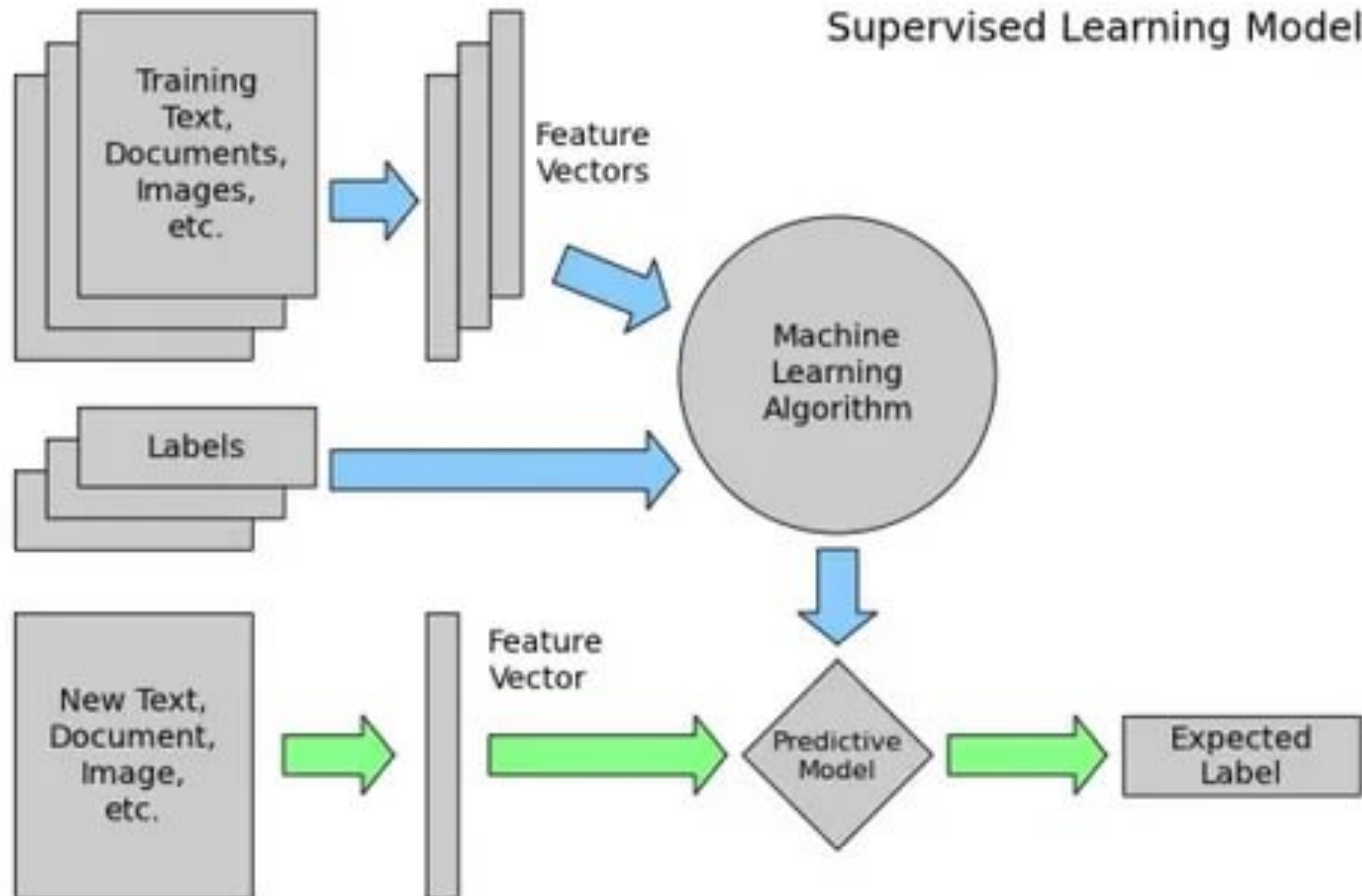


Конволюционные сети для последовательностей

Смирнов Павел БПМИ172

Supervised Learning Model

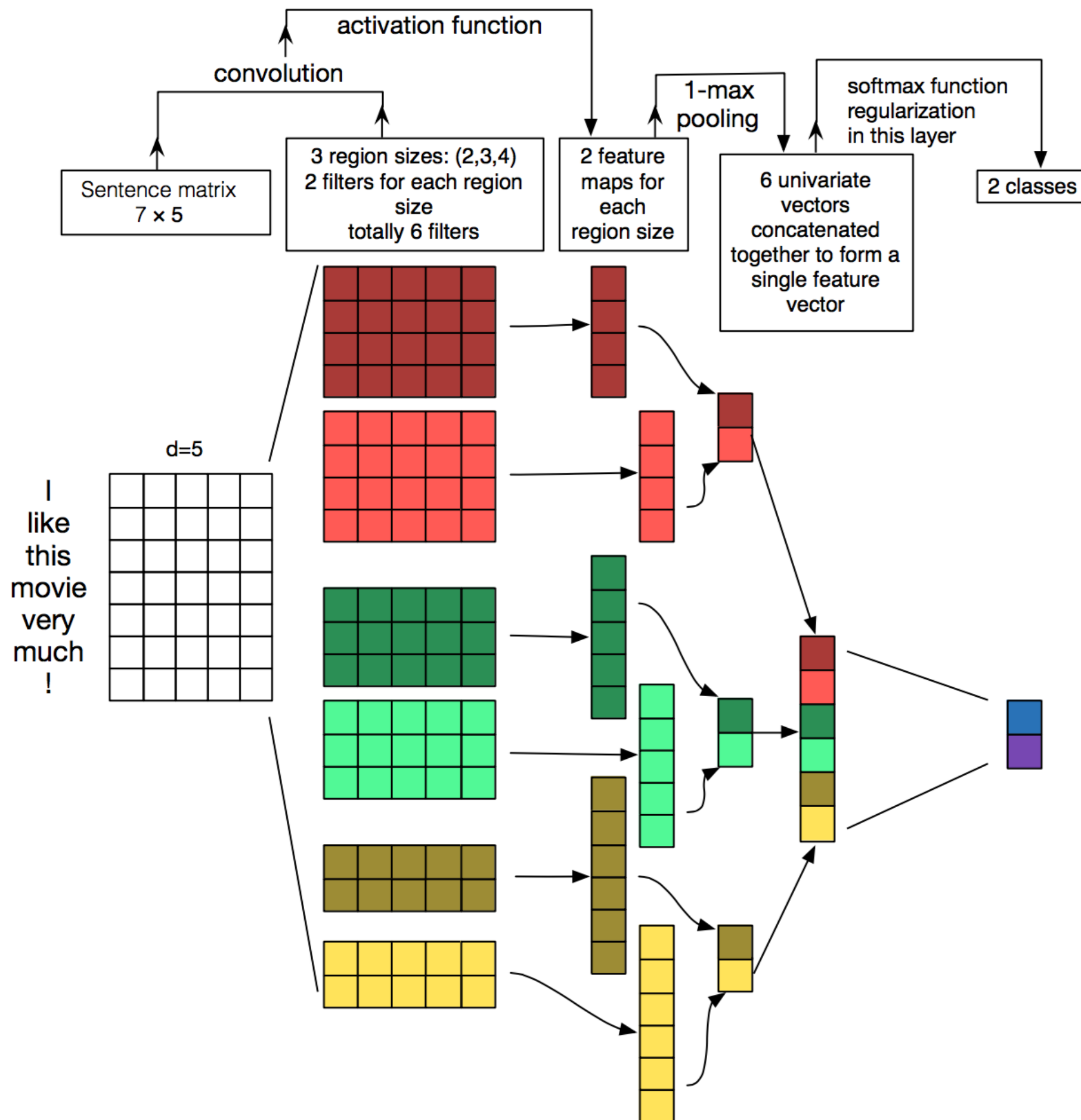


Softmax

- $\sigma(z)_i = \frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}}$ - Функция преобразует вектор z размерности K в вектор σ той же размерности, где каждая координата σ_i полученного вектора представлена вещественным числом в интервале $[0,1]$ и сумма координат равна 1.
- $z = w^T x - \theta$ где x — вектор-столбец признаков объекта размерности $M \times 1$; w^T — транспонированная матрица весовых коэффициентов признаков, имеющая размерность $K \times M$; θ — вектор-столбец с пороговыми значениями размерности $K \times 1$, где K — количество классов объектов, а M — количество признаков объектов.

CNN в текстах

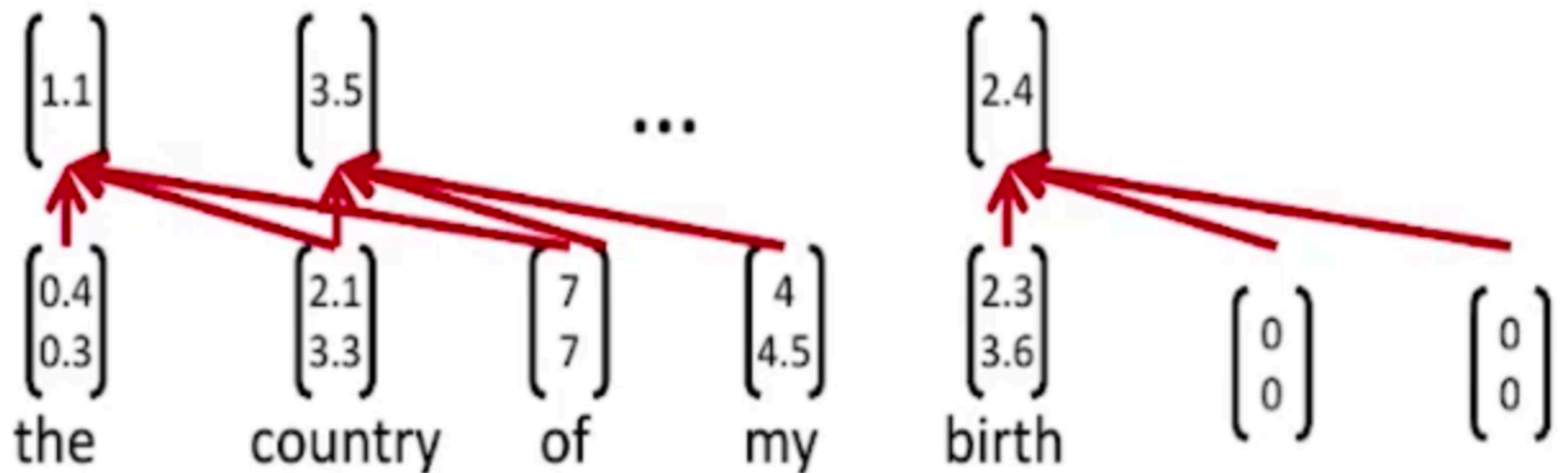
- Съешь ещё этих мягких французских булок:
«Съешь ещё», «ещё этих», ... , «Съешь ещё этих»,
«ещё этих мягких», ..., «Съешь ещё этих мягких», «ещё
этих мягких французских», ... , «Съешь ещё этих
мягких французских булок»



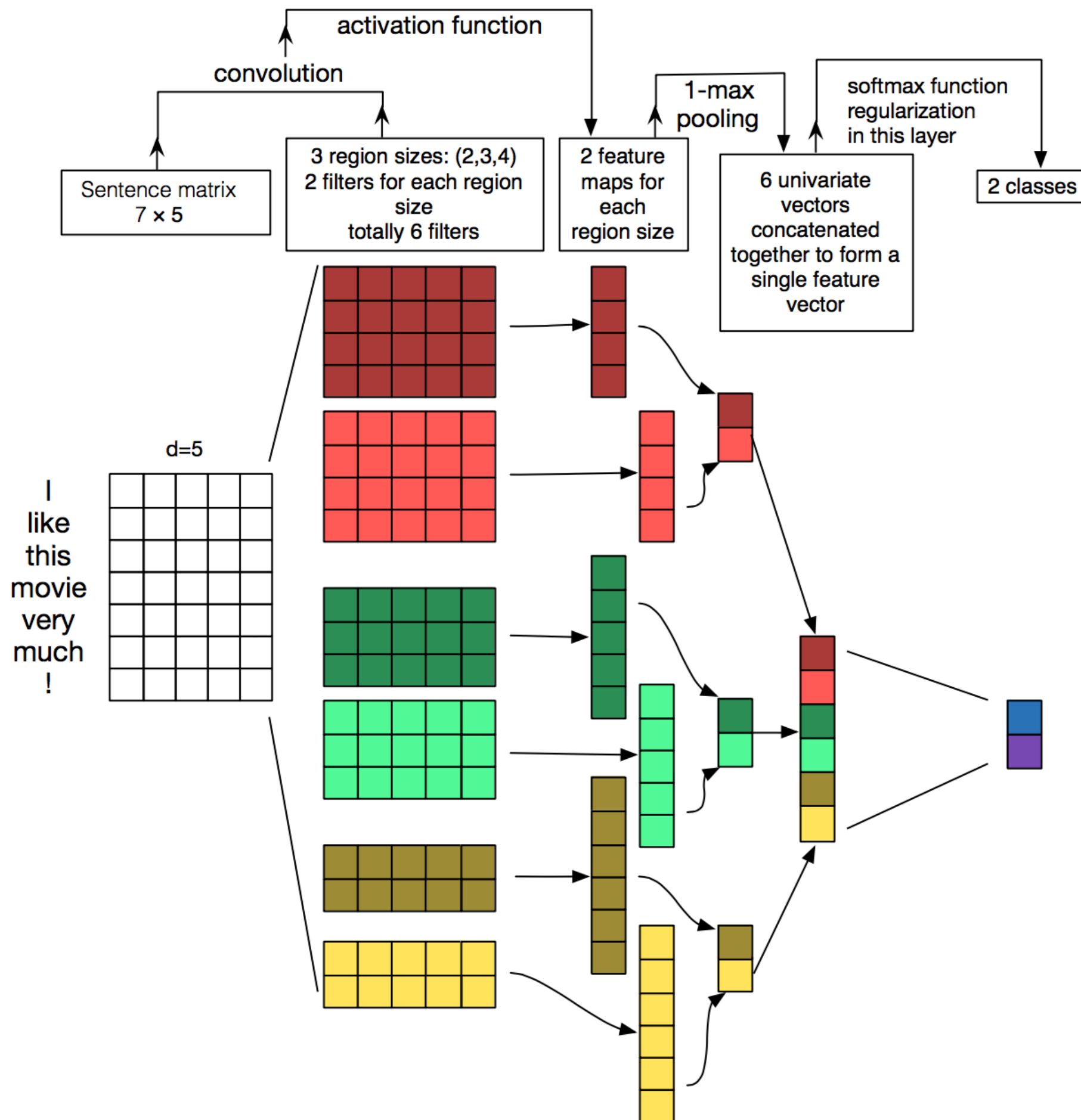
Проблема с фиксированным входом

$$c_i = f(\mathbf{w}^T \mathbf{x}_{i:i+h-1} + b)$$

Result is a feature map: $\mathbf{c} = [c_1, c_2, \dots, c_{n-h+1}] \in \mathbb{R}^{n-h+1}$

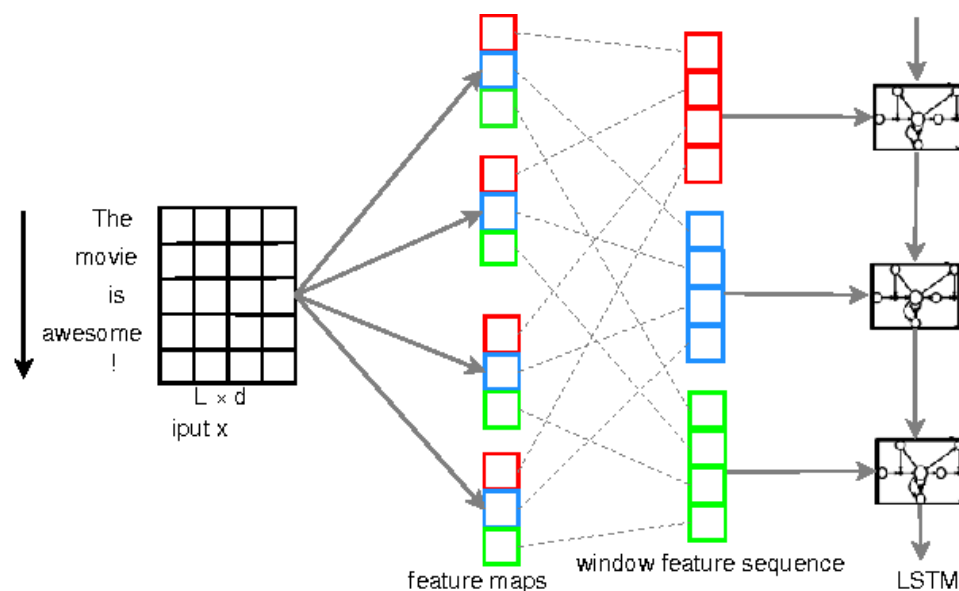


$$\hat{c} = \max_i(c_i)$$

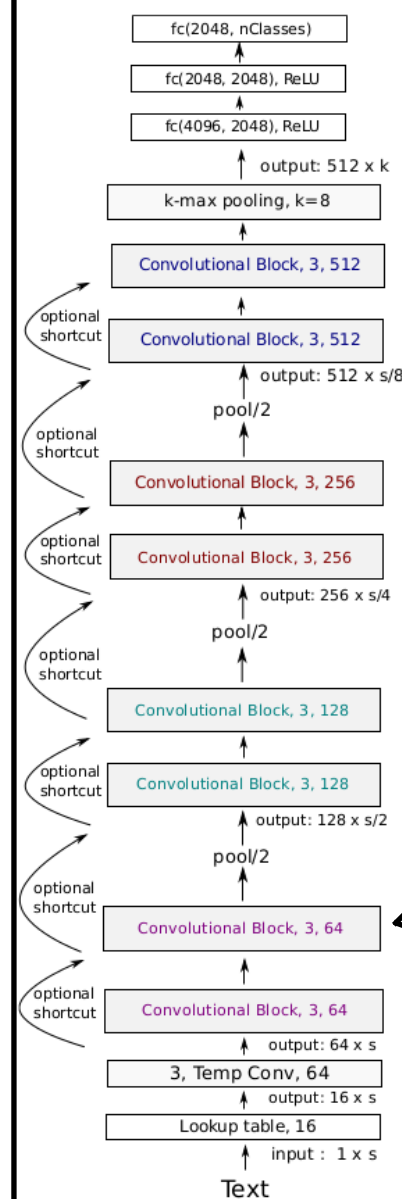


Современные CNN для NLP

- CLSTM



- VDCNN аналог ResNet для текстов (2016) нужны большие input



Каждый convolutional block это:

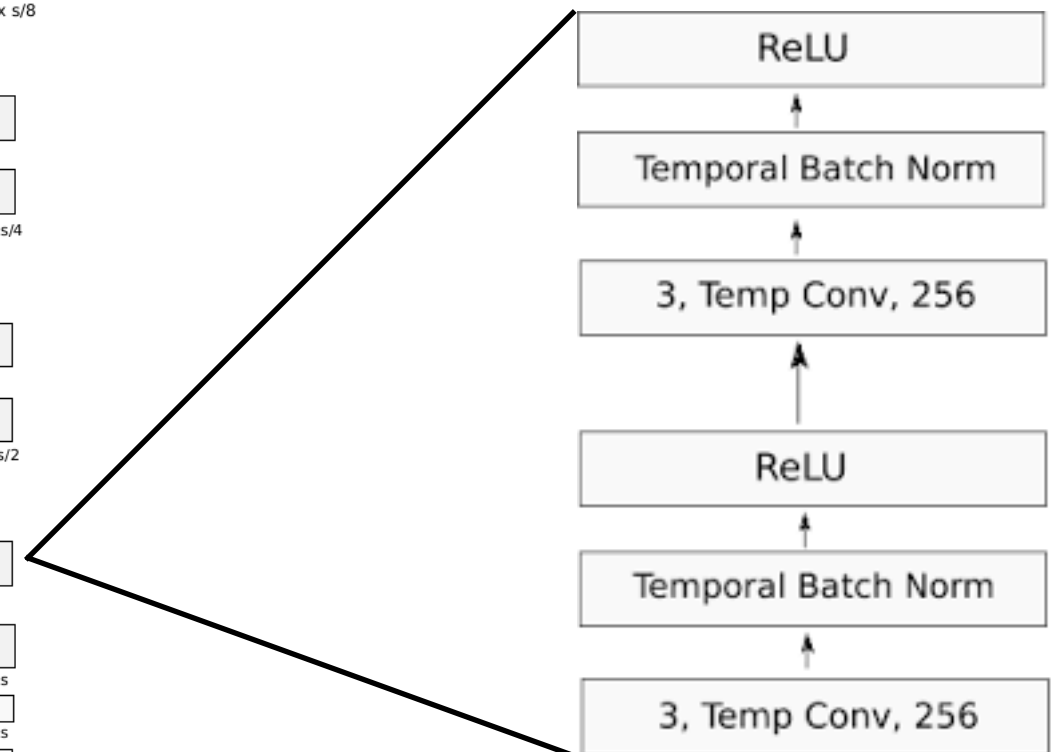
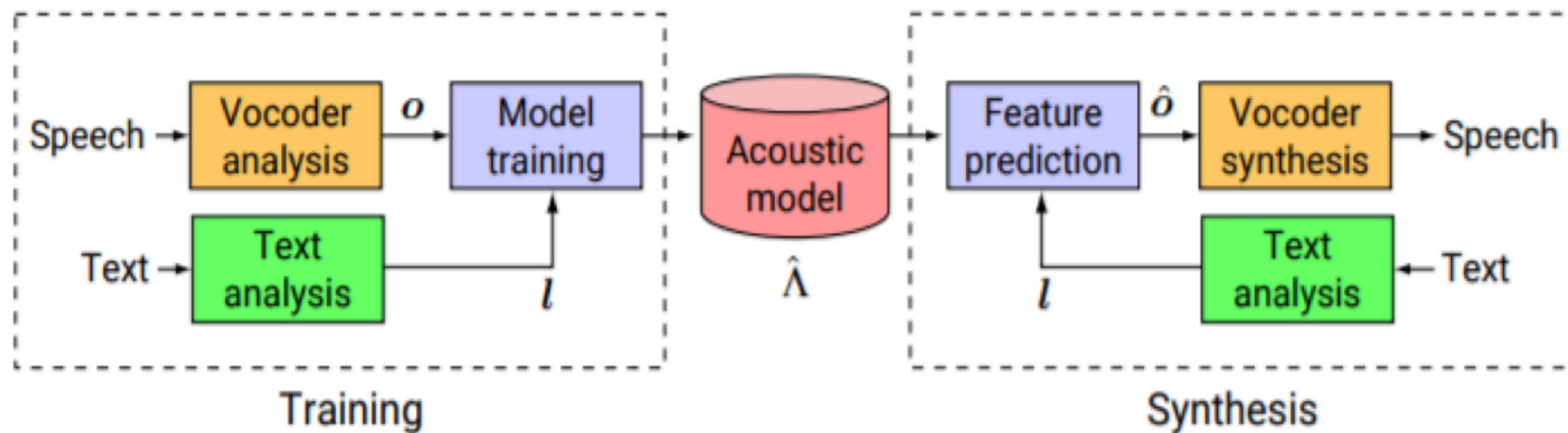


Figure 1: VDCNN architecture.

WaveNet



-

- Совместная вероятность волны $x = (x_1, \dots, x_T)$ описывается как произведение условных вероятностей уравнением:

$$p(x) = \prod_{t=1}^T p(x_t | x_1, \dots, x_{t-1})$$

Свертка

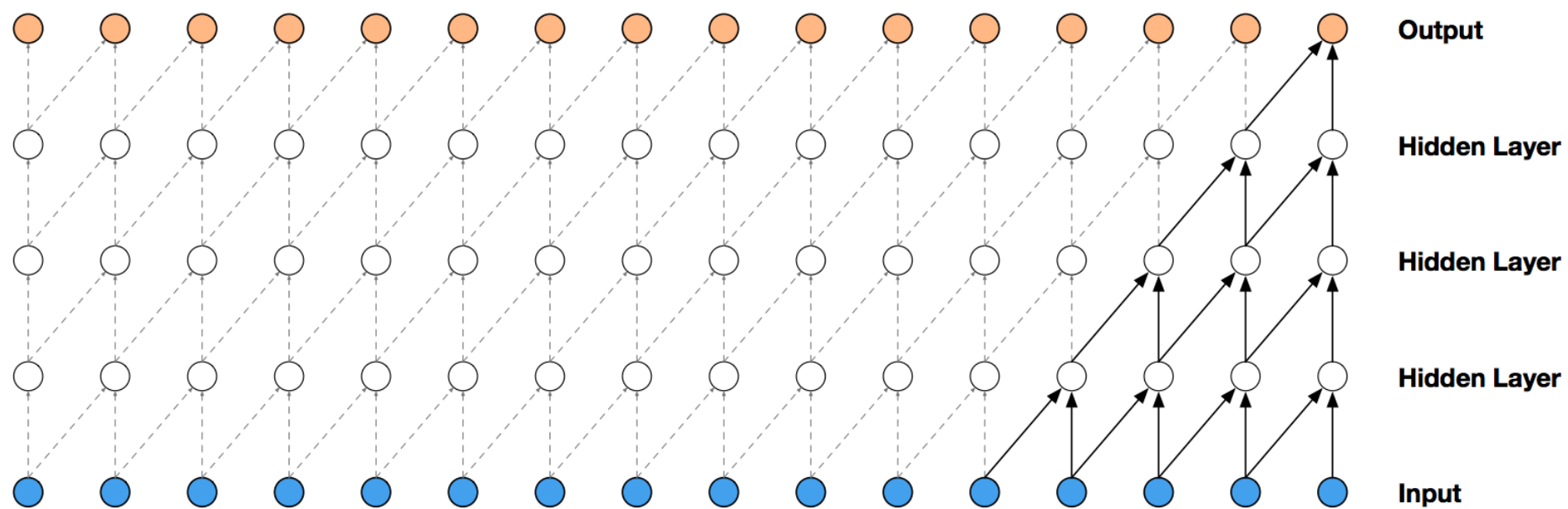
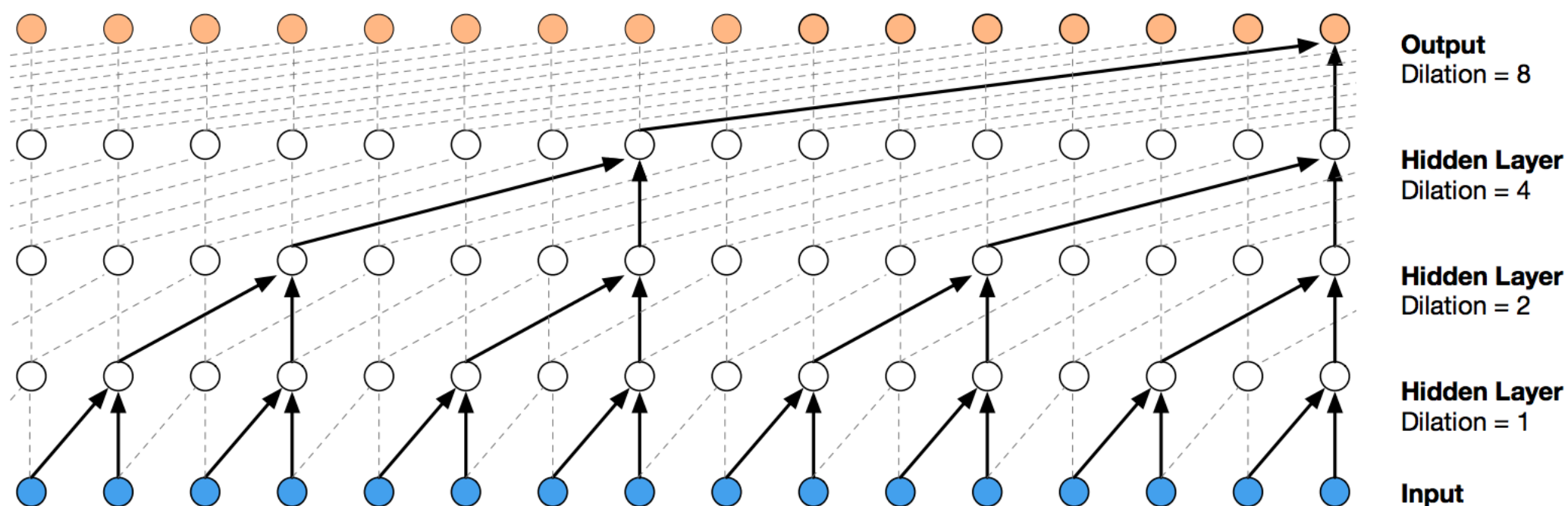


Figure 2: Visualization of a stack of causal convolutional layers.

●

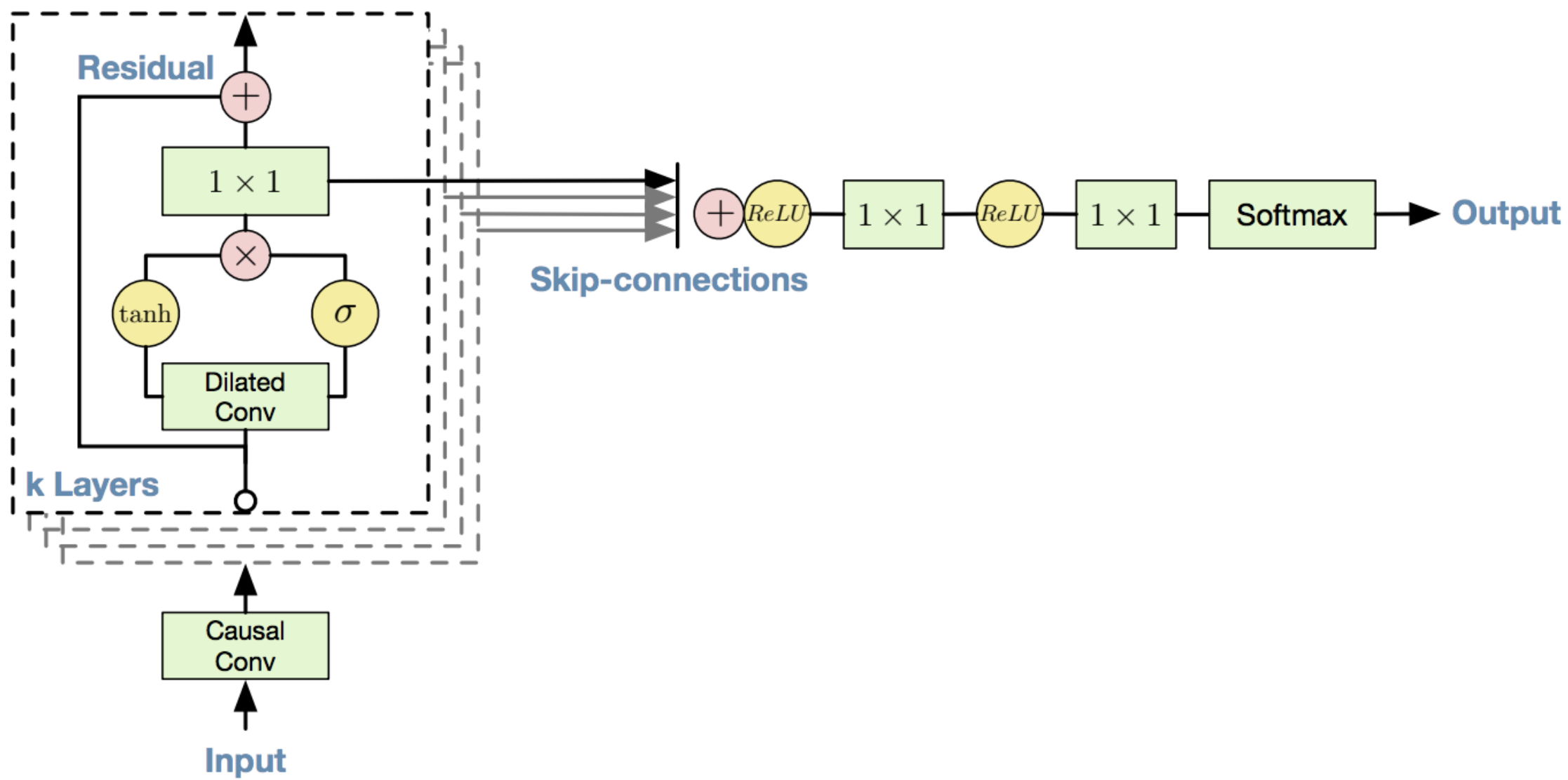
Дырявые свертки



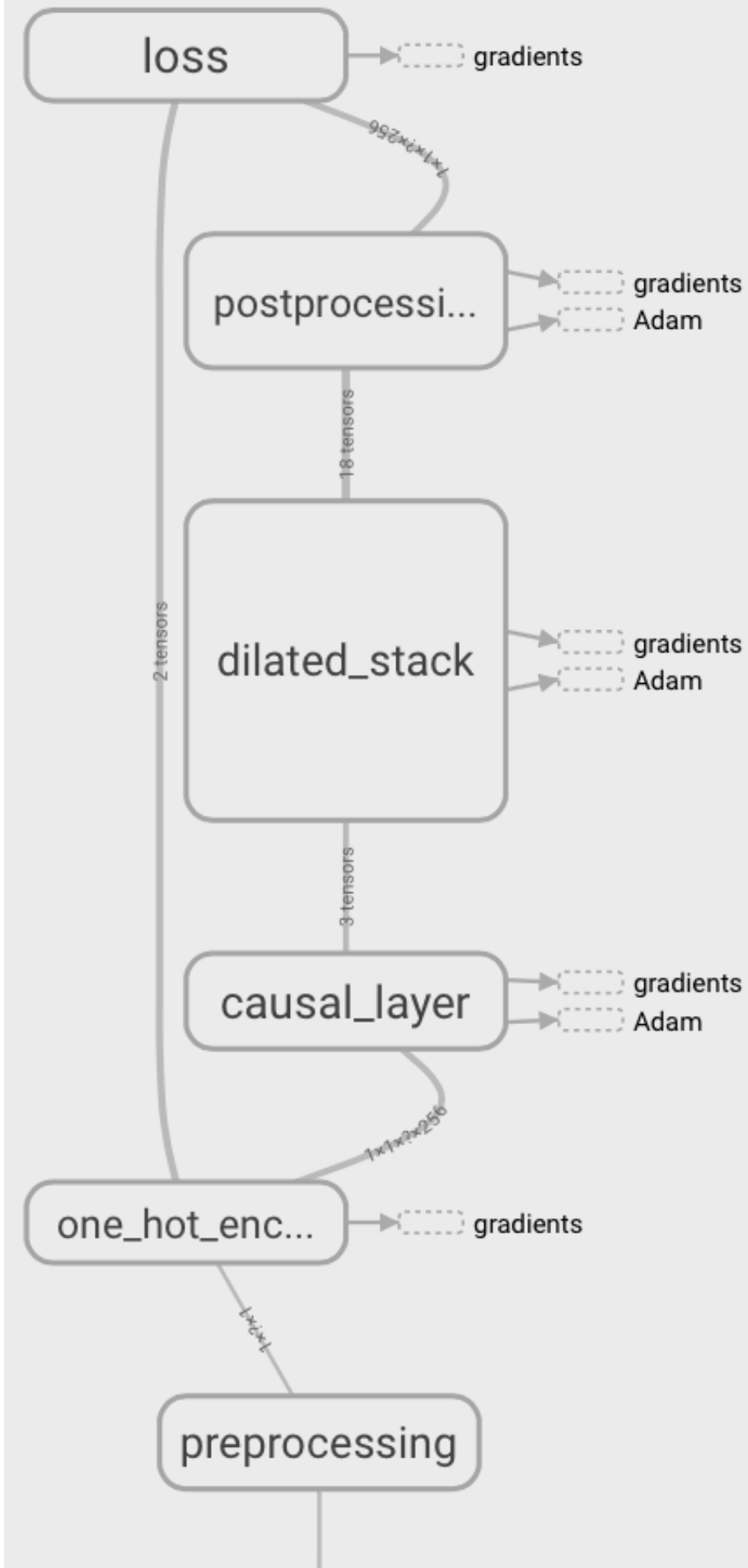
Функция активации

$$z = \tanh(W_{f,k} * x) \odot \sigma(W_{g,k} * x)$$

где $*$ обозначает операцию свёртки, \odot обозначает поэлементное умножение, $\sigma(\cdot)$ сигмоида, k номер слоя, f и g обозначают фильтр и gate соответственно, и W обучаемый свёрточный фильтр.



$$z = \tanh(W_{f,k} * x) \odot \sigma(W_{g,k} * x)$$



WaveNet с условием

- Если дан дополнительный вход h как условие, WaveNet способен моделировать условное распределение $p(x|h) = \prod_{t=1}^T p(x_t|x_1, \dots, x_{t-1}, h)$ аудио по этому входу.
- Функция активации для глобального условия:
$$z = \tanh(W_{f,k} * x + V^T h) \odot \sigma(W_{g,k} * x + V^T h)$$

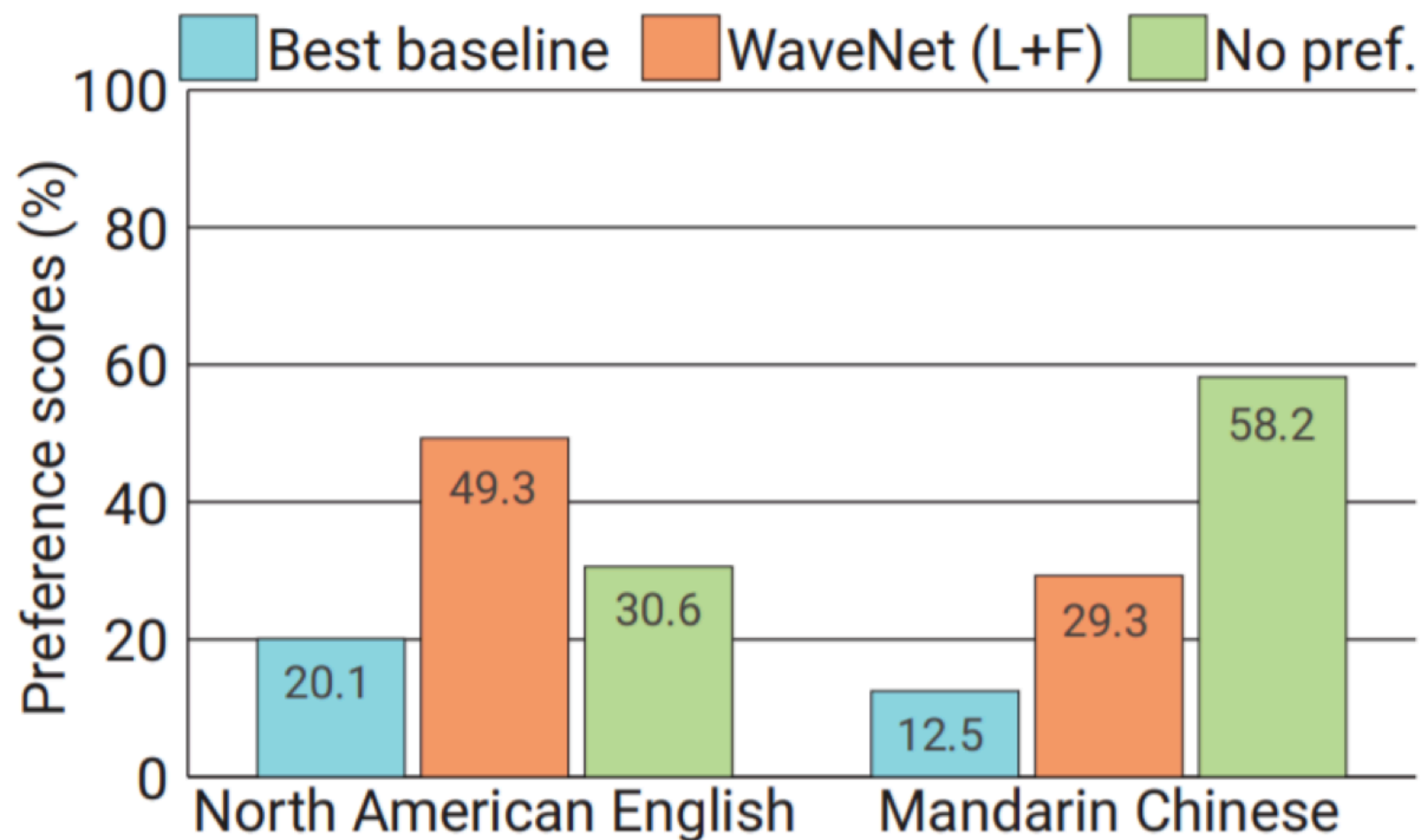
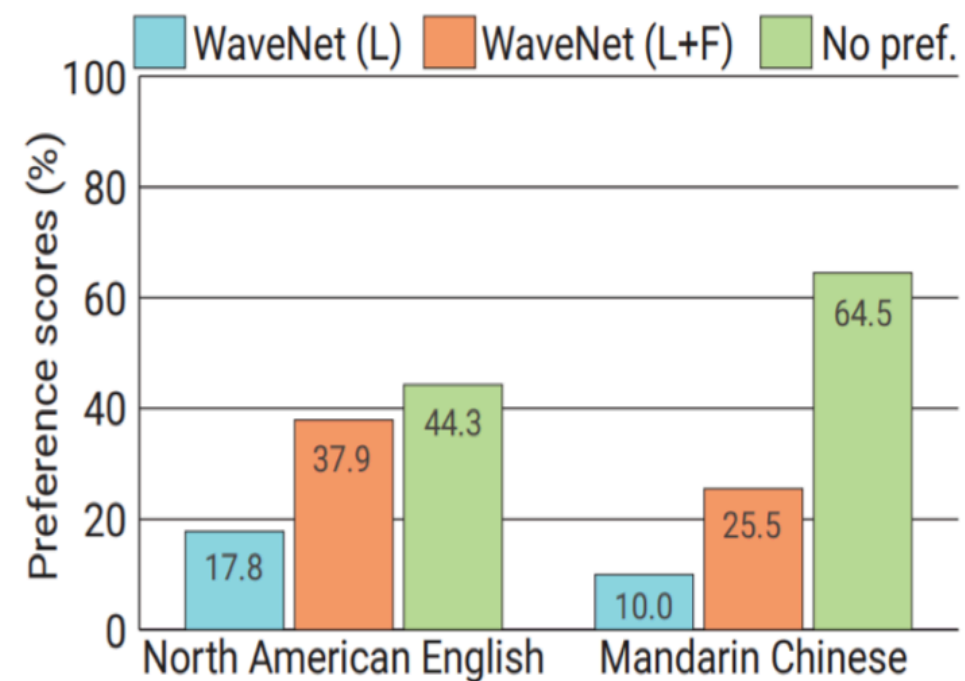
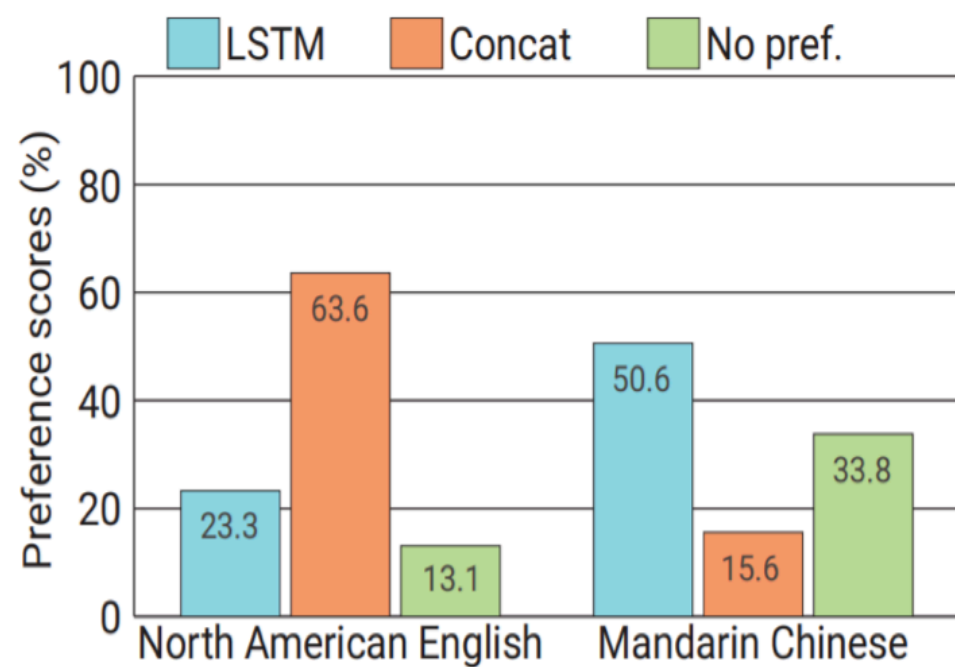
Где V является обучаемой линейной моделью
- Функция активации для локального условия:
$$z = \tanh(W_{f,k} * x + Vh) \odot \sigma(W_{g,k} * x + Vh),$$
 где $V_{*,k}$ это свёртка 1×1

- Локальное условие .
 - Временной ряд, той же длины что и данные. Качество, изменяющееся по времени.
- Глобальное условие.
 - Качество говорящего, не зависящее от времени. Не меняет своего значения в процессе обучения/генерации.

Сравнение качества результатов

Конфигурация WaveNet	Данные	Результат
Без модификаций	весь корпус	речеподобный звук
Без модификаций	одна фраза, много голосов	высокий гул
Без модификаций	один голос, много фраз	шум
Глобальное условие: ID говорящего	весь корпус	речеподобный звук
Глобальное условие: ID говорящего	одна фраза	шум
Глобальное условие: пол говорящего	два голоса, много фраз	шум
Глобальное условие: пол говорящего	два голоса, одна фраза	шум
Локальное условие: текст	весь корпус	шум
Локальное условие: текст	одна фраза	шум
Локальное условие: yandex-speech	весь корпус	речеподобный звук лучшего качества
Локальное условие: yandex-speech	одна фраза	высокий гул
Глобальное условие ID говорящего + локальное условие yandex-speech	весь корпус	шум
Глобальное условие ID говорящего + локальное условие yandex-speech	одна фраза	шум

Speech samples	Subjective 5-scale MOS in naturalness	
	North American English	Mandarin Chinese
LSTM-RNN parametric	3.67 ± 0.098	3.79 ± 0.084
HMM-driven concatenative	3.86 ± 0.137	3.47 ± 0.108
WaveNet (L+F)	4.21 ± 0.081	4.08 ± 0.085
Natural (8-bit μ -law)	4.46 ± 0.067	4.25 ± 0.082
Natural (16-bit linear PCM)	4.55 ± 0.075	4.21 ± 0.071



Производительность

На GPU Tesla K80 с 11 гигабайтами видеопамяти с поправками на виртуализацию.

Чтобы добиться значений функции потерь хотя бы как на рисунке требуется около 4 суток. Генерация пяти секунд аудио занимает около часа.

-

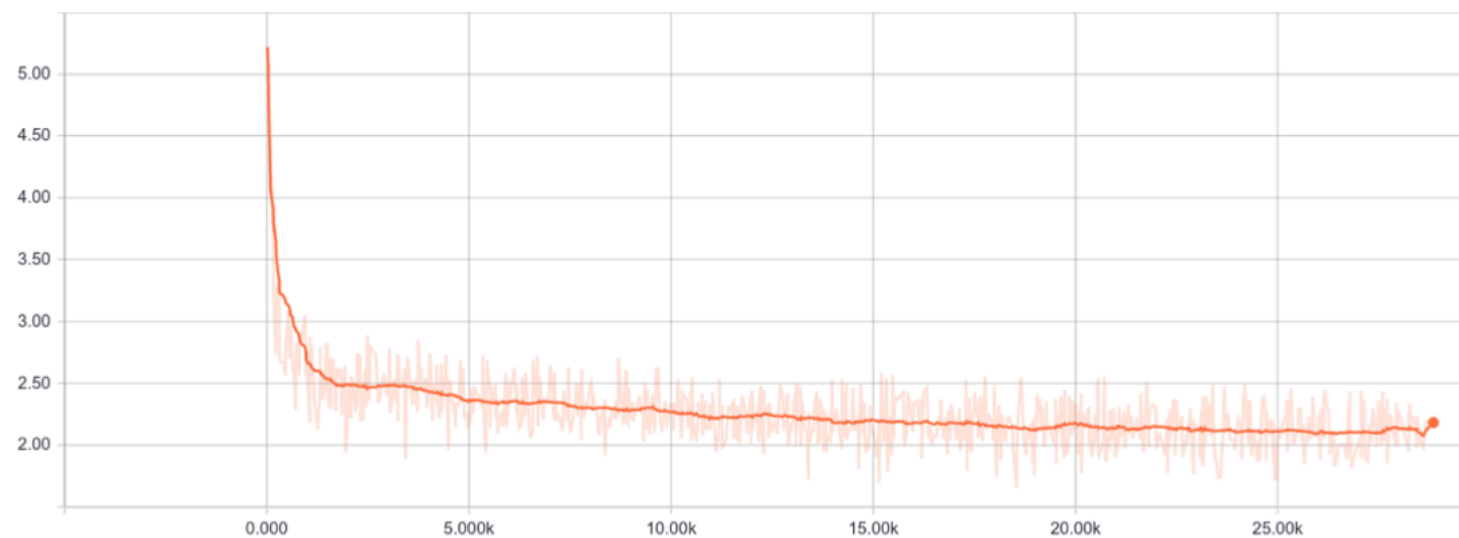


Рис. 11: Изменение функции потерь в процессе обучения

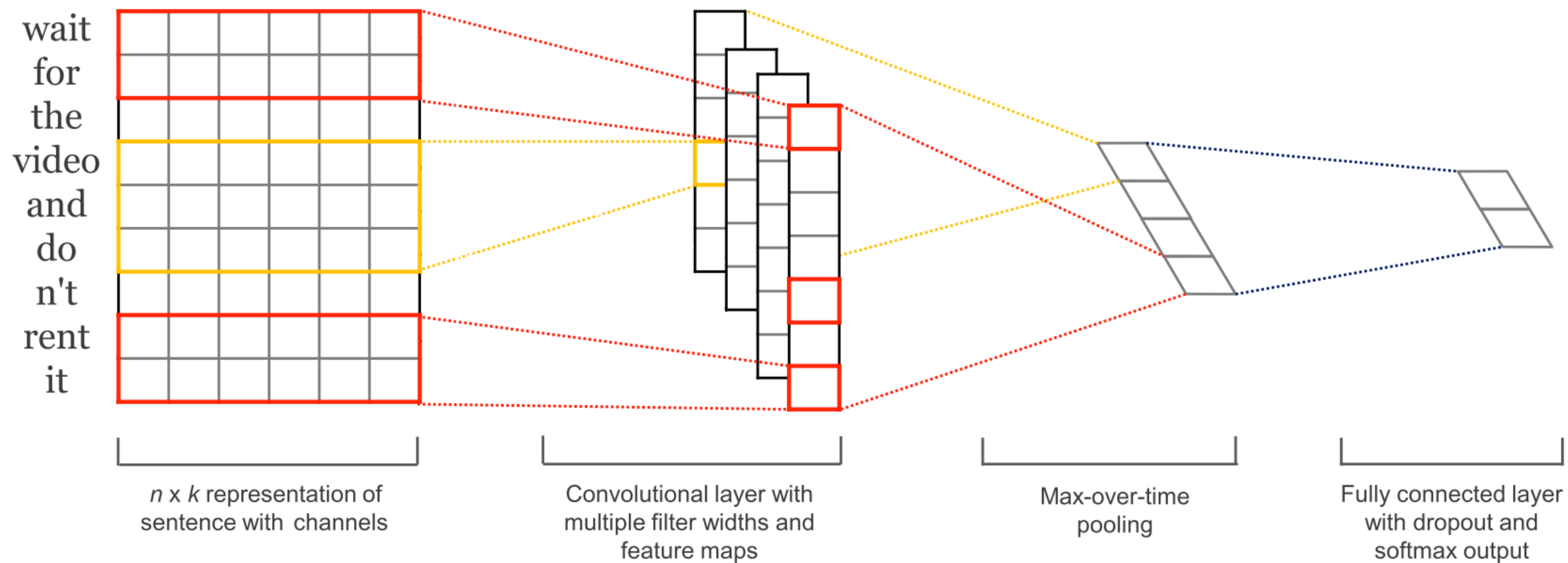
Вопросы

- Напишите функцию активации с условием h ?
- В чем преимущество дырявых сверток (dilated casual convolutions) над casual convolutions?
- Зачем нужны условия в WaveNet?

Ссылки

- http://web.eng.tau.ac.il/deep_learn/wp-content/uploads/2018/01/WaveNet.pdf
- http://mit.spbau.ru/sewiki/images/e/ed/Kurbanov_diploma_master2017.pdf
- <https://habr.com/ru/company/ods/blog/353060/>
- <https://www.youtube.com/watch?v=GYMfcIqMIOU>
- <https://www.youtube.com/watch?v=v-7zVVQqCAs>
- <https://ru.wikipedia.org/wiki/Softmax>

Применение сверток к текстам



- первый шаг - это этап предварительной обработки звука, после квантования входной формы волны до фиксированного целочисленного диапазона. Целочисленные амплитуды затем кодируются для получения тензора формы (num_samples, num_channels).
- Свернутый слой, который получает доступ только к текущим и предыдущим входам, затем уменьшает размер канала.
- Ядро сети построено в виде стека причинно-следственных расширенных слоев, каждый из которых представляет собой расширенную свертку (свертку с дырами), которая обращается только к текущим и прошлым аудиоданным.
- Выходы всех уровней объединяются и расширяются до исходного числа каналов с помощью серии плотных слоев постобработки, за которыми следует функция softmax для преобразования выходов в категориальное распределение.
- Функция потерь представляет собой перекрестную энтропию между выходом для каждого временного интервала и входом на следующем временном интервале.