# Language Models – Basics - 2

``It appears then that a sufficiently complex stochastic process will give a satisfactory representation of a discrete source.''

``A second method is to delete a certain fraction of the letters from a sample of English text and then let someone attempt to restore them. If they can be restored when 50% are deleted the redundancy must be greater than 50%.''
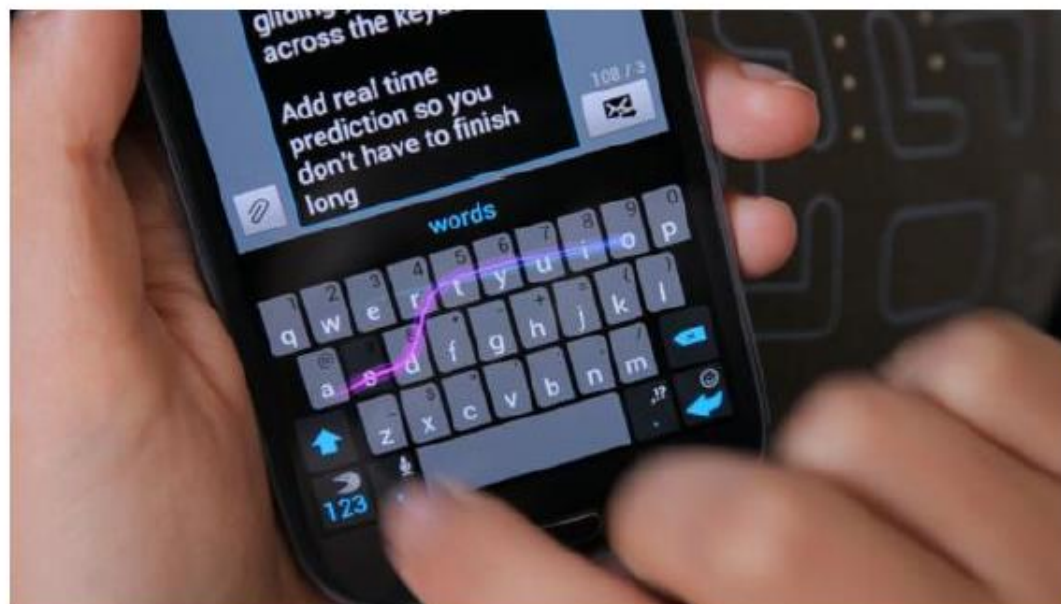
# LMs: "fill in the blank"
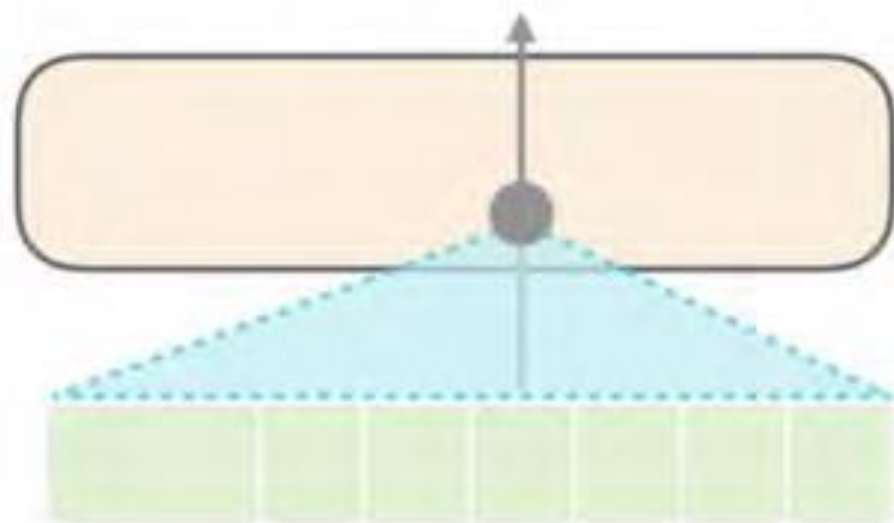
- Think of this as a "fill in the blank" problem.

$$P(w_n | w_1, w_2, w_3, \ldots, w_{n-1})$$
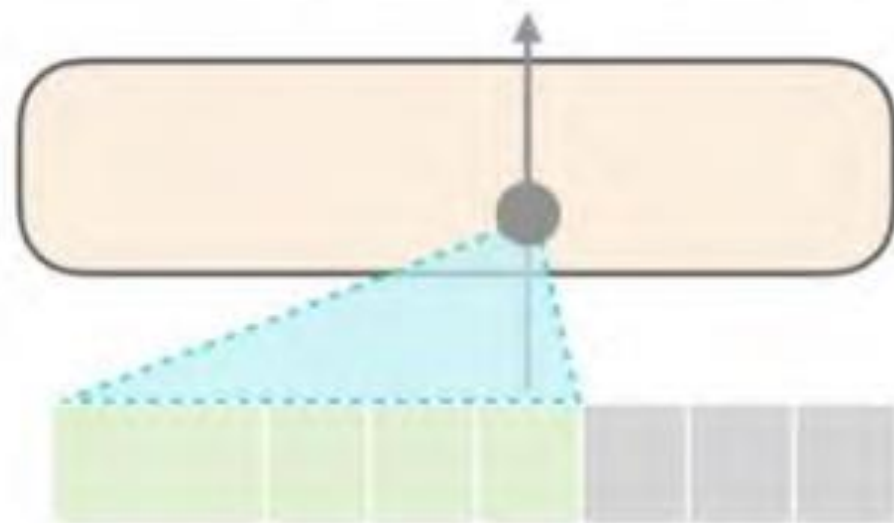
"He picked up the bat and hit the _____"
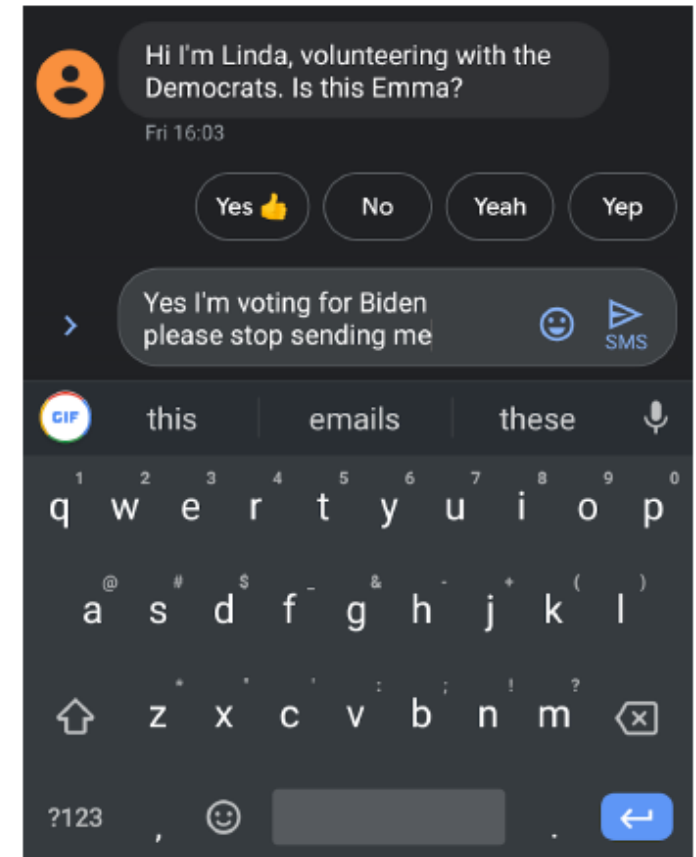
**Ball? Poetry?**

**Self-Attention**

**Masked Self-Attention**

# Probabilistic language models

- Today's goal: assign a probability to a sentence. Why?

# Probabilistic language models

- Today's goal: assign a probability to a sentence. Why?

    - Machine translation:
      P(*high winds tonight*) > P(*large winds tonight*)

    - Spelling correction:
      P(*I'll be five minutes late*) > P(*I'll be five minuets late*)

    - Speech recognition:
      P(*I saw a van*) > P(*eyes awe of an*)

    - Summarization, question answering, …

# Probabilistic language models

- Goal: compute the probability of a sentence (or sequence of words):

$$P(\mathbf{w}) = P(w_1, w_2, w_3, ..., w_n)$$

- Related task: probability of the next word:

$$P(w_5 \mid w_4, w_3, w_2, w_1)$$

- A model that computes either of these is called a **language model** (or **LM**).

# Statistical Language Models (LM)

To calculate $P(\text{word } w \mid \text{history } h)$

or $P(W)$, $W = w_1, w_2, \ldots, w_n$ <u>a 'n-word'</u> <u>sequence</u>

- $P(\text{Random variable } X_i \text{ taking the value "the"})$

$$P(X_i = \text{"the"}) \Rightarrow P(\text{"the"}) \text{ or } P(\text{the})$$

- Sequence of 'n' words $W = w_1, w_2, \ldots, w_n$ or

$$W_{1:n}$$

$$w_{1:1} \Rightarrow w_1$$
$$w_{1:2} \Rightarrow w_1, w_2$$
$$\vdots$$
$$w_{1:n-1} \Rightarrow w_1, w_2, \ldots, w_{n-1}$$
$$w_{1:n} \Rightarrow w_1, w_2, \ldots, w_{n-1}, w_n$$

Joint Probability of each word in a sequence having a particular Value

$$P(X_1 = \omega_1, X_2 = \omega_2 - - \cdot, X_n = \omega_n\}$$

$$\Rightarrow P(\omega_1, \omega_2, - - -, \omega_{n-1}, \omega_n)$$

What is $P(\omega_1, \omega_2, - - -, \omega_n)$? [Prob of Sentence $\overline{W}$]

$$P(X_1 - - - X_n) = P(X_1) P(X_2/X_1) P(X_3 | X_1 X_2) - - - P(X_n | X_{1:n-1})$$

$$= \prod_{k=1}^{n} P(X_k | X_{1:k-1})$$

From

$P(x, y)$

$= P(x/y) \cdot P(y)$

How

By CHAIN RULE OF PROBABILITY

$$P(A, B, C, D) = \underbrace{P(A) \cdot P(B/A)}_{P(A, B)} \cdot \underbrace{P(C/A, B) \cdot P(D/A, B, C)}_{P(A, B, C)}$$

# Applying Chain-Rule to Words

$$P(W) = P(w_{1:n}) = P(w_1) P(w_2/w_1) P(w_3/w_{1:2}) \cdots P(w_n/w_{1:n-1})$$

$$= \prod_{k=1}^{n} P(w_k | w_{1:k-1})$$

**Problem :** How to get estimates (e.g. MLE) of these terms

i.e. Compute the "exact" probability of a word $w_n$ given a long sequence of preceding words

$$w_{1:n-1}$$

i.e. $P(w_n | w_{1:n-1})$

or in general $P(w_k | w_{1:k-1})$

Any particular "Context" (or history "h") might have never occured before. (in the trg corpus)

Instead : " N- gram model "    ④

Approximate the history $w_{1:k-1}$ by just the last few words

e.g Bigram Model ⟹ last `1` word.

$$P\left(w_n \mid w_{1:n-1}\right) \approx P\left(w_n \mid w_{n-1}\right)$$

h: history

Approximate history //

| $w_1$ | $w_2$ | $w_3$ | ... | $w_{n-1}$ | $w_n$ |

e.g (Walden Pond's water is so transparent (that) — the $w_n$

⟹ P(the | history) = P(the | that)

⟹ MARKOV ASSUMPTION of order-1.

Likewise Markov Assumption of Order $N-1$
yields $N$-gram model or approximation

$$P(w_n \mid w_{1:n-1}) = P(w_n \mid w_{n-N+1 : n-1})$$



$w_n$

| 1 | 2 | 3 |
|---|---|---|
| $w_1$ | $w_2$ | $w_3$ |

$n-N+1 \mid - \; - \; - \; \cdot \mid n-2 \mid n-1 \mid n \mid n+1 \mid \cdots \mid - - -$

$N-1$ Context

Note # words in "truncated" Context

$$= n-1 - (n - N+1) + 1$$

$$= n-1 - n + N - n + n$$

$$= N-1 \text{ words} \rightsquigarrow \begin{array}{l}\text{Markov Assumption of order } N-1 \\ \Rightarrow N\text{-gram probabilities}\end{array}$$

| | $P(W) = P(\omega_{1:n}) =$ | MLE Maximum Likelihood Estimate | |
|---|---|---|---|
| Unigram $N=1$ | $\displaystyle\prod_{R=1}^{n} P(\omega_R)$ | $P(\omega_n) = \dfrac{C(\omega_n)}{\displaystyle\sum_{x} C(x)}$ | Total Count of words in the Corpus / Document |
| Bigram. $N=2$ | $\displaystyle\prod_{R=1}^{n} P(\omega_R \mid \omega_{R-1})$ | $P(\omega_n \mid \omega_{n-1}) = \dfrac{C(\omega_{n-1}, \omega_n)}{C(\omega_{n-1})}$ | $\displaystyle\sum_{\omega} C(\omega_{n-1} \omega)$ Marginal Distn. from Joint distn: Integrate out the unwanted variable. |
| Trigram $N=3.$ | $\displaystyle\prod_{R=1}^{n} P(\omega_R \mid \omega_{R-1}, \omega_{R-2})$ $\vdots$ | | |
| N-gram any N! | $\displaystyle\prod_{R=1}^{n} P(\omega_R \mid \omega_{R-N+1:R-1})$ | $P(\omega_n \mid \omega_{n-N+1:n-1})$ $= \dfrac{C(\omega_{n-N+1:n-1}, \omega_n)}{C(\omega_{n-N+1:n-1})}$ | Relative Frequency of Occurrences |

Prefix.

Target word

Wn.

sequence

$$\frac{\text{observed frequency of "sequence"}}{\text{observed frequency of "prefix"}}$$

1) **Problem:** How to deal with "UNKNOWN Words"

⇒ OOVs or Out-of-Vocabulary words.

— Smoothing
— Backoff
— Interpolation
} ⟶

# Backoff :

Use trigram
— If not available — Backoff to (use) Bigram
— If not available — Backoff to Unigram

## Interpolation

$\Rightarrow$ Mix the Probability Estimates

Trigram

$$\hat{P}(\omega_n | \omega_{n-1} \, \omega_{n-2}) = \lambda_1 \, P(\omega_n)$$
$$+ \lambda_2 \, P(\omega_n | \omega_{n-1})$$
$$+ \lambda_3 \, P(\omega_n | \omega_{n-1} \, \omega_{n-1})$$

$$\sum_i \lambda_i = 1$$

# Neural LMs

LMs: Probability Distribution over sequences of $n$ tokens.

$$\{ t_1 \ t_2 \ - \ - \ - \ t_n \}$$

Given such a sequence, an LM assigns a probability

$$P(t_1, t_2 - - - . . t_n)$$

to the whole sequence by modeling the

Prob( token $t_R$ | history $t_1 \ t_2 \ - \ - \ . \ t_{R-1}$)

i.e, $P(t_1 \ t_2 \ - \ - . \ t_n) = \prod_{R=1}^{n} P(t_R \ | \ t_1 \ t_2 \ - - - t_{R-1})$

$$P(t_1, t_2 \dots t_n) = \prod_{k=1}^{n} P(t_k \mid t_1 t_2 \dots t_{k-1})$$

→ Trained by Minimizing the Negative Log-Likelihood,

$$\sum_{k=1}^{n} -\log\left(P\left(t_1, t_2 \dots t_{Rj} \; \theta_t, \theta_{rnn}, \theta_s\right)\right)$$

Parameters to be optimized
$$\theta_t \; , \quad \theta_{rnn} \; , \quad \theta_s$$

$$\Rightarrow$$

$\theta_t$ : Look-up table / Word-to-Vec embedding layer

→ maps each token into a vector of fixed dimension.

→ Word2vec or TBS result on a given vocab of words.
$\hookrightarrow |V|$

$\theta_{rnn}$ : RECURRENT NEURAL NETWORK ( RNN, LSTM, GRU...)

→ Summarizes the sequence of history

e.g. $t_1, t_2 - - - . t_{k-1}$

up to the current time-step (e.g $t_k$)

$\theta_s$: Softmax layer appended at the o/p of EACH

RNN time step for estimating the probability

distribution over the tokens (Posterior vector of dim $|V|$.)

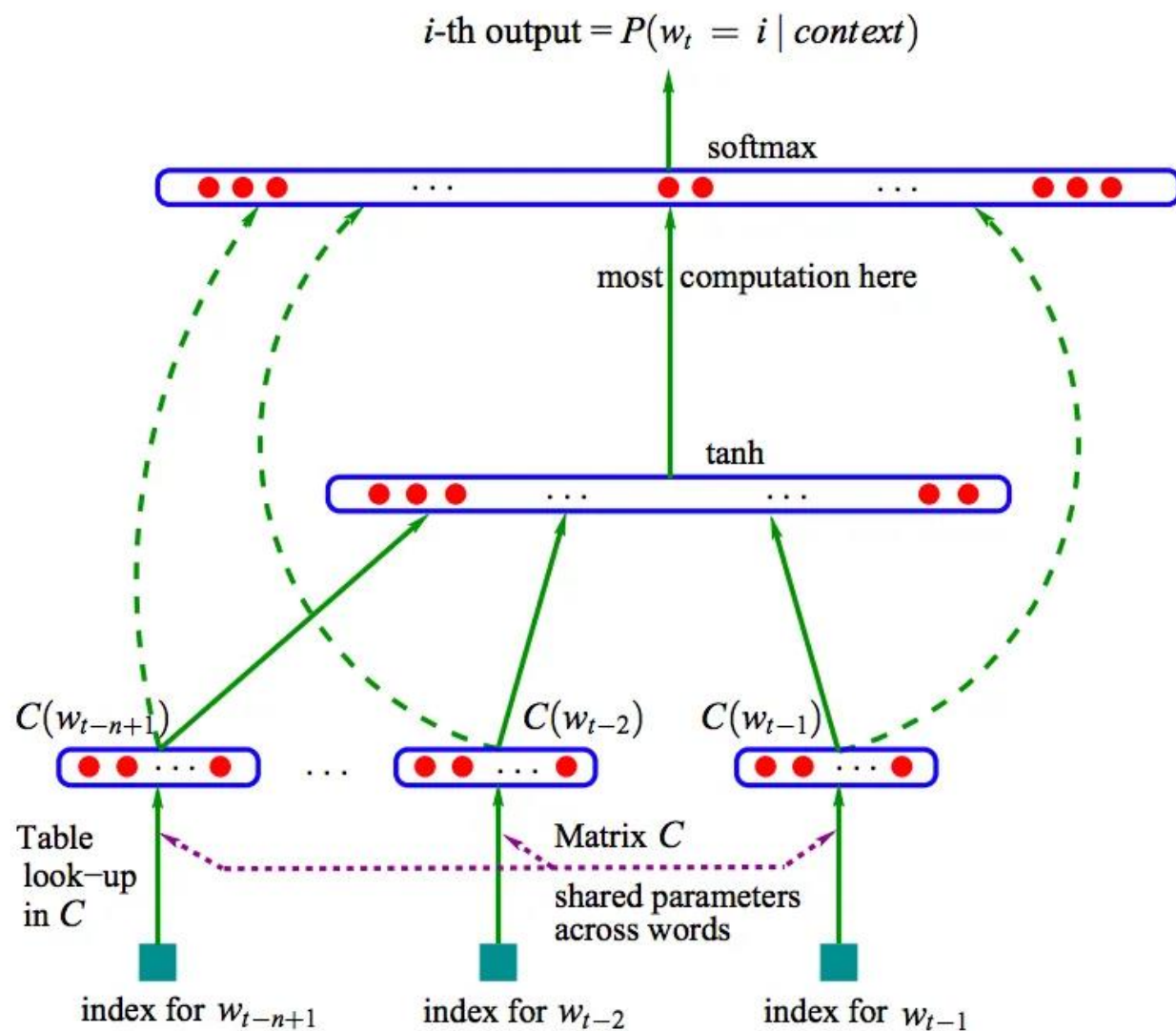@ $k-1$ ⟹ max prob of $t_k$ (ground truth at $k$)

Figure 1: Neural architecture: $f(i, w_{t-1}, \cdots, w_{t-n+1}) = g(i, C(w_{t-1}), \cdots, C(w_{t-n+1}))$ where $g$ is the neural network and $C(i)$ is the $i$-th word feature vector.

# A fixed-window neural Language Model

output distribution

$$\hat{y} = \mathrm{softmax}(\boldsymbol{U}\boldsymbol{h} + \boldsymbol{b}_2) \in \mathbb{R}^{|V|}$$

hidden layer

$$\boldsymbol{h} = f(\boldsymbol{W}\boldsymbol{e} + \boldsymbol{b}_1)$$

concatenated word embeddings

$$\boldsymbol{e} = [\boldsymbol{e}^{(1)}; \boldsymbol{e}^{(2)}; \boldsymbol{e}^{(3)}; \boldsymbol{e}^{(4)}]$$

words / one-hot vectors

$$\boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \boldsymbol{x}^{(3)}, \boldsymbol{x}^{(4)}$$

books

laptops

$\boldsymbol{U}$

$\boldsymbol{W}$

the
$\boldsymbol{x}^{(1)}$

students
$\boldsymbol{x}^{(2)}$

opened
$\boldsymbol{x}^{(3)}$

their
$\boldsymbol{x}^{(4)}$

outputs (optional) $\hat{y}^{(1)}$  $\hat{y}^{(2)}$  $\hat{y}^{(3)}$  $\hat{y}^{(4)}$  ...

hidden states $h^{(1)}$  $W$  $h^{(2)}$  $W$  $h^{(3)}$  $W$  $h^{(4)}$  $W$  ...

input sequence (any length) $x^{(1)}$  $x^{(2)}$  $x^{(3)}$  $x^{(4)}$  ...

# A RNN Language Model

**output distribution**

$$\hat{y}^{(t)} = \text{softmax}\left(Uh^{(t)} + b_2\right) \in \mathbb{R}^{|V|}$$

**hidden states**

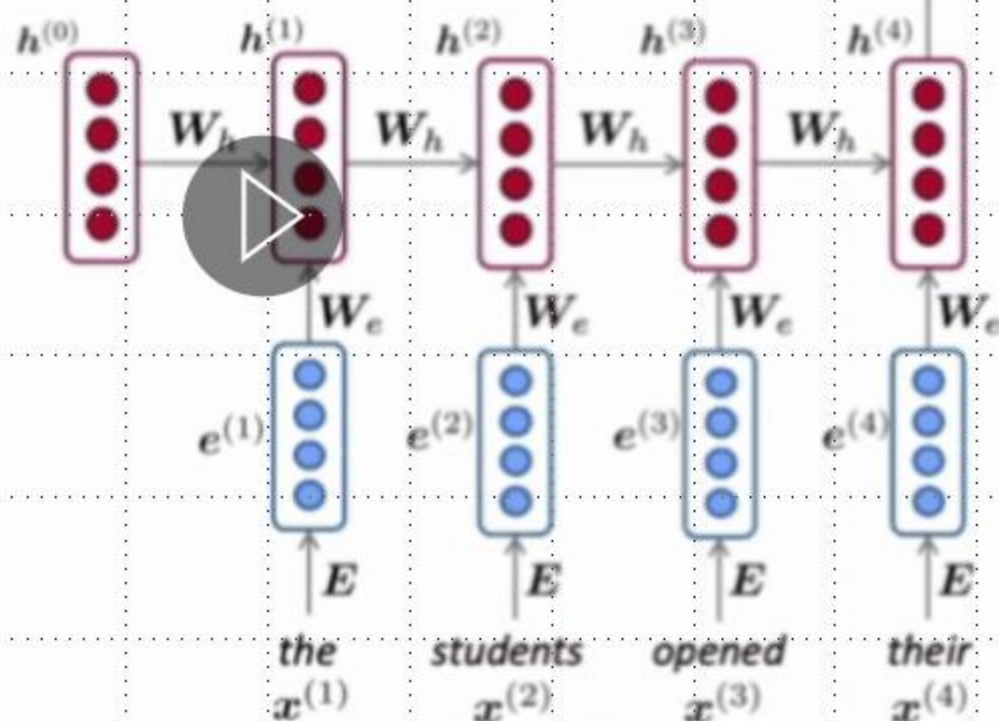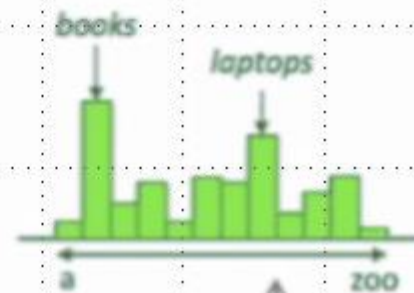$$h^{(t)} = \sigma\left(W_h h^{(t-1)} + W_e e^{(t)} + b_1\right)$$

$h^{(0)}$ is the initial hidden state

**word embeddings**

$$e^{(t)} = E x^{(t)}$$

**words / one-hot vectors**

$$x^{(t)} \in \mathbb{R}^{|V|}$$

23

**Note**: *this input sequence could be much longer, but this slide doesn't have space!*

books    laptops

a    zoo

$U$

$h^{(0)}$    $h^{(1)}$    $h^{(2)}$    $h^{(3)}$    $h^{(4)}$

$W_h$    $W_h$    $W_h$    $W_h$

$W_e$    $W_e$    $W_e$    $W_e$

$e^{(1)}$    $e^{(2)}$    $e^{(3)}$    $e^{(4)}$

$E$    $E$    $E$    $E$

the    students    opened    their

$x^{(1)}$    $x^{(2)}$    $x^{(3)}$    $x^{(4)}$
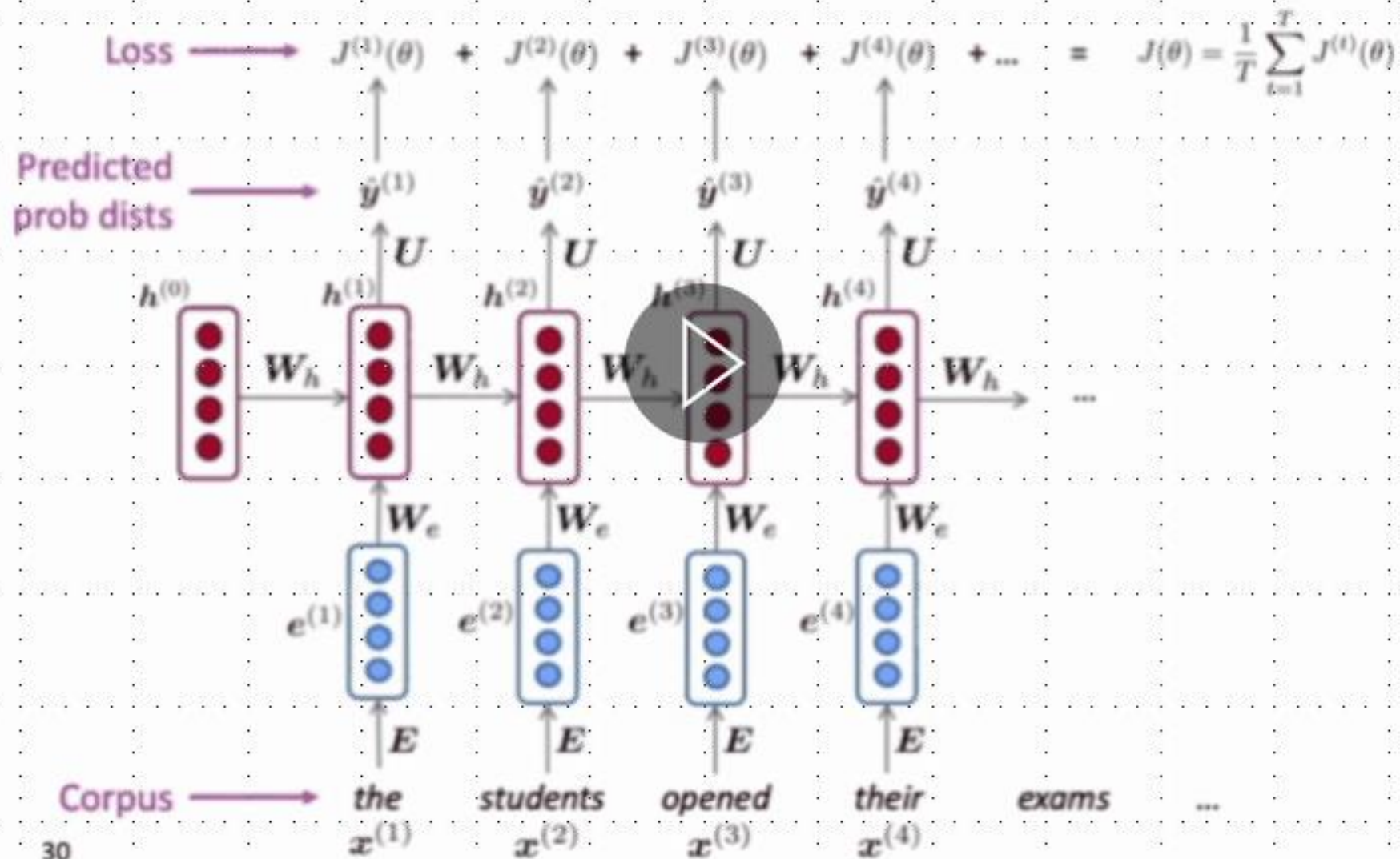
Loss function on step $t$ is cross-entropy between predicted probability distribution $\hat{y}^{(t)}$, and the true next word $y^{(t)}$ (one-hot for $x^{(t+1)}$):

$$J^{(t)}(\theta) = CE(\boldsymbol{y}^{(t)}, \hat{\boldsymbol{y}}^{(t)}) = -\sum_{w \in V} \boldsymbol{y}_w^{(t)} \log \hat{\boldsymbol{y}}_w^{(t)} = -\log \hat{\boldsymbol{y}}_{x_{t+1}}^{(t)}$$

Average this to get overall loss for entire training set:

$$J(\theta) = \frac{1}{T}\sum_{t=1}^{T} J^{(t)}(\theta) = \frac{1}{T}\sum_{t=1}^{T} -\log \hat{\boldsymbol{y}}_{x_{t+1}}^{(t)}$$

# Training a RNN Language Model

Loss $\longrightarrow$ $\quad J^{(1)}(\theta) \quad + \quad J^{(2)}(\theta) \quad + \quad J^{(3)}(\theta) \quad + \quad J^{(4)}(\theta) \quad + \dots \quad = \quad J(\theta) = \frac{1}{T}\sum_{t=1}^{T} J^{(t)}(\theta)$

Predicted $\longrightarrow$ $\quad \hat{y}^{(1)} \qquad \hat{y}^{(2)} \qquad \hat{y}^{(3)} \qquad \hat{y}^{(4)}$
prob dists



Corpus $\longrightarrow$ the $\quad$ students $\quad$ opened $\quad$ their $\quad$ exams $\quad \dots$

$x^{(1)} \qquad x^{(2)} \qquad x^{(3)} \qquad x^{(4)}$

30

Thank you !!