

Научная статья

УДК 517.977.5

URL: <https://trudymai.ru/published.php?ID=186897>

EDN: <https://www.elibrary.ru/AFJZPG>

## ОПТИМИЗАЦИЯ АЛГОРИТМА ПРЕДОТВРАЩЕНИЯ СТОЛКНОВЕНИЙ В ВОЗДУХЕ НА ОСНОВЕ ОБУЧЕНИЯ С ПОДКРЕПЛЕНИЕМ С РЕСУРСНЫМИ ОГРАНИЧЕНИЯМИ

Е.С. НЕРЕТИН<sup>1✉</sup>, Ц. ЛЮ<sup>2</sup>

<sup>1</sup>Филиал ПАО "Яковлев" - Центр комплексирования, Москва, Россия <sup>2</sup>Северо-Западный политехнический университет, Сиань, Китай

✉ [evgeny.neretin@ic.yakovlev.ru](mailto:evgeny.neretin@ic.yakovlev.ru)

---

**Цитирование:** Неретин Е.С., Лю Ц. Оптимизация алгоритма предотвращения столкновений в воздухе на основе обучения с подкреплением с ресурсными ограничениями // Труды МАИ. 2025. № 145. URL: <https://trudymai.ru/published.php?ID=186897>

---

**Аннотация.** С увеличением плотности воздушного трафика возрастает необходимость в эффективных системах предотвращения столкновений в воздухе. Традиционные системы, такие как TCAS, хотя и эффективно поддерживают безопасность, сталкиваются с трудностями в адаптации и оптимизации в современных сложных условиях. Чтобы преодолеть эти ограничения, мы применяем обучение с подкреплением (RL) в рамках марковского процесса принятия решений с ограничениями по ресурсам (RC-MDP), вводя управление виртуальными ресурсами для сокращения числа ложных тревог. Мы предлагаем бонус за время и ресурсы (TRB) для модификации алгоритмов DQN и SAC в DQNTRB и SACTRB, которые поощряют эффективное использование ресурсов при сохранении эффективности предотвращения столкновений. Результаты экспериментов показывают, что эти

модифицированные алгоритмы значительно сокращают количество ложных тревог, достигая почти аналогичной эффективности по сравнению с алгоритмами без ограничений.

**Ключевые слова:** реакция пилота, глубокое обучение с подкреплением, воздушное столкновение, марковский процесс принятия решений, динамическое программирование

.....

## OPTIMIZATION OF AIRBORNE COLLISION AVOIDANCE ALGORITHM BASED ON RESOURCE-CONSTRAINED REINFORCEMENT LEARNING

E.S. NERETIN<sup>1</sup>, L. ZUOCHENG<sup>2</sup>

<sup>1</sup>Branch of PJSC Yakovlev – integration center, Moscow, Russia

<sup>2</sup>Northwestern Polytechnical University, Xi'An, P. R. China

 [evgeny.neretin@ic.yakovlev.ru](mailto:evgeny.neretin@ic.yakovlev.ru)

---

**Citation:** Neretin E.S., Zuocheng L. Optimization of airborne collision avoidance algorithm based on resource-constrained reinforcement learning // Trudy MAI. 2025. No. 145. (In Russ.). URL: <https://trudymai.ru/published.php?ID=186897>

---

**Abstract.** As air traffic density continues to rise with the advancement of aviation technology, the demand for efficient and reliable airborne collision avoidance systems becomes increasingly urgent. Traditional systems, such as the Traffic Collision Avoidance System (TCAS), mainly rely on heuristic rules and parameter settings, which, although effective in maintaining safety, struggle to adapt and optimize under the complexities of modern aviation environments. To address these limitations, we explore the application of reinforcement learning (RL) to optimize the performance of collision avoidance systems. We define the problem within a resource-constrained Markov decision process (RC-MDP) framework, incorporating virtual resource management to control the

frequency of nuisance alerts, which are frequent alarms that do not require actual evasive action. We propose a novel time-resource bonus (TRB) mechanism to modify and enhance two standard RL algorithms, DQN and SAC, into DQNTRB and SACTRB. This approach encourages resource-efficient actions while maintaining collision avoidance performance. Our experimental results demonstrate that these modified algorithms significantly reduce nuisance alerts while achieving near-equivalent collision avoidance performance compared to algorithms without resource constraints.

**Keywords:** pilot respond, deep reinforcement learning, airborne collision avoidance, markov decision-making process, dynamic programming.

---

## Введение

С быстрым развитием авиационных технологий плотность воздушного трафика неуклонно увеличивается, что ведет к необходимости более эффективных и надежных систем предотвращения столкновений в воздухе. Традиционные системы, такие как система предотвращения столкновений в воздухе (TCAS), в основном полагаются на эвристические правила и настройки параметров, которые, хотя и эффективны для обеспечения безопасности, сталкиваются с проблемами при оптимизации и развитии в условиях все более сложных авиационных сред [1, 2]. Недавние исследования сосредоточены на повышении автоматизации и надежности алгоритмов предотвращения столкновений [3, 4, 5], улучшении использования аппаратных средств в авиационных системах [6, 7], решении проблем предотвращения столкновений между несколькими воздушными судами [8], а также использовании технологии автоматического зависящего наблюдения – вещания (ADS-B) для повышения эффективности предотвращения столкновений [9].

Обучение с подкреплением (RL), метод безнадзорного обучения, широко применяется для решения различных задач по уклонению от препятствий и принятия решения [8, 10, 11]. Основная идея RL заключается в постоянной оптимизации политик путем взаимодействия с окружающей средой для

максимизации накопленных вознаграждений . При применении традиционных алгоритмов RL к задаче предотвращения столкновении воздушного движения мы обнаружили, что, хотя RL эффективно решает задачи предотвращения столкновении и планирования трафика, его производительность в уменьшении количества ложных тревог оказалась неудовлетворительной [12]. Для решения этой проблемы мы ввели концепцию виртуальных ресурсов — когда каждое предупреждение о столкновении потребляет определенное количество виртуальных ресурсов, и система запрещает выдачу дальнейших предупреждений при их исчерпании. Эта задача может рассматриваться как проблема RL с ограниченными ресурсами, где ключом к успешному принятию решения является эффективное управление ресурсами при выполнении основной задачи [13, 14].

В данной статье сначала излагаются основные принципы алгоритмов RL, которые мы применили к задаче предотвращения столкновении , включая алгоритмы Deep Q-Network (DQN) и Soft Actor-Critic (SAC). Затем мы устанавливаем основную структуру задачи конфликта между двумя воздушными судами, охватывая пространства состояний и действий , функции вознаграждения для основной и вспомогательной задач, а также показатели оценки производительности алгоритмов. Мы анализируем различия в производительности этих двух алгоритмов при решении задач предотвращения столкновении как с ресурсными ограничениями, так и без них. Мы обнаружили, что в условиях ограничения по ресурсам оба алгоритма демонстрируют слабые результаты и низкую эффективность использования образцов. Дополнительный экспериментальный анализ показал, что исследование агентом среды часто ограничивается количеством ресурсов, и доступные действия зависят от оставшихся ресурсов. Оба алгоритма, DQN и SAC, склонны быстро расходовать ресурсы, ограничивая последующую разведку.

Для решения этой проблемы мы расширили марковский процесс принятия решения (MDP) до MDP с ограниченными ресурсами, введя управление ресурсами [14, 15], и предложили новый механизм бонуса за время и ресурсы (TRB). Мы

модифицировали стандартные алгоритмы SAC и DQN в SACTRB и DQNTRB соответственно. Эти алгоритмы побуждают агента экономить ресурсы, выполняя задачи по предотвращению столкновений и сокращая количество ложных тревог, тем самым исследуя состояния с более широким набором возможных опций. Наконец, мы провели эксперименты для оценки производительности оптимизированных алгоритмов. Результаты показывают, что в условиях ограниченных ресурсов оптимизированные алгоритмы значительно превосходят традиционные алгоритмы как в эффективности предотвращения столкновений, так и в снижении количества ложных тревог. По сравнению с алгоритмами без ограничения по ресурсам оптимизированные алгоритмы достигли почти аналогичной производительности в планировании столкновений, значительно сократив частоту ложных тревог.

### Предварительные сведения

В этом разделе мы представляем модельный каркас и основные алгоритмы, используемые в нашем исследовании.

#### 1) Марковский процесс принятия решений

Марковский процесс принятия решений — это математическая структура, используемая для моделирования последовательных задач принятия решений, которая определяется кортежем  $(S, A, p, r, \gamma)$ , где  $s_t \in S$  — состояние в момент времени  $t$ ,  $a_t \in A$  — действие, совершаемое агентом в момент времени  $t$  в результате процесса принятия решений,  $r_t = R(s_t, a_t, s_{t+1})$  — вознаграждение, получаемое агентом в результате выполнения действия  $a_t$  из состояния  $s_t$  и перехода в состояние  $s_{t+1}$ , а  $p(s_{t+1}, a, s_t)$  — функция переходов, которая отображает вероятность

$p(s_{t+1} | s_t, a_t)$  перехода в состояние  $s_{t+1}$  при выполнении действия  $a_t$  из состояния

$s_t$ .  $\gamma \in [0, 1]$  — это коэффициент дисконтирования, используемый для взвешивания мгновенных вознаграждений по отношению к будущим. Цель MDP заключается в нахождении оптимальной стратегии  $\pi^*$ , которая максимизирует ожидаемую накопленную дисконтированную награду:

$$V(s) = E_{\pi} [ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) | s_0 = s ] \quad (1)$$

## 2) Глубокая Q-сеть (DQN)

Глубокая Q-сеть (DQN) [16, 17] — это алгоритм обучения с подкреплением, основанный на Q-обучении, который использует нейронные сети для аппроксимации функции Q-значений  $Q(s, a)$ , решая проблему многомерных пространств состояний. Функция Q-значений представляет ожидаемую накопленную награду за выполнение действия  $a$  в состоянии  $s$ . DQN обновляет Q-значения по следующей формуле:

$$Q(s, a) \leftarrow \alpha Q(s, a) + (1 - \alpha) [r + \gamma \max_{a'} Q(s', a')] \quad (2)$$

Где  $\alpha$  — это скорость обучения,  $\gamma$  — это коэффициент дисконтирования,  $r$  — это немедленное вознаграждение, а  $s'$  — это следующее состояние после выполнения действия  $a$ . DQN также использует повторное проигрывание опыта для уменьшения корреляции между выборками и целевую сеть для улучшения стабильности обучения.

## 3) Мягкий актор-критик (SAC)

Мягкий актор-критик (SAC) [18] — это алгоритм обучения с подкреплением (off-policy), который объединяет методы градиента политики с регуляризацией энтропии для балансировки исследования и эксплуатации. Цель SAC — максимизировать как ожидаемое вознаграждение, так и энтропию политики, поощряя более исследовательское поведение. Оптимизационная цель:

$$J(\pi) = E_{s \sim p_0} [ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) ] - \lambda H(\pi) \quad (3)$$

Где  $\alpha$  — это параметр температуры, а  $H(\pi|s)$  — это энтропия политики в состоянии  $s_t$ . SAC особенно подходит для пространств с непрерывными действиями, улучшая исследование с помощью регуляризации энтропии. В данном исследовании, для применения алгоритма SAC к задаче предотвращения столкновений, мы внесли соответствующие корректировки как в пространство состояний, так и в пространство действий, чтобы лучше соответствовать операционным требованиям задачи.

### Постановка задачи

Мы формулируем задачу предотвращения столкновений в воздухе между двумя самолетами как MDP и решаем ее с помощью алгоритмов обучения с подкреплением. В этом разделе мы представляем динамическую модель самолетов в конфликтной ситуации и даем подробное описание пространств состояний и действий. Кроме того, мы строим функцию вознаграждения, основанную на целях предотвращения столкновений и уменьшения количества ложных тревог. Также мы устанавливаем метрики оценки для анализа эффективности алгоритмов.

#### 1) Определение пространства состояний и действий

Для поддержания согласованности с TCAS, в данной статье акцент делается на вертикальном избегании столкновений [4]. Пространство состояний для задачи предотвращения столкновений состоит из пяти переменных. Таблица 1 описывает эти переменные и их диапазоны, а на Рисунке 1 представлено их визуальное отображение. Первые три переменные описывают относительные положения и вертикальные скорости собственного и вторгающегося самолета. Четвертая переменная,  $t$ , обобщает горизонтальную геометрию, указывая время до того, как горизонтальное расстояние между двумя самолетами станет менее 500 футов, и это также определяется временем до конфликта (Time to Conflict, TTC). Включение предыдущего совета в пространство состояний позволяет нам штрафовать за изменения или усиления советов, сохраняя марковское свойство.

Переменные пространства состояний

Переменная	Описание	Значения	Единицы измерения
$h$	Относительная высота нарушителя	$[-2500, 2500]$	ft
$\dot{h}_o$	Вертикальная скорость собственного самолета	$[-70, 70]$	ft/s
$\dot{h}_i$	Вертикальная скорость нарушителя	$[-50, 50]$	ft/s
$\tau$	Время до потери горизонтального разделения (TTC)	$[0, 40]$	s
$a_{prev}$	Предыдущий совет	См. Таблица 2	-

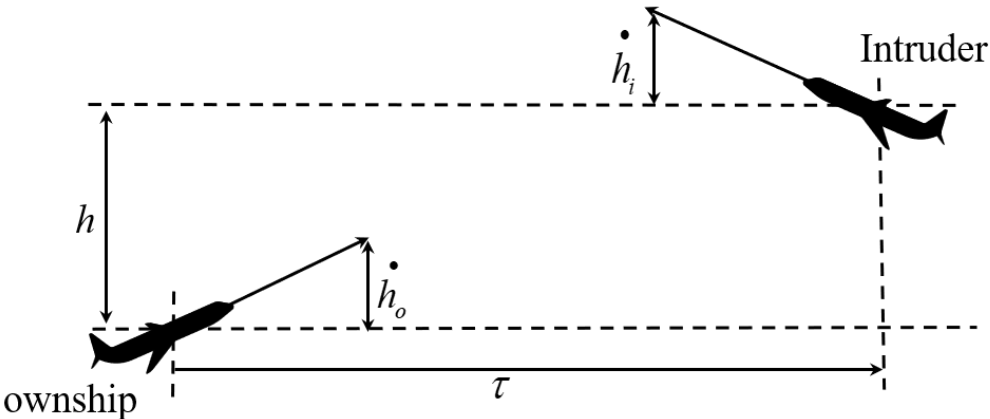


Рисунок 1 - Визуальное представление переменных состояния.

Пространство действий включает советы, которые система предотвращения столкновений может предоставить во время полета, всего 7 возможных советов, как показано в Таблице 2. Все советы, кроме СОС, вызывают предупреждение и направляют самолет в определенный диапазон вертикальных



скорости с соответствующим ускорением. Совет СОС указывает на отсутствие немедленной угрозы столкновения с вторгающимся самолетом.

Таблица 2

Набор советов

Действие	Описание	Ускорение
СОС	Нет конфликта	0
DES1500	Спуск $\leq$ -25ft/s	-g/3
CL1500	Взлет $\geq$ 25ft/s	g/3
SDES1500	Спуск $\leq$ -25ft/s	-g/2.5
SCL1500	Взлет $\geq$ 25ft/s	g/2.5
SDES2500	Спуск $\leq$ -42ft/s	-g/2.5
SCL2500	Взлет $\geq$ 42ft/s	g/2.5

Таблица 3 описывает доступность каждого совета в зависимости от текущего отображаемого совета. Например, СОС может быть выдан в любое время. Однако DES1500 и CL1500, будучи начальными советами, могут быть выданы только если в данный момент пилоту отображается СОС. SDES1500 может быть выдан после CL1500, SCL1500 и SCL2500, выступая как разворот, или после SDES2500, выступая как ослабление. Важно отметить, что SDES1500 не может следовать за СОС и также не может следовать за DES1500 из-за их сходства по своей природе. Таким образом, исходя из доступности каждого совета под текущим советом, агент на самом деле может выбирать только из трех действий в каждом состоянии.

Таблица 3

Доступность советов

Действие	Доступно от
DES1500	СОС
CL1500	СОС
SDES1500	CL1500, SCL1500, SCL2500, SDE2500

время	SCL1500	DES1500, SDE1500, SDES2500, SCL1500	COC	В любое
	SDES2500	DES1500, SDES1500		
	SCL2500	CL1500, SCL1500		

## 2) Динамическая модель самолета

Динамическая модель может быть записана как уравнение (4). Мы предполагаем временной шаг в одну секунду, что означает обновление системы предотвращения столкновений с частотой 1 Гц. Для усложнения симуляционной среды мы ограничиваем диапазон ускорения вторгающегося самолета значениями от  $[-a_{int}, a_{int}]$ , и предполагаем, что скорость вторгающегося самолета  $h_{int}$  изменяется в направлении состояния столкновения на каждом шаге.

$$\begin{aligned}
 & \dot{h} = \begin{cases} h_{int} - 0.5h_{int} - 0.5h_{own} & \text{if } h_{int} > h_{own} \\ h_{own} - 0.5h_{int} - 0.5h_{own} & \text{if } h_{int} < h_{own} \\ 0 & \text{otherwise} \end{cases} \\
 & \dot{h}_{int} = \begin{cases} a_{prev} & \text{if } h_{int} > h_{own} \\ -a_{prev} & \text{if } h_{int} < h_{own} \\ 0 & \text{otherwise} \end{cases} \\
 & \dot{a}_{prev} = \begin{cases} a_{max} & \text{if } h_{int} > h_{own} \\ -a_{max} & \text{if } h_{int} < h_{own} \\ 0 & \text{otherwise} \end{cases}
 \end{aligned}
 \tag{4}$$

## 3) Формирование награды

В задачах предотвращения столкновений функция награды должна балансировать между безопасностью и эффективностью. Поэтому наша цель – предотвратить столкновения, минимизируя при этом выдачу отвлекающих предупреждений и неожиданных действий системы (таких как усиление

рекомендации, изменение рекомендации или создание пересечения высот между самолетами). Кроме того, чтобы минимизировать влияние на другое воздушное пространство во время разрешения конфликта, мы поощряем весь процесс разрешения оставаться в определенном диапазоне высот. Таким образом, дизайн функции награды в первую очередь разделен на две части.

Во-первых, в конце ТТС нам необходимо, чтобы относительная высота самолетов поддерживалась на определенной высоте  $h_{rel}$ , и штрафы налагаются за другие состояния, такие как столкновение, чрезмерно высокая или низкая относительная высота. Первая часть функции награды может быть сформулирована следующим образом:

$$R_1 = \frac{1}{w_{NMAC}} \left( \frac{1}{h_{rel} - h_{col}} \right) \frac{1}{h_{rel} - h_{rel\_s}} \quad (5)$$

$$= \max \left( \frac{1}{2} \frac{1}{h_{rel} - h_{rel\_s}}, 25 \right) \frac{1}{h_{rel} - h_{rel\_s}}$$

Где  $w_{NMAC}$  штрафует за столкновение между собственным самолетом и нарушителем,  $h_{rel}$  представляет относительную высоту между двумя самолетами,  $h_{col}$  указывает высоту, на которой происходит конфликт между двумя самолетами,  $h_{rel\_s}$  обозначает желаемую относительную высоту, которую мы хотим, чтобы два самолета достигли, а  $w_{leav}$  представляет штраф, накладываемый, когда относительная высота между двумя самолетами чрезмерно велика.

В течение периода ТТС нам необходимо обращать внимание на высоту собственного самолета, изменения в статусе предупреждения и ситуацию пересечения высот между двумя самолетами. Функция награды может быть сформулирована следующим образом:

$$R_2 = \frac{1}{w_{leav}} \left( \frac{1}{h_{down} - h_{upper}} \right) \frac{1}{\tau_{coc}} \left( \frac{1}{reversal} + \frac{1}{strength} + \frac{1}{crossing} \right) \exp(-t) \quad (6)$$

$w_{leav}$  представляет штраф, накладываемый, когда высота собственного

самолета ( $h_{own}$ ) превышает заданный верхний предел диапазона высот ( $h_{upper}$ ).  $w_{Alert}$ ,  $w_{stren}$ ,  $w_{reversal}$ ,  $w_{cros}$  штрафуют систему за выдачу предупреждений и ложных тревог. Среди них,  $\tau$  представляет Время до Столкновения (ТТС), и по мере приближения времени к ТТС, штрафы за предупреждения и другие ложные тревоги усиливаются.  $w_{cros}$  предоставляет небольшую награду, когда предупреждение снято.

#### 4) Метрики оценки производительности алгоритма

Для количественной оценки эффективности алгоритма в решении проблемы предотвращения столкновений мы использовали следующие метрики оценки:

- **Результаты сходимости средней награды:** средняя награда за последние 100 эпизодов.
- **Результаты сходимости частоты столкновений:** средняя частота столкновений за последние 100 эпизодов.
- **Результаты сходимости коэффициента успеха координации высоты (ACSR):** средняя доля эпизодов за последние 100 эпизодов, в которых относительная высота самолетов была больше 900 футов, но меньше 1500 футов.
- **Результаты сходимости среднего количества предупреждений:** среднее количество предупреждений, выданных за эпизод, за последние 100 эпизодов.
- **Результаты сходимости среднего количества усиления советов:** среднее количество усиления советов за эпизод за последние 100 эпизодов.
- **Результаты сходимости среднего количества отмен советов:** среднее количество отмен советов за эпизод за последние 100 эпизодов.
- **Результаты сходимости среднего количества пересечений высот:** среднее количество раз, когда высоты обоих самолетов пересекались за эпизод, за последние 100 эпизодов.

## Метод

В этом разделе мы представляем метод оптимизации алгоритма предотвращения воздушных столкновений с использованием ограниченных ресурсов обучения с подкреплением. Мы подробно рассматриваем формулировку RL с ограничением ресурсов. Наши эксперименты показывают, что информация, связанная с ресурсами, имеет ключевое значение для эффективного исследования и реализации основных функций системы предотвращения столкновений. Однако традиционные алгоритмы RL, как правило, игнорируют информацию о ресурсах и предполагают, что среда полностью известна. Чтобы решить эту проблему, мы моделируем задачу предотвращения столкновений с использованием RL с ограничением ресурсов, включая виртуальные ресурсы, и предлагаем подробное объяснение предлагаемого метода Time Resource Bonus (TRB).

### 1) Ресурсоограниченный MDP

В обучении с подкреплением (RL) определённые действия требуют ресурсов, таких как энергия в робототехнике или потребляемые предметы в видеоиграх. Пусть имеется  $d$  типов ресурсов. Обозначим вектор ресурсов как  $\mathbf{r} \in \mathbb{R}^d$ , а множество всех возможных векторов ресурсов — как  $S_r \subseteq \mathbb{R}^d$ . Пространство состояний в проблемах с ограничением ресурсов расширяется до  $S_R \subseteq S \times S_{\mathbf{r}} \subseteq S \times S_o \times S_r$ , что позволяет алгоритмам изучать информацию, связанную с ресурсами.

Для эффективного использования этой информации мы определяем функцию, учитывающую ресурсы,  $I: S \rightarrow \mathbb{R}^d$ , которая отображает состояния в ресурсы. Также мы вводим детерминированную функцию перехода ресурсов  $f: S \times S_r \rightarrow S_r$ , которая считается неизвестной, так как динамика ресурсов часто недоступна в реальных задачах. Распределение вероятности перехода задаётся следующим образом:

$$p_R \left[ s \mid s, a \right], \quad p \left[ s \mid s, a \right]_{o \neq o}, \quad p \left[ I \mid s \right]_r \neq f \left[ s, a \right], \quad (7)$$

где  $s \in \mathcal{S}_{s_o, r}$ ,  $a \in \mathcal{A}_{s \in \mathcal{S}_{s_o, r}}$ . Мы обозначаем распределение вероятностей перехода в задачах с ограниченными ресурсами как  $p$ .

Проблемы с ограничением ресурсов формулируются в виде кортежа  $\langle \mathcal{S}_R, \dots, A, p, r \rangle$ . Набор доступных действий для состояния  $s$  обозначается как  $A(s)$ , и  $\bigcup_{s \in \mathcal{S}_R} A(s) \subseteq \mathcal{A}$ , причём  $A(s)$  зависит от оставшихся ресурсов.

Допустимая политика  $\pi \in \Pi_s$  является функцией плотности вероятности над  $A(s)$ , и мы определяем множество допустимых политик как  $\Pi$ . Подобно традиционным RL, ресурсоограниченное обучение с подкреплением нацелено на решение задачи оптимизации (1) для нахождения оптимальной допустимой политики.

Состояние  $s_2$  доступно из  $s_1$ , если существует стационарная политика  $\pi$  и временной шаг  $t \in \mathbb{N}$ , такой что  $p(s \mid s^t)_{s_1} > 0$ . Если ресурсы невозполнимы и  $\pi \in \Pi_d, I_{s_i} \in \Pi \mid I_{s_i} \in \Pi$ , то состояние  $s_2$  недоступно из  $s_1$ .

Количество доступных ресурсов в состоянии  $s$  определяется как  $I_s \in \Pi \mid I_{s_1} \in \Pi, I_d \in \Pi$ , где  $I_{s_i} \in \Pi$  для всех  $i \in \mathcal{I}_d$ . Ресурсы  $i$ -го типа являются невозполнимыми, если их количество монотонно убывает, то есть  $I_{s_i} \in \Pi \mid I_{s_i} \in \Pi$ . Мы предполагаем, что  $I$  известно заранее, так как информация о текущих ресурсах обычно доступна в реальных задачах.

## 2) Проблема с ресурсными ограничениями в задаче предотвращения столкновений

При применении обучения с подкреплением для решения проблемы предотвращения столкновений в воздухе, несмотря на то что мы штрафovali за ненужные предупреждения, такие как эскалация и отмена рекомендаций через функцию вознаграждения, мы фактически позволяли этим предупреждениям

происходить бесконечно. В результате, хотя система могла завершить задачу по предотвращению столкновений, пилот подвергался множеству ненужных уведомлений. Чтобы решить эту проблему, мы ввели невозстанавливаемый виртуальный ресурс, при этом разные типы действий системы потребляют различные объемы ресурсов. Как только ресурсы исчерпаны, система больше не будет выдавать рекомендации. В рамках этой настройки агент должен не только завершить задачу для достижения высоких вознаграждений, но и минимизировать потребление ресурсов.

На основе условий задачи с ограничением ресурсов мы сначала продемонстрировали, что оба алгоритма RL, DQN и SAC, показали низкую эффективность использования образцов. Как показано на рисунке 2, оба алгоритма сходимости оказались менее успешными при наличии ограничений ресурсов по сравнению с ситуацией, когда таких ограничений не было. Мы также сравнили производительность обоих алгоритмов при различных уровнях ресурсов. Рисунки 3 и 4 иллюстрируют, что по мере увеличения доступности ресурсов производительность обоих алгоритмов улучшалась как в отношении предотвращения столкновений, так и в эффективности планирования. Таблицы 4, 5 и 6 соответствуют дополнительным показателям выполнения задач в ходе экспериментов, представленным на рисунках 2, 3 и 4 соответственно.

Результаты показывают, что после применения ограничений ресурсов количество предупреждений, выданных системой, значительно сократилось, но задачи по предотвращению столкновений и планированию были выполнены не на должном уровне. Рисунок 5 показывает траектории предотвращения столкновений во время обучения при ресурсных ограничениях для обоих алгоритмов. Можно увидеть, что оба алгоритма склонны выдавать непрерывные предупреждения на ранних стадиях сценария столкновения, тем самым быстро исчерпывая ресурсы. Таким образом, хотя эти методы частично решают проблему предотвращения столкновений, они не находят оптимальное решение.

В целом, эти методы страдают от неэффективного исследования в задаче предотвращения столкновений с ограничением ресурсов.

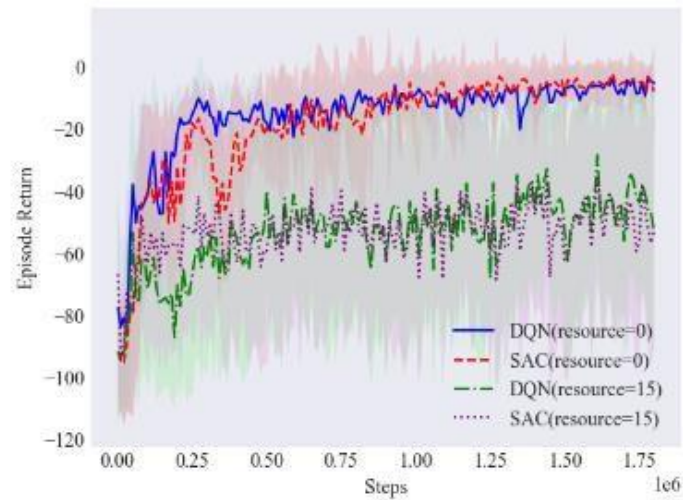


Рисунок 2 - Производительность обучения DQN и SAC в задачах предотвращения столкновений с ограничениями ресурсов и без них.

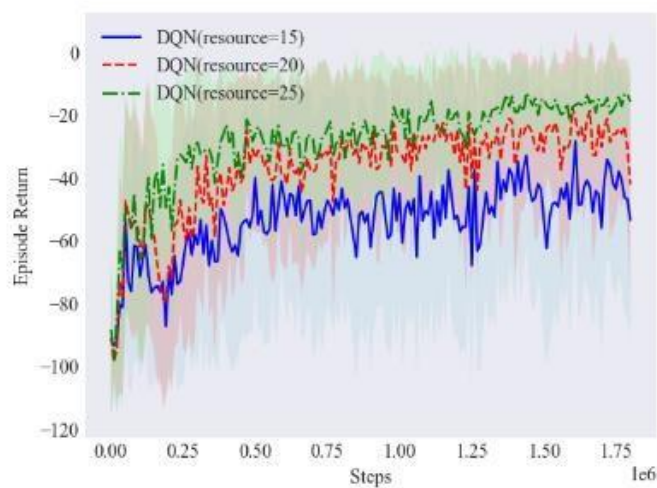


Рисунок 3 - Производительность обучения DQN в задачах предотвращения столкновений при различных ограничениях ресурсов.



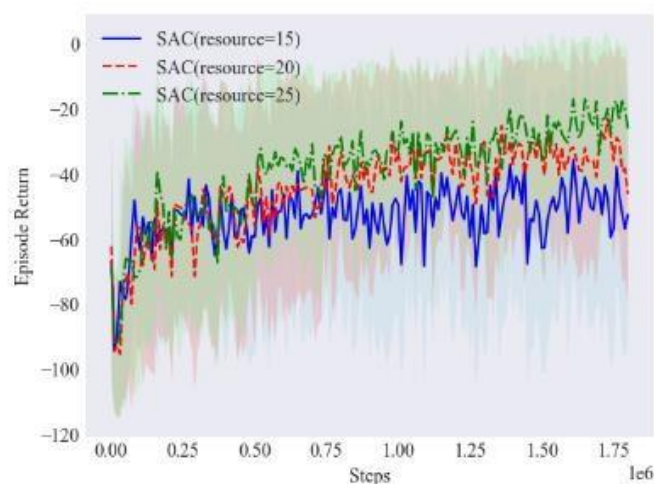


Рисунок 4 - Производительность обучения SAC в задачах предотвращения столкновений при различных ограничениях ресурсов.

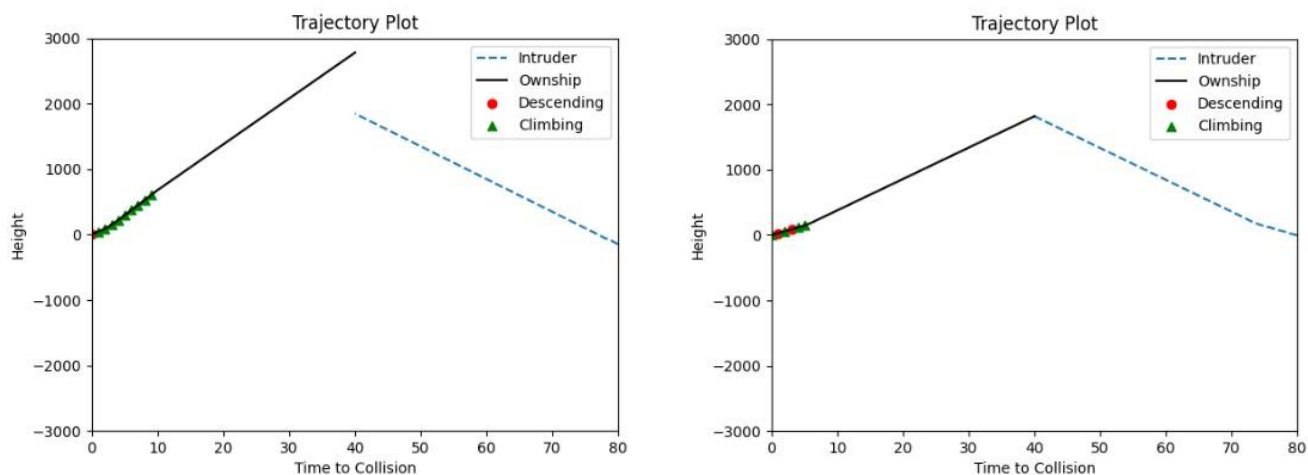


Рисунок 5 - Примеры траекторий предотвращения столкновений во время обучения для DQN (слева) и SAC (справа) при ограничениях ресурсов.

Таблица 4  
Производительность dqn и sac в задачах предотвращения столкновений с ограничениями ресурсов и без них

Алгоритм	Частота столкновений	Уровень успеха координации	Количество оповещений	Количество укреплений	Количество изменений	Количество пересечений
				оповещений	оповещений	высоты
<b>DQN(ресурс=0)</b>	0.021	0.900	6.140	2.011	7.190	0.082
<b>SAC(ресурс=0)</b>	0.045	0.830	4.578	2.734	7.581	0.111
<b>DQN(ресурс=15)</b>	0.292	0.380	0.939	0.474	2.406	0.138

<b>SAC(ресурс=15)</b>	0.262	0.372	0.796	0.427	2.514	0.175
-----------------------	-------	-------	-------	-------	-------	-------

Таблица 5

Производительность dqn в задачах предотвращения столкновений с различными ограничениями ресурсов

Алгоритм	Частота столкновений	Уровень успеха координации	Количество оповещений	Количество		
				укреплений	изменений высоты	пересечений
<b>DQN(ресурс=15)</b>	0.292	0.380	0.939	0.474	2.406	0.138
<b>DQN(ресурс=20)</b>	0.120	0.573	1.104	0.629	3.421	0.118
<b>DQN(ресурс=25)</b>	0.056	0.692	1.372	0.608	3.989	0.108

Таблица 6

Производительность обучения sac в задачах предотвращения столкновений с различными ограничениями ресурсов

Алгоритм	Частота столкновений	Уровень успеха координации	Количество оповещений	Количество		
				укреплений	изменений высоты	пересечений
<b>SAC (ресурс=15)</b>	0.262	0.372	0.796	0.427	2.514	0.175
<b>SAC(ресурс=20)</b>	0.175	0.472	0.980	0.622	4.396	0.163
<b>SAC(ресурс=25)</b>	0.126	0.551	1.060	0.675	5.023	0.177

### 3) Time Resource Bonus

В задачах обучения с подкреплением с ограничениями ресурсов ресурсы играют решающую роль в эффективном исследовании. Мы наблюдали, что в данном состоянии размер доступного множества состояний часто положительно коррелирует с оставшимися доступными ресурсами, так как состояния с высокими ресурсами не могут быть достигнуты из состояний с низкими ресурсами. Ранее

использованные методы исследования на основе расширения [19] показали, что исследование состояний с большими доступными наборами состояний является важным для эффективного исследования. Другими словами, движение к состояниям с высокими ресурсами позволяет агенту достигать большого числа будущих состояний, тем самым улучшая его исследование окружающей среды. Однако при решении проблемы предотвращения столкновений с ограничениями виртуальных ресурсов оба алгоритма, DQN и SAC, не учитывают информацию, связанную с ресурсами, и склонны быстро исчерпывать ресурсы.

Чтобы обеспечить агенту по предотвращению столкновений поддержание высокой эффективности исследования при снижении количества ненужных предупреждений, мы предлагаем Time Resource Bonus (TRB). Этот алгоритм побуждает агента экономить ресурсы при выполнении задач по предотвращению столкновений и снижении количества ненужных предупреждений, таким образом исследуя состояния с более широким набором приемлемых состояний.

Этот метод вводит временной коэффициент дисконтирования в процессе потребления ресурсов. Мы определяем количество ресурсов, потребляемых агентом каждый раз, как функцию  $g$ , связанную с коэффициентом дисконтирования. Мы требуем, чтобы функция  $g$  была возрастающей функцией времени. Поэтому в данной статье мы определяем функцию  $g$  в следующей форме:

$$g(s,a) = \frac{C(s,a)}{T - t} \quad (8)$$

В этой формуле  $\alpha$  — это гиперпараметр, представляющий чувствительность потребления ресурсов к времени. Чем больше значение  $\alpha$ , тем ниже чувствительность потребления ресурсов к времени. В данной статье мы произвольно устанавливаем  $\alpha=T$ , где  $T$  — это максимальное значение интервала времени до конфликта (TTC) в процессе предотвращения столкновений.  $t$  обозначает текущее значение TTC, а  $C(s,a)$  представляет собой количество

потребляемых ресурсов при выполнении действия  $a$  в текущем состоянии  $s$ . Следовательно, мы определяем функцию, учитывающую ресурсы, как

$I_{s,i} = I_{s,i} - g C_{s,a} \cdot \rho_{i,i}$ . Мы используем внутреннее вознаграждение  $r_{int}$ , чтобы представить долю оставшегося виртуального ресурса, что влияет на исследование агента через форму внутренних вознаграждений.  $r_{int}$  принимает следующую форму:

$$r_{int}(s, a, \rho) = \frac{-I_{max} \cdot E_{s,a,g} \cdot \rho_{i,i}}{I_{max} \cdot E_{s,a,g} \cdot \rho_{i,i} + P_{s,i} \cdot T_g} \quad s \in \mathcal{S} \quad (9)$$

## Эксперименты и результаты

Наши эксперименты преследуют две основные цели:

- 1) проверить, превосходит ли алгоритм TRB алгоритмы DQN и SAC в одинаковых условиях ресурсов;
- 2) проанализировать сравнительную производительность алгоритма TRB с ограничениями виртуальных ресурсов и алгоритмов DQN и SAC без ограничений ресурсов.

### 1) Настройки симуляции

Параметры алгоритма установлены следующим образом. В соответствии с определением пространства состояний и действия алгоритм имеет 5 входов и 3 выхода. Архитектура нейронной сети как для актера, так и для критика состоит из двух скрытых слоев с 64 и 32 узлами, с коэффициентом обучения  $6e-5$  и порогом отсечения градиента 3.0.

Размер пакета составляет 64, длина горизонта для исследования — 512, размер буфера для воспроизведения —  $1e6$ , и сети повторно обновляются с использованием буфера воспроизведения, чтобы поддерживать малую потерю критика. Что касается формирования вознаграждения, коэффициент дисконтирования для будущих вознаграждений (гамма) установлен на уровне

0.99, а масштаб вознаграждения равен 1. Процесс обучения прекращается, если общее количество шагов превышает 1 миллион. Операционная система — Windows 11, а GPU — GeForce RTX 3050.

## 2) Экспериментальные результаты и анализ

Сначала мы сравнили производительность TRB с алгоритмами DQN и SAC в одинаковых условиях ресурсов. Мы эмпирически выбрали подходящее количество ресурсов (Ресурс = 20), которое отражает производительность алгоритмов. На рисунке 6 показано, что после включения внутреннего вознаграждения как

DQNTRB, так и SACTRB значительно превосходят стандартные алгоритмы DQN и SAC. Более того, хотя SAC в настоящее время считается современным алгоритмом (state-of-the-art), DQN и DQNTRB достигают лучших результатов сходимости по сравнению с SAC и SACTRB в контексте предупреждения о предотвращении столкновения в воздухе. Конкретные параметры предотвращения столкновения представлены в таблице 7.

Можно наблюдать, что DQNTRB и SACTRB снижают частоту столкновения и повышают уровень успешного согласования высоты по сравнению с DQN и SAC, но также влекут за собой более высокие затраты в отношении предупреждения (с увеличением ненужных предупреждений).

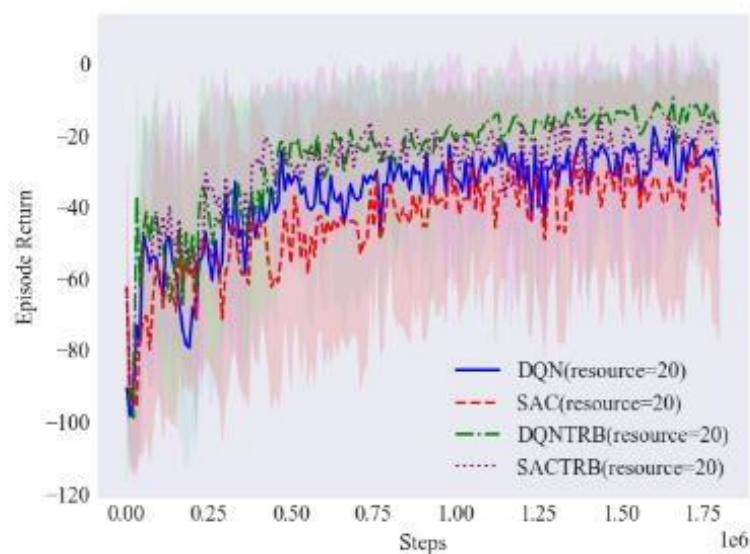


Рисунок 6 - Производительность обучения четырех алгоритмов — DQN, DQNTRB, SAC и SACTRB — в задачах предотвращения столкновений с одинаковыми ограничениями ресурсов.

Таблица 7

Производительность алгоритмов dqn, dqntrb, sac и sactrb в задачах предотвращения столкновений с одинаковыми ограничениями ресурсов

Алгоритм	Количество		Уровень изменений	Количество		Количество	Частота
	успеха	укреплений		пересечений	столкновений		
	координации	оповещений	оповещений	оповещений	высоты	оповещений	
DQN(ресурс=20)	0.120	0.573	1.104	0.628	3.421	0.118	
SAC(ресурс=20)	0.175	0.472	0.979	0.622	4.396	0.164	
DQNTRB(ресурс=20)	0.056	0.717	1.856	0.729	3.436	0.123	
SACTRB(ресурс=20)	0.078	0.657	1.565	0.889	5.352	0.114	

Далее мы сравнили производительность обучения DQNTRB и SACTRB при различных ограничениях ресурсов с производительностью предотвращения столкновении DQN и SAC без ограничении ресурсов.

На рисунке 7 показана производительность обучения DQNTRB при различных ограничениях ресурсов, а также производительность обучения неограниченного алгоритма DQN в среде предупреждения о предотвращении столкновении .

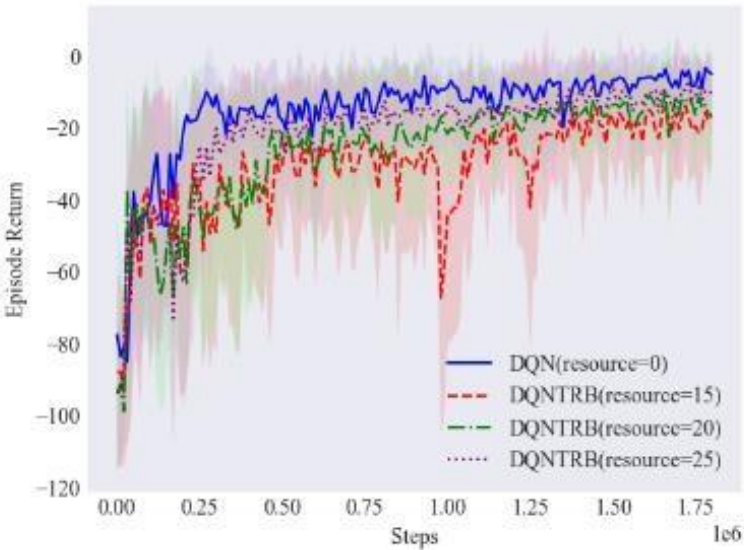


Рисунок 7 - Производительность обучения DQNTRB в задачах предотвращения столкновений при различных ограничениях ресурсов, вместе с производительностью DQN в задачах предотвращения столкновений без ограничений ресурсов.

Явно видно, что без ограничения ресурсов алгоритм DQN сходится быстрее всего и достигает лучших результатов. DQNTRB более чувствителен к изменениям уровня ресурсов. Конкретные параметры предотвращения столкновений в таблице 8 также отражают тот же результат, где алгоритм DQN демонстрирует оптимальную эффективность предотвращения столкновений и уровень успешного согласования высоты. Однако в отношении частоты предупреждения DQNTRB значительно уступает DQN.

Таблица

8 Производительность dqntb в задачах предотвращения столкновений с различными ограничениями ресурсов и dqn в задачах предотвращения столкновений без ограничений ресурсов

Алгоритм	Частота столкновений	Уровень успеха координации	Количество оповещений	Количество укреплений оповещений	Количество изменений оповещений	Количество пересечений высоты
<b>DQN(ресурс=0)</b>	0.021	0.900	6.140	2.011	7.190	0.082
<b>DQNTRB (ресурс=15)</b>	0.082	0.662	1.146	0.682	4.307	0.134
<b>DQNTRB (ресурс=20)</b>	0.056	0.717	1.856	0.729	3.436	0.123
<b>DQNTRB (ресурс=25)</b>	0.041	0.799	2.247	0.889	4.332	0.102

На рисунке 8 показана производительность обучения SACTRB при различных ограничениях ресурсов и производительность алгоритма SAC без ограничения ресурсов в среде предотвращения столкновений. Аналогично рисунку 7, алгоритм SAC демонстрирует хорошие результаты сходимости; однако скорости сходимости различных алгоритмов схожи, и кривые обучения указывают на то, что SACTRB менее чувствителен к изменениям уровня ресурсов.



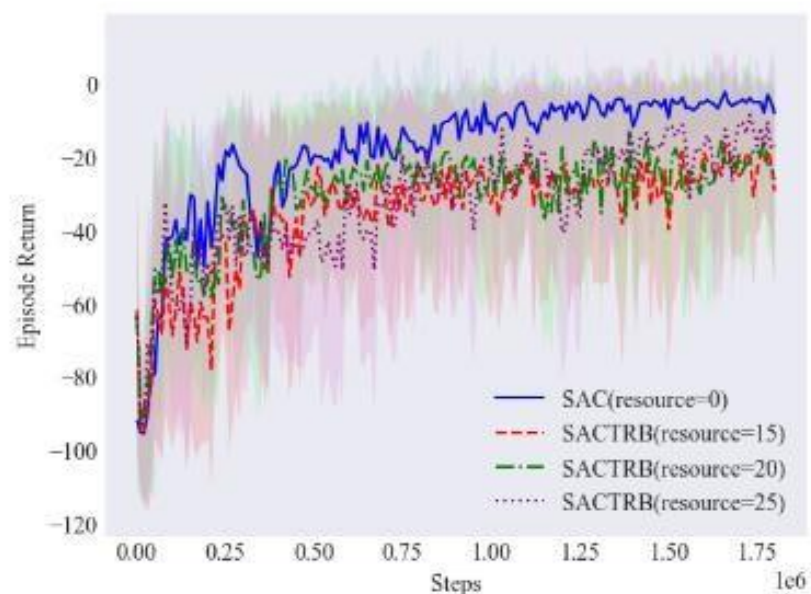


Рисунок 8 - Производительность обучения SACTRB в задачах предотвращения столкновений при различных ограничениях ресурсов, вместе с производительностью SAC в задачах предотвращения столкновений без ограничений ресурсов.

Конкретные параметры предотвращения столкновений в таблице 9 показывают, что при ограничениях ресурсов, в отличие от DQNTRB, алгоритм достигает наилучшей частоты предотвращения столкновений и уровня успешного согласования высоты при количестве ресурсов 20, при этом поддерживая аналогичную эффективность предотвращения столкновений и значительно более низкую частоту предупреждений по сравнению с SAC.



Алгоритм	Частота столкновений	Уровень успеха координации	Количество оповещений	Количество укреплений оповещений	Количество изменений оповещений	Количество пересечений высоты
SAC(ресурс=0)	0.045	0.831	4.578	2.734	7.581	0.111
SACTRB (ресурс=15)	0.111	0.603	0.968	0.766	4.248	0.128
SACTRB (ресурс=20)	0.078	0.657	1.565	0.889	5.352	0.114
SACTRB (ресурс=25)	0.098	0.649	1.621	1.063	6.929	0.181

## Заключение

В данной статье демонстрируется эффективность включения ограничений ресурсов в алгоритмы обучения с подкреплением для предотвращения столкновений в воздухе. Формулируя задачу как задачу МДП с ограничениями ресурсов и вводя механизм бонуса за время ресурсов (TRB), мы успешно оптимизировали как алгоритм DQN, так и SAC, в результате чего были получены DQNTRB и SACTRB.

Эти модифицированные алгоритмы достигли значительного сокращения ненужных предупреждений без ущерба для производительности предотвращения столкновений. Результаты подчеркивают потенциал обучения с подкреплением с учетом ресурсов для повышения эффективности и надежности систем предотвращения столкновений в условиях все более сложной воздушной среды.

## Конфликт интересов

Авторы заявляют об отсутствии конфликта интересов.

## Conflict of interest

The authors declare no conflict of interest.

## Список источников

1. Holland J.E., Kochenderfer M.J., Olson W.A. Optimizing the next generation collision avoidance system for safe, suitable, and acceptable operational performance // Air Traffic Control Quarterly. 2013. Vol. 21, no. 3. P. 275–297.
2. De D., Sahu P.K. A survey on current and next generation aircraft collision avoidance system // International Journal of Systems, Control and Communications. 2018. Vol. 9, iss. 4. P. 306–337.
3. Kochenderfer M.J., Holland J.E., Chryssanthacopoulos J.P. Next generation airborne collision avoidance system // Lincoln Laboratory Journal. 2012. Vol. 19, iss. 1. P. 17–33.
4. Kochenderfer M.J., Chryssanthacopoulos J.P. Robust airborne collision avoidance through dynamic programming : technical report / Massachusetts Institute of Technology, Lincoln Laboratory. 2011, 130 p. Project report ATC-371.
5. Optimized airborne collision avoidance / M.J. Kochenderfer, C. Amato, G. Chowdhary et al. // Decision making under uncertainty: theory and application. MIT Press, 2015. P. 249–276.
6. Julian K.D., Kochenderfer M.J., Owen M.P. Deep neural network compression for aircraft collision avoidance systems // Journal of Guidance, Control, and Dynamics. 2019. Vol. 42, iss. 3. P. 598–608.
7. Julian K.D., Kochenderfer M. Guaranteeing safety for neural network-based aircraft collision avoidance systems // 2019 IEEE/AIAA 38th Digital Avionics Systems Conference (DASC). DOI 10.1109/DASC43569.2019.9081748.
8. Li S, Egorov M., Kochenderfer M. Optimizing collision avoidance in dense airspace using deep reinforcement learning // Thirteenth USA/Europe Air Traffic Management Research and Development Seminar (ATM2019). 2019. DOI 10.48550/arxiv.1912.10146.
9. Online multiple-aircraft collision avoidance method / P. Zhao, W. Wang, L.Ying et al. // Journal of Guidance, Control, and Dynamics. 2020. Vol. 43, iss. 2. P. 1–17.

DOI 10.2514/1.G005161.

10. A partially observable multi-ship collision avoidance decision-making model based on deep reinforcement learning // K. Zheng, X. Zhang, C. Wang et al. // Ocean & Coastal Management. 2023. Vol. 242. Art. 106689.
11. Kormushev P., Calinon S., Caldwell D.G. Reinforcement learning in robotics: Applications and real-world challenges // Robotics. 2013. Vol. 2, iss. 3. P. 122–148.
12. An aircraft collision avoidance method based on deep reinforcement learning / Z. Liu, E. Neretin, X. Gao, et al. // 9th International Conference on Control and Robotics Engineering (ICCRE), IEEE. 2024. P. 241–246.
13. Bhatia A., Varakantham P., Kumar A. Resource constrained deep reinforcement learning // Proceedings of the international conference on automated planning and scheduling. 2019. Vol. 29. P. 610–620.
14. Efficient exploration in resource-restricted reinforcement learning / Z. Wang, T. Pan, Q. Zhou et al. // Proceedings of the AAAI Conference on Artificial Intelligence. 2023 Vol. 37, no. 8. P. 10279–10287.
15. Sutton R.S., Barto A.G. Reinforcement learning: An introduction. Cambridge : MIT press, 2018.
16. Playing atari with deep reinforcement learning / V. Mnih, K. Kavukcuoglu, D. Silver, et al. // ArXiv.org : website / arXiv preprint arXiv:1312.5602: 2013. 9 p.
17. Human-level control through deep reinforcement learning / V. Mnih, K. Kavukcuoglu, D. Silver et al. // Nature. 2015. Vol. 518(7540). P. 529–533.
18. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor / T. Haarnoja, A. Zhou, P. Abbeel et al. // International conference on machine learning. PMLR, 2018. P. 1861–1870.
19. Mohamed S., Jimenez Rezende D. Variational information maximisation for intrinsically motivated reinforcement learning // Advances in neural information processing systems. 2015. 28 p.

## References

1. Holland J.E., Kochenderfer M.J., Olson W.A. Optimizing the next generation collision avoidance system for safe, suitable, and acceptable operational performance // Air Traffic Control Quarterly. 2013. Vol. 21, no. 3. P. 275–297.
2. De D., Sahu P.K. A survey on current and next generation aircraft collision avoidance system // International Journal of Systems, Control and Communications. 2018. Vol. 9, iss. 4. P. 306–337.
3. Kochenderfer M.J., Holland J.E., Chryssanthacopoulos J.P. Next generation airborne collision avoidance system // Lincoln Laboratory Journal. 2012. Vol. 19, iss. 1. P. 17–33.
4. Kochenderfer M.J., Chryssanthacopoulos J.P. Robust airborne collision avoidance through dynamic programming : technical report / Massachusetts Institute of Technology, Lincoln Laboratory. 2011, 130 p. Project report ATC-371.
5. Optimized airborne collision avoidance / M.J. Kochenderfer, C. Amato, G. Chowdhary et al. // Decision making under uncertainty: theory and application. MIT Press, 2015. P. 249–276.
6. Julian K.D., Kochenderfer M.J., Owen M.P. Deep neural network compression for aircraft collision avoidance systems // Journal of Guidance, Control, and Dynamics. 2019. Vol. 42, iss. 3. P. 598–608.
7. Julian K.D., Kochenderfer M. Guaranteeing safety for neural network-based aircraft collision avoidance systems // 2019 IEEE/AIAA 38th Digital Avionics Systems Conference (DASC). DOI 10.1109/DASC43569.2019.9081748.
8. Li S, Egorov M., Kochenderfer M. Optimizing collision avoidance in dense airspace using deep reinforcement learning // Thirteenth USA/Europe Air Traffic Management Research and Development Seminar (ATM2019). 2019. DOI 10.48550/arxiv.1912.10146.
9. Online multiple-aircraft collision avoidance method / P. Zhao, W. Wang, L.Ying et al. // Journal of Guidance, Control, and Dynamics. 2020. Vol. 43, iss. 2. P. 1–17. DOI 10.2514/1.G005161.

10. A partially observable multi-ship collision avoidance decision-making model based on deep reinforcement learning // K. Zheng, X. Zhang, C. Wang et al. // Ocean & Coastal Management. 2023. Vol. 242. Art. 106689.
11. Kormushev P., Calinon S., Caldwell D.G. Reinforcement learning in robotics: Applications and real-world challenges // Robotics. 2013. Vol. 2, iss. 3. P. 122–148.
12. An aircraft collision avoidance method based on deep reinforcement learning / Z. Liu, E. Neretin, X. Gao, et al. // 9th International Conference on Control and Robotics Engineering (ICCRE), IEEE. 2024. P. 241–246.
13. Bhatia A., Varakantham P., Kumar A. Resource constrained deep reinforcement learning // Proceedings of the international conference on automated planning and scheduling. 2019. Vol. 29. P. 610–620.
14. Efficient exploration in resource-restricted reinforcement learning / Z. Wang, T. Pan, Q. Zhou et al. // Proceedings of the AAAI Conference on Artificial Intelligence. 2023 Vol. 37, no. 8. P. 10279–10287.
15. Sutton R.S., Barto A.G. Reinforcement learning: An introduction. Cambridge : MIT press, 2018.
16. Playing atari with deep reinforcement learning / V. Mnih, K. Kavukcuoglu, D. Silver, et al. // ArXiv.org : website / arXiv preprint arXiv:1312.5602: 2013. 9 p.
17. Human-level control through deep reinforcement learning / V. Mnih, K. Kavukcuoglu, D. Silver et al. // Nature. 2015. Vol. 518(7540). P. 529–533.
18. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor / T. Haarnoja, A. Zhou, P. Abbeel et al. // International conference on machine learning. PMLR, 2018. P. 1861–1870.
19. Mohamed S., Jimenez Rezende D. Variational information maximisation for intrinsically motivated reinforcement learning // Advances in neural information processing systems. 2015. 28 p.

