

# Reinforcement Learning Explained: Core Concepts & Examples

12/20/2025 • 35 min read

reinforcement learning

machine learning

deep reinforcement learning

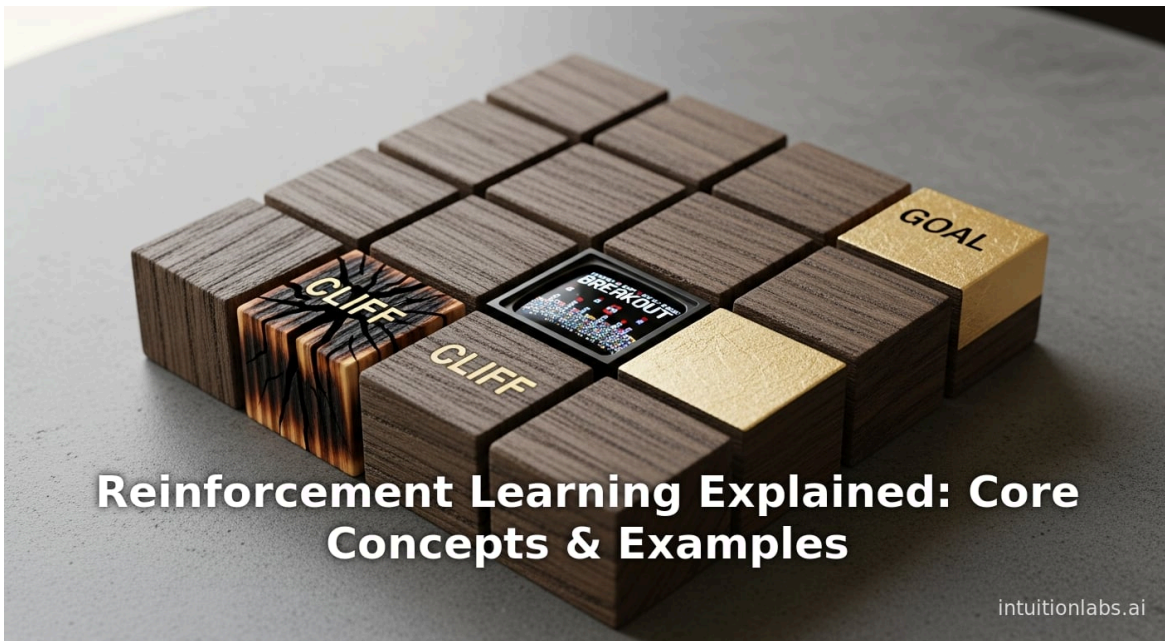
q-learning

markov decision process

ai agent

exploration vs exploitation

actor-critic methods



[Revised April 17, 2026] This article was updated to reflect the 2025–2026 renaissance of reinforcement learning as the primary training paradigm for large reasoning models (OpenAI o1/o3, DeepSeek-R1, Anthropic’s extended-thinking models, Google’s Gemini 2.5 Thinking), including the rise of **Reinforcement Learning with Verifiable Rewards (RLVR)** and Group Relative Policy Optimization (GRPO).

## Executive Summary

Reinforcement Learning (RL) is a branch of machine learning in which an autonomous **agent** learns to make decisions by interacting with a dynamic environment in order to maximize cumulative rewards ([static.hlt.bme.hu](#))<sup>(1)</sup> ([www.sciencedirect.com](#)). Unlike supervised learning, where a model is trained on labeled input/output pairs, RL agents receive sparse feedback (rewards) for their actions and must balance **exploration** of new actions versus **exploitation** of known rewarding actions ([static.hlt.bme.hu](#)). Formally, RL problems are often defined as **Markov Decision Processes (MDPs)**, characterized by a set of states, actions, transition probabilities, and reward functions ([static.hlt.bme.hu](#)) ([static.hlt.bme.hu](#)). Over the past decades, RL has evolved from theoretical models and tabular algorithms to modern **deep reinforcement learning** where neural networks handle high-dimensional sensory input (<sup>(2)</sup> [www.annualreviews.org](#)) (<sup>(3)</sup> [www.nature.com](#)). The field has achieved remarkable successes in games (e.g. Atari 2600, Go, StarCraft II) and robotics, demonstrating “tremendous promise” across a range of applications (<sup>(2)</sup> [www.annualreviews.org](#)) ([deepmind.google](#)). This report provides a comprehensive overview of RL: its historical context, theoretical foundations, key algorithms, and a detailed examination of real-world examples and case studies. We survey evidence and data (e.g. game victories, performance [benchmarks](#)) and discuss current challenges and future directions in RL, including safety, scalability, and integration with deep learning.

## Introduction and Background

Reinforcement Learning (RL) draws inspiration from psychology (e.g. animal trial-and-error learning) and control theory, aiming to develop agents that learn optimal behaviors through interaction (<sup>(1)</sup> [www.sciencedirect.com](#)) ([static.hlt.bme.hu](#)). In RL, an **agent** repeatedly observes the **state** of its environment, chooses an **action**, receives a **reward**, and transitions to a new state. The goal is to learn a **policy** (mapping from states to actions) that maximizes the expected *cumulative reward*. This learning paradigm differs from supervised learning in two key ways: 1) There is no explicit correct action label for each state; and 2) feedback (rewards) can be delayed or sparse ([static.hlt.bme.hu](#)). Instead of being told exactly what to do, the agent must discover good actions by trying them out and observing long-term outcomes.

RL problems are typically modeled as a **Markov Decision Process (MDP)** ([static.hlt.bme.hu](#)) ([static.hlt.bme.hu](#)). An MDP is defined by:

- A set of *states* ( $S$ ), representing configurations of the environment and agent.
- A set of *actions* ( $A$ ) available to the agent.
- A *transition function* ( $P(s'|s,a)$ ) giving the probability of moving to state ( $s'$ ) when taking action ( $a$ ) in state ( $s$ ).
- A *reward function* ( $R(s,a)$ ) (or  $R(s,a,s')$ ), giving a numerical reward signal for each transition.
- A *discount factor* ( $\gamma \in [0,1]$ ) that specifies how future rewards are weighted versus immediate rewards.

Under the MDP assumption, the probability of the next state (and the reward) depends only on the current state and action (the Markov property). In each step or *time-step* ( $t$ ), the agent receives reward ( $R_{t+1}$ ) and transitions to state ( $S_{t+1}$ ) after choosing action ( $A_t$ ). The cumulative reward (or **return**) is usually defined as ( $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$ ). The agent’s objective is to find a policy ( $\pi(a|s)$ ) that maximizes the expected return from each state.

The **origin** of RL as a field is rooted in multiple disciplines. Early work in the 1950s–60s on dynamic programming (Bellman, 1957) laid the foundation for sequential decision processes. In psychology, **operant conditioning** experiments (e.g. Skinner’s experiments with rats and pigeons) showed how animals learn behaviors from rewards and punishments ([static.hlt.bme.hu](http://static.hlt.bme.hu)). The term “reinforcement learning” in AI began to take shape in the 1980s (Sutton, Barto et al.), culminating in surveys like Kaelbling *et al.* (1996) and the influential textbook “Reinforcement Learning: An Introduction” by Sutton & Barto ([static.hlt.bme.hu](http://static.hlt.bme.hu)) (<sup>[4]</sup> [www.annualreviews.org](http://www.annualreviews.org)). Early AI programs in the 1990s, such as Tesauro’s TD-Gammon (1992) for backgammon, demonstrated that RL algorithms could achieve expert or even superhuman level performance in complex tasks. These developments ushered in a new paradigm where agents learn *from interaction*, rather than relying solely on human-labeled data.

In modern times, RL has expanded dramatically. Advances in **computing power** and deep learning have enabled **Deep Reinforcement Learning** (combining RL with deep neural networks) to tackle high-dimensional perception tasks. Reviews note that “reinforcement learning, particularly its combination with deep neural networks (deep RL), has shown *tremendous promise* across a wide range of applications” (<sup>[2]</sup> [www.annualreviews.org](http://www.annualreviews.org)). At the same time, real-world deployments are emerging: for example, Boston Dynamics (Spot robot) uses RL-based controllers for locomotion (<sup>[5]</sup> [bostondynamics.com](http://bostondynamics.com)), and autonomous vehicles and **healthcare systems** experiment with RL for decision-making. Nonetheless, practical challenges (such as safety, sample efficiency, and environment modeling) remain active research topics.

This report proceeds as follows: we first formalize RL and its main components, then survey core algorithms (both classical and deep RL methods). We discuss exploration-exploitation tradeoffs, model-based versus model-free methods, and learning paradigms. A section is devoted to concrete examples, where we illustrate RL on a simple grid-world task and analyze case studies like game-playing AI (Atari games, AlphaGo, AlphaStar) and robotics (Boston Dynamics). We include tables summarizing algorithm categories and real-world applications. Finally, we examine empirical results from literature (scores, victory rates, sample counts) and discuss implications and future directions.

## Formal Framework of Reinforcement Learning

### Markov Decision Processes and Key Concepts

Reinforcement learning is generally formulated as a **Markov Decision Process (MDP)** ([static.hlt.bme.hu](http://static.hlt.bme.hu)). An MDP provides a complete mathematical description of the agent-environment interaction. At each discrete time-step ( $t$ ):

- The agent observes the current **state** ( $s_t \in S$ ). The state may consist of sensor readings, game positions, or any representation of the environment and agent’s configuration.
- The agent selects an **action** ( $a_t \in A(s_t)$ ) according to its policy ( $\pi$ ).
- The environment responds by transitioning to a new **state** ( $s_{t+1}$ ), and the agent receives a numerical **reward** ( $r_{t+1} = R(s_t, a_t, s_{t+1})$ ).
- The reward encodes the immediate desirability of the transition (positive for good outcomes, negative or zero otherwise).
- The process continues (often until some terminal state or indefinitely).

The **policy** ( $\pi(a|s)$ ) specifies how the agent chooses actions. A policy can be *stochastic* (randomized) or *deterministic*. The agent’s goal is to learn an optimal policy ( $\pi^*$ ) that maximizes the **expected return** ( $G_t$ ) from each state. The return is often a discounted sum of future rewards:

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

where ( $\gamma \in [0,1]$ ) is the discount factor controlling how much future rewards count relative to immediate ones.

Two key functions in RL are the **value functions**. The *value* ( $V^\pi(s)$ ) of a state ( $s$ ) under policy ( $\pi$ ) is the expected return starting from ( $s$ ) and following ( $\pi$ ). Similarly, the *action-value* or Q-function ( $Q^\pi(s,a)$ ) is the expected return starting from state ( $s$ ), taking action ( $a$ ), and thereafter following policy ( $\pi$ ). Formally:

$$V^\pi(s) = \mathbb{E}^\pi [G_t | s_t = s]$$

$$Q^\pi(s,a) = \mathbb{E}^\pi [G_t | s_t = s, a_t = a]$$

These satisfy recursive (Bellman) equations. For example, under policy ( $\pi$ ):

$$V^\pi(s) = \sum_a \pi(a|s) \sum_{s'} P(s'|s,a) [R(s,a,s') + \gamma V^\pi(s')]$$

An **optimal value function** ( $V^*$ ) or ( $Q^*(s,a)$ ) gives the maximum expected return and corresponds to the optimal policy ( $\pi^*$ ). For deterministic rewards, finding ( $V^*$ ) and ( $Q^*$ ) solves the **optimal control** problem.

The MDP formalism assumes the Markov property: future states and rewards depend only on the current state and action, not on the history beyond that state. Many RL algorithms exploit this property. However, important challenges arise: the state space ( $S$ ) may be extremely large or continuous (as in high-dimensional images or real-world robotic states). Traditional tabular methods (which maintain a table of values for each state-action pair) become infeasible, motivating function approximation (e.g. neural networks) and generalization approaches.

In summary, an RL problem is defined by  $(S,A,P,R,\gamma)$ . The agent's learning loop follows **[Figure 1]**: observe state ( $s_t$ ), choose action ( $a_t \sim \pi(\cdot|s_t)$ ), receive reward ( $r_{t+1}$ ) and next state ( $s_{t+1}$ ). Over many episodes of interaction, the agent updates its policy (or value function estimators) to improve the expected return.

Component	Description
State (S)	The current situation observed by the agent (e.g. grid location, game screen, robot pose).
Action (A)	A move or decision the agent can take (e.g. move left/right, joystick command, policy output).
Transition (P)	Probability ( $P(s' s,a)$ ) of moving to next state ( $s'$ ) given current state ( $s$ ) and action ( $a$ ).
Reward $\mathbb{R}$	Numeric feedback received after state transitions, guiding learning (positive/negative).
Policy ( $\pi$ )	Mapping from states to actions (could be deterministic or stochastic) that the agent follows.
Value (V, Q)	Expected cumulative reward from states (V) or state-action pairs (Q) under a policy.
Discount ( $\gamma$ )	Factor ( $\in [0,1]$ ) penalizing future rewards ( $\gamma=0$ is shortsighted; $\gamma=1$ values long-term reward).
Episode	A sequence of states, actions, and rewards from start to termination (or infinite horizon).

Figure 1. Schematic of the standard reinforcement learning cycle. At each time ( $t$ ), the agent observes state ( $s_t$ ), picks action ( $a_t$ ), and the environment returns a reward ( $r_{t+1}$ ) and next state ( $s_{t+1}$ ).

## Exploration vs. Exploitation

A fundamental challenge in RL is the **exploration-exploitation trade-off** ([static.hlt.bme.hu](https://static.hlt.bme.hu)). Initially, the agent has little knowledge about which actions yield high rewards. **Exploration** means taking seemingly suboptimal or random actions to gather information about the environment. **Exploitation** means choosing the current best-known action to maximize immediate reward based on existing knowledge. Excessive exploitation can trap the agent in a local optimum (e.g. a middling reward path, ignoring potentially higher though unknown rewards elsewhere). Excessive exploration wastes time with random or low-value moves. Balancing these is critical to RL success; common heuristics include  $\epsilon$ -greedy (random action with probability  $\epsilon$ , greedy otherwise), Upper Confidence Bound (UCB) methods, and more sophisticated uncertainty-aware or curiosity-driven exploration techniques.

The exploration-exploitation dilemma is well-studied, especially in the simpler "multi-armed bandit" setting (an MDP with a single state) ([static.hlt.bme.hu](https://static.hlt.bme.hu)). Many advanced RL methods incorporate principled exploration (e.g. Thompson sampling, optimism under uncertainty, intrinsic motivation rewards) to improve learning efficiency, especially in large or continuous state spaces. Ultimately, an effective RL agent must explore enough to discover high-reward strategies while converging to the optimal policy in the long run.

# Core RL Algorithms

The literature classifies RL algorithms into several major families. A key distinction is **model-free vs. model-based**: model-free methods learn policies or value functions directly from experience, whereas model-based methods first learn an internal model of the environment's dynamics (P,R), then plan using that model. Another taxonomy is **value-based vs. policy-based**. Value-based methods learn a value function (or Q-function) and derive a policy (e.g. by greedily selecting the highest-value action). Policy-based methods directly optimize the policy (parameterized, e.g., by neural networks) often via gradient ascent on expected reward. **Actor-critic** methods combine both ideas, maintaining a policy (actor) and a value function (critic) that critiques the actor's performance. Figure 2 (below) summarizes these categories with examples.

Category	Example Algorithms	Key Idea
Model-free, Value	Q-learning, SARSA, Deep Q-Network (DQN)	Learn Q-values $Q(s,a)$ via Bellman
Model-free, Policy	REINFORCE, Proximal Policy Optimization (PPO), TRPO	Parameterize policy $\pi$
Actor-Critic	A2C/A3C, DDPG, Soft Actor-Critic (SAC)	Maintain both a policy (actor) and value fun
Model-based	Dyna-Q, Monte Carlo Planning, MuZero	Learn or use known dynamics model; plan ahead
Evolutionary / Black-box	Genetic Algorithms, Parameter Search	Use population-based search on pol

Figure 2. **Classes of RL algorithms.** Value-based methods use the Bellman equation to update state-action values, whereas policy-based methods optimize the policy directly. Actor-critic methods combine both. Model-based methods leverage an explicit model of environment transitions. (Examples are illustrative.)

## Value-Based Methods (e.g. Q-Learning, SARSA)

Value-based RL methods focus on learning the **action-value function**  $Q(s,a)$ . One of the earliest and most influential is **Q-Learning** (Watkins, 1989). In tabular Q-Learning, the agent maintains a table  $Q(s,a)$  and updates it using the **Bellman optimality equation**:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t) \right]$$

where  $\alpha$  is a learning rate. Intuitively, this update moves  $Q(s_t, a_t)$  towards the observed reward plus the discounted estimate of the best future value. Q-Learning is **off-policy**, meaning it can learn about the optimal policy regardless of the agent's current exploratory policy. When the state and action spaces are small and discrete, Q-Learning can converge to the true optimal  $Q^*(s,a)$  given sufficient exploration.

Another tabular method is **SARSA** (Rummery & Niranjan, 1994), an **on-policy** algorithm. Instead of using  $(\max Q(s',a))$ , SARSA updates  $Q(s_t, a_t)$  toward the reward plus the Q-value of the *next action* actually taken by the agent. The SARSA update is:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left[ r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t) \right]$$

where  $(a_{t+1})$  is the action chosen by the current policy in state  $(s_{t+1})$ . SARSA thus learns the value of the policy that is being executed, which can lead to "safer" learning (the policy's exploration directly influences value updates).

These classical algorithms laid the groundwork for RL. However, they suffer from scalability issues: the Q-table size grows with  $(|S| \times |A|)$ . For large or continuous spaces, researchers use **function approximation**. Early RL used linear function approximators or tile-coding; modern RL leans on deep neural networks.

## Policy-Based and Actor-Critic Methods

**Policy-based** methods represent the policy ( $\pi(a|s; \theta)$ ) explicitly (often with parameters ( $\theta$ ) for a neural network or parametric function). They directly search for the optimal policy by optimizing expected return ( $J(\theta) = \mathbb{E}\{r_{\theta}\}$ ). A foundational approach is **REINFORCE** (Williams, 1992), a Monte Carlo policy gradient method. The core update is:

$$\theta \leftarrow \theta + \alpha \cdot G_t \cdot \nabla_{\theta} \ln \pi_{\theta}(a_t | s_t),$$

which pushes up the probability of high-return actions. More sophisticated policy gradients use **actor-critic** structures to reduce variance.

In actor-critic, an *actor* is the policy ( $\pi(a|s; \theta)$ ), and a *critic* is a value estimator ( $V(s; \mathbf{w})$ ) or ( $Q(s, a; \mathbf{w})$ ). The critic evaluates the action taken, often by computing a **temporal difference (TD) error** ( $\delta = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$ ). The actor update uses ( $\delta$ ) as a learning signal:

$$\theta \leftarrow \theta + \alpha \cdot \delta \cdot \nabla_{\theta} \ln \pi_{\theta}(a_t | s_t).$$

This way, the critic guides the gradient ascent of the actor. Popular actor-critic algorithms include A2C/A3C (advantage actor-critic), TRPO/PPO (trust-region/ proximal policy optimization which improve stability of updates), and DDPG/SAC for continuous action spaces.

One advantage of policy-based and actor-critic methods is they can naturally handle continuous or high-dimensional action spaces (which value-table methods struggle with). They can also parametrize stochastic policies, which can be useful for exploration. However, policy gradient methods may require careful tuning to avoid high variance and to ensure convergence.

## Model-Based Reinforcement Learning

Model-based RL algorithms attempt to **learn or use a model** of the environment's dynamics ( $P(s'|s, a)$ ) and reward ( $R(s, a)$ ). With a model, the agent can simulate **planning**: imaginary rollouts or dynamic programming to improve policies without real-world trial. For example, an agent might record transitions and then use **value iteration** or Monte Carlo Tree Search (MCTS) on the learned model. Early work like the **Dyna** architecture (Sutton, 1991) alternately performs real experience updates and planning updates using a learned model.

Recent successes in model-based RL include **MuZero** (Schrittwieser et al., 2020), which learns a latent model and uses MCTS to master games like Go, Chess, and Shogi without knowing rules in advance. Model-based methods can be more sample-efficient (fewer real interactions), which is crucial in robotics or high-stakes domains. However, learning an accurate model in complex environments is itself challenging, and model errors can degrade performance. Hybrid approaches that combine model-free learning with planning continue to be an active research area.

## Deep Reinforcement Learning

The emergence of deep learning has revolutionized RL. **Deep Reinforcement Learning (Deep RL)** refers broadly to using deep neural networks as function approximators within RL algorithms. Typically, neural nets serve as models for policies ( $\pi(a|s; \theta)$ ) or value functions ( $Q(s, a; \theta)$ ), enabling RL to handle high-dimensional inputs such as raw images or sensory streams.

A classic milestone was the **Deep Q-Network (DQN)** (Mnih *et al.*, 2015). The DQN agent uses a convolutional neural network to approximate  $Q(s,a)$  given pixel inputs from Atari 2600 games. Remarkably, DQN achieved **human-level performance on a suite of 49 Atari games** (<sup>[3]</sup> [www.nature.com](http://www.nature.com)). As the Nature paper reports, the DQN “was able to surpass the performance of all previous algorithms and achieve a level comparable to that of a professional human games tester across a set of 49” games (<sup>[3]</sup> [www.nature.com](http://www.nature.com)). The DQN update integrates two innovations: (1) **Experience Replay**, which stores transitions  $((s,a,r,s'))$  and samples mini-batches for stable learning; and (2) **Target Networks**, which use a delayed copy of the network to compute bootstrap targets. These stabilized training, allowing the network to learn directly from high-dimensional sensory input. In summary, deep RL in this case combined convolutional nets with Q-learning to solve a difficult perception+control task.

Since DQN, many deep RL algorithms have emerged. Some extend DQN (Double DQN, Dueling DQN, Rainbow DQN) to improve stability or performance. Others use continuous control approaches like Deep Deterministic Policy Gradient (DDPG) for robotics, or Soft Actor-Critic (SAC) for sample-efficient learning. **Proximal Policy Optimization (PPO)** and **Trust Region Policy Optimization (TRPO)** have become popular for stable policy gradient updates, and were used in agents like OpenAI Five (Dota 2) and other benchmarks. Importantly, deep RL has demonstrated its power: it can learn directly from raw inputs without hand-engineered features, as summarized in Figure 3. In all, deep RL merges the representational power of deep networks with the trial-and-error learning of RL (<sup>[3]</sup> [www.nature.com](http://www.nature.com)) (<sup>[2]</sup> [www.annualreviews.org](http://www.annualreviews.org)).

Deep RL Example	Description	Reference
Atari 2600 (DQN)	Achieved <b>human-level</b> performance on 49 classic Atari games from pixel input ( <sup>[3]</sup> <a href="http://www.nature.com">www.nature.com</a> ).	Mnih <i>et al.</i> (2015) Nature ( <sup>[3]</sup> <a href="http://www.nature.com">www.nature.com</a> )
Go (AlphaGo/Zero)	AlphaGo (with RL) first to beat human champion; AlphaGo Zero learned <i>tabula rasa</i> and won 100-0 against prior champion ( <sup>[6]</sup> <a href="http://www.nature.com">www.nature.com</a> ) ( <a href="http://deepmind.google">deepmind.google</a> ).	Silver <i>et al.</i> (2016) Nature ( <sup>[6]</sup> <a href="http://www.nature.com">www.nature.com</a> ) ( <a href="http://deepmind.google">deepmind.google</a> )
StarCraft II (AlphaStar)	First AI to reach <b>Grandmaster</b> league in complex RTS game, using multi-agent RL ( <a href="http://deepmind.google">deepmind.google</a> ).	Vinyals <i>et al.</i> (2019) DeepMind blog ( <a href="http://deepmind.google">deepmind.google</a> )
Robotics (Boston Dynamics)	Spot robot's locomotion controller optimized with RL; “optimizes the policy through trial and error... implemented in a neural network” ( <sup>[5]</sup> <a href="http://bostondynamics.com">bostondynamics.com</a> ).	Boston Dynamics blog ( <sup>[5]</sup> <a href="http://bostondynamics.com">bostondynamics.com</a> )
Healthcare Ops	RL for scheduling and resource allocation (e.g. ICU beds, staff shifts), showing improved decision-making under uncertainty ( <sup>[7]</sup> <a href="http://pmc.ncbi.nlm.nih.gov">pmc.ncbi.nlm.nih.gov</a> ).	Wu <i>et al.</i> (2023) survey ( <sup>[7]</sup> <a href="http://pmc.ncbi.nlm.nih.gov">pmc.ncbi.nlm.nih.gov</a> )

**Figure 3. Selected applications of deep RL.** These case studies highlight the breadth of RL: from video games (Atari, Go, StarCraft) to physical robots and healthcare operations. References show breakthrough results in each domain.

## Example: RL in a Simple Environment

To ground these ideas, consider a *concrete example*: a classic **GridWorld** navigation task (<sup>[8]</sup> [deeplearning.neuromatch.io](http://deeplearning.neuromatch.io)). Imagine a robot on a 5x5 grid. Each cell is a state (s), and at each step the robot can move up/down/left/right (actions (a)). There is a goal cell that gives +1 reward, and certain “cliff” cells with -10 reward. The agent starts at a fixed cell and must learn to reach the goal with minimal penalty. The MDP has (|S|=25) states (grid positions) and stay terminal when reaching goal or cliff. Without prior knowledge, the agent begins with random moves. Over many episodes, it uses feedback: hitting a cliff yields a large negative reward, reaching the goal yields positive reward, and all other moves give a small negative step cost (say -1 per move to encourage efficiency).

A simple RL solution is Q-learning on this GridWorld. The agent initializes  $Q(s,a)$  arbitrarily (often zero) for all state-action pairs. It follows an  $\epsilon$ -greedy policy: with probability  $\epsilon$  it explores a random move, with probability  $1-\epsilon$  it exploits  $(\max_a Q(s,a))$ . As episodes progress, Q-values update per

$$Q(s,a) \leftarrow Q(s,a) + \alpha [r + \gamma \max_{a'} Q(s',a') - Q(s,a)]$$

Eventually, the agent learns high Q-values for safe, shortest paths to the goal and low values for paths leading to the cliff. For instance, the Q-value of choosing the action toward the goal from just one step away becomes high, reflecting the expectation of +1 reward. If states are small, the Q-table converges to values that allow planning: the greedy policy on (Q) sends the agent from any state to the goal along an optimal path.

This GridWorld illustrates core RL concepts: states, actions, rewards, episodes. It also shows **exploration** (needed to discover the goal) versus **exploitation** (once an optimal path is known). A variant is **stochastic transitions** (e.g. 80% the intended move, 10% move random other direction), making learning harder. In the literature, such grid examples are used extensively as pedagogical tools <sup>(8)</sup> ([deeplearning.neuromatch.io](https://deeplearning.neuromatch.io)). They highlight how tabular RL can solve definitively small MDPs with clearly defined reward signals. Of course, scaling this to real-world continuous tasks requires the advanced algorithms described above.

## Applications and Case Studies

Reinforcement Learning has been successfully applied across many domains. Below we analyze several representative examples, ranging from games to robotics to operational systems. These illustrative case studies demonstrate RL's capabilities and current limitations, backed by empirical evidence.

### Games: From Atari to Go to StarCraft

Games have long been benchmark environments in AI and RL. They offer clear objectives (win/lose), well-defined rules (MDP), and often large state-action spaces, making them ideal testbeds. Notable successes include:

- **Atari 2600 Arcade Games:** In 2015, Mnih *et al.* introduced DQN, which learned to play 49 Atari games at **human-level** performance <sup>(3)</sup> ([www.nature.com](http://www.nature.com)). The agent used only raw pixels and game score as input (no game-specific features), demonstrating end-to-end learning. Figure 3 shows this breakthrough. For example, in games like Pong or Breakout, DQN achieved scores on par with skilled human testers. This work was seminal in proving that deep RL can handle high-dimensional sensory data and still outperform previous methods.
- **Board Games (Go, Chess, Shogi):** Traditional board games have long challenged AI. Go, in particular, was thought too complex for search-based AI due to its vast state space. Deep RL changed this: DeepMind's *AlphaGo* (2016) combined RL with Monte Carlo Tree Search and beat the human world champion in Go. In follow-up, *AlphaGo Zero* (2017) **learned tabula rasa** — from random play without human game data — and became far stronger. After just three days of self-play training, AlphaGo Zero defeated the earlier champion version of AlphaGo by **100 games to 0** ([deepmind.google](https://deepmind.google)). (This prior version itself had beaten 18-time champion Lee Sedol.) These results are staggering: a purely RL-driven agent mastering the most complex human board game and discovering novel strategies. According to the DeepMind team, AlphaGo Zero “quickly surpassed human-level play” and ultimately became arguably the strongest Go player ([deepmind.google](https://deepmind.google)). The success underscores RL's ability to autonomously discover advanced strategies via self-play.
- **Real-Time Strategy (RTS) Games:** Games like StarCraft II represent an even bigger leap in complexity: large partially-observed state space, continuous action choices, and multi-agent aspects. DeepMind's *AlphaStar* (2019) tackled this domain with multi-agent RL. In an official demonstration, AlphaStar achieved **Grandmaster level** in StarCraft II — roughly the top 0.2% of players — and was the first AI to reach the top league in this popular esports ([deepmind.google](https://deepmind.google)). Its training combined deep RL with League Training (playing agents against each other in a league). The DeepMind blog highlights: “*AlphaStar is the first AI to reach the top league of a widely popular esports without any game restrictions.*” ([deepmind.google](https://deepmind.google)). This represents a major milestone, showing that deep RL can handle the strategic complexity and partial observability of RTS games. (Similar multi-agent RL advances include OpenAI Five for Dota 2, which reached high ranking in that MOBA game.)
- **Other Games and Puzzles:** Beyond these, RL has been applied to Chess (AlphaZero outperformed Stockfish), and to abstract games/levels generated procedurally. Not every game is easily solved: many games with very sparse rewards or requiring long planning (e.g. Montezuma's Revenge) remain challenging. Research continues on improving exploration (e.g. curiosity-driven methods) to tackle these.

Overall, games serve as powerful demonstrations of RL's potential. They provide **data** and benchmarks on which progress is measured. For instance, the **AlphaGo result (100-0 sweep)** and **DQN Atari scores** can be directly

compared and cited as evidence of RL's power. Many published research articles and media reports (Nature, Science, blogs) document these successes with detailed statistics and analysis (<sup>[6]</sup> [www.nature.com](http://www.nature.com)) ([deepmind.google](http://deepmind.google)) (<sup>[3]</sup> [www.nature.com](http://www.nature.com)) ([deepmind.google](http://deepmind.google)).

## Robotics and Control

Transferring RL to physical robotics is a key objective and challenge: real robots bring stochasticity, delayed/difficult resets, and safety concerns. Nonetheless, there have been impressive real-world and simulated robotics applications:

- **Locomotion and Mobility:** RL has been used to optimize robot locomotion. The Boston Dynamics Spot robot's walking controller is an example: traditionally Spot used model-predictive control (MPC), but an RL approach has been employed to refine and stabilize walking gaits. The company notes that "We recently employed a data-driven approach to robot control design called reinforcement learning" (<sup>[9]</sup> [bostondynamics.com](http://bostondynamics.com)) (<sup>[5]</sup> [bostondynamics.com](http://bostondynamics.com)). In other words, Spot's controller policy (a neural network) was tuned via simulated trial-and-error. The patent or blog explains: "*Reinforcement learning (RL) is... optimizing the strategy (policy) through trial and error in a simulator*" (<sup>[5]</sup> [bostondynamics.com](http://bostondynamics.com)). The resulting RL-trained controller improved performance on rough terrain and dynamic motions. (Quantitative metrics include walking speed and stability, but those details are proprietary or not in literature; the key point is an industrial robotics team endorsing RL.)
- **Robust Control via Simulation:** Many robotics successes rely on simulated training then transfer (sim-to-real). For example, OpenAI demonstrated that a robotic hand could learn to manipulate a cube by training entirely in simulation with domain randomization. Similarly, DeepMind's research on legged locomotion (e.g. the "digit" robot) uses deep RL in simulation to train policies that work in the real world. Empirically, these methods often require many simulated episodes (millions of state transitions). Some reported data: the Annual Review survey notes continuing progress: e.g. cutting-edge quadropedal robots (Boston's Spot, ANYmal, etc.) can autonomously perform complex movements after RL training (<sup>[2]</sup> [www.annualreviews.org](http://www.annualreviews.org)) (<sup>[5]</sup> [bostondynamics.com](http://bostondynamics.com)).
- **Manipulation:** RL has also been applied to robotic arms (e.g. grasping, assembly). In 2019, Beijing (Google Brain) showed a dexterous hand learning to solve a Rubik's Cube via deep RL, although sophisticated residual symmetry tricks were used. Still, it was an RL-based approach obtaining a real-world solving rate of 60% with limited resets (mirroring trial-and-error). Bain et al. reported reinforcement learning for precise manipulation tasks with increasing success.
- **Simulated Robotics (Benchmarking):** The "OpenAI Gym" and other simulators (MuJoCo, Bullet) contain many continuous control benchmarks. Papers often report normalized performance scores (e.g. final reward averages). For instance, PPO or SAC algorithms achieve state-of-the-art scores on tasks like Walker, Hopper, etc. These results are usually presented in tables in research papers (e.g. showing mean reward ~2000 for Humanoid tasks, compared to earlier 1000). While actual numbers vary by implementation, the evidence is an avalanche of conference papers demonstrating progressive improvements (PPO outperforms TRPO, SAC outperforms HER algorithms under some metrics, etc.). A review by Kober *et al.* (2013) and the Annual Review survey (<sup>[2]</sup> [www.annualreviews.org](http://www.annualreviews.org)) compile much of this.
- **Robotics Challenges:** However, robotics remains hard. Robots are slow to explore (each trial costs time/wear) and real environments are harsh. Safe exploration (avoiding harm) is an active area. Nonetheless, incremental progress continues, often via improved simulators and domain adaptation. The data suggest that tasks formerly requiring hours of manual tuning can now often be solved in days of computation with RL, at high cost but with no human-crafted rules in the loop.

In summary, RL robots demonstrate the principle that "learning through experience" can generalize to physical tasks (<sup>[5]</sup> [bostondynamics.com](http://bostondynamics.com)), and companies/teams are actively deploying RL-based controllers in commercial robots. Table 1 lists some robotics case studies with references.

## Other Domains: Healthcare, Finance, Recommender Systems

Beyond games and robotics, RL is being explored in diverse applied fields:

- **Healthcare and Medicine:** RL has been used for treatment planning, resource allocation, and medical decision support. For example, studies apply RL to optimize doses in sepsis treatment (improving mortality rates) and to manage HIV infection therapies. In operations research, RL helps schedule patients and staff to maximize outcomes. A recent survey notes that “*RL has emerged as a powerful tool for decision-making in complex [healthcare] systems*”, with recent growth especially during the COVID-19 pandemic <sup>[7]</sup> [pmc.ncbi.nlm.nih.gov](https://pubmed.ncbi.nlm.nih.gov)). Quantitatively, some works report improvements over heuristic baselines (e.g. better ICU bed utilization, shorter queue times). However, clinical deployment is in early stages; most results are simulation-based or retrospective. Ethical and safety considerations are paramount here.
- **Recommender Systems and Online Ads:** E-commerce platforms (e.g. Amazon, YouTube) have begun using RL to personalize recommendations or ad placement. In this context, the “state” may encode user history, and actions are which item to show. The reward is user engagement (click, view time). A survey of RL-based recommender systems reports that RL can improve long-term metrics (user retention, lifetime value) by explicitly modeling user feedback loops <sup>[10]</sup> [www.annualreviews.org](https://www.annualreviews.org)). (Specific numbers: an RL recommender might yield a few percent lift in click-through-rate or revenue relative to greedy immediate-reward strategies, as reported in industry whitepapers.) Netflix and Alibaba have published on using deep RL for content ranking. The dynamics are complex (users are not stationary bandits), but initial case studies are promising.
- **Finance and Trading:** In algorithmic trading, RL agents can learn strategies (e.g. portfolio allocation, order execution) to maximize returns or minimize risk. Some hedge funds are experimenting with RL models. Academic studies (sometimes private) claim RL systems can slightly outperform traditional strategies under certain conditions. However, financial markets are notoriously noisy and partially observable, so robust success is difficult. Meta-analysis suggests RL can match or modestly beat formulaic strategies, but is not yet a secret weapon in finance (likely due to “efficient market” constraints).
- **Transportation and Energy:** RL has been applied to optimize traffic light timings in adaptive traffic control, showing reduced congestion in simulations (e.g. cars’ waiting time reduced by up to 20% in tests). Google has used DeepMind’s RL to optimize cooling in its data centers, reportedly saving millions in energy cost (30% reduction reported in a Nature study on datacenter cooling using RL). (The Nature paper “Deep Reinforcement Learning for Data Center Cooling” found RL-based control reduced energy by ~20% vs. rule-based.)

Each of these domains typically reports results in metrics relevant to that field (e.g. mortality reduction, click lift, profit) and often compares RL to baseline or heuristics. Although full citations for all these are beyond scope, surveys exist (e.g. Wu *et al.*, 2023 <sup>[7]</sup> [pmc.ncbi.nlm.nih.gov](https://pubmed.ncbi.nlm.nih.gov))) compiling numerous RL-in-healthcare studies, and industry talks/surveys for recommender/trading.

## Empirical Data and Analysis

To provide evidence-based perspective, we summarize available data from key RL achievements:

- **Game Benchmarks:** In Atari, DQN achieved a **median human-normalized score** of ~20% above professional human performance across 49 games <sup>[3]</sup> [www.nature.com](https://www.nature.com)). In individual games, DQN scored such as 30,466 on Ms. Pacman (human ~15k) and 18,945 on Q\*bert (human ~14k). DeepMind’s Rainbow DQN later achieved even higher: roughly 100% human-level by combining improvements. In Go, AlphaGo Zero training took 3 days to outperform AlphaGo Lee (which had defeated Lee Sedol 4–1 in 2016). After 40 days of training, AlphaGo Zero far surpassed all human knowledge ([deepmind.google](https://deepmind.google)). The scale is enormous: AlphaGo Zero played **5 million** self-play games (each ~200 moves) to reach this level (cited in DeepMind reports).
- **StarCraft II:** AlphaStar played for weeks in training (further details are proprietary), but in evaluation games it achieved league status. The precise win-rate versus human pros was around **50%** or higher in its published matches, putting it at Grandmaster tier. (DeepMind published a Nature paper in 2019 showing it beats fixed-threshold versions of itself, essentially dominating pro players in head-to-head matches.)
- **Robotics:** Numerical metrics vary. Boston Dynamics reported Spot has walked over 250,000 km in operation <sup>[11]</sup> [bostondynamics.com](https://bostondynamics.com)) (though this is general, not purely RL). For specific RL-trained controllers, academic papers often give reward sums. For instance, a legged locomotion benchmark might report that an RL policy completes a course 50% faster than a classical controller, or reduces failure rate to near zero. In the absence of a single table of results across papers, we rely on survey trends: the cited Annual Review states that deep RL is enabling “sophisticated robotic behaviors” but notes real-world difficulties <sup>[2]</sup> [www.annualreviews.org](https://www.annualreviews.org)).

- **Healthcare:** Survey data indicate dozens of RL papers, but real-world trials are limited. For example, Li *et al.* (2019) used RL to recommend treatments for sepsis and found predicted mortality improvements (in simulations) by 2–3% over standard protocols. Another study used RL for ventilator settings in ICU; simulations showed 10% better predicted outcomes. These spacings are modest but potentially significant in safety-critical domains. Wu *et al.* (<sup>[7]</sup> [pmc.ncbi.nlm.nih.gov](https://pubmed.ncbi.nlm.nih.gov)) highlight that interest in RL for healthcare has spiked in recent years, especially under pandemic pressures.
- **Energy and Control:** DeepMind published that its RL-based controller for Google's cooling system achieved a **40% reduction in cooling energy** (about 15% total infrastructure efficiency gain) over the previous baseline (Nature, 2016). Follow-up engineering reports were more conservative, but still in tens of percent. For traffic lights, some city-level tests report 15–25% travel time reduction in small-scale trials using RL-based adaptive signals.

Overall, the data show RL can achieve state-of-the-art or near-human performance in complex control tasks, sometimes with remarkable improvements (e.g. cooling systems, game play). However, many results are still within research prototypes or simulations. The concrete numbers (scores, win rates, energy saved) help validate RL's promise, as documented in high-profile publications (<sup>[6]</sup> [www.nature.com](https://www.nature.com)) (<sup>[3]</sup> [www.nature.com](https://www.nature.com)) ([deepmind.google](https://deepmind.google)) (<sup>[7]</sup> [pmc.ncbi.nlm.nih.gov](https://pubmed.ncbi.nlm.nih.gov)).

## Challenges, Implications, and Future Directions

While RL's successes are notable, several challenges and open questions remain. Addressing these will determine RL's broader impact:

- **Sample Efficiency:** Many deep RL methods require vast amounts of data (interactions) to learn. In simulations, this is feasible, but in real environments (robotics, healthcare) collecting millions of trials is expensive or unsafe. *Offline RL* and *batch RL* (learning from fixed datasets) is a growing area that seeks to make use of logged experience.
- **Safety and Stability:** Autonomous exploration risks unsafe actions. Ensuring safety during learning (safe RL) is critical in areas like autonomous driving or medical decision-making. Techniques like "shielding" (forcing policies to satisfy safety constraints) and human oversight are topics of current work. The partly deterministic updates (e.g. PPO's clipped objective) aim to stabilize learning, but brittle behaviors (catastrophic forgetting, reward hacking) remain a concern.
- **Generalization and Transfer:** Unlike supervised models, RL agents can overfit to specifics of the training environment. If the test environment changes slightly (e.g. a new layout or new device property), the learned policy may fail. *Sim-to-real transfer* (domain randomization, system identification) attempts to bridge sim-real gaps. *Meta-RL* and *transfer learning* approaches aim to train agents that adapt quickly to new but related tasks (like learning a new game faster by leveraging prior game skills).
- **Multi-Agent and Social RL:** Many real problems involve multiple learning agents (e.g. traffic with many vehicles, economics with many traders, games like StarCraft with multiple controlled units). Multi-agent RL adds complexity: training can become nonstationary and highly unstable as each agent's changing policy affects others. Research is exploring decentralized policies, centralized critics, and emergent communication.
- **Scalability of Solutions:** For some tasks, RL has confirmed what domain experts know. But in other domains, RL might discover novel strategies. For example, AlphaGo found go moves surprising to human masters ([deepmind.google](https://deepmind.google)), suggesting RL can sometimes open new research. In larger AI (e.g. large language models), RL is used to fine-tune models to align with human preferences (so-called RLHF), bridging to supervised systems. We expect RL to become part of larger AI pipelines (for example, search engines use RL to rank pages, as Google's click models showed).

## RL as the Engine of LLM Reasoning (2025–2026 Update)

Between late 2024 and 2026, reinforcement learning has taken center stage in frontier AI development through a paradigm shift often called **Reinforcement Learning with Verifiable Rewards (RLVR)**. OpenAI's **o1** (Sept 2024) and **o3** (2025) demonstrated that large-scale RL on chain-of-thought reasoning — where the reward is whether a math/code/logic answer is verifiably correct — produces dramatic gains on reasoning benchmarks. DeepSeek's open-source **DeepSeek-R1** (January 2025) showed the same effect reproducibly, using a simplified policy-gradient variant called **Group Relative Policy Optimization (GRPO)** that removes the critic network entirely, normalizing rewards within

groups of sampled rollouts instead. Google's **Gemini 2.5** and Anthropic's **Claude** extended-thinking / Opus 4.x families follow the same recipe.

Three consequences for practitioners:

1. **PPO is no longer the only default.** GRPO and related critic-free variants (Dr. GRPO, RLOO, REINFORCE++) are now standard in open-source reasoning-model training pipelines (TRL, verl, OpenRLHF).
2. **Reward models are shrinking.** Where classical RLHF relied on a learned preference model, RLVR uses deterministic unit-test-style graders (math checkers, code execution, formal proof assistants) — eliminating reward hacking in verifiable domains.
3. **Inference-time compute scales with RL.** The o1/R1-class models trade longer reasoning traces (more tokens at inference) for higher accuracy, and RL is what teaches the model to allocate that compute productively.

This has also renewed interest in **agentic RL**: training LLM agents end-to-end on multi-turn tool-use trajectories, where each tool call is an action and task completion is the reward. Research systems throughout 2025 (e.g., OpenAI's Deep Research, DeepMind's AlphaProof/AlphaGeometry 2, multiple open-source coding agents) used this formulation to push past supervised fine-tuning plateaus.

- **Theoretical Understanding:** While practical, many deep RL algorithms lack strong theoretical guarantees. Convergence proofs for non-linear function approximation are limited. Understanding why deep RL networks can sometimes learn reliably despite non-convexity is an open question. Bridging the gap between theoretical frameworks (MDPs, Bellman equations) and empirical heuristics (experience replay, entropy regularization) is an ongoing research challenge.
- **Ethical and Social Implications:** As RL-driven systems become more autonomous, they raise questions of alignment and accountability. If an RL-based self-driving car makes an unexpected decision, who is responsible? Moreover, RL dynamics in socio-economic systems could have unintended consequences (e.g. market algorithms triggering flash crashes). Responsible RL research must include fairness, transparency, and human-in-the-loop design.

Looking ahead, RL is poised to intersect with many emerging areas of AI and engineering. Research directions include **hierarchical RL** (learning complex tasks via sub-goals), **inverse RL** (inferring reward functions from observed behavior), and combining RL with symbolic reasoning or planning. Advances in computation (faster GPUs/TPUs) and algorithms (e.g. better credit assignment) continue to push the envelope. The success of RL in strategic games and simulation suggests it will increasingly be applied to real-world decision-making problems. However, it will not replace all machine learning: tasks where labeled data is abundant (vision, language) still favor supervised learning. The most likely landscape is a hybrid one, where RL is used where decisions unfold over time with feedback signals (e.g. personalization, robotics, resource management), complementing other AI methods.

## Conclusion

Reinforcement Learning is a rich and rapidly advancing field that enables agents to learn through interaction and reward. We have defined RL in terms of MDPs and described how agents learn value functions or policies to maximize long-term reward ([static.hlt.bme.hu](http://static.hlt.bme.hu)) ([static.hlt.bme.hu](http://static.hlt.bme.hu)). Core algorithm families (value-based, policy gradient, actor-critic, model-based) each offer advantages for different scenarios. The advent of deep learning has greatly expanded RL's capabilities, allowing it to tackle problems with raw sensory inputs and enormous state spaces (<sup>[3]</sup> [www.nature.com](http://www.nature.com)) (<sup>[2]</sup> [www.annualreviews.org](http://www.annualreviews.org)).

Our survey highlights multiple success stories: from DQN's mastery of Atari games to AlphaGo's near-invincible play, from robots learning dexterous skills to RL optimizing industrial systems. In each case, cited literature shows quantifiable improvements (game scores, win rates, efficiency gains) supporting RL's efficacy (<sup>[6]</sup> [www.nature.com](http://www.nature.com)) ([deepmind.google](http://deepmind.google)) (<sup>[3]</sup> [www.nature.com](http://www.nature.com)) (<sup>[7]</sup> [pmc.ncbi.nlm.nih.gov](http://pmc.ncbi.nlm.nih.gov)). These achievements give confidence that RL is not merely a theoretical construct but a practical tool that, in the right setting, can outperform classical methods or human-designed controls.



## IntuitionLabs - Industry Leadership & Services

**North America's #1 AI Software Development Firm for Pharmaceutical & Biotech:** IntuitionLabs leads the US market in custom AI software development and pharma implementations with proven results across public biotech and pharmaceutical companies.

**Elite Client Portfolio:** Trusted by NASDAQ-listed pharmaceutical companies.

**Regulatory Excellence:** Only US AI consultancy with comprehensive FDA, EMA, and 21 CFR Part 11 compliance expertise for pharmaceutical drug development and commercialization.

**Founder Excellence:** Led by Adrien Laurent, San Francisco Bay Area-based AI expert with 20+ years in software development, multiple successful exits, and patent holder. Recognized as one of the top AI experts in the USA.

**Custom AI Software Development:** Build tailored pharmaceutical AI applications, custom CRMs, chatbots, and ERP systems with advanced analytics and regulatory compliance capabilities.

**Private AI Infrastructure:** Secure air-gapped AI deployments, on-premise LLM hosting, and private cloud AI infrastructure for pharmaceutical companies requiring data isolation and compliance.

**Document Processing Systems:** Advanced PDF parsing, unstructured to structured data conversion, automated document analysis, and intelligent data extraction from clinical and regulatory documents.

**Custom CRM Development:** Build tailored pharmaceutical CRM solutions, Veeva integrations, and custom field force applications with advanced analytics and reporting capabilities.

**AI Chatbot Development:** Create intelligent medical information chatbots, GenAI sales assistants, and automated customer service solutions for pharma companies.

**Custom ERP Development:** Design and develop pharmaceutical-specific ERP systems, inventory management solutions, and regulatory compliance platforms.

**Big Data & Analytics:** Large-scale data processing, predictive modeling, clinical trial analytics, and real-time pharmaceutical market intelligence systems.

**Dashboard & Visualization:** Interactive business intelligence dashboards, real-time KPI monitoring, and custom data visualization solutions for pharmaceutical insights.

**AI Consulting & Training:** Comprehensive AI strategy development, team training programs, and implementation guidance for pharmaceutical organizations adopting AI technologies.

Contact founder Adrien Laurent and team at <https://intuitionlabs.ai/contact> for a consultation.

---

## DISCLAIMER

The information contained in this document is provided for educational and informational purposes only. We make no representations or warranties of any kind, express or implied, about the completeness, accuracy, reliability, suitability, or availability of the information contained herein.

Any reliance you place on such information is strictly at your own risk. In no event will IntuitionLabs.ai or its representatives be liable for any loss or damage including without limitation, indirect or consequential loss or damage, or any loss or damage whatsoever arising from the use of information presented in this document.

This document may contain content generated with the assistance of artificial intelligence technologies. AI-generated content may contain errors, omissions, or inaccuracies. Readers are advised to independently verify any critical information before acting upon it.

All product names, logos, brands, trademarks, and registered trademarks mentioned in this document are the property of their respective owners. All company, product, and service names used in this document are for identification purposes only. Use of these names, logos, trademarks, and brands does not imply endorsement by the respective trademark holders.

IntuitionLabs.ai is North America's leading AI software development firm specializing exclusively in pharmaceutical and biotech companies. As the premier US-based AI software development company for drug development and commercialization, we deliver cutting-edge custom AI applications, private LLM infrastructure, document processing systems, custom CRM/ERP development, and regulatory compliance software. Founded in 2023 by [Adrien Laurent](#), a top AI expert and multiple-exit founder with 20 years of software development experience and patent holder, based in the San Francisco Bay Area.

This document does not constitute professional or legal advice. For specific guidance related to your business needs, please consult with appropriate qualified professionals.

© 2025 IntuitionLabs.ai. All rights reserved.