

Д.А. Беляков
Научный руководитель: ст. преп. кафедры ИС, Е.Е. Канунова
Муромский институт (филиал) Владимирского государственного университета
Владимирская обл., г. Муром, ул. Орловская, д.23
E-mail: cwwc@bk.ru

Разработка программы реализации алгоритмов распознавания символов на изображениях архивных текстовых документах

Распознавание – одна из главных задач программирования, затрагивающих машинное зрение. Несмотря на то, что в настоящее время большинство документов составляется на компьютерах, задача создания полностью электронного документооборота ещё далека до полной реализации. Как правило, существующие системы охватывают деятельность отдельных организаций, а обмен данными между организациями осуществляется с помощью традиционных бумажных документов.

Задача перевода информации с бумажных на электронные носители актуальна не только в рамках потребностей, возникающих в системах документооборота. Современные информационные технологии позволяют нам существенно упростить доступ к информационным ресурсам, накопленным человечеством, при условии, что они будут переведены в электронный вид.

Наиболее простым и быстрым является сканирование документов с помощью сканеров. Результат работы является цифровое изображение документа – графический файл. Более предпочтительным, по сравнению с графическим, является текстовое представление информации. Этот вариант позволяет существенно сократить затраты на хранение и передачу информации, а также позволяет реализовать все возможные сценарии использования и анализа электронных документов. Поэтому наибольший интерес с практической точки зрения представляет именно перевод бумажных носителей в текстовый электронный документ.

Процесс формирования структуры представляет собой последовательность итераций. На каждом этапе объемная сцена делится на восемь одинаковых областей и производится анализ каждой области на принадлежность к одному из трех типов, смешанные области подвергаются дальнейшему разбиению.

1. Поступающее на вход системы изображение должно быть очищено от шума и приведено к виду, позволяющему эффективно выделять символы и распознавать их.

2. Система должна разбить изображение на блоки текста, основываясь на особенностях его выравнивания и распределения по нескольким колонкам.

3. Изображение с текстом должно быть разделено на изображения строк, а затем на изображения символов для того, чтобы в дальнейшем обработать каждый символ по отдельности. После данного шага разные системы распознавания работают по своим специфическим алгоритмам.

4. Изображение символа может обрабатываться целиком, для этого оно сравнивается с имеющимися шаблонами. Другим вариантом является выделение характеристик изображаемого символа: отбор характерных признаков, и классификация данных признаков по имеющимся в системе критериям. На выходе четвертого шага появляется возможный вариант буквы. Однако обычно системы на этом не останавливаются и продолжают работу на основе других методов, уточняя полученный результат.

5. Результат распознавания может быть не удовлетворительным. Для получения более хороших результатов в системе может быть встроен блок обучения. С помощью этого блока можно задать системе примеры начертания разных букв в данном шрифте. После процесса обучения предполагается лучшее качество распознавания текста. Система распознавания текста не всегда должна следовать всем описанным шагам, но основные действия процесса распознавания являются общими для любого алгоритма.

Алгоритм работы программы подробно описывается в докладе.

Литература

1. Терехин, А.В. Алгоритм формирования косоугольной проекции трехмерного объекта по модели окто-дерева / А.В. Терехин, С.В. Савичева // Алгоритмы, методы и системы обработки данных. – 2013. – № 3 (25). – С. 74 – 81.

Секция 32. Технологии обработки визуальной информации

2. Садыков, С.С. Технология формирования эталонов трехмерных объектов для их распознавания / С.С. Садыков, А.В. Терехин, А.О. Кравченко // Надежность и качество – 2012:тр. межд. симп. – Пенза: изд. ПГУ. – С. 373 – 376.
3. Терехин, А.В. Алгоритм формирования описания поверхности трехмерного объекта / А.В. Терехин, С.С. Садыков // Распознавание – 2015: сб. мат XII МНТК. – Курск, 2015 – С. 356 – 358.
4. Learning OpenCV // URL: <http://locv.ru/> (Датаобращения 04.02.2016).
5. OpenCV (Open source computer vision) // URL: <http://opencv.org/> (Датаобращения 04.02.2016).
6. Методы и алгоритмы цифровой обработки изображений /под редакцией С.С. Садыкова – Ташкент, УзНпо «Кибернетика» АН РУз, 1992 -296с.
7. Методы компьютерной обработки изображений / Под ред. В.А. Сойфера. – 2-е издание, испр. – М.: Физматлит, 2003. – 784 с.