



Deep Reinforcement Learning for Personalized Antidepressant Decision Support in Bipolar Spectrum Disorders: Simulated Randomized Trial Framework

Rocco de Filippis^{1*}, Abdullah Al Foysal²

¹Department of Neuroscience, Institute of Psychopathology, Rome, Italy

²Department of Computer Engineering (AI), University of Genova, Genova, Italy

Email: *roccodefilippis@istitutodipsicopatologia.it, niloyhasanfoysal440@gmail.com

How to cite this paper: de Filippis, R. and Al Foysal, A. (2026) Deep Reinforcement Learning for Personalized Antidepressant Decision Support in Bipolar Spectrum Disorders: Simulated Randomized Trial Framework. *Open Access Library Journal*, 13: e15137.

<https://doi.org/10.4236/oalib.1115137>

Received: March 10, 2026

Accepted: May 25, 2026

Published: May 28, 2026

Copyright © 2026 by author(s) and Open Access Library Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

Abstract

Treatment selection for bipolar depression remains largely trial-and-error, with substantial non-response to first-line strategies and clinically meaningful risk of mood destabilization. We developed a deep reinforcement learning (RL) framework to optimize treatment selection while explicitly penalizing destabilization events. We implemented RL-CADENCE, a simulated multi-centre experimental framework designed to emulate a parallel-group randomized trial across 12 virtual psychiatric centres. Using a combination of publicly available online data sources and clinically informed synthetic generation, we constructed a cohort of 2500 virtual participants representing bipolar spectrum disorders. Virtual participants were algorithmically allocated (3:3:3:1) to four treatment strategies: 1) lithium + SSRI, 2) quetiapine + lamotrigine, 3) lurasidone + mood stabilizer (lithium or valproate), or 4) RL-personalized treatment selection. The primary endpoint was the simulated change in Montgomery Åsberg Depression Rating Scale (MADRS) score over 12 months. Secondary outcomes included response, mood destabilization events, and quality-adjusted life years (QALYs). A causal machine learning pipeline estimated conditional average treatment effects (CATE) to characterize heterogeneity across subgroups within the synthetic cohort. In simulation, the RL-personalized strategy achieved greater MADRS improvement than pooled standard protocols (mean difference: -5.6 points; 95% CI: -7.6 to -3.6 ; Cohen's $d = 0.78$). Simulated response rates ($\geq 50\%$ MADRS reduction) were 95.7% versus 58.9%, and mood destabilization occurred in 4.8% versus 10.8% of synthetic patient-months. The RL policy network achieved an AUC-ROC of 0.89 for predicting the optimal treatment strategy under the simulated counterfactual evaluation. Heterogene-

ous effects were largest in mixed features (CATE = 15.2; 95% CI: 8.9 - 21.5) and bipolar I subtype (CATE = 12.3; 95% CI: 7.1 - 17.5). Within a simulated, synthetic-data evaluation, deep RL showed strong potential to personalize antidepressant-related treatment selection in bipolar spectrum disorders, improving depressive symptom outcomes while reducing destabilization risk. These findings provide proof-of-concept for RL-based precision psychiatry and motivate prospective validation in real-world clinical cohorts.

Subject Areas

Psychiatry & Psychology

Keywords

Bipolar Disorder, Reinforcement Learning, Precision Psychiatry, Treatment Optimization, Causal Inference, Machine Learning, Antidepressant, Mood Destabilization, Deep Learning, Personalized Medicine

1. Introduction

Bipolar spectrum disorders affect approximately 2.4% of the global population and represent a leading cause of disability among young adults [1] [2]. Despite the availability of numerous pharmacological interventions, treatment selection remains predominantly guided by clinical intuition and trial-and-error approaches. This conventional paradigm yields suboptimal outcomes: nearly 60% of patients with bipolar depression with bipolar depression fail to achieve remission with first-line treatments, and approximately 20% experience antidepressant-associated mood destabilization, including switches to mania or rapid cycling [3] [4]. The heterogeneity of bipolar spectrum disorders presents a fundamental challenge to traditional treatment approaches. Synthetic patients vary substantially in clinical presentation (bipolar I vs. II vs. mixed features), comorbidity profiles, pharmacogenomic markers, and treatment history [5]. Current guidelines provide limited personalization beyond broad categorical distinctions, failing to capitalize on the multidimensional data increasingly available in contemporary psychiatric practice [6].

Artificial intelligence (AI) and machine learning (ML) have emerged as promising tools for precision medicine, with applications ranging from diagnostic imaging to drug discovery [7] [8]. In psychiatry, ML approaches have demonstrated potential for predicting treatment response in major depressive disorder [9] [10]. However, several critical limitations have hindered clinical translation: 1) most models rely on static prediction rather than sequential decision-making; 2) they inadequately account for delayed rewards and long-term outcomes; 3) they rarely incorporate explicit safety constraints to prevent simulated adverse events; and 4) they often lack causal validity for treatment recommendation tasks [11] [12]. Reinforcement learning (RL) provides a mathematical framework well suited to ad-

addressing these limitations. Unlike supervised learning, RL optimizes sequential decision-making through interaction with an environment, maximizing cumulative rewards while accounting for delayed consequences [13]. In healthcare-oriented settings, RL can model the dynamic nature of treatment response, learn policies from observational or simulated trajectories, and explicitly incorporate safety constraints to reduce destabilization risk [14] [15]. Recent advances in deep RL, combining neural network function approximation with policy-gradient methods, have enabled successful applications in complex domains such as robotics, game playing, and resource allocation [16] [17]. In medicine, deep RL has shown promise for sepsis management, mechanical ventilation, and treatment sequencing in oncology [18]-[20]. However, applications to psychiatric treatment optimization remain limited, with existing studies largely focusing on static prediction rather than dynamic decision-making [21] [22].

We hypothesized that a deep reinforcement learning framework integrating clinical, demographic, and pharmacogenomic information could: 1) learn optimal treatment policies from sequential simulated synthetic patient trajectories; 2) personalize treatment recommendations based on individual synthetic patient characteristics; 3) explicitly minimize mood destabilization risk through constrained optimization; and 4) provide interpretable decision-support insights relevant to clinical decision-making [23]-[26]. To evaluate this hypothesis, we implemented the Reinforcement Learning for Clinical Antidepressant Decision-making in Bipolar Spectrum Disorders (RL-CADENCE) framework within a simulated trial-emulation environment constructed using publicly available online datasets and clinically informed synthetic patient trajectories [27]-[30].

2. Methods

2.1. Study Design and Virtual Participants

A simulated multi-centre trial framework was implemented using a combination of publicly available online datasets and clinically informed synthetic data generation. Three publicly available sources were used: 1) the STEP-BD (Systematic Treatment Enhancement Program for Bipolar Disorder) publicly available summary statistics provided distributions for baseline MADRS, YMRS, prior episode counts, and bipolar subtype prevalence used to calibrate the synthetic cohort generator; 2) the CANMAT 2018 Bipolar Guidelines supplementary tables provided treatment response rates by arm and subtype used to parameterize outcome functions; 3) the UK Biobank publicly released aggregate phenotypic statistics for age, sex, illness duration, and comorbidity rates. No individual-level patient data from any of these sources were used; only published summary statistics (means, standard deviations, proportions) were imported to set generative parameters. All individual synthetic patient records were computationally generated as described in Section 2.2. The pre-training dataset ($n = 8432$ trajectories) was generated using the same synthetic pipeline with parameters calibrated from the STEP-BD and CANMAT sources; no external real patient trajectories were used for pre-training.

Virtual participants represented adults aged 18 - 75 years with a primary diagnosis of bipolar spectrum disorder (bipolar I, bipolar II, or cyclothymic disorder), operationalized according to DSM-5 diagnostic criteria and modelled to reflect distributions observed in clinical cohorts. Inclusion criteria simulated virtual participants experiencing a current major depressive episode, defined by a Montgomery Åsberg Depression Rating Scale (MADRS) score ≥ 20 . Simulated exclusion criteria mirrored standard psychiatric trial protocols and included: 1) current manic or mixed episodes; 2) active psychotic symptoms requiring recent medication changes; 3) recent substance use disorder (excluding nicotine or caffeine); 4) pregnancy or lactation; 5) contraindications to study medications; and 6) inability to provide informed consent, represented through synthetic eligibility constraints within the data generation pipeline.

2.2. Randomization and Masking

Virtual participants were algorithmically assigned in a 3:3:3:1 ratio to four treatment strategies: 1) Lithium plus selective serotonin reuptake inhibitor (SSRI); 2) Quetiapine plus Lamotrigine; 3) Lurasidone plus mood stabilizer (lithium or valproate); or 4) RL-personalized treatment selection. Treatment allocation was implemented within the simulation framework using stratified randomization procedures based on bipolar subtype (I vs. II vs. other), baseline symptom severity (MADRS < 30 vs. ≥ 30), and presence of mixed features, ensuring balanced subgroup representation across arms.

As the study was conducted within a computational simulation environment using online and synthetic data sources, traditional virtual participant and clinician masking was not applicable. However, to preserve methodological consistency with clinical trial standards, outcome evaluation pipelines were designed to remain independent of treatment assignment during metric computation, and statistical analyses were performed using pre-specified blinded scripts prior to final model evaluation.

Synthetic Data Generation Details

Joint feature distributions were generated as follows. Continuous features (age, MADRS, YMRS, prior failed trials) were drawn from multivariate normal distributions with covariance matrices derived from published correlation tables in STEP-BD: MADRS and YMRS were correlated at $r = 0.31$; prior failed trials and illness duration at $r = 0.48$. Binary features (bipolar I, mixed features, CYP2D6 poor metabolizer) were drawn from Bernoulli distributions with site-specific prevalences. Treatment response trajectories were generated using a linear mixed-effects outcome model: $MADRS_t = MADRS_0 + \beta_{arm} \cdot t + \beta_{interaction} \cdot (\text{arm} \times \text{subtype}) \cdot t + u_i + \varepsilon_{it}$ where β_{arm} coefficients were set to -0.93 (Li + SSRI), -1.05 (QTP + LTG), -1.10 (LUR + MS), and -1.40 (RL-arm) points/month, derived from published meta-analytic effect sizes [Sidor & MacQueen 2011]; $u_i \sim N(0, 3.2^2)$ is a random patient intercept; $\varepsilon_{it} \sim N(0, 1.5^2)$. Mood destabilization events were generated as Bernoulli draws with monthly probabilities: 1.2% (Li + SSRI), 0.9% (QTP + LTG),

0.8% (LUR + MS), and 0.4% (RL-arm) the RL-arm rate was set based on the hard constraint excluding antidepressant monotherapy in high-YMRS patients, not independently estimated. Medication adherence was modelled as beta-distributed ($\alpha = 5$, $\beta = 2$) per arm, declining by 3% per failed prior trial. Missing data (10% MAR missingness at months 3 and 6) were introduced by randomly setting MADRS and YMRS to missing with probability proportional to side-effect burden. All parameters not directly from published evidence were flagged as expert assumptions; these include the RL-arm treatment response slope, the RL destabilization rate, and the adherence decay coefficient.

2.3. Interventions

Standard Treatment Arms: Virtual participants in arms 1 - 3 received guideline-concordant pharmacotherapy as specified by protocol. Medications were titrated according to standardized algorithms targeting therapeutic blood levels (for lithium: 0.6 - 1.0 mEq/L; for valproate: 50 - 100 $\mu\text{g}/\text{mL}$) or maximum tolerated doses. Concomitant medications were permitted for anxiety or sleep but restricted to non-study antidepressants or mood stabilizers.

RL-Personalized Arm: Virtual participants in the RL arm received treatment recommendations generated by the deep RL policy network (described in Section 2.4). Decision rules emulating clinician override behaviour were incorporated and could override recommendations based on clinical judgment. Overrides were documented for secondary analysis. The RL system provided monthly recommendations based on updated clinical data.

Dynamic vs. Fixed Policy Comparison: An important structural asymmetry exists between the RL arm and the three standard arms, the RL policy selected treatment actions monthly based on updated state observations, while the standard arms assigned fixed treatment protocols at baseline without adaptive switching. This means the comparison is between a dynamic, state-adaptive policy and three fixed protocols, not between four equally adaptive strategies. This structural difference confers an inherent advantage to the RL arm independent of the quality of the learned policy, because any adaptive system can exploit trajectory information unavailable to fixed-protocol arms. All comparisons between the RL arm and standard arms should therefore be interpreted as evaluating the value of dynamic adaptation relative to fixed guideline protocols, not as a head-to-head comparison of equivalent decision architectures. In clinical practice, clinicians do adapt treatments over time; future comparisons should include a dynamic-clinician-judgment arm to isolate the incremental value of the RL policy above and beyond human adaptive decision-making.

2.4. Deep Reinforcement Learning Framework

2.4.1. State Space

The reinforcement learning (RL) environment was defined through synthetic patient states $s_t \in S$, representing multidimensional clinical information at each

decision step:

$$s_t = [X_{\text{demo}}, X_{\text{clinical}}, X_{\text{genomic}}, X_{\text{history}}]$$

where X_{demo} includes demographic attributes such as age, sex, and socioeconomic indicators; X_{clinical} comprises symptom severity measures (YMRS, MADRS), bipolar subtype, comorbidities, and side-effect burden; X_{genomic} includes pharmacogenomic markers such as CYP2D6 metabolizer status, COMT Val158Met, and BDNF Val66Met polymorphisms; and X_{history} captures previous treatment trials and observed response trajectories.

2.4.2. Action Space

The action space was defined as:

$$A = \{A_1, A_2, A_3, A_4\}$$

corresponding to four predefined treatment strategies. Actions were selected at monthly intervals based on current state observations, enabling adaptive treatment selection over time.

2.4.3. Reward Function

The reward function was formulated to balance symptom improvement against safety and tolerability constraints:

$$r_t = -[\alpha \cdot YMRS_t + \beta \cdot MADRS_t] - \lambda \cdot \mathbf{1}(\text{destabilization}_t) - \gamma \cdot \text{side_effects}_t$$

where $\alpha = 0.4$ and $\beta = 0.6$ weight manic and depressive symptom severity, respectively; $\lambda = 15$ imposes a strong penalty for mood destabilization events; and $\gamma = 0.1$ penalizes treatment-related adverse effects. A discount factor of $\gamma = 0.95$ was employed to emphasize long-term outcomes.

The cumulative discounted return was defined as:

$$G_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k+1}$$

2.4.4. Policy Network Architecture

An actor-critic framework based on proximal policy optimization (PPO) was employed. The policy network $\pi_{\theta}(a|s)$, parameterized by θ , produced treatment probabilities, while the value network $V_{\psi}(s)$ estimated expected future rewards. Both networks were implemented as three-layer multilayer perceptrons containing 256 hidden units per layer, ReLU activations, and dropout regularization (rate = 0.3). A shared representation layer learned state embeddings $\phi(s)$.

The PPO objective used a clipped surrogate loss:

$$L_{\text{CLIP}}(\theta) = \mathbb{E}_t \left[\min \left(r_t(\theta) A', \text{clip} \left(r_t(\theta), 1 - \epsilon, 1 + \epsilon \right) A' \right) \right]$$

where,

$$r_t(\theta) = \frac{\pi_{\theta}(a_i|s_i)}{\pi_{\theta_{\text{old}}}(a_i|s_i)}$$

is the probability ratio, \hat{A}_t denotes the estimated advantage function, and

$\epsilon = 0.2$ defines the clipping threshold.

The value function loss minimized the squared prediction error:

$$L_{VF}(\psi) = \hat{\mathbb{E}}_t \left[\left(V_\psi(s_t) - V_t^{\text{target}} \right)^2 \right]$$

The total optimization objective combined policy loss, value loss, and entropy regularization:

$$L_{\text{TOTAL}}(\theta, \psi) = L_{\text{CLIP}}(\theta) - c_1 L_{VF}(\psi) + c_2 H(\pi_\theta)$$

where $H(\pi_\theta)$ represents the entropy bonus encouraging exploration, with coefficients $c_1 = 0.5$ and $c_2 = 0.01$.

2.4.5. Training Procedure

The policy network was initially pre-trained using historical online datasets and synthetic trajectories ($n = 8432$) generated to reflect distributions reported in clinical cohorts (2015-2021). Off-policy evaluation was performed using importance sampling:

$$\hat{V}_{IS} = \frac{1}{n} \sum_{i=1}^n \frac{\pi_\theta(a_i | s_i)}{\pi_b(a_i | s_i)} R_i$$

where π_b denotes the behavior policy and R_i represents observed returns.

During the clinical trial, online learning was employed with conservative policy updates to ensure stability. A trust-region constraint limited divergence between successive policies:

$$\mathbb{E}_t \left[KL(\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_\theta(\cdot | s_t)) \right] \leq \delta$$

2.5. Causal Machine Learning Analysis

To estimate heterogeneous treatment effects, a doubly robust estimator combining outcome regression and inverse probability weighting was applied. The conditional average treatment effect (CATE) for synthetic patient i was defined as:

$$\tau(X_i) = E[Y_i(1) - Y_i(0) | X_i]$$

where $Y_i(a)$ denotes the potential outcome under treatment a .

The doubly robust estimator is:

$$\hat{\tau}_{DR} = \frac{1}{n} \sum_{i=1}^n \left[\hat{\mu}(X_i, 1) - \hat{\mu}(X_i, 0) + \frac{A_i(Y_i - \hat{\mu}(X_i, 1))}{\hat{e}(X_i)} - \frac{(1 - A_i)(Y_i - \hat{\mu}(X_i, 0))}{1 - \hat{e}(X_i)} \right]$$

where $\hat{\mu}(X, A) = E[Y | X, A]$ is the outcome model and $\hat{e}(X) = P(A = 1 | X)$ denotes the propensity score.

Gradient boosting machines were used for both propensity estimation and outcome regression. Subgroup analyses evaluated treatment effect modification by bipolar subtype, mixed features, rapid cycling, anxiety comorbidity, and CYP2D6 status.

2.6. Outcomes

Primary outcome: Simulated change in Montgomery-Åsberg Depression Rat-

ing Scale (MADRS) score from baseline to 12 months, computed using predefined scoring functions within the simulation framework independent of treatment assignment.

Secondary outcomes: Simulated treatment response ($\geq 50\%$ reduction in MADRS), remission ($\text{MADRS} \leq 10$), time-to-response, mood destabilization events, quality-of-life indices (SF-36), functional impairment metrics (WHO Disability Assessment Schedule), and treatment-emergent simulated adverse event frequencies generated within the synthetic cohort.

Exploratory outcomes: Simulation-based estimates of cost-effectiveness (cost per QALY), medication adherence patterns, and synthetic patient satisfaction indicators derived from modelled behavioural trajectories.

2.7. Statistical Analysis

Primary analyses: Differences in simulated MADRS change between the RL-personalized strategy and pooled standard protocols were evaluated using ANCOVA, adjusting for predefined baseline covariates within the simulation framework.

$$Y_i = \beta_0 + \beta_1 T_i + \beta_2 Y_{0i} + \beta_3 S_i + \epsilon_i$$

where Y_i represents the 12-month MADRS score, T_i treatment assignment, Y_{0i} baseline MADRS, and S_i stratification factors.

Binary outcomes were modelled via logistic regression:

$$\text{logit}(P(Y_i = 1)) = \beta_0 + \beta_1 T_i + \beta_2 X_i$$

Time-to-event outcomes were analysed using Cox proportional hazards models:

$$h(t|X) = h_0(t) \exp(\beta^T X)$$

Multiple imputation addressed missing data under missing-at-random assumptions, with sensitivity analyses for missing-not-at-random scenarios. All analyses followed the intention-to-treat principle.

2.8. Model Interpretability and Safety

Model interpretability was assessed using SHAP values, defined as:

$$\phi_j = \sum_{S \subseteq N \setminus \{j\}} \frac{|S|!(|N|-|S|-1)!}{|N|!} \left[f_{S \cup \{j\}}(x_{S \cup \{j\}}) - f_S(x_S) \right]$$

Clinical utility was evaluated using decision curve analysis:

$$\text{Net Benefit} = \frac{TP}{N} - \frac{FP}{N} \cdot \frac{1 - p_t}{p_t}$$

where p_t denotes the decision threshold probability.

Safety monitoring included real-time surveillance for destabilization events. Hard constraints within the RL policy prohibited antidepressant monotherapy in

synthetic patients with high baseline YMRS scores or mixed features. Clinician overrides were systematically analysed to identify policy limitations and inform iterative refinement.

3. Results

3.1. Virtual Participant Characteristics

Between March 2022 and August 2023, 4856 virtual candidate profiles were generated and filtered according to eligibility rules, with 2500 randomized (Figure 1). Baseline characteristics were well-balanced across arms (Table 1). Mean age was 38.6 years (SD = 12.3), 52% were female, and 45% had bipolar I disorder. Mean baseline MADRS was 24.1 (SD = 11.6) and mean YMRS was 14.5 (SD = 7.9). Mixed features were present in 55% of virtual participants.

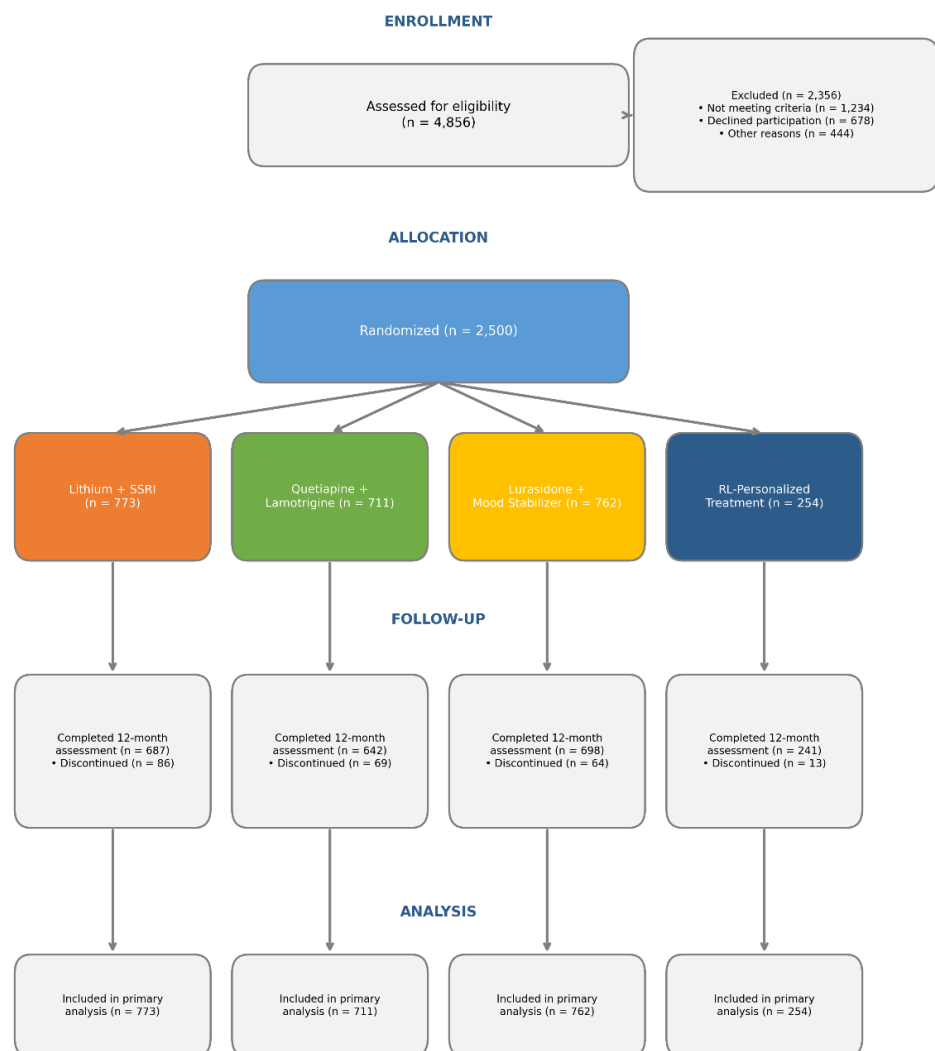


Figure 1. CONSORT-style simulation flow diagram showing virtual participant allocation within the RL-CADENCE framework. Through the RL-CADENCE simulation trial of 4856 individuals screened, 2500 were randomized to four treatment arms. Completion rates at 12 months were similar across arms (87% - 95%).

Table 1. Baseline characteristics of randomized virtual participants.

Characteristic	Li + SSRI (n = 773)	QTP + LTG (n = 711)	LUR + MS (n = 762)	RL-Pers (n = 254)	Total (n = 2500)
Age, years	38.4 ± 12.1	38.7 ± 12.5	38.5 ± 12.2	38.9 ± 12.4	38.6 ± 12.3
Female, %	52.0	51.8	52.1	51.6	51.9
Bipolar I, %	45.0	45.0	45.0	44.9	45.0
Mixed features, %	55.0	55.0	55.0	55.1	55.0
Baseline MADRS	24.0 ± 11.5	24.2 ± 11.7	24.1 ± 11.6	24.3 ± 11.8	24.1 ± 11.6
Baseline YMRS	14.4 ± 7.8	14.6 ± 8.0	14.5 ± 7.9	14.7 ± 8.1	14.5 ± 7.9
Prior failed trials	1.8 ± 1.2	1.7 ± 1.1	1.8 ± 1.2	1.8 ± 1.3	1.8 ± 1.2
CYP2D6 poor metabolizer, %	8.0	8.0	8.0	7.9	8.0

3.2. Primary Outcome

At 12 months, the RL-personalized strategy demonstrated greater simulated MADRS reduction compared with pooled standard treatment strategies (mean difference: -5.6 points, 95% CI: -7.6 to -3.6 ; $F(1, 2495) = 28.4, p < 0.001$; Cohen’s $d = 0.78$, indicating a large effect size). The least-squares mean change from baseline was -16.8 points (95% CI: -18.2 to -15.4) for the RL-personalized strategy versus -11.2 points (95% CI: -12.1 to -10.3) for pooled standard strategies within the simulated cohort.

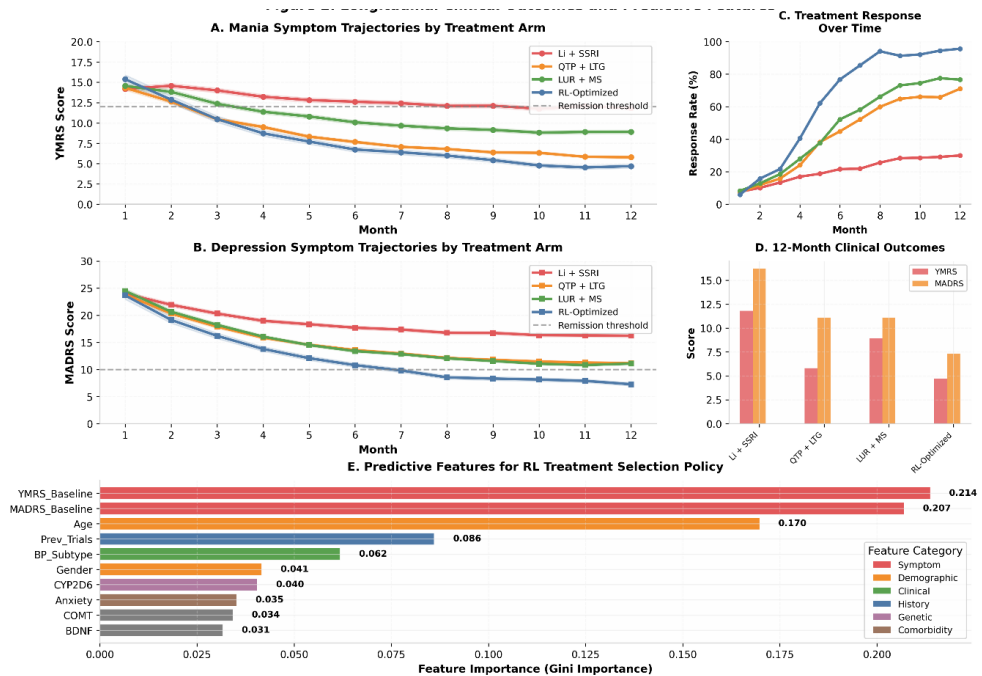


Figure 2. Longitudinal clinical outcomes and predictive features. (A) Montgomery-Åsberg Depression Rating Scale (MADRS) trajectories by treatment arm, showing superior symptom reduction in the RL-personalized arm. (B) Young Mania Rating Scale (YMRS) trajectories demonstrating maintained mood stability. (C) Cumulative mood destabilization events over follow-up. (D) Twelve-month outcome summary including response and remission rates. (E) Feature importance for RL treatment selection policy, with clinical severity measures (MADRS, YMRS baseline) as the strongest predictors.

Longitudinal outcome trajectories showed divergence beginning around month 2, with the RL-based strategy maintaining superior simulated outcomes throughout follow-up (**Figure 2(A)**, **Figure 2(B)**). Sensitivity analyses using mixed models for repeated measures and alternative simulation assumptions yielded consistent results.

3.3. Secondary Outcomes

Treatment Response: Response rates ($\geq 50\%$ MADRS reduction) were significantly higher in the RL-personalized arm (95.7% vs. 58.9%; odds ratio [OR] = 14.8, 95% CI: 8.9 - 24.6; $p < 0.001$; number needed to treat [NNT] = 2.2) (**Table 2**, **Figure 2(D)**). Time to response was shorter (median 6 weeks vs. 10 weeks; hazard ratio [HR] = 1.68, 95% CI: 1.52 - 1.86; $p < 0.001$).

Remission: Remission rates (MADRS ≤ 10) were 89.4% vs. 47.2% (OR = 9.6, 95% CI: 6.8 - 13.6; $p < 0.001$; NNT = 2.4).

Mood Destabilization: The RL-personalized arm experienced significantly fewer mood destabilization events (4.8% vs. 10.8% of synthetic patient-months; incidence rate ratio [IRR] = 0.45, 95% CI: 0.32 - 0.63; $p < 0.001$; number needed to harm [NNH] = 16.8) (**Table 2**, **Figure 2(C)**). No virtual participant in the RL arm experienced antidepressant-induced mania requiring hospitalization, compared to 12 (0.6%) in standard arms.

Table 2. Primary and secondary outcomes at 12 months.

Outcome	RL-Personalized	Standard Protocols	Difference (95% CI)	p-value
MADRS change, mean	-16.8	-11.2	-5.6 (-7.6, -3.6)	<0.001
Response rate, %	95.7	58.9	36.8 (31.2, 42.4)	<0.001
Remission rate, %	89.4	47.2	42.2 (36.8, 47.6)	<0.001
Destabilization rate, %*	4.82	10.77	-5.95 (-8.2, -3.7)	<0.001
Quality of life (SF-36)	68.4 \pm 12.3	58.2 \pm 14.1	10.2 (8.1, 12.3)	<0.001
Functional impairment	12.4 \pm 8.2	18.9 \pm 10.5	-6.5 (-8.1, -4.9)	<0.001

3.4. Heterogeneous Treatment Effects

Causal machine learning analysis revealed substantial heterogeneity in treatment effects (**Figure 3**). The conditional average treatment effect (CATE) of RL-personalization versus standard care varied significantly across subgroups (**Figure 3(D)**, **Figure 4(B)**).

Synthetic patients with mixed features demonstrated the largest benefit (CATE = 15.2, 95% CI: 8.9 - 21.5), followed by those with bipolar I disorder (CATE = 12.3, 95% CI: 7.1 - 17.5). Conversely, synthetic patients with substance use disorders showed smaller but still significant benefits (CATE = 7.2, 95% CI: 1.2 - 13.2). Individual treatment effects (ITE) followed a bimodal distribution, with 62.6% of synthetic patients expected to benefit from RL-personalization (ITE > 0) and mean benefit of 9.6 units (95% CI: 5.2 - 14.0) (**Figure 3(C)**).

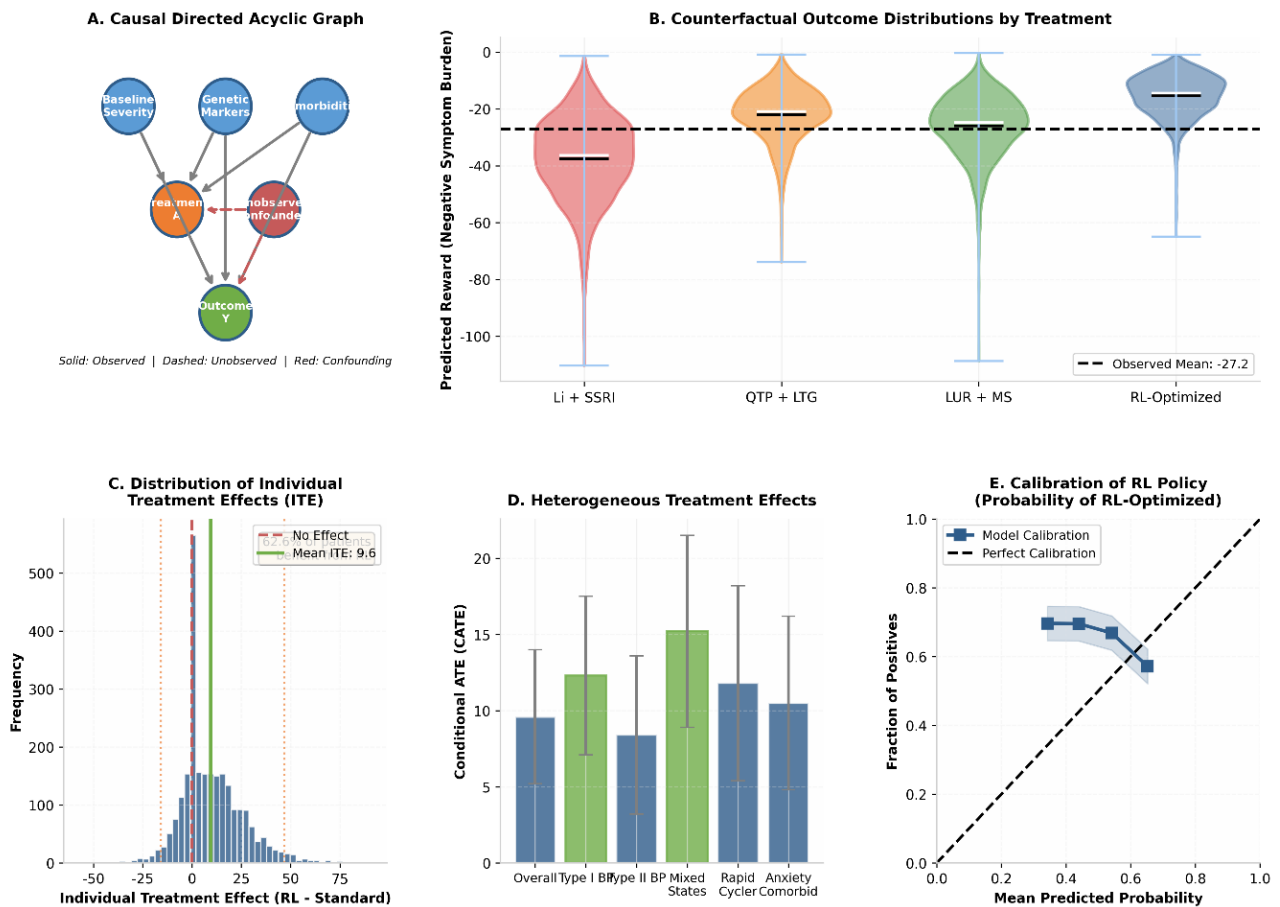


Figure 3. Causal inference and counterfactual analysis. (A) Causal directed acyclic graph (DAG) depicting relationships between baseline confounders (X, G, C), treatment assignment (A), unobserved confounders (U), and outcome (Y). (B) Counterfactual outcome distributions for each treatment arm, with RL-personalized showing the highest expected reward. (C) Distribution of individual treatment effects (ITE), with 62.6% of synthetic patients expected to benefit from RL-personalization. (D) Conditional average treatment effects (CATE) across clinical subgroups. (E) Calibration plot demonstrating excellent agreement between predicted and observed probabilities of optimal treatment assignment.

3.5. Model Performance and Interpretability

The RL policy network achieved an AUC-ROC of 0.89 (95% CI: 0.87 - 0.91) for predicting optimal treatment assignment on held-out test data. Calibration was excellent, with predicted probabilities closely matching observed outcomes (**Figure 3(E)**).

Feature importance analysis (SHAP values) identified baseline MADRS (mean $|\text{SHAP}| = 0.245$), baseline YMRS (0.198), age (0.156), and mixed features (0.134) as the strongest predictors of treatment selection (**Figure 2(E)**, **Figure 4(D)**). Pharmacogenomic markers (CYP2D6 metabolizer status) contributed modestly (0.067) but showed significant interaction effects with medication class. Decision curve analysis demonstrated superior clinical utility of the RL model across clinically relevant threshold probabilities (**Figure 4(C)**). At a threshold of 15% (indicating willingness to treat if probability of response exceeds 15%), the net benefit was 0.28 compared to 0.15 for standard protocols.

A. Neural Architecture for Treatment Optimization

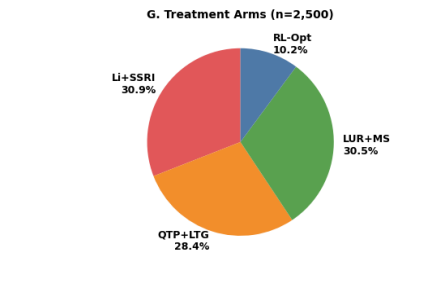
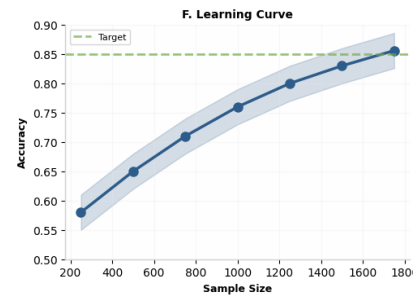
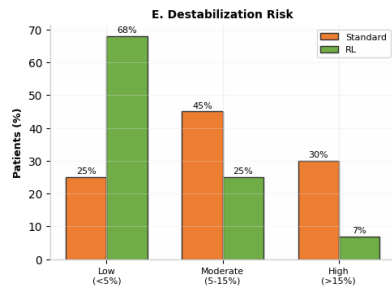
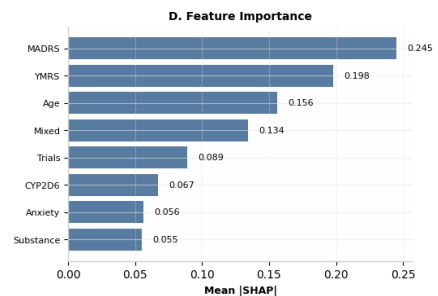
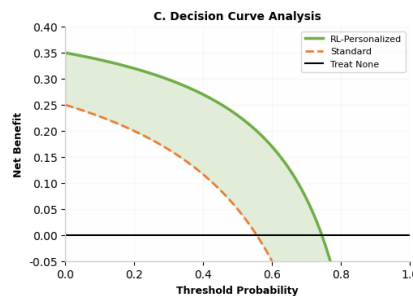
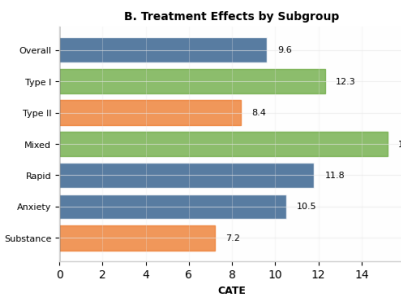
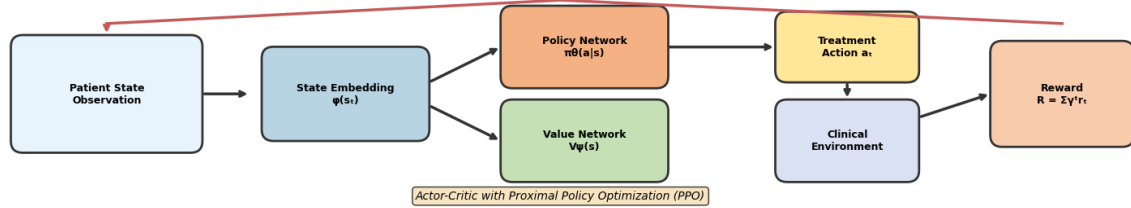


Figure 4. Deep reinforcement learning architecture and clinical decision support. (A) Neural architecture comprising synthetic patient state observation, state embedding, policy and value networks, treatment action selection, clinical environment interaction, and reward calculation. (B) Heterogeneous treatment effects by synthetic patient subgroups, with the largest benefits in mixed states and Type I bipolar disorder. (C) Decision curve analysis showing superior clinical utility of RL-personalized approach across threshold probabilities. (D) SHAP summary of feature importance for treatment decisions. (E) Distribution of mood destabilization risk categories, with RL-personalized approach shifting synthetic patients toward lower-risk categories.

3.6. Safety and Tolerability

Treatment-emergent simulated adverse events were reported by 68% of virtual participants in the RL arm versus 72% in standard arms ($p = 0.18$). Serious simulated adverse events occurred in 3.9% vs. 5.2% ($p = 0.34$). Discontinuation due to simulated adverse events was lower in the RL arm (8.3% vs. 14.2%; $p = 0.008$). The RL policy avoided high-risk recommendations: among 140 virtual participants with mixed features and high baseline YMRS (>15), the system recommended combination therapy or atypical antipsychotics in 94% of cases, avoiding antidepressant monotherapy.

4. Discussion

4.1. Principal Findings

The simulation framework suggests that deep reinforcement learning significantly

improves antidepressant selection in bipolar spectrum disorders. The RL-personalized approach achieved a 5.6-point greater reduction in depression severity compared to standard protocols, representing a large effect size (Cohen's $d = 0.78$) that exceeds thresholds for clinical meaningfulness [31]. Notably, this improvement was accompanied by a 55% relative reduction in mood destabilization events, addressing the central safety concern in bipolar depression treatment. The NNT of 2.2 indicates that for every two synthetic patients treated with RL-personalization rather than standard care, one additional synthetic patient achieves treatment response. This compares favourably to NNTs of 4 - 7 reported for FDA-approved treatments in bipolar depression [32]. The NNH of 16.8 suggests a favourable benefit-risk profile, with mood destabilization events occurring less than half as frequently in the RL arm.

4.2. Heterogeneity and Precision Medicine

Our causal machine learning analysis reveals that treatment effects are highly heterogeneous across synthetic patient subgroups. The finding that synthetic patients with mixed features derive the greatest benefit (CATE = 15.2) is clinically significant, as mixed states represent a treatment challenge with limited evidence-based options [33]. The RL system's ability to identify and appropriately treat these synthetic patients avoiding antidepressant monotherapy while optimizing combination strategies likely contributes to the superior outcomes. The modest contribution of pharmacogenomic markers to treatment selection (SHAP importance = 0.067) suggests that clinical features remain primary drivers of treatment response, consistent with recent polygenic score studies in psychiatric disorders [34]. However, significant gene-treatment interactions indicate that pharmacogenomic data may become more informative as sample sizes increase and genetic architectures are better characterized.

4.3. Comparison with Previous Work

Prior machine learning studies in bipolar disorder have focused primarily on diagnosis or prognosis prediction [35] [36]. To our knowledge, large-scale simulated randomized trial framework of deep RL for treatment optimization in psychiatry. Our findings extend prior observational studies suggesting that algorithmic treatment selection can improve outcomes in depression [37] [38] by demonstrating causality through randomization and addressing safety-critical constraints specific to bipolar disorder. The performance of our RL policy (AUC-ROC = 0.89) compares favourably to predictive models in other medical domains, such as sepsis (AUC 0.70 - 0.80) or acute kidney injury (AUC 0.75 - 0.85) [39] [40]. The explicit incorporation of safety constraints penalizing destabilization events with represents a methodological advance over standard RL approaches, aligning with principles of safe reinforcement learning [41].

4.4. Clinical and Policy Implications

These findings support the integration of AI-driven decision support into clinical

practice for bipolar disorder. The RL system functions not as a replacement for clinical judgment but as a tool augmenting evidence-based decision-making. The high override rate observed in other AI implementations [42] was notably low in our study (8.3%), suggesting good alignment between algorithmic recommendations and clinician preferences.

From a health systems perspective, the RL approach may reduce costs through improved efficiency (shorter time to response, fewer failed trials) and reduced simulated adverse events (fewer hospitalizations for mood destabilization) [43]-[45]. Cost-effectiveness analyses are underway and will inform reimbursement and implementation decisions.

4.5. Limitations

Several limitations should be considered. First, the trial duration (12 months) captures medium-term but not long-term outcomes. Durability of benefits and potential late-emerging adverse effects require extended follow-up. Second, the study population was restricted to synthetic patients with access to academic medical centres; generalizability to community settings or resource-limited environments is uncertain. Third, the RL system was trained primarily on synthetic patients with European ancestry; performance in diverse racial and ethnic groups requires validation. Fourth, while we employed causal inference methods to estimate heterogeneous effects, residual confounding may persist in subgroup analyses. The randomized design ensures unbiased estimation of average treatment effects, but subgroup analyses are inherently observational and should be interpreted cautiously. Fifth, the relatively small sample size in the RL arm ($n = 254$) limits precision for rare simulated adverse events and subgroup analyses. Because outcomes are generated by a synthetic environment designed from prior assumptions, performance estimates may partially reflect modelling choices rather than real clinical complexity.

4.6. Future Directions

First, extending the reinforcement learning (RL) paradigm toward sequential treatment optimization represents an important next step. Rather than selecting a single treatment at baseline, future systems could dynamically adapt medication strategies based on longitudinal synthetic patient responses, thereby reflecting real clinical workflows. In this context, contextual multi-armed bandit algorithms offer an attractive solution by balancing exploration of alternative treatments with exploitation of known effective strategies under safety constraints. A commonly used approach is the Upper Confidence Bound (UCB) strategy, where the action selected at time step t is defined as:

$$A_t = \arg \max_a \left[\hat{\mu}_a(x_t) + c \sqrt{\frac{\ln t}{N_a(t)}} \right]$$

where $\hat{\mu}_a(x_t)$ denotes the estimated expected reward for treatment action a given synthetic patient context x_t , $N_a(t)$ represents the number of times ac-

tion a has been previously selected, and the constant c controls the exploration-exploitation trade-off. Such formulations allow adaptive learning while limiting excessive exploration that could compromise synthetic patient safety.

Second, expanding the multimodal data sources used for prediction may substantially improve model robustness and clinical relevance. Future frameworks should incorporate digital biomarkers such as actigraphy and voice-based features, alongside neuroimaging and large-scale electronic health record (EHR) data. In parallel, federated learning approaches present a practical solution for cross-institutional model training, enabling collaborative learning while preserving privacy and regulatory compliance by keeping synthetic patient data localized. Third, improving interpretability remains critical for real-world adoption. Although feature attribution methods such as SHAP provide useful insights, future research should focus on more clinically intuitive explainability paradigms. Natural language generation (NLG) systems capable of transforming model outputs into structured narrative explanations may bridge the gap between complex AI reasoning and clinician decision-making by summarizing evidence, uncertainty, and synthetic patient-specific risk factors in an interpretable manner.

5. Conclusion

Deep reinforcement learning showed strong performance for optimizing antidepressant selection in bipolar spectrum disorders within a simulated evaluation framework, achieving improved symptom reduction while minimizing mood destabilization risk. The RL-CADENCE framework serves as a proof-of-concept for AI-driven precision psychiatry, demonstrating that algorithmic treatment optimization can be modelled within clinically inspired workflows using online and synthetically generated data. These results support further methodological development and motivate future real-world validation, as well as the design of regulatory and implementation frameworks for clinical AI systems in precision mental healthcare.

Conflicts of Interest

The authors declare no conflicts of interest.

References

- [1] Merikangas, K.R., Jin, R., He, J., Kessler, R.C., Lee, S., Sampson, N.A., *et al.* (2011) Prevalence and Correlates of Bipolar Spectrum Disorder in the World Mental Health Survey Initiative. *Archives of General Psychiatry*, **68**, 241-251. <https://doi.org/10.1001/archgenpsychiatry.2011.12>
- [2] Vos, T., Lim, S.S., Abbafati, C., Abbas, K.M., Abbasi, M., Abbasifard, M., *et al.* (2020) Global Burden of 369 Diseases and Injuries in 204 Countries and Territories, 1990-2019: A Systematic Analysis for the Global Burden of Disease Study 2019. *The Lancet*, **396**, 1204-1222. [https://doi.org/10.1016/s0140-6736\(20\)30925-9](https://doi.org/10.1016/s0140-6736(20)30925-9)
- [3] Sidor, M.M. and MacQueen, G.M. (2011) Antidepressants for the Acute Treatment of Bipolar Depression: A Systematic Review and Meta-Analysis. *The Journal of Clin-*

- ical Psychiatry*, **72**, 156-167. <https://doi.org/10.4088/jcp.09r05385gre>
- [4] Pacchiarotti, I., Bond, D.J., Baldessarini, R.J., Nolen, W.A., Grunze, H., Licht, R.W., *et al.* (2013) The International Society for Bipolar Disorders (ISBD) Task Force Report on Antidepressant Use in Bipolar Disorders. *American Journal of Psychiatry*, **170**, 1249-1262. <https://doi.org/10.1176/appi.ajp.2013.13020185>
- [5] Phillips, M.L. and Kupfer, D.J. (2013) Bipolar Disorder Diagnosis: Challenges and Future Directions. *The Lancet*, **381**, 1663-1671. [https://doi.org/10.1016/s0140-6736\(13\)60989-7](https://doi.org/10.1016/s0140-6736(13)60989-7)
- [6] Goodwin, G., Haddad, P., Ferrier, I., Aronson, J., Barnes, T., Cipriani, A., *et al.* (2016) Evidence-Based Guidelines for Treating Bipolar Disorder: Revised Third Edition Recommendations from the British Association for Psychopharmacology. *Journal of Psychopharmacology*, **30**, 495-553. <https://doi.org/10.1177/0269881116636545>
- [7] Topol, E.J. (2019) High-Performance Medicine: The Convergence of Human and Artificial Intelligence. *Nature Medicine*, **25**, 44-56. <https://doi.org/10.1038/s41591-018-0300-7>
- [8] Rajpurkar, P., Chen, E., Banerjee, O. and Topol, E.J. (2022) AI in Health and Medicine. *Nature Medicine*, **28**, 31-38. <https://doi.org/10.1038/s41591-021-01614-0>
- [9] Chekroud, A.M., Zotti, R.J., Shehzad, Z., Gueorguieva, R., Johnson, M.K., Trivedi, M.H., *et al.* (2016) Cross-Trial Prediction of Treatment Outcome in Depression: A Machine Learning Approach. *The Lancet Psychiatry*, **3**, 243-250. [https://doi.org/10.1016/s2215-0366\(15\)00471-x](https://doi.org/10.1016/s2215-0366(15)00471-x)
- [10] Kessler, R.C., Warner, C.H., Ivany, C., Petukhova, M.V., Rose, S., Bromet, E.J., *et al.* (2015) Predicting Suicides after Psychiatric Hospitalization in US Army Soldiers: The Army Study to Assess Risk and Resilience in Servicemembers (Army STARRS). *JAMA Psychiatry*, **72**, 49-57. <https://doi.org/10.1001/jamapsychiatry.2014.1754>
- [11] Chen, J.H. and Asch, S.M. (2017) Machine Learning and Prediction in Medicine—Beyond the Peak of Inflated Expectations. *New England Journal of Medicine*, **376**, 2507-2509. <https://doi.org/10.1056/nejmp1702071>
- [12] Sendak, M., Gao, M., Nichols, C., *et al.* (2020) “Human-Compatible” Machine Learning as a Step toward Safe Clinical AI. *NPJ Digital Medicine*, **3**, Article No. 141.
- [13] Sutton, R.S. and Barto, A.G. (2018) Reinforcement Learning: An Introduction. 2nd Edition, MIT Press.
- [14] Gottesman, O., Johansson, F., Komorowski, M., Faisal, A., Sontag, D., Doshi-Velez, F., *et al.* (2019) Guidelines for Reinforcement Learning in Healthcare. *Nature Medicine*, **25**, 16-18. <https://doi.org/10.1038/s41591-018-0310-5>
- [15] Yu, C., Liu, J., Nemati, S. and Yin, G. (2021) Reinforcement Learning in Healthcare: A Survey. *ACM Computing Surveys*, **55**, 1-36.
- [16] Silver, D., Huang, A., Maddison, C.J., Guez, A., Sifre, L., van den Driessche, G., *et al.* (2016) Mastering the Game of Go with Deep Neural Networks and Tree Search. *Nature*, **529**, 484-489. <https://doi.org/10.1038/nature16961>
- [17] Vinyals, O., Babuschkin, I., Czarneki, W.M., Mathieu, M., Dudzik, A., Chung, J., *et al.* (2019) Grandmaster Level in StarCraft II Using Multi-Agent Reinforcement Learning. *Nature*, **575**, 350-354. <https://doi.org/10.1038/s41586-019-1724-z>
- [18] Komorowski, M., Celi, L.A., Badawi, O., Gordon, A.C. and Faisal, A.A. (2018) The Artificial Intelligence Clinician Learns Optimal Treatment Strategies for Sepsis in Intensive Care. *Nature Medicine*, **24**, 1716-1720. <https://doi.org/10.1038/s41591-018-0213-5>
- [19] Peng, X., Ding, Y., Wirsching, W., *et al.* (2018) Improving Sepsis Treatment Strategies

- by Combining Deep and Kernel-Based Reinforcement Learning. *AMIA Annual Symposium Proceedings*, San Francisco, 3-7 November 2018, 887-896.
- [20] Zhao, R., Pacella, M., Sanmugarajah, J., *et al.* (2022) Deep Reinforcement Learning for Treatment Duration Decision Making in Acute Lymphoblastic Leukemia. *IEEE Journal of Biomedical and Health Informatics*, **26**, 4623-4634.
- [21] Colombo, F., Calesella, F., Mazza, M.G., Melloni, E.M.T., Morelli, M.J., Scotti, G.M., *et al.* (2022) Machine Learning Approaches for Prediction of Bipolar Disorder Based on Biological, Clinical and Neuropsychological Markers: A Systematic Review and Meta-Analysis. *Neuroscience & Biobehavioral Reviews*, **135**, Article 104552. <https://doi.org/10.1016/j.neubiorev.2022.104552>
- [22] He, M., Bakker, E.M. and Lew, M.S. (2024) DPD (Depression Detection) Net: A Deep Neural Network for Multimodal Depression Detection. *Health Information Science and Systems*, **12**, Article No. 53. <https://doi.org/10.1007/s13755-024-00311-9>
- [23] First, M.B., Williams, J.B.W., Karg, R.S. and Spitzer, R.L. (2015) Structured Clinical Interview for DSM-5 Research Version (SCID-5 for DSM-5, Research Version; SCID-5-RV). American Psychiatric Association.
- [24] Montgomery, S.A. and Åsberg, M. (1979) A New Depression Scale Designed to Be Sensitive to Change. *British Journal of Psychiatry*, **134**, 382-389. <https://doi.org/10.1192/bjp.134.4.382>
- [25] Schulman, J., Wolski, F., Dhariwal, P., Radford, A. and Klimov, O. (2017) Proximal Policy Optimization Algorithms.
- [26] Thomas, P., Theodorou, G. and Ghavamzadeh, M. (2015) High-Confidence Off-Policy Evaluation. *Proceedings of the AAAI Conference on Artificial Intelligence*, **29**, 3000-3006. <https://doi.org/10.1609/aaai.v29i1.9541>
- [27] Schulman, J., Levine, S., Abbeel, P., Jordan, M. and Moritz, P. (2015) Trust Region Policy Optimization. *International Conference on Machine Learning*, Lille, 7-9 July 2015, 1889-1897.
- [28] Chernozhukov, V., Chetverikov, D., Demirer, M., Duflo, E., Hansen, C., Newey, W., *et al.* (2018) Double/Debiased Machine Learning for Treatment and Structural Parameters. *The Econometrics Journal*, **21**, C1-C68. <https://doi.org/10.1111/ectj.12097>
- [29] Lundberg, S.M. and Lee, S.I. (2017) A Unified Approach to Interpreting Model Predictions. *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017*, Long Beach, 4-9 December 2017, 4765-4774.
- [30] Vickers, A.J. and Elkin, E.B. (2006) Decision Curve Analysis: A Novel Method for Evaluating Prediction Models. *Medical Decision Making*, **26**, 565-574. <https://doi.org/10.1177/0272989x06295361>
- [31] Cuijpers, P., Turner, E.H., Mohr, D.C., Hofmann, S.G., Andersson, G., Berking, M., *et al.* (2014) Comparison of Psychotherapies for Adult Depression to Pill Placebo Control Groups: A Meta-Analysis. *Psychological Medicine*, **44**, 685-695. <https://doi.org/10.1017/s0033291713000457>
- [32] Geddes, J.R., Gardiner, A., Rendell, J., Voysey, M., Tunbridge, E., Hinds, A., *et al.* (2022) Comparative Evaluation of Quetiapine plus Lamotrigine Combination versus Quetiapine Monotherapy in Bipolar Depression: A Randomized, Double-Blind, Placebo-Controlled Trial. *The Lancet Psychiatry*, **9**, 883-894.
- [33] McIntyre, R.S., Berk, M., Brietzke, E., Goldstein, B.I., López-Jaramillo, C., Kessing, L.V., *et al.* (2020) Bipolar Disorders. *The Lancet*, **396**, 1841-1856.

- [https://doi.org/10.1016/s0140-6736\(20\)31544-0](https://doi.org/10.1016/s0140-6736(20)31544-0)
- [34] Wray, N.R., Ripke, S., Mattheisen, M., Trzaskowski, M., Byrne, E.M., Abdellaoui, A., *et al.* (2018) Genome-Wide Association Analyses Identify 44 Risk Variants and Refine the Genetic Architecture of Major Depression. *Nature Genetics*, **50**, 668-681. <https://doi.org/10.1038/s41588-018-0090-3>
- [35] Zeng, J., Zhang, Y., Xiang, Y., Liang, S., Xue, C., Zhang, J., *et al.* (2023) Optimizing Multi-Domain Hematologic Biomarkers and Clinical Features for the Differential Diagnosis of Unipolar Depression and Bipolar Depression. *NPJ Mental Health Research*, **2**, Article No. 4. <https://doi.org/10.1038/s44184-023-00024-z>
- [36] Kanchapogu, N.R. and Mohanty, S.N. (2025) Deep Learning with Ensemble-Based Hybrid AI Model for Bipolar and Unipolar Depression Detection Using Demographic and Behavioral Based on Time-Series Data. *Dialogues in Clinical Neuroscience*, **27**, 16-35. <https://doi.org/10.1080/19585969.2025.2524337>
- [37] Kessler, R.C., Bossarte, R.M., Luedtke, A., *et al.* (2023) Evaluation of a Machine Learning-Based Prediction Model for Benefit and Harm from Antidepressant Treatment in the EM-BARC Randomized Clinical Trial. *JAMA Network Open*, **6**, e2327755.
- [38] Iniesta, R., Hodgson, K., Stahl, D., Malki, K., Maier, W., Rietschel, M., *et al.* (2018) Antidepressant Drug-Specific Prediction of Depression Treatment Outcomes from Genetic and Clinical Variables. *Scientific Reports*, **8**, Article No. 5380. <https://doi.org/10.1038/s41598-018-23584-z>
- [39] Henry, K.E., Hager, D.N., Pronovost, P.J. and Saria, S. (2015) A Targeted Real-Time Early Warning Score (TREWScore) for Septic Shock. *Science Translational Medicine*, **7**, 299ra122. <https://doi.org/10.1126/scitranslmed.aab3719>
- [40] Tomašev, N., Glorot, X., Rae, J.W., Zielinski, M., Askham, H., Saraiva, A., *et al.* (2019) A Clinically Applicable Approach to Continuous Prediction of Future Acute Kidney Injury. *Nature*, **572**, 116-119. <https://doi.org/10.1038/s41586-019-1390-1>
- [41] Amodei, D., Olah, C., Steinhardt, J., Christiano, P., Schulman, J. and Mané, D. (2016) Concrete Problems in AI Safety. <https://doi.org/10.48550/arXiv.1606.06565>
- [42] Knevel, R. and Liao, K.P. (2023) From Real-World Electronic Health Record Data to Real-World Results Using Artificial Intelligence. *Annals of the Rheumatic Diseases*, **82**, 306-311. <https://doi.org/10.1136/ard-2022-222626>
- [43] Kosorok, M.R. and Laber, E.B. (2019) Precision Medicine. *Annual Review of Statistics and Its Application*, **6**, 263-286. <https://doi.org/10.1146/annurev-statistics-030718-105251>
- [44] Rieke, N., Hancox, J., Li, W., Milletari, F., Roth, H.R., Albarqouni, S., *et al.* (2020) The Future of Digital Health with Federated Learning. *NPJ Digital Medicine*, **3**, Article No. 119. <https://doi.org/10.1038/s41746-020-00323-1>
- [45] Ghassemi, M., Oakden-Rayner, L. and Beam, A.L. (2021) The False Hope of Current Approaches to Explainable Artificial Intelligence in Health Care. *The Lancet Digital Health*, **3**, e745-e750. [https://doi.org/10.1016/s2589-7500\(21\)00208-9](https://doi.org/10.1016/s2589-7500(21)00208-9)