

MAE IDP

Intelligent Document Processing

PDF Parser • OCR • Invoice Recognition

Автоматическое распознавание счетов и документов с OCR.

Извлекает: вендор, номер счёта, внутренний входящий номер документа, VAT ID

Быстрая установка

Windows

```
git clone https://github.com/mcbile/mae-idp.git  
cd mae-idp  
install.bat
```

macOS / Linux

```
git clone https://github.com/mcbile/mae-idp.git  
cd mae-idp  
chmod +x install.sh run.sh  
../install.sh
```

Требования

Windows

- **Windows 10/11**
- **Python 3.10+** — [скачать](#)
- **Tesseract OCR** — [скачать](#)
- **Poppler** — устанавливается автоматически через `install.bat`

macOS

- **macOS 10.15+ (Catalina или новее)**
- **Python 3.10+** — `brew install python@3.11`
- **Homebrew** — [установка](#)
- **Tesseract & Poppler** — устанавливаются автоматически через `install.sh`

Linux (Ubuntu/Debian/Fedora/Arch)

- **Python 3.10+**
- **Tesseract & Poppler** — устанавливаются автоматически через `install.sh`

Установка зависимостей

Windows — Tesseract

1. Скачай [tesseract-ocr-w64-setup-5.3.3.exe](#)
2. При установке выбери языки: **German, English**
3. Установи в **C:\Program Files\Tesseract-OCR**

macOS — Tesseract

```
brew install tesseract tesseract-lang
```

Linux — Tesseract

```
# Ubuntu/Debian
```

```
sudo apt install tesseract-ocr tesseract-ocr-deu tesseract-ocr-eng poppler-utils
```

```
# Fedora
```

```
sudo dnf install tesseract tesseract-langpack-deu tesseract-langpack-eng poppler-utils
```

```
# Arch
```

```
sudo pacman -S tesseract tesseract-data-deu tesseract-data-eng poppler
```

► Запуск

Windows

```
run.bat
```

Или напрямую:

```
venv\Scripts\python app\mae.py
```

macOS / Linux

```
./run.sh
```

Или напрямую:

```
source venv/bin/activate
```

```
python3 app/mae.py
```

После запуска откройте браузер: <http://127.0.0.1:8766>



Мобильный доступ (iOS / Android)

MAE IDP работает как Progressive Web App (PWA). Вы можете использовать его с мобильных устройств:

Шаг 1: Запустите сервер на компьютере

```
./run.sh # или run.bat на Windows
```

Шаг 2: Узнайте IP-адрес компьютера

```
# macOS/Linux
```

```
ifconfig | grep "inet "
```

```
# Windows
```

```
ipconfig
```

Шаг 3: Откройте в мобильном браузере

Откройте на телефоне: http://ВАШИ_IP:8766

Например: <http://192.168.1.100:8766>

Шаг 4: Установите как приложение (опционально)

iOS (Safari):

1. Нажмите кнопку "Поделиться" (квадрат со стрелкой)
2. Выберите "На экран «Домой»"

Android (Chrome):

1. Нажмите меню (три точки)
2. Выберите "Добавить на главный экран"

Примечание: Компьютер и телефон должны быть в одной Wi-Fi сети.



CLI инструмент (пакетная обработка)

Windows

```
python app\batch_rename.py "D:\Invoices" "D:\Sorted"
```

```
# Только анализ без копирования
```

```
python app\batch_rename.py "D:\Invoices" "D:\Sorted" --dry-run
```

```
# Без Excel отчёта
```

```
python app\batch_rename.py "D:\Invoices" "D:\Sorted" --no-report
```

macOS / Linux

```
python3 app/batch_rename.py ~/Documents/Invoices ~/Documents/Sorted
```

```
# Только анализ без копирования
```

```
python3 app/batch_rename.py ~/Documents/Invoices ~/Documents/Sorted --dry-run
```

```
# Без Excel отчёта
```

```
python3 app/batch_rename.py ~/Documents/Invoices ~/Documents/Sorted --no-report
```



📁 Структура проекта

mae-idp/

 └── app/

 └── mae.py # Веб-приложение (FastAPI + WebView)

 └── core.py # Общая логика OCR (базовый класс)

 └── batch_rename.py # CLI инструмент пакетной обработки

 └── setup_env.py # Настройка окружения (Tesseract, Poppler)

 └── templates/

 └── index.html # Веб-интерфейс

 └── data/

 └── input/ # Входящие документы (временные)

 └── output/ # Excel отчёты

 └── archive/ # Обработанные файлы

 └── install.bat # Автоустановка (Windows)

 └── install.sh # Автоустановка (macOS/Linux)

 └── run.bat # Запуск GUI (Windows)

 └── run.sh # Запуск GUI (macOS/Linux)

 └── requirements.txt # Python зависимости

 └── CHANGELOG.md # История изменений

 └── BACKLOG.md # Планы развития

 └── CLAUDE.md # Инструкции для разработки

 └── README.md

 └── README.pdf

Возможности

GUI (mae.py)

-  **Drag & Drop** — перетащи файлы в окно
-  **Folder Watch** — автоматический мониторинг папки
-  **Cloud Support** — Google Drive Desktop, OneDrive, Dropbox
-  **Excel Export** — выгрузка результатов
-  **Dark/Light Theme** — переключение темы

CLI (batch_rename.py)

-  **Пакетная обработка** — обработка целых папок
-  **Автосортировка** — раскладка по папкам вендоров
-  **Excel отчёт** — автоматический отчёт о результатах
-  **Dry-run режим** — предпросмотр без изменений

Поддерживаемые форматы

- PDF (первая страница)
- JPG / JPEG
- PNG
- TIFF / TIF

Извлекаемые данные

Поле	Описание	Confidence
Vendor	Название вендора (45+ известных компаний)	+30
Invoice Number	Номер счёта/накладной	+30
Internal Number	Внутренний входящий номер (из QR или угла)	+30
VAT ID	Идентификатор НДС (DE, AT, CH)	+10

Порог успеха: 50+ баллов = **success**, иначе **review**

Известные вендоры

Amazon, DHL, UPS, FedEx, Deutsche Telekom, Vodafone, O2, IKEA, MediaMarkt, Saturn, Conrad, Reichelt, RS Components, Mouser, DigiKey, Farnell, Würth, Hoffmann, Grainger, Mercateo, Staples, Office Depot, Viking, Büroshop24 и другие.

Конфигурация

Приложение хранит настройки в `data/config.json`:

```
{
  "watch_path": "G:\\My Drive\\Invoices",
  "output_path": "G:\\My Drive\\Reports"
}
```

API Endpoints

GUI приложение запускает локальный сервер на <http://127.0.0.1:8766>

Метод	Endpoint	Описание
GET	/	HTML интерфейс
GET	/api/status	Статус OCR и Watcher
POST	/api/parse	Загрузить и обработать файл
GET	/api/results	Получить все результаты
DELETE	/api/results	Очистить результаты
POST	/api/export	Экспорт в Excel
POST	/api/watcher/start	Запустить мониторинг папок
POST	/api/watcher/stop	Остановить мониторинг
GET	/api/browse	Диалог выбора папки
GET	/api/detect-gdrive	Поиск Google Drive папок
GET	/api/open/{folder}	Открыть папку в проводнике



Troubleshooting

"Tesseract not found"

Windows: Убедись что Tesseract установлен в `C:\Program Files\Tesseract-OCR`

macOS: Запусти `brew install tesseract tesseract-lang`

Linux: Запусти `sudo apt install tesseract-ocr`

"PDF conversion failed"

Windows: Poppler не установлен. Запусти `install.bat` заново.

macOS: Запусти `brew install poppler`

Linux: Запусти `sudo apt install poppler-utils`

Окно сразу закрывается (Windows)

Запусти через CMD чтобы увидеть ошибку:

`cmd /k run.bat`

Низкий Confidence

- Проверь качество скана (минимум 300 DPI)
- Убедись что документ на немецком или английском
- Проверь что номера не обрезаны

Мобильный доступ не работает

- Убедись что телефон и компьютер в одной Wi-Fi сети
- Проверь что firewall не блокирует порт 8766
- Попробуй отключить VPN