

СОДЕРЖАНИЕ

ВВЕДЕНИЕ	8
1 АНАЛИЗ ПРЕДМЕТНОЙ ОБЛАСТИ	11
1.1 Распознавание объектов на изображениях	11
1.1.1 Исторический обзор	11
1.1.2 Сверточные нейронные сети	12
1.1.3 Методы для распознавания объектов	17
1.2 Архитектура процессоров ARM	23
1.2.1 Особенности	23
1.2.2 Наборы команд	25
СПИСОК ЛИТЕРАТУРЫ	28

ВВЕДЕНИЕ

С каждым днем мы все чаще сталкиваемся с устройствами Интернета вещей (Internet of Things, IoT), которые призваны упростить нашу повседневную жизнь. Концепция IoT подразумевает подключение различных устройств между собой, которые обмениваются данными как в рамках внутренней сети, так и внешней. Классическим примером данной концепции является ”умный” дом, когда, например, мы можем регулировать уровень освещенности в комнате с помощью голосового управления на телефоне.

Однако технологии Интернета вещей внедряются не только в нашу быденную жизнь, но и в промышленности. Так называемый промышленный Интернет вещей (Industrial Internet of Things, IIoT) уже сегодня позволяет повысить эффективность процессов производства. Типичный производственный объект может быть источником большого количества данных, собранных с датчиков, анализ которых позволит обнаружить и устранить проблему на ранних этапах ее появления. В качестве примера внедрения IIoT можно вспомнить компанию Amazon, которая активно применяет роботов и различные датчики на своих складах для более быстрого поиска нужных товаров [1].

На сегодняшний день устройства IoT и IIoT чаще всего применяются в связке с облачной инфраструктурой. То есть данные, собираемые устройствами отправляются в ”облако”, где дальше уже сам анализ происходит на машинах с большими вычислительными ресурсами. Такой подход больше применим для долговременной аналитики, но если результат обработки нужен ”здесь и сейчас”, то задержки при передачи данных от устройств до центров обработки данных (ЦОД) могут быть критичны. В связи с этим начали набирать популярность граничные вычисления (edge computing), когда обработка данных происходит рядом с конечными устройствами. Так как количество данных с каждым годом будет только расти, то можно утверждать,

что граничные вычисления будут находить все более широкое применение. Согласно исследованиям компании Gather, к 2025 году около 75% данных будет обрабатываться на границе, не доходя до ЦОД [2].

Часто для анализа данных собранных с устройств IoT и PoT применяются методы глубокого обучения, которые позволяют установить более сложные зависимости во входных данных. Глубокое обучение подразумевает построение искусственных нейронных сетей для решения прикладных задач. Методы глубокого обучения появились еще в 80-х годах, но начали стремительно набирать популярность около 2010 года. Обусловлено это несколькими причинами. Во-первых, с ростом Интернета увеличивалось количество данных, которые можно было использовать для обучения. Во-вторых, с развитием графических процессоров (Graphics Processing Unit, GPU) процесс обучения нейронных сетей сократился в десятки раз. Впоследствии предобученные модели могут запускаться на устройствах и без GPU, но при этом демонстрировать приемлемое качество работы.

Исторически задача распознавания объектов на изображениях являлась тяжелой, но с развитием методов глубокого обучения удалось повысить эффективность решения данной задачи. Сегодня предобученные модели для распознавания объектов могут быть запущены даже на компьютерах с ограниченными вычислительными ресурсами (КОВР). В данной работе под КОВР подразумевается одноплатный компьютер без каких-либо специальных модулей для аппаратного ускорения методов машинного обучения. Классическим примером КОВР является одноплатный компьютер Raspberry Pi. Чаще всего КОВР собирают данные с различных устройств, проводят их первичную предобработку и далее отправляют их для полного анализа в "облако". Однако при текущем развитии КОВР, их вычислительные ресурсы могут использоваться для куда более сложных вещей, таких как распознавания объектов на видеоизображениях. Вполне ожидаемо, что одному экземпляру КОВР решить такую задачу будет сложно, но при объединении нескольких

экземпляров КОВР в кластер можно добиться приемлемых результатов. Основными преимуществами использования кластера из КОВР являются:

- 1) снижение затрат на анализ видеоизображений. Во-первых, если на объекте уже имеются экземпляры КОВР, то для решения этой задачи может использоваться существующая инфраструктура. Во-вторых, развертывание кластера из КОВР требует меньше финансовых ресурсов по сравнению с долговременным использованием облачной инфраструктуры;
- 2) уменьшение задержек при передаче видеоизображений, что способствует увеличению отзывчивости систем, которые опираются на анализ этих данных;
- 3) повышение конфиденциальности собранных данных за счет их локальной обработки.

Приведенные выше утверждения говорят об **актуальности** данной работы. Главной **проблемой** на сегодняшний день является отсутствие в открытом доступе каких-либо полноценных решений для задачи распределенного распознавания объектов на видеоизображениях на кластере из КОВР. В связи с этим, **целью** данной работы является организация распределенного распознавания объектов на видеоизображениях на кластере из КОВР. Для достижения поставленной цели требуется решить следующие **задачи**:

- 1) провести обзор методов для распознавания объектов на изображениях;
- 2) провести анализ особенностей архитектуры КОВР;
- 3) разработать алгоритм для распределения вычислений, необходимых для распознавания объектов на видеоизображениях, между узлами кластера;
- 4) провести тестирование разработанного алгоритма на кластере из КОВР.

1 АНАЛИЗ ПРЕДМЕТНОЙ ОБЛАСТИ

1.1 Распознавание объектов на изображениях

1.1.1 Исторический обзор

Задачи распознавания объектов на изображениях являются основными в области компьютерного зрения. Разработанные методы могут применяться в различных сферах, начиная от идентификации пользователя по фотографии, заканчивая сложными системами видеонаблюдения и робототехникой. В целом, все существующие методы можно разделить на две категории: методы ручного выделения признаков (handcrafted features) и методы, основанные на нейронных сетях. При использовании методов ручного выделения, какие-то признаки изображений выделяются по заранее определенному алгоритму с учетом специфики данных. Примером таких признаков могут служить границы объектов. Так как с учетом специфики изображений мы знаем, что если значения соседних пикселей сильно различаются, то это может свидетельствовать о том, что в этом месте проходит граница между двумя объектами. В свою очередь, методы с применением нейронных сетей устроены по другому. Нейронная сеть сама определяет набор признаков на основе данных, полученных в процессе обучения.

До 2012 года наибольшей точностью обладали методы с ручным выделением признаков. Тогда пользовались популярностью методы основанные на гистограммах направленных градиентов и признаках Хаара. Однако в 2012 году эта тенденция изменилась после того, как на соревновании ImageNet Large Scale Visual Recognition Challenge (ILSVRC) методы с применением искусственных нейронных сетей показали лучшие результаты.

Участникам соревнования ILSVRC предлагается разработать алгоритмы для задач классификации изображений и обнаружения объектов [3]. Сами соревнования начали проводиться с 2010 года и использовали ImageNet, в качестве базы данных изображений. ImageNet - это база данных, которая

содержит более 14 миллионов изображений, объекты на которых могут относиться к одному из 1000 классов [4].

На рисунке 1.1 представлена диаграмма с лучшими решениями задачи классификации на ILSVRC с 2011 по 2016 год. По горизонтальной оси откладывается время, по вертикальной - процент ошибок. Для каждого из годов в скобках указано название предложенного метода. Как можно заметить, в 2012 году произошел достаточно резкий скачок в точности классификации изображений, почти на 10%. Этот прирост обусловлен тем, что авторы решили отойти от классических подходов и разработали метод, в основе которого лежали сверточные нейронные сети (СНС). В последующие года победителями также становились методы, основанные на СНС, а в 2015 году предложенное решение стало классифицировать изображения даже точнее, чем это может выполнить человек.

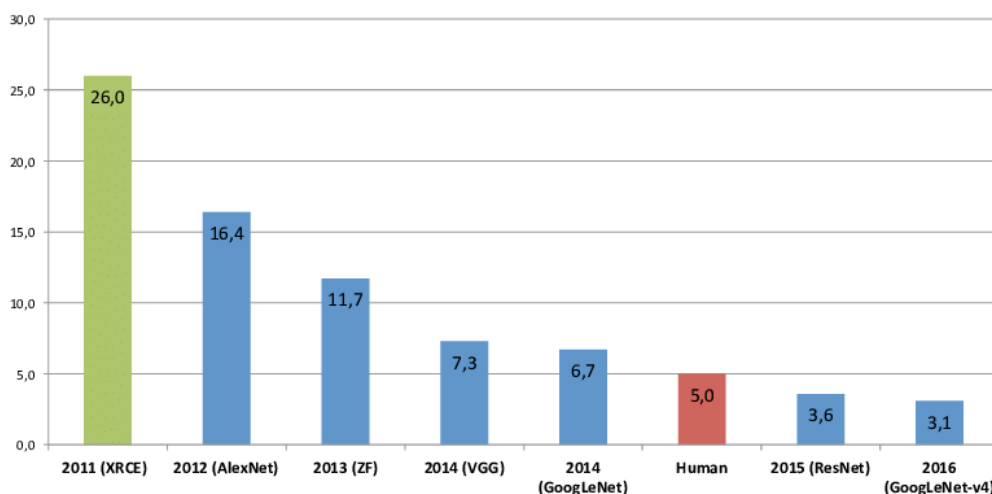


Рисунок 1.1 – Лучшие решения ILSVRC с 2011 по 2016

1.1.2 Сверточные нейронные сети

Идея сверточных нейронных сетей (СНС) появилась достаточно давно, примерно в начале 80-х годов [5]. По принципу своей работы СНС схожи с тем, как функционирует зрительная кора головного мозга человека. Именно из-за этого сходства СНС удается достигать наилучших результатов в различных задачах компьютерного зрения.

СНС получили свое название из-за операции свертки (convolution), которая как раз и выделяет интересные признаки на изображениях. При операциях свертки используется так называемые фильтры (filters) или ядра (kernels). В случаях двумерной свертки это двумерная матрица весов нечетного размера (обычно 3×3 или 5×5). Операцию двумерной свертки для изображения заданного в оттенках серого (то есть каждый пиксель занимает один байт) можно описать следующим образом:

- 1) переместиться в левый-верхний угол входного изображения, таким образом, чтобы каждому весу фильтра соответствовал пиксель изображения;
- 2) выполнить произведение соответствующих пикселей изображения и весов фильтра;
- 3) просуммировать полученные произведения. Эта сумма и будет результатом свертки области изображения под фильтром;
- 4) сдвинуть фильтр вправо на величину шага (stride). При достижении правой границы изображения спуститься на строку ниже;
- 5) вернуться на шаг два, пока полностью не будет обработано все изображение.

Как можно понять из приведенного описания, операция свертки уменьшает размер результирующего изображения, то есть происходит потеря информации на краях изображения. Для решения этой проблемы стали использовать дополнение (padding). Чаще всего применяется нулевое дополнение, когда изображение по краям дополняется нулями. Пример нулевого дополнения для картинке размером 7×7 и размером фильтра 5×5 представлен на рисунке 1.2.

В настоящее время большинство изображений используют цветовую модель RGB, когда каждый пиксель кодируется тремя числами – интенсивностью соответственно красного, зеленого и синего цветов. Поэтому фильтрам

в операции свертки добавляют третью компоненту, чтобы покрывались все каналы входного изображения.



Рисунок 1.2 – Пример нулевого дополнения

Основными составляющими СНС являются сверточные слои, где обычно происходит порядка 90-95% всех вычислительных операций [6]. Каждый из слоев может содержать в себе некоторое количество фильтров, размер и количество которых определяются архитектурой СНС. После применения фильтра к результату предыдущего слоя формируется так называемая карта признаков (feature map). То есть каждый фильтр в сверточном слое выделяет только те признаки, которые он выучил в процессе обучения. В итоге фильтры первого сверточного слоя будут выделять какие-то простые признаки (прямые, наклонные), далее фильтры второго слоя смогут выделять уже полноценные фигуры (квадраты, треугольники) и чем глубже будет находиться слой, тем все более сложные конструкции он сможет выделять.

Так как для всей операции свертки используются одинаковые веса фильтров, то это приводит к более быстрому обучению сверточных слоев по сравнению с теми же полносвязными слоями. Как раз идея разделения весов фильтров является уместной для задач обработки изображений, так как основным свойством всех изображений является их пространственная инвариантность [7]. Иными словами любая форма на изображении должна обрабатываться аналогичным образом, в не зависимости от ее расположения.

За сверточным слоем в СНС следует операция активации. Основной задачей функции активации в нейронных сетях является добавление нелинейности, иначе при отсутствии функции активации (либо при использовании линейной функции активации) увеличение количества слоев не давало бы никакого выигрыша [7]. Иначе говоря при использовании функции активации нейронная сеть может выучить более сложные зависимости во входных данных.

Классическими примерами функций активации являются сигмоида и гиперболический тангенс. На рисунке 1.3 приведены графики данных функций активации. Однако в последнее время все чаще используется выпрямленная линейная функция активации (ReLU). Причин у этого несколько:

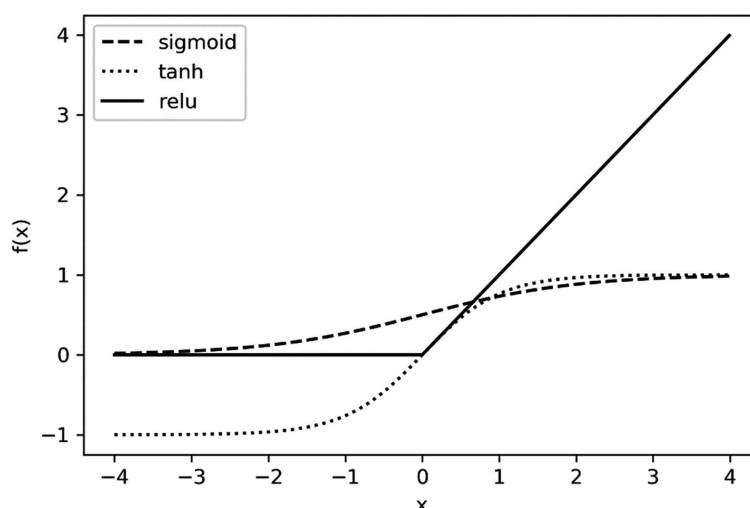


Рисунок 1.3 – Примеры графиков функций активации

- 1) значения функций сигмоиды и гиперболического тангенса изменяются незначительно при стремлении аргумента к $+\infty$ или $-\infty$. Это становится причиной проблемы затухающего градиента (vanishing gradient), что приводит к большому времени обучения для более глубоких слоев. Функция активации ReLU лишена этой проблемы. Именно благодаря этой особенности авторам известной СНС AlexNet удалось в шесть раз ускорить процесс обучения при использовании ReLU по сравнению с гиперболическим тангенсом [8];

- 2) расчет значений функций сигмоиды и гиперболического тангенса более ресурсоемкий с точки зрения центрального процессора, так как необходимо выполнять операции возведения в степень и деления.

После применения функции активации чаще следует операция субдискретизации (pooling). Существует несколько типов данной операции, но на данный момент широко используемой является операция взятия максимума (max-pooling). Сама операция состоит из разбиения результатов работы предыдущего слоя на локальные области размером $w \times h$ и взятия максимума из этой области.

Основной целью операции субдискретизации является уменьшение размеров карты признаков. Это приводит к тому, что последующие операции свертки будут оперировать над все большей областью исходного изображения. Также очевидным следствием данной операции будет уменьшение количества вычислений в последующих слоях. Помимо этого, в результате операции субдискретизации СНС становится более устойчивой к небольшим трансформациям изображения вроде сдвига или поворота [5].

После чередования нескольких сверточных слоев с операциями активации и субдискретизации, завершающим обычно служит полносвязный (fully connected) слой, который уже непосредственно соединяется с выходным слоем СНС. В действительности количество полносвязных слоев может быть больше одного, этот вопрос остается на усмотрение авторам СНС. Главной задачей полносвязного слоя является установления соответствия, какие из высокоуровневых признаков принадлежат к тому или иному классу. Чаще всего полносвязные слои из-за своей структуры содержат порядка 95% всех параметров СНС, но при этом в этих слоях происходит порядка 5-10% от всех вычислительных операций [6].

Обобщая все вышесказанное, типичная архитектура СНС состоит из чередующихся слоев свертки с операциями активации и субдискретизации и завершающим это полносвязным слоем. Структура выходного слоя опреде-

ляется задачей, решаемой СНС. Например, для задач классификации размер выходного слоя будет равен количеству классов. На рисунке 1.4 представлена архитектура известной СНС AlexNet без частей, связанных с разделением слоев между двумя GPU. Современные СНС могут содержать десятки сверточных слоев, поэтому часто операции активации даже не приводят на рисунках, подразумевая, что после каждого сверточного слоя следует ReLU активация. Операция субдискретизации обозначается как MP (сокращение от max-pooling) и можно заметить, что операция субдискретизации может отсутствовать на выходе некоторых сверточных слоев.

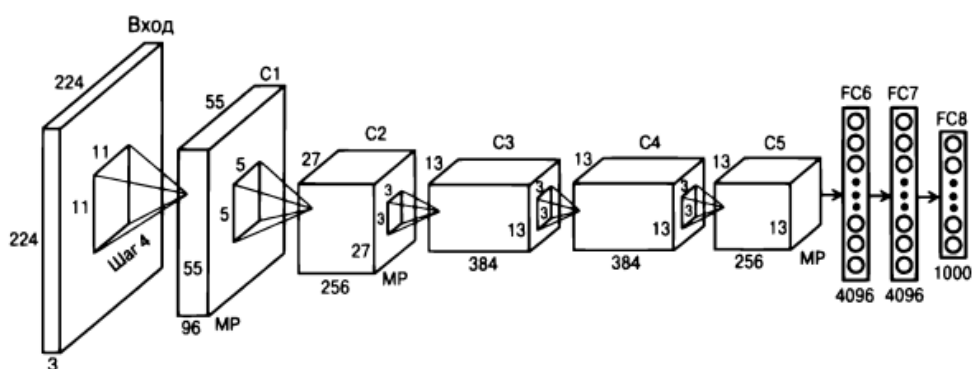


Рисунок 1.4 – Архитектура AlexNet без GPU-разделения [7]

1.1.3 Методы для распознавания объектов

После большого успеха СНС в решении задач классификации, данный тип нейронных сетей также начал широко применяться в методах для распознавания объектов. Под распознаванием объектов будем понимать обнаружение объектов на изображениях, то есть определение координат ограничивающей рамки (bounding box), и определение их класса.

На данный момент существующие методы для распознавания объектов можно разделить на два класса: одношаговые (one-stage) и двухшаговые (two-stage). В двухшаговых методах входное изображение разбивается на регионы (regions of interest, RoI) и далее каждый из RoI обрабатывается отдельно. Также определением координат ограничивающих рамок и классов объектов занимаются разные компоненты. При использовании одношаговых методов

входное изображение полностью обрабатывается нейронной сетью, которая одновременно предсказывает, и координаты ограничивающих рамок, и классы объектов. В общем случае двухшаговые методы требуют на обработку больше времени по сравнению с одношаговыми, но при этом обладают большей точностью распознавания [9].

R-CNN

Первым представителем двухшаговых методов стал метод Regions with CNN features (R-CNN). Тогда авторам удалось повысить среднюю точность (mean average precision) распознавания объектов на 30% по сравнению с предыдущими методами на датасете PASCAL VOC2012 [10]. На рисунке 1.5 представлено описание метода R-CNN. Рассмотрим каждый из шагов отдельно:

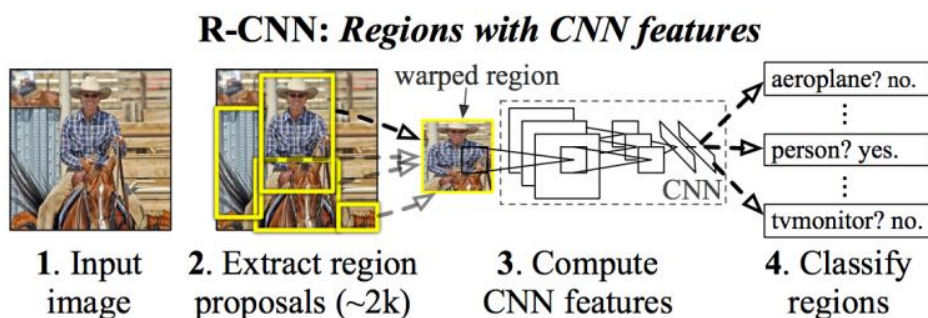


Рисунок 1.5 – Описание метода R-CNN [10]

- 1) изначально из входного изображения выделяются примерно 2000 RoI, на которых могут находиться объекты. Для поиска RoI авторы использовали метод селективного поиска (selective search). Упрощая, можно сказать, что метод селективного поиска определяет границы объектов по изменению цветов и интенсивностей соседних пикселей;
- 2) так как выделенные RoI могут быть разных размеров, то каждый из них необходимо приводить к входным размерам СНС – 227×227 ;
- 3) после каждый из RoI передается на вход СНС. В качестве СНС авторы используют архитектуру AlexNet, у которой удален последний полно-

связный слой. В связи с этим на выходе получается вектор признаков размерностью 4096;

- 4) сама классификация объектов осуществлялась с помощью метода опорных векторов (support-vector machine, SVM). Для каждого из классов был свой бинарный SVM-классификатор;
- 5) координаты объектов, полученные на первом шаге могут быть неточными, поэтому каждый из RoI в конце еще обрабатывается моделью линейной регрессии, отдельной для каждого из классов. По утверждениям авторов данный шаг увеличил среднюю точность на 3-4% [10].

Fast R-CNN

Метод R-CNN обладал хорошей точностью, однако общее время обработки изображения было крайне высоким. Например, при использовании в качестве СНС архитектуру VGG16 обработка одного изображения методом R-CNN занимала 47 секунд (измерения проводились при использовании одной GPU NVIDIA Tesla K40) [11]. У этого было несколько причин. Во-первых, проход 2000 RoI через СНС требовал много вычислительных и временных ресурсов, особенно ситуация ухудшалась для более глубоких СНС, таких как, VGG16. Во-вторых, модель линейной регрессии зависела от результатов классификации, что не давало выполнять эти операции параллельно. Также наличие 3 отдельных моделей требовало много времени и дискового пространства в процессе обучения [11].

Для решения вышеперечисленных проблем авторами R-CNN вскоре был представлен улучшенный метод – Fast R-CNN. На рисунке 1.6 приведено описание данного метода. Рассмотрим обновленный метод подробнее:

- 1) теперь входное изображение сразу поступает на вход СНС. Параллельно с этим из входного изображения выделяются RoI с помощью того же метода селективного поиска;
- 2) далее из карты признаков выделяется вектор признаков для каждого из RoI. Решением этой задачи занимается RoI pooling layer;

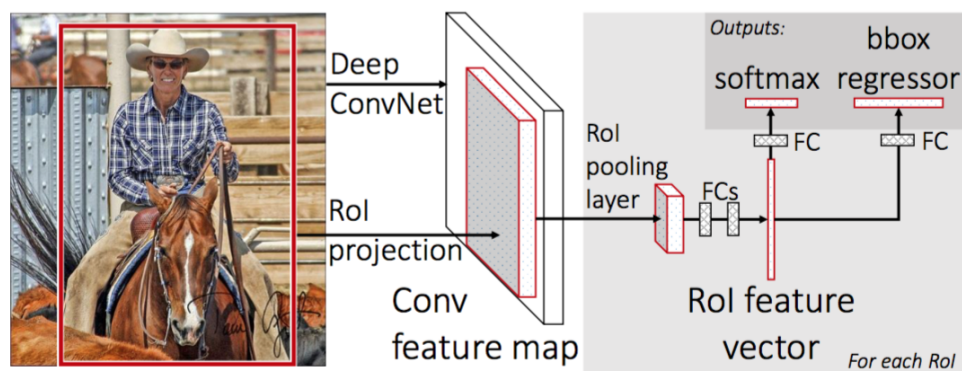


Рисунок 1.6 – Описание метода Fast R-CNN [11]

- 3) после вектор признаков передается в полносвязный слой, а затем следуют два параллельных слоя - Softmax и линейная регрессия. Слой Softmax определяет класс объекта, а слой линейной регрессии уточняет координаты ограничивающей рамки.

Faster R-CNN

Метод Fast R-CNN устранил основные проблемы первоначального метода R-CNN, но все равно еще плохо подходил для тех же задач обработки видео. Например, на обработку одного изображения уходило порядка двух секунд и большая часть этого времени тратилась на поиск RoI с помощью селективного поиска [12]. Для решения этой проблемы авторы метода Faster R-CNN заменили селективный поиск отдельной СНС, которую называли Region Proposal Network (RPN). Все остальные компоненты метода Fast R-CNN остались без изменений.

На рисунке 1.7 иллюстрируется принцип работы RPN. На вход RPN передается карта признаков, ранее выделенных с помощью СНС. Далее по карте признаков проходятся "скользящим окном" (sliding window) размером 3×3 . Для каждого положения окна выделяется k RoI, где k это количество якорей (anchors). Якоря используются далее для определения абсолютных координат RoI. После наложения якорей на текущее положение скользящего окна, из карты признаков выделяется вектор признаков размерностью 256.

Затем полученный вектор передается в два параллельных сверточных слоя с размером фильтра 1×1 . Слой *cls* имеет $2k$ выходов для определения вероятности наличия/отсутствия объекта внутри RoI. В свою очередь слой *reg* имеет $4k$ выходов, каждый из которых соответствует поправкам к координатам якорей.

Проведенные авторами тесты показывают, что предложенная ими RPN тратит порядка 10 мс на выделение RoI [12]. Благодаря этому удалось достичь скорости обработки порядка пяти кадров в секунду, что соответствует 200 мс на кадр. При всем при этом предложенный метод оказался более точным. На датасете PASCAL VOC2007 прирост средней точности составил порядка 9% по сравнению с методом Fast R-CNN.

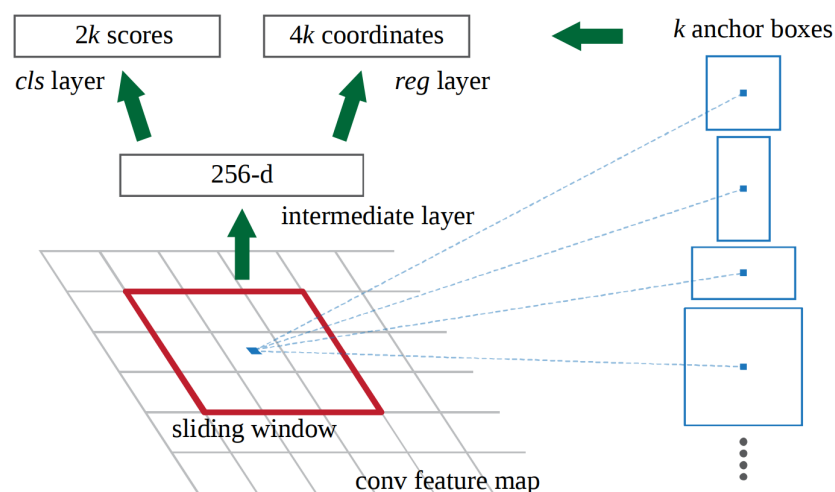


Рисунок 1.7 – Принцип работы RPN [12]

YOLO

Метод You Only Look Once (YOLO) является одним из первых представителей одношаговых методов. По утверждению авторов последняя версия метода YOLOv4 обладает лучшей эффективностью среди существующих методов [13]. Под эффективностью понимается отношение точности к времени обработки одного изображения. На текущий момент представлено четыре версии метода: YOLO [14], YOLOv2 [15], YOLOv3 [16] и YOLOv4 [13]. Рассмотрим обобщенно принцип работы данных методов.

Изначально авторы первой версии YOLO решили рассмотреть проблему распознавания объектов как проблему регрессии, то есть одновременное предсказание координат ограничивающей рамки и класса объекта внутри нее. Сам метод разделяет входное изображение на сетку размером $S \times S$. Каждому из объектов присваивается клетка (grid cell) из этой сетки. Если объект занимает более одной клетки, то ему присваивается та клетка, внутри которой расположен центр этого объекта. Каждая клетка предсказывает B ограничивающих рамок, разных форм и размеров.

Само предсказание для ограничивающих рамок состоит из пяти чисел: (x, y, w, h, a) . Числа (x, y) представляют собой координаты центра ограничивающей рамки, относительно клетки объекта, а числа (w, h) соответственно являются шириной и высотой рамки. Число a представляет собой точность предсказания. Также для каждой из клеток сетки предсказывается вероятность принадлежности к тому или иному классу. Если предположить, что общее количество классов равно C , то на выходе из нейронной сети будет получаться тензор размером $S \times S \times (5B + C)$.

Так как для одного объекта может быть предсказано больше одной ограничивающей рамки, то для выбора наилучшей используют метод Non-Max Suppression (NMS). На рисунке 1.8 представлен результат работы NMS. Рассмотрим этот метод подробнее:

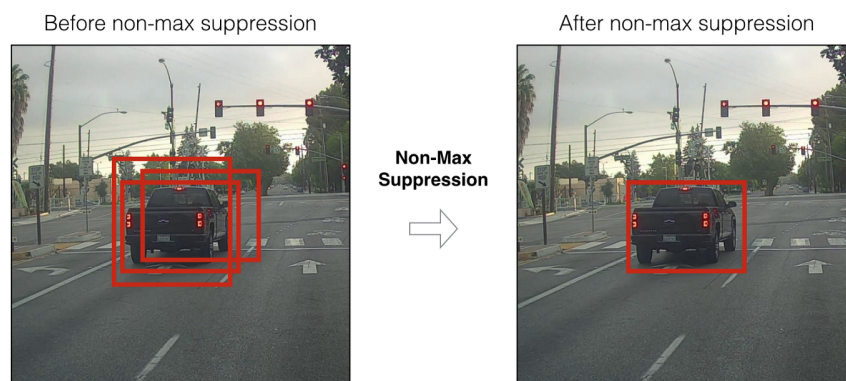


Рисунок 1.8 – Пример работы метода NMS

- 1) выбрать охватывающую рамку, имеющую самую большую точность предсказания;
- 2) отбросить все рамки, которые имеют площадь перекрытия с выбранной рамкой больше определенной границы. Для этого используется значение Intersection over Union, которое рассчитывается как отношение площади пересечения двух рамок к площади их объединения;
- 3) перейти к первому шагу, пока есть лишние рамки.

Также большим преимуществом метода YOLO является тот факт, что авторы предоставляют версии предобученных моделей для запуска в разных окружениях. Например, в последней работе авторы представили легковесную модель YOLOv4-tiny, которая очень хорошо подходит для запуска на устройствах с ограниченными вычислительными ресурсами, при этом средняя точность составляет 42% [17].

1.2 Архитектура процессоров ARM

1.2.1 Особенности

В настоящее время большинство устройств Интернета вещей находятся под управлением именно процессоров архитектуры ARM. Прямым конкурентом ARM является не менее известная архитектура процессоров – x86. Данная архитектура изначально разработана компанией Intel и большинство процессоров настольных компьютеров и серверов реализуют эту архитектуру. Однако со временем повсеместное использование процессоров x86 в перечисленных областях может склониться в сторону ARM, для этого достаточно вспомнить, как в 2020 году компания Apple перешла на использование в своих ноутбуках процессоров архитектуры ARM [18].

Архитектура ARM является типичным представителем архитектуры Reduced Instruction Set Computer (RISC), в то время, как x86 относится к архитектуре Complex Instruction Set Computer (CISC). Рассмотрим отличительные особенности архитектуры RISC и сравним их с CISC:

- 1) меньшее количество инструкций. Так как чаще всего используется только ограниченный набор популярных инструкций, то аппаратные блоки для менее популярных инструкций в архитектуре CISC не задействованы и только занимают место на кристалле. Недостатком этой особенности может быть увеличение размера программы, но при текущем количестве памяти у компьютеров это не является большой проблемой;
- 2) фиксированный размер и простой формат инструкций. Это приводит к тому, что операция декодирования инструкции занимает меньше времени и сама аппаратная структура декодировщика остается простой. В случае CISC операция декодирования будет занимать больше времени, даже для самых простых инструкций [19];
- 3) инструкции выполняются за один такт. Современные процессоры часто используют конвейеризацию инструкций и данная особенность RISC позволяет повысить производительность конвейеров;
- 4) большое количество регистров. Например, процессоры архитектуры ARMv8-a имеют 31 регистр общего назначения [20], в то время как x86_64 предоставляет только 16 регистров общего назначения [21]. В большей степени это связано с тем, что операции обработки данных в RISC в качестве операндов могут использовать только регистры, поэтому перед обработкой данных их нужно загрузить из памяти, а после сохранить результат, используя соответствующие инструкции.

Стоит также упомянуть, что для архитектуры x86 некоторые описанные выше проблемы CISC со временем начали усугубляться из-за необходимости сохранять поддержку не самых популярных инструкций, которые были добавлены десятки лет назад. Поэтому начиная с середины 90-х годов процессоры архитектуры x86 разбивают сложную инструкцию на несколько микроопераций (micro-operations), а после эти инструкции уже выполняются RISC ядром [19].

Также отличительной особенностью процессоров архитектуры ARM является пониженное энергопотребление, по сравнению с процессорами на x86 архитектуре. Например, тот же ARM процессор Apple M1 на данный момент обладает одной из наилучших производительностей на один ватт [18]. Такое достигается за счет использования гетерогенной архитектуры процессоров big.LITTLE, когда используется два типа ядер: высокопроизводительные (big) и энергоэффективные (LITTLE). Высокопроизводительные используются для требовательных задач, в то время как энергоэффективные активизируются для менее ресурсоемких задач, но при этом обладают пониженным энергопотреблением.

1.2.2 Наборы команд

Со временем разброс устройств, для которых применялись процессоры ARM стал увеличиваться и тогда компания решила ввести понятие профиля архитектуры, который характеризует класс решаемых задач. На данный момент определены три профиля [20]:

- 1) профиль приложений (application profile). К данному профилю относятся устройства, требующие высокой производительности, такие как: персональные компьютеры, планшеты, смартфоны. Процессоры данного профиля имеют поддержку виртуальной памяти за счет наличия устройства управления памятью (memory management unit, MMU) [20];
- 2) профиль реального-времени (real-time profile). Процессоры данного профиля применяются для решения задач жесткого реального времени, то есть таких задач, где требуется гарантия, что обработка события займет не более определенного количества времени. Примером таких систем может служить антиблокировочная система тормозов на автомобилях и самолетах. Особенностью данного профиля является отсутствие MMU и наличие тесно связанной памяти (Tightly Coupled Memory, TCM). TCM - это область памяти, которая по своей струк-

туре похожа на кэш процессора, но содержимым этой памяти строго управляет программист и доступ к этой части памяти занимает один такт процессора;

- 3) профиль микроконтроллеров (microcontroller profile). Подобные процессоры обладают пониженным энергопотреблением и позволяют быстрее обрабатывать внешние прерывания [20]. Может использоваться в качестве процессоров для микроконтроллеров или как составная часть системы на кристалле, например, может выступать контроллером ввода-вывода.

После выхода архитектуры ARMv8 у процессоров появилось два режима работы (execution state): aarch64 и aarch32. В режиме работы aarch32 используются 32-битные адреса, есть доступ к 13 регистрам общего назначения размером 32 бита, имеется 32 64-битных регистра для SIMD инструкций и в качестве набора инструкций могут использоваться A32 или T32 [20]. В свою очередь в режиме aarch64 используются 64-битные адреса, количество регистров общего назначения увеличилось до 31, размерность регистров для SIMD инструкций выросла до 128 бит, а в качестве набора команд может использоваться A64 [20]. Процессоры, начиная с архитектуры ARMv8, могут исполнять программы в двух режимах, тем самым сохраняется обратная совместимость с ранее разработанным ПО, для которого использовался набор инструкций A32 или T32.

Как ранее упоминалось, процессоры архитектуры ARM поддерживают SIMD (Single Instruction Multiple Data) инструкции, для этого используется расширение набора команд – NEON. Данный тип инструкций позволяет за одну инструкцию процессора выполнить несколько операций, например, перемножить между собой векторы чисел, каждый из которых содержит по 4 числа. Подобные инструкции могут значительно ускорить операции кодирования/декодирования видео, потому что при решении таких задач часто оперируют векторами чисел. Этими инструкциями можно воспользоваться

напрямую при написании на ассемблере, либо использовать интринсики компиляторов на языке C/C++. Также некоторые библиотеки могут использовать при своей работе SIMD инструкции, например, такую поддержку имеет популярная библиотека OpenCV, разработанная для решения задач компьютерного зрения.

СПИСОК ЛИТЕРАТУРЫ

1. Inside the Amazon Warehouse Where Humans and Machines Become One [Электронный ресурс]. 2019. URL: <https://www.wired.com/story/amazon-warehouse-robots/> (дата обращения: 22.04.2021).
2. What Edge Computing Means for Infrastructure and Operations Leaders [Электронный ресурс]. 2018. URL: <https://www.gartner.com/smarterwithgartner/what-edge-computing-means-for-infrastructure-and-operations-leaders/> (дата обращения: 23.04.2021).
3. ImageNet Large Scale Visual Recognition Challenge [Электронный ресурс]. Официальная страница соревнований ILSVRC. 2021. URL: <http://image-net.org/challenges/LSVRC/> (дата обращения: 13.04.2021).
4. ImageNet [Электронный ресурс]. Официальный сайт проекта ImageNet. 2021. URL: <http://image-net.org/index> (дата обращения: 13.04.2021).
5. Николенко С.И., Кадурын А.А., Архангельская Е.В. Глубокое обучение. Погружение в мир нейронных сетей. СПб: Питер, 2018. 480 с.
6. Krizhevsky A. One weird trick for parallelizing convolutional neural networks // arXiv:1404.5997. 2014.
7. Аггарвал Чару. Нейронные сети и глубокое обучение. Учебный курс. СПб: Диалектика, 2020. 752 с.
8. Krizhevsky A., Sutskever I., Hinton G. E. ImageNet Classification with Deep Convolutional Neural Networks // Advances in Neural Information Processing Systems / Ed. by F. Pereira, C. J. C. Burges, L. Bottou et al. Vol. 25. Curran Associates, Inc., 2012.

9. Deep Learning in Object Detection and Recognition / X. Jiang, A. Hadid, Y. Pang et al. Singapore: Springer, 2019. 240 p.
10. Rich feature hierarchies for accurate object detection and semantic segmentation / R. Girshick, J. Donahue, T. Darrell et al. // arXiv:1311.2524. 2014.
11. Girshick R. Fast R-CNN // arXiv:1504.08083. 2015.
12. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks / S. Ren, K. He, R. Girshick et al. // arXiv:1506.01497. 2016.
13. Bochkovskiy A., Wang C.-Y., Liao H.-Y. M. YOLOv4: Optimal Speed and Accuracy of Object Detection // arXiv:2004.10934. 2020.
14. You Only Look Once: Unified, Real-Time Object Detection / J. Redmon, S. Divvala, R. Girshick et al. // arXiv:1506.02640. 2016.
15. Redmon J., Farhadi A. YOLO9000: Better, Faster, Stronger // arXiv:1612.08242. 2016.
16. Redmon J., Farhadi A. YOLOv3: An Incremental Improvement // arXiv:1804.02767. 2018.
17. Wang C.-Y., Bochkovskiy A., Liao H.-Y. M. Scaled-YOLOv4: Scaling Cross Stage Partial Network // arXiv:2011.08036. 2021.
18. Apple unleashes M1 [Электронный ресурс]. Пресс-релиз компании Apple о процессоре M1. 2020. URL: <https://www.apple.com/newsroom/2020/11/apple-unleashes-m1/> (дата обращения: 09.05.2021).
19. Harris S. L., Harris D. M. Digital Design and Computer Architecture: ARM Edition. Waltham: Morgan Kaufmann, 2016. 584 p.

20. Arm, Cambridge, England. Arm Architecture Reference Manual. Armv8, for Armv8-A architecture profile. 2021.
21. Intel Corporation, Santa Clara, USA. Intel 64 and IA-32 Architectures Software Developer's Manual. 2021.