

Хабр



КАК СТАТЬ АВТОРОМ



Найди стажировку на Битве



empenoso

10 мая в 05:42

Как убрать пустые оборотные страницы из PDF после двухстороннего сканирования

Средний

6 мин

5.8K

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?



+10



37



10



Большой многостраничный документ

▸ TL;DR

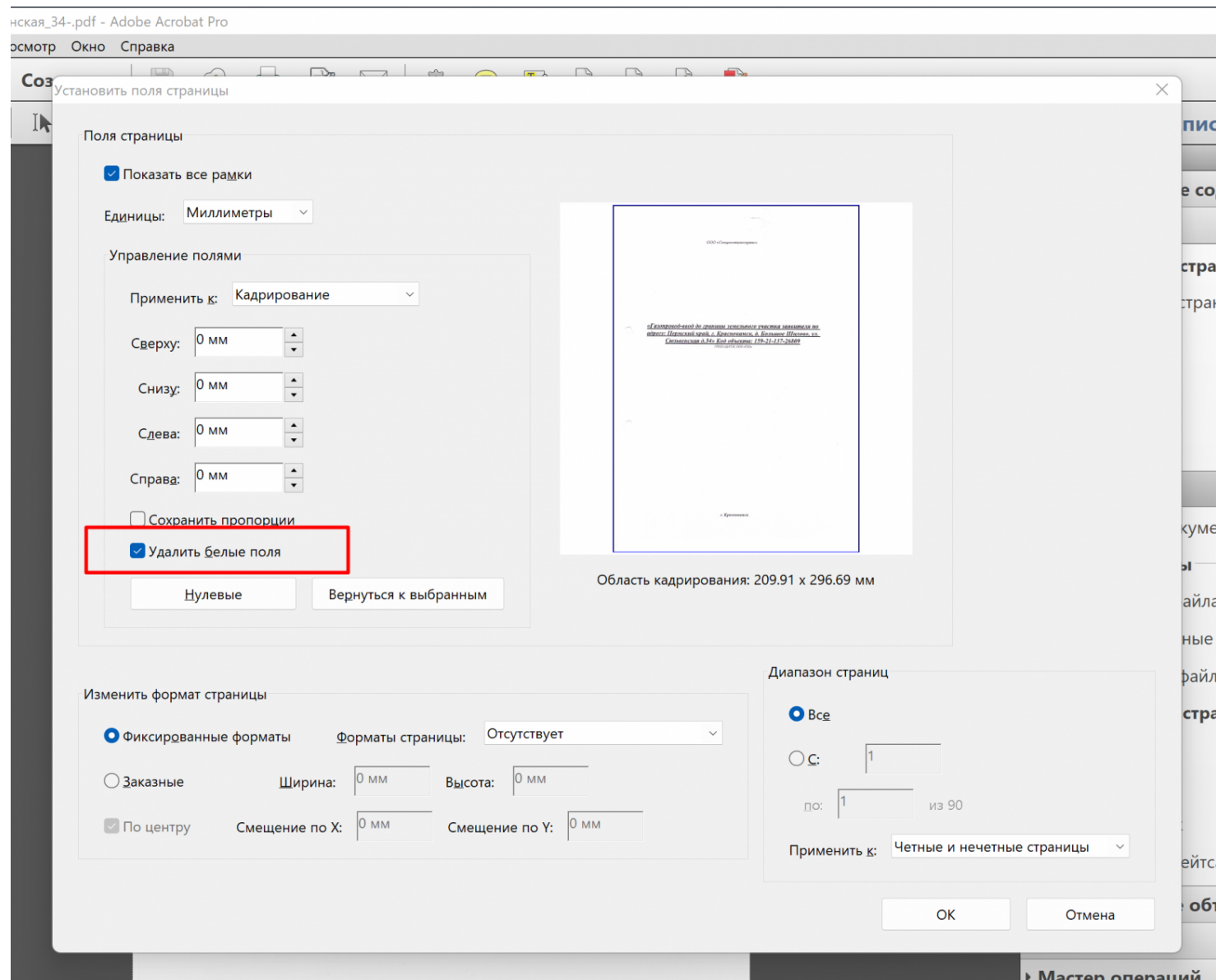
Вариант удаления пустых страниц из pdf при помощи локальной программы

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.

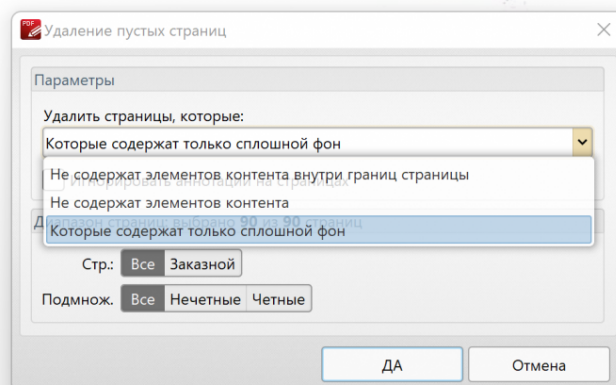
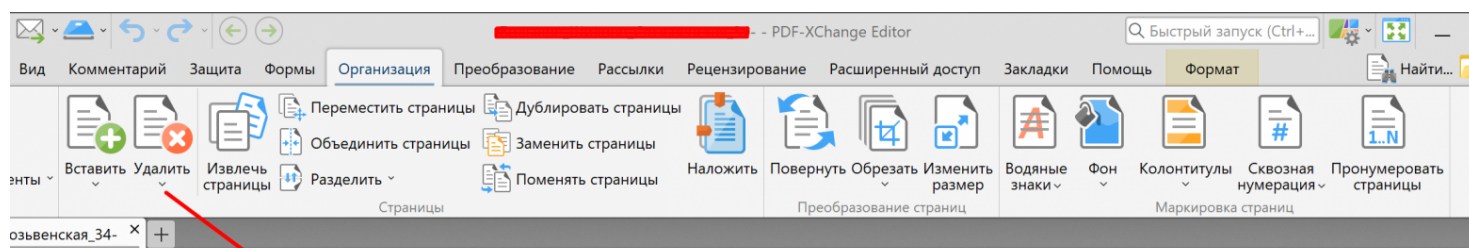


Adobe Acrobat Pro и удаление пустых страниц

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 скан
# https://habr.com/ru/articles/733754/

#!/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_file"
    # 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:

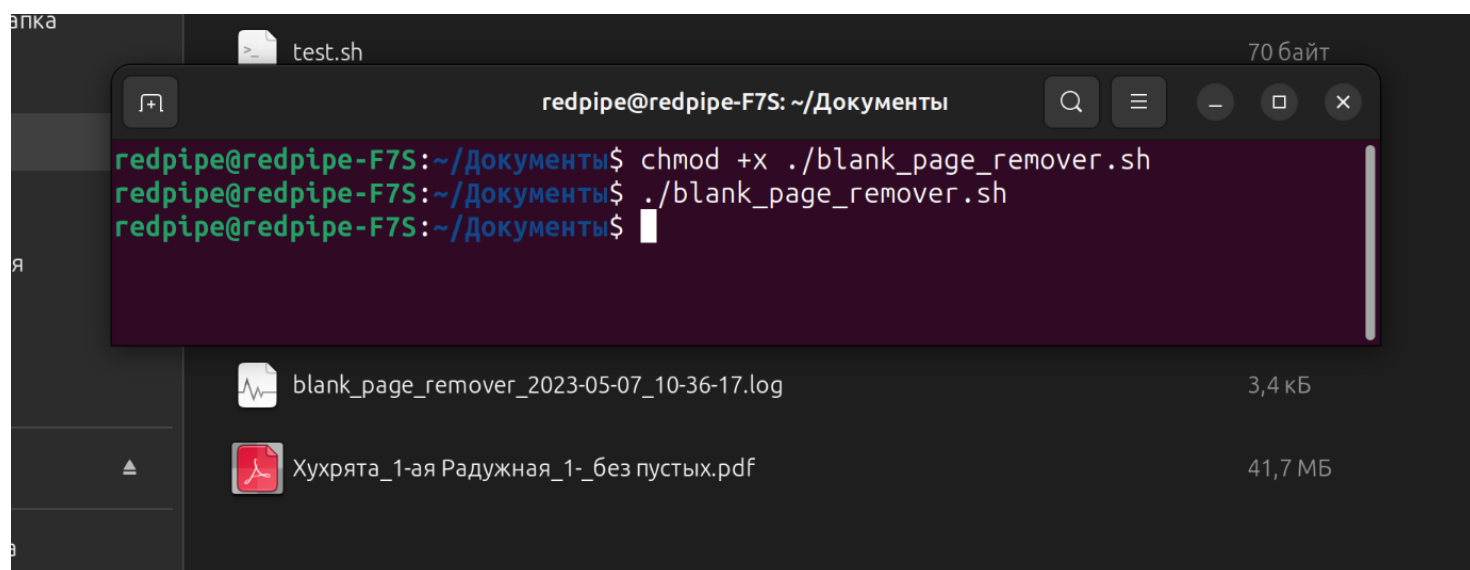
```
chmod +x blank_page_remover.sh
```

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](#) или [подсистему Windows для Linux \(WSL\)](#).
2. Откройте Git Bash или WSL: щелкните правой кнопкой мыши папку, содержащую скрипт, и выберите `GitBash здесь` или `Открыть в WSL`.

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.

Актуальная версия скрипта [всегда доступна на 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](#),

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

**129****2**

Карма

Рейтинг

Михаил Шардин [@empenoso](#)

Разработчик

[Сайт](#)

Реклама

РЕКЛАМА • MEDIASNIPEP

ТИНЬКОФФ

Реклама. Рекламодатель АО «Тинькофф Банк» (ИНН 7710140679), лицензия №2673. Сроки акции с 20.11.2023 по 30.11.2023. Подробнее на [blackfriday23.tinkoff.ru](#).

**ВНИМАНИЕ! СКОРО
ЧЕРНАЯ ПЯТНИЦА
В ТИНЬКОФФ**

**КЭШБЭК
ДО 60%
И БОНУСЫ**

0+

i

Комментарии 10

Публикации

ЛУЧШИЕ ЗА СУТКИ

ПОХОЖИЕ

**yurabeznos**

22 часа назад

Черкаш-код: изобретение и внедрение



Простой



4 мин



9.3K

Тutorial



+72



22



71

**ru_vds**

17 часов назад

Реверс-инжиниринг ячейки регистра процессора Intel 386



Средний



10 мин



3.5K

Кейс

Перевод



+41



40



2

**AlexeyW100**

23 часа назад

Дизайн-система Gravity UI: как легко построить свой интерфейс



Простой



7 мин



7.7K



+40



87



5

**AnaSergeeva**

19 часов назад

97 откликов, 2 тестовых, 3 технических собеседования — и оффер в IT-компанию у меня в кармане



Простой



11 мин



20K



+39



98



140

**ru_vds**

21 час назад

Мобильные суперприложения выгодны корпорациям, но это кошмар для простых людей



Простой



7 мин



6.5K

[Мнение](#)

◆ +36

🔖 25

💬 27

**CatScience**

18 часов назад

Немного про воронье зрение

👉 Простой

🕒 2 мин

👁 8.3K

◆ +35

🔖 27

💬 24

**хепон**

12 часов назад

Может ли быть уязвимость в дизайне, контенте и CSS и разбор такой уязвимости(?) на Госуслугах

👉 Простой

🕒 4 мин

👁 4K

[Обзор](#)

◆ +32

🔖 22

💬 12

**BitterLollipop**

17 часов назад

Сокровища HTML: 7 тегов, которые упростят вам жизнь

👉 Простой

🕒 7 мин

👁 5.3K

[Обзор](#)

◆ +31

🔖 117

💬 10

**Oksana_Nedvigina**

16 часов назад

Машина свободы: как чилийские социалисты придумали компьютер для управления экономикой

🕒 9 мин

👁 2.8K

[Ретроспектива](#)

◆ +26

🔖 23

💬 10



IgnatChuker

18 часов назад

Разговор с одним из основателей «Базальт СПО» Алексеем Смирновым о свободном ПО и взаимодействии в комьюнити

Простой

13 мин

1.6K

Интервью

+25

10

2

Ретушь после генерации: как доработать арты до идеала с нейросетями

Турбо

Показать еще

ИСТОРИИ



Недельный топ-7 статей из блогов компаний

Топ-7 годных статей из блогов компаний



Сколько тратят в IT: сеньор бэкендер

Спрашиваем, сколько зарабатывают и тратят IT-специалисты по категориям бюджета: жилье, еда, хобби и другие

Сколько тратят в IT: сеньор бэкендер



Работаете на удалёнке?

Собрали для вас откровения и лайфхаки об удалённой работе

Мифы и легенды удалёнки

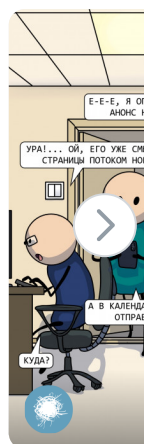


Спасти бизнес за неделю до отпуска



Библиотека айтишника
Подборки новых и просто полезных книг для пополнения библиотеки.

Полезные книги для библиотеки айтишника



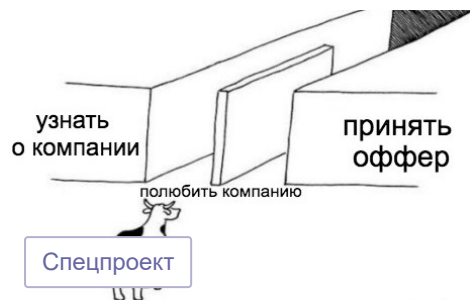
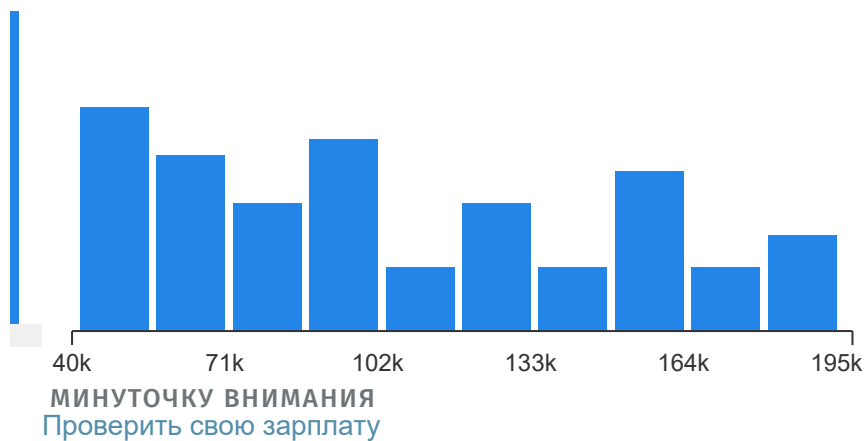
Перевернуть календарь добавить с

СРЕДНЯЯ ЗАРПЛАТА В IT

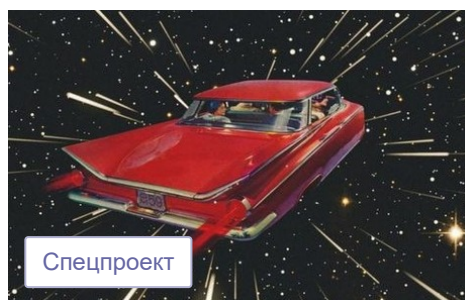
110 667

₽/мес.

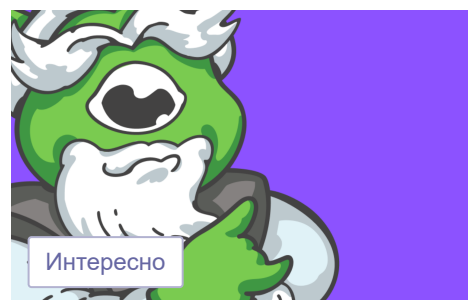
— средняя зарплата во всех IT-специализациях по данным из 27 243 анкет, за 2-ое пол. 2023 года. Проверьте «в рынке» ли ваша зарплата или нет!



Помогаем делать взвешенный выбор с 2020 года



Три карьерные дорожные истории о мягких навыках



Глупым вопросам и ошибкам — быть! IT-менторство на ХК

Хабр



Настройка языка

Техническая поддержка

© 2006–2023, Habr