCoverage + 0.25·MethodTransparency + 0.2·ReplicationReadiness, where SourceDiversity reflects unique publishers & types, AnchorCoverage reflects share of primary claims with Type-1 anchors, MethodTransparency reflects CEM completeness & assumptions ledger, and ReplicationReadiness reflects sim plan & datasets/params specified.

Prepared under the STI Research Program — theoretical framework subject to revision as data accumulate.