



Проблема CSM в частотной области изображения: новый подход к решению

Р. А. Солодуха^а, канд. техн. наук, доцент, /orcid.org/0000-0002-3878-4221, standartal@list.ru

^аВоронежский государственный университет инженерных технологий, Революции пр., 19, Воронеж, 394036, РФ

Введение: одной из проблем, препятствующих применению стеганоанализа в практике цифровой криминалистики, является несоответствие тестируемого контейнера множеству, на котором обучалась модель распознавания. Имеющиеся методы решения данной проблемы не учитывают специфику искажений контейнера, привносимых различными стеганоалгоритмами. При определенной локализации искажений возможно их использование для формирования обучающего множества, соответствующего тестируемому контейнеру. **Цель:** формализация и проверка гипотезы об эффективности формирования обучающего множества на основе упорядочивания расстояния между «калиброванными изображениями» по векторам стеганоаналитических признаков при локализации искажений в частотной области изображений; сравнение точности распознавания при предлагаемом подходе и объединении признаков «калиброванного изображения» и исходного в единый вектор признаков. **Результаты:** показана целесообразность использования стеганоаналитических векторов признаков «калиброванных изображений» для вычисления расстояния между контейнерами. Предложена формализованная процедура формирования обучающего множества на основе расстояния между контейнерами. Создана программная инфраструктура для проведения численного эксперимента. Экспериментально показано, что независимо от качества JPEG-сжатия использование вектора признаков «калиброванного изображения» для формирования обучающего множества эффективнее, чем включение в общий вектор признаков. При этом основные вычисления задействованы на попарный расчет расстояния между контейнерами и могут осуществляться до появления объекта исследования. Для обеспечения воспроизводимости эксперимента наборы данных и программный код представлены в Kaggle. **Практическая значимость:** на примере стеганоалгоритма nsF5 показано преимущество применения стеганоаналитического вектора PEV-274 на соответствующем тестируемому файлу обучающем множестве перед CC-PEV-548 на случайной выборке. Предложенный подход способствует как увеличению точности стеганоанализа, так и уменьшению сроков исследования, что важно в экспертной практике.

Ключевые слова – стеганоанализ, вектор признаков, nsF5, PEV-274, CC-PEV-548, Cover-Source Mismatch, стеганография, машинное обучение, регрессия, расстояние между векторами, экспертиза.

Для цитирования: Солодуха Р. А. Проблема CSM в частотной области изображения: новый подход к решению. *Информационно-управляющие системы*, 2026, № 1, с. 8–18. doi:10.31799/1684-8853-2026-1-8-18, EDN: QIENHD

For citation: Solodukha R. A. A novel approach to solving the CSM problem in the frequency domain of an image. *Informatsionno-upravliaiushchie sistemy* [Information and Control Systems], 2026, no. 1, pp. 8–18 (In Russian). doi:10.31799/1684-8853-2026-1-8-18, EDN: QIENHD

Введение

Тенденции последних конкурсов по стеганоанализу Alaska [1] и Alaska-2 [2] демонстрируют интерес исследователей к наборам данных, содержащих полноцветные изображения, полученные с различных устройств. Основное внимание уделяется точности обнаружения при минимальном количестве ложных срабатываний. Это, а также появление теоретических работ по стеганоаналитической экспертизе [3] свидетельствует о том, что выявление стеганографии в цифровых изображениях находится на пороге качественного скачка. Вероятно, в среднесрочной перспективе стеганоанализ перейдет из исследовательских лабораторий и СТФ в практическую сферу цифровой криминалистики.

За три десятилетия развития стеганоанализа разработаны десятки методов. Использование форматных, сигнатурных и статистических подходов варьируется по вычислительной слож-

ности, точности, надежности, требованиям к наличию дополнительной информации [4–7]. Де-факто стандартом стало применение многомерного вектора признаков с последующей классификацией или регрессией несмотря на то, что статистический стеганоанализ позволяет делать лишь вероятностные выводы [8, 9]. Следует отметить, что соответствующая ему экспертная методика должна включать оценку достоверности полученных результатов.

Рассмотрим ситуацию, когда на экспертизу поступает графический файл I с вопросами:

1. Имеется в представленных файлах/файле стегановложение, выполненное с помощью программы/алгоритма <наименование программы/алгоритма>?

2. Каков размер вложения?

Допустим, в распоряжении эксперта имеется заранее обученный регрессор $R: \mathbf{X} \rightarrow \hat{\mathbf{Y}}$ где

$\mathbf{X} = \{x_{i,j}\}_{i=1, j=1}^{i=n, j=m}$ – матрица реализаций объяс-

няющих переменных (вектора стеганоаналитических признаков), n — количество реализаций, m — количество объясняющих переменных (размер вектора признаков); $\hat{\mathbf{Y}} = \{\hat{y}_i\}_{i=1}^{j=1}$ — вектор-столбец прогнозных значений зависимой переменной.

При этом регрессор обучен на нужной стеганопрограмме/алгоритме, известны метрики качества регрессора. Также эксперту доступен стеганоаналитический алгоритм SA: $\mathbf{I} \rightarrow \mathbf{X}_1$, используемый для обучения, и он может получить вектор признаков.

Вопросы, на которые должен ответить эксперт перед началом исследования:

1. Возможно ли применение регрессора для данного файла?
2. Какова достоверность полученного результата?
3. В какой форме будет сделан вывод?

Первый вопрос относится к известной проблеме Cover-Source Mismatch (CSM), дословно — несоответствие источника контейнера. Изначально под этим понималось, что обучающие и проверяемые изображения получены разными устройствами (image acquisition), но затем понятие CSM распространилось на процессы преобразования (image processing, JPEG compression) и даже семантику изображения [10].

Проблема CSM имеет два диаметрально противоположных решения [11].

Атомистический подход. Если имеется возможность определить параметры изображения, то формируется гомогенный относительно подозрительного контейнера набор данных, т. е. происходит имитация источника подозрительного контейнера [12].

Холистический подход. Обучающее множество формируется из контейнеров, порожденных разнообразными источниками. Задача состоит в получении не столько точной модели, сколько обладающей обобщающей способностью [13].

Отдельно следует отметить адаптационный подход [14]. Идея адаптации заключается в обучении на источнике контейнеров, называемом исходным, и использовании полученных знаний для адаптации к неизвестному источнику контейнеров, называемому целевым. Этот метод позволяет детектору находить пространство, инвариантное к характеристикам, в котором распределения характеристик контейнеров исходного и целевого источников близки.

Настоящая статья посвящена реализации атомистического подхода для частного случая — стеганографии, локализованной в частотной области изображения. В широком смысле статья имеет отношение к формированию обучающего

множества при машинном обучении, чему посвящены работы [15–17].

Идея статьи частично перекликается с высказанной в [18], где предложено перед финальным обнаружением добавить этап предварительной фильтрации контейнеров (отбор «хороших» контейнеров, в которых наличие/отсутствие внедренной информации может быть определено более достоверно, чем во всем множестве).

Статья соответствует направлению по формированию и проверке эффективности векторов признаков с возможностью управления соотношением точность/ресурсоемкость [19–21].

Обоснование идеи исследования

Атомистический подход предполагает формирование обучающей и тестовой выборок из одного множества. Это означает, что характеристики элементов этого множества должны лежать в определенных границах. Под характеристиками в широком смысле можно понимать характеристики трех основных сущностей, участвующих в системе формирования цифрового изображения (DIC): сцены (S), устройства (D), процесса преобразования (P):

$$DIC: S \times D \times P \rightarrow \mathbf{I}.$$

В узком смысле характеристики (CH) — это производные, полученные (DER) от непосредственно изображения \mathbf{I} :

$$DER: \mathbf{I} \rightarrow \mathbf{CH}_1.$$

Значительное количество работ посвящено поиску оптимального набора характеристик изображения, с помощью которых можно отнести изображения к одному множеству [22–24]. Однако данный подход имеет два недостатка:

- 1) реальные изображения, как правило, значительных размеров. Кумулятивные характеристики всего изображения могут не соответствовать характеристикам сегментов изображения;
- 2) априори неизвестно, содержит ли файл, представленный на исследование, вложение. Если содержит, то его характеристики будут искажены. Для нивелирования влияния возможной модификации файлы подвергаются субдискретизации (downsampling) [25] или калибровке [26].

Идея настоящего исследования в следующем. Поскольку решающими характеристиками (функциями от изображения) при стеганоанализе являются реализации стеганоаналитических алгоритмов, то целесообразно сравнивать файлы на принадлежность к одному множеству именно через них, т. е. $\mathbf{CH}_1 = \mathbf{X}_1$.

Однако именно эти характеристики наиболее чувствительны к стегановложению. В условиях априорной неизвестности наличия/размера вложения в исследуемом файле сравнение возможно только при устранении искажений, вызванных вложением. Поскольку это невозможно, остается «пожертвовать» составляющими изображения, где локализованы искажения. Для пространственных областей изображений это возможно, например, для алгоритмов семейства Least Significant Bit Replacement, где искажения затрагивают лишь плоскость младшего бита, и ее можно заполнить нулями, единицами или чередованием нулей и единиц [27]. Однако уже к LSB Matching предлагаемый подход неприменим.

Для частотной области предложено [26] использовать калибровку. «Калиброванное изображение» — изображение, очищенное от стеганографического искажения, но сохранившее семантику. «Калиброванное JPEG-изображение» получается следующим образом. Изображение разворачивается из частотного в пространственное представление, обрезается на несколько пикселей по обоим направлениям, опять сжимается в JPEG с прежними параметрами. «Калиброванное изображение» сохраняет свойства исходного на макроуровне.

Итак, предлагается решить проблему CSM путем формирования обучающего множества для анализируемого изображения из изображений с близкими значениями векторов признаков калиброванных версий.

Одним из вариантов формирования множества из близких векторов является кластеризация. Однако в практической плоскости возникает неопределенность относительно количества как кластеров, так и элементов в кластере. Исследуемый файл может попасть в кластер, мощность которого недостаточна для обучения. В этой связи кластеризация может быть использована для предварительной оценки потенциальной эффективности предложенного подхода и качества выборки.

Второй вариант — определение принадлежности изображений к одному множеству путем непосредственного расчета расстояния между векторами признаков $\text{dist}(\mathbf{X}_I, \mathbf{X}_J)$. При этом метрика должна выбираться из соображений как результативности, так и скорости вычислений.

В статье предложен подход, направленный на решение проблемы CSM для частного случая: при локализации искажений в коэффициентах дискретного косинусного преобразования (ДКП; Discrete Cosine Transform, DCT) и известном алгоритме стегановложения. Приведено формальное описание процесса анализа, архитектура, состав, технологический стек стенда, описание и

результаты численного эксперимента по проверке эффективности предложенного подхода.

Формализация идеи исследования

Рассмотрим частный случай стеганоанализа. На исследование представлено/представлены изображение/изображения (SI — Suspicious Image/Images), известен стеганоалгоритм (СГА) или стеганопрограмма (СГП), с помощью которых могли быть выполнены вложения, искажения привносятся в частотную область изображения (коэффициенты дискретного косинусного преобразования JPEG). Для ответа на вопрос о размере вложения предлагается процедура (схематично изображена на рис. 1), включающая в себя следующие этапы:

1. Формирование обучающего множества (коллекции) изображений (IS — Image Set). Изображения можно получать самостоятельно фотокамерой, загружать из фотохостингов, социальных сетей и пр. В дальнейшем предполагается, что файлы изображений в необходимом количестве и формате имеются в распоряжении аналитика.

2. Формирование множества «калиброванных изображений». Калибровка как изображений IS (ClbrIS — Calibrated IS), так и изображений, переданных на исследование (ClbrSI — Calibrated SI).

3. Формирование множества пустых и заполненных изображений (EPIS — Empty and Payloaded IS). Применение СГА/СГП к коллекции изображений IS, реализация вложений с определенным шагом. Шаг определяется аналитиком исходя из требуемой точности, наличия вычислительных ресурсов и машинной памяти.

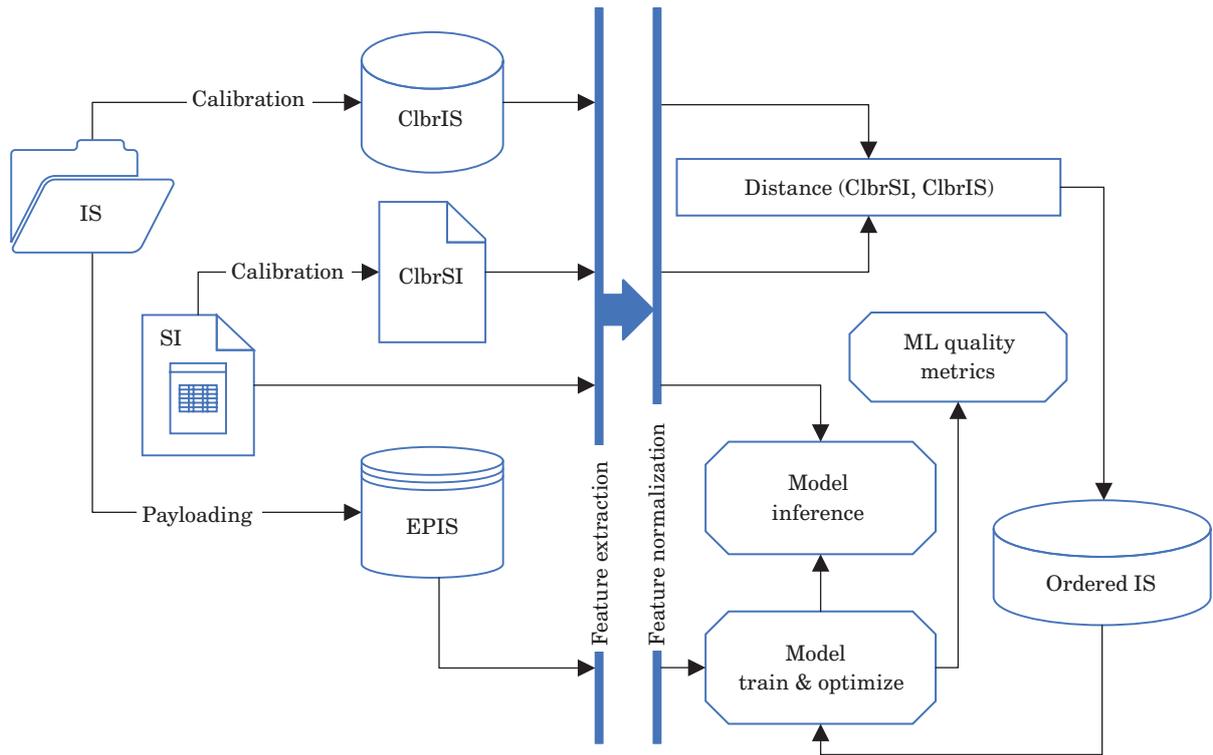
4. Формирование набора данных. Извлечение стеганоаналитических признаков (Feature extraction) из элементов EPIS и SI, нормализация (Feature normalization).

5. Формирование матрицы расстояний — Distance (ClbrSI, ClbrIS). Расчет расстояния между векторами признаков «калиброванных изображений» по выбранной метрике.

6. Упорядочивание изображений по минимуму расстояния от переданного на исследование изображения (Ordered IS). Определение размера набора данных (по мощности, по порогу расстояния, по максимуму точности распознавания).

7. Выбор модели машинного обучения и метрик качества (ML quality metrics).

8. Обучение модели на полученном наборе данных (Model train & optimize), распознавание переданного на исследование изображения (Model inference).



■ **Рис. 1.** Предложенная процедура стеганоанализа
 ■ **Fig. 1.** Scheme of proposed steganalytic technique

Опишем предложенную последовательность действий. Пусть:

- SI (Suspicious Image) – исследуемое изображение;
- IS (Image Set) – коллекция исходных изображений $\{IS_1, IS_2, \dots, IS_N\}$;
- СГА/СГП – $S: S(\mathbf{I}, m, p) \rightarrow \mathbf{I}_{\text{payloaded}}$, где \mathbf{I} – изображение, m – сообщение, p – параметры вложения (опционально);
- ClbrIS = $\{\text{Calibrate}(IS_i) \mid IS_i \in IS\}$ – калиброванные файлы из IS;
- ClbrSI = $\text{Calibrate}(SI)$ – калиброванное исследуемое изображение, калибровка выполняется в соответствии с [26];
- F – функция получения вектора стеганоаналитических признаков;
- $\text{dist}(\mathbf{F}_1, \mathbf{F}_2)$ – метрика расстояния между векторами;
- Q – метрика качества.

Тогда EPIS = $\{IS \cup IS_{\text{payloaded}}\}$ – множество пустых и заполненных изображений (Empty and Payloaded IS), где $IS_{\text{payloaded}} = \{S(IS_i, m_j, p_k) \mid IS_i \in IS, m_j \in M, p_k \in P\}$, M – множество стеганосообщений (файлов), P – набор параметров вложения.

Формируется набор данных относительно SI, для чего необходимо:

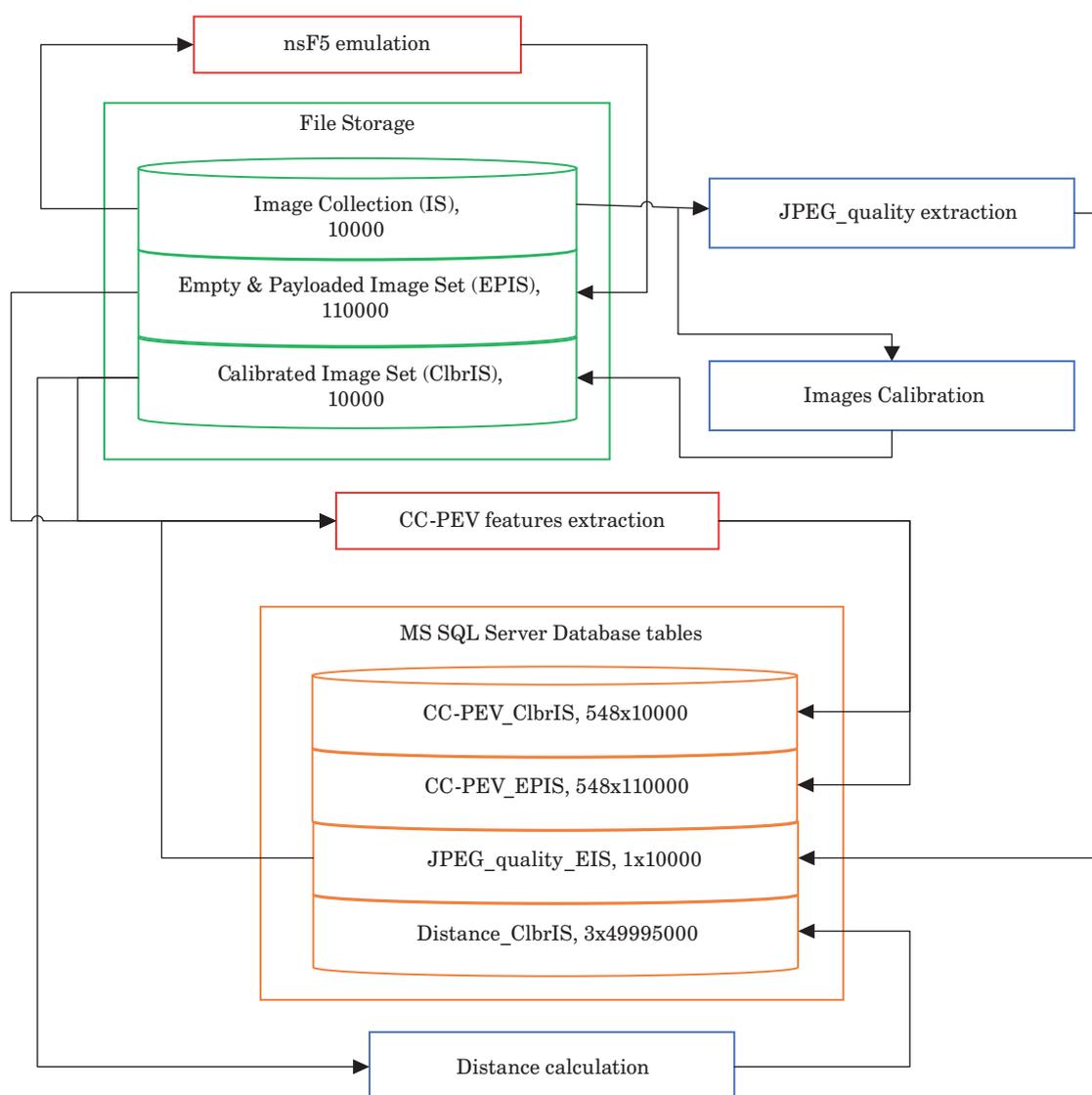
- вычислить векторы признаков $\mathbf{F}_{SI} = F(\text{Clbr}_{SI}), \mathbf{F}_i = F(\text{Clbr}_{IS}_i)$;

- найти вектор расстояний от исследуемого изображения $\mathbf{D} = \{\text{dist}(\mathbf{F}_{SI}, \mathbf{F}_i)\}$;
- упорядочить $\mathbf{D} = (D_i)$ по возрастанию: $D_{\sigma(1)} \leq D_{\sigma(2)} \leq \dots \leq D_{\sigma(N)}$, где $\sigma = \text{argsort}(\mathbf{D})$ – перестановка множества индексов $\text{Ind}(IS)$;
- выбрать подмножество из k наиболее близких к SI изображений $IS_{DS} = \{IS_{\sigma(1)}, IS_{\sigma(2)}, \dots, IS_{\sigma(k)}\}$, где k определяется:
 - фиксированным порогом: $k = \text{const}$;
 - процентилем: $k = \lceil \alpha N \rceil, \alpha \in (0, 1)$;
 - пороговым расстоянием между изображениями: $k = \max\{i \mid D_{\sigma(i)} \leq \varepsilon\}$;
 - максимумом качества распознавания: $k = \max\{i \mid Q(\text{EPIS}_{DS} \mid k = i) > Q(\text{EPIS}_{DS} \mid k = i - 1)\}$, где $\text{EPIS}_{DS} = \{\text{EPIS}_i \mid i \in \text{Ind}(IS_{DS})\}$;
- найти функцию регрессии с обучением по EPIS_{DS}:

$$\hat{f} = \arg \min_{f \in H} \frac{1}{|\text{EPIS}_{DS}|} \times \sum_{i=1}^{|\text{EPIS}_{DS}|} L[f(F(\text{EPIS}_{DS}^i), \text{Pld}(\text{EPIS}_{DS}^i))],$$

где H – множество регрессионных функций; L – функция потерь; $\text{Pld}(\text{EPIS}_{DS}^i)$ – размер стегановложения;

- рассчитать метрики качества машинного обучения;



■ **Рис. 2.** Схема формирования данных для вычислительного эксперимента
 ■ **Fig. 2.** Data processing experiment scheme

– предсказать размер вложения в SI:
 $Pld(SI) = f(F(SI))$.

Экспериментальная часть

Эксперимент по определению эффективности предложенного подхода предполагает не просто вычисление метрик оценки регрессии, но и сравнение с иным подходом [28] к использованию калибровки.

Предложено [26] в состав вектора признака включать признаки как контейнера, так и его калиброванной версии. Эксперимент, проведенный на векторе признаков DCT-23 (признаки извлекаются из коэффициентов ДКП с использованием гистограмм, двумерных гисто-

грамм, матриц совместного появления и иных функционалов), показал эффективность данного подхода.

Указанный подход получил развитие в работе [28], где сформирован вектор признаков CC-PEV-548, состоящий из набора признаков PEV-274, вычисленного по контейнеру и его калиброванной версии (Cartesian Calibration, CC): $CC-PEV-548 = \{CC-PEV-274 \cup PEV-274\}$. В свою очередь PEV-274 [29] представляет комбинацию наборов расширенного DCT (Extended DCT-193) и Markov-81 (разности абсолютных значений коэффициентов ДКП по направлениям агрегированы в матрицы переходных вероятностей марковского процесса 1-го порядка с дальнейшим усреднением по направлениям): $PEV-274 = \{ExtDCT-193 \cup Markov-81\}$.

В рамках эксперимента осуществляется сравнение точности определения размера стегановложения вектором признаков CC-PEV-548 по [28], а также по предложенной процедуре с PEV-274.

Цель эксперимента – проверить эффективность предложенного подхода через наличие/отсутствие эффекта от переноса части распознавательной способности CC-PEV-548 на устранение влияния CSM.

Стенд для проведения эксперимента: Intel i5-12400 2,5 GHz, SSD 500 GB, RAM 32 GB под управлением Windows 10 Pro с установленным программным обеспечением: Python 3, Visual Studio Code, MATLAB R2021, MS SQL Server 2019, MSSS Management Studio. Для данной конфигурации время на обучение по 300 контейнерам и прогноз составляет ~8 с, по 50 контейнерам ~4 с. Время на вычисление попарных расстояний 10 000 контейнеров составило ~14 сут.

Конкретизация моделей, алгоритмов, данных и инструментов, использованных в эксперименте (наборы данных и скрипты доступны в Kaggle: <https://www.kaggle.com/datasets/romansolodukha/clbr-jpeg>):

- коллекция изображений (IS) – первые 10 000 файлов из набора Alaska-2 (<https://www.kaggle.com/c/alaska2-image-steganalysis/data>): файлы JPEG с качеством (QF – Quality Factor) 95, 90, 75 в соотношении 3278/3311/3411. Для определения качества изображения использована функция `get_jpg_quality` (<https://gist.github.com/eddy-geek/c0f01dc5401dc50a49a0a821cdc9b3e8/versions>);

- стеганоалгоритм (S) – использован nsF5 (<https://dde.binghamton.edu/download/nsf5simulator>), так как на нем PEV-274 показал лучшие результаты в рамках конкурса BOSS [30];

- шаг вложения – 10 % от максимально возможного (9, 19, 29...99 %), |EPIS| = 110 000;

- стеганоаналитический вектор признаков (F) – PEV-274 (https://dde.binghamton.edu/download/feature_extractors), данные нормализованы функцией `MinMaxScaler` из библиотеки `sklearn.preprocessing`;

- метрика близости (dist) – расчет осуществлен со следующими метриками из библиотеки `scipy.spatial.distance`: `correlation` (корреляционное расстояние), `euclidean` (евклидово расстояние), `braycurtis` (расстояние Брея – Кертиса), `cosine` (косинусное расстояние). Следует отметить, что значимых различий в результатах при использовании вышеуказанных метрик не наблюдалось. Экспериментальные данные приведены для евклидова расстояния;

- метрика качества (Q) – RMSE и коэффициент детерминации (R2). Используются функции `r2_score`, `mean_squared_error` из библиотеки `sklearn.metrics`;

- регрессионная модель (f) – Ridge. Выбор регрессионной модели осуществлен на основании результатов применения AutoML на базе библиотеки `lazypredict` (<https://lazypredict.readthedocs.io>) к EPIS. Ridge оказался самым быстрым регрессором с приемлемым $R2 \approx 0,89$;

- программный стек:

- MS SQL Server 2019 – база данных с таблицами данных (размерность без учета системных атрибутов приведена на рис. 2), вспомогательные представления и функции;

- MATLAB 2021 – формирование заполненных контейнеров, извлечение PEV-274, CC-PEV-548;

- Python 3.11 + Visual Studio Code – калибровка, определение качества JPEG, вычисление расстояния, формирование обучающего и тестового множеств, обучение, применение модели, вычисление метрик качества.

В эксперименте множество «калиброванных изображений» упорядочивалось по расстоянию относительно каждого файла из тестового множества. Эксперимент проведен как с учетом QF JPEG, так и без учета.

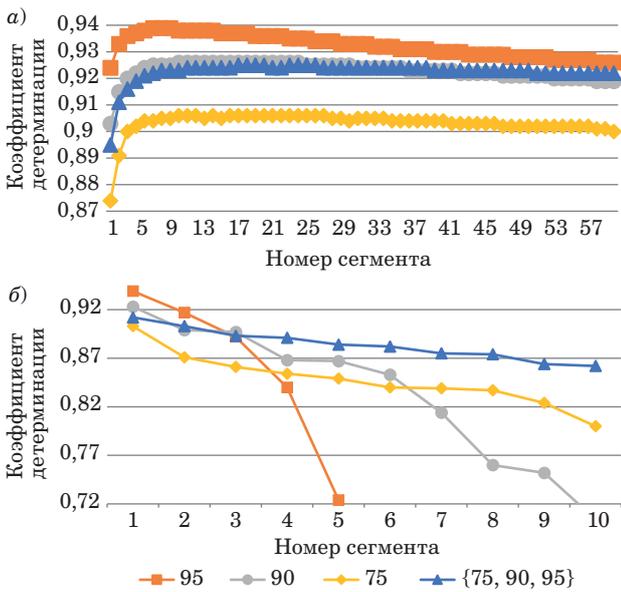
Сначала был определен минимальный размер обучающей выборки, при котором качество регрессии максимально. Для этого обучающее множество из 3000 контейнеров сегментировано по 50 контейнеров (в порядке увеличения dist – «от близких к дальним»), что с учетом вложенный составляет 550 контейнеров (1-й сегмент – 50 контейнеров, наиболее близких к тестируемому файлу). Тестовая выборка – 1000 случайных исходных контейнеров, с учетом вложений – 11 000.

Относительно метрик качества следует отметить, что во всех экспериментах RMSE и R2 коррелированы с коэффициентом в диапазоне (–0,95...–0,98). В этой связи принято решение описывать результаты экспериментов только по R2. Для представления точности прогноза приведем некоторые соответствия RMSE (в процентах от максимально возможного размера стегановложения) и R2: (14; 0,7), (11,6; 0,8), (9,1; 0,9), (7,6; 0,94).

На рис. 3, а приведено усредненное значение R2 тестового множества в зависимости от сегмента обучающего множества с накоплением. Точка N на оси абсцисс означает, что обучение осуществлено на контейнерах, составляющих сегменты $n \leq N$.

Анализ графиков рис. 3, а показывает, что R2 достигает максимума в зависимости от QF при $|IS_{DS}| = 300-500$, $|EPIS_{DS}| = 3300$, затем плавно уменьшается.

Дальнейшие эксперименты проводились с сегментами по 300 контейнеров (нижняя граница оптимального размера сегмента выбрана



■ **Рис. 3.** Зависимость усредненного коэффициента детерминации от сегмента обучающей выборки с накоплением (а) и без накопления (б) (обучающее множество упорядочено по возрастанию расстояния)

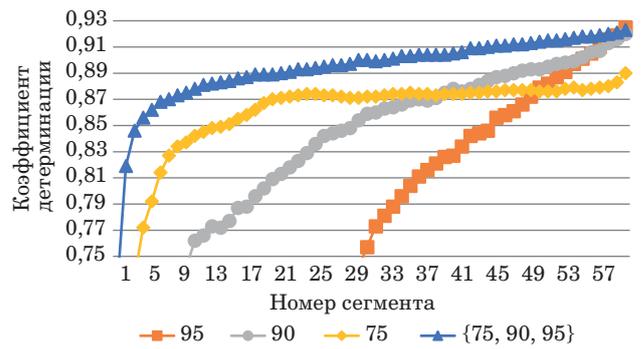
■ **Fig. 3.** Dependence of the averaged determination coefficient on the training sample segment with accumulation (a) and without accumulation (b) (the training set is ordered by increasing distance)

из соображений минимизации времени на вычисления). На графике рис. 3, б приведен R2 в зависимости от сегмента (больше номер сегмента — дальше от тестового изображения). Наблюдается влияние качества изображения на скорость убывания качества распознавания. Чем выше качество JPEG, тем сильнее влияние сегментации.

Интерес вызывает поведение кривой, описывающей вариант с изображениями произвольного качества, — $Q_{\{75,90,95\}}$. Ожидалось, что эти значения должны являться усреднением значений Q_{75} , Q_{90} , Q_{95} при соответствующих сегментах (аналогично кривым рис. 3), однако, начиная с 3-го сегмента: $Q_{\{75,90,95\}} > (Q_{75} + Q_{90} + Q_{95})/3$. Это означает, что расстояние между «калиброванными изображениями» влияет на R2 сильнее, чем идентичность QF элементов обучающего и тестового множества.

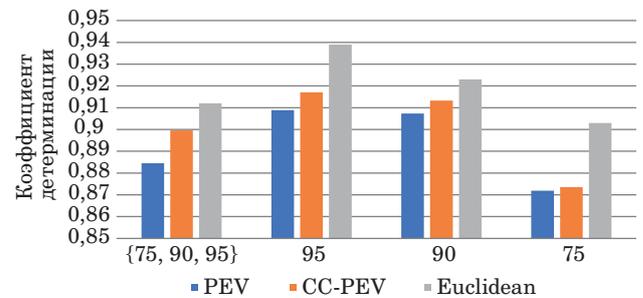
Для понимания данного явления построен график рис. 4, аналогичный рис. 3, а, но IS_{DS} упорядочено «от дальних к близким». Кривая $Q_{\{75,90,95\}}$ находится выше Q_{75} , Q_{90} , Q_{95} вплоть до сегмента с самыми «близкими» контейнерами.

На качественном уровне это можно объяснить тем, что в сегментах, где расстояние между контейнерами мало, идентичность QF положительно влияет на объясняющую способность модели. По мере увеличения расстояния между тестируемым контейнером и обучающим множеством отфильтрованные по QF выборки проигрывают



■ **Рис. 4.** Зависимость усредненного коэффициента детерминации от сегмента обучающей выборки с накоплением (обучающее множество упорядочено по убыванию расстояния)

■ **Fig. 4.** Dependence of the averaged determination coefficient on the training sample segment with accumulation (the training set is ordered by decreasing distance)



■ **Рис. 5.** Эффективность распознавания в зависимости от подхода

■ **Fig. 5.** Prediction accuracy depends on the approach

смешанной. Это связано с тем, что на первый план выходит именно расстояние между контейнерами.

Обобщить проведенный эксперимент можно гистограммой рис. 5, где представлена сводная информация по сравнению определения размера вложения с помощью PEV-274, CC-PEV-548 и предложенного подхода в комбинации с PEV-274 (на гистограмме PEV, CC-PEV, Euclidean соответственно). Для Euclidean использован сегмент из 300 самых «близких» контейнеров. Для PEV-274, CC-PEV-548 приведено усреднение по 10 наборам данных из 300 случайных контейнеров.

Для контейнеров с $QF = \{90, 95\}$ $R2_{PEV} \approx 0,91$, при $QF = 75$ наблюдается уменьшение коэффициента детерминации $R2_{PEV} \approx 0,87$. Использование CC-PEV-548 незначительно улучшает результаты для $QF = \{90, 95\}$, при этом для $QF = 75$ $R2_{PEV} \approx R2_{CC-PEV}$.

Применение предложенного подхода увеличило точность прогноза во всех группах контейнеров, что экспериментально подтвердило эффективность использования признаков, из-

влеченных из «калиброванных изображений», не в векторе признаков, а для определения расстояния между контейнерами. Другими словами, перенос информации о «калиброванном изображении» из признакового пространства в процедуру формирования релевантной обучающей выборки является более эффективным способом борьбы с CSM.

Заключение

В работе предложен, формализован и экспериментально проверен один из вариантов атомистического подхода к решению проблемы CSM в частотной области изображений. Обучающее множество формируется из контейнеров, наиболее «близких» к исследуемому файлу по расстоянию между векторами признаков их «калиброванных» версий.

Эксперимент с алгоритмом nsF5 показал, что предложенный подход обеспечивает более точную оценку размера стегановложения по сравнению с использованием как базового вектора PEV-274, так и расширенного SS-PEV-548 на случайной выборке. Установлено, что расстояние между векторами признаков по влиянию на точность прогноза преобладает над идентичностью качества JPEG элементов выборки. Предложенная процедура допускает вынесение наиболее ресурсоемкого этапа (расчет попарных расстояний) на этап предварительной под-

готовки, что потенциально содействует уменьшению сроков анализа.

Ограничение применения полученных результатов обусловлено частным характером валидации, так как эксперимент проведен на одном стеганоалгоритме и одном наборе данных. Обобщение результатов на другие алгоритмы и наборы изображений требует дополнительных исследований.

Таким образом, в перспективе планируется изучение границ и методики применения предложенного подхода с определением статистической достоверности результатов для каждого размера вложения. Выбор регрессионной модели обусловлен, в том числе, алгоритмической простотой и теоретической проработанностью оценки статистической значимости параметров линейной регрессии. В комбинации с оценкой достоверности приведенный в статье подход можно будет рассматривать в качестве прообраза экспертной методики.

В инфраструктурном плане для хранения векторов целесообразно осуществить миграцию с реляционной системы управления базами данных на векторную. Это связано не только с быстрой скоростью, но и с тем, что реляционные базы данных имеют ограничения по количеству столбцов в таблицах, что приводит к невозможности нативного (файлы в строках, признаки в столбцах) представления в них многомерных стеганоаналитических векторов признаков, например Milvus или Qdrant.

Литература

1. **Cogranne R., Giboulot Q., Bas P.** ALASKA#2: Challenging academic research on steganalysis with realistic images. *IEEE International Workshop on Information Forensics and Security (WIFS)*, 2020, pp. 1–5. doi:10.1109/WIFS49906.2020.9360896
2. **Cogranne R., Giboulot Q., Bas P.** The ALASKA steganalysis challenge: A first step towards steganalysis. *Proceedings of the ACM Workshop on Information Hiding and Multimedia Security*, 2019. doi:10.1145/3335203.3335726
3. **Bobok I. I., Koboziyeva A. A.** Theoretical foundations of digital content integrity expertise. *Problemele Energeticii Regionale*, 2025, no. 1 (65), pp. 105–121. doi:10.52254/1857-0070.2025.1-65.08
4. **Солодуха Р. А., Атласов И. В., Кубасов И. А.** *Стеганализ цифровых изображений: технологии, алгоритмы, программная реализация: монография.* Воронеж, Воронежский институт МВД России, 2022. 172 с.
5. **Michaylov K., Sarmah D.** Steganography and steganalysis for digital image enhanced forensic analysis and recommendations. *Journal of Cyber Security Technology*, 2025, vol. 9(1), pp. 1–27. doi:10.1080/23742917.2024.2304441
6. **Сирота А. А., Дрюченко М. А., Иванков А. Ю.** Стеганализ цифровых изображений с использованием методов поверхностного и глубокого машинного обучения: известные подходы и новые решения. *Вестник ВГУ. Серия: Системный анализ и информационные технологии*, 2021, № 1, с. 33–52. doi:10.17308/sait.2021.1/3369
7. **Полунин А. А., Яндашевская Э. А.** Использование аппарата сверточных нейронных сетей для стеганализа цифровых изображений. *Труды ИСП РАН*, 2020, т. 32, № 4, с. 155–163. doi:10.15514/ISPRAS-2020-32(4)-11
8. **Лубин А. Ф.** О допустимости вероятностных выводов экспертного заключения в уголовном судопроизводстве. *Юридическая наука и практика: Вестник Нижегородской академии МВД России*, 2019, т. 3, № 47, с. 138–142. doi:10.36511/2078-5356-2019-3-138-142
9. **Овсянников И. В.** К вопросу о вероятном заключении эксперта. *Российская юстиция*, 2014, № 11, с. 56–59.
10. **Giboulot Q., Cogranne R., Borghys D., Bas P.** Effects and solutions of cover-source mismatch in image

- steganalysis. *Signal Processing: Image Communication*, Elsevier, 2020, vol. 86. doi:10.1016/j.image.2020.115888
11. **Mallet A., Benes M., Cogramne R.** Cover-source mismatch in steganalysis: Systematic review. *EURASIP Journal on Information Security*, 2024, no. 26. doi:10.1186/s13635-024-00171-6
 12. **Hou X., Zhang T., Xiong G., Wan B.** Forensics aided steganalysis of heterogeneous bitmap images with different compression history. *KSII Transactions on Internet and Information Systems*, 2012, vol. 6, no. 8, pp. 1926–1945. doi:10.3837/tiis.2012.08.003
 13. **Pasquet J., Bringay S., Chaumont M.** Steganalysis with cover-source mismatch and a small learning database. *EUSIPCO: European Signal Processing Conference*, 2014, pp. 2425–2429. doi:10.5281/ZENODO.43792
 14. **Abecidan R., Itier V., Boulanger J., Bas P.** Unsupervised JPEG domain adaptation for practical digital image forensics. *WIFS 2021: IEEE International Workshop on Information Forensics and Security*, 2021. doi:10.1109/WIFS53200.2021.9648397
 15. **Кафтаников И. Л., Парасич А. В.** Проблемы формирования обучающей выборки в задачах машинного обучения. *Вестник Южно-Уральского государственного университета. Серия: Компьютерные технологии, управление, радиоэлектроника*, 2016, т. 16, № 3, с. 15–24. doi:10.14529/ctcr160302
 16. **Парасич А. В., Парасич В. А., Парасич И. В.** Формирование обучающей выборки в задачах машинного обучения. Обзор. *Информационно-управляющие системы*, 2021, № 4, с. 61–70. doi:10.31799/1684-8853-2021-4-61-70
 17. **Лебедев И. С.** Адаптивное применение моделей машинного обучения на отдельных сегментах выборки в задачах регрессии и классификации. *Информационно-управляющие системы*, 2022, № 3, с. 20–30. doi:10.31799/1684-8853-2022-3-20-30
 18. **Монарев В. А., Пестунов А. И.** Повышение эффективности методов стегоанализа при помощи предварительной фильтрации контейнеров. *Прикладная дискретная математика*, 2016, № 2(32), с. 87–99. doi:10.17223/20710410/32/6
 19. **Солодуха Р. А.** Стеганоанализ изображений, модифицированных алгоритмом Bit Plane Complexity Segmentation. *Информационно-управляющие системы*, 2023, № 2, с. 27–38. doi:10.31799/1684-8853-2023-2-27-38, EDN: DXURBZ
 20. **Солодуха Р. А., Перминов Г. В., Атласов И. В.** Редукция набора детекторов LSB с заданной достоверностью. *Научно-технический вестник информационных технологий, механики и оптики*, 2022, т. 22, № 1, с. 74–81. doi:10.17586/2226-1494-2022-22-1-74-81
 21. **Солодуха Р. А.** Повышение точности стеганоанализа пространственной области изображений за счет дополнительных стегановложений. *Информационно-управляющие системы*, 2024, № 3, с. 2–10. doi:10.31799/1684-8853-2024-3-2-10, EDN: FOOKRY
 22. **Donghui Hu, Zhongjin Ma, Yuqi Fan, Lina Wang.** A study of the two-way effects of cover source mismatch and texture complexity in steganalysis. *Lecture Notes in Computer Science*, 2017, 10082, pp. 601–615. doi:10.1007/978-3-319-53465-7_45
 23. **Alkhalidi M., Abu-Elnasr O., Elarif T.** A robust steganalysis method for detecting the steganography in images. *International Journal of Intelligent Computing and Information Sciences*, 2017. doi:10.21608/ijicis.2017.19819
 24. **Hammad B., Ahmed I., Jamil N.** A steganalysis classification algorithm based on distinctive texture features. *Symmetry*, 2022, vol. 14, iss. 2, Art. 236. doi:10.3390/sym14020236
 25. **Kato H., Osuge K., Haruta S., Sasase I.** A preprocessing by using multiple steganography for intentional image downsampling on CNN-based steganalysis. *IEEE Access*, 2020, vol. 8, pp. 195578–195593. doi:10.1109/ACCESS.2020.3033814
 26. **Fridrich J.** Feature-based steganalysis for JPEG images and its implications for future design of steganographic schemes. *Information Hiding 2004. Lecture Notes in Computer Science*, 2004, vol. 3200, Springer. doi:10.1007/978-3-540-30114-1_6
 27. **Solodukha R.** The CSM problem solving technique for LSBR steganography in image spatial domain. *6th International Conference on Control Systems, Mathematical Modeling, Automation and Energy Efficiency*, 2024, pp. 521–525. doi:10.1109/SUMMA64428.2024.10803776
 28. **Kodovsky J., Fridrich J.** Calibration revisited. *Proceedings of the 11th ACM Multimedia and Security Workshop*, 2009. doi:10.1145/1597817.1597830
 29. **Pevny T., Fridrich J.** Merging Markov and DCT features for multi-class JPEG steganalysis. *Proceedings SPIE, Electronic Imaging, Security, Steganography, and Watermarking of Multimedia Contents IX*, 2007, vol. 6505. doi:10.1117/12.696774
 30. **Bas P., Filler T., Pevny T.** Break our steganographic system: The ins and outs of organizing BOSS. *Information Hiding. Lecture Notes in Computer Science*, 2011, vol. 6958. doi:10.1007/978-3-642-24178-9_5

UDC 519.6

doi:10.31799/1684-8853-2026-1-8-18

EDN: QIENHD

A novel approach to solving the CSM problem in the frequency domain of an imageR. A. Solodukha^a, PhD, Tech., Associate Professor, orcid.org/0000-0002-3878-4221, standartal@list.ru^aVoronezh State University of Engineering Technologies, 19, Revolucii Ave., 394036, Voronezh, Russian Federation

Introduction: One of the problems hindering the use of steganalysis in digital forensics concerns the discrepancy between the tested container and the training set. Currently available methods for solving this problem do not take into account the specifics of container distortions introduced by various steganographic algorithms. Nevertheless, a certain location of distortions allows using them to form a training set fitting the tested container. **Purpose:** To formalize and check the hypothesis on the effectiveness of training set formation based on distance ordering for “calibrated images” by feature steganalytic vectors while determining the locations of distortions in the frequency domain of images. To compare the recognition accuracy while using the proposed approach and while combining the features of a “calibrated image” and the original image into a single feature vector. **Results:** We demonstrate the feasibility of using steganalytic feature vectors of calibrated images to calculate the distance between containers. We propose the formalized procedure for the formation of a training set based on the distance between containers. We have created the software infrastructure for a numerical experiment which has shown that, regardless of the JPEG compression quality, the use of the feature vector of a calibrated image to form a training set is more effective than including it in the general feature vector. The main calculations are done to compute pairwise distances and can be carried out before the object of study appears. To ensure reproducibility of the experiment, we have presented datasets and program code in Kaggle. **Practical relevance:** Taking nsF5 steganographic algorithm as an example, we have demonstrated the advantage of applying for a random sample the PEV-274 steganalytic vector with the training set corresponding to the test file over CC-PEV-548. The proposed approach helps to increase the accuracy of steganalysis and to reduce the research time which is important in digital forensic practice.

Keywords — steganalysis, feature vector, nsF5, PEV-274, CC-PEV-548, steganography, Cover-Source Mismatch, machine learning, regression, distance between vectors, forensics.

For citation: Solodukha R. A. A novel approach to solving the CSM problem in the frequency domain of an image. *Informatsionno-upravliaiushchie sistemy* [Information and Control Systems], 2026, no. 1, pp. 8–18 (In Russian). doi:10.31799/1684-8853-2026-1-8-18, EDN: QIENHD

Reference

- Cogranne R., Giboulot Q., Bas P. ALASKA#2: Challenging academic research on steganalysis with realistic images. *IEEE International Workshop on Information Forensics and Security (WIFS)*, 2020, pp. 1–5. doi:10.1109/WIFS49906.2020.9360896
- Cogranne R., Giboulot Q., Bas P. The ALASKA steganalysis challenge: A first step towards steganalysis. *Proceedings of the ACM Workshop on Information Hiding and Multimedia Security*, 2019. doi:10.1145/3335203.3335726
- Bobok I. I., Koboziyeva A. A. Theoretical foundations of digital content integrity expertise. *Problemele Energeticii Regionale*, 2025, no. 1 (65), pp. 105–121. doi:10.52254/1857-0070.2025.1-65.08
- Solodukha R. A., Atlasov I. V., Kubasov I. V. *Steganoanaliz cifrovyyh izobrazheniy: tekhnologii, algoritmy, programmnaya realizatsiya* [Steganalysis of digital images: technologies, algorithms, software implementation]. Voronezh, Voronezhskij institut Ministerstva vnutrennih del Rossii Publ., 2022. 172 p. (In Russian)
- Michaylov K., Sarmah D. Steganography and steganalysis for digital image enhanced forensic analysis and recommendations. *Journal of Cyber Security Technology*, 2025, vol. 9(1), pp. 1–27. doi:10.1080/23742917.2024.2304441
- Sirota A. A., Dryuchenko M. A., Ivankov A. Y. Steganalysis of digital images by means of shallow and deep machine learning: existing approaches and new solutions. *Proceedings of Voronezh State University. Series: Systems Analysis and Information Technologies*, 2021, no. 1, pp. 33–52 (In Russian). doi:10.17308/sait.2021.1/3369
- Polunin A. A., Yandashevskaya E. A. Using of convolutional neural networks for steganalysis of digital images. *Proceedings of the Institute for System Programming of the RAS*, 2020, vol. 32, iss. 4, pp. 155–164 (In Russian). doi:10.15514/ISPRAS-2020-32(4)-11
- Lubin A. F. About admission of probabilistic conclusions of expert conclusion in a criminal trial. *Legal Science and Practice: Journal of Nizhny Novgorod Academy of the Ministry of Internal Affairs of Russia*, 2019, vol. 3, no. 47, pp. 138–142 (In Russian). doi:10.36511/2078-5356-2019-3-138-142
- Ovsyannikov I. V. To question about probabilistic conclusion of expert. *Rossiyskaya yustitsiya*, 2014, no. 11, pp. 56–59 (In Russian)
- Giboulot Q., Cogranne R., Borghys D., Bas P. Effects and solutions of cover-source mismatch in image steganalysis. *Signal Processing: Image Communication*, Elsevier, 2020, vol. 86. doi:10.1016/j.image.2020.115888
- Mallet A., Benes M., Cogranne R. Cover-source mismatch in steganalysis: Systematic review. *EURASIP Journal on Information Security*, 2024, no. 26. doi:10.1186/s13635-024-00171-6
- Hou X., Zhang T., Xiong G., Wan B. Forensics aided steganalysis of heterogeneous bitmap images with different compression history. *KSI Transactions on Internet and Information Systems*, 2012, vol. 6, no. 8, pp. 1926–1945. doi:10.3837/tiis.2012.08.003
- Pasquet J., Bringay S., Chaumont M. Steganalysis with cover-source mismatch and a small learning database. *EUSIPCO: European Signal Processing Conference*, 2014, pp. 2425–2429. doi:10.5281/ZENODO.43792
- Abecidan R., Itier V., Boulanger J., Bas P. Unsupervised JPEG domain adaptation for practical digital image forensics. *WIFS 2021: IEEE International Workshop on Information Forensics and Security*, 2021. doi:10.1109/WIFS53200.2021.9648397
- Kaftannikov I. L., Parasich A. V. Problems of training set's formation in machine learning tasks. *Bulletin of the South Ural State University. Ser. Computer Technologies, Automatic Control, Radio Electronics*, 2016, vol. 16, no. 3, pp. 15–24. (In Russian). doi:10.14529/ctcr160302
- Parasich A. V., Parasich V. A., Parasich I. V. Training set formation in machine learning tasks. Survey. *Informatsionno-upravliaiushchie sistemy* [Information and Control Systems], 2021, no. 4, pp. 61–70 (In Russian). doi:10.31799/1684-8853-2021-4-61-70
- Lebedev I. S. Adaptive application of machine learning models on separate segments of a data sample in regression and classification problems. *Informatsionno-upravliaiushchie sistemy* [Information and Control Systems], 2022, no. 3, pp. 20–30 (In Russian). doi:10.31799/1684-8853-2022-3-20-30
- Monarev V. A., Pestunov A. I. Enhancing steganalysis accuracy via tentative filtering of stego-containers. *Applied Discrete Mathematics*, 2016, no. 2 (32), pp. 87–99 (In Russian). doi:10.17223/20710410/32/6
- Solodukha R. Steganalysis of Bit Plane Complexity Segmentation algorithm. *Informatsionno-upravliaiushchie sistemy* [Information and Control Systems], 2023, no. 2, pp. 27–38 (In Russian). doi:10.31799/1684-8853-2023-2-27-38, EDN: DXURBZ
- Solodukha R. A., Perminov G. V., Atlasov I. V. Reduction of LSB detectors set with definite reliability. *Scientific and Technical Journal of Information Technologies, Mechanics and Optics*, 2022, vol. 22, no. 1, pp. 74–81 (In Russian). doi:10.17586/2226-1494-2022-22-1-74-81
- Solodukha R. A. Increasing the accuracy of spatial domain steganalysis through additional embeddings. *Informatsionno-upravliaiushchie sistemy* [Information and Control Sys-

- tems], 2024, no. 3, pp. 2–10 (In Russian). doi:10.31799/1684-8853-2024-3-2-10, EDN: FOOKRY
22. Donghui Hu, Zhongjin Ma, Yuqi Fan, Lina Wang. A study of the two-way effects of cover source mismatch and texture complexity in steganalysis. *Lecture Notes in Computer Science*, 2017, 10082, pp. 601–615. doi:10.1007/978-3-319-53465-7_45
 23. Alkhalidi M., Abu-Elnasr O., Elarif T. A robust steganalysis method for detecting the steganography in images. *International Journal of Intelligent Computing and Information Sciences*, 2017. doi:10.21608/ijicis.2017.19819
 24. Hammad B., Ahmed I., Jamil N. A steganalysis classification algorithm based on distinctive texture features. *Symmetry*, 2022, vol. 14, iss. 2, Art. 236. doi:10.3390/sym14020236
 25. Kato H., Osuge K., Haruta S., Sasase I. A preprocessing by using multiple steganography for intentional image down-sampling on CNN-based steganalysis. *IEEE Access*, 2020, vol. 8, pp. 195578–195593, doi:10.1109/ACCESS.2020.3033814
 26. Fridrich J. Feature-based steganalysis for JPEG images and its implications for future design of steganographic schemes. *Information Hiding 2004. Lecture Notes in Computer Science*, 2004, vol. 3200, Springer. doi:10.1007/978-3-540-30114-1_6
 27. Solodukha R. The CSM problem solving technique for LSBR steganography in image spatial domain. *6th International Conference on Control Systems, Mathematical Modeling, Automation and Energy Efficiency*, 2024, pp. 521–525. doi:10.1109/SUMMA64428.2024.10803776
 28. Kodovsky J., Fridrich J. Calibration revisited. *Proceedings of the 11th ACM Multimedia and Security Workshop*, 2009. doi:10.1145/1597817.1597830
 29. Pevny T., Fridrich J. Merging Markov and DCT features for multi-class JPEG steganalysis. *Proceedings SPIE, Electronic Imaging, Security, Steganography, and Watermarking of Multimedia Contents IX*, 2007, vol. 6505. doi:10.1117/12.696774
 30. Bas P., Filler T., Pevny T. Break our steganographic system: The ins and outs of organizing BOSS. *Information Hiding. Lecture Notes in Computer Science*, 2011, vol. 6958. doi:10.1007/978-3-642-24178-9_5

УВАЖАЕМЫЕ АВТОРЫ!

Научные базы данных, включая Scopus и Web of Science, обрабатывают данные автоматически. С одной стороны, это ускоряет процесс обработки данных, с другой — различия в транслитерации ФИО, неточные данные о месте работы, области научного знания и т. д. приводят к тому, что в базах оказывается несколько авторских страниц для одного и того же человека. В результате для всех по отдельности считаются индексы цитирования, что снижает рейтинг ученого.

Для идентификации авторов в сетях Thomson Reuters проводит регистрацию с присвоением уникального индекса (ID) для каждого из авторов научных публикаций.

Процедура получения ID бесплатна и очень проста, есть возможность провести регистрацию на 12 языках, включая русский (чтобы выбрать язык, кликните на зеленое поле вверху справа на стартовой странице): <https://orcid.org>