

Методы индексирования и поиска изображений и видеоданных на основании визуального содержания

Байгарова Н.С., Бухштаб Ю.А., Горный А.А., Евтеева Н.Н.,
Лялин В.Ю., Монастырский А.В., Стрелков А.Ю.
Институт прикладной математики им. М.В.Келдыша РАН
E-mail: kikom@glasnet.ru

1. Введение

Исследования направлены на развитие информационных технологий, применяемых в электронных библиотеках, ориентированных на работу с изображениями и видеоматериалами, а именно на разработку и реализацию методов анализа, индексирования и поиска изображений и видеоданных на основании визуальных атрибутов.

Оцифровка и хранение больших объемов изображений и видеоданных в настоящее время не представляет проблему с технической точки зрения, чего нельзя пока сказать про поиск релевантной визуальной информации. До последнего времени традиционным способом являлся поиск визуальной информации, опирающийся на индексирование текстовых строк, ассоциированных с изображением или фильмом. Авторы доклада разработали и используют для представления коллекции изображений специализированную полнотекстную поисковую машину, которая эффективно функционирует как на CD, так и в Интернет. Однако поиск по названию, авторам, теме, словам описания содержания и по другой текстовой информации, ассоциированной с изображениями коллекции, представляется недостаточным. Неоднозначность соответствия между визуальным содержанием и текстовым описанием снижает показатели точности и полноты поиска. Некоторые изображения вообще трудно описать словами (очевидный пример - абстрактные картины).

В связи с этим решается проблема обеспечения доступа к современным электронным коллекциям изображений с использованием различных средств - как текстовых описаний, так и характеристик визуального содержания, простейших типа цветовой гаммы, и более сложных, связанных с распознаванием образов, наиболее интересных для предметной области.

2. Визуальные примитивы и механизм поиска по образцу

Для организации электронных библиотек, связанных с визуальными данными, требуются методы создания и использования поисковых образов, отражающих визуальное содержание изображений. Методы распознавания образов и понимания сцены в настоящее время из-за отсутствия эффективных универсальных алгоритмов применяются в узких предметных областях. Современная универсальная технология доступа к коллекциям изображений связана с сопоставлением изображению набора визуальных примитивов (характеристик цвета, формы, текстуры, а для видео еще и параметров движения сцены и объектов) и определением количественной оценки близости изображений по значениям примитивов [6, 12, 13, 14].

Визуальные примитивы - это характеристики изображения, которые автоматически вычисляются по оцифрованным визуальным данным, позволяют эффективно индексировать их и обрабатывать запросы с использованием визуальных свойств изображения. Поисковый образ изображения, сгенерированный из визуальных примитивов, невелик по размеру в сравнении с самим изображением и удобен для организации поиска. Вычисление подобия изображений заменяет принятую в традиционных СУБД операцию установления соответствия запросу. Хотя запросом в такой системе может быть описание набора примитивов, более удобен запросный механизм поиска по образцу, когда система отыскивает изображения, визуально похожие на предоставленный образец. Система анализирует образец аналогично тому, как это делается при составлении поисковых образов изображений базы. Вычисление подобия изображения-образца изображениям коллекции осуществляется на основании сравнения значений отдельных визуальных примитивов, при этом система определяет меру их отличия, а затем сортирует изображения базы в соответствии с близостью к образцу по всем параметрам, с учетом указываемой в запросе степени важности каждого параметра. Поиск на таком уровне абстракции не предпо-

© Вторая Всероссийская научная конференция
ЭЛЕКТРОННЫЕ БИБЛИОТЕКИ:
ПЕРСПЕКТИВНЫЕ МЕТОДЫ И ТЕХНОЛОГИИ,
ЭЛЕКТРОННЫЕ КОЛЛЕКЦИИ
26-28 сентября 2000г., Протвино

лагает идентификацию объектов. Скажем, если в качестве образца взято изображение собаки, то система будет искать изображения, похожие на образец по цветовой гамме, композиции, наличию определенных форм и т.п., но нет никакой гарантии, что среди них окажется изображение именно этого животного. Тем не менее, метод поиска по образцу на основании визуальных примитивов представляется на сегодняшний день достаточно эффективным и универсальным средством доступа к коллекциям оцифрованных изображений.

3. Методы анализа изображений

Различными группами исследователей уже накоплен определенный опыт реализации алгоритмов, позволяющих автоматически описывать изображения в терминах простых вычислимых визуальных свойств, а также определять меру их отличия. Авторами доклада был подготовлен обзор этих алгоритмов [7].

Наши текущие исследования в этой области направлены на дальнейшее развитие методов вычисления и сравнения визуальных примитивов. Реализован метод количественной оценки близости статичных изображений по их цветовым гистограммам. Решена задача пространственного сегментирования изображения. Разработан и реализован алгоритм, осуществляющий вычисление параметров форм для выделенных объектов картинке и сравнение форм по их параметрам. Проводятся работы и имеются результаты, которые позволят выполнять локальное индексирование, отражающее распределение на изображении цветовых множеств. С целью вычисления измерений текстур исследуются возможности использования метода функций Габора и характеристик матрицы взаиморасположения оттенков серого цвета.

3.1. Цветовые гистограммы

Метод цветовых гистограмм – наиболее популярный из методов, использующих цветовые характеристики для индексирования изображений. Возможно также использование таких показателей, как средний или основной цвета, а также множества цветов; эти характеристики имеет смысл использовать для локального индексирования областей изображения [9, 17].

Идея метода цветовых гистограмм для индексирования и сравнения изображений сводится к следующему. Все множество цветов разбивается на набор непересекающихся, полностью покрывающих его подмножеств V_i , $0 \leq i < N$. Будем называть такое разбиение множества цветов базовой палитрой. Для изображения формируется гистограмма, отражающая долю каждого подмножества цветов в общей цветовой гамме изображения - массив $H[i] = N[i] / \sum N[i]$, где $N[i]$ - число точек с цветом из множества V_i . Для сравнения гистограмм вводится понятие расстояния между ними. Известны различные способы построения и сравнения цветовых гистограмм [1, 2, 7, 17], отличающиеся между собой изначальной цветовой схемой (RGB, CMY, HSV, grayscale и т. д.), размерностью гистограммы и определением расстояния между гистограммами.

В данной работе реализовано несколько модификаций метода, применяющих разные способы квантования множества цветов и вычисления расстояния между гистограммами. Используются две базовые палитры и, следовательно, два метода построения гистограммы.

1) Разбиение RGB-цветов по яркости.

В базовой палитре V_i ($0 \leq i < N$) определяется как множество цветов C , интенсивность которых $I(C)$: $256i/N \leq I(C) < 256(i+1)/N$. Интенсивность вычисляется по классической формуле: $I(C) = 0.3 R(C) + 0.59 G(C) + 0.11 B(C)$, где R , G и B – красная, зеленая и синяя компоненты цвета C , $0 \leq I(C) < 256$.

В частности, для черно-белых полутоновых изображений на N подмножеств разбивается исходное множество оттенков. Значение N выбиралось практически произвольно, сейчас установлено $N=16$.

В качестве расстояния между гистограммами используется сумма модулей разности соответствующих элементов гистограмм. Некоторое усовершенствование метода достигается при вычислении расстояния на основании поэлементного сравнения гистограмм с учетом соседних элементов. Для каждого элемента гистограммы первого изображения вычисляется не одна, а три разности:

$$R1[i] = |H1[i] - H2[i-1]|$$

$$R2[i] = |H1[i] - H2[i]|$$

$$R3[i] = |H1[i] - H2[i+1]| \text{ (для } i=0 \text{ и } i=N \text{ вместо невычислимых разностей подставляются заведомо большие значения), итоговое же расстояние равно:}$$

$$S = \sum_{i=0}^{N-1} \min_{1 \leq k \leq 3} R_k[i]$$

Этот способ не годится для произвольной базовой палитры, т. к. предполагает строгую упорядоченность множества цветов, как в случае с разбиением по яркости. Заметим, что так определенное S не является расстоянием в математическом смысле из-за несимметричности (нельзя гарантировать, что $S(H1,H2)=S(H2,H1)$). Основное преимущество алгоритма состоит в том, что он слабо чувствителен к изменению освещенности, что ощутимо улучшает результаты его применения на широком классе изображений.

Этот метод построения гистограмм наиболее эффективен для черно-белых полутоновых изображений. Для цветных RGB-изображений лучшие результаты дает другой способ.

2) Разбиение RGB-цветов по прямоугольным параллелепипедам.

Цветовое RGB-пространство рассматривается как трехмерный куб, каждая ось которого соответствует одному из трех основных цветов (красному, зеленому или синему), деления на осях пронумерованы от 0 до 255 (большее значение соответствует большей интенсивности цвета). При таком рассмотрении любой цвет RGB-изображения может быть представлен точкой куба. Для построения цветовой гистограммы каждая сторона делится на $n=4$ равных интервалов, соответственно RGB-куб делится на $N=64$ прямоугольных параллелепипедов. V_i – множество цветов, все компоненты которых попадают в определенные интервалы. Гистограмма изображения отражает распределение точек RGB-пространства, соответствующих цветам пикселей изображения, по параллелепипедам.

В качестве расстояния между гистограммами используется покомпонентная сумма модулей разности между ними. Несмотря на предельную простоту подхода, он показывает довольно стабильные результаты. Распознаются схожие по цветовой гамме серии картинок, если они имеются в базе. (В качестве тестовой базы была использована коллекция абсолютно разных фотографий, представленная на CD РИА «Новости».)

Более точное сравнение изображений достигается с помощью техники квадродеревьев, когда методы вычисления и сравнения цветных гистограмм применяются не ко всему изображению, а к его четверти (одной шестнадцатой и т. д.). Сейчас программа позволяет работать не только с полными изображениями, но и с их разбиением на четверти. Для реализации этой возможности при построении гистограммы автоматически считаются и гистограммы всех четырех квадрантов изображения. Сравнение изображений основывается на расстоянии, определенном как Евклидово в пространстве расстояний между гистограммами их частей - вместо вычисления расстояния между полными гистограммами, рассчитываются расстояния между четвертями, итоговым результатом считается корень из суммы их квадратов. Этот метод дает результат, семантически отличный от других вариантов: изображения, различающиеся только по взаимному расположению идентичных по цвету объектов, считаются непохожими, в то время как могли быть определены как близкие без использования этой техники. Целесообразность ее применения определяется значением для пользователя расположения на картинке-образце определенных цветных областей.

3.2. Объекты изображения

Пространственное сегментирование изображения может осуществляться автоматически, когда выделяются области с некоторыми общими свойствами - одинаковыми или сильно схожими значениями того или иного примитива. Полученные в результате области характеризуются расположением на картинке, размерами, значениями примитивов. Например, определение объекта изображения на основании близких значений интенсивности соседних точек позволяет довольно точно характеризовать выделенную область с помощью такого показателя, как средний цвет ее точек. Таким образом, для выделенных объектов могут быть определены и включены в индекс такие характеристики, как координаты на изображении, размеры, характеристики цвета, измерения формы и текстуры.

Контур - граница объекта - представляет собой замкнутую последовательность точек (x_s, y_s) , где $1 \leq s \leq N$. Удобно считать, что $(x_{s+N}, y_{s+N}) = (x_s, y_s)$. Задача выявления контуров связана с локализацией на изображении резких перепадов яркости цвета или изменений параметров, характеризующих текстуру.

Определение границ объектов изображения выполняется нами по следующей схеме: цветное изображение переводится в черно-белое полутоновое и сглаживается, осуществляется пространственное дифференцирование - вычисляется градиент функции интенсивности в каждой точке изображения и, наконец, подавляются значения меньше установленного порога. За основу взят метод Собеля [18], использующий для вычисления градиента интенсивности специальные ядра, известные как «операторы Собеля».

Ядра применяются к каждому пикселу изображения: он помещается в центр ядра, и значения интенсивности в соседних точках умножаются на соответствующие коэффициенты ядра, после чего полученные значения суммируются. X- оператор Собеля, примененный к 3×3 матрице исходного изображения, дает величину горизонтальной составляющей градиента интенсивности в центральной точке этой матрицы, а Y-оператор Собеля дает величину вертикальной составляющей градиента. Коэффициенты ядра выбраны так, чтобы при его применении одновременно выполнялось сглаживание в одном направлении и вычисление пространственной производной - в другом.

Величина градиента определяется как квадратный корень из суммы квадратов значений горизонтальной и вертикальной составляющих градиента.

В результате образуется массив чисел $G(i, j)$, характеризующих изменения яркости в различных точках изображения. Затем выполняется операция сравнения с порогом и определяется положение элементов изображения с наиболее сильными перепадами яркости. Выбор порога является одним из ключевых вопросов выделения перепадов. В нашей реализации он отличается от оригинального метода Собеля. В качестве основного порога берется средняя для изображения величина градиента - $Smid$. Для большого изображения с малым числом точек, обладающих сильным перепадом яркости, данной пороговой величины недостаточно, оказывается весьма сильным влияние шума. Для ликвидации этой проблемы для каждой точки изображения считается величина $Slocal$, равная средней величине градиента в области 3×3 вокруг анализируемой точки. Пороговое условие выглядит так:

$$(G(i, j) \geq Smid) \text{ AND } (Slocal \geq Smid)$$

В результате обработки получается бинарная матрица, где единицам соответствуют точки со значительным перепадом яркости, нулям – все остальные. В качестве дополнительной меры в борьбе с шумом и ликвидации возможных разрывов в контурах применяются морфологические операции.

Следующий этап – сегментирование изображения. Целью является выделение на изображении контуров объектов. В бинарной матрице единицами представлены точки, принадлежащие искусственно утолщенным на предыдущем этапе границам объектов. Для выделения границы одного объекта в матрице по определенному алгоритму ищется элемент, равный единице, не отнесенный ранее ни к какому другому объекту; далее считается, что все соседние элементы, равные единице, также принадлежат этому объекту; и т. д. Для выделения точек внешнего контура используется обход полученного объекта по внешней его стороне, начиная с нижней левой точки объекта и заканчивая ею же. Обход точек ведется последовательно против часовой стрелки. В результате получаем массив точек, образующих замкнутый контур объекта. Небольшие объекты при этом исключаются из рассмотрения.

3.3. Характеристики формы

Существует практика использования формы объектов для индексирования изображений с целью их дальнейшего сравнения [2, 7, 20].

В данной работе для вычисления предназначенных для индексирования характеристик формы из контура объекта выбирается 128 точек (x_s, y_s) , $1 \leq s \leq 128$.

Вводятся две функции:

- 1) Функция расстояния от точек контура до центра фигуры

$$R(s)^2 = (x_c - x_s)^2 + (y_c - y_s)^2$$

где (x_c, y_c) - центр масс контура:

$$x_c = \sum_{s=1}^{128} x_s / 128, \quad y_c = \sum_{s=1}^{128} y_s / 128$$

- 2) Углы поворота:

$$Y(s) = \arccos((a^2 + b^2 - c^2) / 2ab)$$

$$a^2(s) = (x_{s-1} - x_s)^2 + (y_{s-1} - y_s)^2$$

где

$$b^2(s) = (x_{s+1} - x_s)^2 + (y_{s+1} - y_s)^2$$

$$c^2(s) = (x_{s-1} - x_{s+1})^2 + (y_{s-1} - y_{s+1})^2$$

Для обеспечения инвариантности относительно поворотов и масштаба, выполняется нормирование величин, а в качестве начальной точки контура берется та, расстояние до центра от которой наименьшее, соответственно упорядочиваются элементы векторов.

Предлагается способ сравнения форм объектов на основании вычисления общего расстояния между двумя парами соответствующих векторов - используется покомпонентная сумма модулей разности.

Реализующая данный метод программа показала приемлемые результаты при поиске изображений базы, содержащих объекты с формами, похожими на форму объекта изображения-образца. (Тестирование проводилось на примерах, взятых из дистрибутива пакета Macromedia Director.)

4. Методы анализа видеоданных

4.1. Временное сегментирование видеофильма

В связи с большим объемом видео-файлов для организации эффективного поиска данных с удовлетворительными показателями полноты и точности, а также для обеспечения быстрого предоставления пользователю релевантной информации имеет смысл индексировать каждый фильм не как единое целое, а как последовательность логически самостоятельных частей — видеофрагментов [13]. Задача сводится к определению границ видеофрагментов, они могут быть связаны с точками монтажа, изменением положения снимающей камеры и т.п. Формально задачу можно поставить так: на вход подается упорядоченный набор кадров, необходимо выделить из них последовательность номеров, каждый из которых соответствует началу нового фрагмента.

Временное сегментирование может выполняться путем автоматического анализа изображения, соответствующие приемы известны [3, 5, 10]. Достаточно эффективны для выделения кадров, на которых происходит значительное изменение видеоизображения, методы, базирующиеся на вычислении низкоуровневых характеристик изображения.

Предлагается алгоритм, основанный на сравнении цветовых гистограмм соседних кадров. Система вычисляет цветовую гистограмму очередного кадра и сравнивает с гистограммой предыдущего кадра. Построение и сравнение гистограмм осуществляются идентично работе со статичными изображениями. Есть дополнительная возможность “выравнивания” гистограмм, применение которой имеет смысл только в случае разбиения множества цветов по яркости. Целью ее является приведение гистограммы к виду, при котором верно:

Потребность такой обработки объясняется тем, что во многих реальных образцах видео (особенно черно-белых и/или плохого качества) часто встречаются практически идентичные соседние кадры, отличающиеся только по яркости среднего освещения. Семантически они должны попадать в один фрагмент, но их первоначальные гистограммы

$$\sum_{i=0}^{N/2-1} H[i] = \sum_{i=N/2}^{N-1} H[i]$$

могут существенно различаться, что и приводит к необходимости выравнивания. Прием не дает полного решения проблемы, но, по крайней мере, существенно улучшает результаты на широком классе изображений.

Результатом сравнения гистограмм последовательных кадров является массив чисел, где i -ый компонент — расстояние между гистограммами i -ого и $(i+1)$ -ого кадров. Опираясь на эти данные, фильм разбивается на фрагменты. Граница фрагментов считается обнаруженной, если разница гистограмм между рассматриваемыми кадрами выше некоторого абсолютного порога L и одновременно в K раз превышает среднее значение разницы гистограмм соседних кадров, посчитанное от начала выделяемого фрагмента до текущего кадра (относительный порог).

Попытка ограничиться абсолютным порогом не привела к успеху. Значение порога, дающее желаемые результаты, существенно зависит от качества записи фильма, динамики его связанных фрагментов и средней освещенности. Эти характеристики могут быть различны в разных фрагментах одного и того же фильма (т. е. для их вычисления потребуются уже готовые результаты временной сегментации); кроме того, определение, например, качества записи — трудная задача, и зависимость от нее неизбежно внесет дополнительную погрешность в итоговый результат. Аналогично, не удалось решить задачу только с помощью относительного порога, используя среднее значение разницы гистограмм соседних кадров, посчитанное по кадрам от начала выделяемого фрагмента. Проблема заключается в возможности практически полного совпадения первых нескольких гистограмм фрагмента и минимального отличия от них следующей — условие порога будет выполнено, и программа выдаст лишний фрагмент. Наиболее эффективной оказалась комбинация этих подходов, т. е. $(i+1)$ -й кадр считается началом нового фрагмента, если расстояние между гистограммами i -ого и $(i+1)$ -ого кадров превышает оба порога.

Результаты сегментирования, разумеется, сильно зависят от выбора параметров. Они установлены эмпирически для достижения приемлемых результатов с точки зрения минимизации числа ошибок, связанных с обнаружением ложной границы и пропуском действительной. (При уменьшении вероятности одной из ошибок, неизбежно повышается вероятность другой.) Текущие значения $L=0.15$ и $K=3$ подобраны так, чтобы ложные обнаружения встречались примерно на порядок чаще, чем пропуск переходов. Вызвано это тем, что на “двойном” видеофрагменте невозможно корректно вычислить оптический поток, что является одной из главных целей временного сегментирования, в то время как два фрагмента вместо одного дают лишь некоторое увеличение требуемых для дальнейшей обработки ресурсов.

Для тестирования программы использовались видеофильмы, взятые из различных источников, разного качества. Все фильмы разбиты на несколько непересекающихся групп, результаты внутри которых были более или менее одинаковы. Приводим усредненные результаты, отражающие процент ошибок, допускаемых различными методами временной сегментации на различных типах видеофрагментов.

Ложные обнаружения

Метод	(1)	(2)	(3)	(4)	(5)
Мультфильмы	10%	10%	31%	18%	0%
Цветное видео	33%	11%	33%	33%	27%
Черно-белое видео	23%	-	25%	-	12%

Пропущенные границы фрагментов

Метод	(1)	(2)	(3)	(4)	(5)
Мультфильмы	0%	0%	0%	0%	0%
Цветное видео	0%	0%	0%	0%	0%
Черно-белое видео	4%	-	2%	-	10%

Методы:

1 - Палитра разбивается по интенсивности.

2 - Палитра разбивается на RGB-параллелепипеды.

3 – Квадродерево & палитра разбивается по интенсивности.

4 – Квадродерево & палитра разбивается на RGB-параллелепипеды.

5 - Палитра разбивается по интенсивности, применяется усложненная формула расстояния и выравнивание гистограмм.

(Процент ошибок – отношение числа ошибок к сумме ошибок и верно обнаруженных переходов, умноженное на 100%.)

4.2. Индексирование видеофрагментов

После того как видеопоток разбивается на фрагменты, из них выделяются для исследования ключевые стоп-кадры. Стратегия извлечения представительных стоп-кадров из каждого выделенного фрагмента может быть, например, такой [1]: если фрагмент короче секунды, берется один центральный кадр, для более длинных фрагментов берется по одному в секунду. Для каждого выделенного кадра вычисляются с целью индексирования визуальные примитивы: цветовые гистограммы, характеристики формы и цвета объектов изображения, измерения текстуры; для этого применяются те же методы, что и для анализа статичных изображений. Кроме того, представляется важным индексировать фрагмент также характеристиками движения камеры/сцены и движения объектов, определяемыми на основании совокупности кадров видеофрагмента [3, 5].

4.3. Вычисление оптического потока

Для индексирования видеоданных по движению применяется метод оптического потока. Он основан на том, что для видеофрагмента, содержащего некоторые объекты в движении, можно вычислить направление и величину скорости движения в каждой точке видеокadra. Известны разные алгоритмы вычисления оптического потока [8].

В данной работе реализован дифференциальный метод расчета оптического потока, который предполагает вычисление пространственно-временных производных интенсивности. Для того чтобы повысить точность вычислений, все кадры видеофрагмента предварительно сглаживаются с помощью фильтра Гаусса. Предварительно осуществляется выделение кадров из видеофрагмента, цветные изображения преобразуются в черно-белые полутоновые. Авторы работы [8] предлагают применять, кроме пространственного, также временное сглаживание, а именно использовать для расчета интенсивности в точке кадра значения в близких точках соседних кадров. К сожалению, эта техника дает приемлемые результаты только в том случае, если скорость движения объектов на видеофрагменте не превышает одного-двух пикселей на кадр, а это условие далеко не всегда выполняется для реальных видеофрагментов.

Дифференциальная техника вычисления скорости в каждой точке опирается на простое правило: при движении объекта интенсивность составляющих его точек не изменяется:

$$I(\mathbf{x} + \mathbf{v}dt, t + dt) = I(\mathbf{x}, t), \text{ или } \frac{dI(\mathbf{x}, t)}{dt} = 0,$$

где $\mathbf{v} = (u, v)^t$ — скорость точки $\mathbf{x} = (x, y)^t$, $I(\mathbf{x}, t)$ — интенсивность в точке \mathbf{x} в момент времени t .

Это правило дает одно линейное уравнение для двухкомпонентного вектора скорости:
 $\nabla I(\mathbf{x}, t) \cdot \mathbf{v} + I_t(\mathbf{x}, t) = 0,$

где $\nabla I(\mathbf{x}, t) = \left(\frac{\partial I(\mathbf{x}, t)}{\partial x}, \frac{\partial I(\mathbf{x}, t)}{\partial y} \right)^t$, $\nabla I(\mathbf{x}, t) \cdot \mathbf{v}$ — скалярное произведение векторов.

Дополнительные ограничения могут быть получены различными способами. Например, реализованный в настоящей работе метод минимума градиента (метод Хорна) [11], предполагает гладкое изменение значения скорости от точки к точке: $\|\nabla u(\mathbf{x})\|^2 + \|\nabla v(\mathbf{x})\|^2 = 0$, или $\left(\frac{\partial u}{\partial x} \right)^2 + \left(\frac{\partial u}{\partial y} \right)^2 + \left(\frac{\partial v}{\partial x} \right)^2 + \left(\frac{\partial v}{\partial y} \right)^2 = 0$.

Итак, для определения вектора скорости необходимо минимизировать интеграл

$$\iint_D \left[(\nabla I(\mathbf{x}, t) \cdot \mathbf{v} + I_t(\mathbf{x}, t))^2 + \alpha^2 (\|\nabla u(\mathbf{x})\|^2 + \|\nabla v(\mathbf{x})\|^2) \right] d\mathbf{x}, \quad (1)$$

где D — множество точек изображения, а α — параметр, определяющий вес слагаемого, отвечающего за гладкость оптического потока. В настоящей реализации используется значение $\alpha = 100$, которое в процессе тестирования оказалось оптимальным.

В работе [11] показано, что задача минимизации интеграла (1) сводится к решению следующей системы равенств:

$$\begin{cases} (\alpha^2 + I_x^2 + I_y^2)(u - \bar{u}) = -I_x(I_x \bar{u} + I_y \bar{v} + I_t) \\ (\alpha^2 + I_x^2 + I_y^2)(v - \bar{v}) = -I_y(I_x \bar{u} + I_y \bar{v} + I_t) \end{cases} \quad (2)$$

где \bar{u} и \bar{v} — локальное среднее, соответственно, для u , v .

Авторы метода предлагают решать полученную систему линейных уравнений итерационным методом, например методом Гаусса — Зейделя. Тогда

$$\begin{cases} u_{n+1}(j, i) = \bar{u}_n(j, i) - I_x(j, i) \frac{(I_x(j, i)\bar{u}_n(j, i) + I_y(j, i)\bar{v}_n(j, i) + I_t(j, i))}{(\alpha^2 + I_x^2(j, i) + I_y^2(j, i))}, \\ v_{n+1}(j, i) = \bar{v}_n(j, i) - I_y(j, i) \frac{(I_x(j, i)\bar{u}_n(j, i) + I_y(j, i)\bar{v}_n(j, i) + I_t(j, i))}{(\alpha^2 + I_x^2(j, i) + I_y^2(j, i))}, \end{cases} \quad (3)$$

где $\mathbf{v}_k(j, i) = (u_k(j, i), v_k(j, i))$ — значение вектора скорости в соответствующей точке кадра на k -й итерации; $\mathbf{v}_0 = 0$. Итерация связана с временным шагом. Итеративное вычисление оптического потока выполняется на основании всех кадров фрагмента для более точного определения скоростей. В отличие от оригинального метода [11], для вычисления производных интенсивности по времени и по каждой из координат, вместо двухточечной схемы, применялась центральная 5-точечная разностная схема с коэффициентами $\frac{1}{12}(-1, 8, 0, -8, 1)$.

Полученный оптический поток используется затем для определения более высокоуровневых характеристик, связанных с движением, которые предназначены для индексирования и поиска видеоданных.

4.4. Выделение движущихся объектов

Применяемый для вычисления оптического потока метод позволяет, выполнив лишь одну итерацию, правильно определить нормальную составляющую вектора скорости на границе объекта изображения. Затем, по мере увеличения числа итераций, значение на границе приближается к реальному значению скорости. С другой стороны, согласно формуле (3) в точках изображения с нулевым градиентом интенсивности скорость будет определяться как среднее значение скоростей соседних точек. Следовательно, с увеличением числа итераций размер и форма объектов будут искажаться. Кроме того, ненулевые значения вектора скорости получают все точки изображения, принадлежавшие объекту хотя бы на одном кадре последовательности.

Значит, чем больше число итераций, тем более правильно определяются величины скоростей, но в то же время неточно воспроизводится форма объекта. Поэтому для выделения движущихся объектов осуществляется вычисление потока отдельно по последним пяти кадрам без учета предыдущих, с последующей подстановкой в точках с ненулевыми скоростями результатов итеративно вычисленного по всем кадрам оптического потока, как более адекватных. В полученном оптическом потоке будут представлены как форма движущихся объектов (если таковые присутствуют на видеофрагменте), так и достаточно точные значения скоростей.

Информация об оптическом потоке используется для пространственного сегментирования изображения. Группу расположенных близко друг от друга точек, движущихся с примерно одинаковыми скоростями (или хотя бы приблизительно однонаправленными), можно считать движущимся объектом. Выделение областей происходит по схеме: две точки считаются принадлежащими одному объекту, если они отстоят друг от друга не более чем на 3 пиксела и направления скоростей в них отличаются не более чем на 45° ; не принимаются во внимание области, размер которых пренебрежимо мал ($< 0.1\%$) либо слишком велик ($> 30\%$) относительно размера кадра (в последнем случае движение области учитывается при определении глобальных характеристик движения сцены/камеры).

Для выделенных областей рассматриваются минимальные охватывающие их прямоугольники. Помимо их размеров и расположения, определяется тип движения, по той же схеме, что и для картинки в целом (см. следующий раздел). Для этого вычисляются средние значения скорости в четырех квадрантах этого прямоугольника.

Затем вычисляется средний модуль скорости по всем точкам, принадлежащим объекту, для обеспечения возможности поиска видеофрагментов с требуемой интенсивностью движения объекта. Кроме того, вычислив оптический поток не только для последних кадров, но также для первых и для некоторых промежуточных, в случае поступательного движения можно определить характеристики траектории движения объектов. В качестве дальнейшего этапа исследований рассматривается задача определения происходящих с объектами событий.

После обработки объектов исследуются глобальные характеристики движения, для чего вычисляются средние значения вектора скорости в квадрантах изображения и средняя интенсивность движения сцены (без учета скоростей выделенных объектов).

4.5. Характеристики движения

После вычисления оптического потока в каждой точке видеофрагмента, возникает задача привести эти сложные данные к простой и пригодной для индексирования и поиска форме. В работе [1] излагается возможный способ решения проблемы.

В данной работе предлагаются новые характеристики, вычисляемые исходя из средних скоростей четвертей видеофрагмента. Алгоритмы могут применяться для вычисления как глобальных, так и локальных характеристик видео: они определяются отдельно для каждой прямоугольной области, содержащей движущийся объект, а также для всего изображения.

Создана многоуровневая классификация видеофрагментов по типу движения:

1) Идентификатор схемы движения

На основании средних скоростей в квадрантах анализируемого изображения выбирается наиболее близкая схема движения, определяющая для каждого квадранта одно из 8 основных направлений движения (с точностью до 45° градусов) или же отсутствие существенного движения. При определении схемы учитываются не только направления средних скоростей, но и соотношение их модулей. Так как в разных фрагментах интенсивность может сильно различаться, порог, значения ниже которого считаются нулевыми, не может быть установлен изначально. Сейчас принят следующий алгоритм его вычисления. Значения скоростей разбиты на интервалы экспоненциально растущего размера. Определяется интервал, в который попадает наибольшее значение модуля скорости, нулевыми же считаются все значения из остальных интервалов. В результате скорости каждого квадранта переводятся в целые числа от 0 до 8, что для четырех квадрантов дает 6561 комбинацию. Если назвать комбинацию идентификатором схемы, то и получится первый вид классификации: схожими будут считаться фрагменты с одинаковыми идентификаторами их схемы движения. Несмотря на простоту, данная методика соответствует семантике многих реальных запросов. В частности, распознаваемы характерные функции камеры (приближение, удаление, сдвиг). Например, постепенному переходу к крупному плану будет соответствовать движение от центра к углам во всех четвертях.

2) Доминирующее направление

Идея – разбиение всего множества фильмов на два класса: с выраженным общим направлением движения и без него. С человеческой точки зрения, в первый тип попадают фрагменты с крупным планом и движущимся центральным объектом, а также эпизоды, снятые движущейся камерой, или с движением фона. Во второй – все остальные фрагменты. Естественно, для фрагментов с доминирующим движением целесообразно хранить не только сам факт его наличия, но и направление, а также число квадрантов с этим направлением (мощность доминанты).

Вычисление производится по идентификатору схемы. Предварительно создается массив $D[i]$ ($1 \leq i \leq 8$), такой, что для каждого i $D[i]$ равно числу квадрантов с направлением i . Определение доминирующего направления

после этого сводится к нахождению такого i , что $D[i] \geq 3$ или $D[i] = 2$ и $\forall j \neq i \Rightarrow D[j] \leq 1$. (Значение $D[i]$ - мощность доминанты). Если удовлетворяющего этим условиям i не найдено, то доминирующего направления нет.

3) Мощность схемы определяется как количество квадрантов с не близкими к нулю скоростями.

4) Эквивалентность схем с точностью до поворотов

Поддерживается разбиение на набор классов, каждый из которых образован поворотом схемы с базовым идентификатором вокруг своей оси на 0 (базовая схема), $\pi/2$, π и $3\pi/2$. Семантически класс соответствует некому целостному движению, показанному с разных сторон. Полностью сохраняются все типичные видеоэффекты.

Реализация, организованная с использованием библиотек классов, обеспечивает возможность настройки и расширения классификации.

Предложенная классификация обеспечит разносторонний поиск видеофрагментов. При таком подходе запрос сможет задавать для искомого видеофрагмента частично или полностью определенную схему движения, наличие некоторого доминирующего направления, количество квадрантов с ненулевыми скоростями, а также степень интенсивности движения. Понятие эквивалентности схем с точностью до поворотов позволит определять в запросе относительную схему движения.

5. Распознавание лица

Пользователю электронной библиотеки изображений должна быть предоставлена возможность строить запросы с использованием различных визуальных средств - в терминах не только визуальных примитивов, но и высокоуровневых объектов. Для этого в поисковом образе должен отражаться факт присутствия на изображении объектов наиболее интересных классов, а также размеры и расположение на кадре этих объектов. Задача нахождения на изображении объектов в настоящее время не ставится глобально. Как правило, речь идет об объектах определенного класса, особенно интересных для рассматриваемой предметной области. В рамках данной работы решается задача локализации фронтального вида лица человека на неподвижных изображениях / стоп-кадрах с помощью нейронной сети. Использование большого количества положительных и отрицательных примеров для обучения классифицирующего механизма позволяет автоматически получить достаточно точную модель объекта. Примеры систем распознавания лица, использующих контролируемое обучение - [16, 19].

Разрабатываемая система применяется к черно-белым полутоновым изображениям.

Та часть системы, которая непосредственно определяет наличие или отсутствие лица на картинке, применяется к небольшой области изображения, размеры которой выбираются так, чтобы в этой области можно было бы свободно поместить неискаженное лицо, и в то же время чтобы вычисления не занимали много времени (20 на 20 пикселей). На выходе выдается число, близкое к "1", если эта часть картинки содержит лицо, и к "-1" в противном случае.

Чтобы определить наличие лица на изображении, описанный фильтр применяется к каждому его участку. Для определения лиц, размеры которых превосходят размеры входного изображения для фильтра, картинку уменьшают в размерах и снова к каждому участку полученного изображения применяют фильтр и так далее, пока размеры картинки не уменьшатся до размеров входного изображения фильтра. Механизм дает возможность находить лица разного размера.

Перед непосредственным применением фильтра его входная картинка проходит предварительную обработку, которая позволяет сети работать с изображениями независимо от их средней интенсивности и от влияния источника света на изображение.

Результатом выполнения этих шагов является множество областей, где обнаружены лица. Последний этап состоит в том, чтобы отбросить те области, где ошибочно обнаружены лица, и объединить те, в которых обнаружено одно и то же лицо несколько раз, учитывая процесс многократного масштабирования.

Таким образом, систему можно разделить на 2 подсистемы: 1) предварительная обработка и нейросетевой фильтр, который для каждой области исходного изображения размером 20 на 20 пикселей выдает ответ, подтверждающий или опровергающий факт наличия лица, и 2) арбитр, который отбрасывает ошибочно обнаруженные лица. На данный момент реализована первая подсистема, а также механизм обучения сети.

Нейронная сеть данной системы состоит из трех слоев - входного и двух уровней нейронов. Вектор примитивов $\bar{x} \in [-1; 1]^{400}$, поступающий на вход нейронной сети, является значением функции интенсивности пикселей изображения размером 20 на 20. Внутренний слой состоит из 26 узлов, чувствительных к определенным областям входного изображения. Области выбраны для облегчения нахождения характерных черт лица. Внешний уровень состоит из одного узла, выходное значение которого интерпретируется как наличие или отсутствие лица. Используются биполярные функции активации нейронов [15].

Обучение сети на примерах позволяет автоматически настроить ее параметры. На вход системы подается набор изображений с известным выходом (1/-1). Система настраивается так, чтобы данное изображение отображать в заданное число. Процесс обучения основан на градиентном методе (метод обратного распространения ошибки), ми-

нимизирующем функцию ошибки, которая характеризует степень отклонения реального выхода сети от желаемого, по всему обучающему множеству [15]. Обучающее множество состоит из двух подмножеств картинок. Одно содержит картинки с лицами, другое - картинки без лиц. Первое подмножество генерируется из реальных изображений, лица на которых различаются размерами, расположением, наклоном, интенсивностью освещения. На каждой картинке глаза помечаются вручную. Эти метки используются для приведения каждого лица к одному размеру, расположению и наклону. Нормализация выполняется путем поворота, масштабирования и сдвига исходного изображения так, чтобы метки глаз заняли предопределенное положение в окне 20 на 20 пикселей. В обучающий набор включаются как сами преобразованные изображения, так и изображения, полученные из них путем поворота на угол 5 градусов по часовой и против часовой стрелки, также включаются зеркальные отображения и повернутые зеркальные отображения на тот же угол. Это делает систему распознавания инвариантной к небольшим наклонам головы. Обучающий набор проходит предварительную обработку, для устранения влияния источника света на изображение. Набор отрицательных примеров генерируется случайным образом.

Чтобы определить момент окончания обучения, используется проверочное множество (25% от общего количества обучающих примеров). Несколько итераций система обучается на 75% обучающего множества, затем вычисляется суммарная квадратичная ошибка на проверочном множестве. Такая процедура выполняется до тех пор, пока ошибка проверочного множества не начнет расти - тогда обучение заканчивается. Такая процедура позволяет избежать настройки системы исключительно на обучающее множество.

Для обучения системы было подготовлено 500 изображений с лицами, из которых было сгенерировано 3000 примеров лиц. Из них 2250 использовались для обучения и 750 для проверки. Случайным образом было сгенерировано 375 изображений, не содержащих лица, которые использовались для обучения, и 125, которые составили проверочное множество. Предусмотрено повторное обучение системы на примерах, не содержащих лиц, если система ошибочно их выявляет.

Качество обучения проверяется на тестовом наборе задач, не пересекающемся с обучающей выборкой. После проведенного на данный момент обучения система показывает высокие показатели распознавания лиц, однако пока велик процент ошибочного их обнаружения. Это связано со сложностью подготовки всеобъемлющего множества отрицательных примеров, необходимостью использования для обучения большего количества примеров и с тем, что пока не применяется арбитражная система.

6. Заключение

Создается прототип системы, позволяющей осуществлять поиск изображений и видеоданных в электронных коллекциях на основании визуальных атрибутов. Нами рассматриваются вопросы реализации такой системы в контексте проекта создания электронной библиотеки «Кинолетопись России» на базе собрания документальных кинофильмов Российского государственного архива кинофотодокументов [4].

Реализация разработанных алгоритмов осуществлена для изображений в формате Windows device-independent bitmap - DIB (BMP). Видео-файлы обрабатываются в формате *audio-video interleaved* (AVI).

Список литературы

1. Ardizzone, E., La Cascia, M., and Molinelli, D., Motion and Color Based Video Indexing and Retrieval, Proc. Int. Conf. on Pattern Recognition, (ICPR-96), Wien, Austria, Aug. 1996.
<http://www.cs.bu.edu/associates/marco/publications.html>
2. Ardizzone, E., La Cascia, M., Vito di Gesu, and Valenti, C., Content Based Indexing of Image and Video Databases by Global and Shape Features, 1996.
<http://www.cs.edu/associates/marco/publications.html>
3. Н.С. Байгарова, Ю.А. Бухштаб
Некоторые принципы организации поиска видеоданных
Программирование, N 3, 1999, стр. 165-170

4. Н.С. Байгарова, Ю.А. Бухштаб
Проект «Кинолетопись России» : представление и поиск видеоинформации
I Всероссийская конференция «Электронные библиотеки», Санкт-Петербург, 1999, стр. 209-215
5. Н.С. Байгарова, Ю.А. Бухштаб, Н.Н. Евтеева
Организация электронной библиотеки видеоматериалов
Препринт Института прикладной математики им. М.В. Келдыша РАН, 2000, N 5
6. Н.С. Байгарова, Ю.А. Бухштаб, А.А. Воробьев, А.А. Горный
Организация управления базами визуальных данных
Препринт Института прикладной математики им. М.В. Келдыша РАН, 2000, N 6
7. Н.С. Байгарова, Ю.А. Бухштаб, А.А. Горный
Методы индексирования и поиска визуальных данных
Препринт Института прикладной математики им. М.В. Келдыша РАН, 2000, N 7
8. Baron, J. L., Fleet, D. J., and Beauchemin, S. S., Performances of optical flow techniques.
Int. Journal of Computer Vision, 12:1, pp.43—77, 1994
9. Carson, C. and Ogle, V.E.,
Storage and Retrieval of Feature Data for a Very Large Online Image Collection. 1996.
<http://elib.cs.berkeley.edu/papers/>
10. Chrictel, M., Stevens, S., Kanade, T., Mauldin, M., Reddy, R., and Wactlar, H.,
Techniques for the Creation and Exploration of Digital Video Libraries,
Multimedia Tools and Applications, Boston: Kluwer, 1996, vol. 2.
11. Horn, B.K.P. and B.G.Schunk, Determining optical flow, Artificial intelligence, 17,1981.
12. Jain, R. and Gupta, A., Computer Vision and Visual Information Retrieval, 1996
<http://vision.ucsd.edu/papers/rosenfeld/>
13. Jain, R. and Gupta, A., Visual Information Retrieval,
Communications of the ACM, 1997, vol. 40, no. 5.
14. Jain, R., Pentland, A.P., Petkovic, D.,
Workshop Report: NSF – ASPA Workshop on Visual Information Management Systems, 1995.
<http://www.virage.com/vim/vimsreport95.html>
15. Looney, C.G., Pattern Recognition Using Neural Networks,
Theory and Algorithms for Engineers and Scientists, Oxford University Press, 1997.
16. Rowley, H.A., Baluja, S., and Kanade, T., Neural Network-Based Face Detection,
IEEE Transactions on Pattern Analysis and Machine Intelligence, 1998.
and In Proceedings of International Conference on Computer Vision
and Pattern Recognition, pp. 203-208, San Francisco, CA.
17. Smith, J.R. and Shih-Fu Chang. Tools and Techniques for Color Image Retrieval. 1996.
<http://www.ctr.columbia.edu/~jrsmith/html/pubs/>
18. Sobel, I., An isotropic image gradient operator.
Machine Vision for Three-Dimensional Scenes, pp.376-379. Academic Press, 1990
19. Sung, K-K. and Poggio, T., Example-Based Learning for View-based Human Face Detection.
A.I. Memo No. 1521, December 1994.
20. Wei-Ying Ma, NETRA: A Toolbox for Navigating Large Image Databases. 1997.
<http://vivaldi.ece.ucsb.edu/Netra/>