

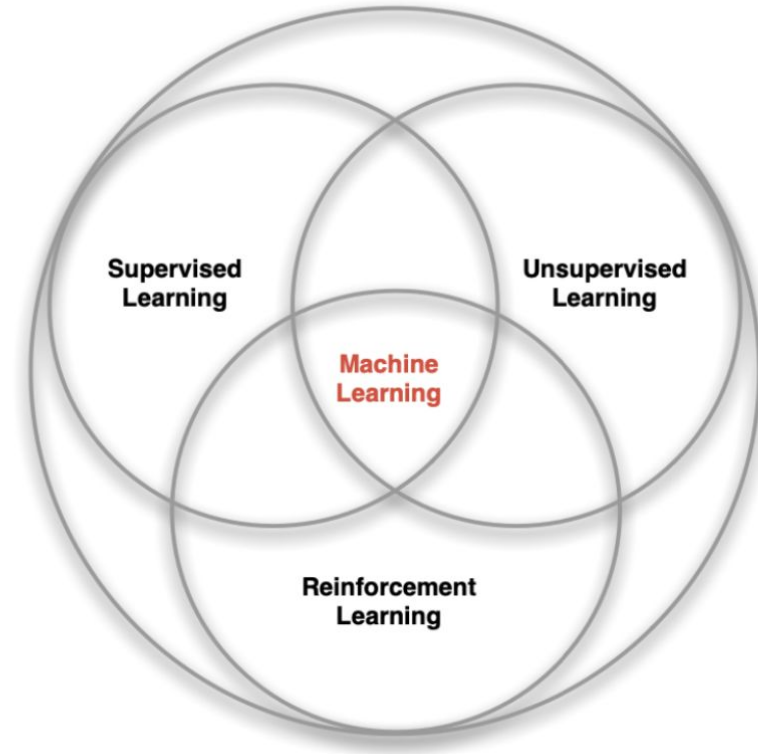
Разработка интеллектуальных агентов компьютерных игр

Веселко Никита Игоревич

- CEO и Lead Product Owner студии Винторог (ex IThub games) и студии Contrast Games
- CPO и Co-Founder в Current Agency Development
- CEO в HashCats
- Аспирант ФКН НИУ ВШЭ
- Преподаватель НИУ ВШЭ и школы Нарраторика
- Соавтор книги “Программирование в Unreal Engine 5”



Reinforcement Learning in Machine Learning



Supervised Learning

- Train sample: X, Y
- Approximation $Y \approx \hat{y} = f(X)$, f - семейство алгоритмов
- Loss minimization: $L(Y, \hat{y}) \rightarrow \min$ по f

Sequential Decision Making

Baby learning



Sequential Decision Making

Self-driving car



Sequential Decision Making

Self-driving car



RL

- Отображений ситуации на действия, максимизируя численный сигнал - вознаграждения
- Агент сам должен сам понять, какое действие в какой ситуации приносит максимальную награду
 - В сложных моделях понять, как действие повлияет и на награды за следующие действия
- Восприятие, действие и цель

RL vs (Un)Supervised Learning

- Supervised Learning

- Есть обучающая размеченная выборка с метками правильного действия под каждую ситуацию
- Модель экстраполирует или обобщает свою реакцию на ситуации, которые не предъявлены в обучающей выборке
- Не подходит для обучения с помощью взаимодействия, так как в интерактивных задачах нельзя получить примеры желаемого поведения

- Upsupervised Learning

- Цель: обнаружить структуру, скрытую в наборе данных, а не максимизировать награду

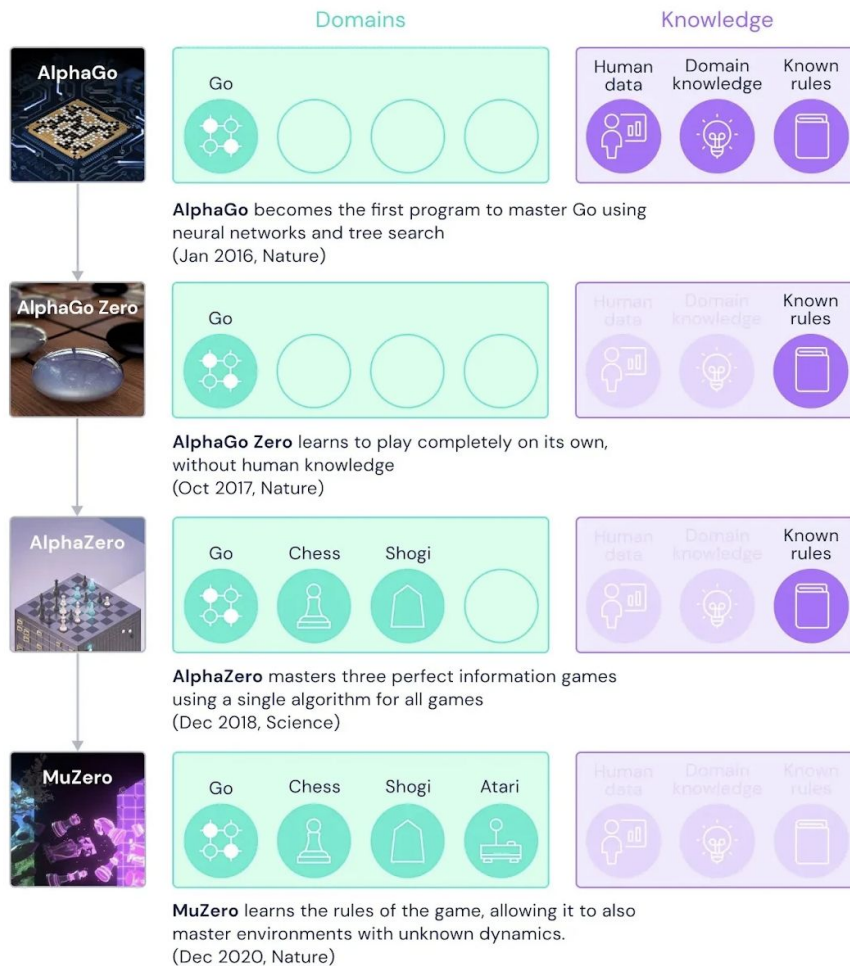
- RL

- Компромисс между исследованием и использованием
- Рассмотрение целостной проблемы
 - В Sup не ставится вопрос о конечной пользе приобретенных способностей обобщения
-

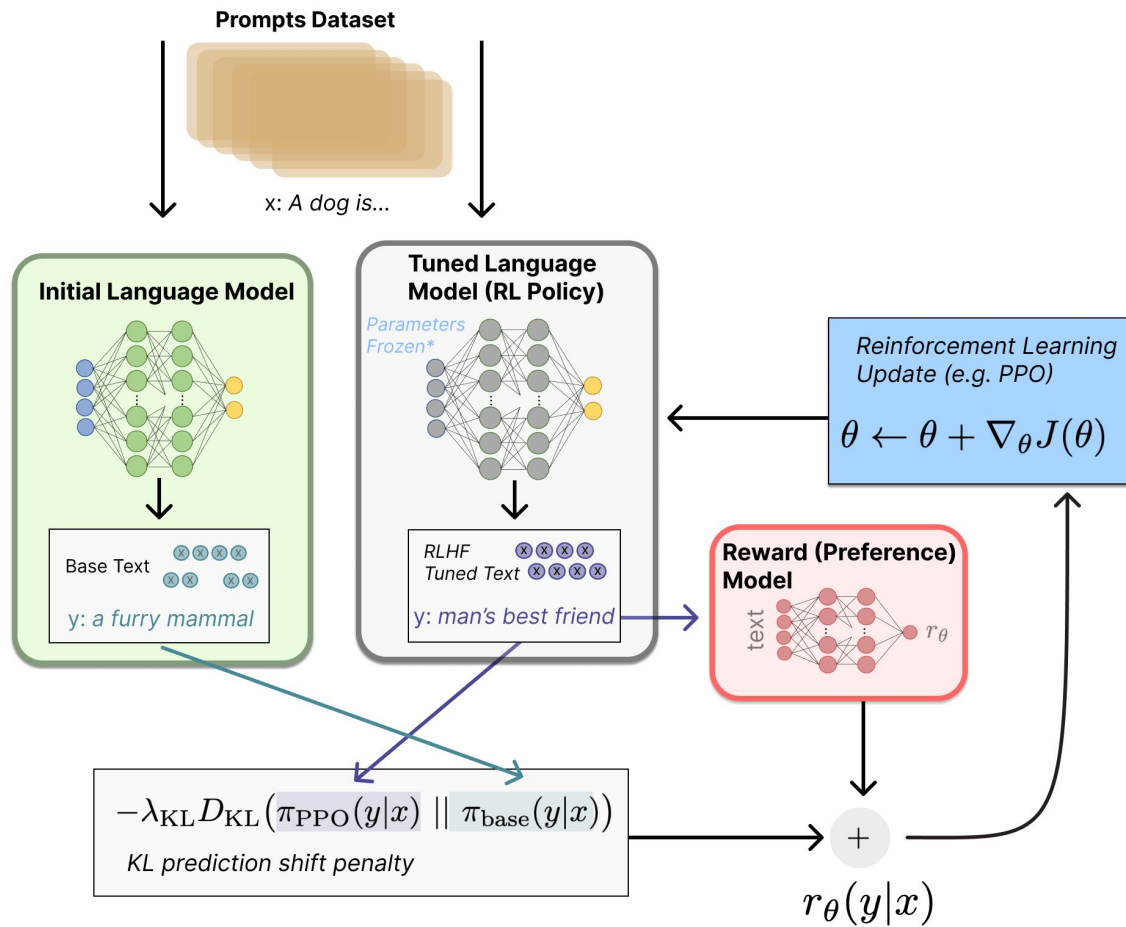
Характеристики RL

- Нет учителя, только награды
- Фидбек отложенный, а не мгновенный
- Время имеет значение (последовательность)
- Действия агента влияют на последующие данные

Пример



Пример

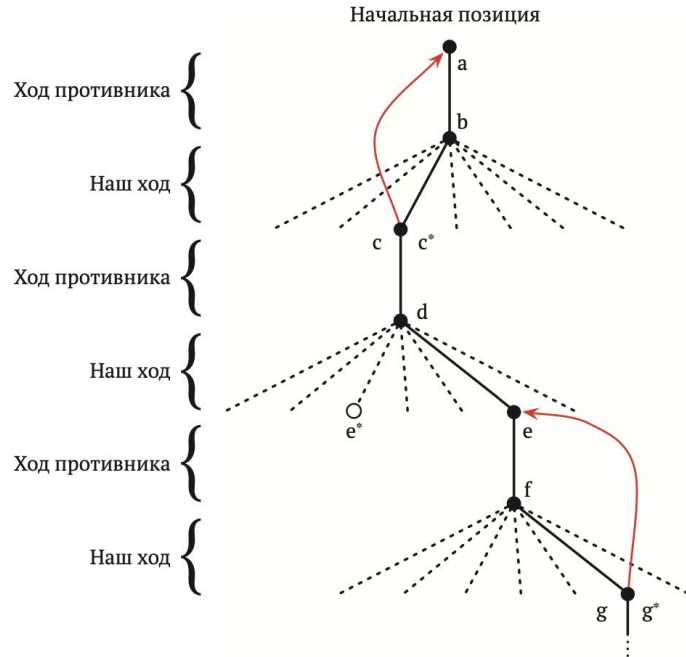


Элементы RL

- Агент
- Окружающая среда
- Стратегия – отображение множества воспринимаемых состояний среды на действия, предпринимаемые в этих состояниях
- Сигнал вознаграждения – определяет цель в задаче обучения с подкреплением
- Функция ценности – это полное вознаграждение, которого агент может ожидать в будущем, если начнет работу в этом состоянии
- *Модель окружающей среды – используются для планирования

Пример “Крестики-Нолики”

- Таблица чисел для каждого состояния игры - последняя оценка вероятности выиграть, начав с этого состояния - функция ценности



На подумать

- Предположим, что вместо игры со случайным противником описанный выше алгоритм обучения с подкреплением играет против себя самого, и обе стороны обучаются. Как вы думаете, что произойдет в таком случае? Обучится ли игрок другой стратегии выбора ходов?

Награды

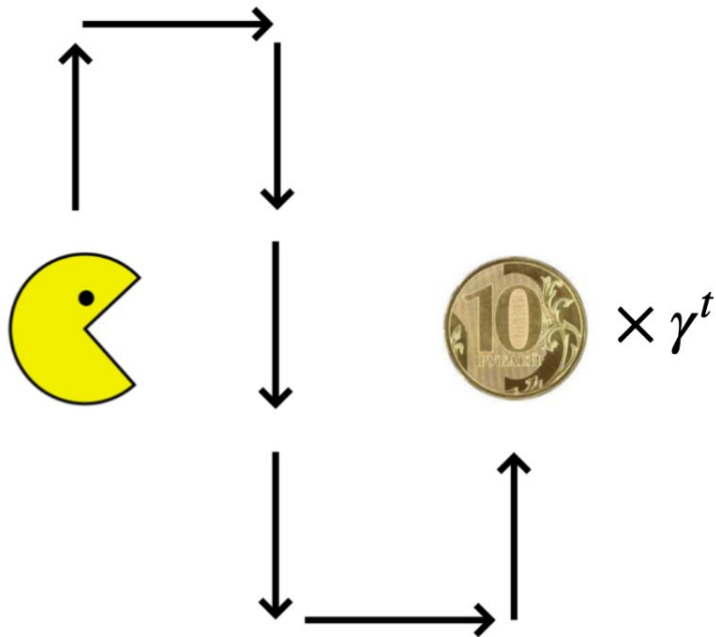
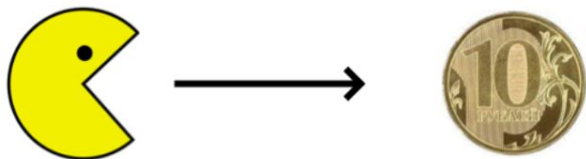
- Награда R_t - скалярный сигнал вознаграждения
- Показывает, на сколько хорошее действие принял агент на шаге t
- Задача агента - максимизировать кумулятивную сумму наград (функцию ценности)
- The reward hypothesis
 - All goals can be described by the maximisation of expected cumulative reward

Дизайн наград



Дисконтирование

$\gamma \in [0,1]$ is a discount factor

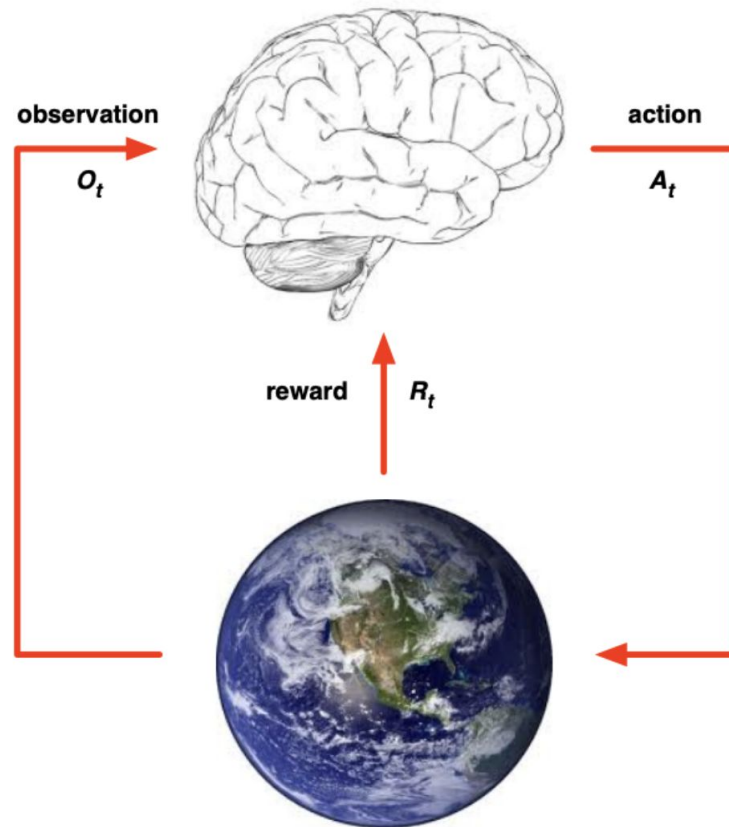


Зачем дисконтировать

- Избегать бесконечные награды
- Неопределенность в отношении будущего
- Поведение людей демонстрирует предпочтение немедленным наградам

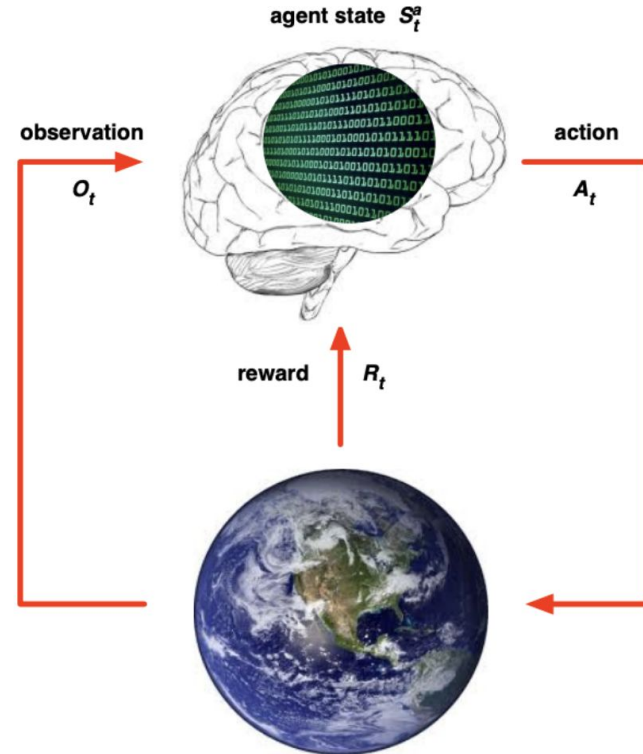
Агент и среда

- В каждом шаге t агент начинает в O_t
 - Выполняет действие A_t
 - Получает награду R_t
 - Получает следующее наблюдение O_{t+1}
- На каждом шаге t среда
 - Получает действие A_t
 - Передает награду R_t
 - Передает наблюдение O_t



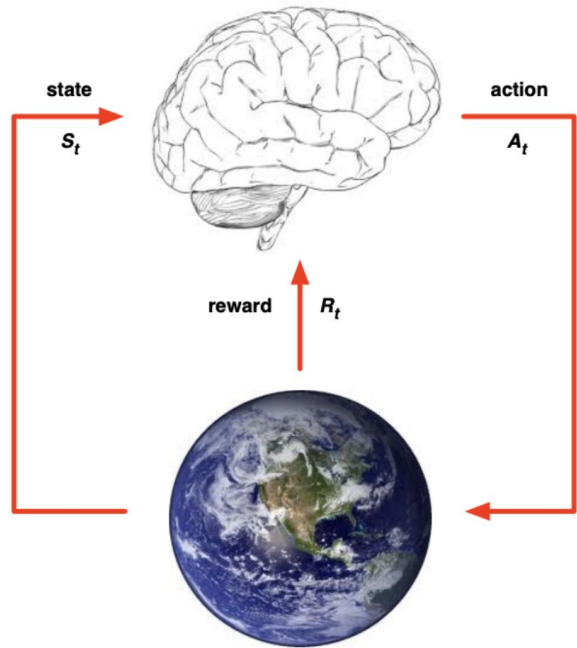
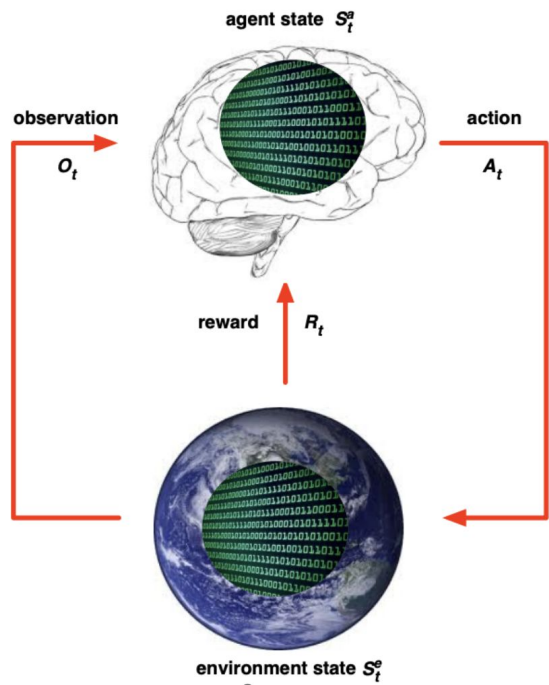
Состояние (state)

- Состояние среды S_t - закрытое представление среды
- Состояние агента - внутреннее представление агента



Observability (наблюдаемость)

- Частично наблюдаемы: Состояние среды и агента не равны
- Полностью наблюдаемы: Состояние среды и агента равны



Политика (Policy)

- Политика полностью определяет поведение агента (стратегия)
- Отображение из состояний в действия
- Детерминированная политика $a = \pi(s)$
- Стохастическая политика $\pi(a | s) = P[A_t = a | S_t = s]$

В итоге цель

$$\mathbb{E}_{\pi} \left[\sum_{t=0}^T \gamma^t R_t \right] \rightarrow \max_{\pi}$$