# Метрики качества

# WER, CER, SER

какого дьявола ты здесь шумишь

какого дявола ты здесь шумишь

$$CWER = \frac{S + D + I}{N}$$

where...
S = number of substitutions
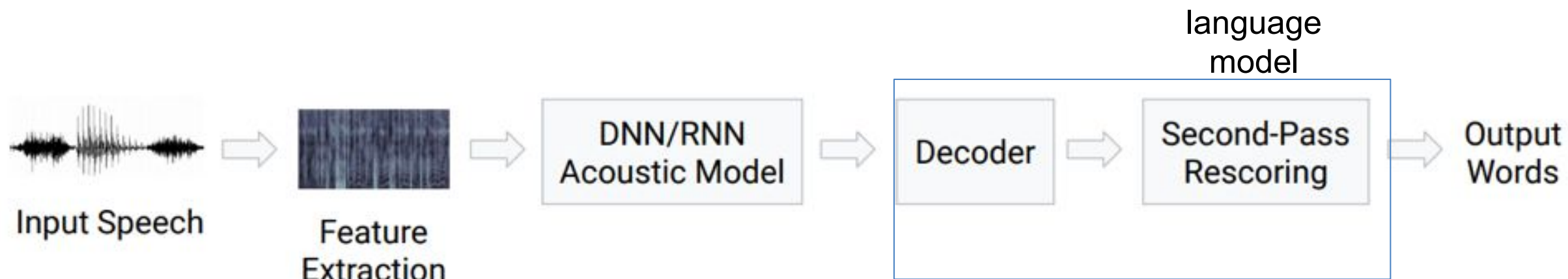D = number of deletions
I = number of insertions
N = number of words in the reference
chars
sentences
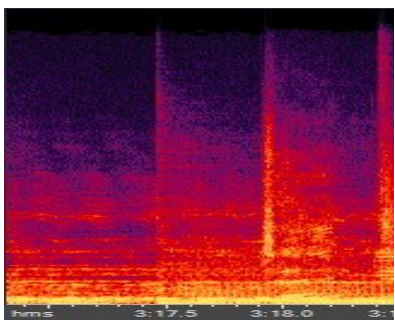
# ASR pipeline

(Hybrid ASR)

# Conventional ASR

Pipeline



Input Speech → Feature Extraction → DNN/RNN Acoustic Model → Decoder → language model / Second-Pass Rescoring → Output Words

## Пайплайн

wav -> melspectrogram (wav2vec) -> Acoustic Model -> Decoding with language model -> (punc) -> words
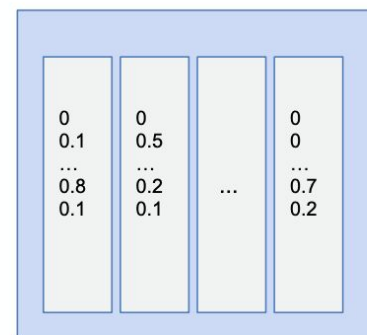

## Метрика

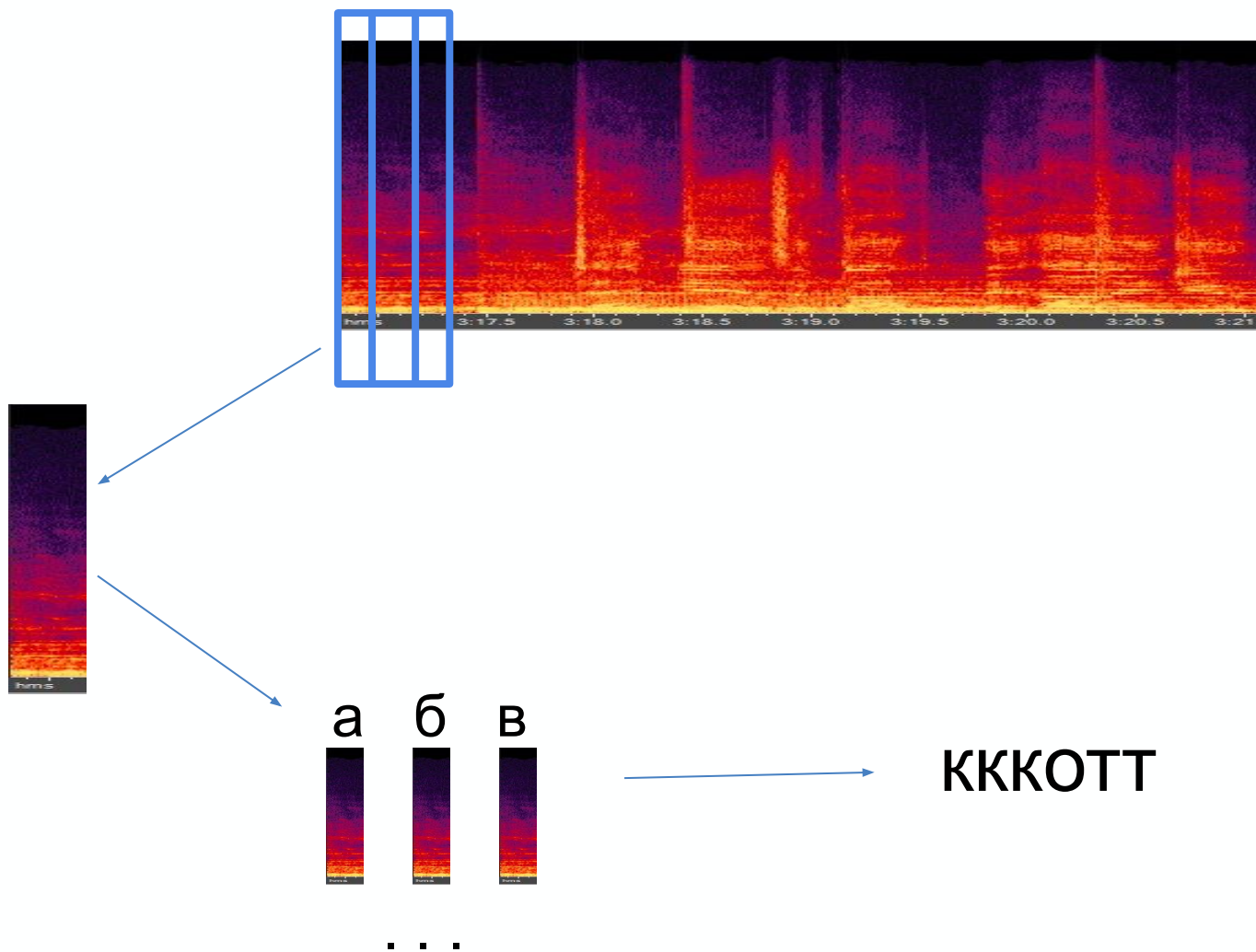wer (word error rate)

# Акустическая модель

# frame-level prediction?
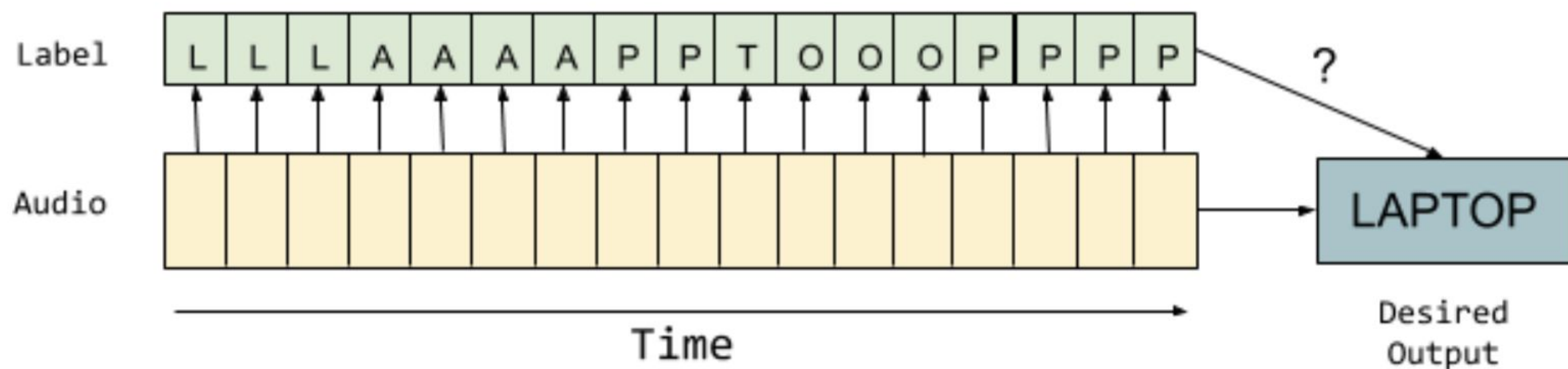


а б в

ККОТТ

# CTC-loss

Figure 1. Here, we have aligned audio data, where the audio is chopped up into time slices and each is labeled with a letter. But it's very difficult to go from those labels to the correct transcript, especially considering words with repeated letters (such as "book").

log (Pr (output: "BOOK" | audio)) = log (Pr ( BOO-OOO - КК | audio)) +  log (Pr (ВBO - OO-ККК | audio)) + ...).
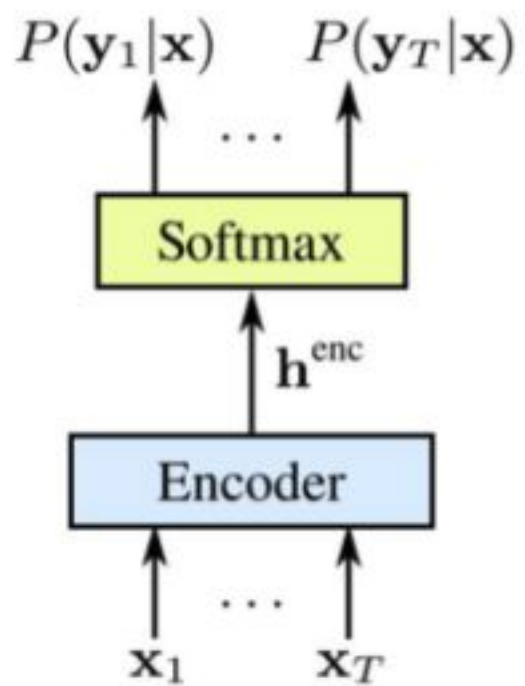
На практике мы можем использовать подход динамического программирования, чтобы рассчитать это, накапливая наши логарифмические вероятности по разным «путям» через выходы softmax на каждом шаге.

коллапсирование

ккклллллаасс_с-с

класс

$$P(\mathbf{y}_1|\mathbf{x}) \quad P(\mathbf{y}_T|\mathbf{x})$$

Softmax

$$\mathbf{h}^{enc}$$

Encoder

$$\mathbf{x}_1 \quad \mathbf{x}_T$$
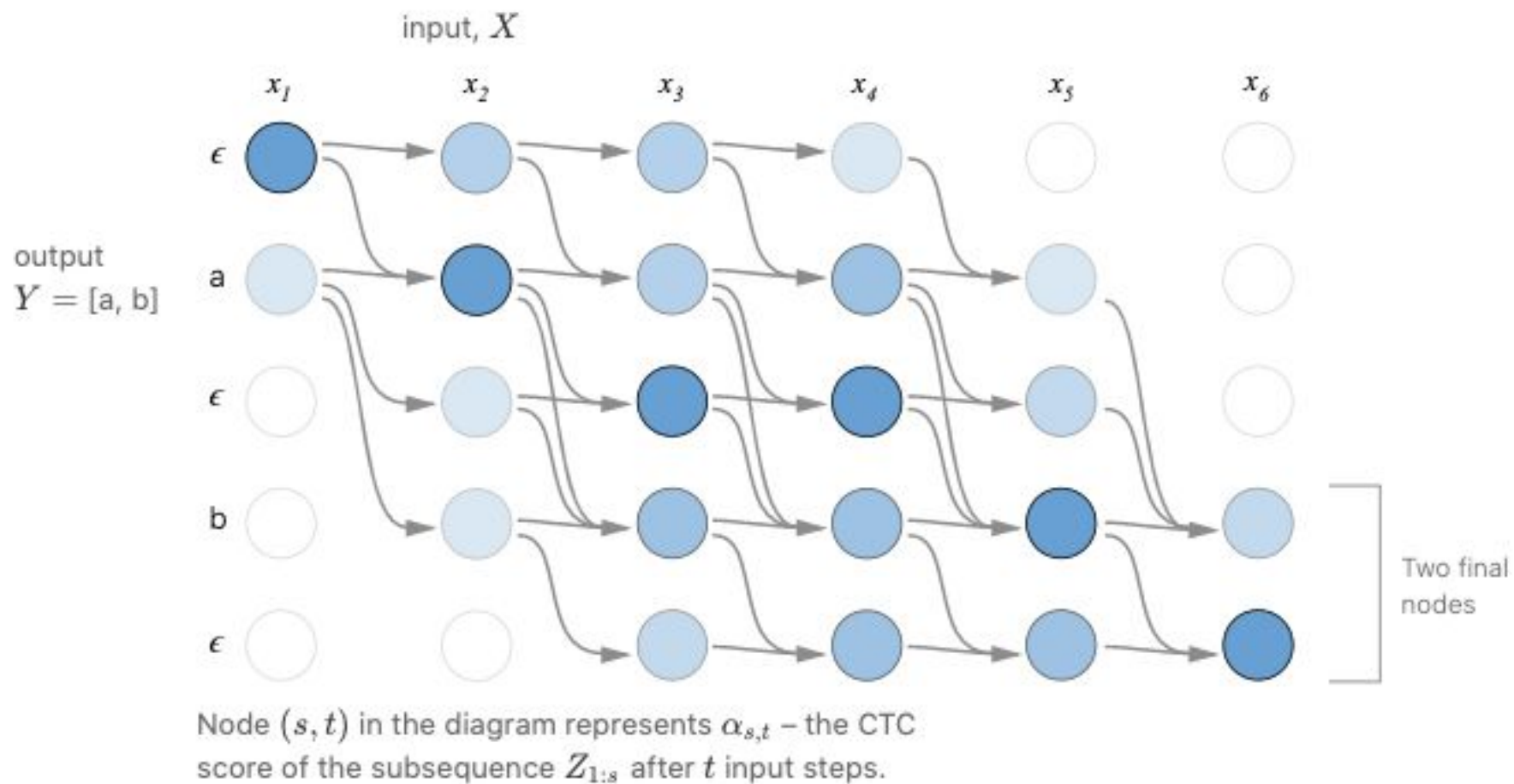
```
B  B  c  B  B  a  a  B  B  t
B  c  c  B  a  B  B  B  B  t
              ...
B  c  B  B  a  B  B  t  t  B
```

$$P(\mathbf{y}|\mathbf{x}) = \sum_{\hat{\mathbf{y}} \in \mathcal{B}(\mathbf{y},\mathbf{x})} \prod_{t=1}^{T} P(\hat{y}_t|\mathbf{x})$$

# CTC loss



input, $X$

$x_1$    $x_2$    $x_3$    $x_4$    $x_5$    $x_6$

output $Y = [a, b]$

Two final nodes

Node $(s, t)$ in the diagram represents $\alpha_{s,t}$ – the CTC score of the subsequence $Z_{1:s}$ after $t$ input steps.

**Handwriting recognition:** The input can be $(x, y)$ coordinates of a pen stroke or pixels in an image.

**Speech recognition:** The input can be a spectrogram or some other frequency based feature extractor.

# Listen-Attent-Spell (2015)

- RNN
- Autoregressive
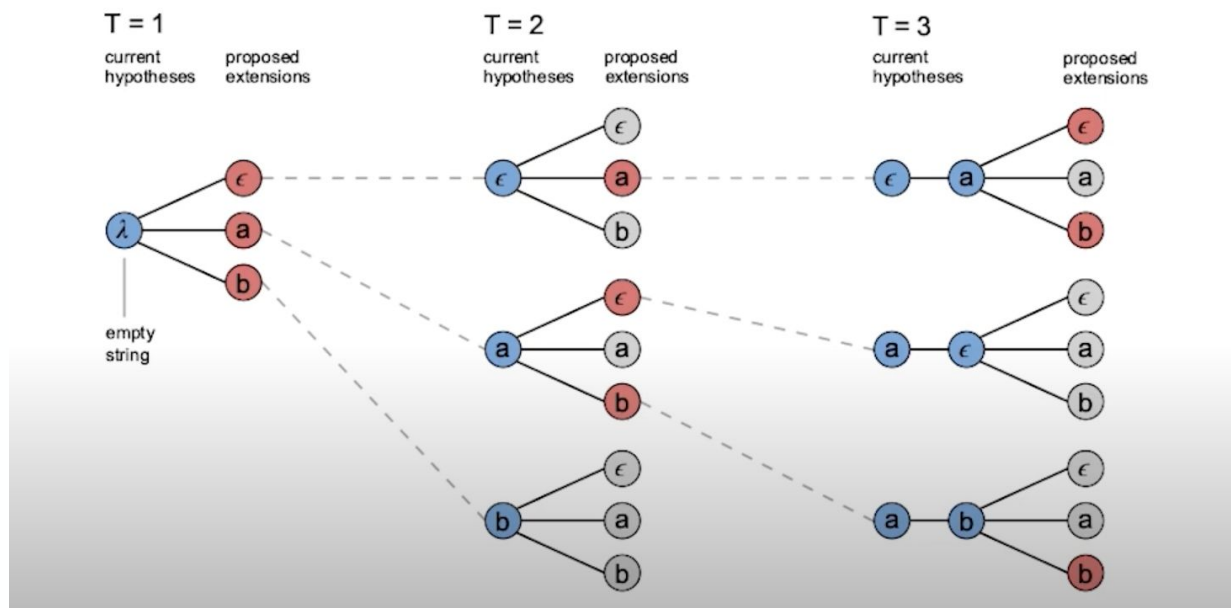- No need beam search & LM
- Cross-Entropy

characters: $\{a, b, c, \cdots, z, 0, \cdots, 9, \langle space \rangle, \langle comma \rangle, \langle period \rangle, \langle apostrophe \rangle, \langle unk \rangle\}$

Speller

$P(y_i | \mathbf{x}, y_{<i}) = \text{CharacterDistribution}(s_i, c_i)$

$s_i = \text{RNN}(s_{i-1}, y_{i-1}, c_{i-1})$

historical information

$c_i = \text{AttentionContext}(s_i, \mathbf{h})$

current context

acoustic model encoder

$h = (h_1, \ldots, h_U)$ — dense representation of x

Listener

input sequence of filter bank spectra features

# DeepSpeech 2 (2015)

- RNN & Conv
- Non-Autoregressive
- Need LM beam search & LM
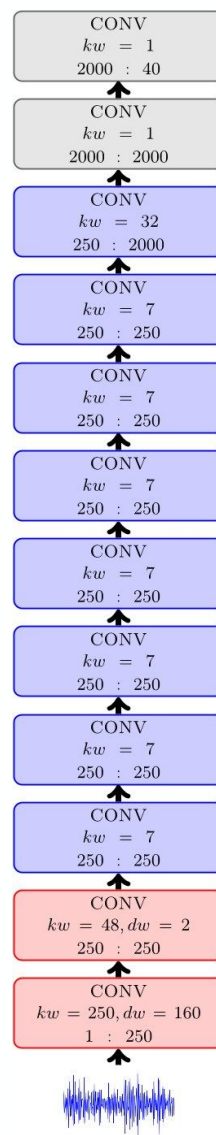- CTC

# Beam Search & LM



BEAM SEARCH

Only beam search

$$y^* = \arg\max_{y} \; \log p(y|x)$$

Beam search & LM (shallow fusion)

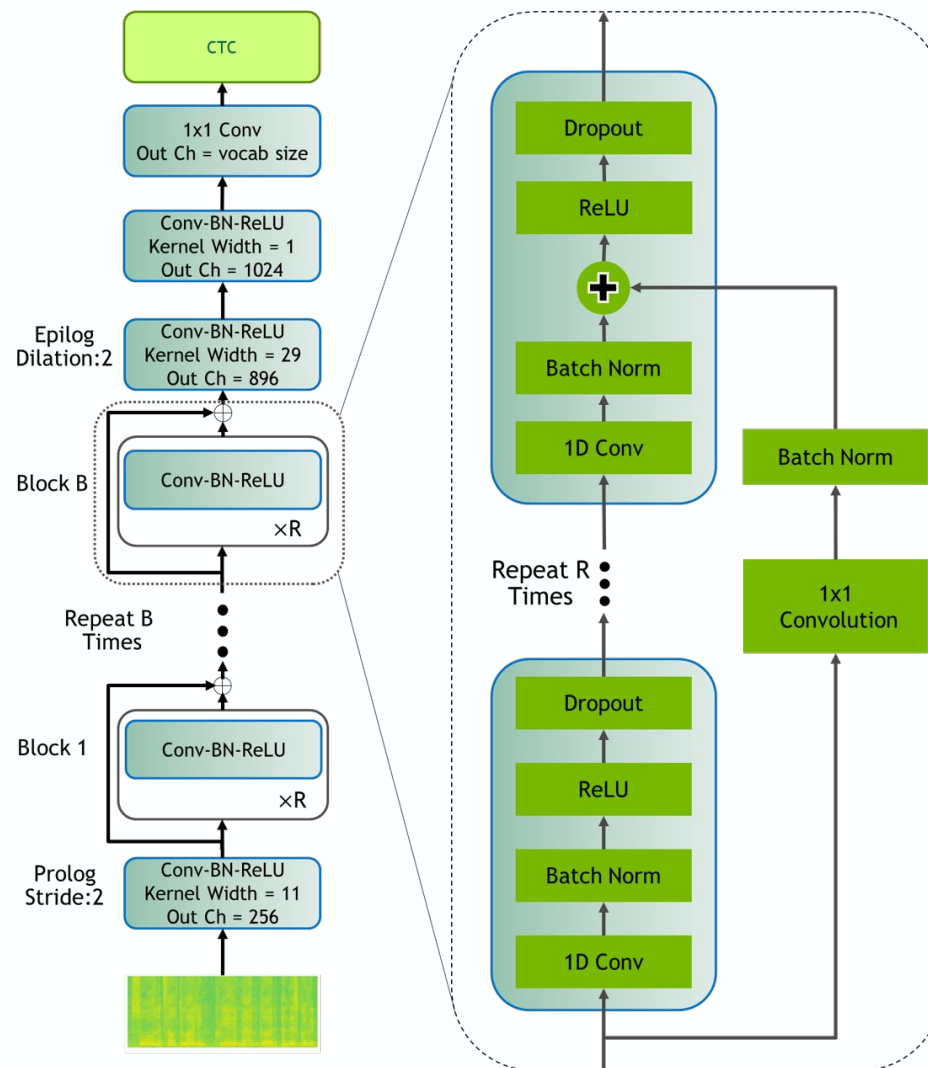$$y^* = \arg\max_{y} \; \log p(y|x) + \lambda \log p_{LM}(y)$$

# Wav2Letter (2016)

- Conv
- Non-Autoregressive
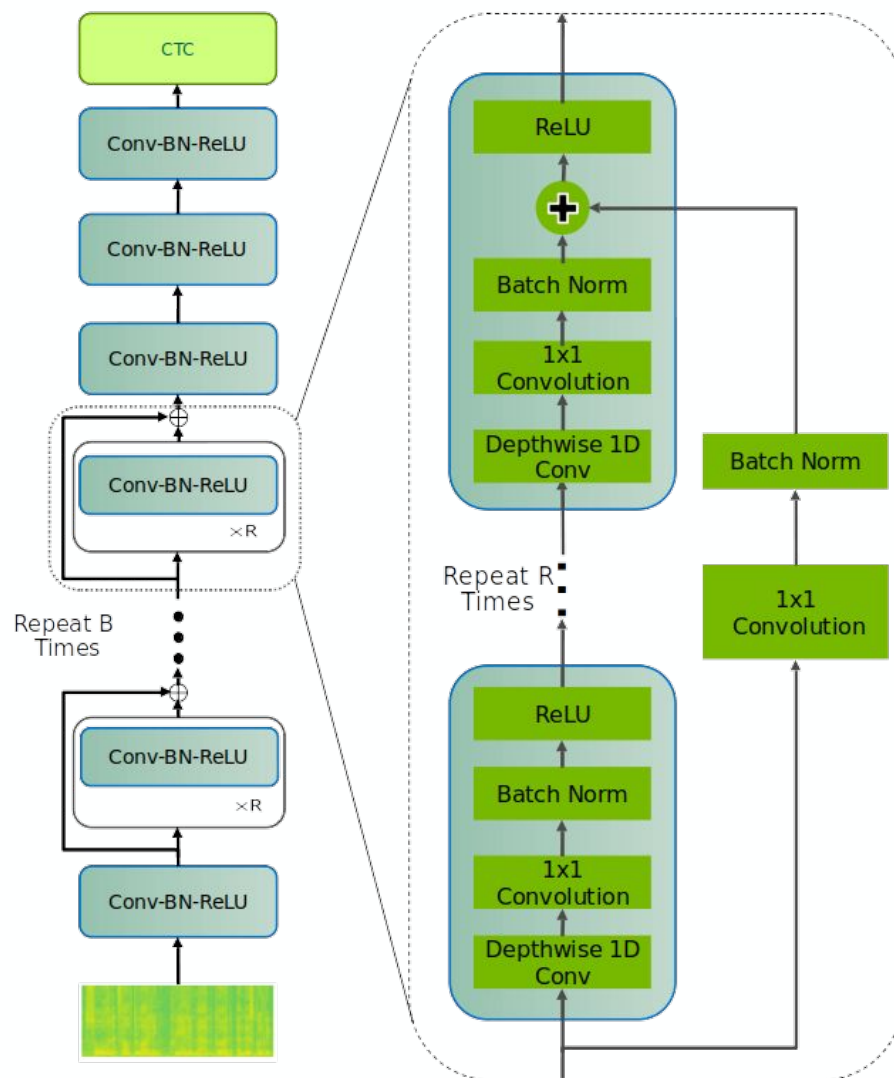- Need beam-search & LM
- CTC

# Jasper (2019)

- Conv
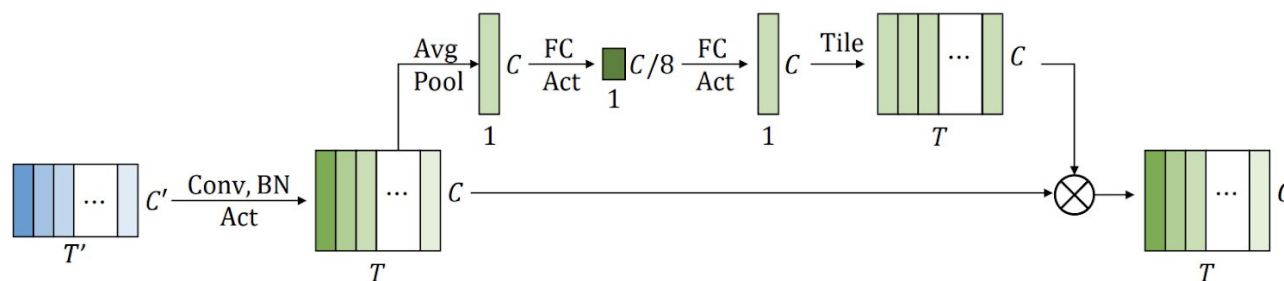- Non-Autoregressive
- Need beam-search & LM
- CTC

# QuartzNet (2019)

- Conv
- Non-Autoregressive
- Need beam-search & LM
- CTC

# ContextNet (2020)

- RNN & Conv
- Autoregressive
- Better with beam-search & LM, but can work without it
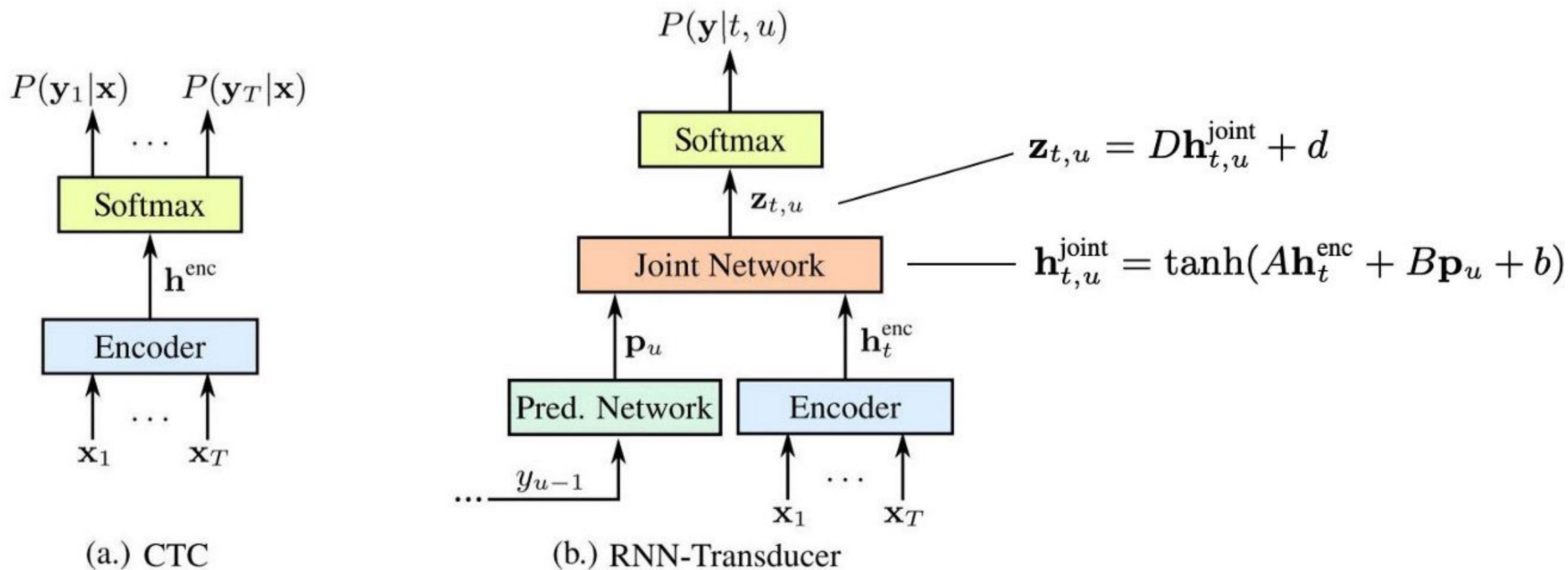- RNN-T loss

# RNN-Transducer



$$\mathbf{z}_{t,u} = D\mathbf{h}^{\text{joint}}_{t,u} + d$$

$$\mathbf{h}^{\text{joint}}_{t,u} = \tanh(A\mathbf{h}^{\text{enc}}_t + B\mathbf{p}_u + b)$$

(a.) CTC

(b.) RNN-Transducer

# Conformer (2020)

- RNN, Conv & Transformer
- Autoregressive
- Better with beam-search & LM, but can work without it
- RNN-T loss

| Method | #Params (M) | WER Without LM | | WER With LM | |
|---|---|---|---|---|---|
| | | testclean | testother | testclean | testother |
| **Hybrid** | | | | | |
| Transformer [33] | - | - | - | 2.26 | 4.85 |
| **CTC** | | | | | |
| QuartzNet [9] | 19 | 3.90 | 11.28 | 2.69 | 7.25 |
| **LAS** | | | | | |
| Transformer [34] | 270 | 2.89 | 6.98 | 2.33 | 5.17 |
| Transformer [19] | - | 2.2 | 5.6 | 2.6 | 5.7 |
| LSTM | 360 | 2.6 | 6.0 | 2.2 | 5.2 |
| **Transducer** | | | | | |
| Transformer [7] | 139 | 2.4 | 5.6 | 2.0 | 4.6 |
| ContextNet(S) [10] | 10.8 | 2.9 | 7.0 | 2.3 | 5.5 |
| ContextNet(M) [10] | 31.4 | 2.4 | 5.4 | **2.0** | 4.5 |
| ContextNet(L) [10] | 112.7 | **2.1** | 4.6 | **1.9** | 4.1 |
| **Conformer (Ours)** | | | | | |
| Conformer(S) | 10.3 | **2.7** | **6.3** | **2.1** | **5.0** |
| Conformer(M) | 30.7 | **2.3** | **5.0** | **2.0** | **4.3** |
| Conformer(L) | 118.8 | **2.1** | **4.3** | **1.9** | **3.9** |