

LegalSim: Multi-Agent Simulation of Legal Systems for Discovering Procedural Exploits

Sanket Badhe

Rutgers University

sanketbadhe1611@gmail.com

Abstract

We present LEGALSIM, a modular multi-agent simulation of adversarial legal proceedings that explores how AI systems can exploit procedural weaknesses in codified rules. Plaintiff and defendant agents choose from a constrained action space (for example, discovery requests, motions, meet-and-confer, sanctions) governed by a JSON rules engine, while a stochastic judge model with calibrated grant rates, cost allocations, and sanction tendencies resolves outcomes. We compare four policies: PPO, a contextual bandit with an LLM, a direct LLM policy, and a hand-crafted heuristic; Instead of optimizing binary case outcomes, agents are trained and evaluated using effective win rate and a composite exploit score that combines opponent-cost inflation, calendar pressure, settlement pressure at low merit, and a rule-compliance margin. Across configurable regimes (e.g., bankruptcy stays, inter partes review, tax procedures) and heterogeneous judges, we observe emergent “exploit chains”, such as cost-inflating discovery sequences and calendar-pressure tactics that remain procedurally valid yet systemically harmful. Evaluation via cross-play and Bradley-Terry ratings shows, PPO wins more often, the bandit is the most consistently competitive across opponents, the LLM trails them, and the heuristic is weakest. The results are stable in judge settings, and the simulation reveals emergent exploit chains, motivating red-teaming of legal rule systems in addition to model-level testing.

1 Introduction

The legal system is an adversarial process guided by dense procedural rules that shape how disputes unfold. Litigants do not only argue substance; they sequence filings, exploit timing, and impose tactical costs to influence outcomes. As AI enters legal practice, these tactics may be amplified: learning agents can search large procedural spaces, probe edge cases at scale, and coordinate strategies with

speed and persistence beyond human capacity. This possibility raises questions at the intersection of natural legal language processing, multi-agent reinforcement learning, and AI safety (Amodei et al., 2016).

Most work in legal NLP treats models as assistive tools that classify, summarize, retrieve, or predict (Chalkidis et al., 2020, 2022; Zhong et al., 2019). These settings assume a largely passive role for AI within human workflows. Far less is known about what happens when AI agents interact directly with codified procedure and with each other. In complex systems, agents trained to optimize rewards often uncover loopholes that remain technically compliant while socially harmful (Amodei et al., 2016). The legal process, with its motion practice, deadlines, and rule-based gates, is a natural domain where such behavior may emerge.

We argue that studying these dynamics requires a simulation environment that treats litigation as strategic interaction under rules. Our approach frames procedure as a structured action space with observable state, limited information, and stochastic judicial response. Agents learn over repeated play to pursue objectives that extend beyond win or loss, including cost imposition, delay, and settlement leverage under sanction risk. By varying rule sets across domains, the same environment can reveal how different procedural regimes encourage or deter exploitative behavior.

We introduce LEGALSIM, a modular multi-agent framework for adversarial legal proceedings. Plaintiff and defendant agents select structured actions validated by a JSON rule engine that encodes domain-specific procedural gates; a stochastic judge mediates outcomes via calibrated grant rates, cost allocations, and sanction tendencies. Policies include a hand-crafted heuristic baseline, a contextual bandit over tactic families, a PPO policy trained in self-play, and a direct LLM policy (Schulman et al., 2017; Silver et al., 2016; Lowe

et al., 2020; Silver et al., 2017; OpenAI et al., 2019; Vinyals et al., 2019). Rather than optimizing a binary case outcome, agents receive a composite exploit score that aggregates opponent-cost inflation, calendar pressure, settlement pressure conditional on low merits, and a rule-compliance margin.

Contributions.

1. **Formalization of litigation as a MARL environment.** We model adversarial legal proceedings as a multi-agent environment with a structured token space, machine-readable procedural gates, and a calibrated judge model, enabling regime-agnostic studies.
2. **Discovery of emergent legal exploits.** Through self-play, interacting agents discover strategies that were not pre-programmed, including tactics observed in practice and novel exploit chains that expose systemic vulnerabilities (cost inflation, calendar pressure, settlement leverage under sanction risk).
3. **Evaluation protocol and artifacts.** We evaluate with head-to-head and all-against-all cross-play and fit role-symmetric Bradley-Terry-Luce ratings and robustness sweeps across judges, enabling comparison of heuristic, contextual bandit, PPO, and LLM-guided policies.
4. **AI-safety perspective on law.** We argue for red-teaming codified legal systems themselves rather than only individual models, offering a testbed for measuring and mitigating AI-amplified procedural abuse.

Our findings suggest an AI-safety perspective that red-teams not only models but the legal rule systems themselves. LEGALSIM offers a testbed for measuring and mitigating procedural exploitation, linking methods from legal NLP, MARL, and robustness analysis (Balduzzi et al., 2019; Omidshafiei et al., 2019).

2 Background and Related Work

Our research sits at the intersection of four areas: AI for legal reasoning, AI safety in legal contexts, multi-agent systems for strategic discovery, and the formalization of law as code. Each area is established, but their synthesis to red-team legal frameworks is new.

2.1 AI in Legal Reasoning and Prediction

A substantial body of Natural Legal Language Processing (NLLP) focuses on analytical and predictive tasks. Early work showed that machine learn-

ing can predict judicial outcomes from case facts (Aletras et al., 2016). More recent approaches leverage large language models and legal-specific pre-training, such as LEGAL-BERT (Chalkidis et al., 2020), achieving strong results on legal judgment prediction, document classification, and argument mining. These tools reason about or predict outcomes in a static setting; they are not designed to act as strategic agents within a procedural process. Our work shifts the focus from passive prediction to active, strategic participation.

2.2 AI Safety and Fairness in Law

As capabilities grow, concerns about safety and fairness have intensified. The dominant paradigm is to identify and mitigate model-level flaws, including demographic bias in predictive justice systems, as highlighted by the COMPAS investigation (Angwin et al., 2016). Additional lines address robustness of legal text classifiers and the explainability of black-box models in service of due process (Richmond et al., 2023). This work is essential, but it primarily addresses harms from systems that are wrong or biased. We study a complementary risk: harms produced by agents that are competent and strategically exploit codified rules to achieve unfair or inefficient outcomes.

2.3 Multi-Agent Systems and Emergent Strategy

Outside law, multi-agent reinforcement learning (MARL) has uncovered novel strategies in complex adversarial settings. Self-play has yielded superhuman policies in Go (Silver et al., 2016) and StarCraft II (Vinyals et al., 2019). Work on emergent tool use demonstrates autocurricula in competitive environments (Baker et al., 2020). Classical game-theoretic analyses of litigation exist (Baird et al., 1994), but modern MARL for discovering procedural strategies from a blank slate remains underexplored. Recent agent systems that blend planning with language interaction, such as C-CERO for Diplomacy (FAIR et al., 2022), suggest feasibility for mixed-motive negotiation akin to litigation. Concurrently, agentic-risk studies examine malicious or deceptive uses of LLM agents; for example, ScamAgents demonstrates how autonomous AI agents can be architected to simulate and execute complex, human-level scam calls (Badhe, 2025). Our work brings these ideas into a rules-constrained, legally grounded domain.

2.4 Computational Law and Rules-as-Code

Formalizing legal rules in machine-readable form (*rules-as-code*) is a prerequisite for procedural simulation. Foundational visions in computational law aim to represent statutes, regulations, and contracts with logical precision (Genesereth, 2005; Surden, 2012). Prior applications emphasize compliance checking, digital advisory tools, and expert systems. We build directly on this foundation but use the formalized ruleset as the “physics engine” for an adversarial simulation, enabling stress tests where intelligent agents interact strategically. While rules-as-code focuses on encoding law as written, our objective is to surface unintended and exploitative consequences that can emerge in practice.

3 Problem Formulation

We model litigation as a multi-agent adversarial game governed by codified procedure. Two agents, plaintiff and defendant, act in a structured environment mediated by a judge. The objective is to allow learning agents to discover strategies from the environment dynamics rather than rely on hand-coded heuristics.

3.1 State

At time t the environment state is

$$s_t = \{P^{\text{pl}}, P^{\text{df}}, C, H\},$$

where P^{pl} and P^{df} are party states, C encodes court attributes including active procedural gates and judicial tendencies, and H is a structured history of filings, rulings, and citations. Party states track budgets, accumulated burden, fees, sanctions, and merits, enabling decisions conditioned on posture and history. Judicial tendencies are parameterized by a profile with a grant rate and a sanction tendency that shape probabilistic rulings.

3.2 Actions and Gates

At each step an agent chooses $a_t \in \mathcal{A}(s_t)$, with availability constrained by active gates:

$$a_t \in \mathcal{A}(s_t) \quad \text{iff no gate blocks } a_t.$$

The action space covers core procedural moves such as filing a proceeding, referencing authority, requesting discovery, moving for sanctions, changing venue, and making settlement offers. In implementation these map to structured action tokens. Gates implement rule-based blocks that delay or

nullify certain actions until expiry, matching the formal constraint above. The full inventory of abstract action tokens is listed in App. B.

3.3 Transition Dynamics

A rule engine evaluates actions and updates state with deterministic and stochastic effects: impose or lift gates, allocate costs and burdens, apply sanctions, and progress the case toward termination. This supports faithful procedural interaction without prescribing strategies.

3.4 Rewards

Each agent’s reward combines competing litigation objectives:

$$\begin{aligned} R_t = & w_1 \cdot \text{OpponentCost}_t + w_2 \cdot \text{DelayCredit}_t \\ & + w_3 \cdot \text{OutcomeBonus}_t \\ & - w_4 \cdot \text{SanctionPenalty}_t \end{aligned} \quad (1)$$

with weights w_i tuning strategic preferences. Plaintiffs seek favorable outcomes with cost control, while defendants emphasize dismissal or delay with minimal sanction exposure. One instantiated shaping in code rewards increases in opponent burden while penalizing own cost and burden, with a terminal bonus or penalty at resolution.

3.5 Learning Objective

We optimize policies for both sides under discounted returns:

$$\pi^{\text{pl}}, \pi^{\text{df}} = \arg \max_{\pi} \mathbb{E} \left[\sum_{t=0}^T \gamma^t R_t(s_t, a_t) \right],$$

where $\gamma \in (0, 1)$ is the discount factor.

3.6 Illustrative Gate Scenario

If a defendant files for Chapter 11¹, an automatic-stay gate activates and blocks motions and discovery for sixty timesteps. This induces a delay exploit for the plaintiff. After expiry, a calibrated discovery sequence can inflate the defendant’s costs if judicial sanctions are unlikely, yielding a cost-inflation exploit while each step remains procedurally valid.

This formulation specifies the state, action, gating, transition, reward, and objective needed to study emergent procedural exploits within a principled multi-agent reinforcement learning setting.

¹Chapter 11 of the U.S. Bankruptcy Code

4 System Architecture

Our system’s architecture is organized into a three-tiered structure: the LEGALSIM environment, the agents that interact within it, and a training and analysis harness. This design separates the core simulation logic from the agent policies and experiment orchestration, which enables modularity and the easy substitution of components without changing core interfaces.

4.1 Environment Layer

The environment layer is a domain-agnostic litigation simulator. It is driven by a rules-as-code engine that processes an abstract action space and updates the simulation state based on a set of JSON rules. Each environment instance maintains the state of two adversarial parties, tracking their budgets, accumulated burdens, and merits. It also incorporates a stochastic judge profile that influences probabilistic rulings, such as the granting of motions or the imposition of sanctions.

4.1.1 Rules-as-code Engine

The core of the environment is a JSON rule engine (a finite set of state–action predicates with effect handlers) that loads procedural rules and gates (temporary action blocks) from JSON files (Governatori et al., 2011). Rules are defined by when conditions (e.g., a specific action is taken) and effects that modify the environment state. Effects can include applying costs, transferring burdens, and, most importantly, activating procedural gates. The RulesOracle component provides an approximation of legal principles like proportionality and sanction risk, which are parameterized to allow for policy studies rather than strict encoding of legal doctrine. We include rule example at App. A.

4.1.2 Procedural regimes.

The environment swaps procedural regimes by loading different JSON rule files without code changes. In our experiments we use a default bankruptcy regime and also include domain-specific sets for patent, tax collection, immigration, and corporate disputes. Each regime defines gates (temporary blocks on actions) and effects that shape costs, burdens, delay, and sanctions.

4.1.3 Action Interface and State

Agents interact with the environment by emitting abstract, tokenized actions. In total, there

are 13 possible action tokens, including representative examples such as REQUEST_DOCS, FILE_PROCEEDING, and SETTLEMENT_OFFER. The environment validates these tokens against a pre-defined schema, ensuring the action space remains structured and the simulation runs smoothly. The environment state includes party-specific metrics (budget, burden, merits) and global context from the judge, such as grant rate and sanction tendency.

4.2 Agent Layer

The agent layer supports multiple policy families that are swappable through a common interface. This allows for a mix-and-match approach in experiments, where different agent types can be pitted against each other. The policies include:

1. **Heuristic Policy:** A hand-coded, rule-based baseline that makes decisions based on simple, pre-defined logic related to costs and burdens.
2. **LLM-driven Policy:** A policy that queries a large language model (LLM) for a reasoned action. It uses few-shot prompting and enforces a strict JSON output contract to ensure the LLM’s free-form reasoning can be translated into a valid action token. Throughout this paper, all LLM calls use OpenAI’s GPT-4o.
3. **Contextual Bandit Policy:** A hybrid policy that first uses a contextual bandit to select a high-level "tactic" (e.g., DELAY, BURDEN_OPP), and then uses the LLM to propose a specific action consistent with that tactic. (Li et al., 2010)
4. **PPO Policy:** A policy based on Proximal Policy Optimization (PPO), a reinforcement learning algorithm that learns to select actions from the environment’s observations (Schulman et al., 2017; Silver et al., 2016).

4.3 Training and Evaluation Harness

The harness coordinates self-play experiments, enforcing role alternation and judge rotation, scheduling learning updates for PPO and the contextual bandit, and validating all emitted action tokens.

4.3.1 Episode Flow

A single episode unfolds as follows: the harness initializes the environment, agents observe the state and propose actions, the environment validates and executes these actions, and the state advances. This process repeats until a termination condition is met (e.g., budget exhaustion, settlement, or maximum steps). At termination, composite exploit metrics

are calculated, and learning updates are applied to the agents’ policies if enabled.

5 Experiments and Evaluation

We evaluate LEGALSIM under a controlled protocol that alternates roles each game, sweeps ten random seeds, and rotates between two judge profiles: *permissive* (grant_rate 0.65, sanction_tendency 0.25, calendar_load 0.55) and *strict* (0.35, 0.70, 0.60). Domains are loaded from JSON rule files; unless noted, we use the default regime.

Policies and training. We evaluate the four policy families introduced in Sec. 4 (Heuristic, LLM, Contextual Bandit, PPO). The *Heuristic* is non-learning. The *LLM* policy uses a single API model at inference time and emits JSON-constrained tokens without any fine-tuning. The *Contextual Bandit* selects a high-level tactic via an ϵ -greedy linear contextual bandit with a bias term ($\epsilon=0.1$, learning rate 0.05), then asks the same LLM to instantiate a concrete token; it performs one SGD update per episode on the terminal composite reward. The *PPO* agent is an actor critic over the discrete token set with a 13-D observation (budgets, burdens, judge features, merits, progress, gate summaries); both actor and critic are two-layer MLPs (64 Tanh units each) trained with Adam (3×10^{-4}), $\gamma=0.99$, GAE $\lambda=0.95$, clip $\epsilon=0.2$, and entropy coefficient 0.005. PPO optimizes the shaped reward

$$r_t = 0.20 \Delta(\text{opponent burden}) - 0.01 \Delta(\text{own cost}) - 0.01 \Delta(\text{own burden}). \quad (2)$$

with a terminal bonus of +5 (plaintiff win) and −5 (defendant win), is trained for 300 episodes against the Heuristic while alternating judges by episode, and is then frozen for evaluation.

Environment and rules. Litigation is modeled as a turn-based process with a rules-as-code core. Agents emit tokens such as REQUEST_DOCS, FILE_MOTION, MOVE_SANCTIONS, MEET_CONFER, SETTLEMENT_OFFER, and FILE_PROCEEDING. A JSON rule engine maps state–action conditions to cost transfers, burden updates, temporary gates that block actions, and judge-sensitive sanction events; a RulesOracle provides proportionality and sanction-risk proxies. Each episode tracks budgets, burdens, merits, fees, sanctions, active gates, and the judge profile.

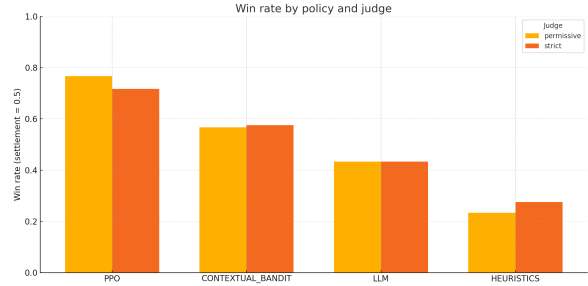


Figure 1: Win rate by policy & judge (settlement = 0.5). Bars show mean effective win rate across ten seeds under *permissive* and *strict* judges; higher is better for the policy.

Protocols. We use two complementary designs. (i) *Head-to-head*: selected pairs play under both judges with alternating roles; we log token sequences, rulings, and per-role metrics. (ii) *Cross-play league*: all policies play all others across seeds and judges. From these games we build a role-symmetric payoff matrix A whose entry

$$A_{ij} = \mathbb{E}[(\text{plaintiff composite}) - (\text{defendant composite})] \quad (3)$$

when policy i faces j , flipping the sign when roles swap so cells are comparable. We also fit role-symmetric Bradley-Terry-Luce (BTL) ratings with a sum-to-zero constraint (Bradley and Terry, 1952):

$$\Pr(i \succ j) = \sigma(s_i - s_j), \quad \sigma(x) = \frac{1}{1 + e^{-x}},$$

and report 95% bootstrap confidence intervals over 500 resamples. Prompt templates used during evaluation are reproduced in App. C

Overall effectiveness. The win-rate analysis in Figure 1 (Win-rate by policy & judge) shows a consistent ordering: PPO attains the highest effective win rate overall, followed by the contextual bandit, then the LLM policy, with the heuristic trailing. The same ordering holds within each judge profile. This indicates that learning a direct policy over the token space (PPO) converts the observable state into match wins more reliably than tactic selection with LLM instantiation (contextual bandit) or purely generative action proposals (LLM). Table 1 quantifies these differences alongside mean composite exploit scores for plaintiff/defendant roles.

Who beats whom (meta-game structure). Figure 2 summarizes pairwise performance using two

Policy	Win rate _{eff}	Flag rate	\bar{C}_{pl}	\bar{C}_{df}	BTL	BTL CI _{low}	BTL CI _{high}
contextual_bandit	0.571	0.958	1.130	0.948	99.6	74.3	124.7
llm	0.433	1.000	15.789	5.509	95.9	70.3	120.9
ppo	0.742	0.958	1.280	0.845	-97.6	-123.0	-72.6
heuristic	0.254	1.000	39.995	8.337	-97.9	-122.6	-71.6

Table 1: Evaluation summary by policy. Win rate_{eff} treats settlements as 0.5. \bar{C}_{pl} and \bar{C}_{df} are mean composite exploit scores for plaintiff/defendant roles. BTL and 95% bootstrap CIs are from role-symmetric Bradley-Terry-Luce fits on the cross-play league.

role-symmetric metrics: (i) win rate with settlements counted as 0.5, averaged across both role assignments for each policy pair (row policy i vs. column policy j); and (ii) composite margin, the mean difference in exploit score (plaintiff composite – defendant composite) with the sign flipped when roles swap, so positive values indicate that the row policy systematically exerts more procedural pressure than the column policy. The heatmaps yield a consistent ordering: the Contextual Bandit dominates win rates ≥ 0.56 against all opponents and positive margins (largest vs. Heuristic, modest vs. LLM); the LLM policy is second, clearly ahead of PPO and Heuristic; PPO shows advantage only over the Heuristic; and the Heuristic is uniformly weakest. This ranking holds in both outcome space (win rate) and pressure space (composite margin), indicating the bandit’s broad competitiveness across the meta-game. An example episode underlying a high-margin cell is unpacked in App. D.

Exploitiveness metrics. Each episode produces per-role components already defined in the environment: (i) *opponent cost inflation* (opponent fees divided by own fees), (ii) *calendar pressure* (opponent burden divided by $1 +$ own burden), (iii) *settlement pressure at low merit* (settlement offers times $1 -$ own merits, clipped), and (iv) *rule-compliance margin* (a penalty for self-sanctions). The composite exploit score is a fixed weighted sum of these components (0.35/0.25/0.25/0.15). We summarize these per policy with means and standard errors and also report the *flag rate*, the fraction of episodes with composite ≥ 0.6 .

Exploitiveness results. Applying these definitions, Table 2 reports per-policy \times judge means, standard errors, flag rates, and episode counts. Two patterns emerge. First, the Heuristic and LLM policies produce very high composite scores with near-ubiquitous flagging across judges, indicating heavy procedural pressure. Second, PPO and the Contextual

Bandit maintain composites near 1 with non-maximal flag rates, and show judge sensitivity (the bandit declines under the strict judge, whereas PPO ticks up slightly). Together with Figure 1, this confirms that effectiveness (win rate) and exploitiveness are related but not identical: PPO converts state to wins while applying less extreme procedural pressure than the Heuristic or LLM, and the bandit sits between these extremes.

Judge effects. Breaking out the bars in Figure 1 by judge shows that absolute win rates shift with judicial temperament, but the relative ordering of policies remains stable. On the permissive judge, motion-driven strategies benefit more; on the strict judge, margins compress but the ranking persists. This mirrors the sanction and grant-rate sensitivities in Table 1 and the stratified means in Table 3.

Robustness. We stress-test the policies in two simple ways: (i) we make the judge more likely to impose sanctions, and (ii) we add random $\pm 10\%$ – 20% perturbations to cost and burden parameters to mimic modeling noise. For each policy and judge profile we then recompute two summaries, the mean composite exploit score and the fraction of episodes that are flagged ($C \geq 0.6$), and present them as a policy-by-stress matrix (Table 3).

Across all stress tests, the qualitative ordering of policies does not change: PPO remains strongest on outcomes, the contextual bandit is generally second, the LLM trails, and the heuristic is consistently weakest. Making sanctions stricter reliably lowers exploit scores and flag rates, with the biggest reductions for strategies that lean on filing volume and burden (LLM, then bandit), while PPO is least affected. Injecting cost/burden noise increases variability but does not reverse pairwise rankings. In short, the effects we report are stable to reasonable procedural and parameter changes; stricter sanction regimes act as a partial brake on exploit-heavy behavior without reshuffling the policy hierarchy.

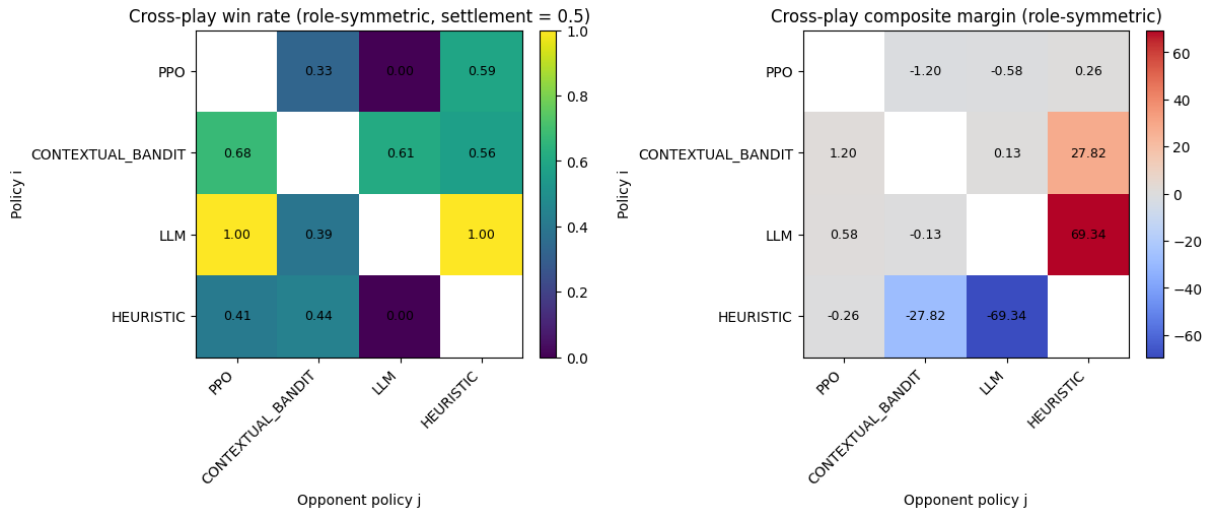


Figure 2: **Cross-play performance heatmaps (role-symmetric)**. Left: win rate with settlements counted as 0.5, averaged over both role assignments. Right: composite margin (plaintiff composite – defendant composite) with the sign flipped when roles swap. Rows index the *row* policy i and columns the *opponent* policy j ; numbers are cell means across seeds and judges. Higher (warmer) values indicate that the row policy systematically outperforms (or exerts more procedural pressure than) the column policy.

Policy	Judge	Mean Composite	SE Composite	Flag Rate	N Episodes
Contextual Bandit	Permissive	1.18	0.145	0.83	60
Contextual Bandit	Strict	0.90	0.098	0.70	60
Heuristic	Permissive	23.24	5.353	1.00	60
Heuristic	Strict	25.09	5.893	0.98	60
LLM	Permissive	10.92	2.558	1.00	60
LLM	Strict	10.38	2.650	1.00	60
PPO	Permissive	1.01	0.054	0.78	60
PPO	Strict	1.12	0.065	0.83	60

Table 2: Exploiteness summary by policy and judge. *Mean Composite* is the average composite exploit score C ; *SE Composite* is the standard error; *Flag Rate* is the fraction of episodes with $C \geq 0.6$; and *N Episodes* is the number of episodes summarized.

Reconciling win rate with BTL ratings. BTL summarizes global competitiveness from the full cross-play, not just wins. It can rank a policy higher when it draws fewer severe losses and plays most opponents close, even if its raw win rate is slightly lower. In our run, PPO tops Table 2 for win rate; BTL places the contextual bandit and LLM closer in the middle of the meta-game (with overlapping confidence intervals), while the heuristic sits clearly below. This is consistent with Figure 2 showing near-zero margins between the mid-tier policies and large negative margins concentrated in the heuristic row and column.

6 Defense, and Mitigation

Risks

Optimizing agents can find strategies that are legal but harmful, turning procedural gaps into cost

and delay weapons (Amodei et al., 2016). Law is especially exposed because it is highly codified, adversarial, and variably adjudicated. This risk of “reward hacking,” where an agent satisfies the literal specification of a reward function in an unintended way, is a fundamental challenge in aligning AI with complex, real-world objectives (Leike et al., 2018).

Beyond accuracy. The problem is not only wrong predictions but system-level exploits that emerge when agents play the rules. League and cross-play results echo findings in open multi-agent games: small tactical gains can snowball into undesirable equilibria (Balduzzi et al., 2019). Evaluating performance in such ecologies requires methods beyond simple win-rates, such as AlphaRank or Bradley-Terry models, to capture the full matrix of strategic interactions (Omidshafiei et al., 2019; Bradley and Terry, 1952).

Sanction tendency sweep			
Sanction	Mean C	95% low	95% high
0.10	24.63	16.61	32.52
0.25	24.62	16.84	32.68
0.50	24.19	16.86	32.04
0.75	24.44	17.40	32.07
0.90	25.69	17.54	33.91
Parameter noise sweep			
Noise	Mean C	95% low	95% high
-0.20	32.22	18.54	46.43
-0.10	27.15	15.56	39.12
0.00	26.49	15.20	38.26
+0.10	24.48	13.85	35.26
+0.20	22.02	12.31	32.01

Table 3: Robustness summary (aggregate). Mean composite exploit score C and bootstrap 95% CIs under (top) sanction-tendency sweep and (bottom) multiplicative parameter-noise sweep applied to costs/burdens. Values are aggregated across policies and both judges, with $n_{\text{episodes}} = 60$ simulation runs per configuration.

Design-time defenses. Harden procedures before deployment: introduce light randomization in scheduling, add rule linting to detect long burden-inflating chains, and tie cost shifting to burden ratios so exploit-heavy sequences become expensive.

Governance. Require pre-deployment red-teaming in a rules-as-code sandbox, disclosure of agent capabilities, and auditable reports with cross-play matrices, BTL ratings, and exploit dashboards. This aligns with emerging AI governance standards, such as the NIST AI Risk Management Framework, which emphasizes continuous testing, evaluation, and risk mitigation throughout the AI lifecycle

The core risk is not that agents can win, but that they can steer rule-driven systems toward exploit-heavy equilibria. Simulation-first analysis, league benchmarks, and targeted procedural guardrails provide a practical path to measure and mitigate these effects.

7 Ethical Considerations

This work examines how agentic AI might exploit codified procedure with the goal of improving safety and fairness in legal automation by identifying failure modes before real-world use. By surfacing and quantifying “exploit chains,” we aim to support due-process values and reduce risks. We acknowledge dual-use concerns: the same insights could enable misuse. To mitigate

this, we confine our analysis to simulation, avoid jurisdiction-specific guidance, and emphasize safeguards. LEGALSIM is a research simulator and not legal advice.

8 Limitations

Our simulator necessarily abstracts complex laws and institutional practice: rules are encoded at a coarse level, the judge model is parametric and stationary, and strategy discovery is constrained by a tokenized action space. The exploit metrics and their weights, though motivated, are ultimately design choices that may not capture the full spectrum of welfare-relevant harms. Similarly, policy coverage remains narrow (four families) and training horizons modest, such that stronger or more sample-efficient methods could shift the observed rankings. These results should therefore not be assumed to generalize across jurisdictions, case types, or institutional settings, especially since the environment omits strategic behavior by non-party actors such as regulators or multi-judge panels.

A central limitation is that our findings have not yet been grounded in real-world judicial data or case law. While the experiments reveal how artificial agents may exploit procedural rules in silico, we have not examined whether comparable exploitative dynamics occur in practice, nor how judicial actors (e.g., judges, clerks, regulators) adapt to mitigate such behaviors.

Finally, we emphasize that simulations cannot substitute for legal or ethical judgment. Insights derived here should inform, but never replace, human governance and procedural safeguards.

9 Conclusion

We introduced LEGALSIM, a modular multi-agent simulation that treats procedure as rules-as-code and measures how AI-driven strategies can exploit legal process. Across head-to-head and cross-play evaluations, we observed consistent ordering among policies, documented emergent “exploit chains,” and quantified exploitiveness with outcome and pressure centric metrics. These results frame procedural robustness as an AI-safety problem: not only how models behave, but how codified rules can be gamed. The framework provides a controlled setting to study defenses, e.g. randomized gates, human review for high-impact actions, and system-level red-teaming, before deployment in real practice. We hope this work catalyzes collab-

oration between NLP, MARL, and legal communities on measuring and mitigating AI-amplified procedural abuse.

References

- Nikolaos Aletras, Dimitrios Tsarapatsanis, Daniel Preoŕiuc-Pietro, and Vasileios Lampos. 2016. [Predicting judicial decisions of the european court of human rights: a natural language processing perspective](#). *PeerJ Computer Science*, 2:e93.
- Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. 2016. [Concrete problems in ai safety](#). *Preprint*, arXiv:1606.06565.
- Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2016. Machine bias: There’s software used across the country to predict future criminals. and it’s biased against blacks. <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>.
- Sanket Badhe. 2025. [Scamagents: How ai agents can simulate human-level scam calls](#). *Preprint*, arXiv:2508.06457.
- Douglas G Baird, Robert H Gertner, and Randal C Picker. 1994. *Game theory and the law*. Harvard University Press.
- Bowen Baker, Ingmar Kanitscheider, Todor Markov, Yi Wu, Glenn Powell, Bob McGrew, and Igor Mordatch. 2020. [Emergent tool use from multi-agent autotutorials](#). *Preprint*, arXiv:1909.07528.
- David Balduzzi, Marta Garnelo, Yoram Bachrach, Wojciech Czarnecki, Julien Perolat, Max Jaderberg, and Thore Graepel. 2019. [Open-ended learning in symmetric zero-sum games](#). In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 434–443. PMLR.
- Ralph Allan Bradley and Milton E. Terry. 1952. [Rank analysis of incomplete block designs: I. the method of paired comparisons](#). *Biometrika*, 39(3/4):324–345.
- Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2020. [LEGAL-BERT: The muppets straight out of law school](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2898–2904, Online. Association for Computational Linguistics.
- Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael Bommarito, Ion Androutsopoulos, Daniel Katz, and Nikolaos Aletras. 2022. [LexGLUE: A benchmark dataset for legal language understanding in English](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4310–4330, Dublin, Ireland. Association for Computational Linguistics.
- FAIR, Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, Andrew Goff, Jonathan Gray, Hengyuan Hu, Athul Paul Jacob, Mojtaba Komeili, Karthik Konath, Minae Kwon, Adam Lerer, Mike Lewis, Alexander H. Miller, Sasha Mitts, Adithya Renduchintala, and 8 others. 2022. [Human-level play in the game of <i>diplomacy</i> by combining language models with strategic reasoning](#). *Science*, 378(6624):1067–1074.
- Michael Genesereth. 2005. [Computational law](#). In *Proceedings of the 10th International Conference on Artificial Intelligence and Law (ICAIL ’05)*, pages 12–13, New York, NY, USA. ACM.
- Guido Governatori, Francesco Olivieri, Simone Scanapieco, and Matteo Cristani. 2011. [Designing for compliance: Norms and goals](#). In *Rule-Based Modeling and Computing on the Semantic Web*, pages 282–297, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Jan Leike, David Krueger, Tom Everitt, Miljan Martic, Vishal Maini, and Shane Legg. 2018. [Scalable agent alignment via reward modeling: a research direction](#). *Preprint*, arXiv:1811.07871.
- Lihong Li, Wei Chu, John Langford, and Robert E. Schapire. 2010. [A contextual-bandit approach to personalized news article recommendation](#). In *Proceedings of the 19th international conference on World wide web*, WWW ’10, page 661–670. ACM.
- Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. 2020. [Multi-agent actor-critic for mixed cooperative-competitive environments](#). *Preprint*, arXiv:1706.02275.
- Shayegan Omidshafiei, Christos Papadimitriou, Georgios Piliouras, Karl Tuyls, Mark Rowland, Jean-Baptiste Lespiau, Wojciech M. Czarnecki, Marc Lanctot, Julien Perolat, and Remi Munos. 2019. [\$\alpha\$ -rank: Multi-agent evaluation by evolution](#). *Preprint*, arXiv:1903.01373.
- OpenAI, :, Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Dębiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, Rafal Józefowicz, Scott Gray, Catherine Olsson, Jakub Pachocki, Michael Petrov, Henrique P. d. O. Pinto, Jonathan Raiman, and 8 others. 2019. [Dota 2 with large scale deep reinforcement learning](#). *Preprint*, arXiv:1912.06680.
- Karen Richmond, Satya Muddamsetty, Thomas Gammeltoft-Hansen, Henrik Olsen, and Thomas Moeslund. 2023. [Explainable ai and law: An evidential survey](#). *Digital Society*, 3.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. [Proximal policy optimization algorithms](#). *Preprint*, arXiv:1707.06347.

David Silver, Aja Huang, Christopher J. Maddison, Arthur Guez, Laurent Sifre, George van den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. 2016. [Mastering the game of go with deep neural networks and tree search](#). *Nature*, 529:484–503.

David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Simonyan, and Demis Hassabis. 2017. [Mastering chess and shogi by self-play with a general reinforcement learning algorithm](#). *Preprint*, arXiv:1712.01815.

Harry Surden. 2012. Computable contracts. *UC Davis Law Review*, 46:629.

Oriol Vinyals and 1 others. 2019. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575:350–354.

Haoxi Zhong, Chaojun Xiao, Cunchao Tu, Tianyang Zhang, Zhiyuan Liu, and Maosong Sun. 2019. [Jecqa: A legal-domain question answering dataset](#). *Preprint*, arXiv:1911.12011.

Appendix

A Rules-as-Code Examples

Tax regime excerpt. Gates and effects that suspend collection-related actions and extend a stay when certain sections are cited:

```
{
  "gates": {
    "collection_stay": {"blocks_actions": ["REQUEST_DOCS", "MOVE_COMPEL", "MOVE_SANCTIONS"]},
    "offshore_complexity": {"blocks_actions": ["MOVE_SANCTIONS"]}
  },
  "rules": [
    {
      "name": "tax_collection_stay",
      "when": {"action": "FILE_PROCEEDING", "conditions": [
        {"param": "proceeding_type", "op": "eq", "value": "tax_petition"},
        {"param": "section", "op": "in", "value": ["26 USC 6213", "26 USC 6330"], "optional": true}
      ]},
      "effects": [
        {"type": "set_gate", "gate": "collection_stay", "duration": 20},
        {"type": "add_cost", "who": "self", "amount": 2.0},
        {"type": "add_delay_credit", "k": 1}
      ]
    },
    {
      "name": "reference_tax_code_enforcement",
      "when": {"action": "REFERENCE_AUTHORITY", "conditions": [
        {"param": "code", "op": "in", "value": ["26 USC 6331", "26 U.S.C. 6331"]}
      ]},
      "effects": [
        {"type": "add_citation", "code_from_params": "code"},
        {"type": "extend_gate", "gate": "collection_stay", "inc": 3}
      ]
    }
  ]
}
```

Listing 1: Tax regime JSON excerpt

B Action Tokens and Observation Vector

Action token vocabulary (13 tokens).

```
NOOP, REQUEST_DOCS, OBJECT_REQUEST, FILE_MOTION, RESPOND_MOTION,
MOVE_COMPEL, MOVE_SANCTIONS, MEET_CONFER, SETTLEMENT_OFFER,
CHANGE_VENUE, WITHDRAW, FILE_PROCEEDING, REFERENCE_AUTHORITY
```

Listing 2: Action token set

C LLM Prompting Templates

JSON-constrained action proposal (few-shot).

System: You are an adversarial litigation strategy assistant. Return ONLY valid JSON.

User: ROLE: <plaintiff|defendant>

TACTIC_HINT: <optional>

OBS: [... 13-d observation ...]

Examples of good and bad moves (abstract tokens):

OBS:[...], TACTIC:SEEK_DISMISSAL -> GOOD: {"action":{"type":"FILE_MOTION","params":{"aggr":0.3}}}

OBS:[...], TACTIC:DELAY -> GOOD:

{"action":{"type":"FILE_PROCEEDING","params":{"proceeding_type":"bankruptcy","chapter":11,"forum":"BK"}}}

OBS:[...], TACTIC:TAX_STAY -> GOOD: {"action":{"type":"REFERENCE_AUTHORITY","params":{"code":"26 USC 6331","weight":0.7}}}

Return ONLY a JSON object with keys:

- reasoning: string

- action: {type: one of ACTION_TOKENS, params: {...}}

Listing 3: JSON-constrained LLM prompt

Free-text suggestion prompt.

System: You are an adversarial litigation strategy assistant. Respond concisely.

User: Given the legal discovery game, propose ONE next action for the indicated role.

Include the action token name from: NOOP, REQUEST_DOCS, ..., REFERENCE_AUTHORITY.

Optionally include params like custodians=10, complexity=0.7, aggr=0.4, amount=25.

Listing 4: Free-text LLM prompt

D Example Exploit Chain (Episode Trace)

Cost-inflating discovery loop (with procedural replies). Moderate, repeated REQUEST_DOCS raises the opponent's burden/fees; interleaved SETTLEMENT_OFFERS add leverage. The opponent's FILE_MOTION (e.g., protective) and RESPOND_MOTION replies keep the exchange inside ordinary procedure and below sanction/proportionality gates.

```
plaintiff_seq: [  
  {"type": "MEET_CONFER"},  
  {"type": "REQUEST_DOCS", "params": {"custodians": 10, "complexity": 0.6}},  
  {"type": "SETTLEMENT_OFFER", "params": {"amount": 100000, "importance": 0.8}},  
  {"type": "REQUEST_DOCS", "params": {"custodians": 12, "complexity": 0.6}},  
  {"type": "REQUEST_DOCS", "params": {"custodians": 8, "complexity": 0.55}}  
]  
  
defendant_seq: [  
  {"type": "FILE_MOTION", "params": {"kind": "protective", "aggr": 0.2}},  
  {"type": "RESPOND_MOTION"},  
  {"type": "RESPOND_MOTION"}  
]
```

Listing 5: Exploit chain with procedural replies