

Хабр



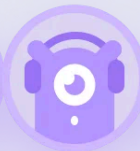
КАК СТАТЬ АВТОРОМ



Войти

1000+

вакансий с удалёнкой



Хабр Карьера



empenoso

10 мая 2023 в 03:42

Как убрать пустые оборотные страницы из PDF после двухстороннего сканирования

Средний

6 мин

9.1K

Open source*, PDF, Софт, Лайфхаки для гиков

Кейс

Около двух месяцев назад я написал статью [как сканировать многостраничные двухсторонние документы, если под рукой только обычный сканер с автоподачей](#), в которой затронул проблему того, что МФУ часто имеют дуплексную двустороннюю печать, но односторонний сканер.

Однако после решения проблемы быстрого сканирования больших двухсторонних документов, была обнаружена ещё одна проблема — некоторое количество страниц могут оказаться односторонними. И это означает, что PDF будет иметь белые страницы, например, со сканами перфораций или отверстий под кольца.

Конечно, можно удалить несколько страниц из PDF вручную, но что если таких файлов сотни, а сами документы имеют несколько десятков или даже сотен страниц как на фотографии?



Большой многостраничный документ

▸ TL;DR

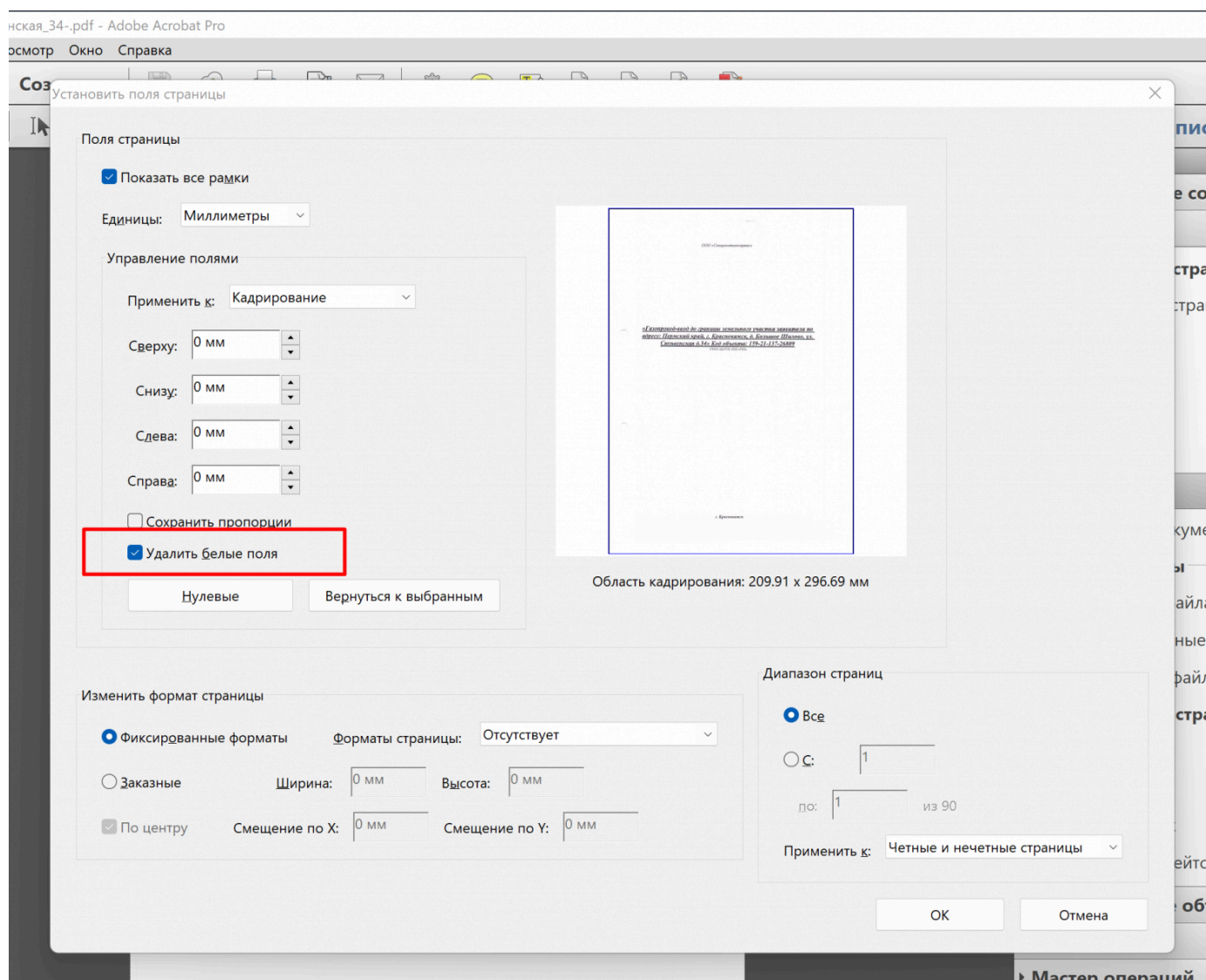
Вариант удаления пустых страниц из pdf при помощи локальной программы

Перед тем как начать писать свой скрипт я честно пытался разобраться как удалить пустые страницы из пдф при помощи штатных средств какой-нибудь программы:

1. Пытался сделать это при помощи бесплатной открытой [PDFsam Basic](#), которая доступна под Linux и Windows, и MacOS, потому что в интернете нашёл инструкции, но они оказались устаревшими.
2. Пытался сделать это при помощи Adobe Acrobat Pro, но у меня не получилось. Делал по инструкции:
 1. Откройте файл PDF в Adobe Acrobat.

2. Нажмите на вкладку «Инструменты» в верхней строке меню.
3. Выберите «Страницы» из списка инструментов справа.
4. Нажмите «Обрезать» в меню инструментов «Страницы».
5. В диалоговом окне «Обрезка страниц» выберите параметры «Удалить белые поля» и «Удалить белые поля для всех страниц».
6. Нажмите «ОК», чтобы применить изменения.

Эти действия должны были автоматически удалить все пустые страницы из файла PDF, но у меня этого не произошло.

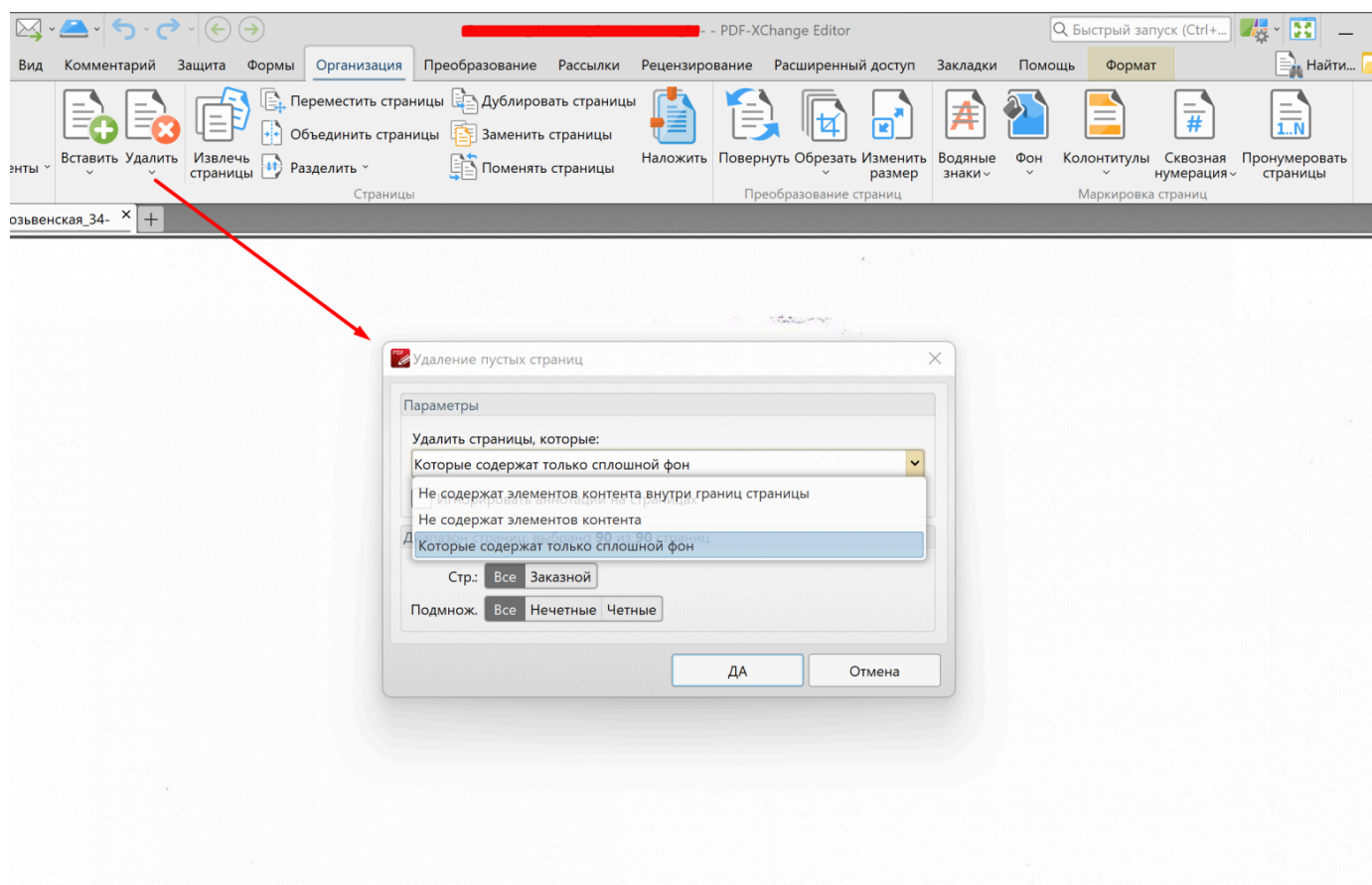


Adobe Acrobat Pro и удаление пустых страниц

3. Попытка сделать это при помощи PDF-XChange Editor, но у меня тоже не получилось.
У меня была инструкция:

1. Загрузите файл PDF: выберите «Файл» > «Открыть» или нажмите Ctrl + O на клавиатуре, затем найдите и выберите файл PDF, из которого вы хотите удалить пустые страницы.
2. После загрузки PDF-файла щелкните вкладку «Организация» на верхней панели инструментов.
3. Выбрав все страницы, нажмите кнопку «Удалить пустые страницы».

Прогресс пробежал, но пустые страницы оставались на месте для любых из трех вариантов.



PDF-XChange Editor

Использование локальной программы, конечно, было бы лучшим вариантом, потому что это гарантировало, что PDF-файлы останутся на компьютере, обеспечивая конфиденциальность и безопасность по сравнению с использованием онлайн-инструментов.

Вариант удаления пустых страниц из pdf при помощи онлайн-инструментов

Но раз с локальными инструментами у меня не пошло, решил попробовать онлайн сервисы.

Я смог найти несколько доступных онлайн-инструментов, которые могли бы помочь автоматически удалить пустые страницы из PDF-файла:

1. Sejda (<https://www.sejda.com/delete-pdf-pages>)
2. Smallpdf (<https://smallpdf.com/delete-pages-from-pdf>)
3. DeftPDF (<https://deftpdf.com/delete-pdf-pages>)

Ни в одном из них я не смог найти опцию автоматического распознавания пустых страниц, хотя в поисковике попадались ссылки на несуществующие сейчас страницы (pdf remove blank pages) этих сервисов.

Ну и конечно использование онлайн-инструментов может поставить под угрозу конфиденциальность и безопасность ваших документов.

Вариант удаления пустых страниц из pdf при помощи локального bash скрипта и консольной программы PDFtk

После постигшей неудачи решил написать свой собственный скрипт который удалит пустые страницы из всех pdf файлов в текущем каталоге.

При изучении вопроса наткнулся на [большую дискуссию](#), где обсуждался вопрос как лучше [удалить пустые страницы из pdf при помощи командной строки](#). Предлагались разные методы, но у меня были все документы сканированные и это значит, что даже на пустом листе какая-то информация всё равно была — сканы отверстий под перешивку или просто грязь со сканера.

Решил что будет следующий алгоритм:

1. Разделяю PDF документ на отдельные файлы.
2. Страницы меньше определенного размера удаляю.
3. Склеиваю оставшиеся страницы обратно.
4. Повторяю столько раз, сколько PDF файлов в текущей папке.
5. PROFIT

После нехитрых манипуляций получился файл `blank_page_remover.sh` :

```
# Подробнее в статье Как убрать пустые оборотные страницы из PDF после двухстороннего с
# https://habr.com/ru/articles/733754/
# Михаил Шардин https://shardin.name/

#!/bin/bash
datetime=$(date +"%Y-%m-%d_%H-%M-%S")
# Создаём единый лог файл для всех действий и папку куда перемещаем вырезанные страницы
log_file="blank_page_remover_${datetime}.log"
touch $log_file
mkdir removed
# Перебираем все PDF файлы в текущем каталоге
for file in *.pdf; do
    echo "Работаем с $file..." >> "$log_file"
    # Разделяем PDF файл на отдельные страницы
    echo "Разделяем $file на отдельные страницы..." >> "$log_file"
    pdftk "$file" burst output "${file%.*}_pg_%04d.pdf" >> "$log_file" 2>&1
    # Удаляем файлы страниц, размер которых меньше чем XX килобайт
    echo "Удаляем файлы страниц, размер которых меньше чем 35 килобайт..." >> "$log_file"
    for page in "${file%.*}_pg_*.pdf; do
        size=$(wc -c < "$page")
        if [[ $size -lt 35000 ]]; then
            echo "Удаляем $page (размер: $size байт)..." >> "$log_file"
            mv "$page" "removed/"
            #rm "$page"
        fi
    done
    # Склеиваем оставшиеся страницы в новый файл
    echo "Склеиваем оставшиеся страницы в новый файл..." >> "$log_file"
    pdftk "${file%.*}_pg_*.pdf cat output "${file%.*}_без пустых.pdf" compress >> "$log_
    # Удаляем временные файлы
    echo -e "Удаляем временные файлы...\n" >> "$log_file"
    rm "${file%.*}_pg_*.pdf
done
```

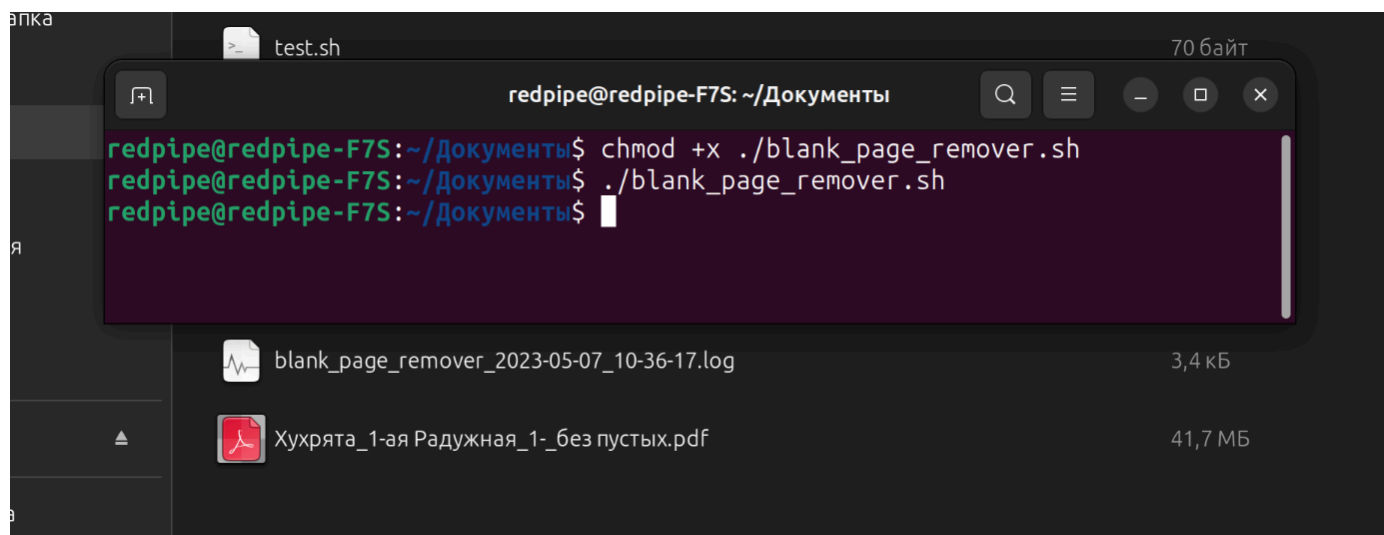
Для работы скрипта понадобится Pdftk (сокращение от PDF Toolkit) — это инструмент командной строки для работы с PDF-файлами. Как его установить для разных операционных систем можно узнать [в предыдущей статье](#).

Как воспользоваться скриптом удаления пустых страниц из PDF документа

Чтобы выполнить сценарий `bash` на компьютере, выполните следующие действия в зависимости от операционной системы:

Для Linux и macOS:

1. Откройте Терминал: нажмите `Ctrl + Alt + T` в Linux или откройте **Терминал** из папки **Приложения > Утилиты** в macOS.
2. Перейдите в каталог, где находится скрипт: используйте команду `cd`, за которой следует путь к каталогу. Например:
`cd /путь/к/скрипту`
3. Сделайте скрипт исполняемым:
`chmod +x blank_page_remover.sh`
4. Выполните этот сценарий. Запустите сценарий, введя `./`, а затем имя сценария:
`./blank_page_remover.sh`
5. PROFIT!
Скрипт создаст новые pdf файлы без пустых страниц и подробный лог действий.



Терминал в Ubuntu и результат выполнения скрипта `blank_page_remover.sh`

Для Windows (используя GitBash или WSL):

1. Установите GitBash или WSL: если вы еще этого не сделали, установите [GitBash](https://gitbash.com/) или [подсистему Windows для Linux \(WSL\)](https://docs.microsoft.com/en-us/windows/wsl/).

- Откройте Git Bash или WSL: щелкните правой кнопкой мыши папку, содержащую скрипт, и выберите `GitBash здесь` или `Открыть в WSL`.
- Сделайте скрипт исполняемым:

```
chmod +x blank_page_remover.sh
```
- Выполните этот сценарий. Запустите сценарий, введя `./`, а затем имя сценария:

```
./blank_page_remover.sh
```
- PROFIT!
Скрипт создаст новые pdf файлы без пустых страниц и подробный лог действий.

Актуальная версия скрипта [всегда доступна на гитхабе](#).

Заключение

Удаление пустых страниц из PDF-файлов после двустороннего сканирования может оказаться непростой задачей, особенно при работе с большими объемами документов. Тем не менее, эта статья предоставила вам решение в виде использования автоматического локального сценария bash с консольной программой PDFtk.

Следуя подробным инструкциям вы сможете эффективно избавиться от пустых страниц и поддерживать чистый профессиональный вид отсканированных PDF-документов.

Независимо от объема или сложности ваших файлов, это решение упростит ваш рабочий процесс и сэкономит ваше время и усилия.

Автор: Михаил Шардин

 [Моя онлайн-визитка](#)

 [Telegram «Умный Дом Инвестора»](#)

10 мая 2023 г.

Теги: [bash](#), [pdftk](#), [сканирование](#), [документы](#)

Хабы: [Open source](#), [PDF](#), [Софт](#), [Лайфхаки для гиков](#)

Редакторский дайджест

Присылаем лучшие статьи раз в месяц



**183****87.1**

Карма

Рейтинг

Михаил Шардин @empenoso

Автоматизация / Данные / Финансы / Умные дома

[Подписаться](#)[Сайт](#) [Сайт](#) [Github](#)

Комментарии 10

Публикации

[ЛУЧШИЕ ЗА СУТКИ](#)[ПОХОЖИЕ](#)**rssdev10**

23 часа назад

Почему въехав по «визе талантов» в США я с радостью вернулся в Россию

**Средний**

32 мин



35K

[Мнение](#) **+174**

112

512

**melnik909**

19 часов назад

Вы не знаете CSS. Мои вопросы о CSS с ответами. Часть 2

**Средний**

7 мин



2K

[Обзор](#) **+39**

32

1

**DAN_SEA**

17 часов назад

Генерация случайных чисел

**Средний**

10 мин



2.4K

[Обзор](#) **+32**

25

34

**OrkBiotechnologist**

23 часа назад

VPS за 139 рублей — дом для вашего резюме на основе Hugo

**Простой**

7 мин



8.4K

[Тutorial](#)

 +27 46 12**PatientZero**

2 часа назад

Пишем стек TCP/IP с нуля: Ethernet, ARP, IPv4 и ICMPv4



Простой



13 мин



879

Тutorial

Перевод

 +19 30 1**tertiumnon**

17 часов назад

Минимум книг, которые нужно прочитать начинающему или продолжающему свою кривую обучения программисту



Простой



3 мин



8.8K

Обзор

 +19 193 21**Ibkanter**

1 час назад

Бэкдор Auto-color: разбор угрозы, технический анализ и способы защиты



Средний



4 мин



281

Обзор

 +13 5 2**FlatSpike**

18 часов назад

Создаём многомодульную библиотеку на Android: как же собрать fat-aar?



Средний



19 мин



531

Кейс

+13

16

0



alexander-shustanov

19 часов назад

В поисках идеального Database-клиента для IDE: Amplicode выбирает DBeaver



Простой



6 мин



2.6K

+13

13

7



ptsecurity

21 час назад

Безопасность без боли: плагины, которые упрощают жизнь разработчикам



7 мин



930

Кейс

+12

13

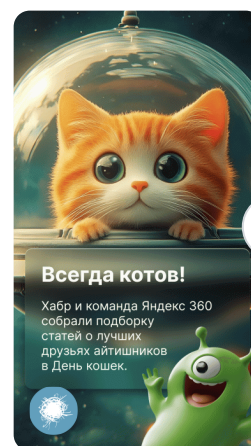
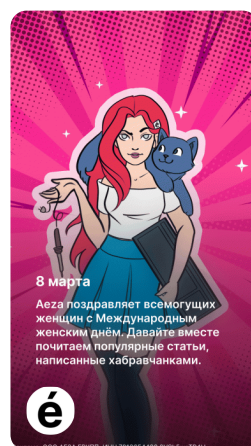
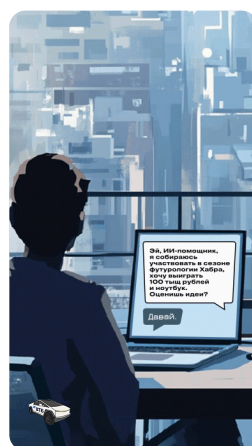
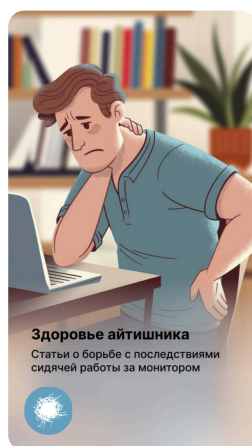
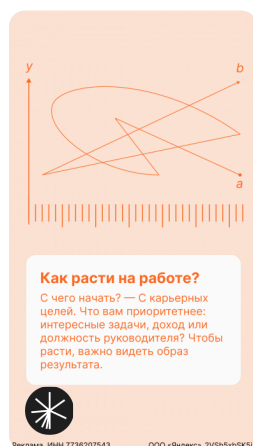
3

Huawei Mate 70Pro: макрофото на 10/10

Промо

Показать еще

ИСТОРИИ



[Как расти на работе?](#)[Здоровье айтишника](#)[Угадайте будущее в новом сезоне](#)[С праздником весны!](#)[Всегда котов!](#)

ВОПРОСЫ И ОТВЕТЫ

Как правильно настроить сканирование в папку на МФУ Ricoh Aficio MP 5002SP?

Сканирование · Простой · 1 ответ

Как отредактировать nginx с помощью ansible?

Nginx · Простой · 7 ответов

Быстро перенести файлы с 1 хостинга на другой в потоке, нужны идеи для скрипта?

bash · Простой · 0 ответов

Какие есть инструменты минификации shell кода или bash-скриптов?

bash · Простой · 3 ответа

Почему bash-скрипт неправильно копирует папку на macOS?

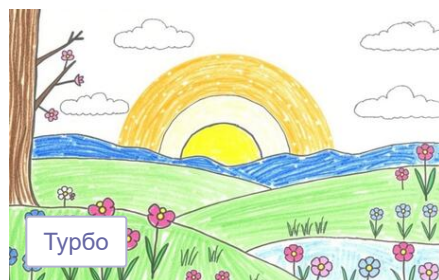
bash · Простой · 2 ответа

[Больше вопросов на Хабр Q&A](#)

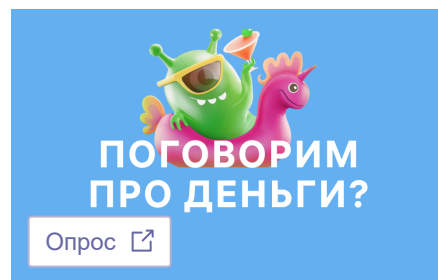
МИНУТОЧКУ ВНИМАНИЯ



Девушка с розовыми волосами и Слизень на планете Рбах

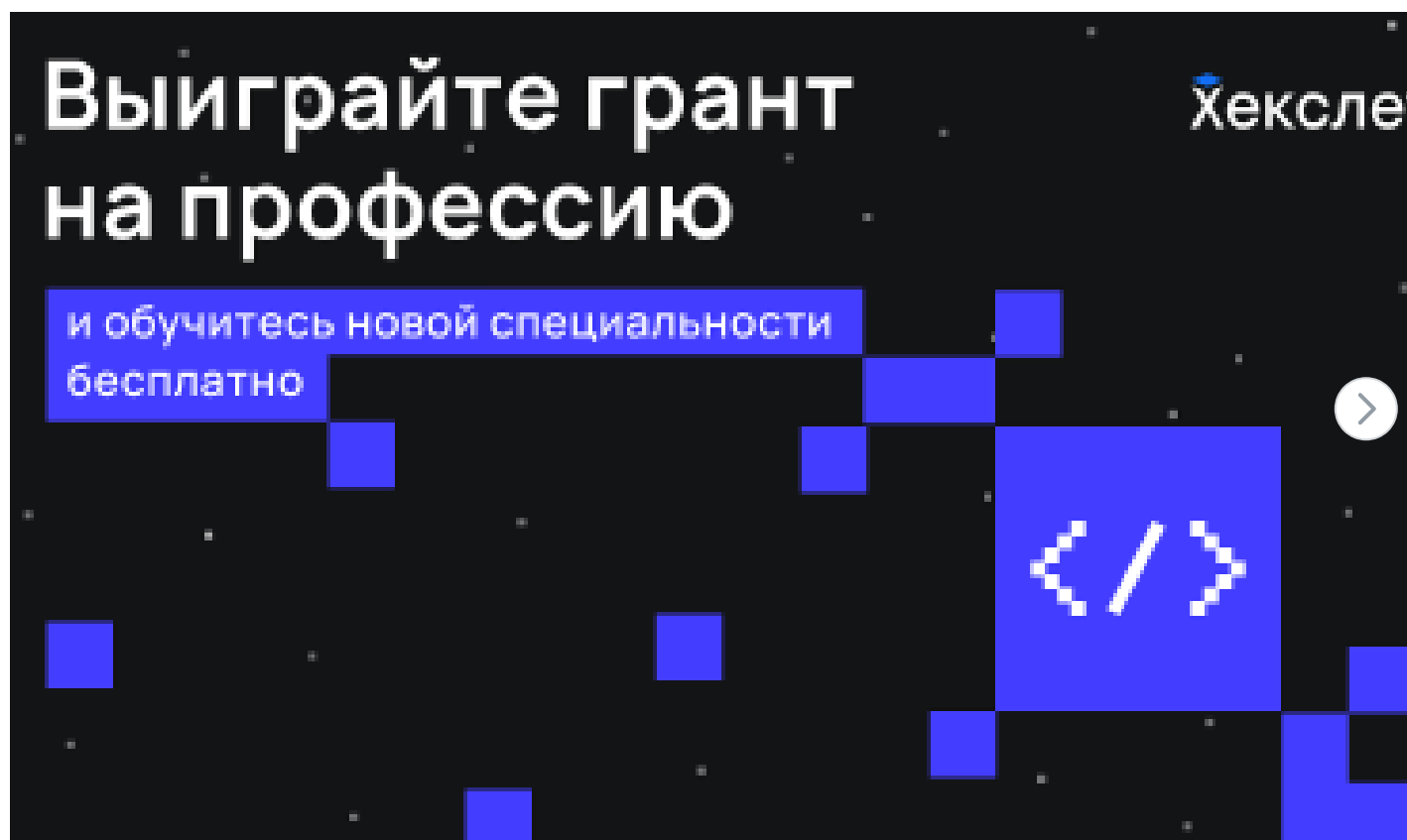


Весеннее пробуждение IT-мероприятий в Календаре



Как айтишники подходят к финансовому планированию?

БЛИЖАЙШИЕ СОБЫТИЯ



17 февраля – 24 марта

Конкурс «Снежный код» от Хекслета. Три гранта на бесплатное 10-месячное обучение

Онлайн

Разработка

[Больше событий в календаре](#)

Хабр



🌐 [Настройка языка](#)

[Техническая поддержка](#)

© 2006–2025, Habr