

# Распознавание речи обзор задачи: СТС подход на примере QuartzNet

На основе статьи:  
QUARTZNET: DEEP AUTOMATIC SPEECH RECOGNITION WITH 1D  
TIME-CHANNEL SEPARABLE CONVOLUTIONS, 2019 г.

Щербакова Светлана, ММОВС21\_3

**Высшая школа экономики, 2022**

# План

Voice technologies

История

ASR

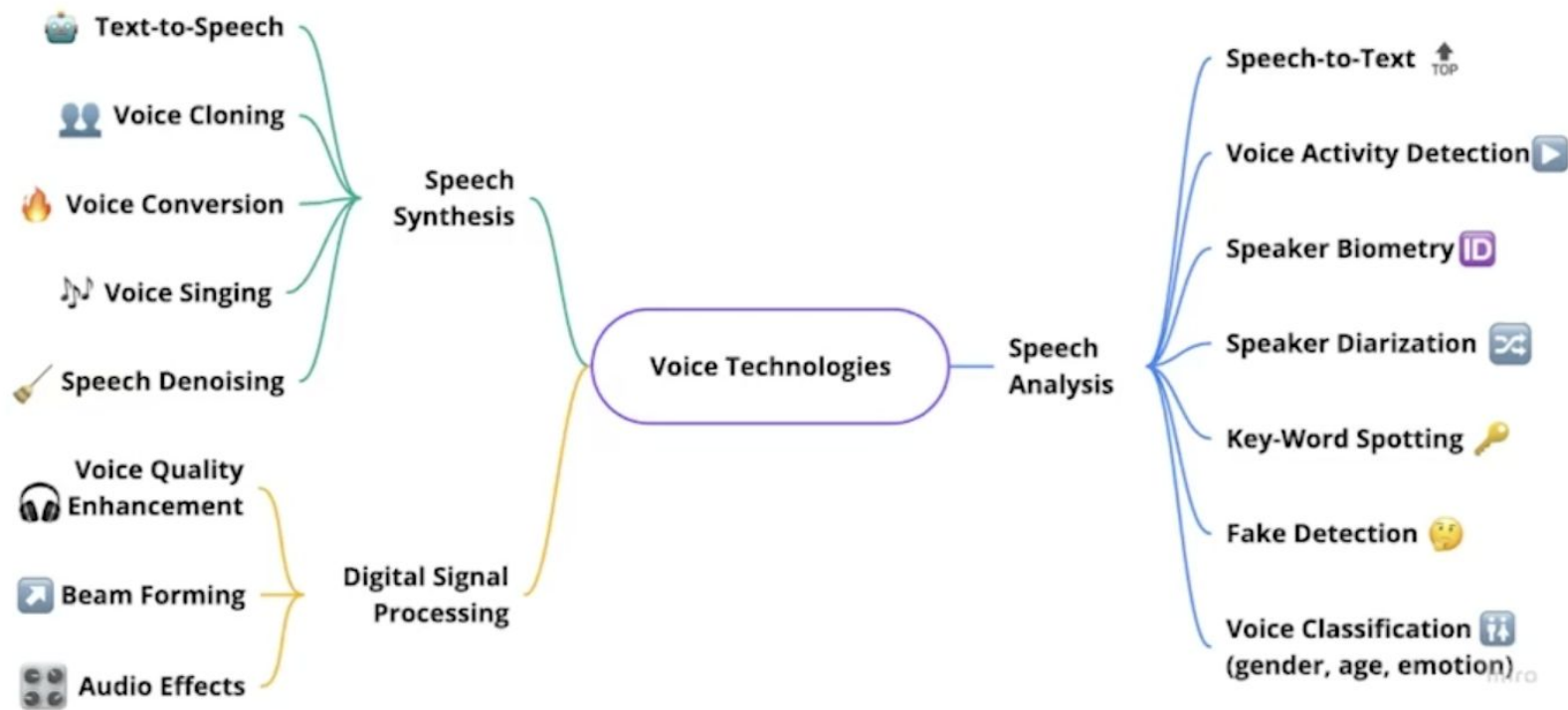
Преобразование сигнала

Преобразование Фурье

СТС

QuartzNet и сепарабельные свертки

# Voice technologies



# HISTORY



1784

Wolfgang von Kempelen creates the **Acoustic-Mechanical Speech Machine** in Vienna

**Thomas Edison** invents the first dictation machine

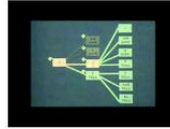


1952

Bell Labs releases **Audrey**, capable of recognizing spoken digits with **90% accuracy** - but only when spoken by its inventor



**IBM Shoebox** can understand 16 English words



1971

**Harpy**, created at Carnegie Mellon University, can comprehend **1,011 words** - and some phrases



**IBM Tangora**, using the Hidden Markov Model, predicts upcoming phonemes in speech



1986



2006

The **National Security Agency (NSA)** starts using speech recognition to isolate key words in recorded speech

**Google** launches a voice search app, bringing speech recognition to mobile devices



2011

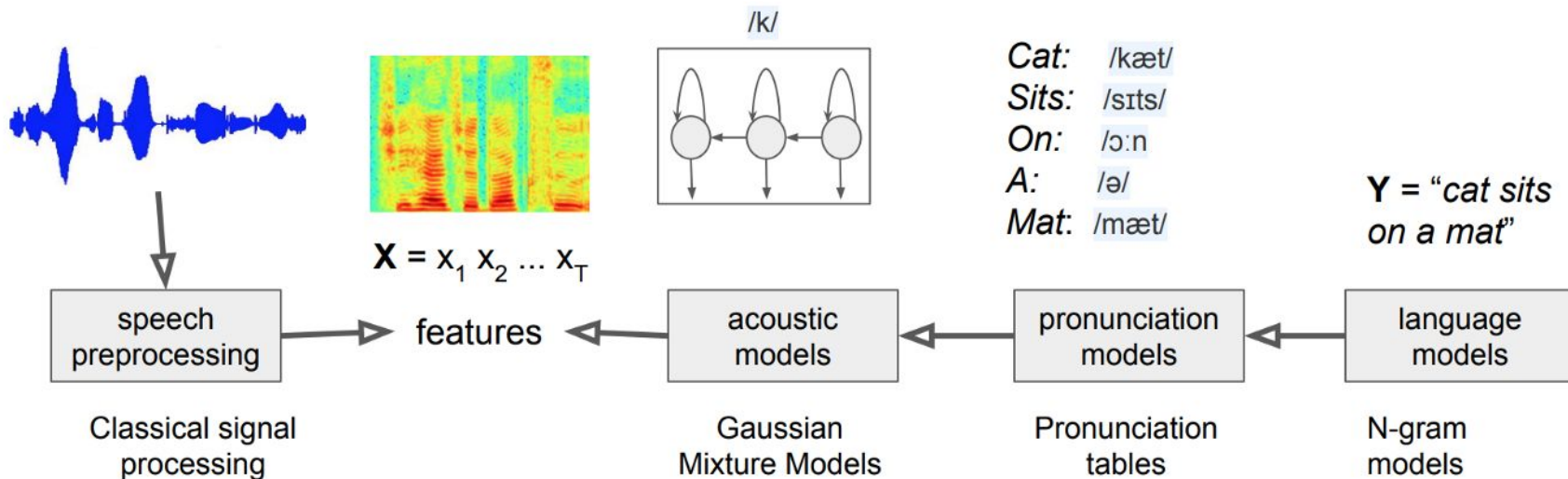
2008



**Apple** announces **Siri**, ushering in the age of the voice-enabled digital assistant

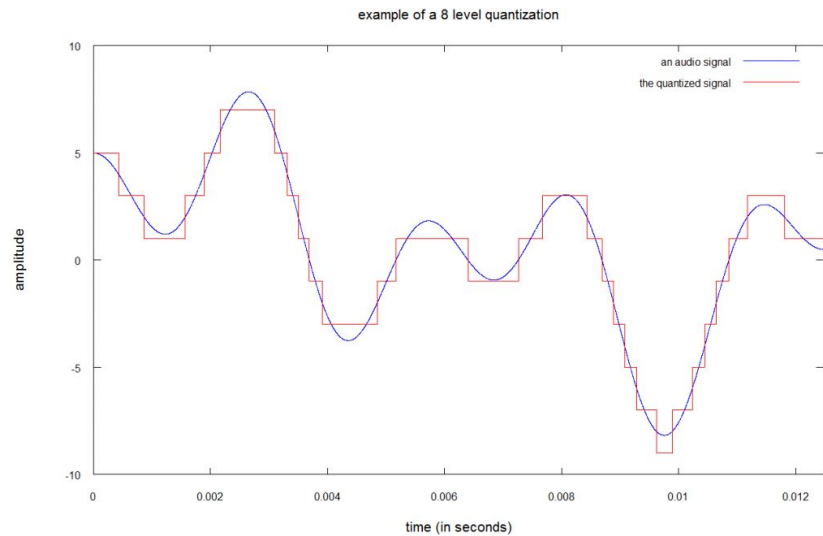
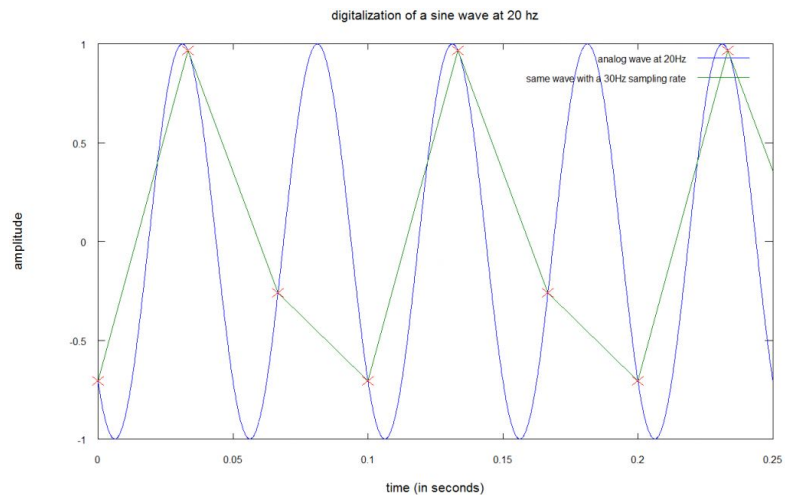


# ASR

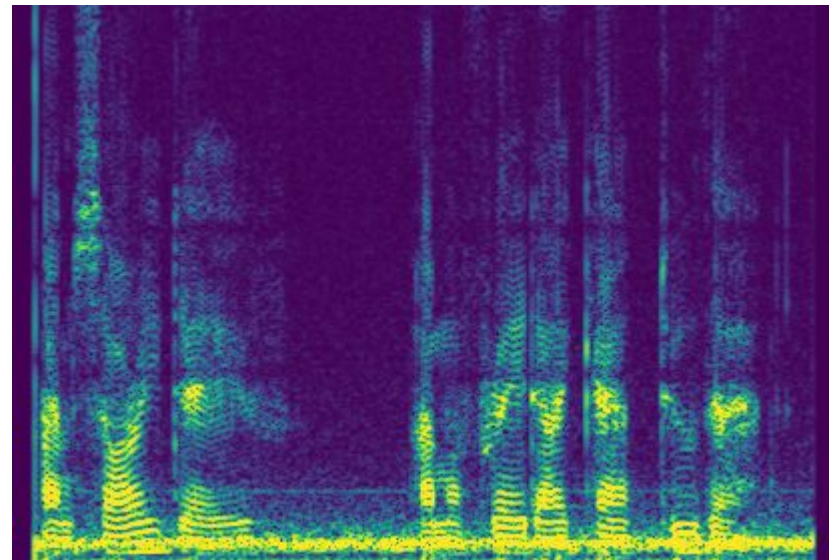
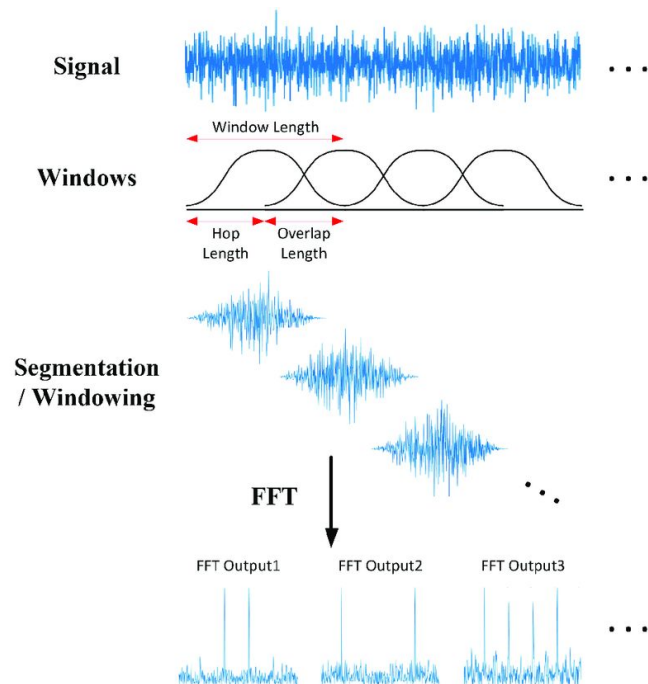


$$P(W_n | W_1, W_2, \dots, W_{n-1}) \approx P(W_n | W_{n-1}) \quad b_j(\mathbf{x}) = p(\mathbf{x} | S=j) = \sum_{m=1}^M c_{jm} \mathcal{N}(\mathbf{x}; \mu_{jm}, \Sigma_{jm})$$

# Преобразование сигнала



# Преобразование Фурье





# МЕЛ Спектрограмма

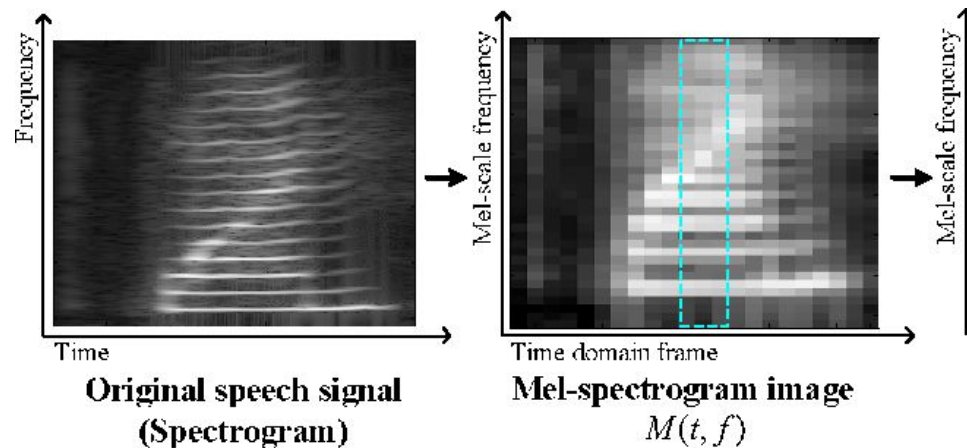
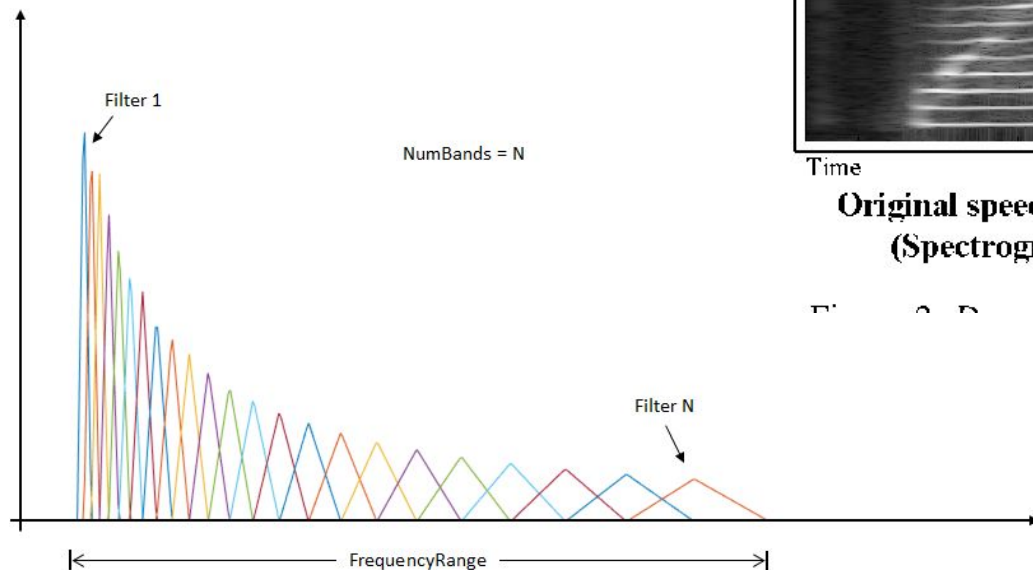
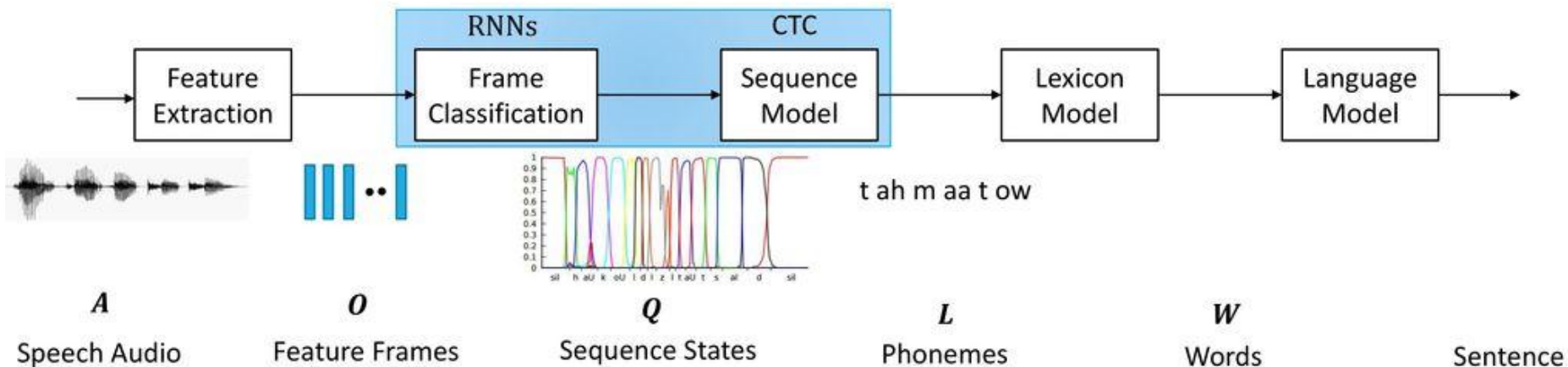


Figure 2.2.1. Generating Mel-spectrogram image

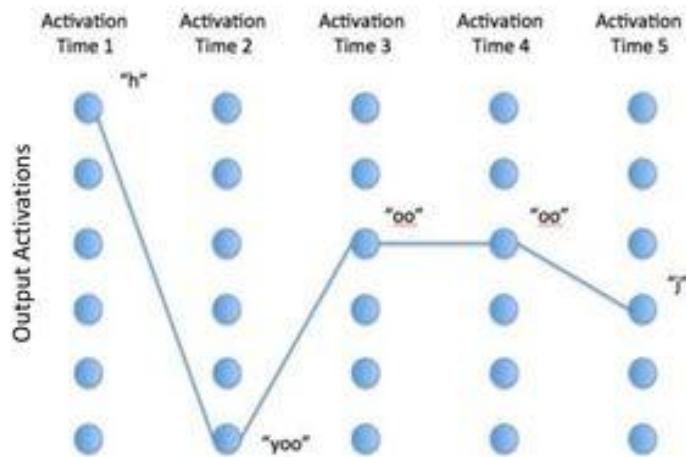
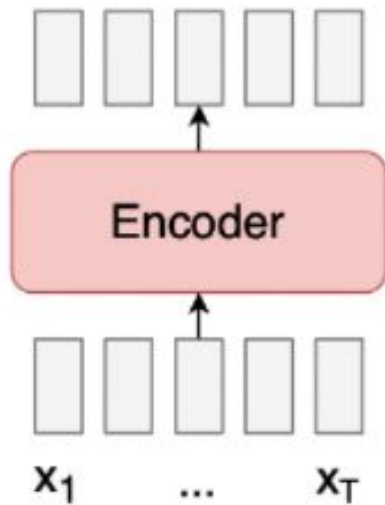
# CTC

RNN: Recurrent Neural Networks

CTC: Connectionist Temporal Classification



# CTC



h h e  $\epsilon$   $\epsilon$  l l l  $\epsilon$  l l o

h e  $\epsilon$  l  $\epsilon$  l o

h e l l o

h e l l o

First, merge repeat characters.

Then, remove any  $\epsilon$  tokens.

The remaining characters are the output.

# LOSS CTC

the ground truth of the word sequence

acoustic frames

$$\mathcal{L}_{CTC} = -\log P(\mathbf{S}|\mathbf{X})$$

$$P(\mathbf{S}|\mathbf{X}) = \sum_{\mathbf{c} \in A(\mathbf{S})} P(\mathbf{C}|\mathbf{X})$$

sum over all possible paths

(e.g. cceaaett, cccεaett, cεaetttt, ...)

$$P(\mathbf{C}|\mathbf{X}) = \prod_{t=1}^T y(c_t, t)$$

joint probability of a path

(e.g. cceaaett)

# LOSS CTC

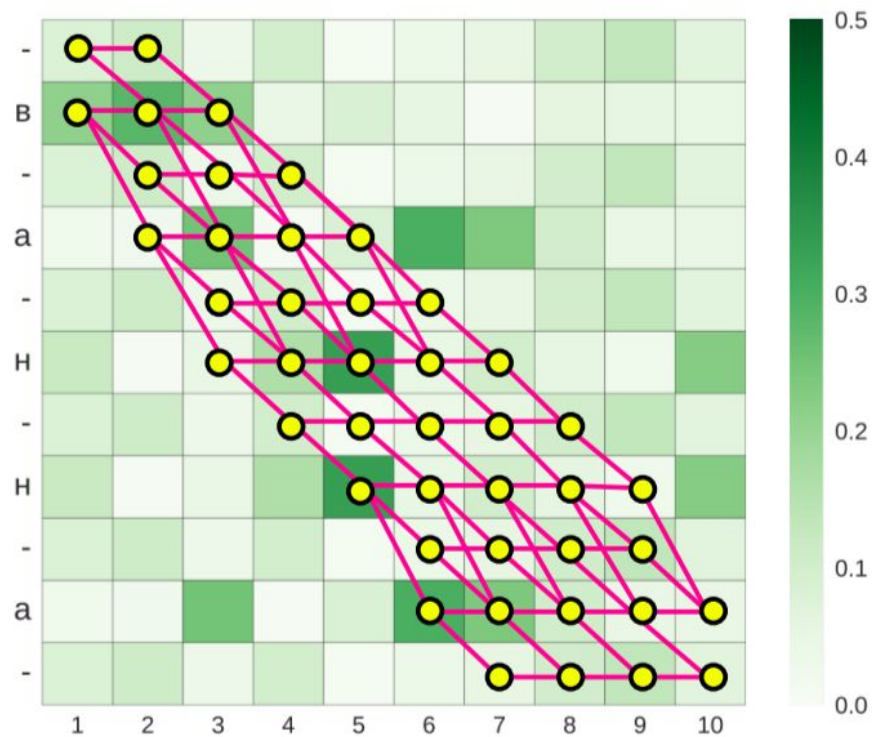
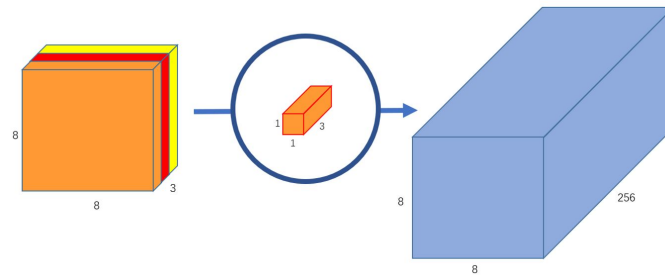
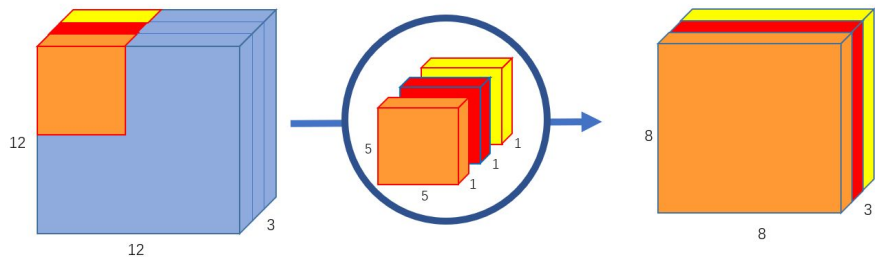
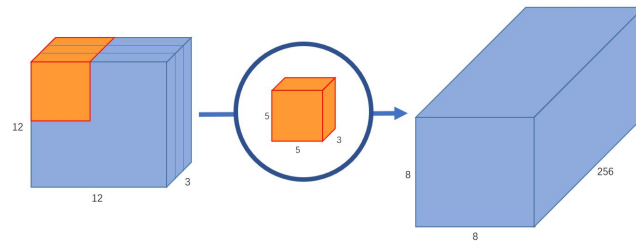


Табл. 2.3 — Всевозможные пути, приводящие за время  $T$  к результату  $l$  при преобразовании их функцией  $B$ .

# Сепарабельные свертки

$$\begin{bmatrix} 3 & 6 & 9 \\ 4 & 8 & 12 \\ 5 & 10 & 15 \end{bmatrix} = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix} \times \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}$$



# QuartzNet

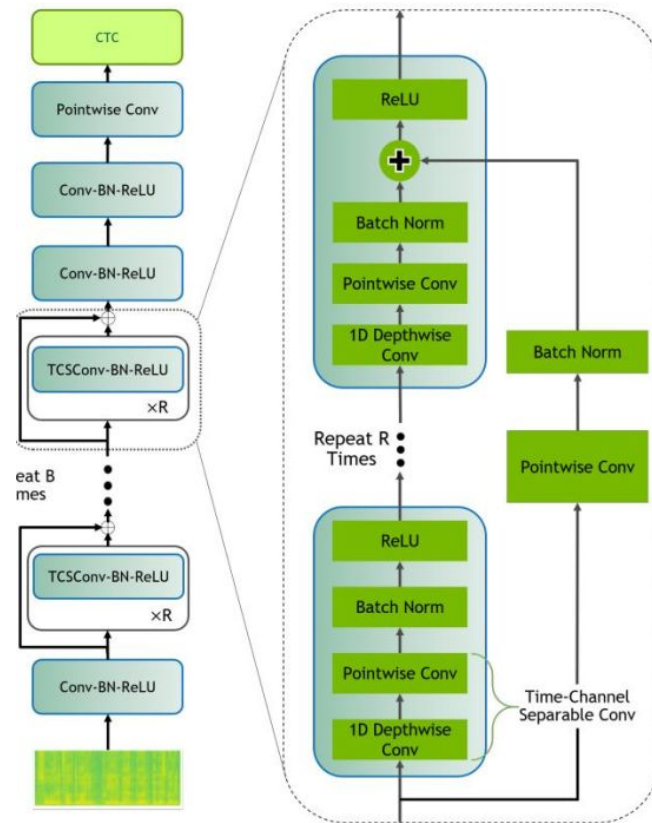
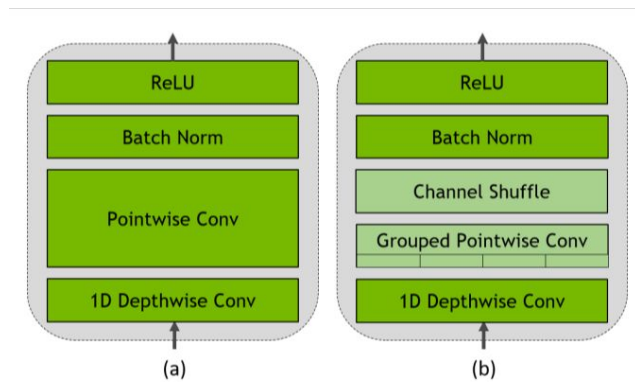


Рис. 1. Архитектура QuartzNet BxR

# Вопросы

Что такое ASR

Чем полезна CTC в архитектуре модели распознавания речи

Зачем нужно все остальное если есть архитектуры основанные на Attention