

Розробка алгоритму пошуку та розпізнавання зони машинного зчитування на паспорті

В. І. Уманський

Машинозчитувана зона (МЗЗ) присутня на будь-якому сучасному паспорті громадянина більшості країн світу та на закордонних паспортах України. Це два алфавітно-цифрових рядки, розташовані на нижній частині головної сторінки документа, в яких дублюється основна інформація про власника паспорта з метою її зчитування комп'ютером. В даній роботі пропонується метод локалізації та розпізнавання МЗЗ на сторінці паспорту та реалізація, що може бути застосована у системах мобільної ідентифікації та реєстрації.

Вступ

Міжнародна асоціація авіаперевезників (ICAO) запропонувала стандарт [1], що був пізніше затверджений комітетом ISO (ISO/IEC 7501-1) щодо машинозчитуваних подорожніх документів, в тому числі паспортів. Стандарт покликаний поліпшити безпеку та спростити формальні процедури перетину кордонів.



Рис. 1. Приклад головної сторінки машинозчитуваного паспорта.

Головна сторінка паспорту призначена як для візуального, так і машинного зчитування. Вона складається з декількох зон (див. Рис. 2), деякі з них є необов'язковими, деякі мають нефіксоване розташування. Зона, призначена для машинного зчитування є обов'язковим елементом та завжди займає нижню частину документа.

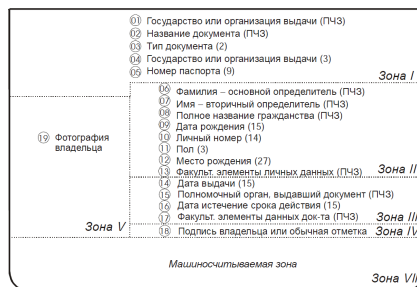


Рис. 2. Приклад розташування зон на паспорті [1].

МЗЗ складається з двох рядків по 44 символи. Стандартом зафіксовані фізичні розміри сторінок паспорту, розташування зони, розмір та гарнітура шрифту. Також обмежено набір допустимих символів: дозволеними лише є великі символи латиниці, цифри та знак < (див. Рис. 3).

МЗЗ містить основну інформацію про власника паспорта: номер паспорта, прізвище та ім'я, стать, дату народження, термін дії документу, код країни видачі паспорту та громадянство, додаткові дані. Найбільш важливі дані як, наприклад, номер паспорту та

0123456789
ABCDEFGHI
JKLMNOPQR
STUVWXYZ <

Рис. 3. Набір символів, дозволений в зоні машинного зчитування.

його термін дії захищені контрольними сумами – таким чином, в процесі розпізнавання, є можливість перевірити відсутність помилок.

ІСАО лише пропонує стандарт, що описує вимоги до машинозчитуваних документів. Уточнення специфікацій, розробка засобів друку документів (або придбання відповідних технологій) та програмного забезпечення для роботи з ними залишається за державами.

У випадку мобільних застосувань для отримання зображення паспорту припускається використання веб-камери або протяжних сканерів замість стаціонарних. Це призводить до значно більших допусків на зміщення та поворот паспорту на зображенні, а також більш низьку якість зображення. Крім цього, нерівномірність освітлення та інші фактори значно поскладнюють процес локалізації та розпізнавання зони машинного зчитування. Враховуючи зростаючі вимоги до безпеки кордонів держави, поставлена задача є досить актуальною.

Пошук машинозчитуваної зони на зображенні

На вході задається кольорове зображення або зображення в тонах сірого головної сторінки паспорта. Оскільки маркерів для пошуку зони машинного зчитування не передбачено, орієнтація та зміщення паспорта на зображенні є невідомими, єдиним надійним способом пошуку МЗЗ є виявлення структури, що відповідає тексту машинозчитуваної зони. Це можливо реалізувати, якщо врахувати наступні факти: кількість символів в рядку та кількість рядків зафіксовано в стандарті, так само як і відстань між символами рядка та відстань між рядками. Крім того, проаналізувавши набір допустимих символів (Рис. 3), можна помітити, що всі символи складаються з єдиної компоненти зв'язності.

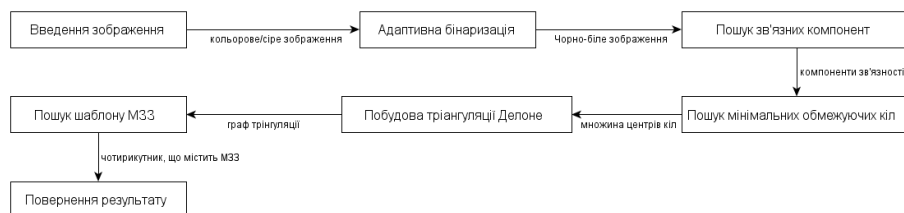


Рис. 4. Схема процесу виділення зони машинного зчитування.

На Рис. 4 зображено загальну схему алгоритму. Першим етапом роботи є бінаризація вхідного зображення: з кольорового зображення утворюється чорно-біле – таким чином текст відділяється від фонового візерунка паспорта. Оскільки сторінка паспорта може бути освітлена нерівномірно, звичайна порогова бінаризація є неприйнятною, тому використовується адаптивна бінаризація, яка визначає поріг для кожної точки окремо на основі інтенсивності точок з певного околу.

На бінаризованому зображенні виконується пошук зв'язних компонент [2]. Деякі зв'язні компоненти відповідають символам машинозчитуваної зони, деякі – іншим символам на сторінці та, грубо кажучи, шумам (несуттєві для даної задачі дані, що можуть завадити процесу її розв'язання). Для кожної компоненти зв'язності виконується пошук мінімального охоплюючого кола [3]. Оскільки розміри символів МЗЗ є відомими, відомими є приблизні розміри відповідних кіл, завдяки чому можна провести попередню фільтрацію знайдених компонент зв'язності.

Для пошуку структури, що відповідає рядкам МЗЗ, будується триангуляція Делоне [4] на центрах знайдених кіл. Серед ребер триангуляції присутні ті, що сполучають кола компонент зв'язності МЗЗ. Оскільки відстані між символами є фіксованими, відстань

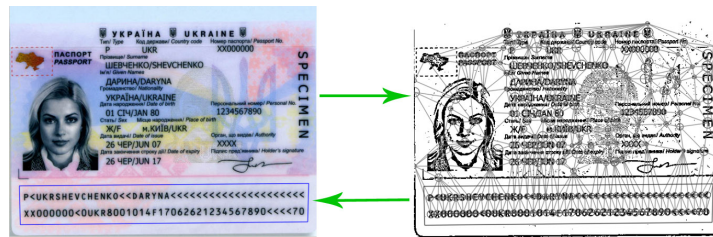


Рис. 5. Вхідне та оброблене зображення. Для відфільтрованих компонент зв'язності зображені їх мінімальні охоплюючі кола. На їх центрах побудовано триангуляцію Делоне.

між центрами кіл знаходиться в певних межах. Серед всіх ребер графа триангуляції відбираються лише ті, довжина яких знаходиться в певних межах. В наслідок цього граф триангуляції розпадається на компоненти зв'язності. Експериментально встановлено, що вершини графа, які відповідають символам МЗЗ, знаходяться в одній компоненті зв'язності (з високою ймовірністю) та утворюють ланцюг. Таким чином, задача зводиться до пошуку двох таких ланцюгів довжиною 44 елементи (кількість символів в рядку МЗЗ), що робиться тривіальним чином. Координати чотирикутника, що містить символи МЗЗ визначається координатами крайніх символів ланцюгів.

Розпізнавання символів машинозчитуваної зони

Після знаходження чотирикутника, що містить МЗЗ, відповідна частина оригінального зображення копіюється. Перед розпізнаванням символів, скопійована область бінаризується та з неї видаляються шуми шляхом пошуку та видалення зв'язних компонент невідповідного розміру. Для безпосереднього розпізнавання тексту на зображенні використано вільну бібліотеку Tesseract, що була додатково навчена на обмеженому наборі символів набраних шрифтом, зазначеним в [1]. Після розпізнавання, у випадку, якщо всі контрольні суми співпадають, виконується розбір МЗЗ на складові компоненти за допомогою регулярних виразів згідно того ж стандарту.

Висновки

Було розроблено та реалізовано мовою C++ алгоритм локалізації зони машинного зчитування паспорту з використанням примітивів вільної бібліотеки OpenCV та подальше розпізнавання за допомогою вільної OCR бібліотеки Tesseract. Алгоритм працює в режимі реального часу та є досить стійким щодо нерівномірного освітлення та зашумленості зображення. Якість розпізнавання можна ще дещо покращити шляхом навчання бібліотеки Tesseract на більшій кількості еталонних символів з необхідного набору.

Список літератури

- [1] Машиночитываемые проездные документы. Часть 1. Том 1. Издание шестое. — Международная ассоциация гражданской авиации, 2006.
- [2] Gonzales R.C., Woods R. E. Digital image processing, 2nd edition — Prentice Hall, 2001.
- [3] Bradsky G., Kaehler A. Learning OpenCV, 1st edition — O'Reilly, 2001.
- [4] de Berg M., van Creveld M., Overmars M., Schwarzkopf O. Computational Geometry. Algorithms and Applications. Second, revised edition — Springer, 2000

Автори

Віталій Ігорович Уманський — магістр 2-го року навчання, кафедра математичної інформатики, факультет кібернетики, Київський національний університет імені Тараса Шевченка, Київ, Україна; E-mail: vietal@list.ru