

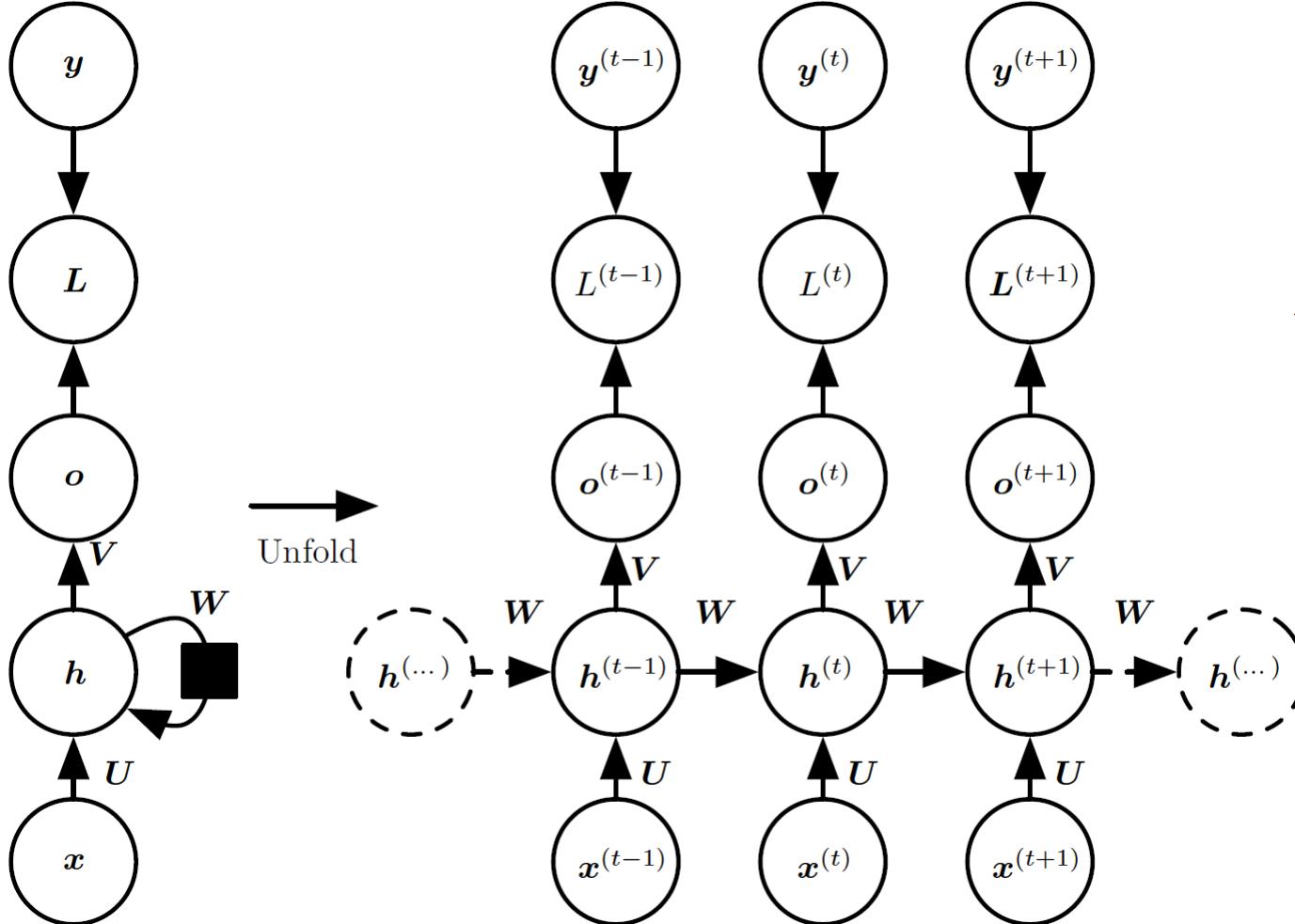


RNN architectures & applications

stats403_deep_learning
spring_2025
lecture_6

6.1 RNN architectures

backpropagation through time (BPTT)



$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)},$$

$$h^{(t)} = \tanh(a^{(t)}),$$

$$o^{(t)} = c + Vh^{(t)},$$

$$\hat{y}^{(t)} = \text{softmax}(o^{(t)}),$$

$$L(\{x^{(1)}, \dots, x^{(\tau)}\}, \{y^{(1)}, \dots, y^{(\tau)}\})$$

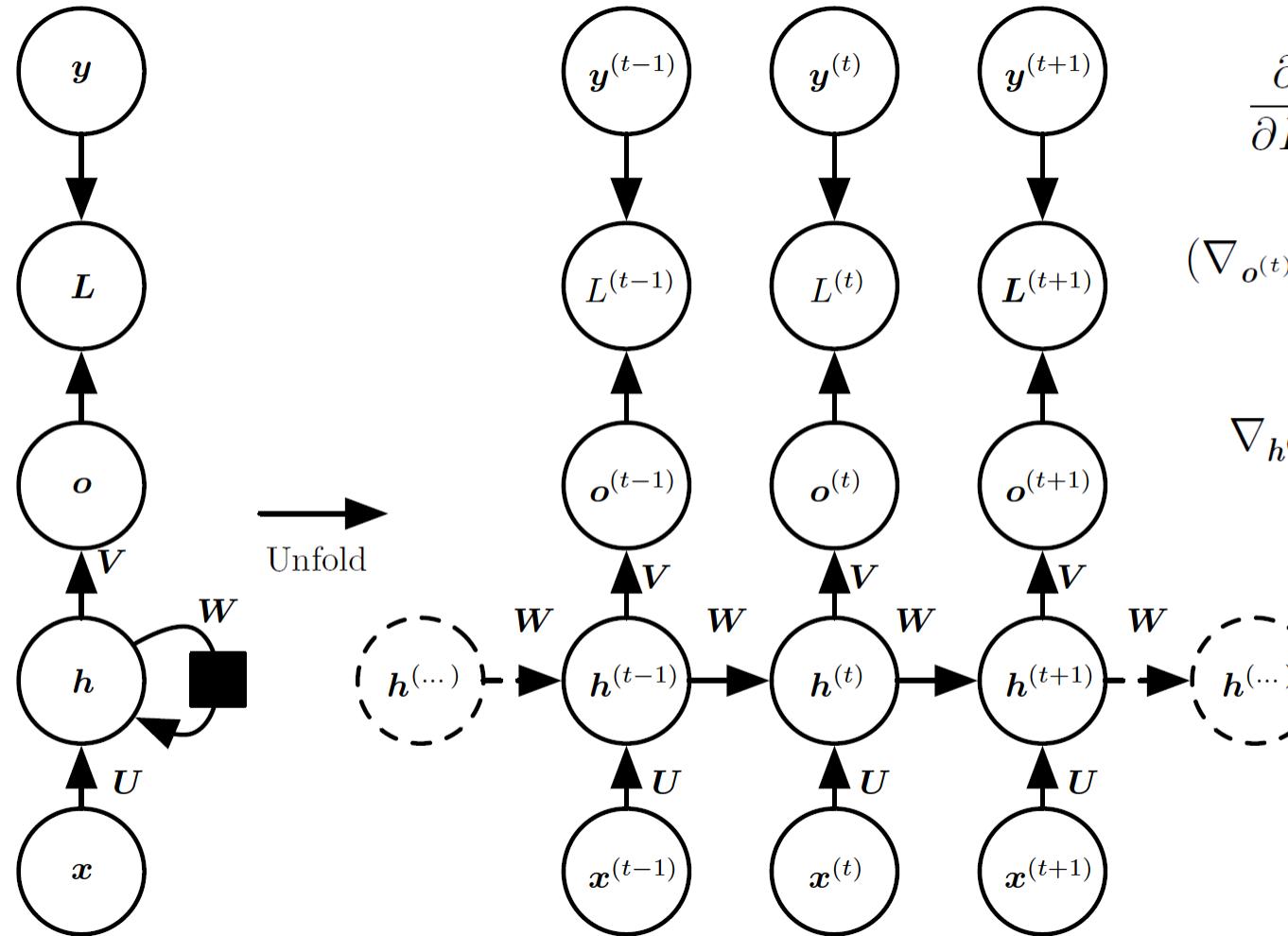
$$= \sum_t L^{(t)}$$

$$= - \sum_t \log p_{\text{model}}(y^{(t)} \mid \{x^{(1)}, \dots, x^{(t)}\})$$

BPTT

$$\begin{aligned} \mathbf{a}^{(t)} &= \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)}, \\ \mathbf{h}^{(t)} &= \tanh(\mathbf{a}^{(t)}), \\ \mathbf{o}^{(t)} &= \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)}, \\ \hat{\mathbf{y}}^{(t)} &= \text{softmax}(\mathbf{o}^{(t)}), \end{aligned}$$

$$\begin{aligned} L(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}\}, \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\}) \\ = \sum_t L^{(t)} \\ = - \sum_t \log p_{\text{model}}(y^{(t)} \mid \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\}) \end{aligned}$$

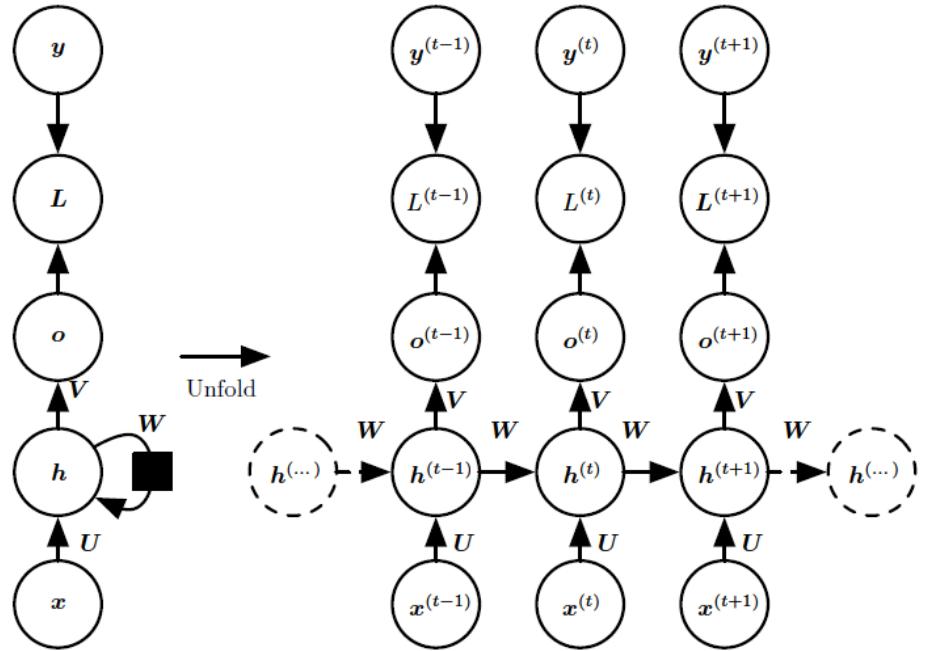


$$\frac{\partial L}{\partial L^{(t)}} = 1$$

$$(\nabla_{\mathbf{o}^{(t)}} L)_i = \frac{\partial L}{\partial o_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial o_i^{(t)}} = \hat{y}_i^{(t)} - \mathbf{1}_{i,y^{(t)}}$$

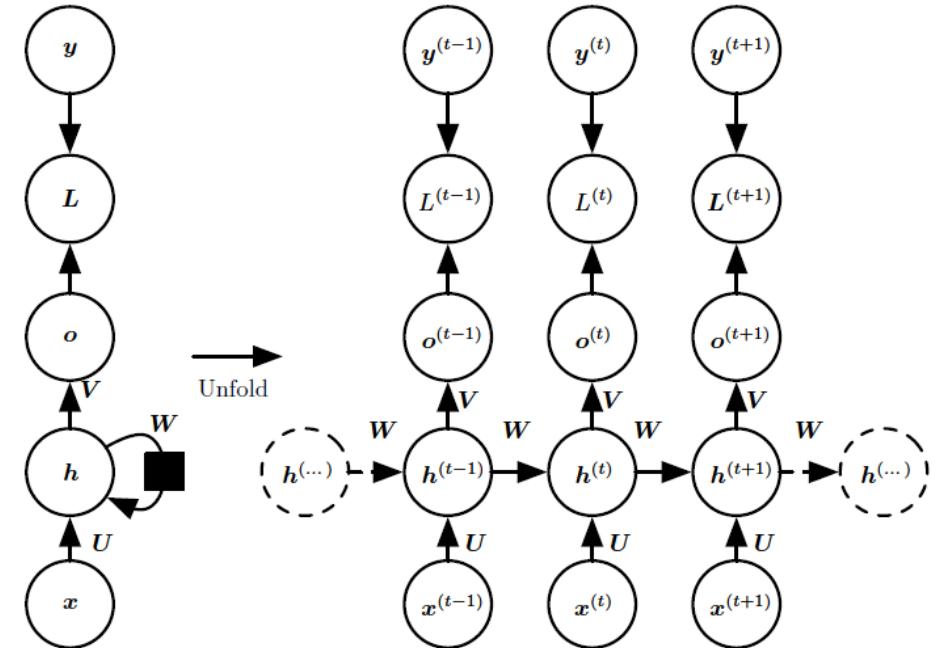
$$\nabla_{\mathbf{h}^{(\tau)}} L = \mathbf{V}^\top \nabla_{\mathbf{o}^{(\tau)}} L$$

BPTT



$$\begin{aligned}
 \text{for } t \leq \tau - 1, \quad \nabla_{\mathbf{h}^{(t)}} L &= \left(\frac{\partial \mathbf{h}^{(t+1)}}{\partial \mathbf{h}^{(t)}} \right)^\top (\nabla_{\mathbf{h}^{(t+1)}} L) + \left(\frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{h}^{(t)}} \right)^\top (\nabla_{\mathbf{o}^{(t)}} L) \\
 &= \mathbf{W}^\top (\nabla_{\mathbf{h}^{(t+1)}} L) \text{diag} \left(1 - (\mathbf{h}^{(t+1)})^2 \right) + \mathbf{V}^\top (\nabla_{\mathbf{o}^{(t)}} L),
 \end{aligned}$$

BPTT



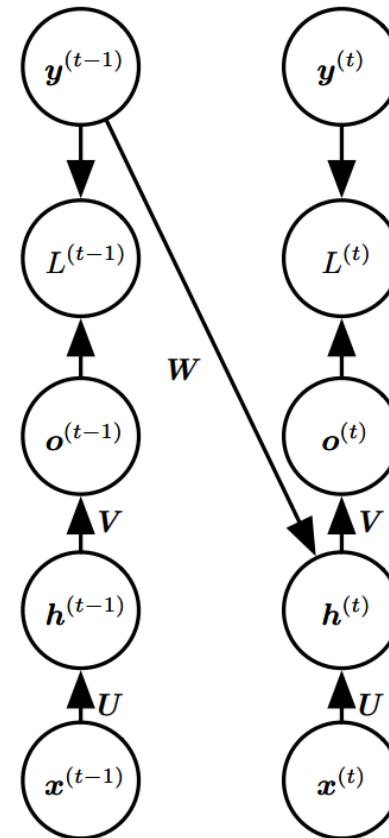
Once we get $\nabla_{\mathbf{h}^{(t)}} L$ for each t , we can calculate the gradient w.r.t. the parameters. For example,

$$\begin{aligned}\nabla_{\mathbf{W}} L &= \sum_t \sum_i \left(\frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{W}^{(t)}} h_i^{(t)} \\ &= \sum_t \text{diag} \left(1 - (\mathbf{h}^{(t)})^2 \right) (\nabla_{\mathbf{h}^{(t)}} L) \mathbf{h}^{(t-1)\top}\end{aligned}$$

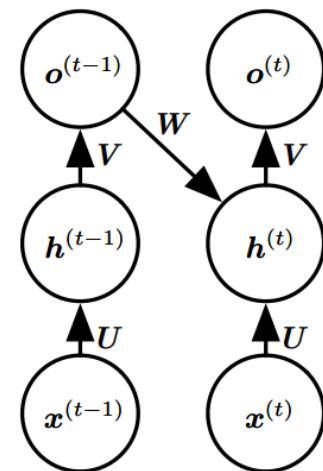
teacher forcing

$$\begin{aligned} \log p(\mathbf{y}^{(1)}, \mathbf{y}^{(2)} \mid \mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \\ = \log p(\mathbf{y}^{(2)} \mid \mathbf{y}^{(1)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}) + \log p(\mathbf{y}^{(1)} \mid \mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \end{aligned}$$

depends on
previous output



Train time



Test time

exploding and vanishing gradients

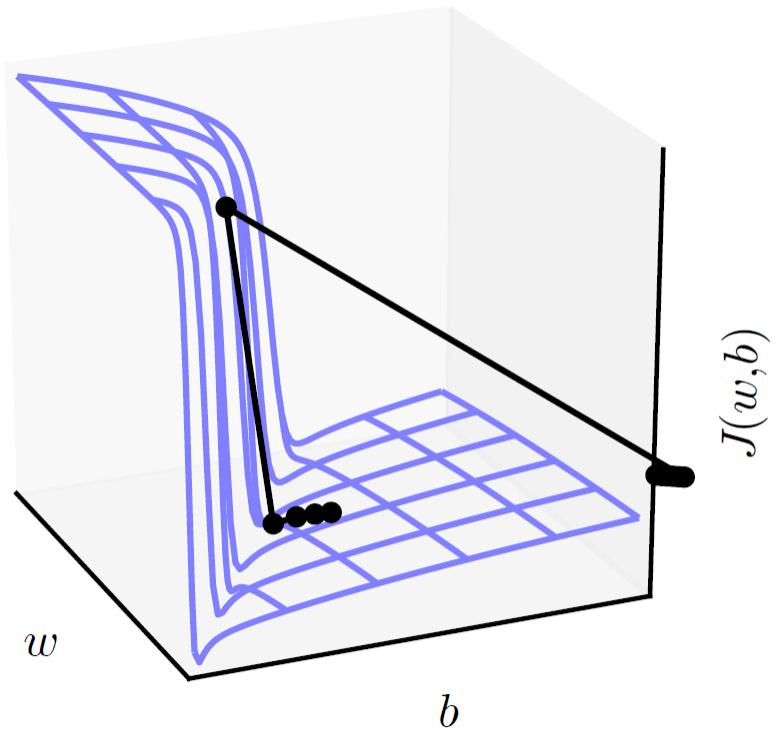
- consider a simple recurrence relation

$$\mathbf{h}^{(t)} = \mathbf{W}^\top \mathbf{h}^{(t-1)}$$

- we have $\mathbf{h}^{(t)} = (\mathbf{W}^t)^\top \mathbf{h}^{(0)}$
- assume $\mathbf{W} = \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^\top$
- we can further simplify it to

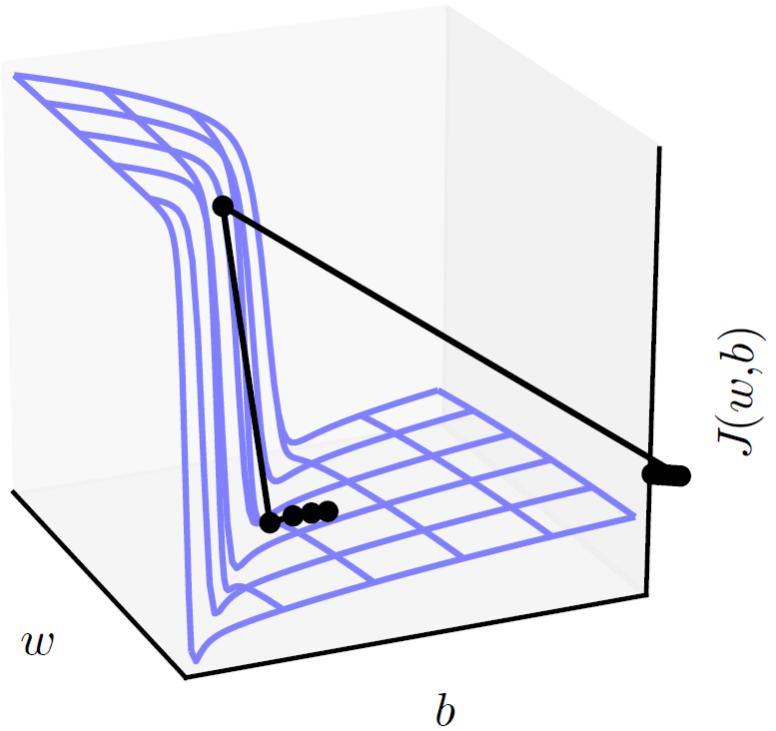
$$\mathbf{h}^{(t)} = \mathbf{Q}^\top \boldsymbol{\Lambda}^t \mathbf{Q} \mathbf{h}^{(0)}$$

exploding gradients

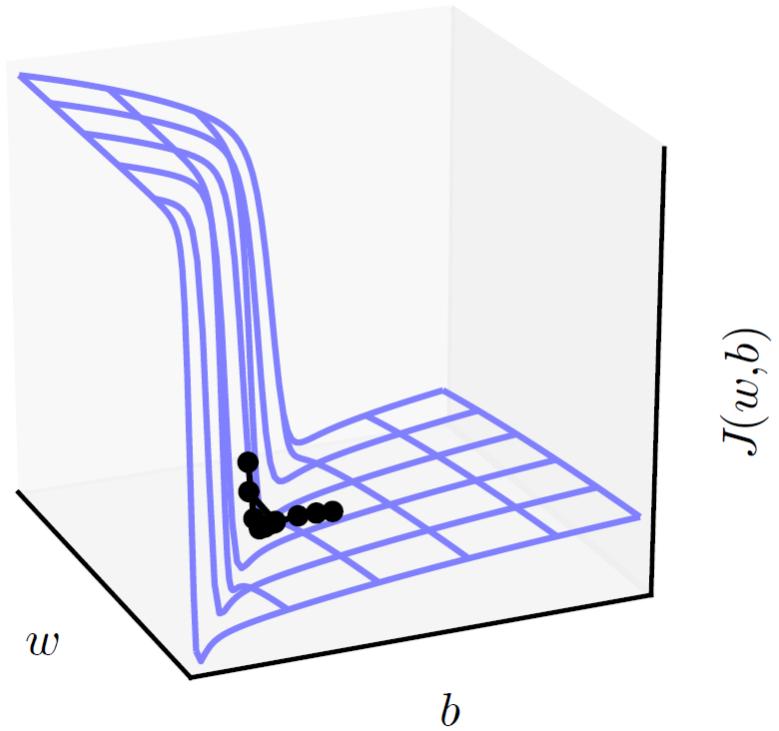


gradient clipping

Without clipping



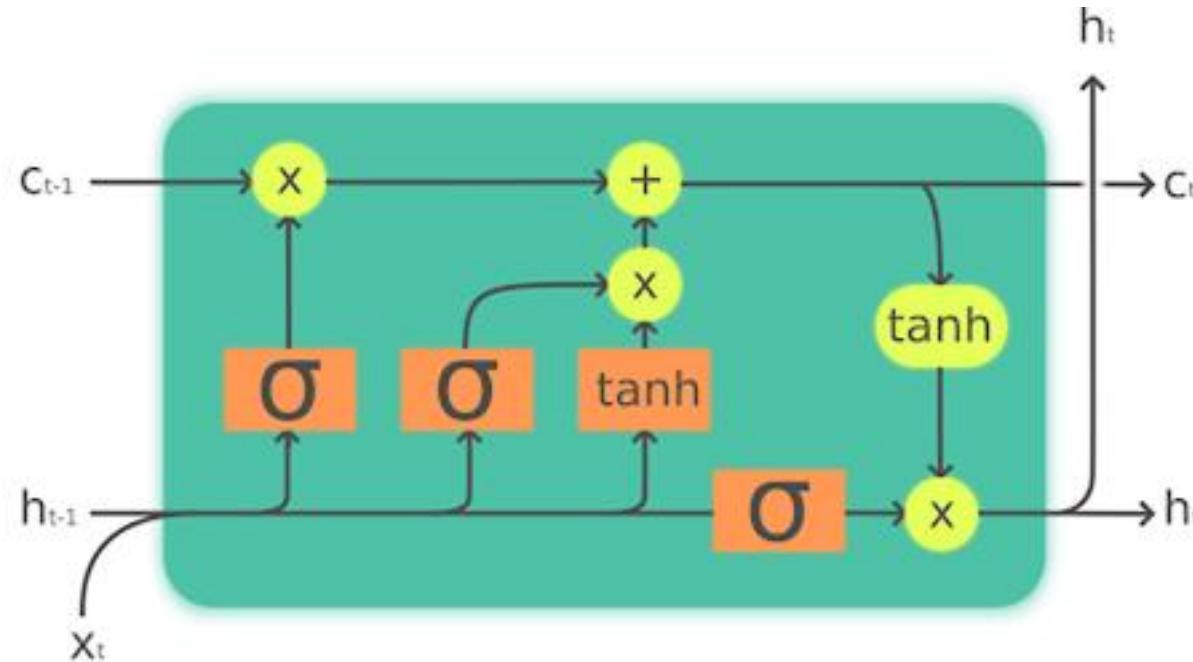
With clipping



if $\|g\| > v$
$$g \leftarrow \frac{gv}{\|g\|}$$

Long Short-Term Memory (LSTM)

Hochreiter & Schmidhuber (1997)



$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f)$$

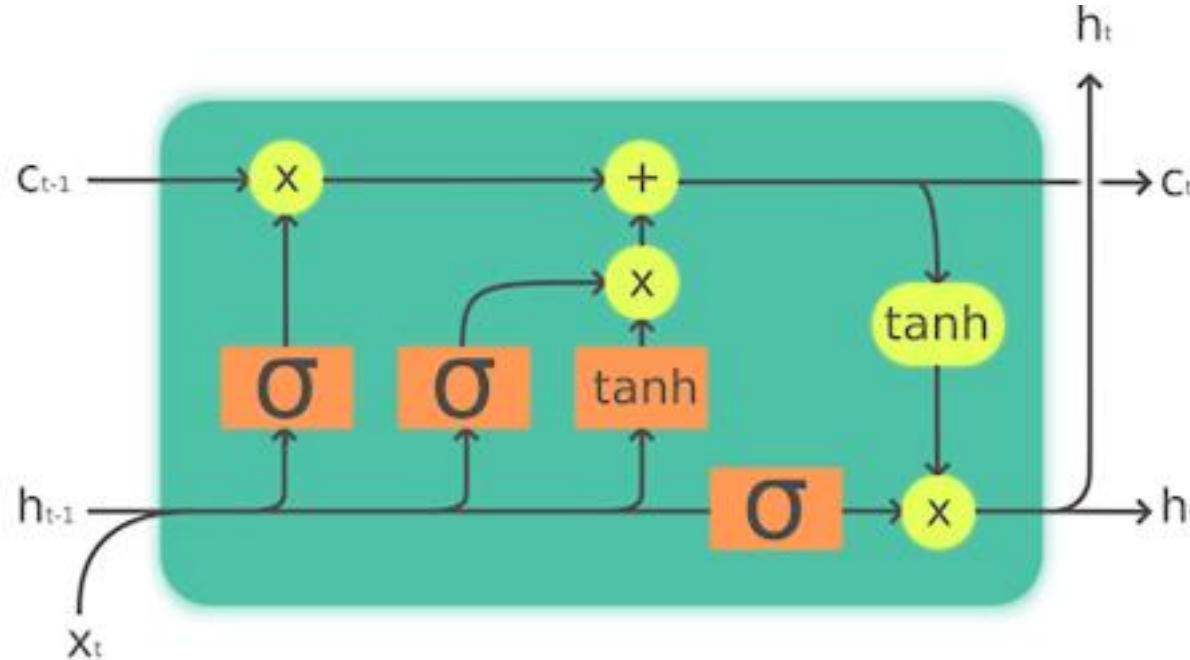
$$i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i)$$

$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o)$$

$$c_t = f_t \circ c_{t-1} + i_t \circ \tanh(W_c x_t + U_c h_{t-1} + b_c)$$

$$h_t = o_t \circ \tanh(c_t)$$

Long Short-Term Memory (LSTM)



Hochreiter & Schmidhuber (1997)

forget gate

$$f_t = \sigma_g(W_f x_t + U_f h_{t-1} + b_f)$$

input gate

$$i_t = \sigma_g(W_i x_t + U_i h_{t-1} + b_i)$$

output gate

$$o_t = \sigma_g(W_o x_t + U_o h_{t-1} + b_o)$$

internal state

$$c_t = f_t \circ c_{t-1} + i_t \circ \tanh(W_c x_t + U_c h_{t-1} + b_c)$$

output

$$h_t = o_t \circ \tanh(c_t)$$

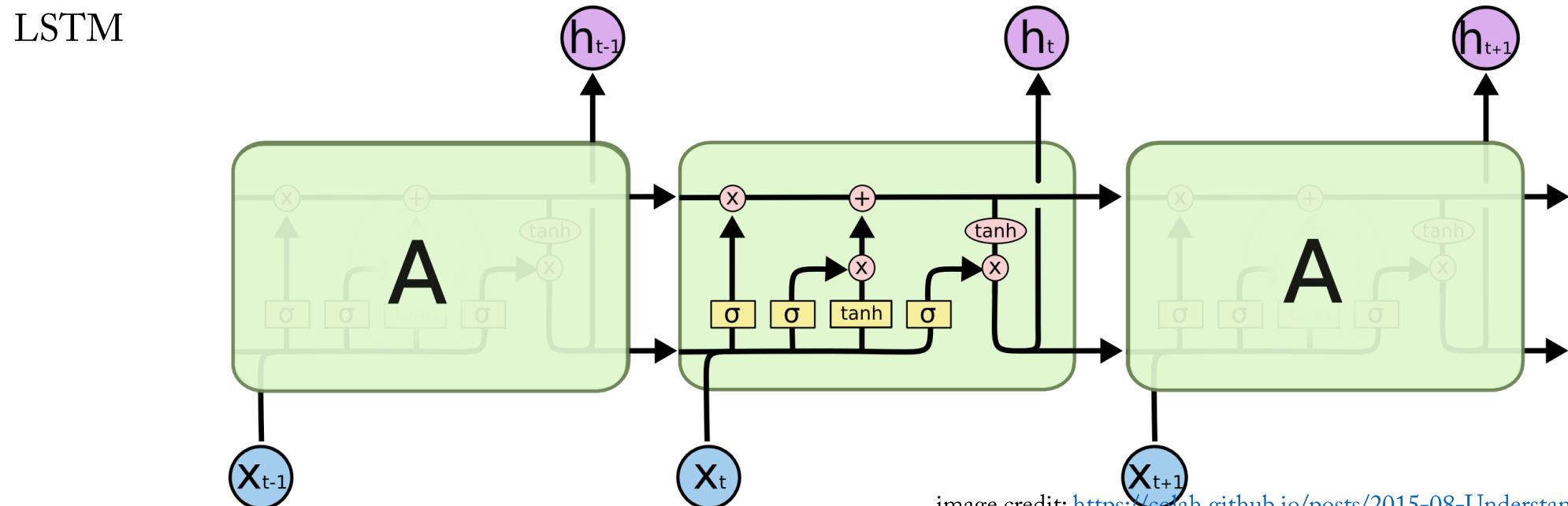
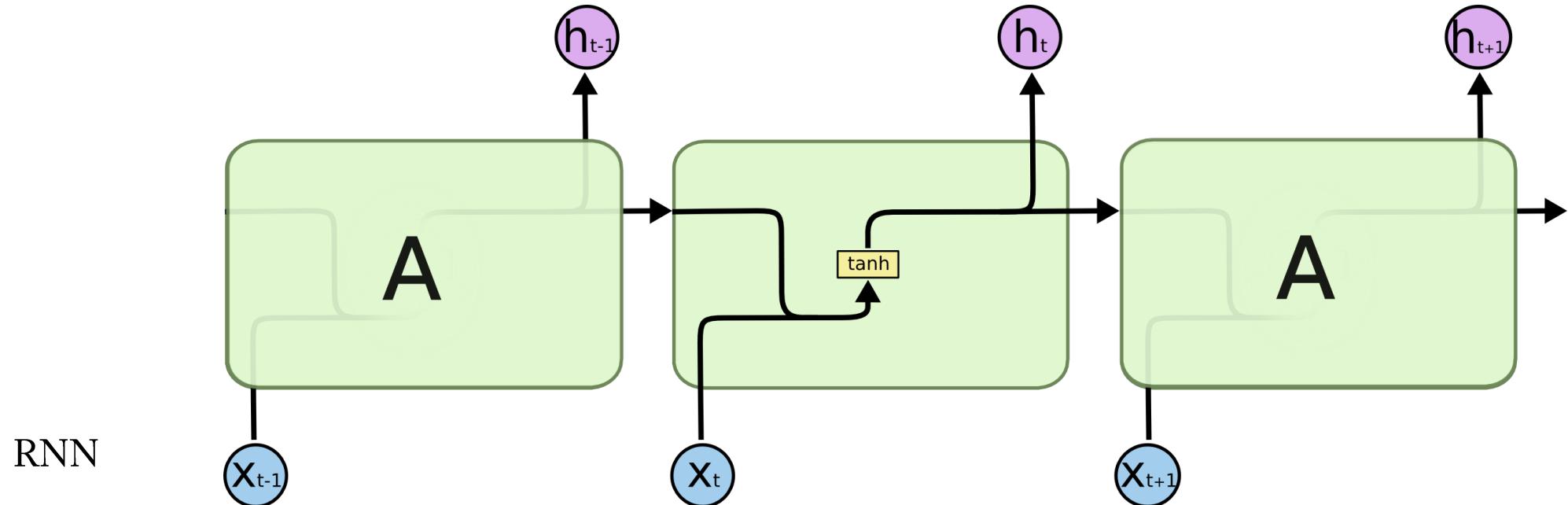
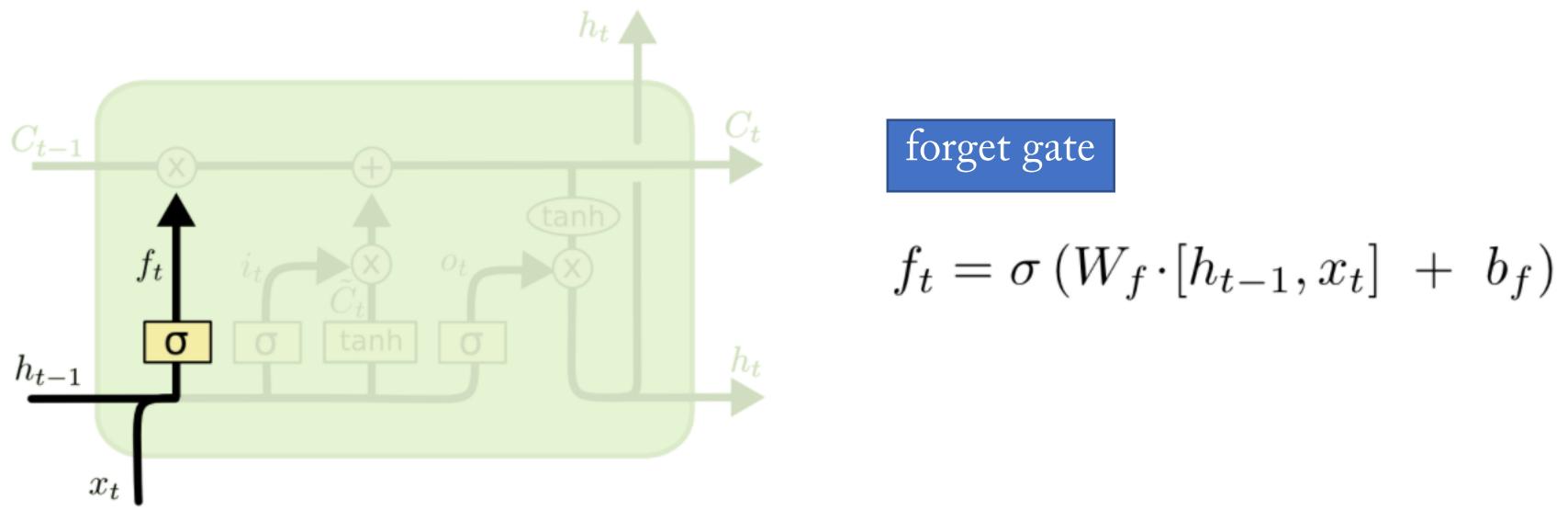
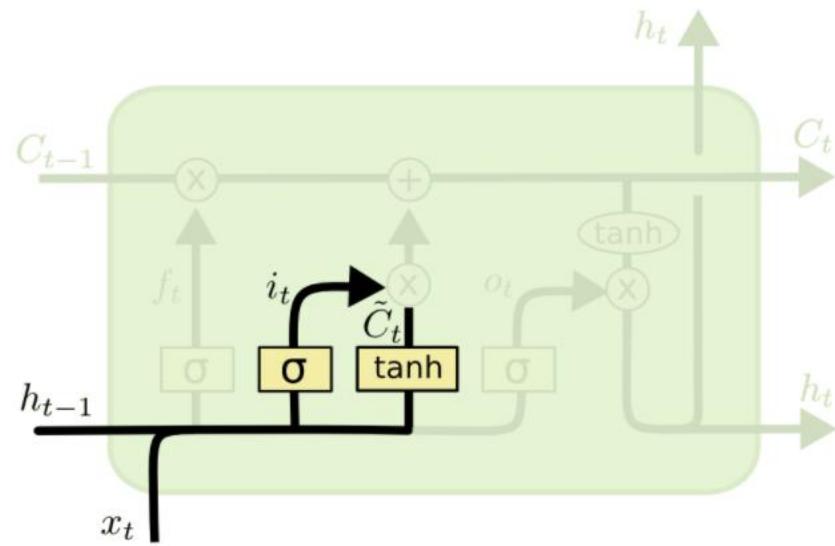


image credit: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

LSTM



LSTM

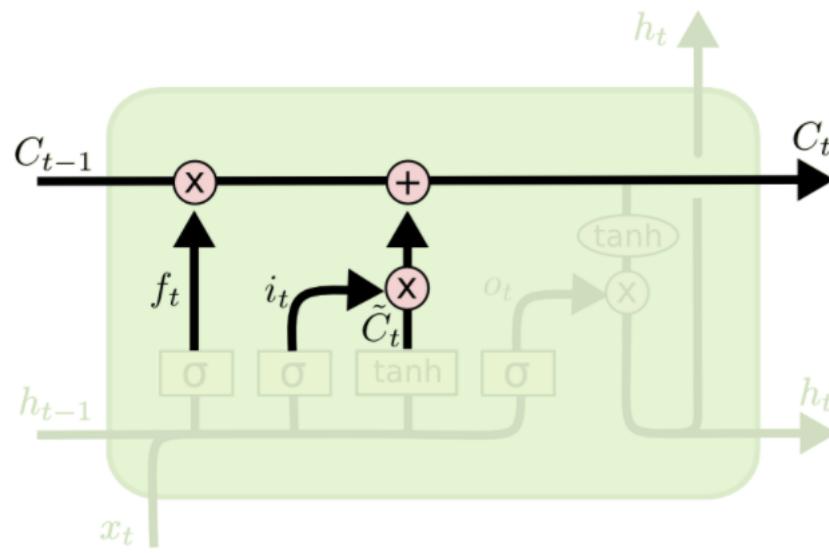


input gate

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

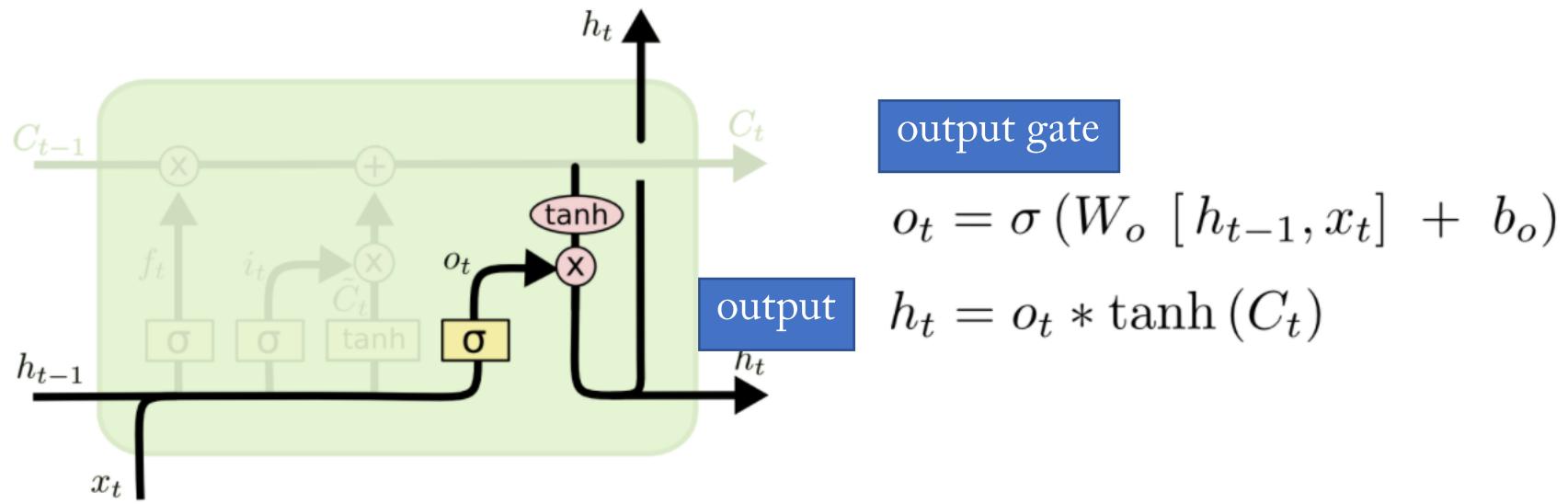
LSTM



internal state

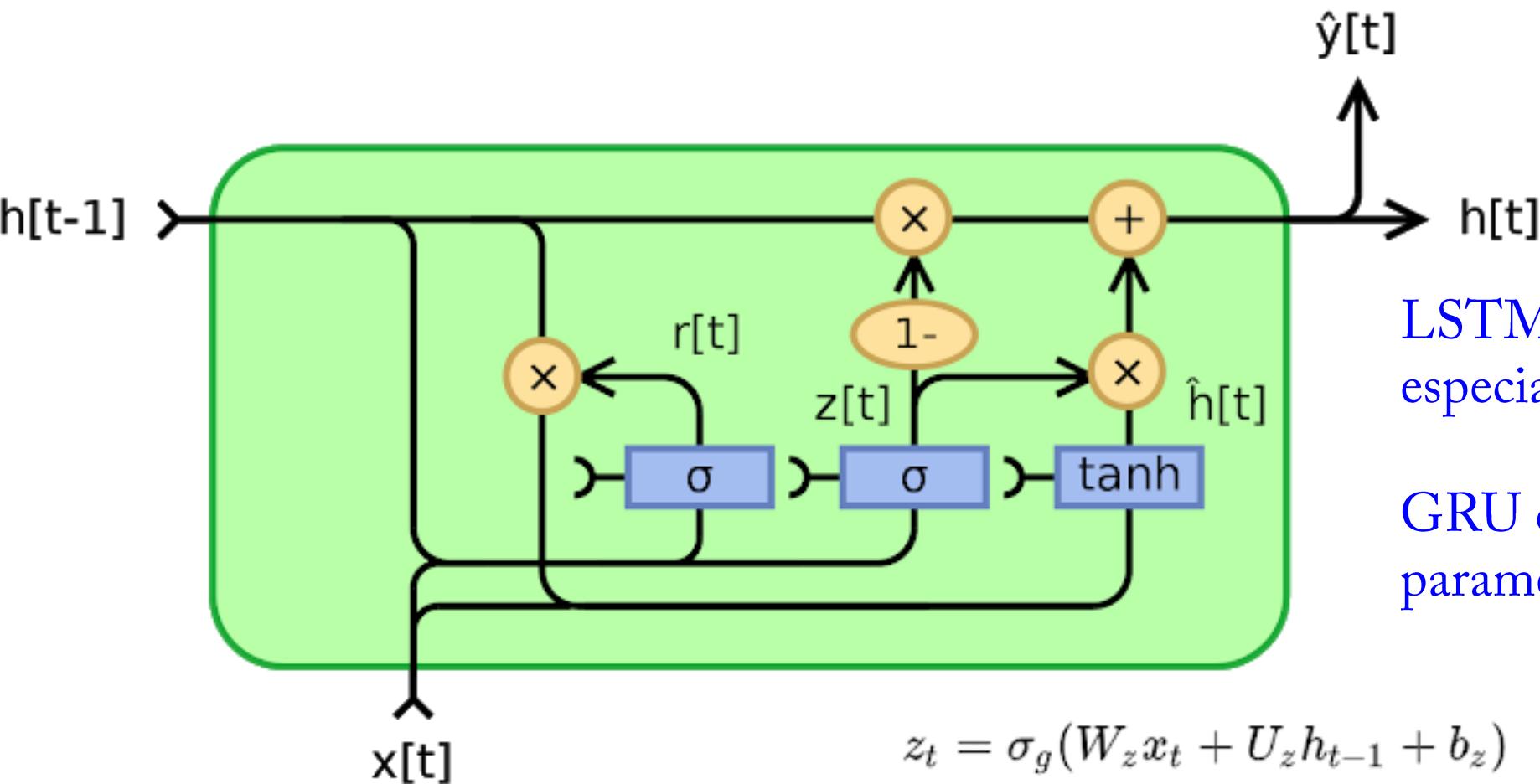
$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

LSTM



Gated Recurrent Unit (GRU)

Cho (2004)



LSTM maybe the first to try
especially for long dependence

GRU can be faster and fewer
parameters

$$z_t = \sigma_g(W_z x_t + U_z h_{t-1} + b_z)$$

$$r_t = \sigma_g(W_r x_t + U_r h_{t-1} + b_r)$$

$$h_t = z_t \odot h_{t-1} + (1 - z_t) \odot \phi_h(W_h x_t + U_h (r_t \odot h_{t-1}) + b_h)$$

6.2 part-of-speech (POS) tagging

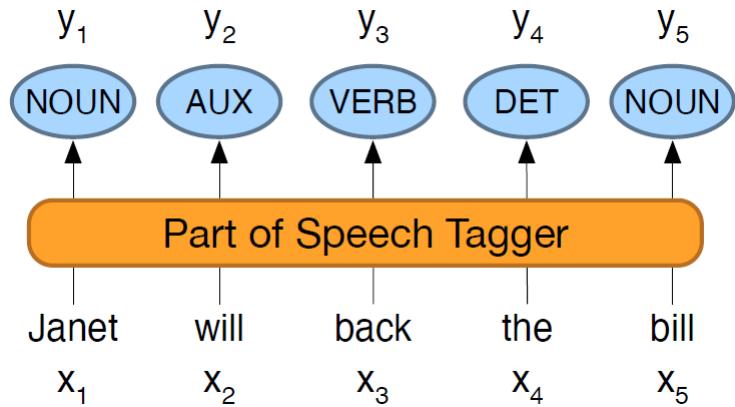
parts of speech

	Tag	Description	Example
Open Class	ADJ	Adjective: noun modifiers describing properties	<i>red, young, awesome</i>
	ADV	Adverb: verb modifiers of time, place, manner	<i>very, slowly, home, yesterday</i>
	NOUN	words for persons, places, things, etc.	<i>algorithm, cat, mango, beauty</i>
	VERB	words for actions and processes	<i>draw, provide, go</i>
	PROPN	Proper noun: name of a person, organization, place, etc..	<i>Regina, IBM, Colorado</i>
Closed Class Words	INTJ	Interjection: exclamation, greeting, yes/no response, etc.	<i>oh, um, yes, hello</i>
	ADP	Adposition (Preposition/Postposition): marks a noun's spacial, temporal, or other relation	<i>in, on, by, under</i>
	AUX	Auxiliary: helping verb marking tense, aspect, mood, etc.,	<i>can, may, should, are</i>
	CCONJ	Coordinating Conjunction: joins two phrases/clauses	<i>and, or, but</i>
	DET	Determiner: marks noun phrase properties	<i>a, an, the, this</i>
	NUM	Numeral	<i>one, two, first, second</i>
	PART	Particle: a function word that must be associated with another word	<i>'s, not, (infinitive) to</i>
	PRON	Pronoun: a shorthand for referring to an entity or event	<i>she, who, I, others</i>
Other	SCONJ	Subordinating Conjunction: joins a main clause with a subordinate clause such as a sentential complement	<i>that, which</i>
	PUNCT	Punctuation	<i>; , ()</i>
	SYM	Symbols like \$ or emoji	<i>\$, %</i>
	X	Other	<i>asdf, qwfg</i>

parts of speech: Penn Treebank tags

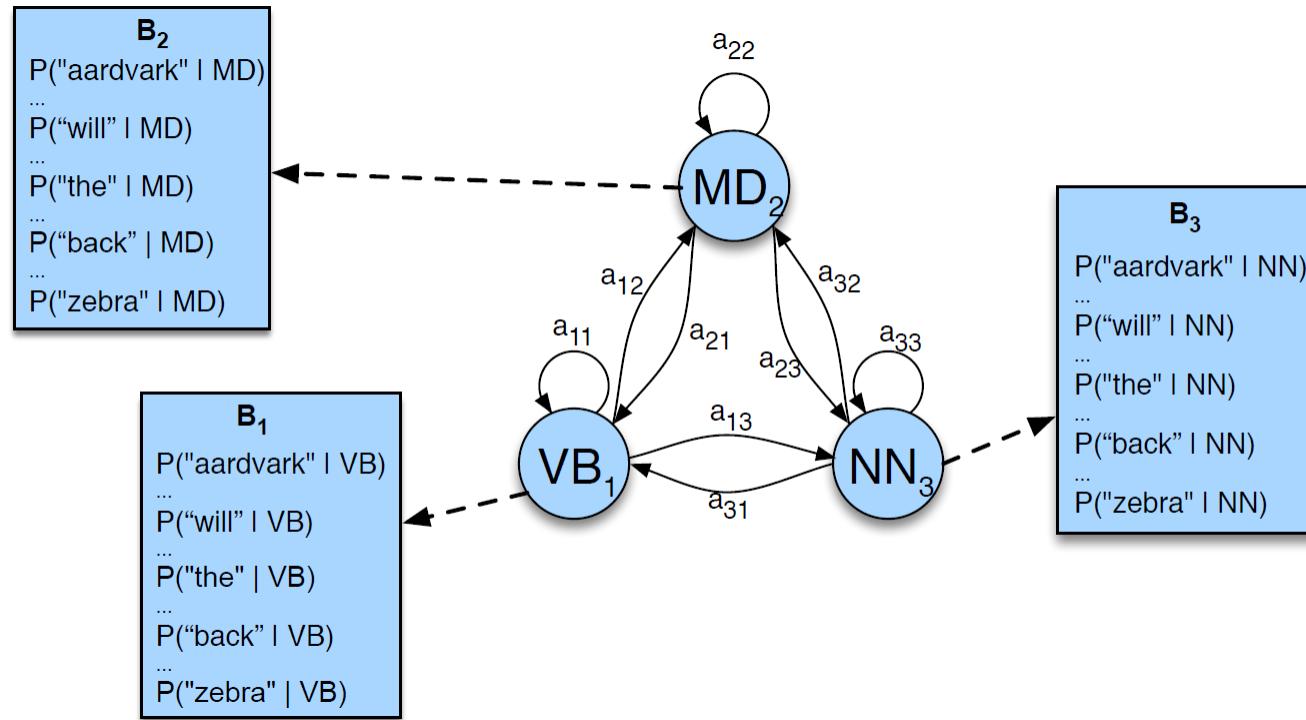
Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coord. conj.	<i>and, but, or</i>	NNP	proper noun, sing.	<i>IBM</i>	TO	“to”	<i>to</i>
CD	cardinal number	<i>one, two</i>	NNPS	proper noun, plu.	<i>Carolinas</i>	UH	interjection	<i>ah, oops</i>
DT	determiner	<i>a, the</i>	NNS	noun, plural	<i>llamas</i>	VB	verb base	<i>eat</i>
EX	existential ‘there’	<i>there</i>	PDT	predeterminer	<i>all, both</i>	VBD	verb past tense	<i>ate</i>
FW	foreign word	<i>mea culpa</i>	POS	possessive ending	<i>'s</i>	VBG	verb gerund	<i>eating</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	PRP	personal pronoun	<i>I, you, he</i>	VBN	verb past participle	<i>eaten</i>
JJ	adjective	<i>yellow</i>	PRP\$	possess. pronoun	<i>your, one's</i>	VBP	verb non-3sg-pr	<i>eat</i>
JJR	comparative adj	<i>bigger</i>	RB	adverb	<i>quickly</i>	VBZ	verb 3sg pres	<i>eats</i>
JJS	superlative adj	<i>wildest</i>	RBR	comparative adv	<i>faster</i>	WDT	wh-determ.	<i>which, that</i>
LS	list item marker	<i>1, 2, One</i>	RBS	superlatv. adv	<i>fastest</i>	WP	wh-pronoun	<i>what, who</i>
MD	modal	<i>can, should</i>	RP	particle	<i>up, off</i>	WP\$	wh-possess.	<i>whose</i>
NN	sing or mass noun	<i>llama</i>	SYM	symbol	<i>+, %, &</i>	WRB	wh-adverb	<i>how, where</i>

tagging



- ❑ Naïve idea: check the dictionary for tagging
- ❑ Problem: not unique
- ✓ Sequential labeling: tagging according to the context

Hidden-Markov Model (HMM)



Hidden-Markov Model (HMM)

$$Q = q_1 q_2 \dots q_N$$

a set of N **states**

$$A = a_{11} \dots a_{ij} \dots a_{NN}$$

a **transition probability matrix** A , each a_{ij} representing the probability of moving from state i to state j , s.t. $\sum_{j=1}^N a_{ij} = 1 \quad \forall i$

$$O = o_1 o_2 \dots o_T$$

a sequence of T **observations**, each one drawn from a vocabulary $V = v_1, v_2, \dots, v_V$

$$B = b_i(o_t)$$

a sequence of **observation likelihoods**, also called **emission probabilities**, each expressing the probability of an observation o_t being generated from a state q_i

$$\pi = \pi_1, \pi_2, \dots, \pi_N$$

an **initial probability distribution** over states. π_i is the probability that the Markov chain will start in state i . Some states j may have $\pi_j = 0$, meaning that they cannot be initial states. Also, $\sum_{i=1}^n \pi_i = 1$

Hidden-Markov Model (HMM)

- Markov assumption:

$$P(q_i|q_1, \dots, q_{i-1}) = P(q_i|q_{i-1})$$

- output independence:

$$P(o_i|q_1, \dots, q_i, \dots, q_T, o_1, \dots, o_i, \dots, o_T) = P(o_i|q_i)$$

Hidden-Markov Model (HMM)

- transition probability
 - ❑ e.g. modal verbs [will] are very likely to be followed by a verb in the base form [eat]
 - ❑ In the WSJ corpus, we can calculate

$$P(VB|MD) = \frac{C(MD, VB)}{C(MD)} = \frac{10471}{13124} = .80$$

Hidden-Markov Model (HMM)

- emission probability
 - In the WSJ corpus, we can calculate

$$P(will|MD) = \frac{C(MD, will)}{C(MD)} = \frac{4046}{13124} = .31$$

Hidden-Markov Model (HMM)

- decoding: given as input an HMM and a sequence of observations, find the most probable sequence of states

$$\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} P(t_1 \dots t_n | w_1 \dots w_n)$$

Hidden-Markov Model (HMM)

$$\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} P(t_1 \dots t_n | w_1 \dots w_n)$$



$$\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} \frac{P(w_1 \dots w_n | t_1 \dots t_n) P(t_1 \dots t_n)}{P(w_1 \dots w_n)}$$



$$\hat{t}_{1:n} = \operatorname{argmax}_{t_1 \dots t_n} P(w_1 \dots w_n | t_1 \dots t_n) P(t_1 \dots t_n)$$

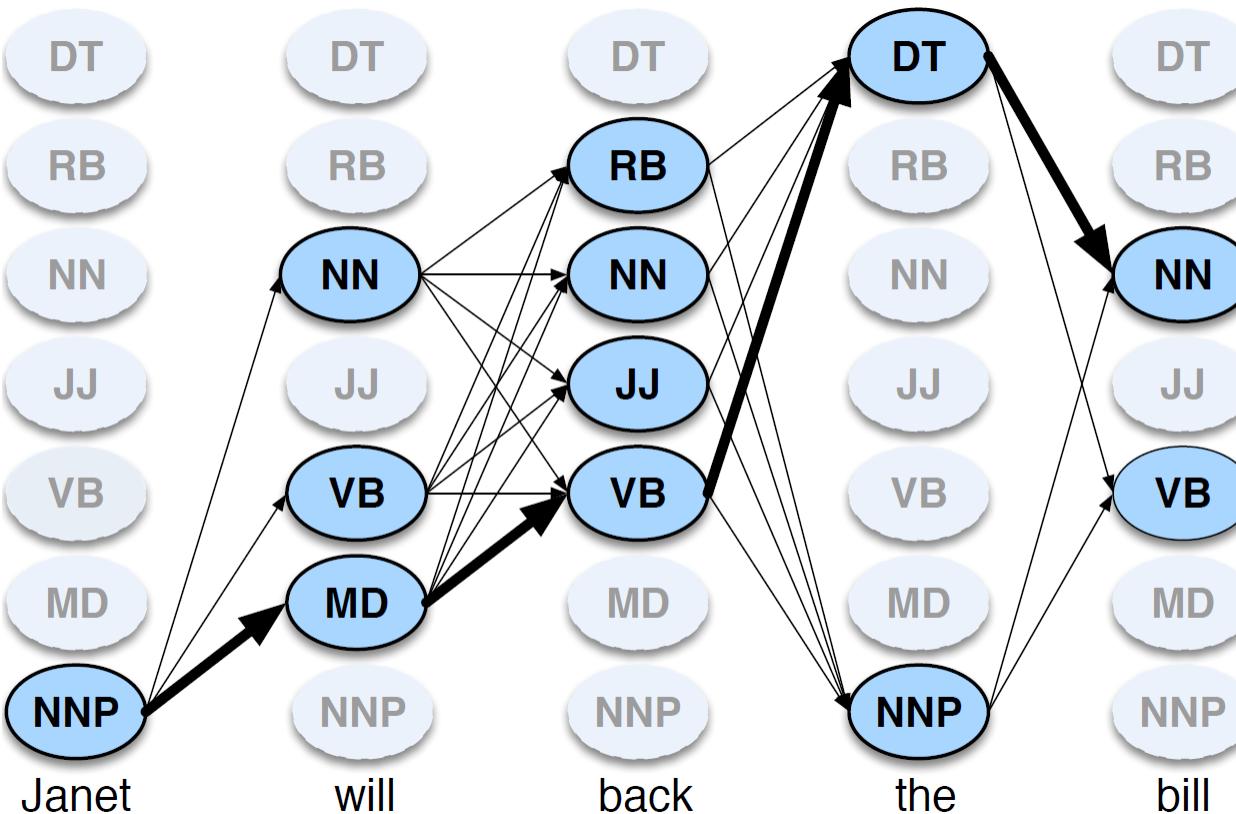
$$\prod_{i=1}^n P(w_i | t_i)$$

$$\prod_{i=1}^n P(t_i | t_{i-1})$$

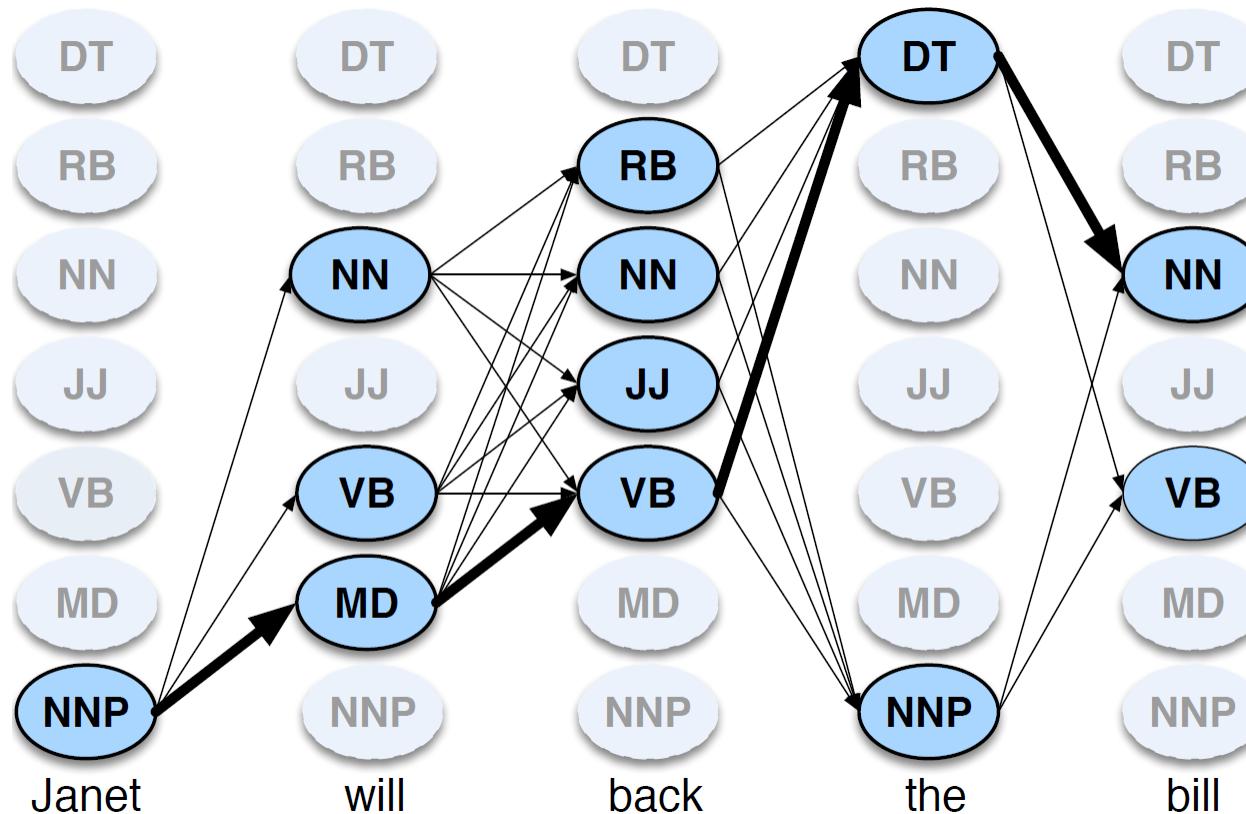
Hidden-Markov Model (HMM)

$$\hat{t}_{1:n} = \underset{t_1 \dots t_n}{\operatorname{argmax}} P(t_1 \dots t_n | w_1 \dots w_n) \approx \underset{t_1 \dots t_n}{\operatorname{argmax}} \prod_{i=1}^n \overbrace{P(w_i | t_i)}^{\text{emission transition}} \overbrace{P(t_i | t_{i-1})}^{\text{transition}}$$

Viterbi algorithm (used on top of HMM in this case)



Viterbi algorithm



each cell $v_t(j) := \max_{q_1, \dots, q_{t-1}} P(q_1, \dots, q_{t-1}; o_1, \dots, o_t, q_t = j | \text{HMM}) = \max_{i=1, \dots, N} v_{t-1}(i) a_{ij} b_j(o_t)$

Viterbi algorithm

transition probability a_{ij}

	NNP	MD	VB	JJ	NN	RB	DT
<S>	0.2767	0.0006	0.0031	0.0453	0.0449	0.0510	0.2026
NNP	0.3777	0.0110	0.0009	0.0084	0.0584	0.0090	0.0025
MD	0.0008	0.0002	0.7968	0.0005	0.0008	0.1698	0.0041
VB	0.0322	0.0005	0.0050	0.0837	0.0615	0.0514	0.2231
JJ	0.0366	0.0004	0.0001	0.0733	0.4509	0.0036	0.0036
NN	0.0096	0.0176	0.0014	0.0086	0.1216	0.0177	0.0068
RB	0.0068	0.0102	0.1011	0.1012	0.0120	0.0728	0.0479
DT	0.1147	0.0021	0.0002	0.2157	0.4744	0.0102	0.0017

Viterbi algorithm

state observation likelihood $b_j(o)$

	Janet	will	back	the	bill
NNP	0.000032	0	0	0.000048	0
MD	0	0.308431	0	0	0
VB	0	0.000028	0.000672	0	0.000028
JJ	0	0	0.000340	0	0
NN	0	0.000200	0.000223	0	0.002337
RB	0	0	0.010446	0	0
DT	0	0	0	0.506099	0

Viterbi algorithm

- Let's illustrate the algorithm using a simpler and unrealistic example:

“Cats cannot fly.”

Initial distribution and transition probabilities

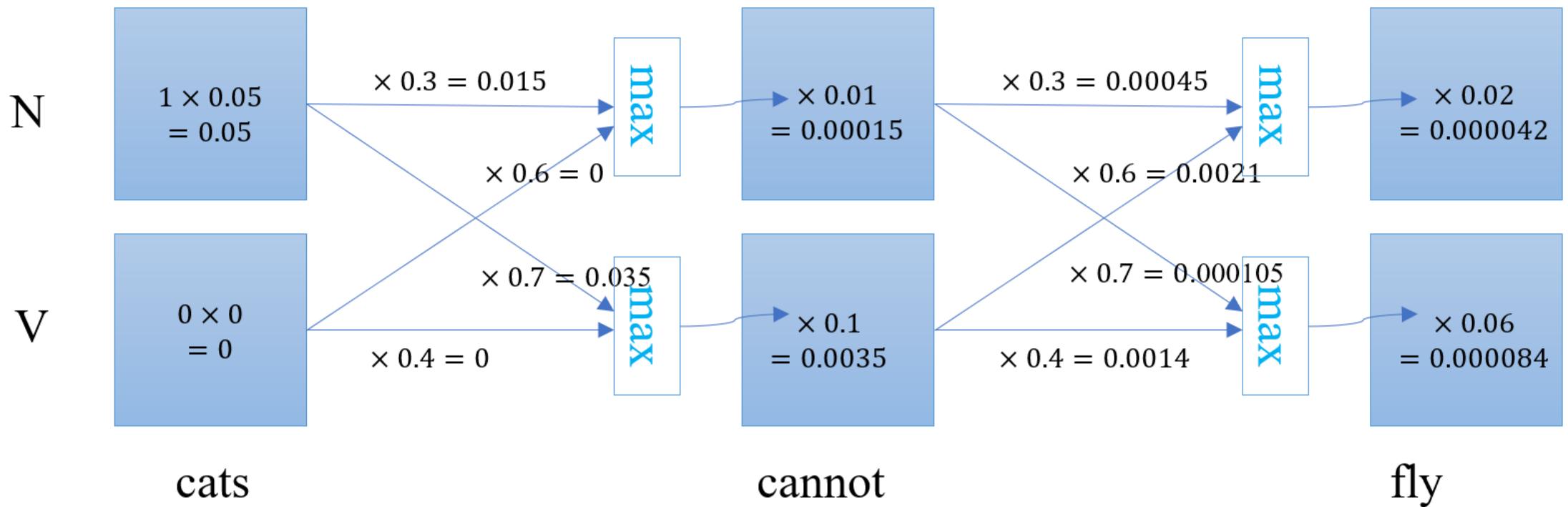
	N	V
initial	1	0
N	0.3	0.7
V	0.6	0.4

Emission probabilities

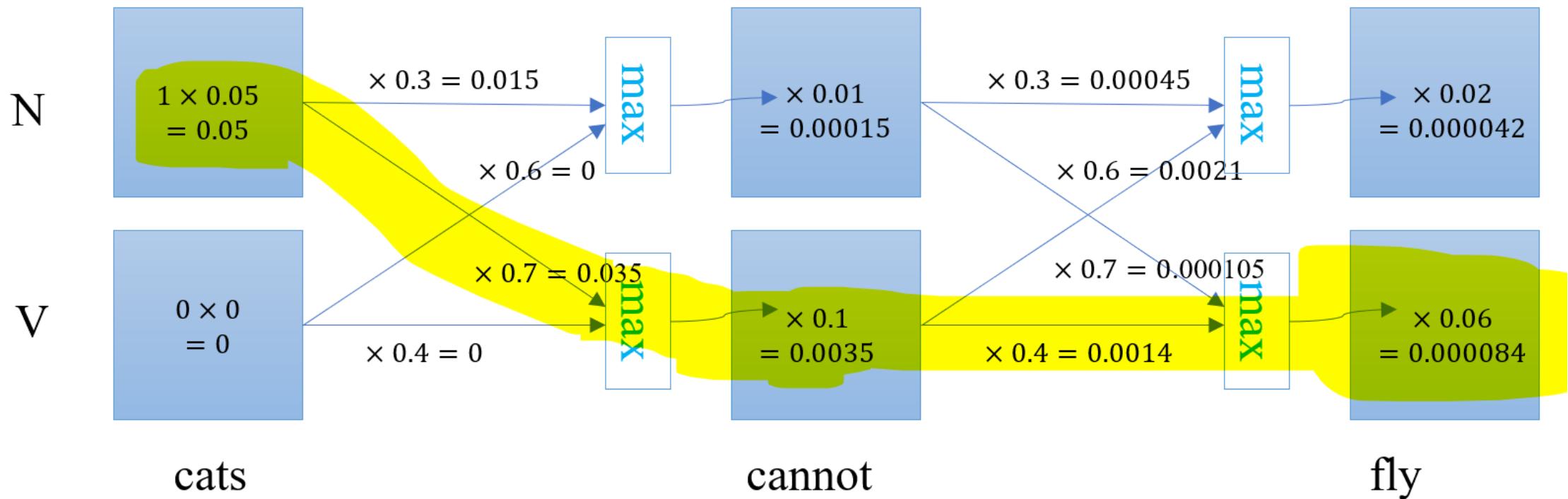
	cats	cannot	fly
N	0.05	0.01	0.02
V	0	0.1	0.06

Viterbi algorithm

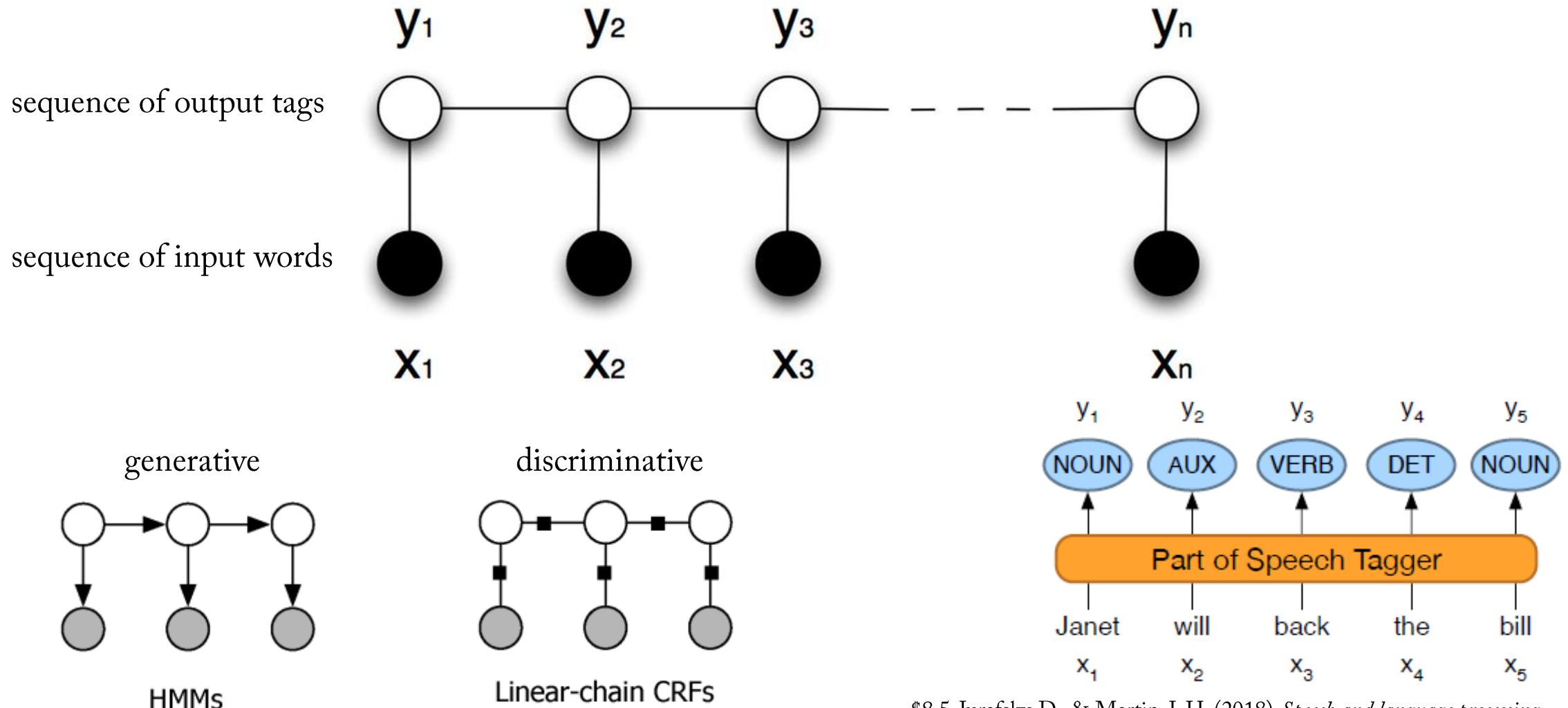
$$\text{each cell } v_t(j) := \max_{q_1, \dots, q_{t-1}} P(q_1, \dots, q_{t-1}; o_1, \dots, o_t, q_t = j | \text{HMM}) = \max_{i=1, \dots, N} v_{t-1}(i) a_{ij} b_j(o_t)$$



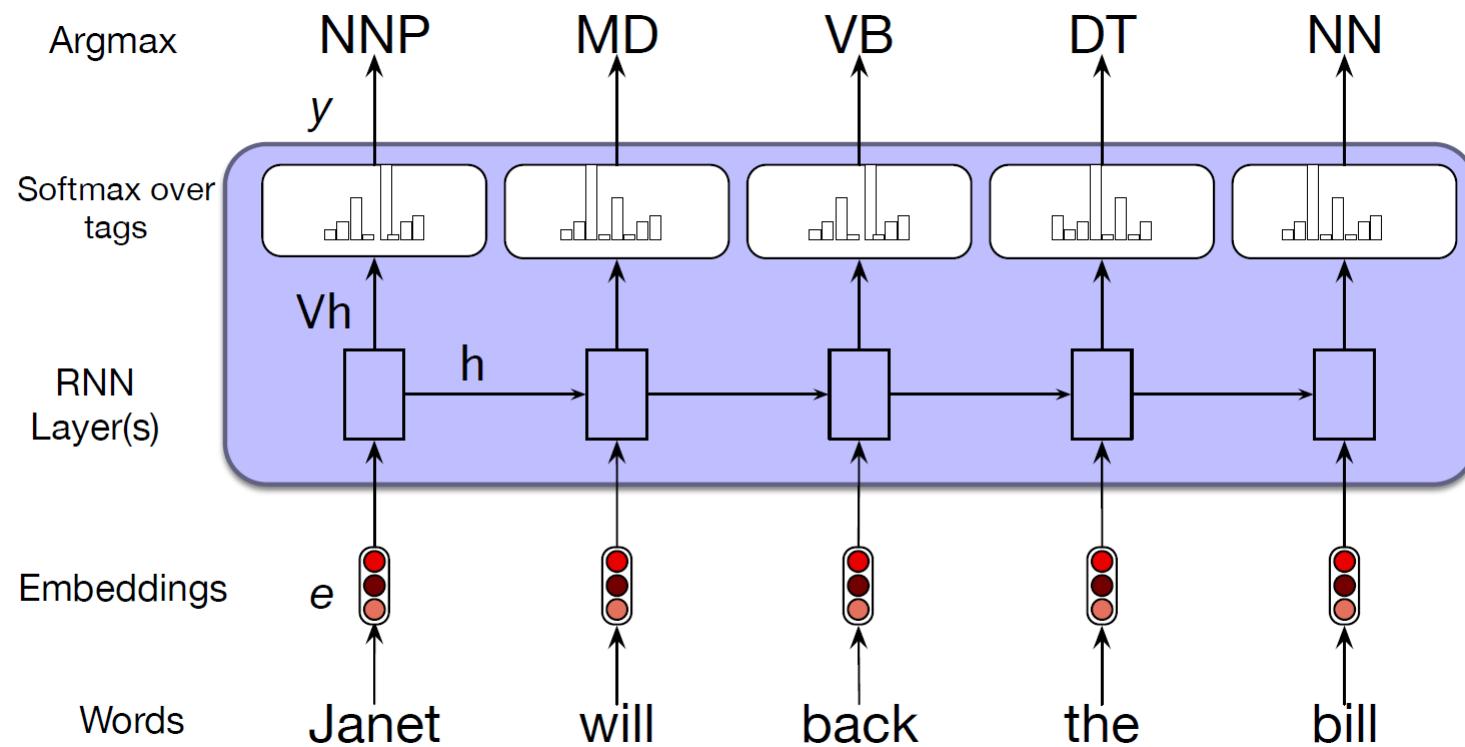
Viterbi algorithm



conditional random field (CRF)



RNN for POS tagging



Let's use LSTM for POS tagging!

- click [here](#)

6.3 RNN for forecasting

RNN for forecasting

Uber Engineering

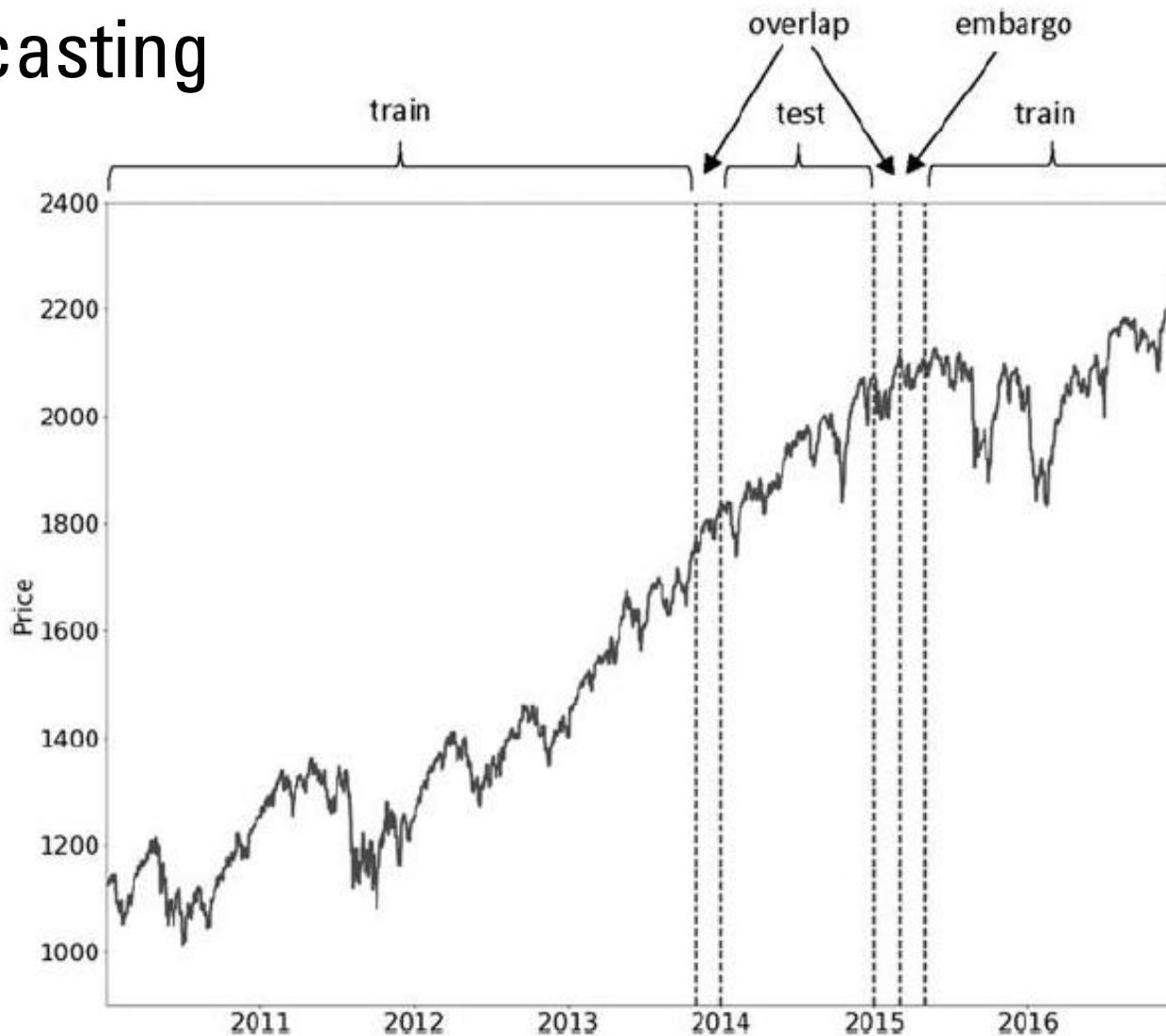
Uber Data

Uber Data

M4 Forecasting Competition: Introducing a New Hybrid ES-RNN Model

Uber's business depends on accurate forecasting. For instance, we use forecasting to predict the expected supply of drivers and demands of riders in the 600+ cities we operate in, to identify when our systems are having outages, to ensure we always have enough **customer obsession** agents managing our support systems, and of course, to plan our business expenditures.

RNN for forecasting



6.4 machine translation

machine translation: difficulty

English: *He wrote a letter to a friend*

Japanese: *tomodachi ni tegami-o kaita*
friend to letter wrote

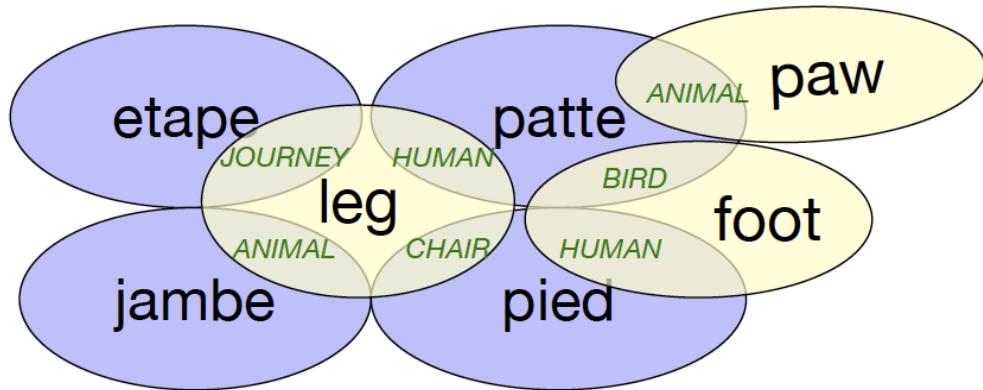
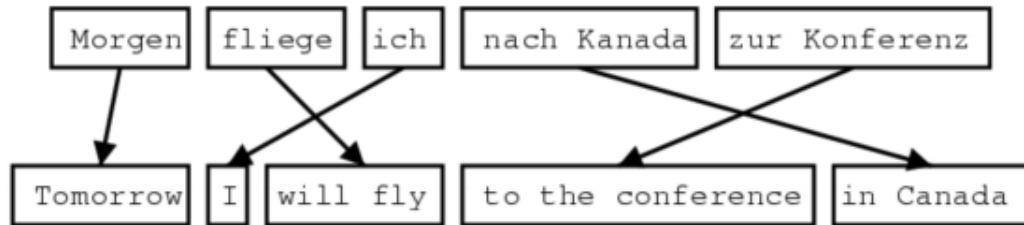
machine translation: difficulty

大会/General Assembly 在/on 1982年/1982 12月/December 10日/10 通过了/adopted 第37号/37th 决议/resolution , 核准了/approved 第二次/second 探索/exploration 及/and 和平peaceful 利用/using 外层空间/outer space 会议/conference 的/of 各项/various 建议/suggestions 。

On 10 December 1982 , the General Assembly adopted resolution 37 in which it endorsed the recommendations of the Second United Nations Conference on the Exploration and Peaceful Uses of Outer Space .

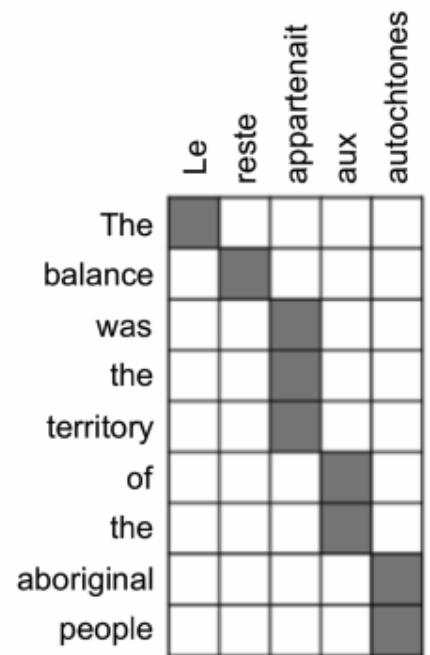
encoder-decoder handles well this sort of difficulty

machine translation: difficulty



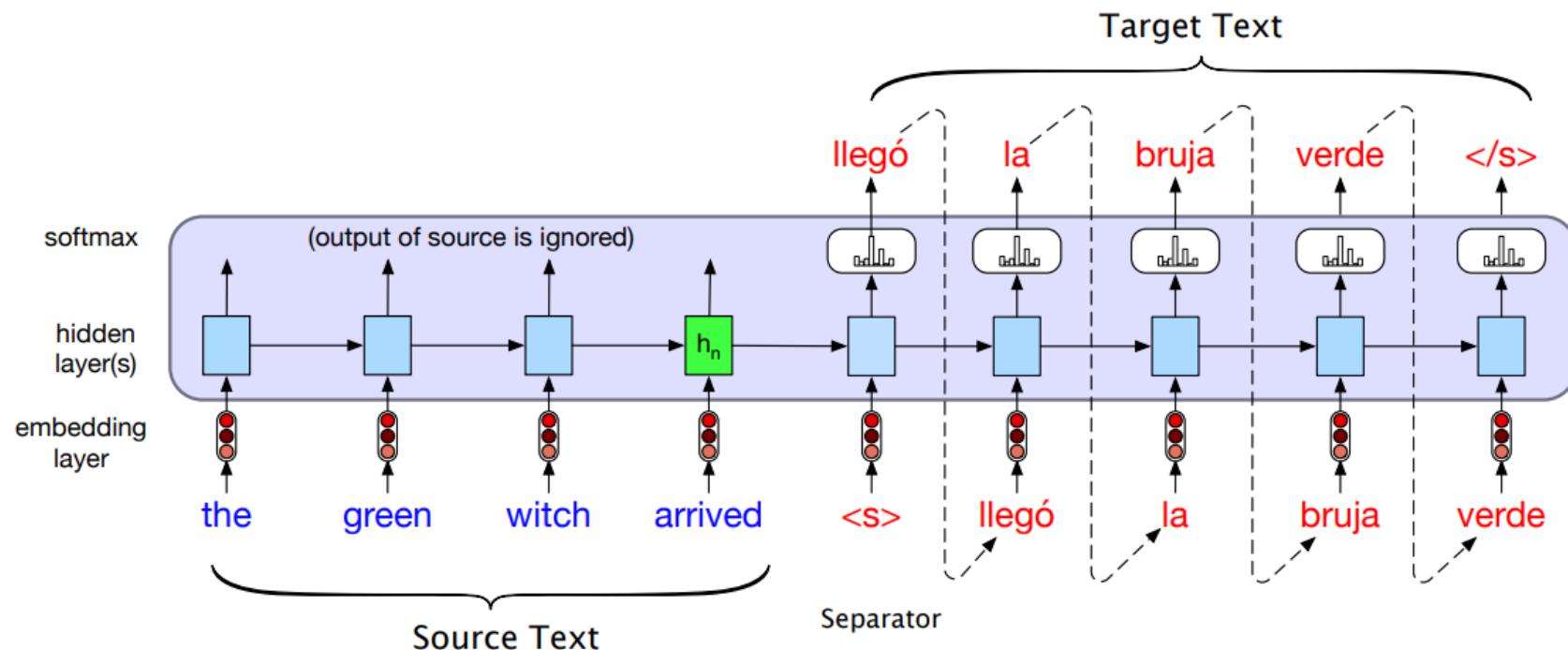
The ————— Le
balance ————— reste
was —————
the ————— appartenait
territory —————
of ————— aux
the —————
aboriginal ————— autochtones
people —————

many-to-one
alignments

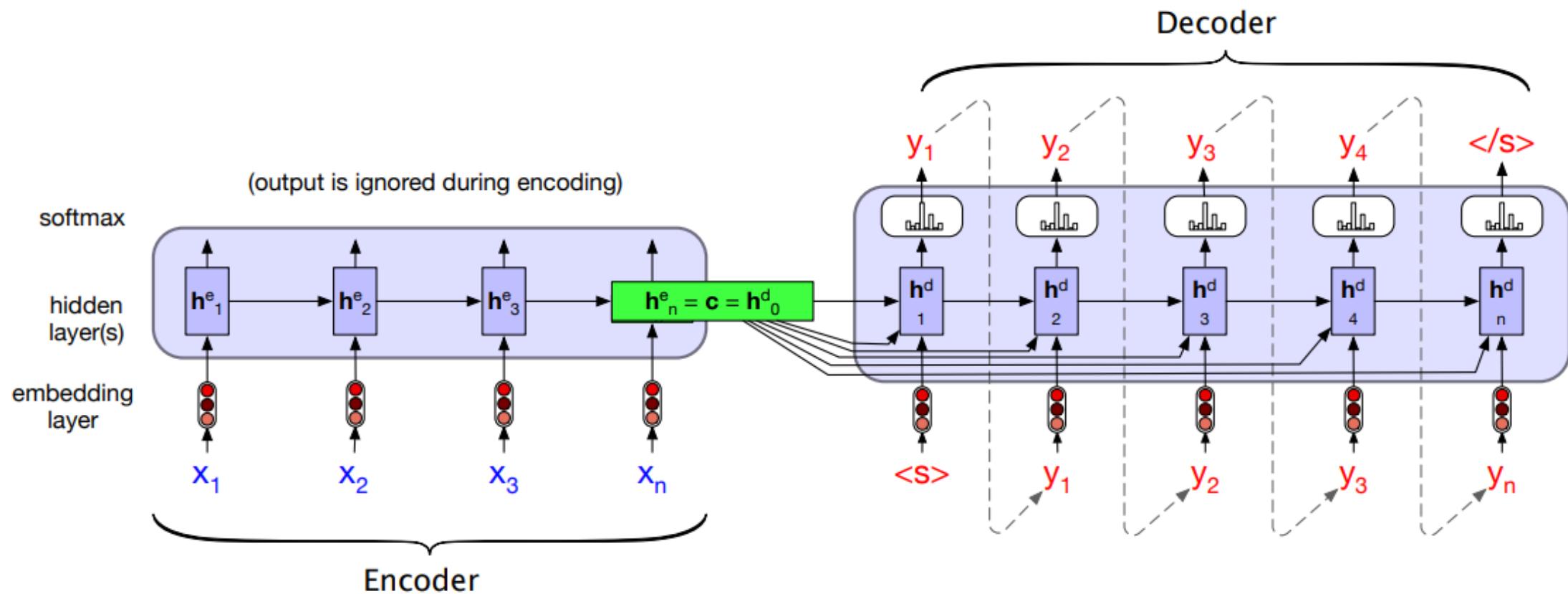


machine translation

Translating a single sentence (inference time) in the basic RNN version of encoder-decoder (sequence-to-sequence) approach to machine translation. Source and target sentences are concatenated with a separator token in between, and the decoder uses context information from the encoder's last hidden state.



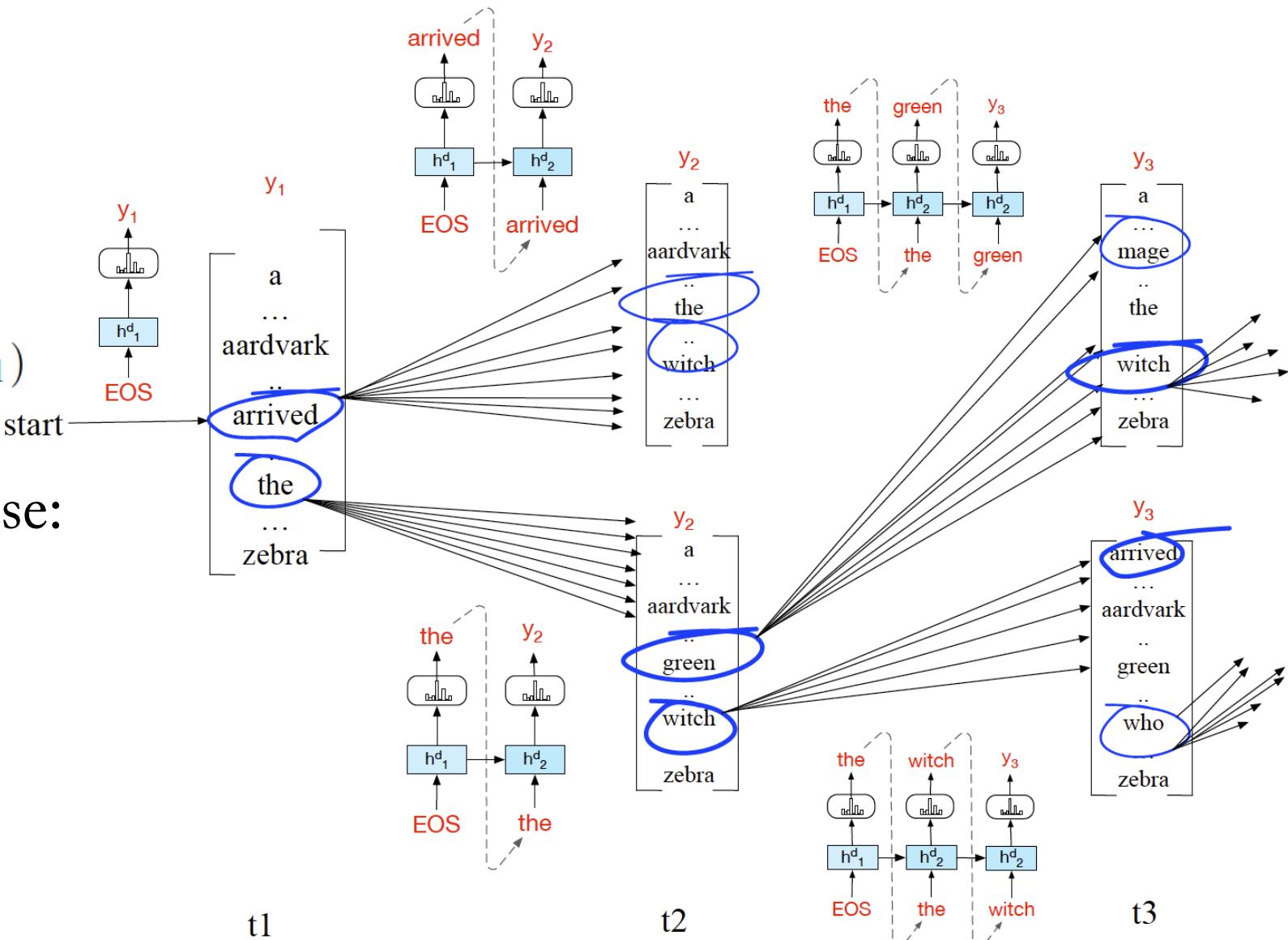
machine translation



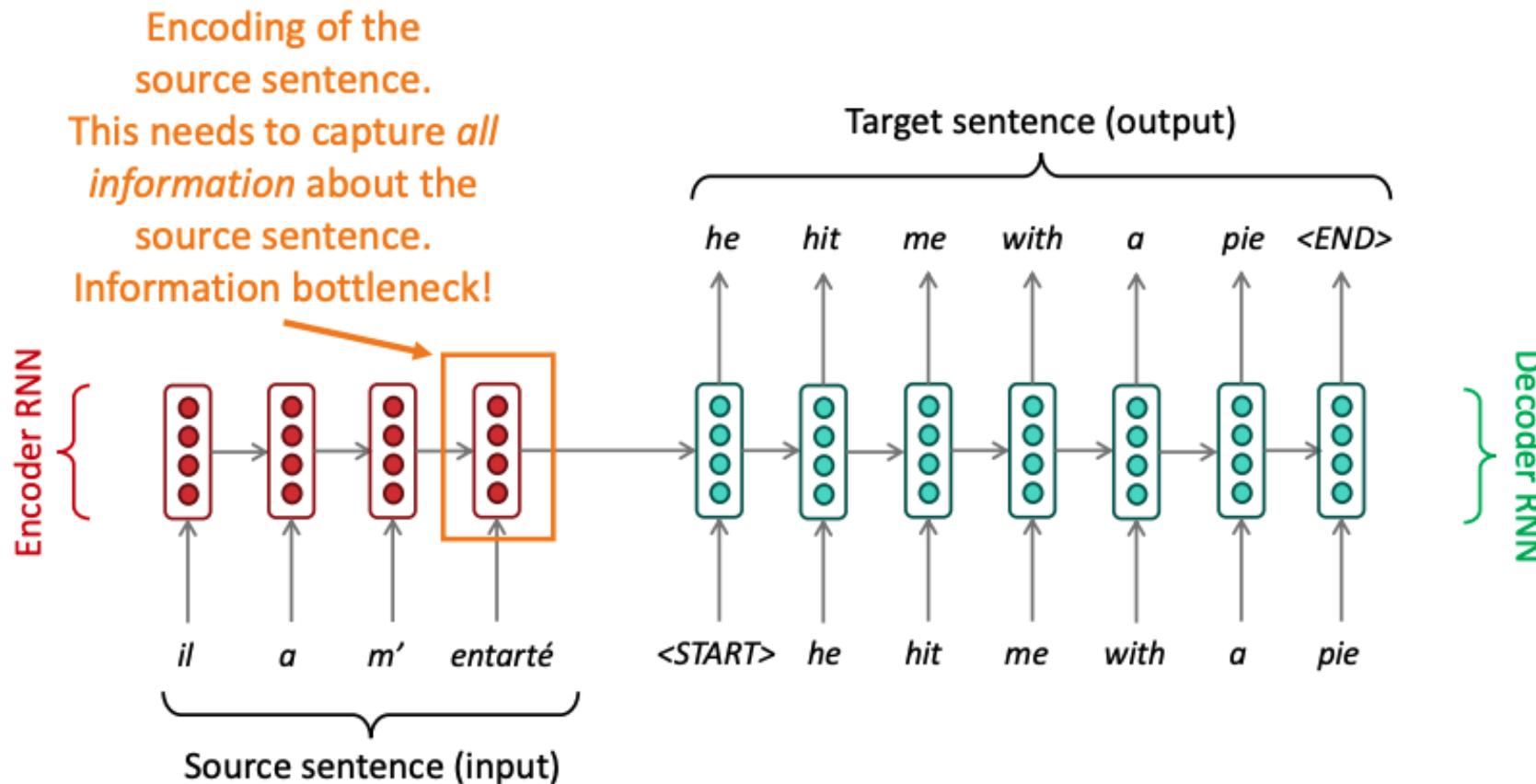
decoding

- greedy:

$$\hat{y}_t = \operatorname{argmax}_{w \in V} P(w|x, y_1 \dots y_{t-1})$$

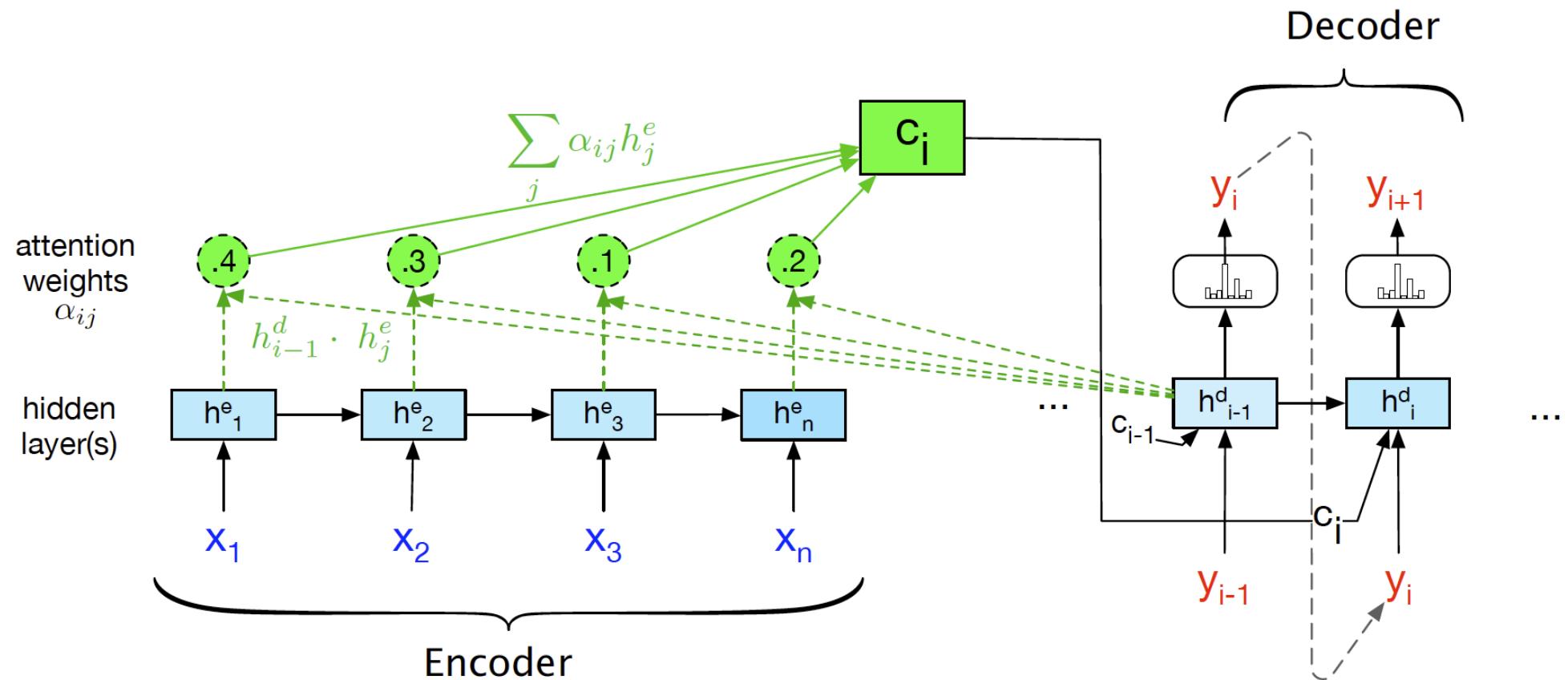


drawback of encoder-decoder



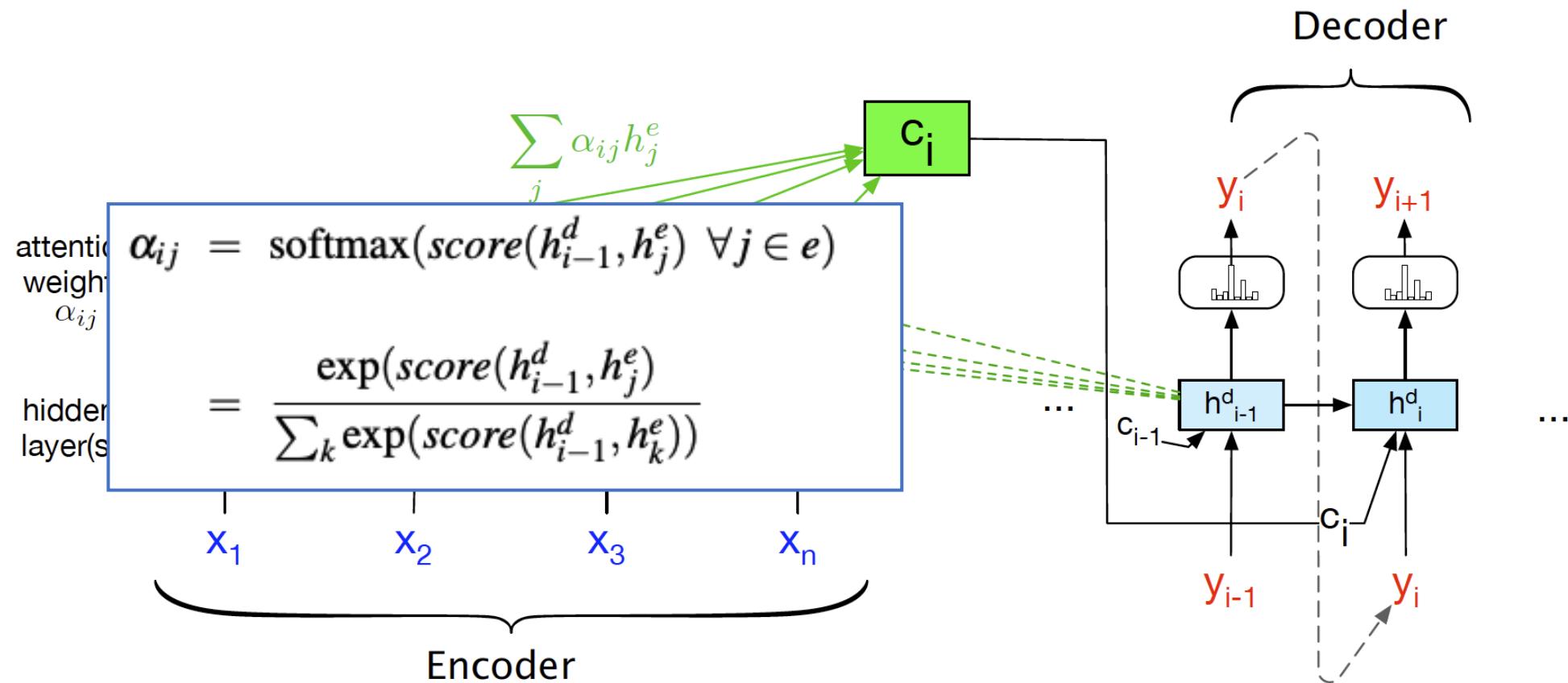
machine translation: bottleneck

- solved by “attention” – to be covered

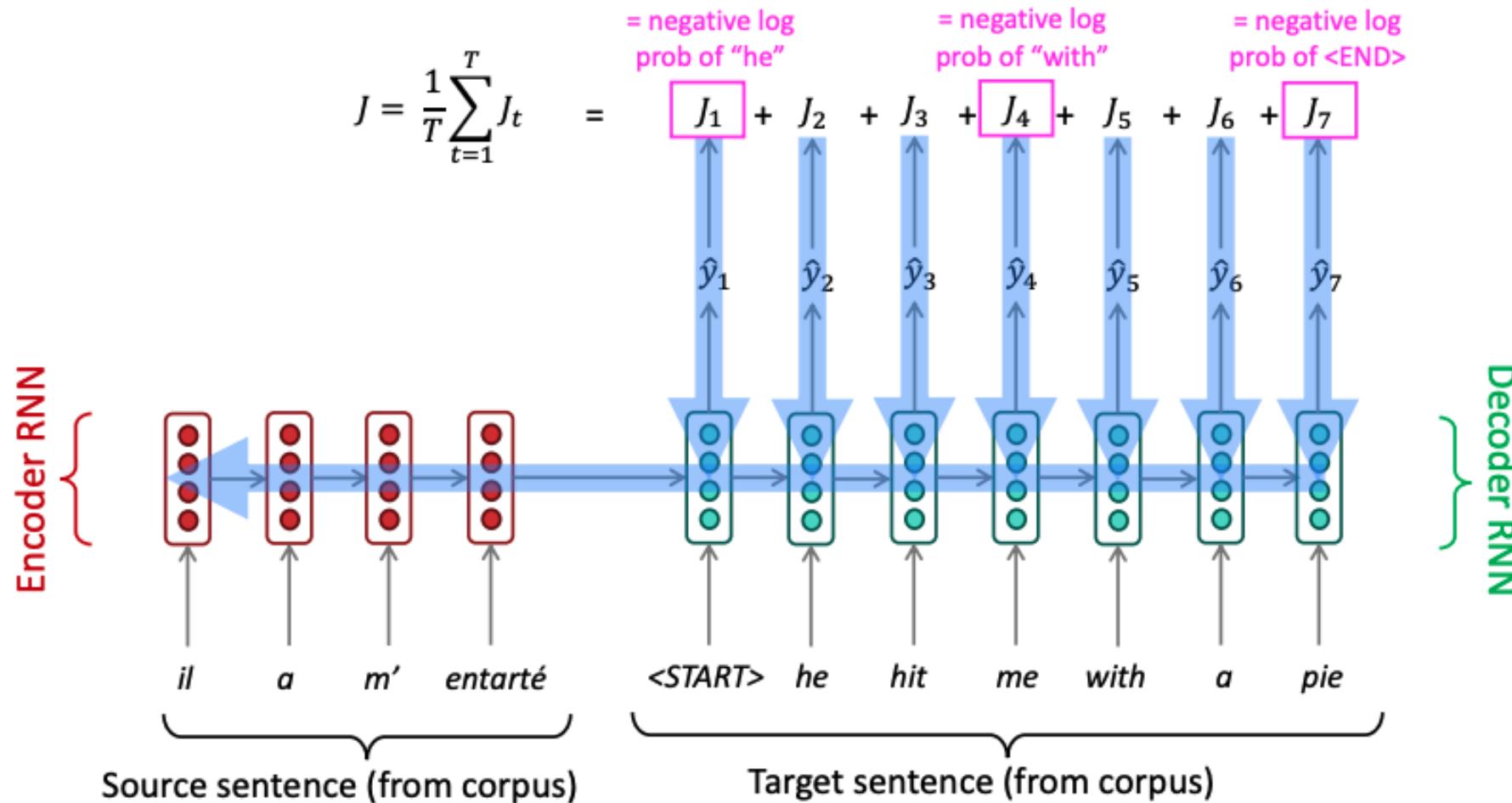


machine translation: bottleneck

- solved by “attention” – to be covered



machine translation: training



machine translation: evaluation

BLEU (bilingual evaluation understudy)

- compare the machine translation with human translation, and compute a similarity score
 - n-gram precision
 - penalty for too-short translation

Source

la verdad, cuya madre es la historia, émula del tiempo, depósito de las acciones, testigo de lo pasado, ejemplo y aviso de lo presente, advertencia de lo por venir.

Reference

truth, whose mother is history, rival of time, storehouse of deeds, witness for the past, example and counsel for the present, and warning for the future.

Candidate 1

truth, whose mother is history, voice of time, deposit of actions, witness for the past, example and warning for the present, and warning for the future

Candidate 2

the truth, which mother is the history, émula of the time, deposition of the shares, witness of the past, example and notice of the present, warning of it for coming

machine translation: bias

The screenshot shows a machine translation interface with Malay input and English output. The Malay input is "Dia bekerja sebagai jururawat." and "Dia bekerja sebagai pengaturcara." The English output is "She works as a nurse." and "He works as a programmer." A pink arrow points from the text "Didn't specify gender" at the bottom to the gendered English outputs.

Malay - detected ▾

English ▾

Dia bekerja sebagai jururawat.

Dia bekerja sebagai pengaturcara. Edit

She works as a nurse.

He works as a programmer.

Didn't specify gender

machine translation: bias

The image displays four separate machine translation results from Google Translate, each showing a different interpretation of the Chinese phrase '信达雅' (Xìndáyǎ), which translates to 'faithfulness, expressiveness, and elegance'.

Top Left: Input '信达雅' (Xìndáyǎ) is translated to 'Sindaya'. The interface shows 'Chinese - detected' and 'English' as the source and target languages respectively. There are microphone and speaker icons below the translation.

Top Right: Input '信达雅' (Xìndáyǎ) is translated to 'Faith, Express, Elegant'. The interface shows 'Chinese - detected' and 'English' as the source and target languages respectively. There are microphone and speaker icons below the translation.

Bottom Left: Input '信达雅' (Xìndáyǎ) is translated to 'Faithfulness, Expressiveness, Elegance'. The interface shows 'Chinese - detected' and 'English' as the source and target languages respectively. There are microphone and speaker icons below the translation.

Bottom Right: A note stating: 'This example was taken in 2021. Google translate has been evolving and you may see a different result.'

machine translation: non-interpretable

Somali	↔	English
Translate from Irish		
ag ag ag ag ag ag ag ag ag ag ag ag ag ag	Edit	As the name of the LORD was written in the Hebrew language, it was written in the language of the Hebrew Nation

6.5 speech recognition and text to speech

Speech recognition and text to speech



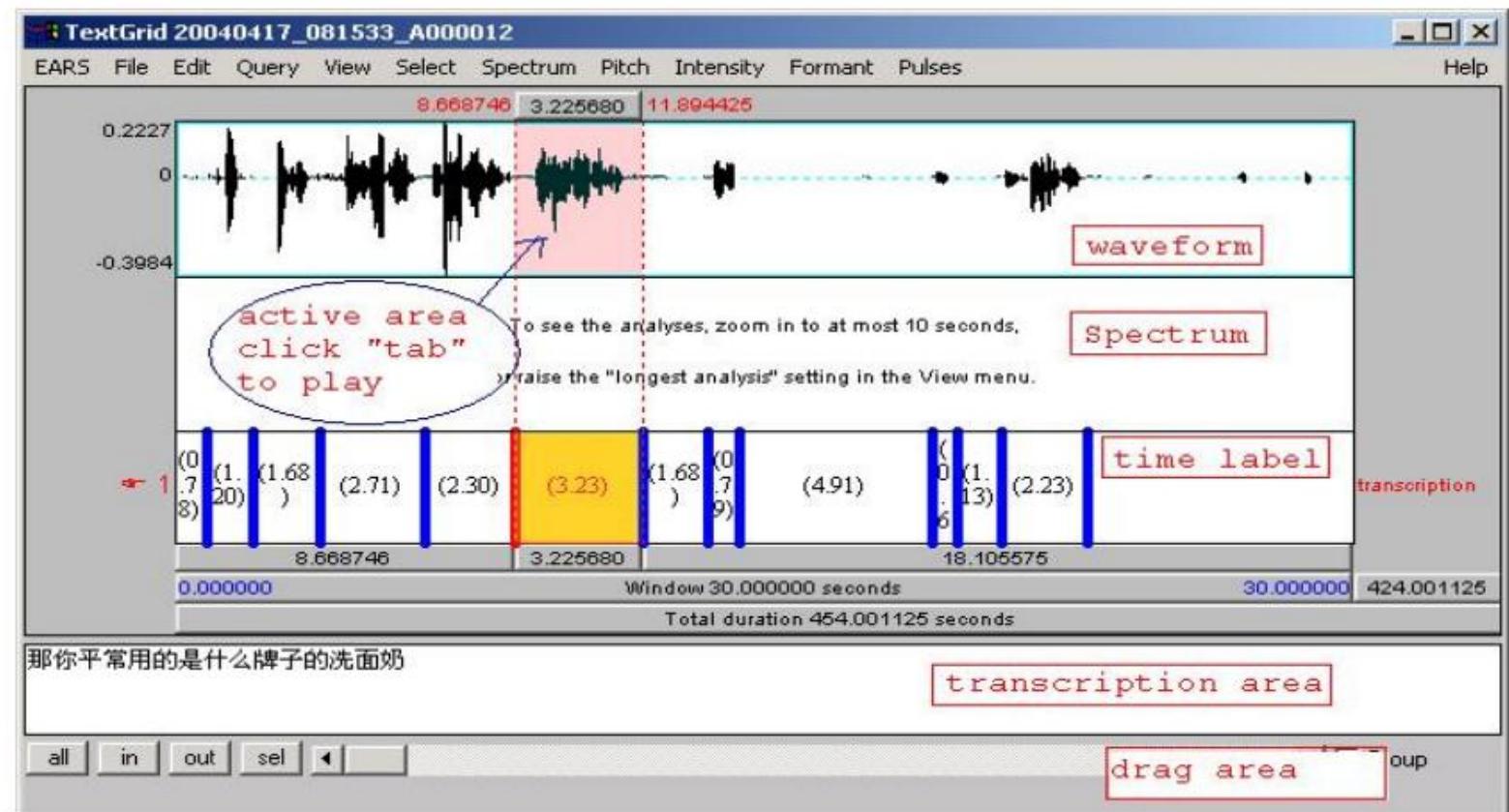
datasets

- **Open Speech and Language Resources** <https://www.openslr.org/resources.php>
 - LibriSpeech: audiobooks
 - Switchboard: telephone conversation
- CORAAL: African American speakers

datasets

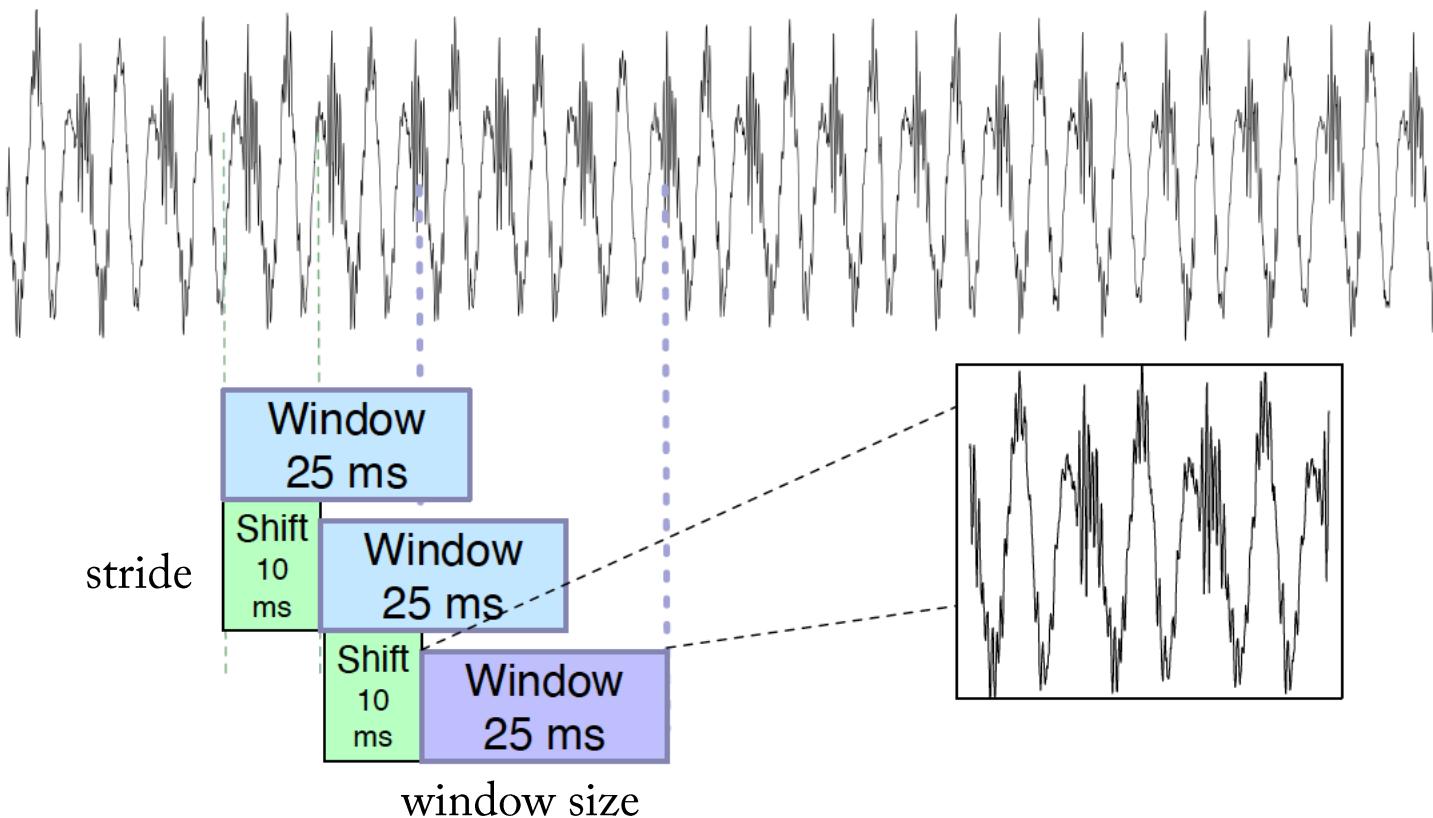
- HKUST Mandarin Telephone Speech

官话 Mandarin	闽 Min	湘 Xiang	赣 Gan	粤 Yue	客家 Kejia	吴 Wu
北东 N.E.	闽南 Minnan	土话 Tuhua	南昌 Nanchang	广州 G.Zh.	梅县 Meixian	太湖 Taihu
冀鲁 Jilu	蒲仙 Puxian	新湘语 N. X	鹰潭 Yingtan	五邑 Wuyi		台州 Taizhou
胶辽 Jiaoliao	闽东 Mindong	老湘语 O. X	抚州 Fuzhou			婺州 Wuzhou
北京 Beijing	闽北 Minbei		宜春 Yichun			处衢 Chuqu
中原 Central	闽中 Minzhong			吉安 Ji'an		瓯江 Oujiang
兰银 Lanyin	琼文 Qiongwen					宣州 Xuanzhou
西南 S.W.	邵将 Shaojiang					
江淮 Jianghuai						
桂柳片 Gui-Liu						

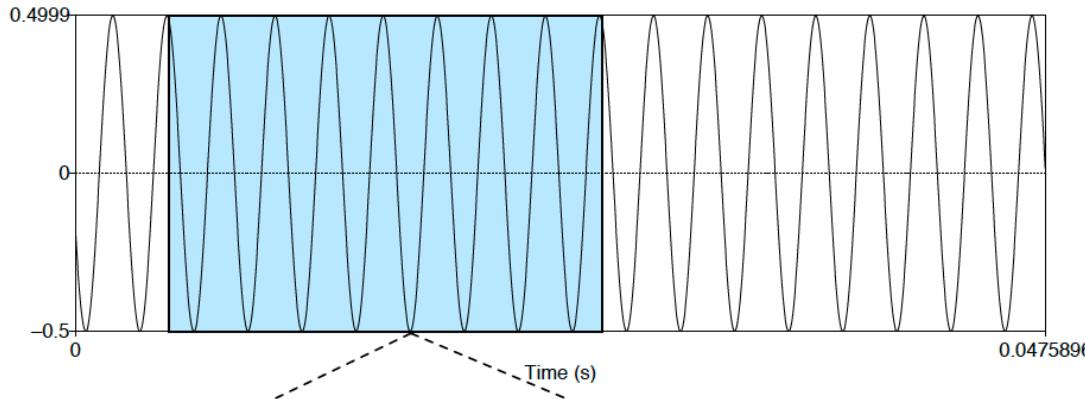


windowing

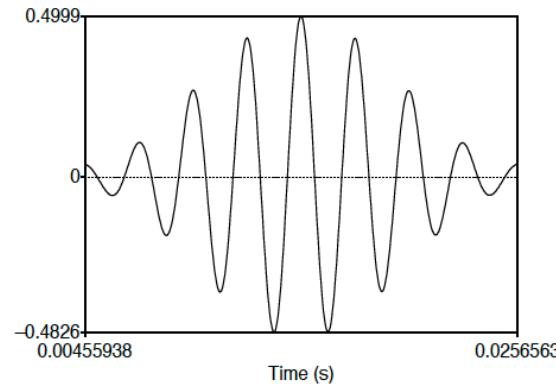
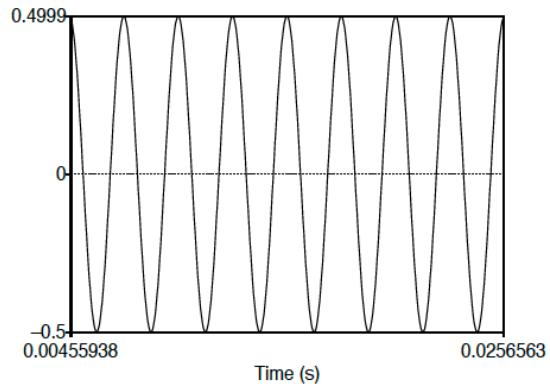
- need to extract features (frames) from small windows



windowing



Rectangular window

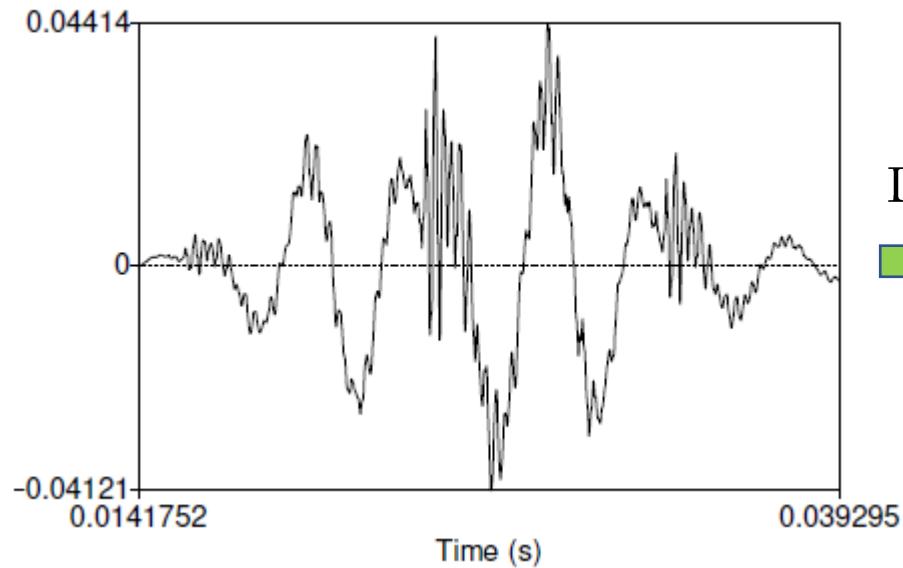


$$y[n] = w[n] s[n]$$

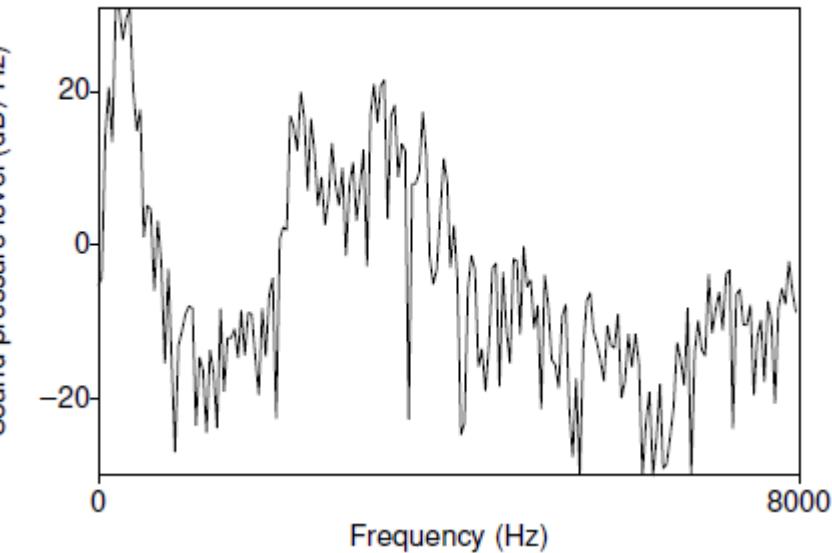
rectangular
Hamming

$$w[n] = \begin{cases} 1 & 0 \leq n \leq L-1 \\ 0 & \text{otherwise} \end{cases}$$
$$w[n] = \begin{cases} 0.54 - 0.46 \cos\left(\frac{2\pi n}{L}\right) & 0 \leq n \leq L-1 \\ 0 & \text{otherwise} \end{cases}$$

spectrum



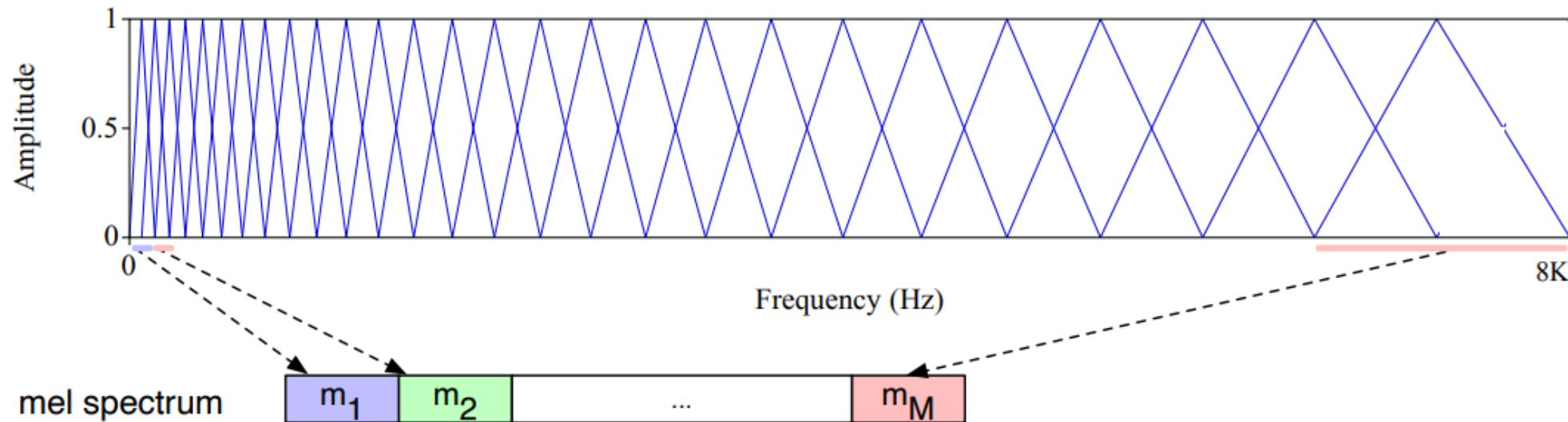
DFT
→



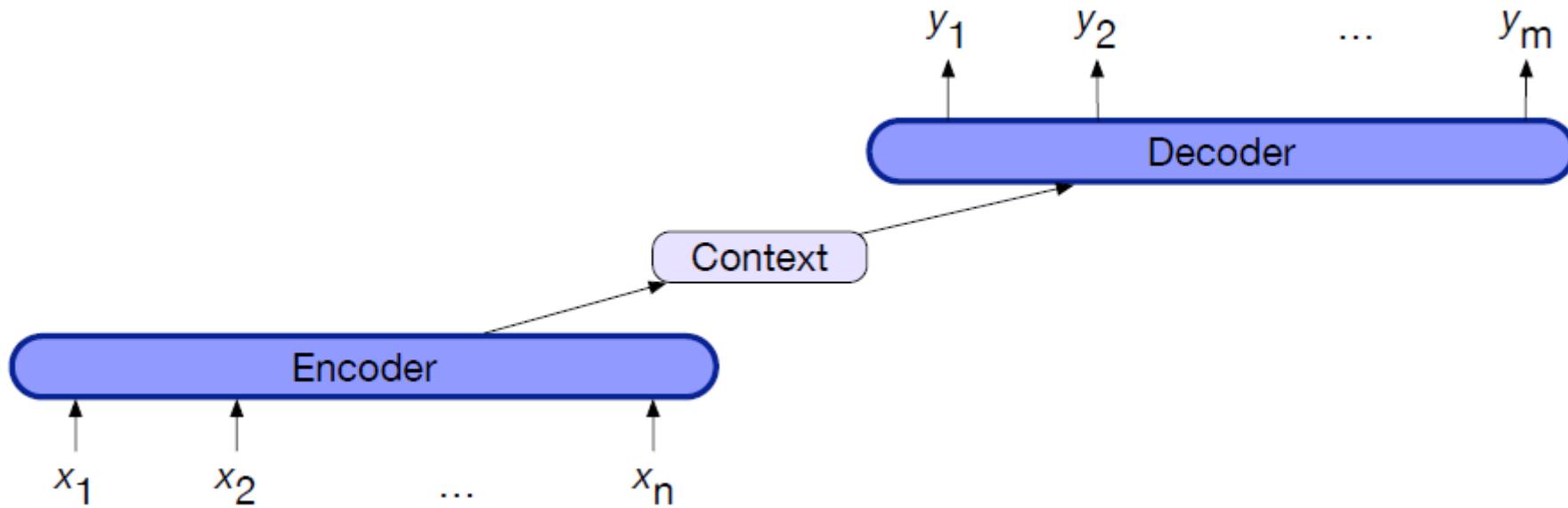
spectrum

- mel scale

$$\text{mel}(f) = 1127 \ln\left(1 + \frac{f}{700}\right)$$



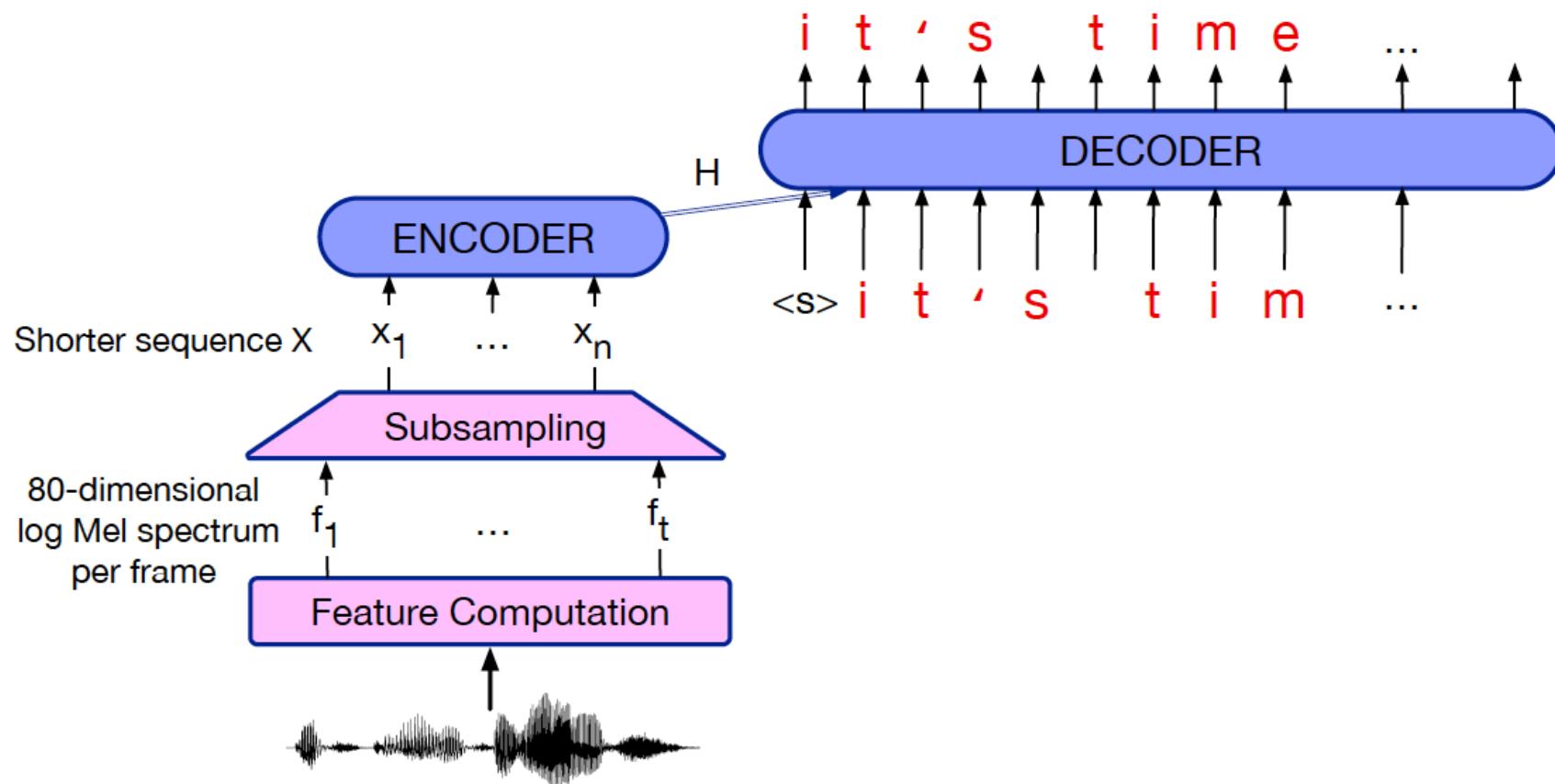
encoder-decoder RNN



encoder-decoder speech recognizer

$$Y = (\langle \text{SOS} \rangle, y_1, \dots, y_m \langle \text{EOS} \rangle)$$

$y_i \in \{a, b, c, \dots, z, 0, \dots, 9, \langle \text{space} \rangle, \langle \text{comma} \rangle, \langle \text{period} \rangle, \langle \text{apostrophe} \rangle, \langle \text{unk} \rangle\}$



training

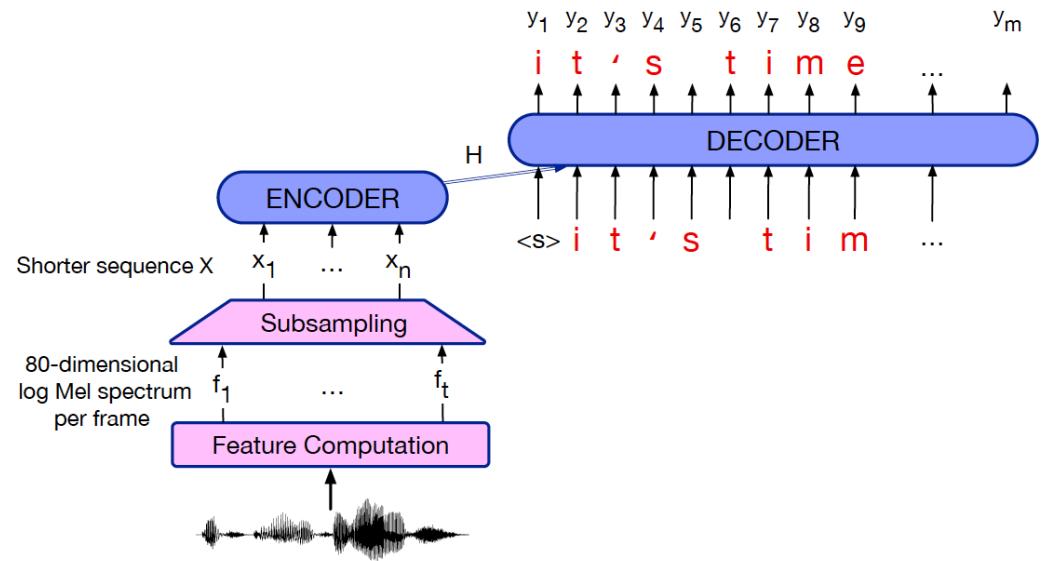
- inference

$$p(y_1, \dots, y_n) = \prod_{i=1}^n p(y_i | y_1, \dots, y_{i-1}, X)$$

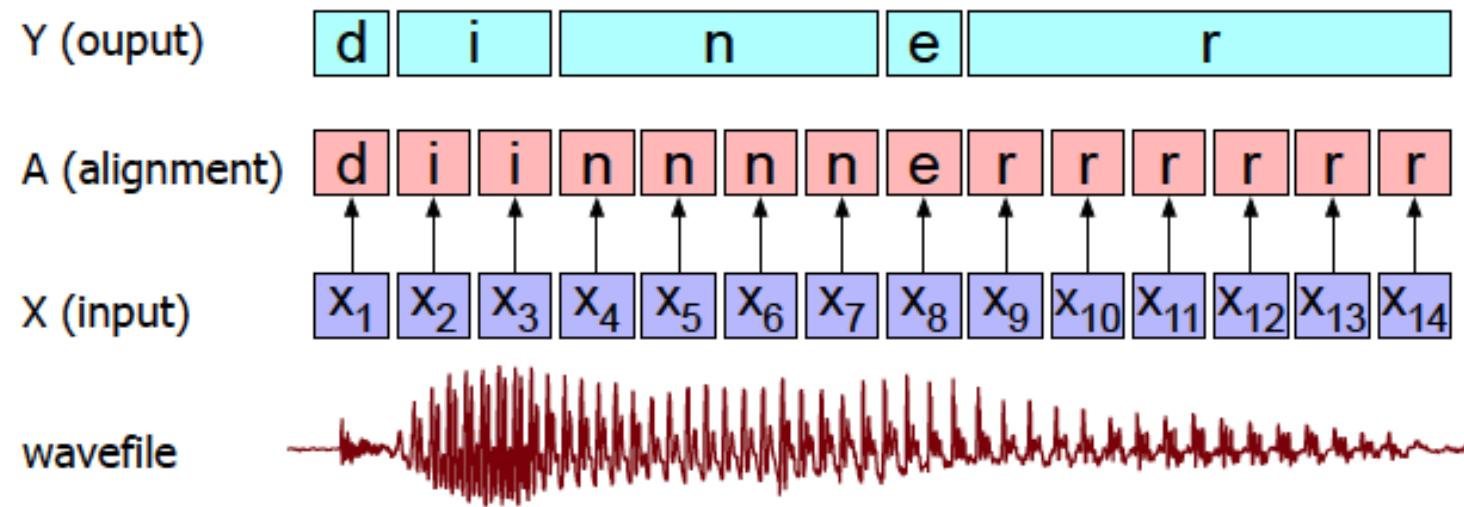
$$\hat{y}_i = \operatorname{argmax}_{\text{char} \in \text{Alphabet}} P(\text{char} | y_1 \dots y_{i-1}, X)$$

- learning: cross-entropy loss
- insufficient data: use a large language model for improved training,
e.g.,

$$score(Y|X) = \frac{1}{|Y|_c} \log P(Y|X) + \lambda \log P_{LM}(Y)$$

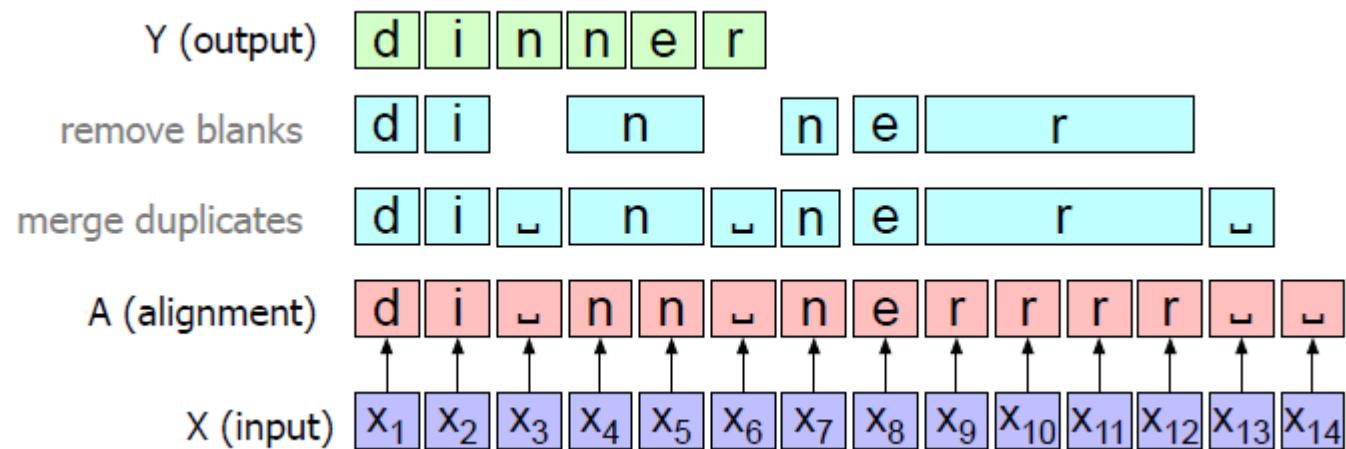


CTC speech recognizer

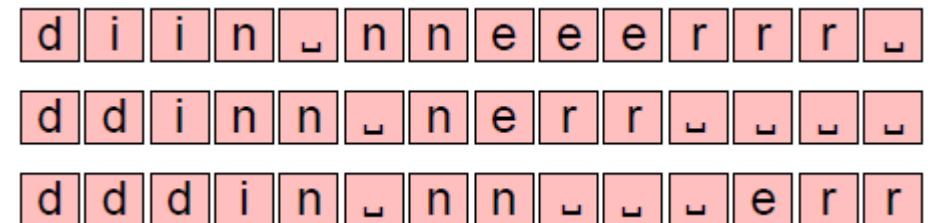


Connectionist Temporal Classification: an alternative to
encoder-decoder

CTC speech recognizer



non-one-to-one-ness



evaluation

$$\text{Word Error Rate} = 100 \times \frac{\text{Insertions + Substitutions + Deletions}}{\text{Total Words in Correct Transcript}}$$

REF:	i *** ** UM the PHONE IS	i LEFT THE portable **** PHONE UPSTAIRS last night
HYP:	i GOT IT TO the ***** FULLEST i LOVE TO portable FORM OF STORES last night	
Eval:	I I S D S S I S S	

evaluation

English Tasks	WER%
LibriSpeech audiobooks 960hour clean	1.4
LibriSpeech audiobooks 960hour other	2.6
Switchboard telephone conversations between strangers	5.8
CALLHOME telephone conversations between family	11.0
Sociolinguistic interviews, CORAAL (AAL)	27.0
CHiMe5 dinner parties with body-worn microphones	47.9
CHiMe5 dinner parties with distant microphones	81.3
Chinese (Mandarin) Tasks	CER%
AISHELL-1 Mandarin read speech corpus	6.7
HKUST Mandarin Chinese telephone conversations	23.5

TTS text-to-speech

- spectrum prediction

It's time for lunch!



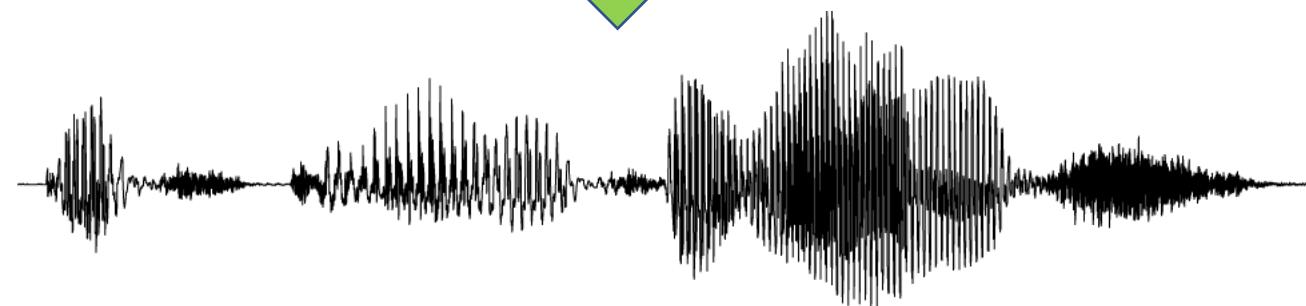
TTS text-to-speech

It's time for lunch!

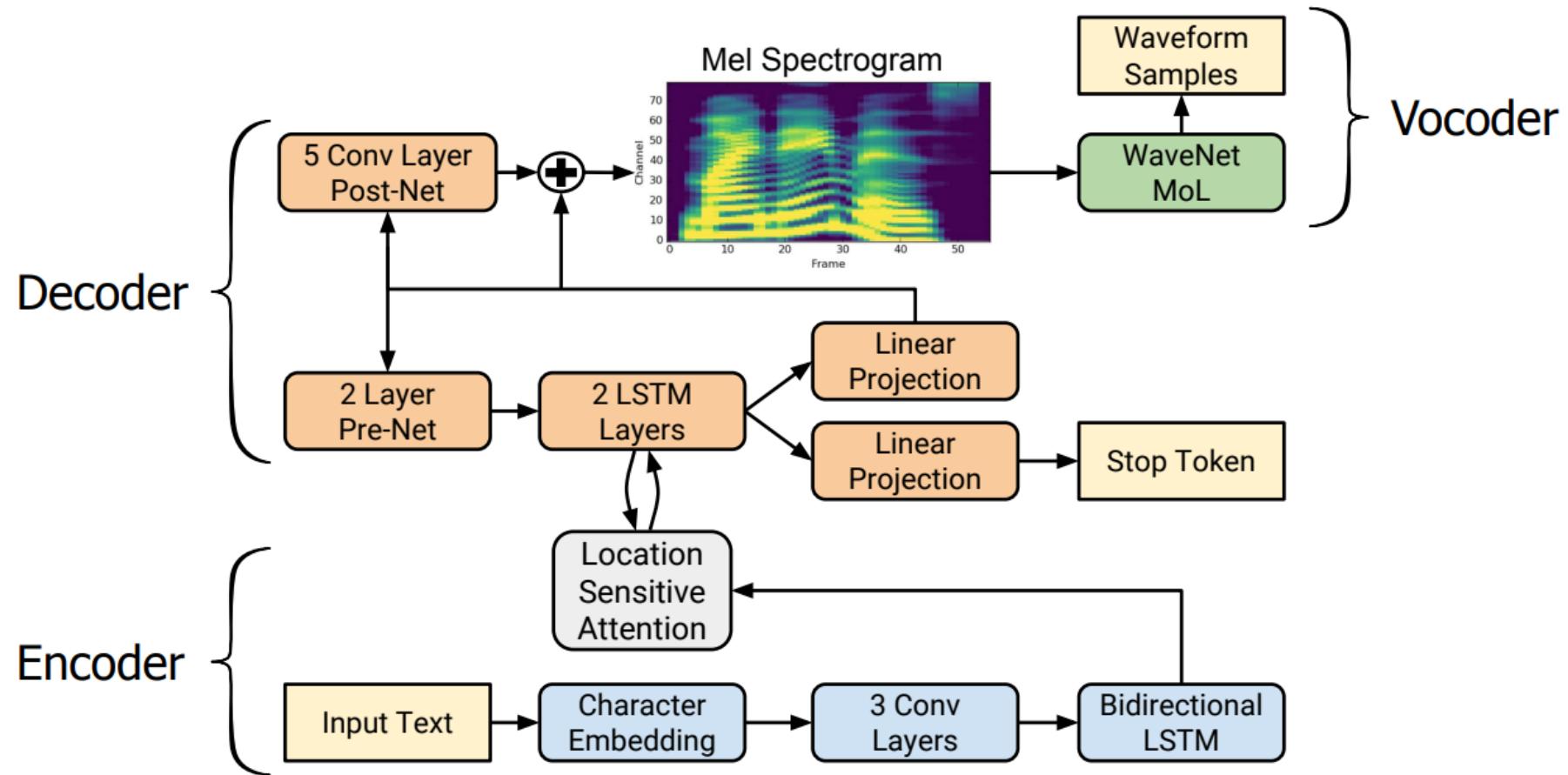
encoder-decoder



vocoder



Tacotron2



visit https://pytorch.org/hub/nvidia_deeplearningexamples_tacotron2/
or play [here](#)

Thank you!

Reference

- Ch 10, *Deep Learning*.
- Pascanu, R., Mikolov, T., & Bengio, Y. (2013, May). On the difficulty of training recurrent neural networks. In ICML (pp. 1310-1318).
- Ch 8, Jurafsky, D., & Martin, J. H. (2018). *Speech and language processing (draft)*. <https://web.stanford.edu/~jurafsky/slp3>.
- Goldberg, Y. (2016). A primer on neural network models for natural language processing. *Journal of Artificial Intelligence Research*, 57, 345-420.



Thank you!

Reference

- Ch 11 & 26, Jurafsky, D., & Martin, J. H. (2018). *Speech and language processing (draft)*. <https://web.stanford.edu/~jurafsky/slp3>.
- Goldberg, Y. (2016). A primer on neural network models for natural language processing. *Journal of Artificial Intelligence Research*, 57, 345-420.

