

Санкт-Петербургский государственный университет

Выпускная квалификационная работа
общей образовательной программы бакалавриата
«Прикладная математика и информатика»

студента 4 курса, группы 18.Б04-мм
Ефима Сергеевича Бутакова

на тему

«Методы машинного обучения в задаче распознавания аудиосигнала»

Научный руководитель:
к.ф.-м.н., Михаил Сергеевич Ананьевский
кафедра теоретической кибернетики

Рецензент:
д.т.н., Игорь Борисович Фуртат
в.н.с. лаб. УСС ИПМаш РАН

г. Санкт-Петербург

2022

Saint Petersburg State University

Graduation Project

Applied Mathematics and Computer Science

Control and Processing of Information in Cybernetical and Robotic Systems

Butakov Efim

Machine Learning Techniques for Audio Recognition

Scientific Supervisor:

Mikhail Ananyevskiy

Candidate of Physico-Mathematical Sciences

Reviewer:

Igor Furtat

Doctor of Engineering Sciences

Saint Petersburg

2022

Оглавление

1.	Введение	2
1.1.	Постановка задачи.	3
2.	Обзор существующих методов.	4
2.1.	Предобработка звука.	4
	Подготовка звука.	4
	Создание спектограммы.	4
2.2.	Сверточная нейронная сеть (CNN)	6
2.3.	Рекуррентная нейронная сеть (RNN)	6
2.4.	Сети с долгой кратковременной памятью (LSTM)	7
2.5.	Закрытый рекуррентный блок GRU	9
2.6.	DeepSpeech	10
2.7.	DeepSpeech 2	11
2.8.	QuartzNet	14
3.	Эксперименты по обучению моделей.	17
3.1.	Данные для распознавания речи LJ Speech Dataset.	17
3.2.	Функция ошибки.	17
3.3.	Тестирование и сравнение моделей.	18
3.4.	Возможные улучшения полученных результатов	19
4.	Результаты.	20
5.	Список литературы.	21
6.	Приложение	23

1. Введение

Автоматическое распознавание речи (Automatic speech recognition - ASR) – сфера, которая активно исследуется последние пять десятилетий, эта тема считается важным пунктом улучшения взаимодействия в сценариях человек - человек и человек-машина. Однако, в недалёком прошлом речь не являлась ключевым инструментом во взаимодействии человека и машины. Отчасти это связано с тем, что инструменты распознавания речи в то время были недостаточно хороши, чтобы преодолеть необходимый уровень точности и использоваться в реальных условиях. В то же время, альтернативные способы взаимодействия, например клавиатура и мышь, значительно превосходили речь по эффективности, ограничениям и точности.

В последние несколько лет наблюдается новый всплеск интереса к ASR. Это может быть связано с повышением требований к задаче автоматического распознавания речи в мобильных устройствах и успех новых виртуальных речевых помощников (например, Apple's Siri, Google Now и Microsoft's Cortana). Не менее важным пунктом является развитие методов глубокого обучения и увеличение вычислительных возможностей. Комбинированное использование методов глубокого обучения позволило уменьшить коэффициент ошибок в распознавании речи больше чем на треть относительно известных ранее классических методов GMM-НММ. Например, на данный момент точность распознавания слова для Английского языка достигает 95% и более.

Таким образом, объектом данного исследования является: задача ASR, предметом исследования является: использование актуальных архитектур нейронных сетей для задачи автоматического распознавания речи и сравнение их между собой.

1.1. Постановка задачи.

В рамках данной работы были поставлены следующие задачи:

- Исследовать современных подходы к решению задачи автоматического распознавания речи.
- Рассмотреть детали реализации этих моделей.
- Обучить несколько вариаций современных нейронных сетей.
- Сравнить обученные модели и выявить наиболее подходящий метод к решению задачи автоматического распознавания речи.

2. Обзор существующих методов.

2.1. Предобработка звука.

Подготовка звука.

Как известно звук — это волна. Но как превратить волну в набор чисел? Звуковые волны одномерны. В каждый момент времени у них есть одно значение, зависящее от амплитуды волны. Чтобы превратить звуковую волну в числа, мы просто записываем значения амплитуды волны в равноотстоящих точках. Этот процесс называется дискретизацией. Мы считываем данные тысячи раз в секунду и записываем числа, соответствующие амплитуде звуковой волны в этот момент времени. Чтобы не потерять

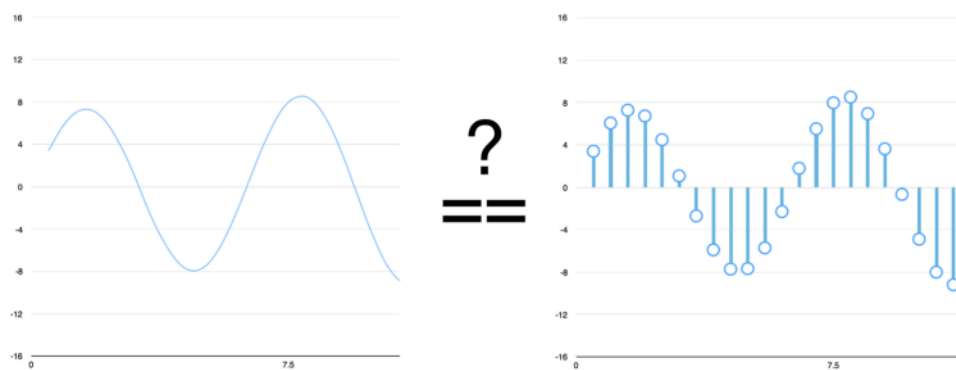


Рис. 1. Дискретизация.

часть информации мы воспользуемся теоремой Котельникова, которая говорит, что для идеального воссоздания исходной звуковой волны достаточно использовать частоту дискретизации, вдвое превышающую самую высокую частоту записываемого звука. Таким образом из звуковой волны мы получили большой массив чисел.

Создание спектограммы.

Мы могли бы просто обучить нейросеть на этих числах, но распознавание речевых моделей путем обработки этих чисел напрямую затруднительно. Вместо этого мы можем облегчить задачу, проведя некоторую предварительную обработку аудиоданных. Для работы с цифровым звуком необходимо представить итоговую звуковую волну в виде суммы синусоидальных функций с помощью метода прямого преобразования Фурье, то есть найти спектр звукового сигнала.

Спектр звукового сигнала это совокупность синусоидальных звуковых волн, в результате суперпозиции образующих звуковую волну. Спектр звукового сигнала является самым важным инструментом анализа и обработки любого звука. Процесс преобразования звуковой волны в набор отдельных гармоник называется прямым преобразованием Фурье. Обратный процесс называется обратным преобразованием Фурье

Ряд Фурье для временной периодической функции $y(t)$:

$$y(t) = \frac{a_0}{2} + \sum_{i=1}^{\infty} (a_i \cos(i\omega t) + b_i \sin(i\omega t))$$

где

$$a_0 = \frac{2}{T} \int_0^T y(t) dt$$

$$a_i = \frac{2}{T} \int_0^T (y(t) \cos(i\omega t)) dt$$

$$b_i = \frac{2}{T} \int_0^T (y(t) \sin(i\omega t)) dt$$

На примере произвольного фрагмента записи человеческой речи, записанного с помощью компьютера приведу пример спектра звукового сигнала и амплитудного графика. На следующем рисунке представлена некоторая зависимость амплитуды звукового сигнала от времени, которая характеризует определенный фрагмент человеческой речи. Данное представление ее в виде дв речи. Данный способ представления звука называется спек-

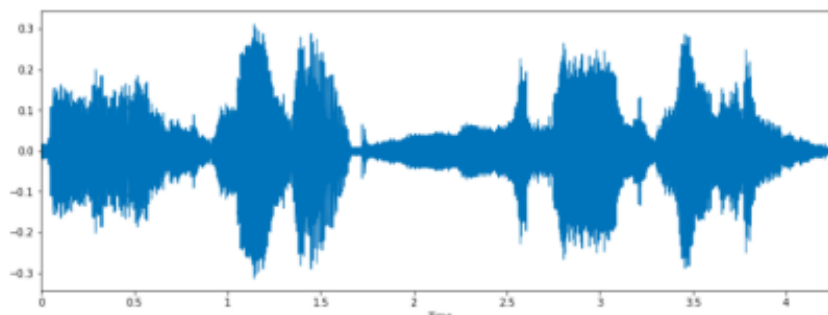


Рис. 2. Амплитудный график фрагмента человеческой речи..

трограммой звука или сигналограммой. Наиболее распространенным представлением спектрограммы является двумерной диаграммы:

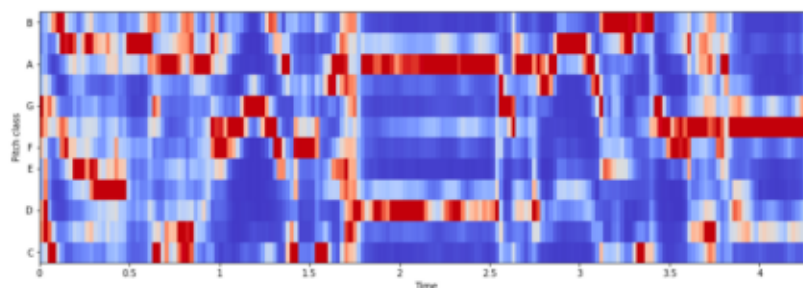


Рис. 3. Спектрограмма.

В данной работе для обучения и сравнения методов обучения нейронных сетей мы будем использовать спектрограммы.

2.2. Сверточная нейронная сеть (CNN)

Сверточная нейронная сеть (англ. convolutional neural network, CNN) — специальная архитектура нейронных сетей, предложенная Яном Лекуном, изначально нацеленная на эффективное распознавание изображений. В сверточной нейронной сети выходы промежуточных слоев образуют матрицу (изображение) или набор матриц (несколько слоев изображения). Так, например, на вход сверточной нейронной сети можно подавать три слоя изображения (R-, G-, B-каналы изображения). Основными видами слоев в сверточной нейронной сети являются сверточные слои (англ. convolutional layer), пулинговые слои (англ. pooling layer).

2.3. Рекуррентная нейронная сеть (RNN)

Начнем с основ, а именно с простого RNN.

Весь смысл обработки связной последовательности данных заключается в том, чтобы кроме выделения отклика для каждого элемента, суметь учесть и связь элементов. Например, можно представить изображение в виде набора векторов, каждый из которых содержит пиксели одной колонки. Будем последовательно подавать такие вектора на вход сети, если мы хотим, чтобы сеть классифицировала изображение, необходимо сделать так, чтобы сеть умела “запоминать” информацию об уже просмотренных колонках. Аналогично, если мы хотим научить сеть классифицировать предложения естественного языка (например, определять эмоциональный окрас предложения), и подаём в сеть одно слово за другим, нам желательно, чтобы сеть “помнила” уже переданные

слова. Если мы хотим, чтобы сеть переводила предложение с одного языка на другой, то тоже было бы не плохо учитывать начало предложение при переводе середины и конца.

Как раз задачу “запоминания” уже рассмотренных элементов последовательности и предполагается решать при помощи рекуррентной сети. Для этого, кроме выходного вектора, сеть должна иметь еще и некоторый вектор, который описывает текущее внутреннее состояние сети, т.е. в нем содержатся воспоминания о всех уже рассмотренных сетью элементах. Более формально это выглядит так.

Допустим у нас есть набор входных векторов $(X^{(1)}, X^{(2)}, \dots, X^{(n)})$, мы их последовательно преобразуем:

$$H^t = \sigma(W^{hx} X^t + W^{hh} H^{(t-1)} + b_h)$$

$$Y^{(t)} = W^{yh} H^{(t)} + b_y$$

При этом кроме выхода $Y(t)$ (который нам возможно и не нужен на каждом шаге), мы имеем еще вектор $H(t)$, описывающий текущее состояние.

Таким образом сеть выглядит следующим образом:

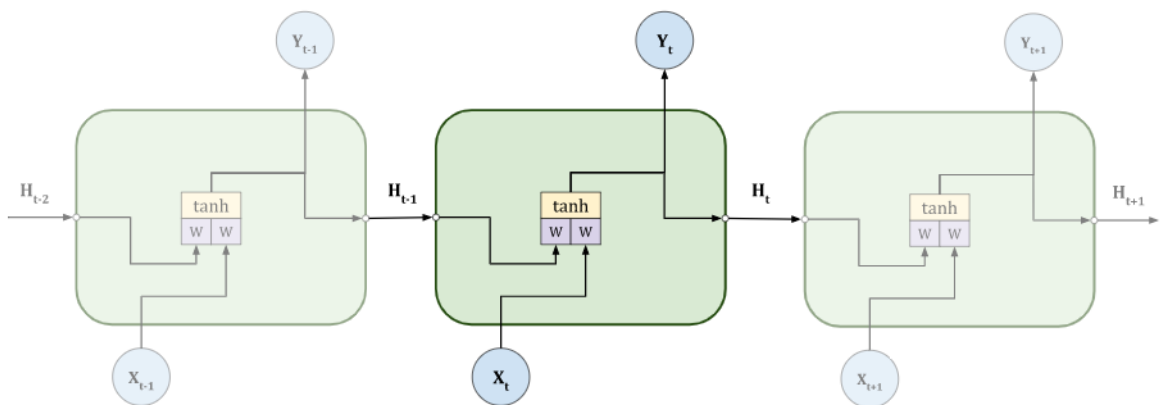


Рис. 3. RNN.

2.4. Сети с долгой кратковременной памятью (LSTM)

Проблема с обычными RNN ячейками, описанными выше, заключается в том, что они не могут “удержать в памяти” слишком длинные последовательности. Связано это

с тем, что когда мы прокидываем градиенты через, достаточно, длинную последовательность, то мы сталкиваемся с одной из двух проблем: либо градиенты уменьшаются настолько, что ошибки на конце последовательности перестают влиять на ее начало, либо градиенты увеличиваются, и процесс расходится.

Чтобы победить эту проблему было предложено заменить обычные RNN ячейки на более продвинутый вариант LSTM. В базовом варианте LSTM было добавлено еще одно внутреннее состояние ячейки S_t , которое на каждом шаге преобразовывается, и за счет этого гасится проблема исчезающих или взрывающихся градиентов.

Так же в LSTM ячейку были добавлены input и output gate, которые должны были воздействовать на то, что из входных данных окажет влияние на внутреннее состояние ячейки и что из внутреннего состояния должно отправиться на выход. Также после этого начали добавлять forget gate, которое позволяет занулять некоторые компоненты внутреннего состояния $S^{(t-1)}$ прежде чем их передать дальше. В итоге получаем следующую ячейку LSTM сети:

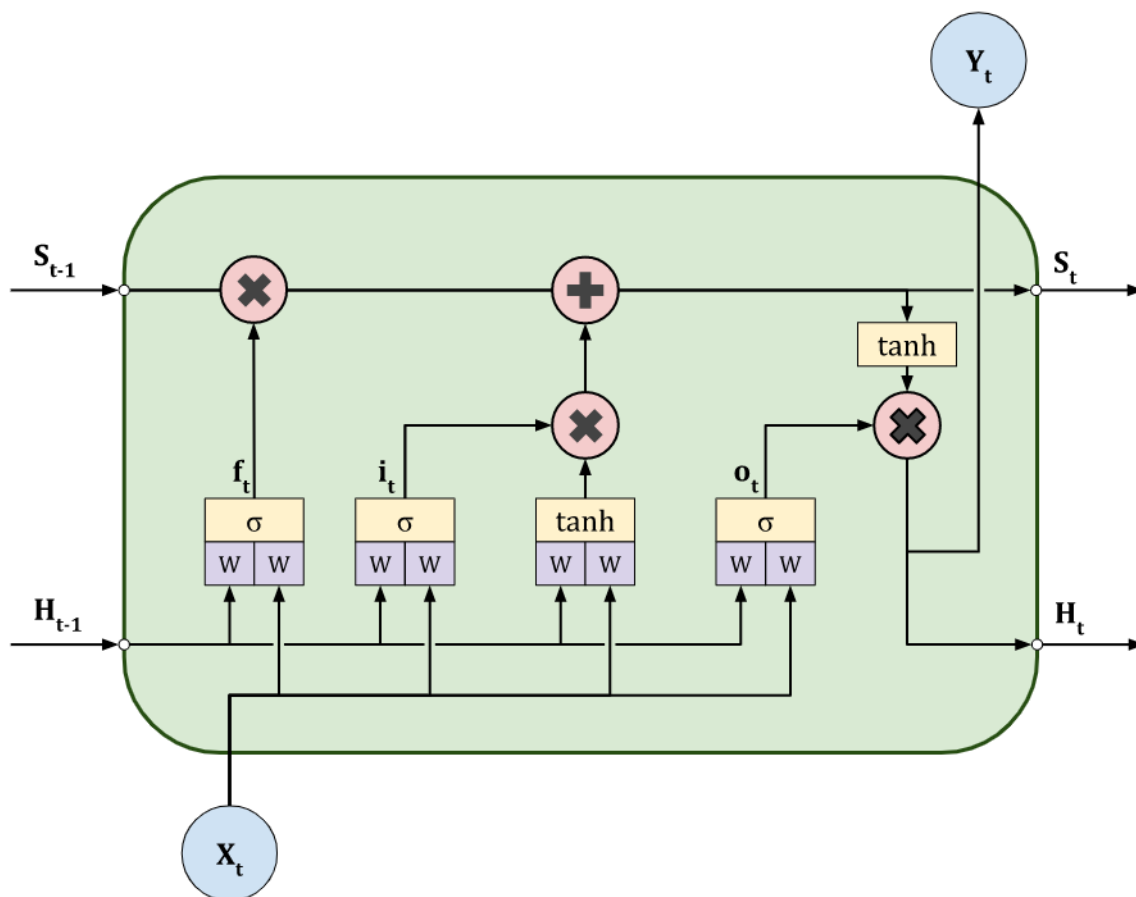


Рис. 4. Ячейка LSTM сети.

Формально обучение происходит следующим образом:

$$g^{(t)} = \phi(W^{gx} X^{(t)} + W^{gh} H^{(t-1)} + b^g)$$

$$i^{(t)} = \sigma(W^{ix} X^{(t)} + W^{ih} H^{(t-1)} + b^i)$$

$$o^{(t)} = \sigma(W^{ox} X^{(t)} + W^{oh} H^{(t-1)} + b^o)$$

$$f^{(t)} = \sigma(W^{fx} X^{(t)} + W^{fh} H^{(t-1)} + b^f)$$

$$S^{(t)} = g^{(t)} * i^{(t)} + S^{(t-1)} * f^{(t)}$$

$$H^{(t)} = \phi(S^{(t)}) * o^{(t)}$$

где $i^{(t)}$, $o^{(t)}$ это соответственно input и output gate,

2.5. Закрытый рекуррентный блок GRU

Еще один вариант ячейки, очень похожий на LSTM, но с меньшим количеством gate (двумя), а значит и параметров (качество при этом обещают сравнимое с LSTM с таким же количеством параметров):

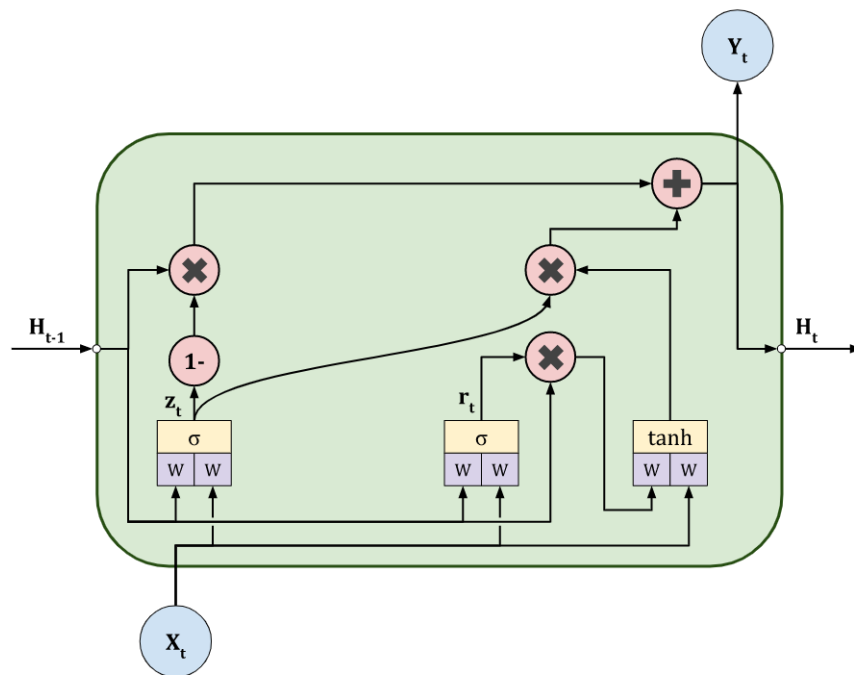


Рис. 5. Ячейка GRU.

Формально обучение происходит следующим образом:

$$\begin{aligned}
z^{(t)} &= \sigma(W^{zx}X^{(t)} + W^{zh}H^{(t-1)+b^z}) \\
r^{(t)} &= \sigma(W^{rx}X^{(t)} + W^{rh}H^{(t-1)+b^r}) \\
\tilde{H}^{(t)} &= \phi(W^{\tilde{h}x}X^t + W^{\tilde{h}h}(r^t * H^{(t-1)}) + b_{\tilde{w}}) \\
H^{(t)} &= (1 - z^{(t)}) * H^{(t-1)} + z^{(t)} * \tilde{h}^{(t)}
\end{aligned}$$

2.6. DeepSpeech

Следующим этапом в рассмотрении известных подходов к решению задачи автоматического распознавания речи является обзор наиболее знаковых систем и архитектур, с помощью которых производится наиболее эффективное распознавание речи в последние годы. Без сомнения к таким архитектурам могут быть отнесены DeepSpeech, DeepSpeech2 и QuartzNet.

DeepSpeech представляет собой сквозную (end-to-end) систему для автоматического распознавания речи, которая является более простой относительно классических подходов (GMM и HMM) и достигает лучших результатов. Основа архитектуры представляет собой рекуррентную нейронную сеть (RNN), которая обучается напрямую из данных и не требует настройки отдельных компонент для адаптации спикера или фильтрации шума. DeepSpeech (2014) опережает по точности ранее опубликованные методы на корпусе *SwitchboardHub5'00* и достигает 16.0% ошибок. Рекуррентная нейронная сеть (RNN) обучена на восприятие спектрограмм и её целью является генерация текстовых транскрипций ко входящим аудио. Пусть одна запись x и её метка (транскрипция) y взяты из тренировочного набора $X = (x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots$. Каждая запись $x^{(i)}$ - это временной ряд длины $T^{(i)}$, где каждый временной срез - это вектор аудио признаков, $x_t^{(i)}$, $t = 1, \dots, T^{(i)}$. В качестве аудиопризнаков используются спектрограммы, таким образом $x_{(t,p)}^{(i)}$ обозначает мощность p -ого элемента по частоте в момент времени t в аудио записи. Целью RNN является преобразование входной последовательности x в вероятностную последовательность символов для транскрипции y ,

$$\hat{y}_t = P(c_t|x), c_t \in a, b, c, \dots, z, space, apostrophe, blank$$

Модель RNN состоит из 5 скрытых состояний. Для входа x скрытые состояния на слое l обозначаются как $h^{(l)}$, то есть что $h^{(0)}$ это входное состояние. Первые три

слоя RNN не являются рекуррентными. Для первого слоя, в каждый момент времени t выходные данные зависят от кадра спектрограммы x_t вместе с контекстом из $C \in \{5, 7, 9\}$ кадров. Второй и третий слои работают с независимыми данными для каждого шага (без использования контекста).

Таким образом, для каждого момента времени t первые три слоя определяются следующим образом:

$$h_t^{(l)} = g(W^{(l)}h_t^{(l-1)} + b^{(l)})$$

где $g(z) = \min(\max(0, z), 20)$ - это ограниченная сверху функция активации ReLu (rectified-linear unit) и $W^{(l)}, b^{(l)}$ - матрица весов и параметр смещения для слоя l . Четвертый слой - это двунаправленный рекуррентный слой. Этот слой содержит два подмножества скрытых состояний: с прямой рекуррентностью $h^{(f)}$ и с обратной $h^{(b)}$:

$$h_t^{(f)} = g(W^{(4)}h_t^{(3)} + W_r^{(f)}h_t^{(f)}t - 1) + b^{(4)}$$

$$h_t^{(b)} = g(W^{(4)}h_t^{(3)} + W_r^{(b)}h_t^{(b)}t + 1) + b^{(4)}$$

Стоит заметить, что h^f должен быть вычислен последовательно с $t = 1$ до $t = T^{(i)}$ для i -ой записи, в то время как $h^{(b)}$ должен быть вычислен последовательно в инверсированном порядке с $t = 1$ до $t = T^{(i)}$. Пятый (не рекуррентный) слой принимает на вход скрытые состояния с прямой и обратной рекуррентностью $h_t^{(4)} = g(W^{(5)}h_t^{(4)} + b^{(5)})$, где $h_t^{(4)} = h_t^{(f)} + h_t^{(b)}$. Результатом данного слоя является стандартная функция softmax, которая вычисляет вероятности предсказанных символов для каждого временного интервала t и символа k из алфавита:

$$h_{t,k}^{(4)} = \hat{y}_{t,k} \equiv P(c_t = k|x) = \frac{\exp(W^{(6)}h_t^{(5)} + b^{(6)})}{\sum_j W_j^{(5)}h_t^{(5)} + b_j^{(6)}}$$

где $W_k^{(6)}$ и $b_k^{(6)}$ обозначают k -ый столбец матрицы весов и k -е смещение (bias), соответственно. Полная модель RNN, использованной в архитектуре DeepSpeech проиллюстрирована на рис. 6.

2.7. DeepSpeech 2

Следующей ступенью эволюции DeepSpeech стала архитектура DeepSpeech 2, которая в своей основе так же имела RNN, кроме того, авторами была рассмотрена

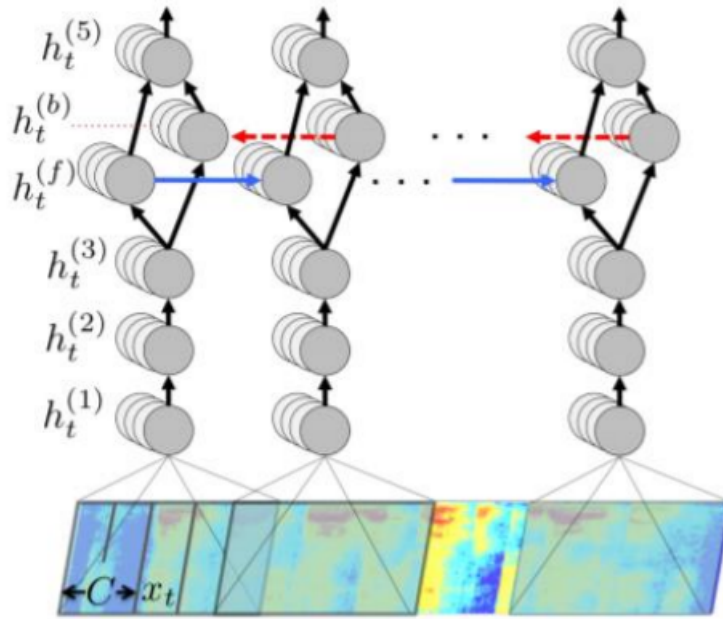


Рис. 6. Архитектура DeepSpeech.

возможность адаптации модели под несколько языков (Английский, Китайский). На рис. 7 представлена архитектура DeepSpeech 2, нетрудно заметить отличия от архитектуры её предшественника – DeepSpeech: добавлены нормализация по батчам (Batch Normalization), инвариантные свёртки к входным спектрограммам и один полносвязный слой перед применением softmax-слоя. Данная модель по-прежнему тренируется с помощью CTC функции потерь, которая позволяет напрямую предсказывать последовательность символов из входного аудио. Кроме того, в качестве основного блока сети были рассмотрены GRU (Gated Recurrent Units) и LSTM (Long Short-Term Memory).

В качестве дальнейшего исследования авторами была выбрана модель GRU, так как в экспериментах с малыми наборами данных GRU и LSTM достигают одинаковой точности со схожим количеством параметров, при этом GRU имеет более высокую скорость обучения и меньшую вероятность расхождения. Также авторы обращают внимание на использование пакетной нормализации (batch normalization) при обучении глубоких рекуррентных нейронных сетей (Deep RNNs). При увеличении числа рекуррентных слоёв в RNN увеличивается время тренировки сети, так как глубина и размер модели увеличивается. Авторы DeepSpeech 2 также провели эксперименты для исследования возможности ускорения тренировки основного блока системы. Результаты тестирования DeepSpeech2 демонстрируют, что система хорошо справляется с аудио, которые представляют собой прочитанный вслух текст (строка Read) и для $\frac{3}{4}$ наборов

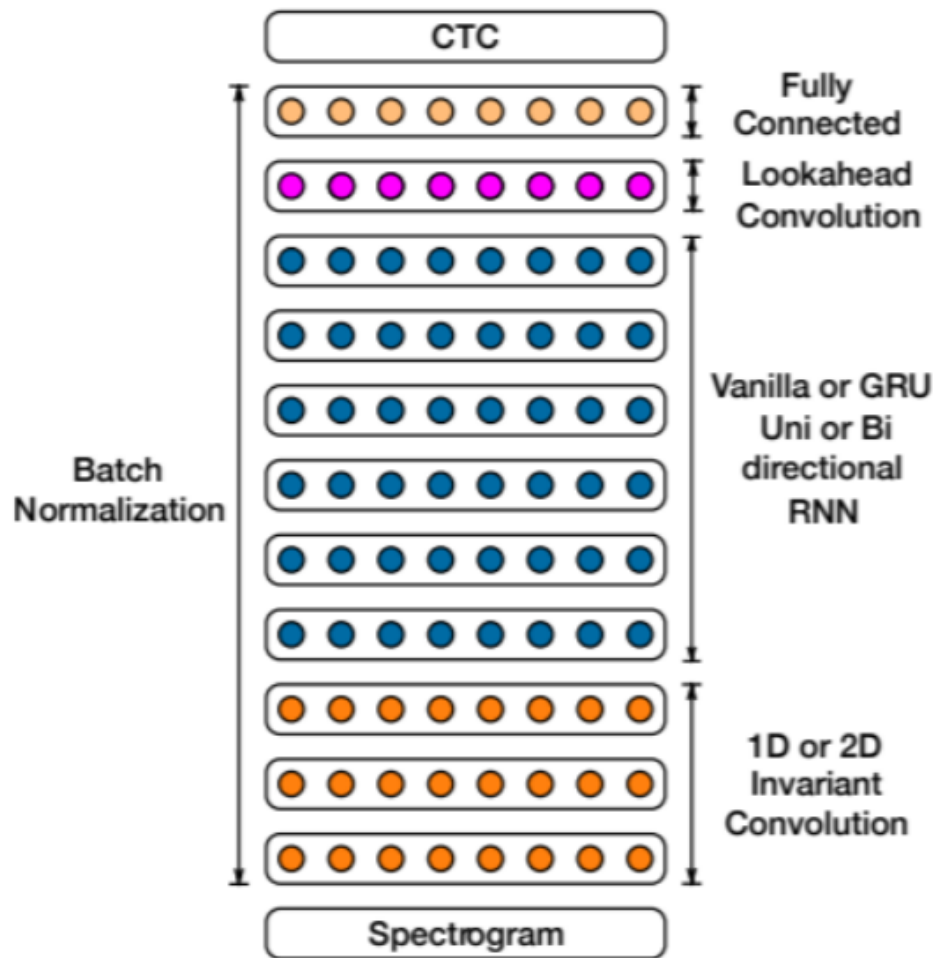


Рис. 7. Архитектура DeepSpeech 2.

данных опережает по точности человека. В то же время, на шумных данных (строка Noisy) модель показывает существенное отставание от результатов показанных человеком, в среднем, на 10% - это говорит о том, что исследование темы автоматического распознавания на шумных данных актуально и на данный момент, так как DeepSpeech 2 по прежнему активно используется для коммерческого распознавания речи и имеет множество адаптаций под разные языки.

2.8. QuartzNet

Следующая архитектура, которая будет рассмотрена имеет название QuartzNet была представлена в 2019 году. QuartzNet представляет собой сквозную (end-to-end) нейро-акустическую модель для автоматического распознавания речи, эта архитектура представляет собой множество блоков, которые соединены остаточными связями (residual connections). Каждый блок состоит из одного или нескольких модулей с одномерными свёртками, разделёнными по времени (1D time-channel separable convolutional layer), пакетной нормализацией и ReLU (rectified-linear unit) слоями. Модель обучена с помощью CTC (Connectionist Temporal Classification) функции ошибок. Основу архитектуры QuartzNet легла свёрточная модель JasperNet (2019). Основным нововведением архитектуры QuartzNet является замена обычных одномерных свёрток на одномерные свёртки разделяемые по времени с реализацией разделения по глубине (depthwise separable convolutions). Польза этой модификации объясняется тем, что одномерные свёртки с разделением на временные каналы могут быть разделены на одномерный свёрточный слой по глубине с длиной ядра K который работает на каждом канале индивидуально, но через K временных кадров, и точечный свёрточный слой, который работает на каждом временном кадре независимо, но по всем каналам.

Модель QuartzNet имеет следующую структуру: свёрточный одномерный слой C_1 , за которым следует последовательность блоков B .

Каждый блок B_i повторяется S_i и имеет остаточные связи (residual connections) между блоками. Каждый блок B_i состоит из одинаковых базовых модулей, которые повторяются R_i раз и включают в себя четыре слоя:

- Глубинный свёрточный слой размера K с каналами c_{out}
- Поточечная свёртка

- Слой нормализации
- ReLU

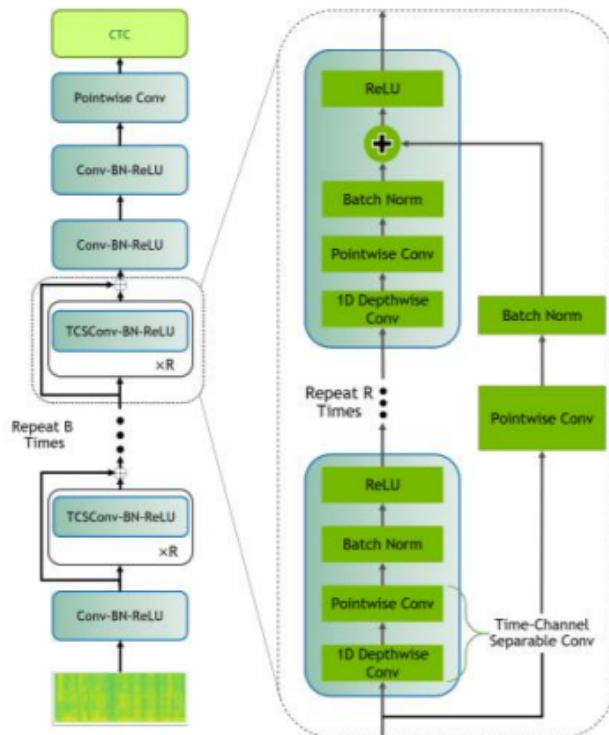


Рис. 8. Архитектура модели QuartzNet

Для нейросетевой акустической модели одним из определяющих параметров является количество параметров, эта характеристика влияет на скорость обучения, время распознавания и адаптивность модели к аппаратному обеспечению с различными ограничениями. Авторы модели QuartzNet предложили модифицировать один из блоков, пронаблюдать за метрикой точности и выяснить: как уменьшение числа параметров влияет на точность акустической модели.

Одномерный свёрточный (1D convolutional layer) слой имеет количество параметров равное произведению $K \times c_{in} \times c_{out}$, где K – размер ядра свёртки, c_{in} – число входных каналов и c_{out} – число выходных каналов. Разделяемые по временным каналам свёртки имеют количество параметров равное $K \times c_{in} + c_{in} \times c_{out}$. Поскольку K в большинстве случаев в несколько раз меньше, чем c_{out} рассматривается возможность использования групповых свёрток для этого модуля и последующий обмен информацией между группами свёрток (рис. 9).

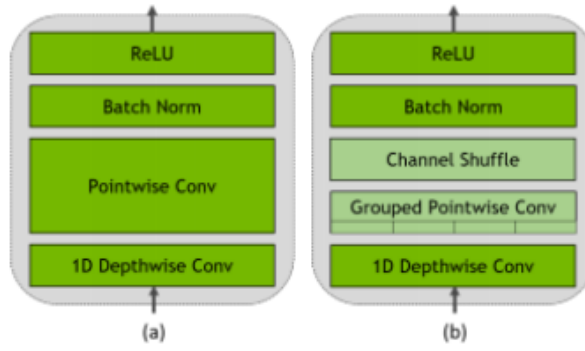


Рис. 9. (a) Разделяемые по временным каналам свертки (b) Разделяемые по временным каналам свертки с разделением по группам и дальнейшим обменом информацией

Применение данного метода модификации для свёрточного модуля позволило уменьшить количество параметров модели более чем в два раза (с 18.9 до 8.7), при этом ухудшение метрики WER для максимального числа групп составило менее 1%. Лучшие

# Groups	dev-clean	dev-other	Params, M
1	3.98	11.58	18.9
2	4.29	12.52	12.1
4	4.51	13.48	8.70

Рис. 10. Результаты тестирования (WER метрика) модели QuartzNet-15×5 с применением групповых свёрток (число групп: 1, 2, 4), для обучения использовался набор данных LibriSpeech, 300 эпох обучения. Тестирование на подвыборках dev (clean and other)

результаты на наборе данных LibriSpeech достигнуты с моделью QuartzNet-15x5, которая состоит из пятнадцати блоков с пятью свёрточными модулями на каждый блок. Модель с разделяемыми по временным каналам свёртками, описанная выше, имеет гораздо меньшее число параметров, чем модель с обыкновенными свёртками и менее склонна к чрезмерному переобучению, поэтому в качестве регуляризации авторы используют только расширение набора данных (data augmentation) и сокращения веса (weight decay).

3. Эксперименты по обучению моделей.

3.1. Данные для распознавания речи LJ Speech Dataset.

LJ Speech Dataset - это общедоступный набор речевых данных, состоящий из 13 100 коротких аудиоклипов, в которых один спикер читает отрывки из 7 научно-популярных книг. Транскрипция предоставляется для каждого клипа. Клипы различаются по длине от 1 до 10 секунд и имеют общую продолжительность примерно 24 часа.

Тексты были опубликованы между 1884 и 1964 годами. Аудио было записано в 2016-17 годах проектом LibriVox и также находится в открытом доступе.

Статистика:

- Всего клипов - 13,100
- Всего слов - 225,715
- Всего персонажей - 1,308,678
- Общая продолжительность - 23:55:17
- Средняя продолжительность клипа - 1,11 сек.
- Отдельные слова - 13,821

3.2. Функция ошибки.

В системах автоматического распознавания речи основным показателем качества является точность распознавания, которая определяется как процент правильно распознанных слов (WRR — Word Recognition Rate) или, наоборот, неправильно распознанных слов (WER — Word Error Rate). Иногда также используется показатель ошибок распознавания фраз/предложений (SER — Sentence Error Rate), который является важным в диалоговых системах, где корректировка гипотезы распознавания невозможна в отличие от задачи диктовки текста. В последнее время в качестве основного показателя точности работы систем распознавания речи используется показатель WER, а именно, его абсолютное значение или относительное, если сравниваются различные модели/системы. Поскольку с развитием речевых технологий показатель WER все более приближается к нулю, то улучшение его значения более наглядно, чем повышение точности распознавания слов. Метод определения показателя WER состоит в выравнивании двух

текстовых строк (первая — это результат распознавания, а вторая — запись того, что было сказано в действительности) с помощью алгоритма динамического программирования с вычислением расстояния Левенштейна. Расстояние Левенштейна представляет собой „стоимость“ редактирования данных (минимальное количество или взвешенная сумма операций редактирования) для преобразования первой строки во вторую с наименьшим числом операций ручной замены (S), удаления (D) и вставки (I) слов:

$$WER = \frac{S + D + I}{T}, WRR = 1 - WER$$

где T — количество слов в распознаваемой фразе. В данной работе мы будем сравнивать модели с помощью WER .

3.3. Тестирование и сравнение моделей.

В качестве первой модели была выбрана модель с двумя сверточными слоями и одним полносвязным слоем. После 50ти эпох обучений с помощью CTC (Connectionist Temporal Classification) функции ошибок модель показала $WER = 0.6$ Что является очень плохим результатом.

Следующим этапом были выбраны модели на основе DeepSpeech:

- 2 сверточных слоя нейронной сети, 2 слоя GRU , обучаемая с помощью CTC функции ошибок и Dropout 50%. Модель уже имеет порядка 5.6 миллиона параметров и после 50ой эпохи имела $WER = 0.37$
- 2 сверточных слоя нейронной сети, 2 слоя $LSTM$, обучаемая с помощью CTC функции ошибок и Dropout 50%. Здесь в качестве эксперимента 2слоя GRU были заменены на слои $LSTM$. Что в теории при большем количестве параметров и более длительном обучении, модель должна была показать результат лучше, чем у предыдущей модели. Модель уже имеет порядка 7.6 миллиона параметров и после 50ой эпохи имела $WER = 0.29$

Позже были обучены современные методы автоматического распознавания речи на основе DeepSpeech 2 и QuartzNet. А именно:

- DeepSpeech 2 с 5ю слоями GRU внутри, обучаемая с помощью CTC функции ошибок, Dropout 50%, а также наложенным шумом. Модель имеет порядка 26 миллионов параметров и после 50ой эпохи имела $WER = 0.17$

- DeepSpeech 2 с 5ю слоями LSTM внутри, обучаемая с помощью CTC функции ошибок, Dropout 50%, а также наложенным шумом. Здесь в качестве эксперимента 5 слоев GRU были заменены на слои LSTM. Модель уже имеет порядка 35 миллиона параметров и после 50ой эпохи имела $WER = 0.13$
- Нейросеть jasper_10x5 с 330 миллионами параметров и $WER = 0.087$
- Нейросеть quartznet15x5 с 19 миллионами параметров и $WER = 0.074$

Модель	Количество параметров	WER
2 Слоя CNN + Dense + CTC Loss + Drop(0.5)	1.9M	0.6
DeepSpeech: 2 Слоя CNN + 2 слоя GRU + CTC Loss + Drop(0.5)	5.6M	0.37
DeepSpeech: 2 Слоя CNN + 2 слоя LSTM + CTC Loss + Drop(0.5)	7.6M	0.29
DeepSpeech 2: 2 Слоя CNN + 5 слоев LSTM + CTC Loss + Drop(0.5) + Noise + shifter + stretch	35M	0.13
DeepSpeech 2: 2 Слоя CNN + 5 слоев GRU + CTC Loss + Drop(0.5) + Noise + shifter + stretch	26M	0.17
quartznet15x5	19M	0.074
jasper_10x5	330M	0.087

Рис. 11. Результаты тестирования моделей.

Результаты тестирования моделей представленные на рис. 11 демонстрируют способность модели QuartzNet достигать современных и релевантных результатов в области автоматического распознавания речи с использованием модели с меньшим числом параметров.

3.4. Возможные улучшения полученных результатов

В качестве улучшения могут быть рассмотрены несколько вариантов: увеличение количества обучающих данных, увеличение вычислительных мощностей, использование нейросетевых моделей в качестве лингвистической модели, например Transformer-XL.

4. Результаты.

Проведённый обзор современных подходов позволяет оценить подходы к решению задачи автоматического распознавания речи с помощью использования нейросетевых акустических моделей. Анализ архитектур рассмотренных моделей (GRU, LSTM, DeepSpeech, DeepSpeech 2 и QuartzNet) и результаты, полученные при их тестировании, позволяют сделать вывод: модель QuartzNet является наиболее подходящей архитектурой для решения задачи автоматического распознавания речи ввиду меньшего числа параметров и при этом высоких показателей при распознавании (*WER*).

В ходе выполнения работы были исследованы современные подходы к решению задачи автоматического распознавания речи, рассмотрены детали реализации этих моделей, обучены несколько вариаций нейронных сетей с последующим сравнением между собой.

Полученные модели могут быть использованы для распознавания речи в реальной среде и доработаны для решения реальных задач в сфере автоматического распознавания речи.

5. Список литературы.

1. Levenshtein V. I. Binary codes capable of correcting deletions, insertions and reversals // *Sov. Phys. Dokl.* 1966. Vol. 6. P. 707–710.
2. Wang, Y.; Acero, A.; Chelba, C. (2003). Is Word Error Rate a Good Indicator for Spoken Language Understanding Accuracy. *IEEE Workshop on Automatic Speech Recognition and Understanding*. St. Thomas, US Virgin Islands. CiteSeerX 10.1.1.89.424.
3. ImageNet Classification with Deep Convolutional Neural Networks / Alex Krizhevsky [et al.]. –November 2013, 9 S.
4. Juang, B. H., Rabiner, L. R. (1991). Hidden Markov models for speech recognition. *Technometrics*, 33(3), 251-272.
5. Hannun, A., Case, C., Casper, J., Catanzaro, B., Diamos, G., Elsen, E., ... Ng, A. Y. (2014). Deep speech: Scaling up end-to-end speech recognition. *arXiv preprint arXiv:1412.5567*.
6. Amodei, D., Ananthanarayanan, S., Anubhai, R., Bai, J., Battenberg, E., Case, C., ... Zhu, Z. (2016, June). Deep speech 2: End-to-end speech recognition in english and mandarin. In *International conference on machine learning* (pp. 173-182). PMLR.
7. Krیمان, S., Beliaev, S., Ginsburg, B., Huang, J., Kuchaiev, O., Lavrukhin, V., ... Zhang, Y. (2020, May). Quartznet: Deep automatic speech recognition with 1d time-channel separable convolutions. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 6124-6128). IEEE.
8. M. Schuster and K. K. Paliwal. Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing*, 45(11):2673–2681, 1997
9. Graves, A., Fernández, S., Gomez, F., and Schmidhuber, J. Connectionist 36 temporal classification: Labelling unsegmented sequence data with recurrent neural networks. In *ICML*, pp. 369–376. ACM, 2006.
10. Hochreiter, Sepp and Schmidhuber, Jürgen. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.

11. A. Graves, S. Fernandez, F. Gomez, and J. Schmidhuber, “Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks,” in ICML, 2006.
12. J. Li, V. Lavrukhin, B. Ginsburg, R. Leary, O. Kuchaiev, J.M. Cohen, H. Nguyen, and R.T. Gadde, “Jasper: An end-to-end convolutional neural acoustic model,” arXiv:1904.03288, 2019
13. Kuchaiev, O., Li, J., Nguyen, H., Hrinchuk, O., Leary, R., Ginsburg, B., ... Cohen, J. M. (2019). Nemo: a toolkit for building ai applications using neural modules. arXiv preprint arXiv:1909.09577.
14. NVIDIA, QuartzNet15x5Base-En, - <https://ngc.nvidia.com/catalog/models/>
15. Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Rafal Jozefowicz, Yangqing Jia, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Mike Schuster, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
16. Bracewell, R. N., Bracewell, R. N. (1986). The Fourier transform and its applications(Vol. 31999, pp. 267-272). New York: McGraw-Hill.

6. Приложение

```
1 gpu_info = !nvidia-smi
2 gpu_info = '\n'.join(gpu_info)
3 if gpu_info.find('failed') >= 0:
4     print('Not connected to a GPU')
5 else:
6     print(gpu_info)
7
8 from psutil import virtual_memory
9 ram_gb = virtual_memory().total / 1e9
10 print('Your runtime has {:.1f} gigabytes of available RAM\n'.format(ram_gb))
11
12 if ram_gb < 20:
13     print('Not using a high-RAM runtime')
14 else:
15     print('You are using a high-RAM runtime!')
16
17 !pip install jiwer
18
19 import pandas as pd
20 import numpy as np
21 import tensorflow as tf
22 from tensorflow import keras
23 from tensorflow.keras import layers
24 import matplotlib.pyplot as plt
25 from IPython import display
26 from jiwer import wer
27 from tensorflow.keras import backend as K
28 from tqdm.notebook import tqdm
29
30 data_url = "https://data.keithito.com/data/speech/LJSpeech-1.1.tar.bz2"
31 data_path = keras.utils.get_file("LJSpeech-1.1", data_url, untar=True)
32 wavs_path = data_path + "/wavs/"
33 metadata_path = data_path + "/metadata.csv"
34
35
36 # Read metadata file and parse it
37 metadata_df = pd.read_csv(metadata_path, sep="|", header=None, quoting=3)
38 metadata_df.columns = ["file_name", "transcription", "normalized_transcription"]
39 metadata_df = metadata_df[["file_name", "normalized_transcription"]]
40 metadata_df = metadata_df.sample(frac=1, random_state = 228).reset_index(drop=True)
41 metadata_df.head(3)
42
43 split = int(len(metadata_df) * 0.85)
44 df_train = metadata_df[:split]
45 df_val = metadata_df[split:]
```

```

46
47 print(f"Size of the training set: {len(df_train)}")
48 print(f"Size of the training set: {len(df_val)}")
49
50 for x in tqdm(df_val['file_name']):
51     !cp /root/.keras/datasets/LJSpeech-1.1/wavs/{x}.wav /content/OpenSeq2Seq/{x}.wav
52
53 wer()
54
55 test = pd.DataFrame()
56 test['wav_filename'] = df_val['file_name'].tolist()
57 test['wav_filename'] = test['wav_filename'] + '.wav'
58 test['wav_filesize'] = 'UNUSED'
59 test['transcript'] = 'UNUSED'
60 test.to_csv('/content/OpenSeq2Seq/test.csv', index = None)
61
62 test
63
64 %tensorflow_version 1.x
65
66 import os
67 from os.path import exists, join, basename, splitext
68
69 git_repo_url = 'https://github.com/NVIDIA/OpenSeq2Seq.git'
70 project_name = splitext(basename(git_repo_url))[0]
71 if not exists(project_name):
72     # clone and install dependencies
73     !git clone -q --depth 1 {git_repo_url}
74     !git checkout e958b7d
75     !pip uninstall -y -q pymc3
76     !pip install --upgrade joblib
77     #!cd {project_name} && pip install -q -r requirements.txt
78     !pip install -q youtube-dl librosa python_speech_features sentencepiece
79     !pip install -q --upgrade gdown
80
81     # create eval config
82     !cp {project_name}/example_configs/speech2text/jasper10x5_LibriSpeech_nvgrad.py {project_name}/
83     !sed -i -e 's/\/data\/librispeech\/librivox-test-clean/test/' {project_name}/conf.py
84     #!sed -i -e 's/# "use_lang/"use_lang/' {project_name}/conf.py
85     !echo 'backend = "librosa"' >> {project_name}/conf.py
86     #!cat {project_name}/conf.py
87     !echo "wav_filename, wav_filesize, transcript" > {project_name}/test.csv
88     !echo "test.wav, UNUSED, UNUSED" >> {project_name}/test.csv
89
90 import sys
91 sys.path.append(project_name)
92 from IPython.display import YouTubeVideo

```

```

93
94 if not exists(join(project_name, 'w2l_log_folder')):
95     import gdown
96     gdown.download('https://drive.google.com/uc?id=1gzGT8HoVNKY1i5HNQTKaSoCu7JHV4siR', 'jasper_10x5
97     !unzip jasper_10x5_dr_sp_nvgrad.zip
98     !mv checkpoint {project_name}/jasper_log_folder
99
100 !cd /content/{project_name} && python run.py --config_file conf.py --mode=infer --infer_output_fi
101
102 answer = pd.read_csv('/content/OpenSeq2Seq/output.txt')
103
104 !unzip jasper_10x5_dr_sp_nvgrad.zip
105
106 import torch
107 quartz_net = torch.load('/content/quartznet15x5_multidataset/JasperEncoder-STEP-243800.pt') 18972
108
109 !wget https://api.ngc.nvidia.com/v2/models/nvidia/jasper_pytorch_ckpt_amp/versions/20.10.0/files/nvidi
110
111 jasper = torch.load('/content/nvidia_jasper_210205.pt')
112 sum(p.numel() for p in jasper['ema_state_dict'].values())
113
114 jasper['ema_state_dict']
115
116 import string
117 df_val['predicted_transcript'] = answer['predicted_transcript'].tolist()
118 wer(df_val['normalized_transcription']).apply(lambda x:x.lower().translate(str.maketrans('', '', s
119
120 #@title
121 import os
122 from os.path import exists, join, basename, splitext
123 from IPython.display import YouTubeVideo
124
125 if not exists('quartznet15x5_multidataset'):
126     # download the pretrained weights
127     !wget -nc -q --show-progress -O quartznet15x5.zip https://api.ngc.nvidia.com/v2/models/nvidia/m
128     !unzip quartznet15x5.zip && mkdir quartznet15x5_multidataset && mv Jasper* quartznet15x5.yaml q
129
130 import json
131 from ruamel.yaml import YAML
132 import nemo
133 import nemo_asr
134
135 WORK_DIR = "/content/quartznet15x5_multidataset"
136 MODEL_YAML = "/content/quartznet15x5_multidataset/quartznet15x5.yaml"
137 CHECKPOINT_ENCODER = "/content/quartznet15x5_multidataset/JasperEncoder-STEP-243800.pt"
138 CHECKPOINT_DECODER = "/content/quartznet15x5_multidataset/JasperDecoderForCTC-STEP-243800.pt"
139 # Set this to True to enable beam search decoder
140 ENABLE_NGRAM = False

```

```

141 # This is only necessary if ENABLE_NGRAM = True. Otherwise, set to empty string
142 LM_PATH = "<PATH_TO_KENLM_BINARY>"
143
144 # Read model YAML
145 yaml = YAML(typ="safe")
146 with open(MODEL_YAML) as f:
147     jasper_model_definition = yaml.load(f)
148 labels = jasper_model_definition['labels']
149
150 # Instantiate necessary Neural Modules
151 # Note that data layer is missing from here
152 neural_factory = nemo.core.NeuralModuleFactory(
153     placement=nemo.core.DeviceType.GPU,
154     backend=nemo.core.Backend.PyTorch)
155 data_preprocessor = nemo_asr.AudioToMelSpectrogramPreprocessor(factory=neural_factory)
156 jasper_encoder = nemo_asr.JasperEncoder(
157     jasper=jasper_model_definition['JasperEncoder']['jasper'],
158     activation=jasper_model_definition['JasperEncoder']['activation'],
159     feat_in=jasper_model_definition['AudioToMelSpectrogramPreprocessor']['features'])
160 jasper_encoder.restore_from(CHECKPOINT_ENCODER, local_rank=0)
161 jasper_decoder = nemo_asr.JasperDecoderForCTC(
162     feat_in=1024,
163     num_classes=len(labels))
164 jasper_decoder.restore_from(CHECKPOINT_DECODER, local_rank=0)
165 greedy_decoder = nemo_asr.GreedyCTCDecoder()
166
167 def wav_to_text(manifest, greedy=True):
168     from ruamel.yaml import YAML
169     yaml = YAML(typ="safe")
170     with open(MODEL_YAML) as f:
171         jasper_model_definition = yaml.load(f)
172     labels = jasper_model_definition['labels']
173
174     # Instantiate necessary neural modules
175     data_layer = nemo_asr.AudioToTextDataLayer(
176         shuffle=False,
177         manifest_filepath=manifest,
178         labels=labels, batch_size=1)
179
180     # Define inference DAG
181     audio_signal, audio_signal_len, _, _ = data_layer()
182     processed_signal, processed_signal_len = data_preprocessor(
183         input_signal=audio_signal,
184         length=audio_signal_len)
185     encoded, encoded_len = jasper_encoder(audio_signal=processed_signal,
186                                         length=processed_signal_len)
187     log_probs = jasper_decoder(encoder_output=encoded)
188     predictions = greedy_decoder(log_probs=log_probs)

```

```

189
190     if ENABLE_NGRAM:
191         print('Running with beam search')
192         beam_predictions = beam_search_with_lm(
193             log_probs=log_probs, log_probs_length=encoded_len)
194         eval_tensors = [beam_predictions]
195
196     if greedy:
197         eval_tensors = [predictions]
198
199     tensors = neural_factory.infer(tensors=eval_tensors)
200     if greedy:
201         from nemo_asr.helpers import post_process_predictions
202         prediction = post_process_predictions(tensors[0], labels)
203     else:
204         prediction = tensors[0][0][0][0][1]
205     return prediction
206
207 def create_manifest(file_path):
208     # create manifest
209     manifest = dict()
210     manifest['audio_filepath'] = file_path
211     manifest['duration'] = 18000
212     manifest['text'] = 'todo'
213     with open(file_path+".json", 'w') as fout:
214         fout.write(json.dumps(manifest))
215     return file_path+".json"
216
217 !cd /content/{project_name} && python run.py --config_file conf.py --mode=infer --infer_output_file
218
219 from tqdm.notebook import tqdm
220 predict = [wav_to_text(create_manifest(f'/content/OpenSeq2Seq/{x}')) for x in tqdm(answer['wav_filepaths'])]
221
222 answer = pd.read_csv('/content/OpenSeq2Seq/output.txt')
223
224 import string
225 df_val['predicted_transcript'] = [ x[0] for x in predict]
226 wer(df_val['normalized_transcription']).apply(lambda x:x.lower().translate(str.maketrans(' ', '', string.punctuation)))
227
228 # The set of characters accepted in the transcription.
229 characters = [x for x in "abcdefghijklmnopqrstuvwxy'?! "]
230 # Mapping characters to integers
231 char_to_num = keras.layers.StringLookup(vocabulary=characters, oov_token="")
232 # Mapping integers back to original characters
233 num_to_char = keras.layers.StringLookup(
234     vocabulary=char_to_num.get_vocabulary(), oov_token="", invert=True
235 )
236

```

```

237 print(
238     f"The vocabulary is: {char_to_num.get_vocabulary()} "
239     f"(size ={char_to_num.vocabulary_size()})"
240 )
241
242 # An integer scalar Tensor. The window length in samples.
243 frame_length = 256
244 # An integer scalar Tensor. The number of samples to step.
245 frame_step = 160
246 # An integer scalar Tensor. The size of the FFT to apply.
247 # If not provided, uses the smallest power of 2 enclosing frame_length.
248 fft_length = 384
249
250
251 def encode_single_sample(wav_file, label):
252     #####
253     ## Process the Audio
254     #####
255     # 1. Read wav file
256     file = tf.io.read_file(wavs_path + wav_file + ".wav")
257     # 2. Decode the wav file
258     audio, _ = tf.audio.decode_wav(file)
259     audio = tf.squeeze(audio, axis=-1)
260     # 3. Change type to float
261     audio = tf.cast(audio, tf.float32)
262     # 4. Get the spectrogram
263     spectrogram = tf.signal.stft(
264         audio, frame_length=frame_length, frame_step=frame_step, fft_length=fft_length
265     )
266     # 5. We only need the magnitude, which can be derived by applying tf.abs
267     spectrogram = tf.abs(spectrogram)
268     spectrogram = tf.math.pow(spectrogram, 0.5)
269     # 6. normalisation
270     means = tf.math.reduce_mean(spectrogram, 1, keepdims=True)
271     stddevs = tf.math.reduce_std(spectrogram, 1, keepdims=True)
272     spectrogram = (spectrogram - means) / (stddevs + 1e-10)
273     #####
274     ## Process the label
275     #####
276     # 7. Convert label to Lower case
277     label = tf.strings.lower(label)
278     # 8. Split the label
279     label = tf.strings.unicode_split(label, input_encoding="UTF-8")
280     # 9. Map the characters in label to numbers
281     label = char_to_num(label)
282     # 10. Return a dict as our model is expecting two inputs
283     return spectrogram, label
284

```

```

285 for i in range(10):
286     sp, label = encode_single_sample(metadata_df.loc[i, 'file_name'], metadata_df.loc[i, 'normali
287     print(sp.shape)
288
289 batch_size = 32
290 # Define the training dataset
291 train_dataset = tf.data.Dataset.from_tensor_slices(
292     (list(df_train["file_name"]), list(df_train["normalized_transcription"])))
293 )
294 train_dataset = (
295     train_dataset.map(encode_single_sample, num_parallel_calls=tf.data.AUTOTUNE)
296     .padded_batch(batch_size)
297     .prefetch(buffer_size=tf.data.AUTOTUNE)
298 )
299 train_dataset = train_dataset.cache()
300
301 # Define the validation dataset
302 validation_dataset = tf.data.Dataset.from_tensor_slices(
303     (list(df_val["file_name"]), list(df_val["normalized_transcription"])))
304 )
305 validation_dataset = (
306     validation_dataset.map(encode_single_sample, num_parallel_calls=tf.data.AUTOTUNE)
307     .padded_batch(batch_size)
308     .prefetch(buffer_size=tf.data.AUTOTUNE)
309 )
310 validation_dataset = validation_dataset.cache()
311
312
313 fig = plt.figure(figsize=(8, 5))
314 for batch in train_dataset.take(1):
315     spectrogram = batch[0][0].numpy()
316     spectrogram = np.array([np.trim_zeros(x) for x in np.transpose(spectrogram)])
317     label = batch[1][0]
318     # Spectrogram
319     label = tf.strings.reduce_join(num_to_char(label)).numpy().decode("utf-8")
320     ax = plt.subplot(2, 1, 1)
321     ax.imshow(spectrogram, vmax=1)
322     ax.set_title(label)
323     ax.axis("off")
324     # Wav
325     file = tf.io.read_file(wavs_path + list(df_train["file_name"])[0] + ".wav")
326     audio, _ = tf.audio.decode_wav(file)
327     audio = audio.numpy()
328     ax = plt.subplot(2, 1, 2)
329     plt.plot(audio)
330     ax.set_title("Signal Wave")
331     ax.set_xlim(0, len(audio))
332     display.display(display.Audio(np.transpose(audio), rate=16000))

```

```

333 plt.show()
334
335 def CTCLoss(y_true, y_pred):
336     # Compute the training-time loss value
337     batch_len = tf.cast(tf.shape(y_true)[0], dtype="int64")
338     input_length = tf.cast(tf.shape(y_pred)[1], dtype="int64")
339     label_length = tf.cast(tf.shape(y_true)[1], dtype="int64")
340
341     input_length = input_length * tf.ones(shape=(batch_len, 1), dtype="int64")
342     label_length = label_length * tf.ones(shape=(batch_len, 1), dtype="int64")
343
344     loss = keras.backend.ctc_batch_cost(y_true, y_pred, input_length, label_length)
345     return loss
346
347 from tensorflow.keras.layers import Conv1D, MaxPooling1D, Dropout, \
348 BatchNormalization, GlobalAveragePooling1D, Softmax, Add, Dense, Activation
349
350 def build_model(input_dim, output_dim, lstm_layers=5, lstm_units=128):
351     """Model similar to 2CNN + Dence"""
352     # Model's input
353     input_spectrogram = layers.Input((None, input_dim), name="input")
354     # Expand the dimension to use 2D CNN.
355     x = layers.Reshape((-1, input_dim, 1), name="expand_dim")(input_spectrogram)
356     # Convolution layer 1
357     x = layers.Conv2D(
358         filters=32,
359         kernel_size=[11, 41],
360         strides=[2, 2],
361         padding="same",
362         use_bias=False,
363         name="conv_1",
364     )(x)
365     x = layers.BatchNormalization(name="conv_1_bn")(x)
366     x = layers.ReLU(name="conv_1_relu")(x)
367     # Convolution layer 2
368     x = layers.Conv2D(
369         filters=32,
370         kernel_size=[11, 21],
371         strides=[1, 2],
372         padding="same",
373         use_bias=False,
374         name="conv_2",
375     )(x)
376     x = layers.BatchNormalization(name="conv_2_bn")(x)
377     x = layers.ReLU(name="conv_2_relu")(x)
378     # Reshape the resulted volume to feed the RNNs layers
379     x = layers.Reshape((-1, x.shape[-2] * x.shape[-1]))(x)
380     # RNN layers

```



```

381     for i in range(1, lstm_layers + 1):
382         recurrent = layers.GRU(
383             units=lstm_units,
384             activation="tanh",
385             recurrent_activation="sigmoid",
386             use_bias=True,
387             return_sequences=True,
388             # reset_after=True,
389             name=f"lstm_{i}",
390         )
391         # x = layers.Bidirectional(
392             # recurrent, name=f"bidirectional_{i}", merge_mode="concat"
393         # )(x)
394         # if i < lstm_layers:
395             # x = layers.Dropout(rate=0.5)(x)
396         # Dense layer
397         x = layers.Dense(units=lstm_units * 2, name="dense_1")(x)
398         x = layers.ReLU(name="dense_1_relu")(x)
399         x = layers.Dropout(rate=0.5)(x)
400         # Classification layer
401         output = layers.Dense(units=output_dim + 1, activation="softmax")(x)
402         # Model
403         model = keras.Model(input_spectrogram, output, name="2CNN + Dence")
404         # Optimizer
405         opt = keras.optimizers.Adam(learning_rate=1e-4)
406         # Compile the model and return
407         model.compile(optimizer=opt, loss=CTCLoss)
408         return model
409
410
411     # Get the model
412     model = build_model(
413         input_dim=fft_length // 2 + 1,
414         output_dim=char_to_num.vocabulary_size(),
415         lstm_units=512,
416     )
417     model.summary(line_length=110)
418
419
420     # A utility function to decode the output of the network
421     def decode_batch_predictions(pred):
422         input_len = np.ones(pred.shape[0]) * pred.shape[1]
423         # Use greedy search. For complex tasks, you can use beam search
424         results = keras.backend.ctc_decode(pred, input_length=input_len, greedy=True)[0][0]
425         # Iterate over the results and get back the text
426         output_text = []
427         for result in results:
428             result = tf.strings.reduce_join(num_to_char(result)).numpy().decode("utf-8")

```

```

429         output_text.append(result)
430     return output_text
431
432
433 # A callback class to output a few transcriptions during training
434 class CallbackEval(keras.callbacks.Callback):
435     """Displays a batch of outputs after every epoch."""
436
437     def __init__(self, dataset):
438         super().__init__()
439         self.dataset = dataset
440
441     def on_epoch_end(self, epoch: int, logs=None):
442         predictions = []
443         targets = []
444         for batch in self.dataset:
445             X, y = batch
446             batch_predictions = model.predict(X)
447             batch_predictions = decode_batch_predictions(batch_predictions)
448             predictions.extend(batch_predictions)
449             for label in y:
450                 label = (
451                     tf.strings.reduce_join(num_to_char(label)).numpy().decode("utf-8")
452                 )
453                 targets.append(label)
454         wer_score = wer(targets, predictions)
455         print("-" * 100)
456         print(f"Word Error Rate: {wer_score:.4f}")
457         print("-" * 100)
458         for i in np.random.randint(0, len(predictions), 2):
459             print(f"Target      : {targets[i]}")
460             print(f"Prediction: {predictions[i]}")
461             print("-" * 100)
462
463
464 # Define the number of epochs.
465 epochs = 50
466 # Callback function to check transcription on the val set.
467 validation_callback = CallbackEval(validation_dataset)
468 # Train the model
469 history = model.fit(
470     train_dataset,
471     validation_data=validation_dataset,
472     epochs=epochs,
473     callbacks=[validation_callback],
474 )
475
476 def build_model(input_dim, output_dim, lstm_layers=2, lstm_units=128):

```

```

477 """Model similar to DeepSpeech. GRU"""
478 # Model's input
479 input_spectrogram = layers.Input((None, input_dim), name="input")
480 # Expand the dimension to use 2D CNN.
481 x = layers.Reshape((-1, input_dim, 1), name="expand_dim")(input_spectrogram)
482 # Convolution layer 1
483 x = layers.Conv2D(
484     filters=32,
485     kernel_size=[11, 41],
486     strides=[2, 2],
487     padding="same",
488     use_bias=False,
489     name="conv_1",
490 )(x)
491 x = layers.BatchNormalization(name="conv_1_bn")(x)
492 x = layers.ReLU(name="conv_1_relu")(x)
493 # Convolution layer 2
494 x = layers.Conv2D(
495     filters=32,
496     kernel_size=[11, 21],
497     strides=[1, 2],
498     padding="same",
499     use_bias=False,
500     name="conv_2",
501 )(x)
502 x = layers.BatchNormalization(name="conv_2_bn")(x)
503 x = layers.ReLU(name="conv_2_relu")(x)
504 # Reshape the resulted volume to feed the RNNs layers
505 x = layers.Reshape((-1, x.shape[-2] * x.shape[-1]))(x)
506 # RNN layers
507 for i in range(1, lstm_layers + 1):
508     x = layers.GRU(
509         units=lstm_units,
510         activation="tanh",
511         recurrent_activation="sigmoid",
512         use_bias=True,
513         return_sequences=True,
514         # reset_after=True,
515         name=f"lstm_{i}",
516     )(x)
517     if i < lstm_layers:
518         x = layers.Dropout(rate=0.5)(x)
519 # Dense layer
520 x = layers.Dense(units=lstm_units * 2, name="dense_1")(x)
521 x = layers.ReLU(name="dense_1_relu")(x)
522 x = layers.Dropout(rate=0.5)(x)
523 # Classification layer
524 output = layers.Dense(units=output_dim + 1, activation="softmax")(x)

```

```

525     # Model
526     model = keras.Model(input_spectrogram, output, name="DeepSpeech. GRU")
527     # Optimizer
528     opt = keras.optimizers.Adam(learning_rate=1e-4)
529     # Compile the model and return
530     model.compile(optimizer=opt, loss=CTCLoss)
531     return model
532
533
534     # Get the model
535     model = build_model(
536         input_dim=fft_length // 2 + 1,
537         output_dim=char_to_num.vocabulary_size(),
538         lstm_units=512,
539     )
540     model.summary(line_length=110)
541
542     # Define the number of epochs.
543     epochs = 50
544     # Callback function to check transcription on the val set.
545     validation_callback = CallbackEval(validation_dataset)
546     # Train the model
547     history = model.fit(
548         train_dataset,
549         validation_data=validation_dataset,
550         epochs=epochs,
551         callbacks=[validation_callback],
552     )
553
554     def build_model(input_dim, output_dim, lstm_layers=2, lstm_units=128):
555         """Model similar to DeepSpeech. LSTM"""
556         # Model's input
557         input_spectrogram = layers.Input((None, input_dim), name="input")
558         # Expand the dimension to use 2D CNN.
559         x = layers.Reshape((-1, input_dim, 1), name="expand_dim")(input_spectrogram)
560         # Convolution layer 1
561         x = layers.Conv2D(
562             filters=32,
563             kernel_size=[11, 41],
564             strides=[2, 2],
565             padding="same",
566             use_bias=False,
567             name="conv_1",
568         )(x)
569         x = layers.BatchNormalization(name="conv_1_bn")(x)
570         x = layers.ReLU(name="conv_1_relu")(x)
571         # Convolution layer 2
572         x = layers.Conv2D(

```

```

573     filters=32,
574     kernel_size=[11, 21],
575     strides=[1, 2],
576     padding="same",
577     use_bias=False,
578     name="conv_2",
579 ) (x)
580 x = layers.BatchNormalization(name="conv_2_bn")(x)
581 x = layers.ReLU(name="conv_2_relu")(x)
582 # Reshape the resulted volume to feed the RNNs layers
583 x = layers.Reshape((-1, x.shape[-2] * x.shape[-1]))(x)
584 # RNN layers
585 x = layers.Dense(units=lstm_units * 2, name="dense_1")(x)
586 x = layers.ReLU(name="dense_1_relu")(x)
587 x = layers.Dropout(rate=0.5)(x)
588 for i in range(1, lstm_layers + 1):
589     x = layers.LSTM(
590         units=lstm_units,
591         activation="tanh",
592         recurrent_activation="sigmoid",
593         use_bias=True,
594         return_sequences=True,
595         # reset_after=True,
596         name=f"lstm_{i}",
597     )(x)
598     # x = layers.Bidirectional(
599         # recurrent, name=f"bidirectional_{i}", merge_mode="concat"
600     # )(x)
601     if i < lstm_layers:
602         x = layers.Dropout(rate=0.5)(x)
603     # Dense layer
604     x = layers.Dense(units=lstm_units * 2, name="dense_2")(x)
605     x = layers.ReLU(name="dense_2_relu")(x)
606     x = layers.Dropout(rate=0.5)(x)
607     # Classification layer
608     output = layers.Dense(units=output_dim + 1, activation="softmax")(x)
609     # Model
610     model = keras.Model(input_spectrogram, output, name="DeepSpeech. LSTM")
611     # Optimizer
612     opt = keras.optimizers.Adam(learning_rate=1e-4)
613     # Compile the model and return
614     model.compile(optimizer=opt, loss=CTCLoss)
615     return model
616
617
618 # Get the model
619 model = build_model(
620     input_dim=fft_length // 2 + 1,

```

```

621     output_dim=char_to_num.vocabulary_size(),
622     lstm_units=512,
623 )
624 model.summary(line_length=110)
625
626 # Define the number of epochs.
627 epochs = 50
628 # Callback function to check transcription on the val set.
629 validation_callback = CallbackEval(validation_dataset)
630 # Train the model
631 history = model.fit(
632     train_dataset,
633     validation_data=validation_dataset,
634     epochs=epochs,
635     callbacks=[validation_callback],
636 )
637
638 def build_model(input_dim, output_dim, rnn_layers=5, rnn_units=128):
639     """Model similar to DeepSpeech2. LSTM"""
640     # Model's input
641     input_spectrogram = layers.Input((None, input_dim), name="input")
642     # Expand the dimension to use 2D CNN.
643     x = layers.Reshape((-1, input_dim, 1), name="expand_dim")(input_spectrogram)
644     # Convolution layer 1
645     x = layers.Conv2D(
646         filters=32,
647         kernel_size=[11, 41],
648         strides=[2, 2],
649         padding="same",
650         use_bias=False,
651         name="conv_1",
652     )(x)
653     x = layers.BatchNormalization(name="conv_1_bn")(x)
654     x = layers.ReLU(name="conv_1_relu")(x)
655     # Convolution layer 2
656     x = layers.Conv2D(
657         filters=32,
658         kernel_size=[11, 21],
659         strides=[1, 2],
660         padding="same",
661         use_bias=False,
662         name="conv_2",
663     )(x)
664     x = layers.BatchNormalization(name="conv_2_bn")(x)
665     x = layers.ReLU(name="conv_2_relu")(x)
666     # Reshape the resulted volume to feed the RNNs layers
667     x = layers.Reshape((-1, x.shape[-2] * x.shape[-1]))(x)
668     # RNN layers

```

```

669     for i in range(1, rnn_layers + 1):
670         recurrent = layers.LSTM(
671             units=rnn_units,
672             activation="tanh",
673             recurrent_activation="sigmoid",
674             use_bias=True,
675             return_sequences=True,
676             reset_after=True,
677             name=f"gru_{i}",
678         )
679         x = layers.Bidirectional(
680             recurrent, name=f"bidirectional_{i}", merge_mode="concat"
681         )(x)
682         if i < rnn_layers:
683             x = layers.Dropout(rate=0.5)(x)
684         # Dense layer
685         x = layers.Dense(units=rnn_units * 2, name="dense_1")(x)
686         x = layers.ReLU(name="dense_1_relu")(x)
687         x = layers.Dropout(rate=0.5)(x)
688         # Classification layer
689         output = layers.Dense(units=output_dim + 1, activation="softmax")(x)
690         # Model
691         model = keras.Model(input_spectrogram, output, name="DeepSpeech2. LSTM")
692         # Optimizer
693         opt = keras.optimizers.Adam(learning_rate=1e-4)
694         # Compile the model and return
695         model.compile(optimizer=opt, loss=CTCLoss)
696         return model
697
698 K.clear_session()
699
700 # Get the model
701 model = build_model(
702     input_dim=fft_length // 2 + 1,
703     output_dim=char_to_num.vocabulary_size(),
704     rnn_units=256,
705     rnn_layers = 2
706 )
707
708 model.summary(line_length=110)
709
710 # Define the number of epochs.
711 epochs = 50
712 # Callback function to check transcription on the val set.
713 validation_callback = CallbackEval(validation_dataset)
714 # Train the model
715 history = model.fit(
716     train_dataset,

```

```

717     validation_data=validation_dataset,
718     epochs=epochs,
719     callbacks=[validation_callback],
720     verbose = 1
721 )
722
723 def build_model(input_dim, output_dim, rnn_layers=5, rnn_units=128):
724     """Model similar to DeepSpeech2. GRU"""
725     # Model's input
726     input_spectrogram = layers.Input((None, input_dim), name="input")
727     # Expand the dimension to use 2D CNN.
728     x = layers.Reshape((-1, input_dim, 1), name="expand_dim")(input_spectrogram)
729     # Convolution layer 1
730     x = layers.Conv2D(
731         filters=32,
732         kernel_size=[11, 41],
733         strides=[2, 2],
734         padding="same",
735         use_bias=False,
736         name="conv_1",
737     )(x)
738     x = layers.BatchNormalization(name="conv_1_bn")(x)
739     x = layers.ReLU(name="conv_1_relu")(x)
740     # Convolution layer 2
741     x = layers.Conv2D(
742         filters=32,
743         kernel_size=[11, 21],
744         strides=[1, 2],
745         padding="same",
746         use_bias=False,
747         name="conv_2",
748     )(x)
749     x = layers.BatchNormalization(name="conv_2_bn")(x)
750     x = layers.ReLU(name="conv_2_relu")(x)
751     # Reshape the resulted volume to feed the RNNs layers
752     x = layers.Reshape((-1, x.shape[-2] * x.shape[-1]))(x)
753     # RNN layers
754     for i in range(1, rnn_layers + 1):
755         recurrent = layers.GRU(
756             units=rnn_units,
757             activation="tanh",
758             recurrent_activation="sigmoid",
759             use_bias=True,
760             return_sequences=True,
761             reset_after=True,
762             name=f"gru_{i}",
763         )
764     x = layers.Bidirectional(

```



```

765         recurrent, name=f"bidirectional_{i}", merge_mode="concat"
766     )(x)
767     if i < rnn_layers:
768         x = layers.Dropout(rate=0.5)(x)
769     # Dense layer
770     x = layers.Dense(units=rnn_units * 2, name="dense_1")(x)
771     x = layers.ReLU(name="dense_1_relu")(x)
772     x = layers.Dropout(rate=0.5)(x)
773     # Classification layer
774     output = layers.Dense(units=output_dim + 1, activation="softmax")(x)
775     # Model
776     model = keras.Model(input_spectrogram, output, name="DeepSpeech2. GRU")
777     # Optimizer
778     opt = keras.optimizers.Adam(learning_rate=1e-4)
779     # Compile the model and return
780     model.compile(optimizer=opt, loss=CTCLoss)
781     return model
782
783 K.clear_session()
784
785 # Get the model
786 model = build_model(
787     input_dim=fft_length // 2 + 1,
788     output_dim=char_to_num.vocabulary_size(),
789     rnn_units=256,
790     rnn_layers = 2
791 )
792
793 model.summary(line_length=110)
794
795 # Define the number of epochs.
796 epochs = 50
797 # Callback function to check transcription on the val set.
798 validation_callback = CallbackEval(validation_dataset)
799 # Train the model
800 history = model.fit(
801     train_dataset,
802     validation_data=validation_dataset,
803     epochs=epochs,
804     callbacks=[validation_callback],
805     verbose = 1
806 )

```