# Automatic Differentiation
# &
# Transformers

Matt Gormley
Lecture 28
Nov. 1, 2023

# Reminders

- **Homework 6: Learning Theory & Generative Models**
  - **Out: Fri, Oct 27**
  - **Due: Fri, Nov 3 at 11:59pm**
- **Exam 2 Practice Problems**
  - **Out: Fri, Nov 3**
- **Exam 2: Thu, Nov 9**

# MODULE-BASED AUTOMATIC DIFFERENTIATION

# Backpropagation

## Automatic Differentiation – Reverse Mode (aka. Backpropagation)

<u>Forward Computation</u>
1. Write an **algorithm** for evaluating the function y = f(**x**). The algorithm defines a **directed acyclic graph,** where each variable is a node (i.e. the "**computation graph**")
2. Visit each node in **topological order.**
   For variable $u_i$ with inputs $v_1, \ldots, v_N$
   a. Compute $u_i = g_i(v_1, \ldots, v_N)$
   b. Store the result at the node

<u>Backward Computation (Version A)</u>
1. **Initialize** dy/dy = 1.
2. Visit each node $v_j$ in **reverse topological order.**
   Let $u_1, \ldots, u_M$ denote all the nodes with $v_j$ as an input
   Assuming that y = h(**u**) = h($u_1, \ldots, u_M$)
   and **u** = g(**v**) or equivalently $u_i = g_i(v_1, \ldots, v_j, \ldots, v_N)$ for all i
   a. We already know dy/du$_i$ for all i
   b. Compute dy/dv$_j$ as below (Choice of algorithm ensures computing (du$_i$/dv$_j$) is easy)

$$\frac{dy}{dv_j} = \sum_{i=1}^{M} \frac{dy}{du_i}\frac{du_i}{dv_j}$$

**Return** partial derivatives dy/du$_i$ for all variables

# Backpropagation

## Automatic Differentiation – Reverse Mode (aka. Backpropagation)

Forward Computation
1. Write an **algorithm** for evaluating the function y = f(**x**). The algorithm defines a **directed acyclic graph,** where each variable is a node (i.e. the "**computation graph**")
2. Visit each node in **topological order.**
   For variable $u_i$ with inputs $v_1,\ldots, v_N$
   a. Compute $u_i = g_i(v_1,\ldots, v_N)$
   b. Store the result at the node

Backward Computation (Version B)
1. **Initialize** all partial derivatives $dy/du_j$ to 0 and $dy/dy = 1$.
2. Visit each node in **reverse topological order.**
   For variable $u_i = g_i(v_1,\ldots, v_N)$
   a. We already know $dy/du_i$
   b. Increment $dy/dv_j$ by $(dy/du_i)(du_i/dv_j)$
      (Choice of algorithm ensures computing $(du_i/dv_j)$ is easy)

**Return** partial derivatives $dy/du_i$ for all variables

# Backpropagation

*Why is the backpropagation algorithm efficient?*

1. Reuses **computation from the forward pass** in the backward pass

2. Reuses **partial derivatives** throughout the backward pass (*but only if the algorithm reuses shared computation in the forward pass*)

    (Key idea: partial derivatives in the backward pass should be thought of as variables stored for reuse)

A Recipe for

## Gradients

1. Given training dat

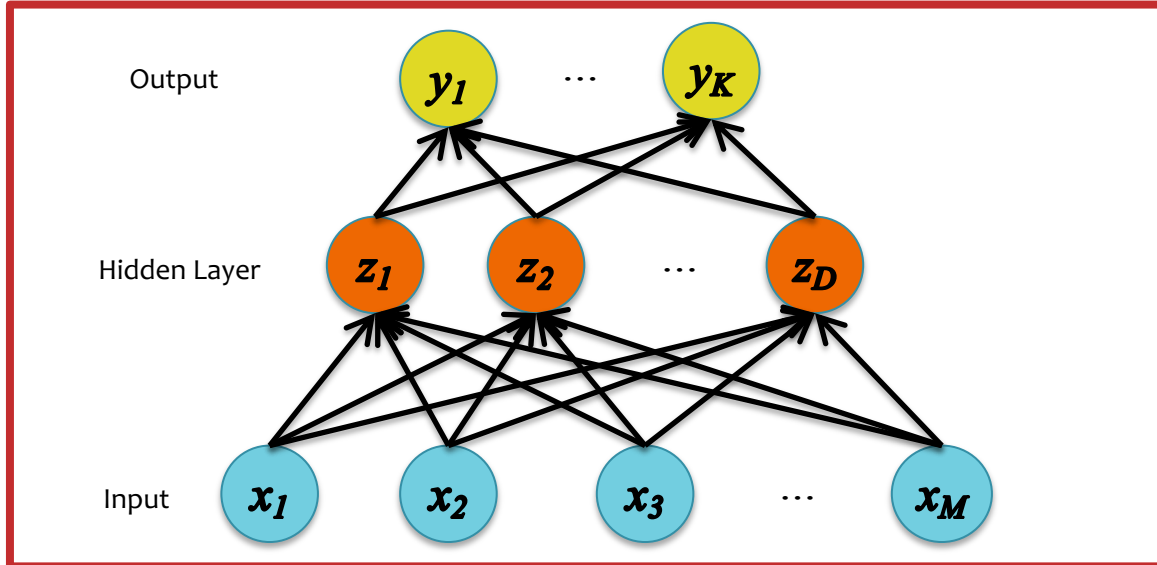$$\{\boldsymbol{x}_i, \boldsymbol{y}_i\}_{i=1}^{N}$$

**Backpropagation** can compute this gradient!

And it's a **special case of a more general algorithm** called reverse-mode automatic differentiation that can compute the gradient of any differentiable function efficiently!
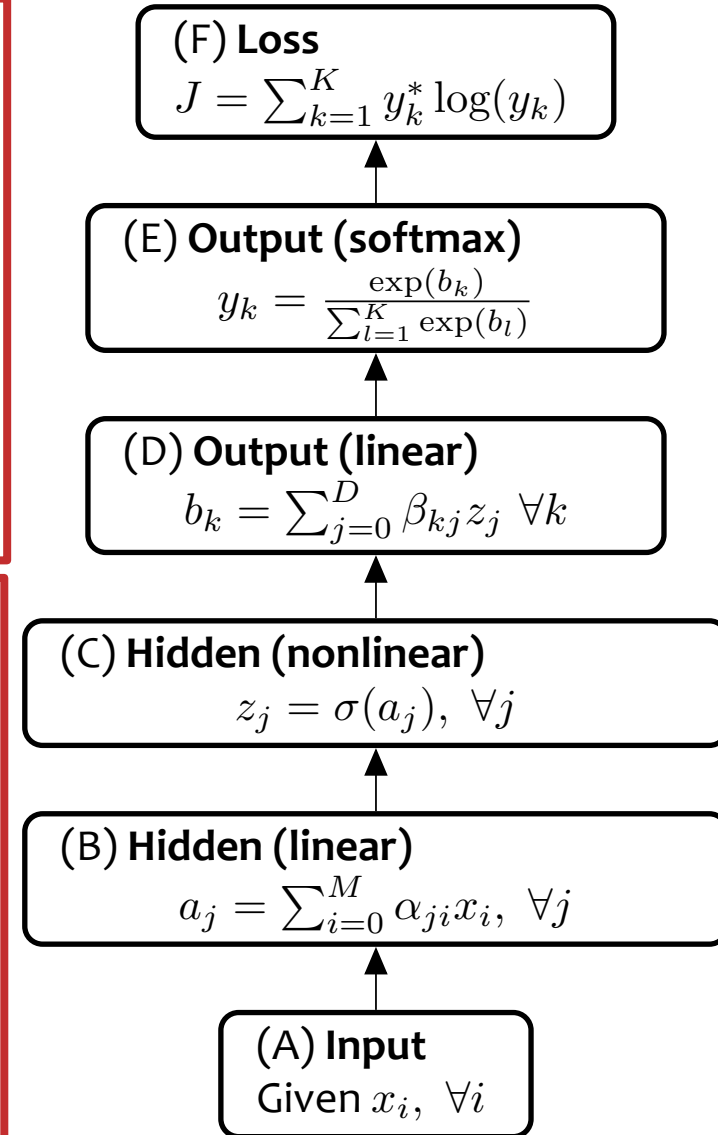
2. Choose each of t
   – Decision functio

$$\hat{y} = f_{\boldsymbol{\theta}}(\boldsymbol{x}_i)$$

   – Loss function

$$\ell(\hat{\boldsymbol{y}}, \boldsymbol{y}_i) \in \mathbb{R}$$

opposite the gradient)

$$\boldsymbol{\theta}^{(t)} \quad {}^{(t)} - \eta_t \nabla \ell(f_{\boldsymbol{\theta}}(\boldsymbol{x}_i), \boldsymbol{y}_i)$$

# Backpropagation: Abstract Picture



Output

$y_1$ ... $y_K$

Hidden Layer

$z_1$ $z_2$ ... $z_D$

Input

$x_1$ $x_2$ $x_3$ ... $x_M$

## Forward

5. $J = -\mathbf{y}^T \log \hat{\mathbf{y}}$

4. $\hat{\mathbf{y}} = \text{softmax}(\mathbf{b})$

3. $\mathbf{b} = \boldsymbol{\beta}\mathbf{z}$

2. $\mathbf{z} = \sigma(\mathbf{a})$

1. $\mathbf{a} = \boldsymbol{\alpha}\mathbf{x}$

## Backward

6. $\mathbf{g_{\hat{y}}} = -\mathbf{y} \div \hat{\mathbf{y}}$

7. $\mathbf{g_b} = \mathbf{g_{\hat{y}}^T}\left(\text{diag}(\hat{\mathbf{y}}) - \hat{\mathbf{y}}\hat{\mathbf{y}}^T\right)$

8. $\mathbf{g_\beta} = \mathbf{g_b^T}\mathbf{z}^T$

$\mathbf{g_z} = \boldsymbol{\beta}^T\mathbf{g_b^T}$

10. $\mathbf{g_a} = \mathbf{g_z} \odot \mathbf{z} \odot (1 - \mathbf{z})$

11. $\mathbf{g_\alpha} = \mathbf{g_a}\mathbf{x}^T$

(F) **Loss**
$J = \sum_{k=1}^{K} y_k^* \log(y_k)$

(E) **Output (softmax)**
$y_k = \frac{\exp(b_k)}{\sum_{l=1}^{K}\exp(b_l)}$

(D) **Output (linear)**
$b_k = \sum_{j=0}^{D}\beta_{kj}z_j \ \forall k$

(C) **Hidden (nonlinear)**
$z_j = \sigma(a_j), \ \forall j$

(B) **Hidden (linear)**
$a_j = \sum_{i=0}^{M}\alpha_{ji}x_i, \ \forall j$

(A) **Input**
Given $x_i, \ \forall i$

8

# Backpropagation: Procedural Method

**Algorithm 1** Forward Computation

  1:  **procedure** NNFORWARD(Training example $(\mathbf{x}, \mathbf{y})$, Params $\alpha, \beta$)

  2:       $\mathbf{a} = \alpha \mathbf{x}$

  3:       $\mathbf{z} = \sigma(\mathbf{a})$

  4:       $\mathbf{b} = \beta \mathbf{z}$

  5:       $\hat{\mathbf{y}} = \text{softmax}(\mathbf{b})$

  6:       $J = -\mathbf{y}^T \log \hat{\mathbf{y}}$

  7:       $\mathbf{o} = \texttt{object}(\mathbf{x}, \mathbf{a}, \mathbf{z}, \mathbf{b}, \hat{\mathbf{y}}, J)$

  8:       **return** intermediate quantities $\mathbf{o}$

**Algorithm 2** Backpropagation

  1:  **procedure** NNBACKWARD(Training example $(\mathbf{x}, \mathbf{y})$, Params $\alpha, \beta$, Intermediates $\mathbf{o}$)

  2:       Place intermediate quantities $\mathbf{x}, \mathbf{a}, \mathbf{z}, \mathbf{b}, \hat{\mathbf{y}}, J$ in $\mathbf{o}$ in scope

  3:       $\mathbf{g}_{\hat{\mathbf{y}}} = -\mathbf{y} \div \hat{\mathbf{y}}$

  4:       $\mathbf{g}_{\mathbf{b}} = \mathbf{g}_{\hat{\mathbf{y}}}^T \left( \text{diag}(\hat{\mathbf{y}}) - \hat{\mathbf{y}}\hat{\mathbf{y}}^T \right)$

  5:       $\mathbf{g}_{\beta} = \mathbf{g}_{\mathbf{b}}^T \mathbf{z}^T$

  6:       $\mathbf{g}_{\mathbf{z}} = \beta^T \mathbf{g}_{\mathbf{b}}^T$

  7:       $\mathbf{g}_{\mathbf{a}} = \mathbf{g}_{\mathbf{z}} \odot \mathbf{z} \odot (1 - \mathbf{z})$

  8:       $\mathbf{g}_{\alpha} = \mathbf{g}_{\mathbf{a}} \mathbf{x}^T$

  9:       **return** parameter gradients $\mathbf{g}_{\alpha}, \mathbf{g}_{\beta}$

**Drawbacks of Procedural Method**

1. Hard to reuse / adapt for other models

2. (Possibly) harder to make individual steps more efficient

3. Hard to find source of error if finite-difference check reports an error (since it tells you only that there is an error somewhere in those 17 lines of code)
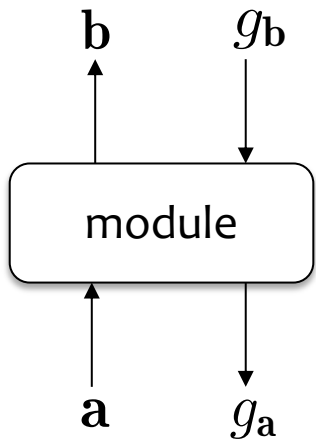
9

# Module-based AutoDiff

Module-based automatic differentiation (AD / Autodiff) is a technique that has long been used to develop libraries for deep learning

- **Dynamic neural network packages** allow a specification of the computation graph dynamically at runtime
  - PyTorch http://pytorch.org
  - Torch http://torch.ch
  - DyNet https://dynet.readthedocs.io
  - TensorFlow with Eager Execution https://www.tensorflow.org
- **Static neural network packages** require a static specification of a computation graph which is subsequently compiled into code
  - TensorFlow with Graph Execution https://www.tensorflow.org
  - Aesara (and Theano) https://aesara.readthedocs.io
  - *(These libraries are also module-based, but herein by "module-based AD" we mean the dynamic approach)*

# Module-based AutoDiff

- **Key Idea:**
  - – componentize the computation of the neural-network into layers
  - – each layer consolidates multiple **real-valued nodes** in the computation graph (a subset of them) into one **vector-valued node** (aka. a **module**)
- Each **module** is capable of two actions:

  1. Forward computation of output $\mathbf{b} = [b_1, \ldots, b_B]$ given input $\mathbf{a} = [a_1, \ldots, a_A]$ via some differentiable function $f$. That is $\mathbf{b} = f(\mathbf{a})$.

  2. Backward computation of the gradient of the input $\mathbf{g_a} = \nabla_{\mathbf{a}} J = [\frac{\partial J}{\partial a_1}, \ldots, \frac{\partial J}{\partial a_A}]$ given the gradient of output $\mathbf{g_b} = \nabla_{\mathbf{b}} J = [\frac{\partial J}{\partial b_1}, \ldots, \frac{\partial J}{\partial b_B}]$, where $J$ is the final real-valued output of the entire computation graph. This is done via the chain rule $\frac{\partial J}{\partial a_i} = \sum_{j=1}^{J} \frac{\partial J}{\partial b_j} \frac{db_j}{da_i}$ for all $i \in \{1, \ldots, A\}$.

$\mathbf{b}$    $g_{\mathbf{b}}$

module

$\mathbf{a}$    $g_{\mathbf{a}}$

# Module-based AutoDiff

**Dimensions:** input $\mathbf{a} \in \mathbb{R}^A$, output $\mathbf{b} \in \mathbb{R}^B$, gradient of output $\mathbf{g_a} \triangleq \nabla_{\mathbf{a}} J \in \mathbb{R}^A$, and gradient of input $\mathbf{g_b} \triangleq \nabla_{\mathbf{b}} J \in \mathbb{R}^B$.

**Sigmoid Module** The sigmoid layer has only one input vector $\mathbf{a}$. Below $\sigma$ is the sigmoid applied element-wise, and $\odot$ is element-wise multiplication s.t. $\mathbf{u} \odot \mathbf{v} = [u_1 v_1, \ldots, u_M v_M]$.

1: **procedure** SIGMOIDFORWARD($\mathbf{a}$)
2:    $\mathbf{b} = \sigma(\mathbf{a})$
3:    **return** $\mathbf{b}$
4: **procedure** SIGMOIDBACKWARD($\mathbf{a}, \mathbf{b}, \mathbf{g_b}$)
5:    $\mathbf{g_a} = \mathbf{g_b} \odot \mathbf{b} \odot (1 - \mathbf{b})$
6:    **return** $\mathbf{g_a}$

**Softmax Module** The softmax layer has only one input vector $\mathbf{a}$. For any vector $\mathbf{v} \in \mathbb{R}^D$, we have that $\text{diag}(\mathbf{v})$ returns a $D \times D$ diagonal matrix whose diagonal entries are $v_1, v_2, \ldots, v_D$ and whose non-diagonal entries are zero.

1: **procedure** SOFTMAXFORWARD($\mathbf{a}$)
2:    $\mathbf{b} = \text{softmax}(\mathbf{a})$
3:    **return** $\mathbf{b}$
4: **procedure** SOFTMAXBACKWARD($\mathbf{a}, \mathbf{b}, \mathbf{g_b}$)
5:    $\mathbf{g_a} = \mathbf{g_b}^T \left( \text{diag}(\mathbf{b}) - \mathbf{b}\mathbf{b}^T \right)$
6:    **return** $\mathbf{g_a}$

**Linear Module** The linear layer has two inputs: a vector $\mathbf{a}$ and parameters $\omega \in \mathbb{R}^{B \times A}$. The output $\mathbf{b}$ is not used by LINEARBACKWARD, but we pass it in for consistency of form.

1: **procedure** LINEARFORWARD($\mathbf{a}, \omega$)
2:    $\mathbf{b} = \omega \mathbf{a}$
3:    **return** $\mathbf{b}$
4: **procedure** LINEARBACKWARD($\mathbf{a}, \omega, \mathbf{b}, \mathbf{g_b}$)
5:    $\mathbf{g_\omega} = \mathbf{g_b} \mathbf{a}^T$
6:    $\mathbf{g_a} = \omega^T \mathbf{g_b}$
7:    **return** $\mathbf{g_\omega}, \mathbf{g_a}$

**Cross-Entropy Module** The cross-entropy layer has two inputs: a gold one-hot vector $\mathbf{a}$ and a predicted probability distribution $\hat{\mathbf{a}}$. It's output $b \in \mathbb{R}$ is a scalar. Below $\div$ is element-wise division. The output $b$ is not used by CROSSENTROPYBACKWARD, but we pass it in for consistency of form.

1: **procedure** CROSSENTROPYFORWARD($\mathbf{a}, \hat{\mathbf{a}}$)
2:    $b = -\mathbf{a}^T \log \hat{\mathbf{a}}$
3:    **return** $\mathbf{b}$
4: **procedure** CROSSENTROPYBACKWARD($\mathbf{a}, \hat{\mathbf{a}}, b, g_b$)
5:    $\mathbf{g_{\hat{a}}} = -g_b(\mathbf{a} \div \hat{\mathbf{a}})$
6:    **return** $\mathbf{g_a}$

# Module-based AutoDiff

**Algorithm 1** Forward Computation

1: **procedure** NNFORWARD(Training example $(\mathbf{x}, \mathbf{y})$, Parameters $\alpha$, $\beta$)
2:     $\mathbf{a} = \text{LINEARFORWARD}(\mathbf{x}, \alpha)$
3:     $\mathbf{z} = \text{SIGMOIDFORWARD}(\mathbf{a})$
4:     $\mathbf{b} = \text{LINEARFORWARD}(\mathbf{z}, \beta)$
5:     $\hat{\mathbf{y}} = \text{SOFTMAXFORWARD}(\mathbf{b})$
6:     $J = \text{CROSSENTROPYFORWARD}(\mathbf{y}, \hat{\mathbf{y}})$
7:     $\mathbf{o} = \text{object}(\mathbf{x}, \mathbf{a}, \mathbf{z}, \mathbf{b}, \hat{\mathbf{y}}, J)$
8:     **return** intermediate quantities $\mathbf{o}$

**Algorithm 2** Backpropagation

1: **procedure** NNBACKWARD(Training example $(\mathbf{x}, \mathbf{y})$, Parameters $\alpha, \beta$, Intermediates $\mathbf{o}$)
2:     Place intermediate quantities $\mathbf{x}, \mathbf{a}, \mathbf{z}, \mathbf{b}, \hat{\mathbf{y}}, J$ in $\mathbf{o}$ in scope
3:     $g_J = \frac{dJ}{dJ} = 1$                                        ▷ Base case
4:     $\mathbf{g_{\hat{y}}} = \text{CROSSENTROPYBACKWARD}(\mathbf{y}, \hat{\mathbf{y}}, J, g_J)$
5:     $\mathbf{g_b} = \text{SOFTMAXBACKWARD}(\mathbf{b}, \hat{\mathbf{y}}, \mathbf{g_{\hat{y}}})$
6:     $\mathbf{g_\beta}, \mathbf{g_z} = \text{LINEARBACKWARD}(\mathbf{z}, \mathbf{b}, \mathbf{g_b})$
7:     $\mathbf{g_a} = \text{SIGMOIDBACKWARD}(\mathbf{a}, \mathbf{z}, \mathbf{g_z})$
8:     $\mathbf{g_\alpha}, \mathbf{g_x} = \text{LINEARBACKWARD}(\mathbf{x}, \mathbf{a}, \mathbf{g_a})$       ▷ We discard $\mathbf{g_x}$
9:     **return** parameter gradients $\mathbf{g_\alpha}, \mathbf{g_\beta}$

**Advantages of Module-based AutoDiff**

1. Easy to reuse / adapt for other models
2. Encapsulated layers are easier to optimize (e.g. implement in C++ or CUDA)
3. Easier to find bugs because we can run a finite-difference check on each layer separately

# Module-based AutoDiff (OOP Version)

Object-Oriented Implementation:
- Let each module be an **object**
- Then allow the **control flow** dictate the creation of the **computation graph**
- No longer need to implement NNBackward($\cdot$), just follow the computation graph in **reverse topological order**

```
1   class  Sigmoid (Module)
2       method  forward(a)
3           b = σ(a)
4           return  b
5       method  backward(a , b , g_b)
6           g_a = g_b ⊙ b ⊙ (1 − b)
7           return  g_a
```

```
1   class  Linear (Module)
2       method  forward(a , ω)
3           b = ωa
4           return  b
5       method  backward(a , ω , b , g_b)
6           g_ω = g_b a^T
7           g_a = ω^T g_b
8           return  g_ω, g_a
```

```
1   class  Softmax (Module)
2       method  forward(a)
3           b = softmax(a)
4           return  b
5       method  backward(a , b , g_b)
6           g_a = g_b^T (diag(b) − bb^T)
7           return  g_a
```

```
1   class  CrossEntropy (Module)
2       method  forward(a , â)
3           b = −a^T log â
4           return  b
5       method  backward(a , â , b , g_b)
6           g_â = −g_b(a ÷ â)
7           return  g_a
```

# Module-based AutoDiff (OOP Version)

```
1   class NeuralNetwork(Module):
2
3       method init()
4           lin1_layer = Linear()
5           sig_layer = Sigmoid()
6           lin2_layer = Linear()
7           soft_layer = Softmax()
8           ce_layer = CrossEntropy()
9
10      method forward(Tensor x, Tensor y, Tensor α, Tensor β)
11          a = lin1_layer.apply_fwd(x, α)
12          z = sig_layer.apply_fwd(a)
13          b = lin1_layer.apply_fwd(z, β)
14          ŷ = soft_layer.apply_fwd(b)
15          J = ce_layer.apply_fwd(y, ŷ)
16          return J.out_tensor
17
18      method backward(Tensor x, Tensor y, Tensor α, Tensor β)
19          tape_bwd()
20          return lin1_layer.in_gradients[1], lin2_layer.in_gradients[1]
```

# Module-based AutoDiff (OOP Version)

```
1    class NeuralNetwork(Module):
2
3        method init()
4            lin1_layer = Linear()
5            sig_layer = Sigmoid()
6            lin2_layer = Linear()
7            soft_layer = Softmax()
8            ce_layer = CrossEntropy()
9
10       method forward(Tensor x, Tensor y, Tensor c
11           a = lin1_layer.apply_fwd(x, α)
12           z = sig_layer.apply_fwd(a)
13           b = lin1_layer.apply_fwd(z, β)
14           ŷ = soft_layer.apply_fwd(b)
15           J = ce_layer.apply_fwd(y, ŷ)
16           return J.out_tensor
17
18       method backward(Tensor x, Tensor y, Tensor
19           tape_bwd()
20           return lin1_layer.in_gradients[1], lin2_layer
```

```
1    global tape = stack()
2
3    class Module:
4
5        method init()
6            out_tensor = null
7            out_gradient = 1
8
9        method apply_fwd(List in_modules)
10           in_tensors = [x.out_tensor for x in in_modules]
11           out_tensor = forward(in_tensors)
12           tape.push(self)
13           return self
14
15       method apply_bwd():
16           in_gradients = backward(in_tensors, out_tensor, out_gradient)
17           for i in 1,..., len(in_modules):
18               in_modules[i].out_gradient += in_gradients[i]
19           return self
20
21   function tape_bwd():
22       while len(tape) > 0
23           m = tape.pop()
24           m.apply_bwd()
```

# Module-based AutoDiff (OOP Version)

```
1    global tape = stack()
2
3    class Module:
4
5        method init()
6            out_tensor = null
7            out_gradient = 1
8
9        method apply_fwd(List in_modules)
10           in_tensors = [x.out_tensor for x in in_modules]
11           out_tensor = forward(in_tensors)
12           tape.push(self)
13           return self
14
15       method apply_bwd():
16           in_gradients = backward(in_tensors, out_tensor, out_gradient)
17           for i in 1,..., len(in_modules):
18               in_modules[i].out_gradient += in_gradients[i]
19           return self
20
21   function tape_bwd():
22       while len(tape) > 0
23           m = tape.pop()
24           m.apply_bwd()
```

# PyTorch

The same simple neural network we defined in pseudocode can also be defined in PyTorch.

```python
1  # Define model
2  class NeuralNetwork(nn.Module):
3      def __init__(self):
4          super(NeuralNetwork, self).__init__()
5          self.flatten = nn.Flatten()
6          self.linear1 = nn.Linear(28*28, 512)
7          self.sigmoid = nn.Sigmoid()
8          self.linear2 = nn.Linear(512,512)
9
10     def forward(self, x):
11         x = self.flatten(x)
12         a = self.linear1(x)
13         z = self.sigmoid(a)
14         b = self.linear2(z)
15         return b
16
17 # Take one step of SGD
18 def one_step_of_sgd(X, y):
19     loss_fn = nn.CrossEntropyLoss()
20     optimizer = torch.optim.SGD(model.parameters(), lr=1e-3)
21
22     # Compute prediction error
23     pred = model(X)
24     loss = loss_fn(pred, y)
25
26     # Backpropagation
27     optimizer.zero_grad()
28     loss.backward()
29     optimizer.step()
```

Example adapted from https://pytorch.org/tutorials/beginner/basics/quickstart_tutorial.html

# PyTorch

**Q:** Why don't we call linear.forward() in PyTorch?

**A:** This is just syntactic sugar. There's a special method in Python `__call__` that allows you to define what happens when you treat an object as if it were a function.

In other words, running the following:
```
linear(x)
```
is equivalent to running:
```
linear.__call__(x)
```
which in PyTorch is (nearly) the same as running:
```
linear.forward(x)
```

This is because PyTorch defines every Module's `__call__` method to be something like this:
```
def __call__(self):
    self.forward()
```

# PyTorch

**Q:** Why don't we pass in the parameters to a PyTorch Module?

**A:** This just makes your code cleaner.

In PyTorch, you store the parameters inside the Module and "mark" them as parameters that should contribute to the eventual gradient used by an optimizer

```
0      method forward(Tensor x , Tensor y , Tensor α , Tensor β)
11         a =lin1_layer.apply_fwd(x, α)
12         z =sig_layer.apply_fwd(a)
3          b =lin1_layer.apply_fwd(z, β)
4          ŷ =soft_layer.apply_fwd(b)
5          J =ce_layer.apply_fwd(y, ŷ)
6          return J.out_tensor
```

```python
9
10    def forward(self, x):
11        x = self.flatten(x)
12        a = self.linear1(x)
13        z = self.sigmoid(a)
14        b = self.linear2(z)
15        return b
```

# LARGE LANGUAGE MODELS

# What is ChatGPT?

- ChatGPT is a large (in the sense of having many parameters) language model, fine-tuned to be a dialogue agent
- The base language model is GPT-3.5 which was trained on a large quantity of text
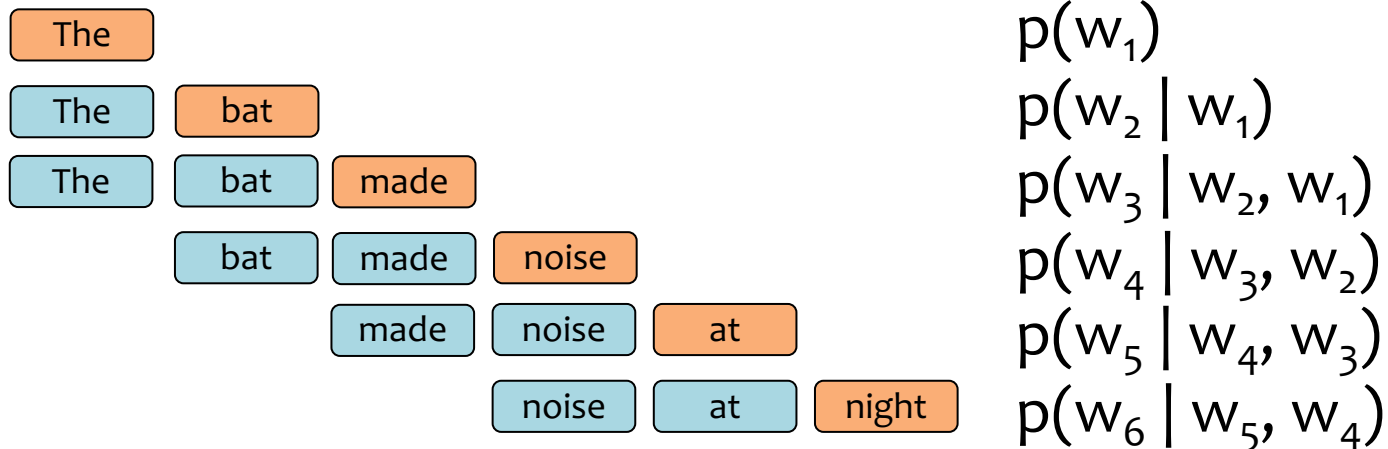
# TASK: LANGUAGE MODELING

# n-Gram Language Model

*Question*: How can we **define** a probability distribution over a sequence of length T?
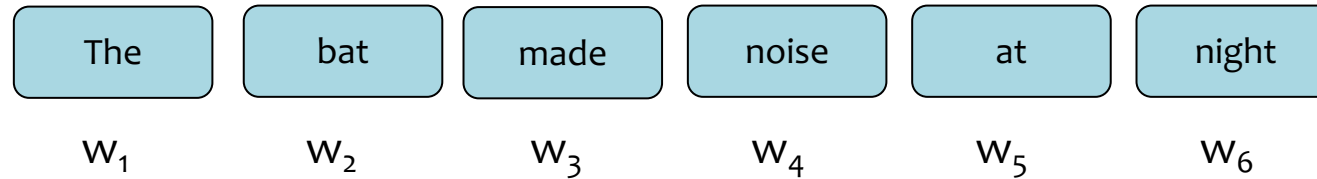
| The | bat | made | noise | at | night |
|-----|-----|------|-------|-----|-------|
| $w_1$ | $w_2$ | $w_3$ | $w_4$ | $w_5$ | $w_6$ |

**n-Gram Model (n=3)**

$$p(w_1, w_2, \ldots, w_T) = \prod_{t=1}^{T} p(w_t \mid w_{t-1}, w_{t-2})$$
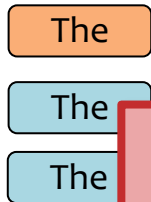
$$p(w_1, w_2, w_3, \ldots, w_6) =$$

| The | | | $p(w_1)$ |
| The | bat | | $p(w_2 \mid w_1)$ |
| The | bat | made | $p(w_3 \mid w_2, w_1)$ |
| bat | made | noise | $p(w_4 \mid w_3, w_2)$ |
| made | noise | at | $p(w_5 \mid w_4, w_3)$ |
| noise | at | night | $p(w_6 \mid w_5, w_4)$ |

# n-Gram Language Model

*Question*: How can we **define** a probability distribution over a sequence of length T?

| The | bat | made | noise | at | night |
|-----|-----|------|-------|----|----|
| $w_1$ | $w_2$ | $w_3$ | $w_4$ | $w_5$ | $w_6$ |

**n-Gram Model (n=3)**
$$p(w_1, w_2, \ldots, w_T) = \prod_{t=1}^{T} p(w_t \mid w_{t-1}, w_{t-2})$$

$p(w_1, w_2, w_3, \ldots, w_6) =$

The

$p(w_1)$

The

The

*Note*: This is called a **model** because we made some **assumptions** about how many previous words to condition on
(i.e. only n-1 words)

# Learning an n-Gram Model

*Question*: How do we **learn** the probabilities for the n-Gram Model?

$p(w_t \mid w_{t-2} = \text{The}, w_{t-1} = \text{bat})$

| $w_t$ | $p(\cdot \mid \cdot, \cdot)$ |
|---|---|
| ate | 0.015 |
| … | |
| flies | 0.046 |
| … | |
| zebra | 0.000 |

$p(w_t \mid w_{t-2} = \text{made}, w_{t-1} = \text{noise})$

| $w_t$ | $p(\cdot \mid \cdot, \cdot)$ |
|---|---|
| at | 0.020 |
| … | |
| pollution | 0.030 |
| … | |
| zebra | 0.000 |

$p(w_t \mid w_{t-2} = \text{cows}, w_{t-1} = \text{eat})$

| $w_t$ | $p(\cdot \mid \cdot, \cdot)$ |
|---|---|
| corn | 0.420 |
| … | |
| grass | 0.510 |
| … | |
| zebra | 0.000 |

# Learning an n-Gram Model

*Question*: How do we **learn** the probabilities for the n-Gram Model?

*Answer*: From data! Just **count** n-gram frequencies

$$p(w_t \mid w_{t-2} = \text{cows}, w_{t-1} = \text{eat})$$

…the **cows eat grass**…
…our **cows eat hay** daily…
…factory-farm **cows eat corn**…
…on an organic farm, **cows eat hay** and…
…do your **cows eat grass** or corn?…
…what do **cows eat if** they have…
…**cows eat corn** when there is no…
…which **cows eat which** foods depends…
…if **cows eat grass**…
…when **cows eat corn** their stomachs…
…should we let **cows eat corn**?…
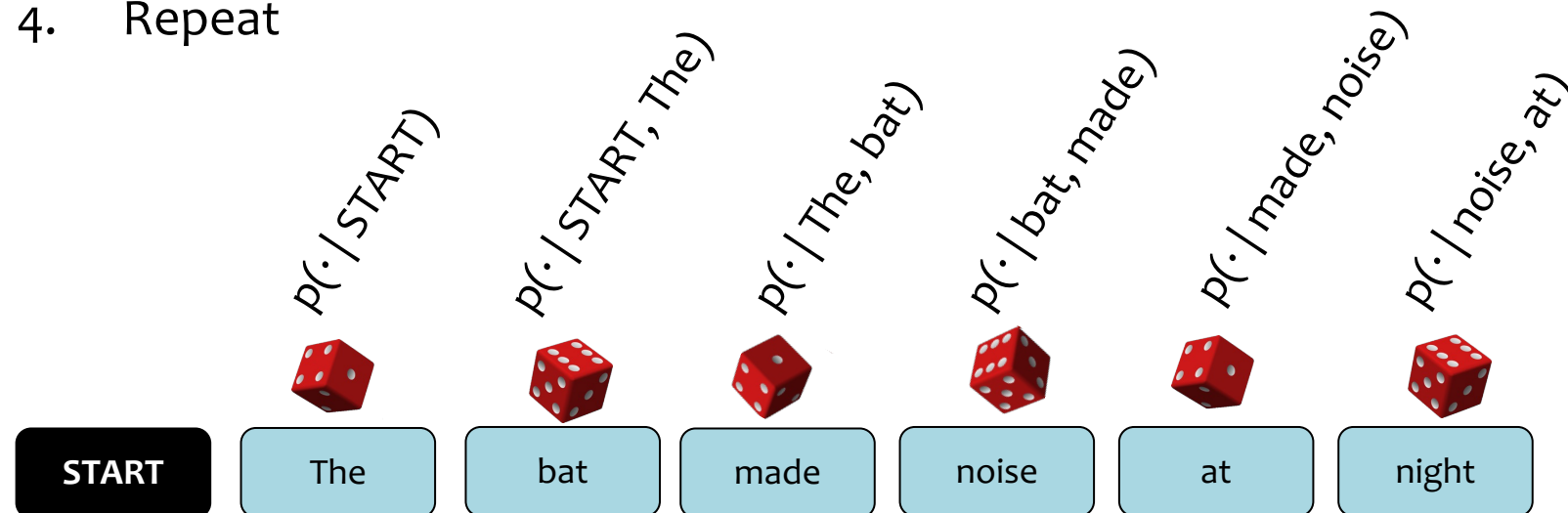
| $w_t$ | $p(\cdot \mid \cdot, \cdot)$ |
|---|---|
| corn | 4/11 |
| grass | 3/11 |
| hay | 2/11 |
| if | 1/11 |
| which | 1/11 |

# Sampling from a Language Model

*Question*: How do we sample from a Language Model?

*Answer*:

1. Treat each probability distribution like a (50k-sided) weighted die
2. Pick the die corresponding to $p(w_t | w_{t-2}, w_{t-1})$
3. Roll that die and generate whichever word $w_t$ lands face up
4. Repeat



| $p(\cdot | \text{START})$ | $p(\cdot | \text{START, The})$ | $p(\cdot | \text{The, bat})$ | $p(\cdot | \text{bat, made})$ | $p(\cdot | \text{made, noise})$ | $p(\cdot | \text{noise, at})$ |

| START | The | bat | made | noise | at | night |

# Noisy Channel Models

- Prior to 2017, two tasks relied heavily on language models:
  - speech recognition
  - machine translation
- Definition: a **noisy channel model** combines a *transduction model* (probability of converting **y** to **x**) with a *language model* (probability of **y**)

$$\hat{\mathbf{y}} = \operatorname*{argmax}_{\mathbf{y}} p(\mathbf{y} \mid \mathbf{x}) = \operatorname*{argmax}_{\mathbf{y}} p(\mathbf{x} \mid \mathbf{y})p(\mathbf{y})$$

<span style="color:red">transduction model</span>    <span style="color:red">language model</span>

- **Goal:** to recover **y** from **x**
  - *For speech:* **x** is acoustic signal, **y** is transcription
  - *For machine translation:* **x** is sentence in source language, **y** is sentence in target language

# Large (n-Gram) Language Models

- The earliest (truly) large language models were n-gram models

- Google n-Grams:
  - 2006: first release, English n-grams
    - trained on **1 trillion tokens** of web text (95 billion sentences)
    - included 1-grams, 2-grams, 3-grams, 4-grams, and 5-grams
  - 2009 – 2010: n-grams in Japanese, Chinese, Swedish, Spanish, Romanian, Portuguese, Polish, Dutch, Italian, French, German, Czech

English n-gram model is **~3 billion parameters**

| Number of unigrams: | 13,588,391 |
|---|---|
| Number of bigrams: | 314,843,401 |
| Number of trigrams: | 977,069,902 |
| Number of fourgrams: | 1,313,818,354 |
| Number of fivegrams: | 1,176,470,663 |

```
serve as the incoming 92
serve as the incubator 99
serve as the independent 794
serve as the index 223
serve as the indication 72
serve as the indicator 120
serve as the indicators 45
serve as the indispensable 111
serve as the indispensible 40
serve as the individual 234
serve as the industrial 52
serve as the industry 607
```

```
accessoire Accessoires </S> 515
accessoire Accord i-CTDi 65
accessoire Accra accu 312
accessoire Acheter cet 1402
accessoire Ajouter au 160
accessoire Amour Beauté 112
accessoire Annuaire LOEIL 49
accessoire Architecture artiste 531
accessoire Attention : 44
```

```
惯例 为 电影 创作      52
惯例 为 的 是          95
惯例 为 目标 职位      49
惯例 为 确保 合作      69
惯例 为 确保 重组      213
惯例 为 科研 和        55
惯例 为 统称 </s>     183
惯例 为 维 和          50
惯例 为 自己 的        43
惯例 为 艺术类 学院    44
惯例 为 避免 侵权     148
```

# Large (n-Gram) Language Models

- The earliest (truly) large language models were n-gram models

- Google n-Grams:
  - 2006: first release, English n-grams
    - trained on **1 trillion tokens** of web text (95 billion sentences)
    - included 1-grams, 2-grams, 3-grams, 4-grams, and 5-grams
  - 2009 – 2010: n-grams in Japanese, Chinese, Swedish, Spanish, Romanian, Portuguese, Polish, Dutch, Italian, French, German, Czech

English n-gram model is **~3 billion parameters**

| | |
|---|---|
| Number of unigrams: | 13,588,391 |
| Number of bigrams: | 314,843,401 |
| Number of trigrams: | 977,069,902 |
| Number of fourgrams: | 1,313,818,354 |
| Number of fivegrams: | 1,176,470,663 |

Q: Is this a large training set?

A: Yes!

oire Accessoir
oire Accord i—
oire Accra acc
oire Acheter
oire Ajouter a
oire Amour Bea
oire Annuaire
oire Architect
oire Attention

Q: Is this a large model?

A: Yes!

serve as the individual 234
serve as the industrial 52
serve as the industry 607

惯例 为 艺术类 学院    44
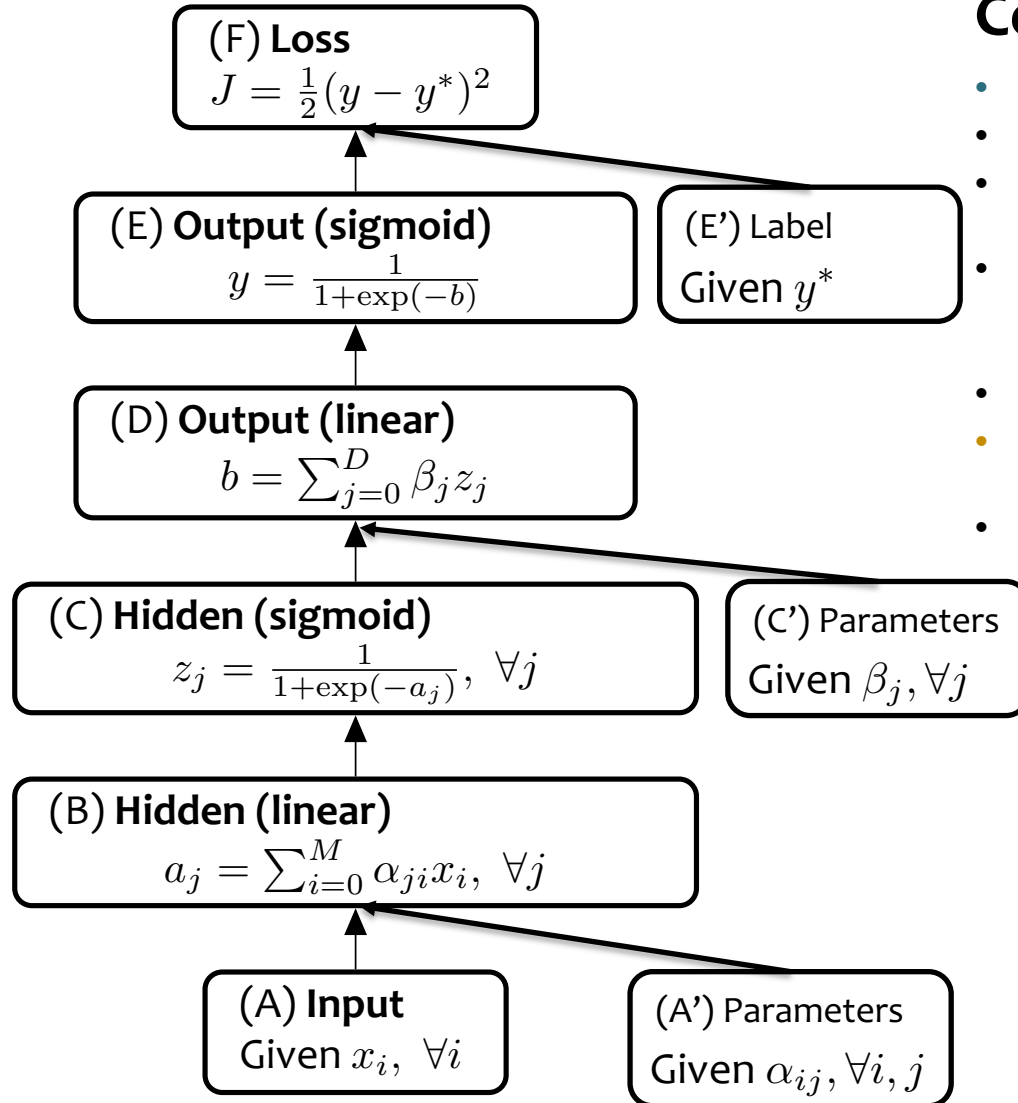惯例 为 避免 侵权    148

# How large are LLMs?

Comparison of some recent **large language models** (LLMs)

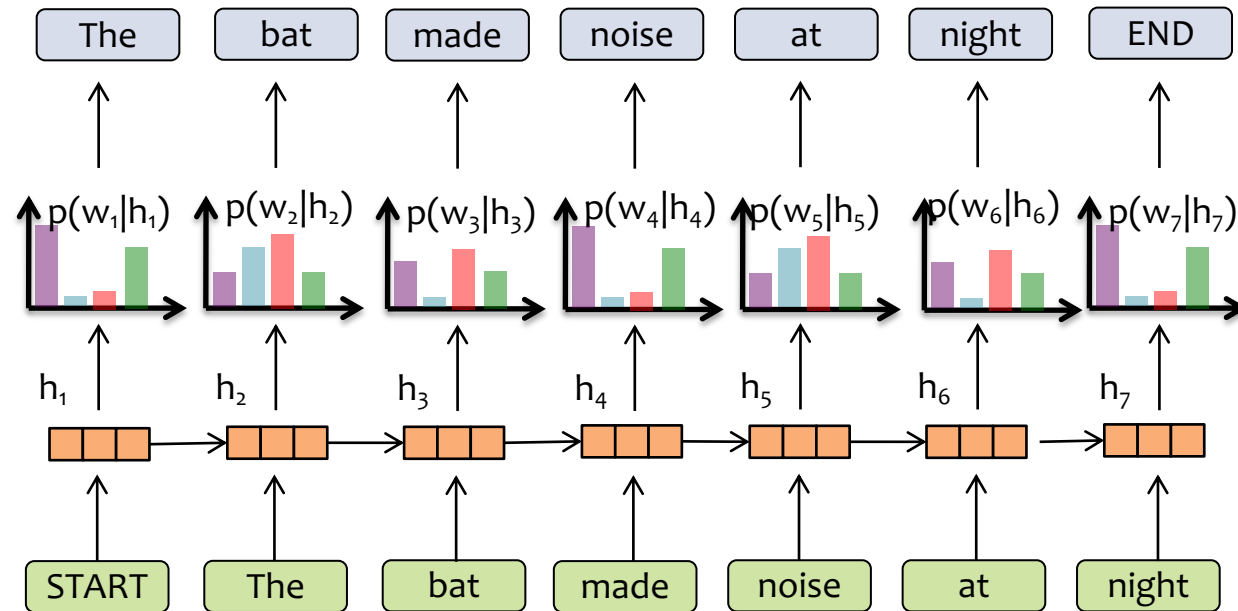| Model | Creators | Year of release | Training Data (# tokens) | Model Size (# parameters) |
|---|---|---|---|---|
| GPT-2 | OpenAI | 2019 | ~10 billion (40Gb) | 1.5 billion |
| GPT-3 (cf. ChatGPT) | OpenAI | 2020 | 300 billion | 175 billion |
| PaLM | Google | 2022 | 780 billion | 540 billion |
| Chinchilla | DeepMind | 2022 | 1.4 trillion | 70 billion |
| LaMDA (cf. Bard) | Google | 2022 | 1.56 trillion | 137 billion |
| LLaMA | Meta | 2023 | 1.4 trillion | 65 billion |
| GPT-4 | OpenAI | 2023 | ? | ? |

Transformer Language Models

# MODEL: GPT

# Ways of Drawing Neural Networks

(F) **Loss**
$$J = \frac{1}{2}(y - y^*)^2$$

(E) **Output (sigmoid)**
$$y = \frac{1}{1+\exp(-b)}$$

(E') Label
Given $y^*$

(D) **Output (linear)**
$$b = \sum_{j=0}^{D} \beta_j z_j$$

(C) **Hidden (sigmoid)**
$$z_j = \frac{1}{1+\exp(-a_j)}, \quad \forall j$$

(C') Parameters
Given $\beta_j, \forall j$

(B) **Hidden (linear)**
$$a_j = \sum_{i=0}^{M} \alpha_{ji} x_i, \quad \forall j$$

(A) **Input**
Given $x_i, \forall i$

(A') Parameters
Given $\alpha_{ij}, \forall i, j$

## Computation Graph

- The diagram represents an algorithm
- Nodes are **rectangles**
- One node per intermediate variable in the algorithm
- Node is labeled with the function that it computes (inside the box) and also the variable name (outside the box)
- *Edges are directed*
- Edges do not have labels (since they don't need them)
- For neural networks:
  - Each intercept term should appear as a node (if it's not folded in somewhere)
  - Each parameter should appear as a node
  - Each constant, e.g. a true label or a feature vector should appear in the graph
  - It's perfectly fine to include the loss
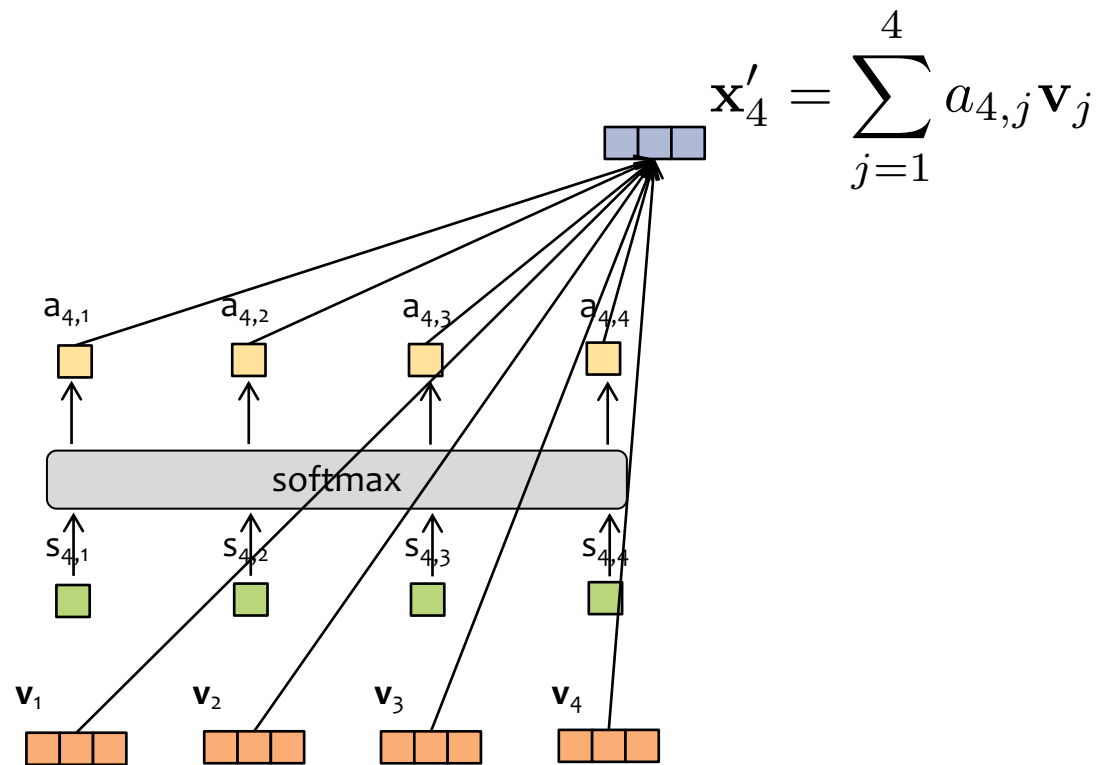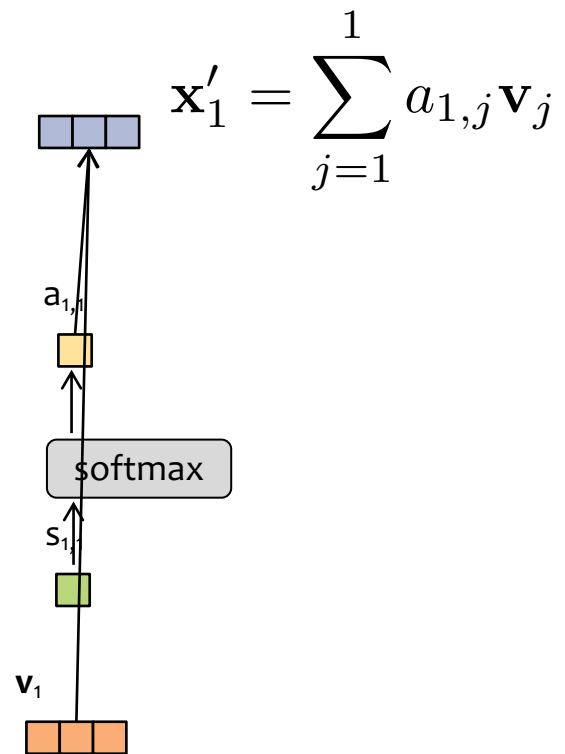
37

# RNN Language Model

_Key Idea_:
(1) convert all previous words to a **fixed length vector**
(2) define distribution $p(w_t \mid f_\theta(w_{t-1}, \ldots, w_1))$ that conditions on the vector $\mathbf{h}_t = f_\theta(w_{t-1}, \ldots, w_1)$
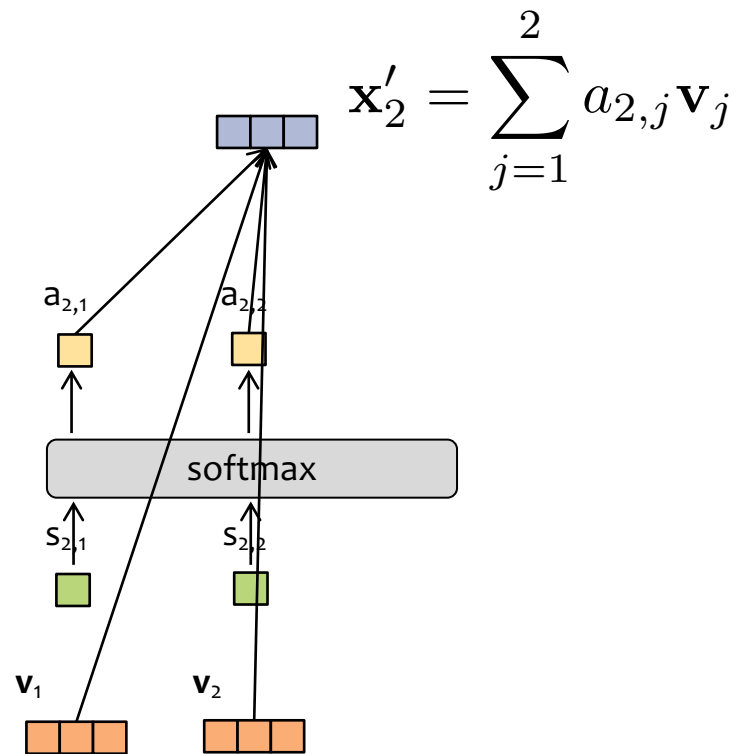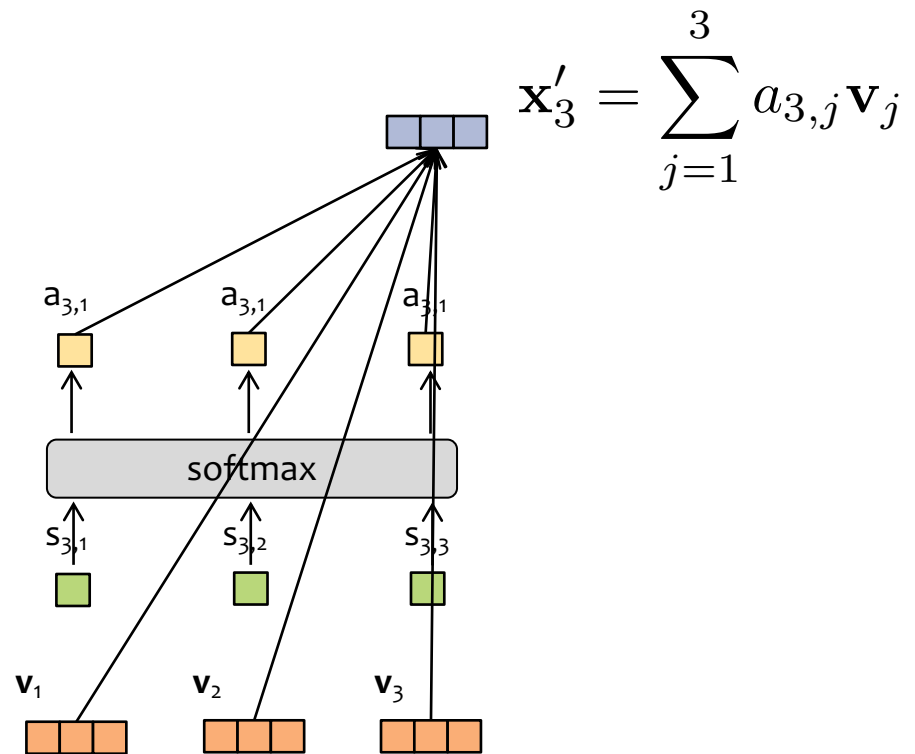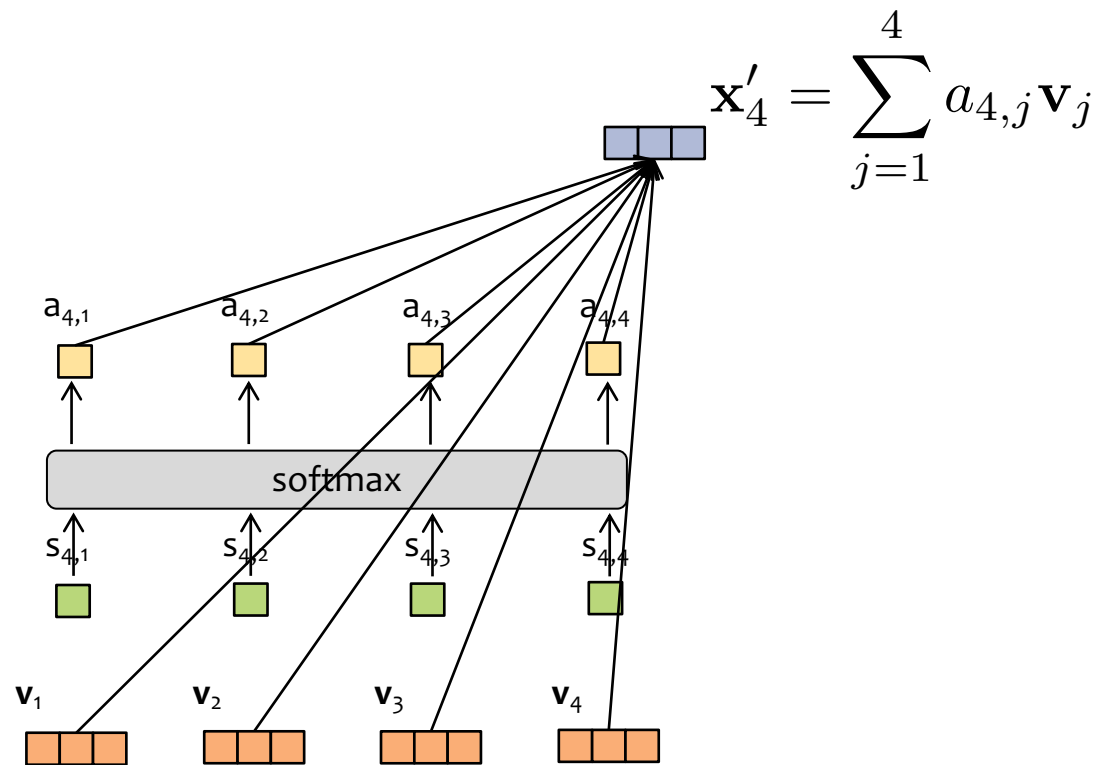
38

# RNNs and Forgetting

# Attention

$$\mathbf{x}'_4 = \sum_{j=1}^{4} a_{4,j} \mathbf{v}_j$$

$a_{4,1}$    $a_{4,2}$    $a_{4,3}$    $a_{4,4}$

softmax

$s_{4,1}$    $s_{4,2}$    $s_{4,3}$    $s_{4,4}$

$\mathbf{v}_1$    $\mathbf{v}_2$    $\mathbf{v}_3$    $\mathbf{v}_4$

# Attention

$$\mathbf{x}_1' = \sum_{j=1}^{1} a_{1,j} \mathbf{v}_j$$

$a_{1,1}$

softmax

$s_{1,1}$

$\mathbf{v}_1$

# Attention

$$\mathbf{x}'_2 = \sum_{j=1}^{2} a_{2,j} \mathbf{v}_j$$

$a_{2,1}$

$a_{2,2}$

softmax

$s_{2,1}$

$s_{2,2}$

$\mathbf{v}_1$

$\mathbf{v}_2$

# Attention

$$\mathbf{x}'_3 = \sum_{j=1}^{3} a_{3,j} \mathbf{v}_j$$

$a_{3,1}$   $a_{3,1}$   $a_{3,1}$

softmax

$s_{3,1}$   $s_{3,2}$   $s_{3,3}$

$\mathbf{v}_1$   $\mathbf{v}_2$   $\mathbf{v}_3$

# Attention



$$\mathbf{x}'_4 = \sum_{j=1}^{4} a_{4,j} \mathbf{v}_j$$

# Attention



$$\mathbf{x}_t' = \sum_{j=1}^{t} a_{t,j} \mathbf{v}_j$$

attention weights

scores

values

# Scaled Dot-Product Attention



$$\mathbf{x}'_4 = \sum_{j=1}^{4} a_{4,j} \mathbf{v}_j$$

$a_{4,1}$  $a_{4,2}$  $a_{4,3}$  $a_{4,4}$

softmax

$s_{4,1}$  $s_{4,2}$  $s_{4,3}$  $s_{4,4}$

$\mathbf{W}_v$

$\mathbf{v}_1$  $\mathbf{v}_2$  $\mathbf{v}_3$  $\mathbf{v}_4$

$$\mathbf{v}_j = \mathbf{W}_v^T \mathbf{x}_j \qquad \text{values}$$

$\mathbf{x}_1$  $\mathbf{x}_2$  $\mathbf{x}_3$  $\mathbf{x}_4$

# Scaled Dot-Product Attention

$$\mathbf{x}'_4 = \sum_{j=1}^{4} a_{4,j} \mathbf{v}_j$$



$a_{4,1}$   $a_{4,2}$   $a_{4,3}$   $a_{4,4}$

softmax

$s_{4,1}$   $s_{4,2}$   $s_{4,3}$   $s_{4,4}$

$\mathbf{W}_k$

$\mathbf{k}_1$   $\mathbf{k}_2$   $\mathbf{k}_3$   $\mathbf{k}_4$

$\mathbf{W}_v$

$\mathbf{v}_1$   $\mathbf{v}_2$   $\mathbf{v}_3$   $\mathbf{v}_4$

$\mathbf{x}_1$   $\mathbf{x}_2$   $\mathbf{x}_3$   $\mathbf{x}_4$

$$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j \qquad \text{keys}$$

$$\mathbf{v}_j = \mathbf{W}_v^T \mathbf{x}_j \qquad \text{values}$$

# Scaled Dot-Product Attention



$$\mathbf{x}'_4 = \sum_{j=1}^{4} a_{4,j} \mathbf{v}_j$$

$$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j \quad \text{queries}$$

$$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j \quad \text{keys}$$

$$\mathbf{v}_j = \mathbf{W}_v^T \mathbf{x}_j \quad \text{values}$$

# Scaled Dot-Product Attention
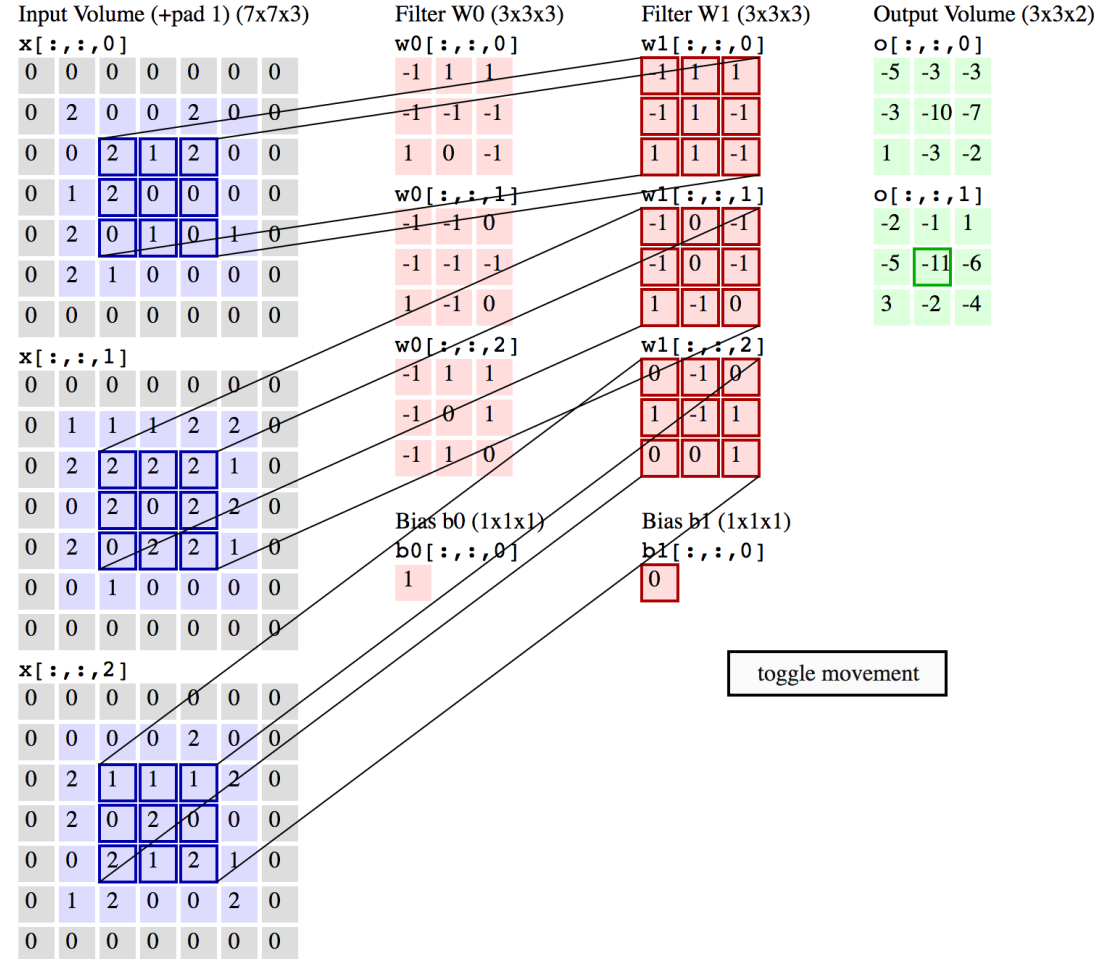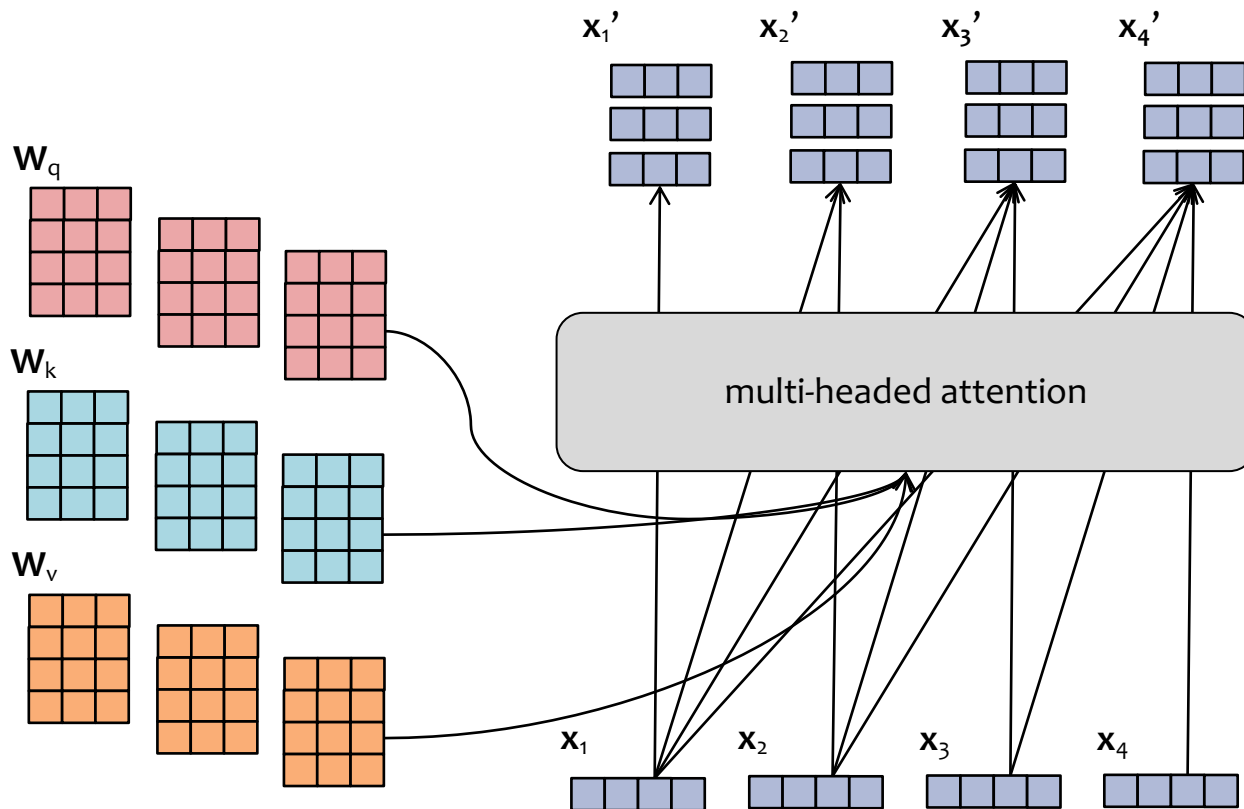


$$\mathbf{x}'_4 = \sum_{j=1}^{4} a_{4,j} \mathbf{v}_j$$

$$s_{4,j} = \mathbf{k}_j^T \mathbf{q}_4 / \sqrt{d_k} \quad \text{scores}$$

$$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j \quad \text{queries}$$

$$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j \quad \text{keys}$$

$$\mathbf{v}_j = \mathbf{W}_v^T \mathbf{x}_j \quad \text{values}$$

# Scaled Dot-Product Attention

$$\mathbf{x}'_4 = \sum_{j=1}^{4} a_{4,j} \mathbf{v}_j$$



$\mathbf{a}_4 = \text{softmax}(\mathbf{s}_4)$ attention weights

$s_{4,j} = \mathbf{k}_j^T \mathbf{q}_4 / \sqrt{d_k}$ scores

$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j$     queries

$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j$     keys

$\mathbf{v}_j = \mathbf{W}_v^T \mathbf{x}_j$     values

# Scaled Dot-Product Attention



$$\mathbf{x}'_4 = \sum_{j=1}^{4} a_{4,j} \mathbf{v}_j$$

$\mathbf{a}_4 = \text{softmax}(\mathbf{s}_4)$ attention weights

$s_{4,j} = \mathbf{k}_j^T \mathbf{q}_4 / \sqrt{d_k}$ scores

$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j$     queries

$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j$     keys

$\mathbf{v}_j = \mathbf{W}_v^T \mathbf{x}_j$     values

# Scaled Dot-Product Attention



$$\mathbf{x}'_t = \sum_{j=1}^{t} a_{t,j} \mathbf{v}_j$$

$\mathbf{a}_t = \text{softmax}(\mathbf{s}_t)$ attention weights

$s_{t,j} = \mathbf{k}_j^T \mathbf{q}_t / \sqrt{d_k}$ scores

$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j$     queries

$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j$     keys

$\mathbf{v}_j = \mathbf{W}_v^T \mathbf{x}_j$     values

# Animation of 3D Convolution

http://cs231n.github.io/convolutional-networks/

Figure from Fei-Fei Li & Andrej Karpathy & Justin Johnson (CS231N)

# Multi-headed Attention



- Just as we can have **multiple channels** in a **convolution** layer, we can use **multiple heads** in an **attention** layer
- Each head gets **its own parameters**
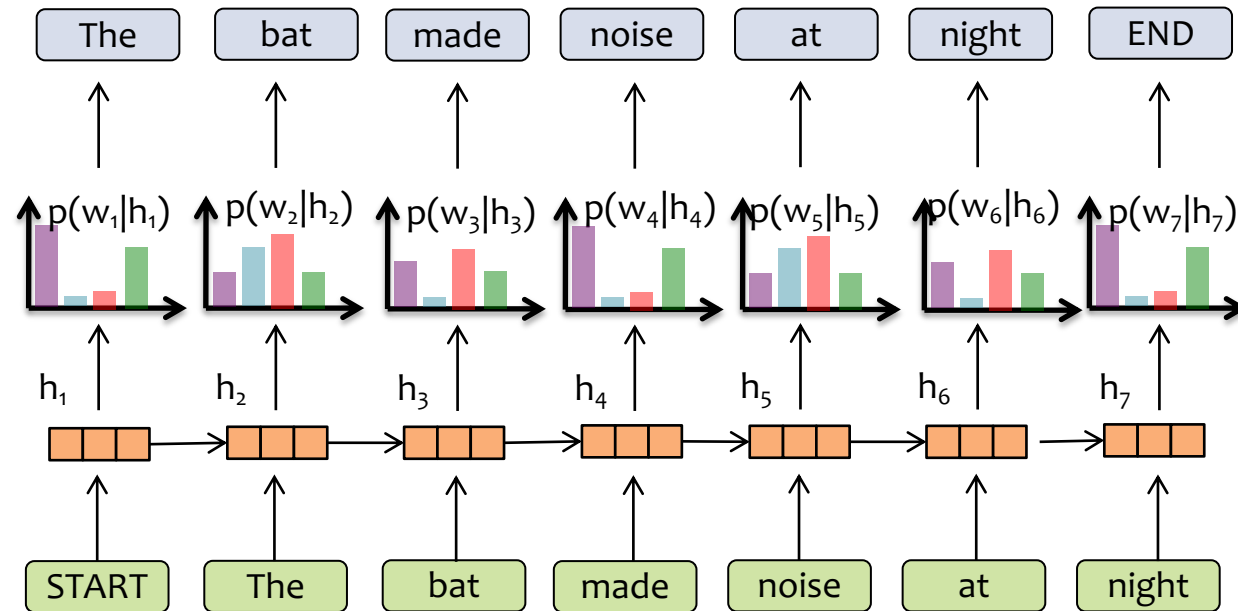- We can **concatenate** all the outputs to get a single vector for each time step

# Multi-headed Attention

- To ensure the dimension of the **input** embedding $x_t$ is the same as the **output** embedding $x_t'$, Transformers usually choose the embedding sizes and number of heads appropriately:
  - $d_{model}$ = dim. of inputs
  - $d_k$ = dim. of each output
  - h = # of heads
  - Choose $d_k = d_{model} / h$
- Then concatenate the outputs



- Just as we can have **multiple channels** in a **convolution** layer, we can use **multiple heads** in an **attention** layer
- Each head gets **its own parameters**
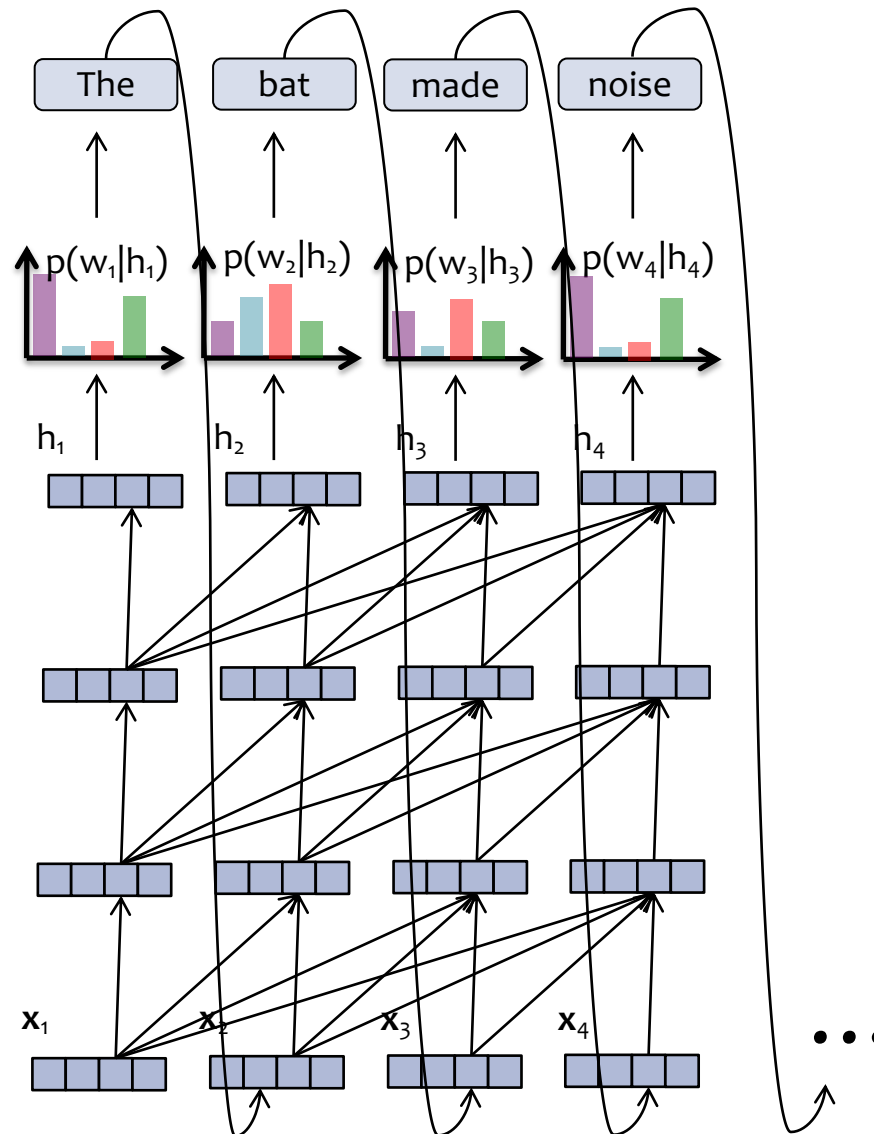- We can **concatenate** all the outputs to get a single vector for each time step

# Multi-headed Attention

- To ensure the dimension of the **input** embedding $x_t$ is the same as the **output** embedding $x_t'$, Transformers usually choose the embedding sizes and number of heads appropriately:
  - $d_{model}$ = dim. of inputs
  - $d_k$ = dim. of each output
  - h = # of heads
  - Choose $d_k = d_{model} / h$
- Then concatenate the outputs



- Just as we can have **multiple channels** in a **convolution** layer, we can use **multiple heads** in an **attention** layer
- Each head gets **its own parameters**
- We can **concatenate** all the outputs to get a single vector for each time step

# RNN Language Model

*Key Idea:*
(1) convert all previous words to a **fixed length vector**
(2) define distribution $p(w_t \mid f_\theta(w_{t-1}, \ldots, w_1))$ that conditions on the vector $\mathbf{h}_t = f_\theta(w_{t-1}, \ldots, w_1)$

# Transformer Language Model



**Important!**

- RNN computation graph grows **linearly** with the number of input tokens
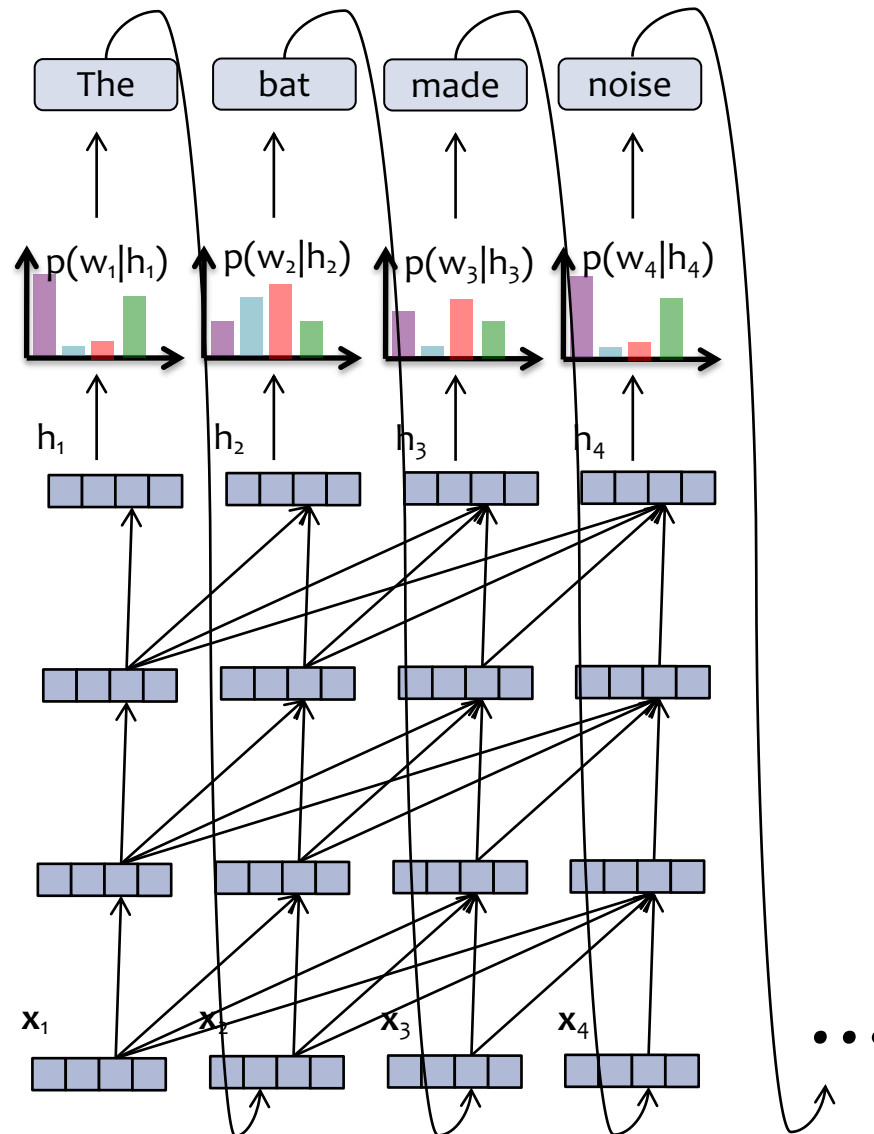- Transformer-LM computation graph grows **quadratically** with the number of input tokens

Each hidden vector looks back at the hidden vectors of the **current and previous timesteps in the previous layer.**

The language model part is just like an RNN-LM!

$p(w_1|h_1)$  $p(w_2|h_2)$  $p(w_3|h_3)$  $p(w_4|h_4)$

The   bat   made   noise

$h_1$   $h_2$   $h_3$   $h_4$

$x_1$   $x_2$   $x_3$   $x_4$

# Transformer Language Model



Important!
- RNN computation graph grows **linearly** with the number of input tokens
- Transformer-LM computation graph grows **quadratically** with the number of input tokens

**Each layer** of a Transformer LM consists of several **sublayers**:
1. attention
2. feed-forward neural network
3. layer normalization
4. residual connections

Each hidden vector looks back at the hidden vectors of the **current and previous timesteps in the previous layer.**

The language model part is just like an RNN-LM!

# Layer Normalization

- *The Problem:* **internal covariate shift** occurs during training of a deep network when a small change in the low layers amplifies into a large change in the high layers

- *One Solution:* **Layer normalization** normalizes each layer and learns elementwise gain/bias

- Such normalization allows for higher learning rates (for **faster convergence**) without issues of diverging gradients

Given input $\mathbf{a} \in \mathbb{R}^K$, LayerNorm computes output $\mathbf{b} \in \mathbb{R}^K$:

$$\mathbf{b} = \boldsymbol{\gamma} \odot \frac{\mathbf{a} - \mu}{\sigma} \oplus \boldsymbol{\beta}$$

where we have mean $\mu = \frac{1}{K} \sum_{k=1}^{K} a_k$,

standard deviation $\sigma = \sqrt{\frac{1}{K} \sum_{k=1}^{K} (a_k - \mu)^2}$,
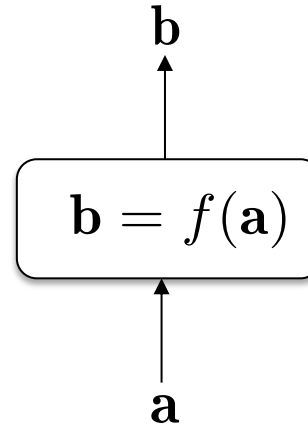
and parameters $\boldsymbol{\gamma} \in \mathbb{R}^K, \boldsymbol{\beta} \in \mathbb{R}^K$.

$\odot$ and $\oplus$ denote elementwise multiplication and addition.
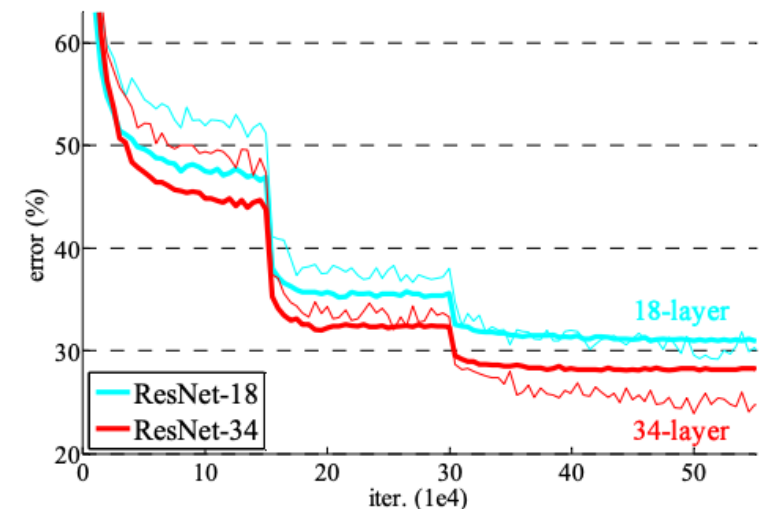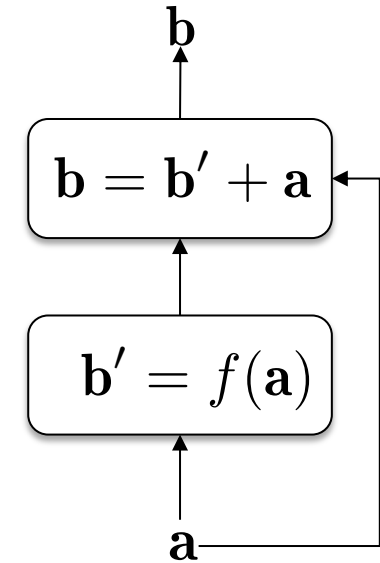
Figure from https://arxiv.org/pdf/1607.06450.pdf

# Residual Connections

- *The Problem:* as network depth grows very large, a **performance degradation** occurs that is not explained by overfitting (i.e. train / test error both worsen)

- *One Solution:* **Residual connections** pass a copy of the input alongsidethe

- These residual connections allow for **effective training of very deep networks** that perform better than their shallower (though still deep) counterparts
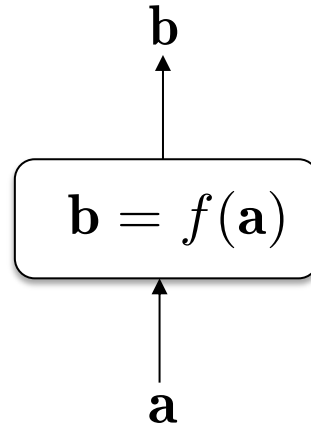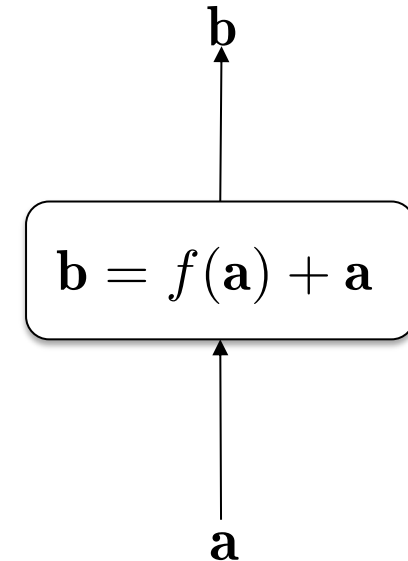
Plain Connection

$$b$$

$$b = f(a)$$

$$a$$

Residual Connection

$$b$$

$$b = b' + a$$

$$b' = f(a)$$

$$a$$

# Residual Connections

- *The Problem:* as network depth grows very large, a **performance degradation** occurs that is not explained by overfitting (i.e. train / test error both worsen)

- *One Solution:* **Residual connections** pass a copy of the input alongsidethe

- These residual connections allow for **effective training of very deep networks** that perform better than their shallower (though still deep) counterparts
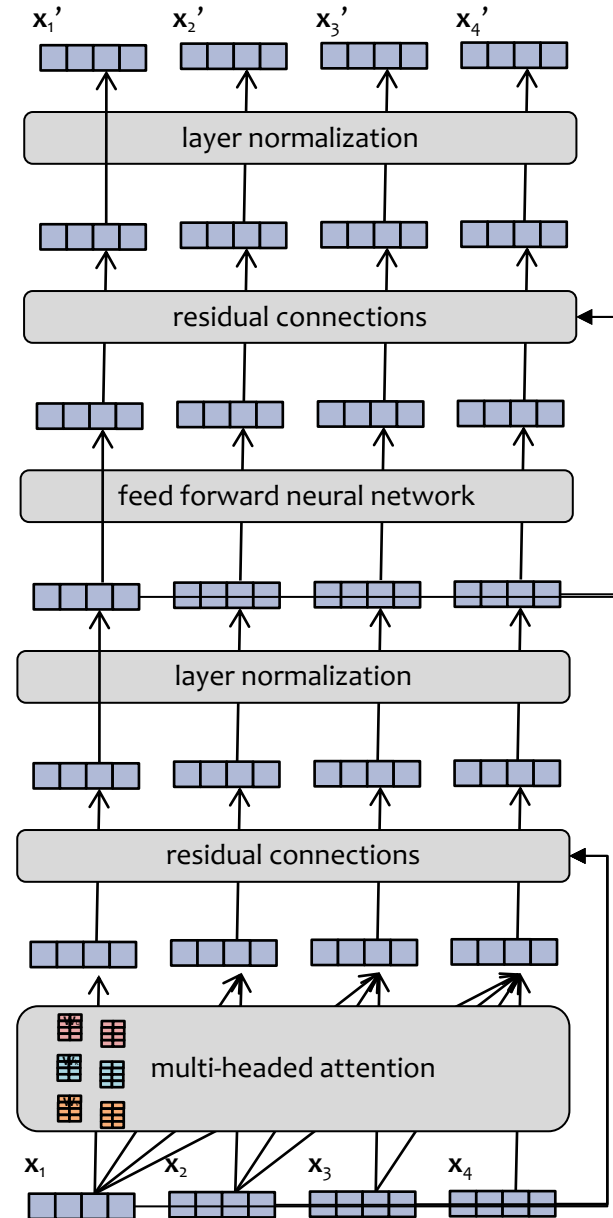
Plain Connection

$$\mathbf{b}$$

$$\mathbf{b} = f(\mathbf{a})$$

$$\mathbf{a}$$

Residual Connection

$$\mathbf{b}$$

$$\mathbf{b} = f(\mathbf{a}) + \mathbf{a}$$
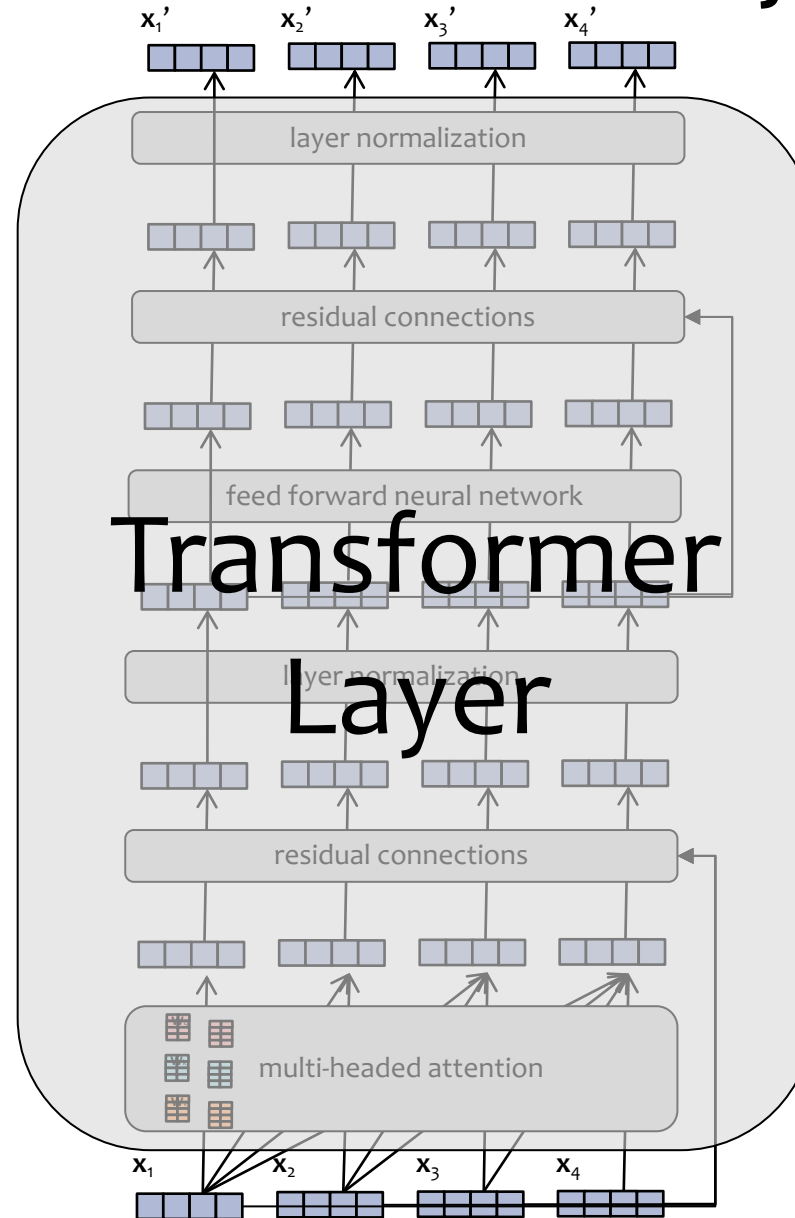
$$\mathbf{a}$$

**Why are residual connections helpful?**

Instead of f(a) having to learn a full transformation of a, f(a) only needs to learn an additive modification of a (i.e. the residual).

# Transformer Layer



**Each layer** of a Transformer LM consists of several **sublayers:**
1. attention
2. feed-forward neural network
3. layer normalization
4. residual connections

# Transformer Layer

$x_1'$  $x_2'$  $x_3'$  $x_4'$

layer normalization

residual connections

feed forward neural network

# Transformer Layer

layer normalization

residual connections

multi-headed attention
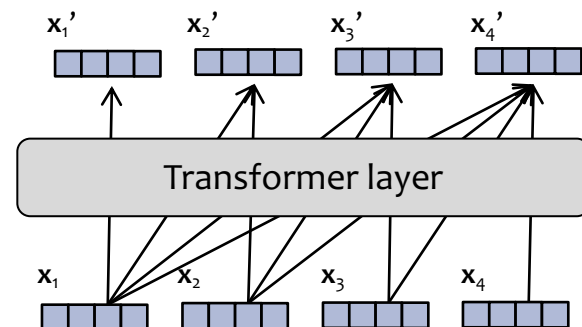
$x_1$  $x_2$  $x_3$  $x_4$

**Each layer** of a Transformer LM consists of several **sublayers:**
1. attention
2. feed-forward neural network
3. layer normalization
4. residual connections

# Transformer Layer

$x_1'$ $x_2'$ $x_3'$ $x_4'$

## Transformer Layer

$x_1$ $x_2$ $x_3$ $x_4$

**Each layer** of a Transformer LM consists of several **sublayers:**

1. attention
2. feed-forward neural network
3. layer normalization
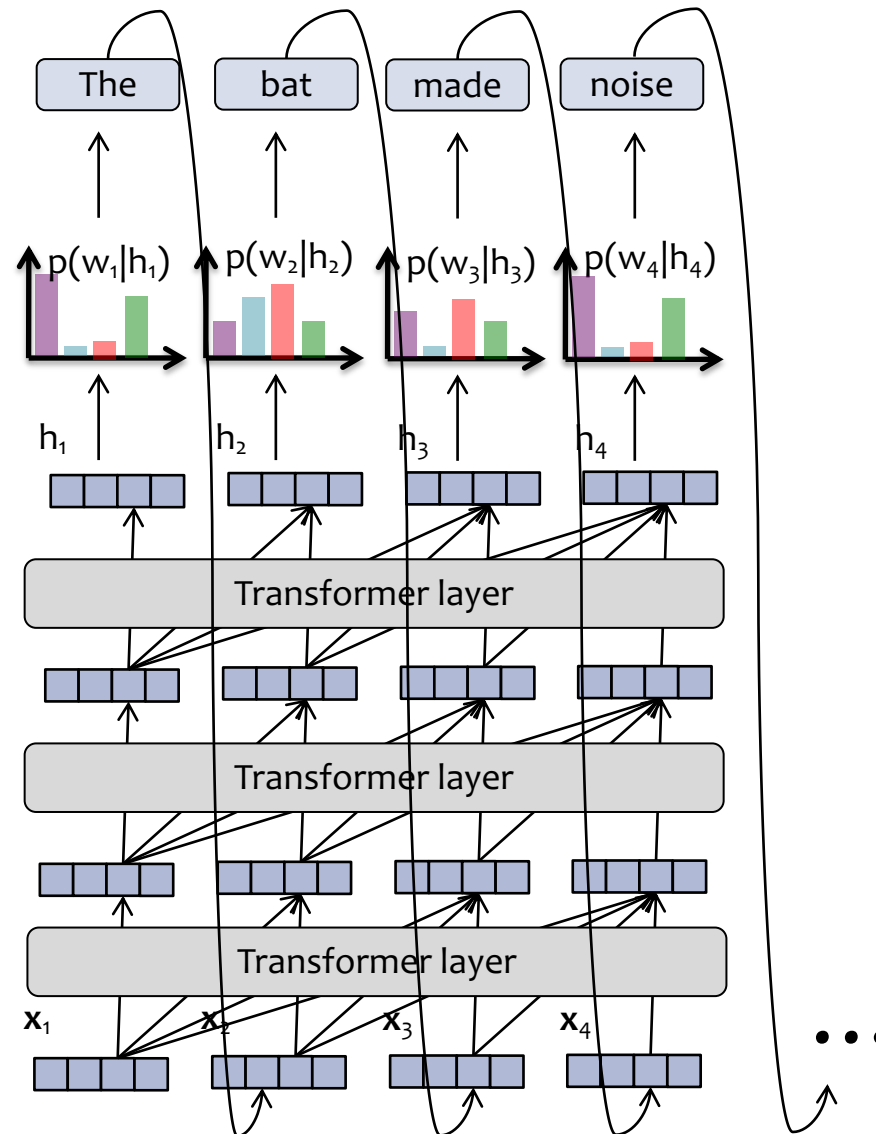4. residual connections

# Transformer Layer

Each layer of a Transformer LM consists of several **sublayers:**
1. attention
2. feed-forward neural network
3. layer normalization
4. residual connections

# Transformer Language Model



**Each layer** of a Transformer LM consists of several **sublayers:**
1. attention
2. feed-forward neural network
3. layer normalization
4. residual connections

Each hidden vector looks back at the hidden vectors of the **current and previous timesteps in the previous layer.**

The language model part is just like an RNN-LM.

# In-Class Poll

$$\mathbf{x}'_4 = \sum_{j=1}^{4} a_{4,j} \mathbf{v}_j$$

**Question:**

Suppose we have the following input embeddings and attention weights:
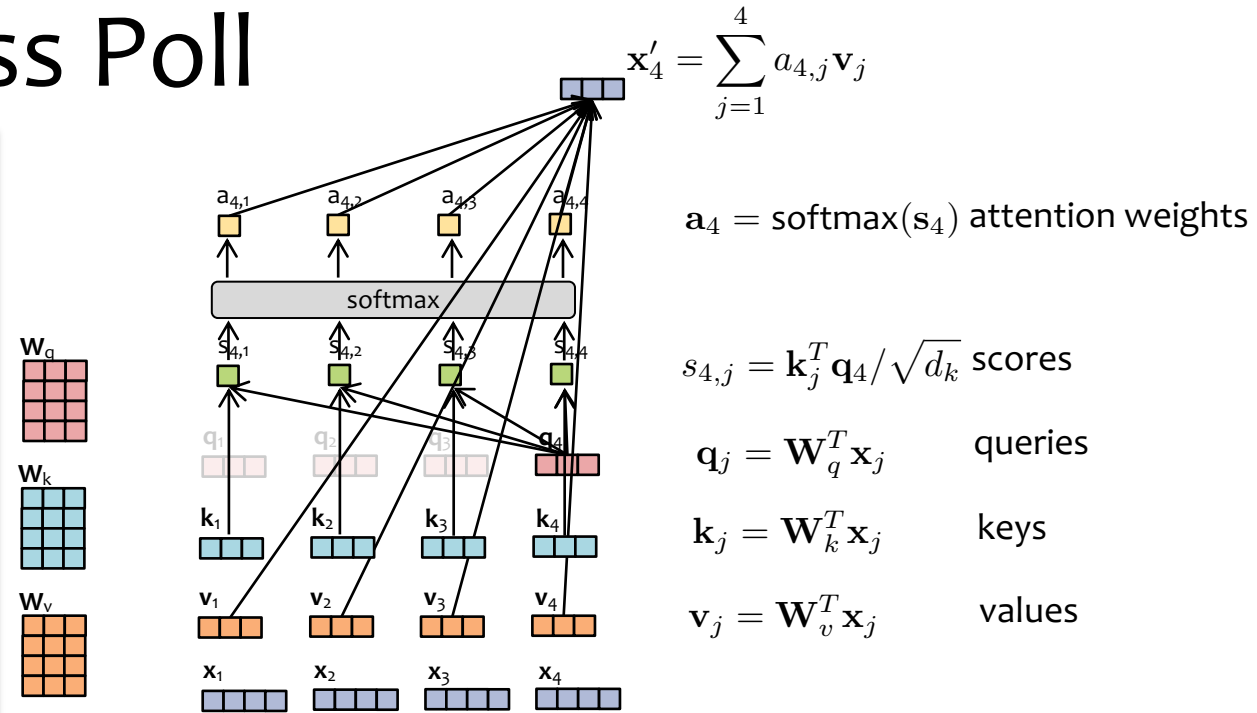
- $x_1 = [1,0,0,0]$  $a_{4,1} = 0.1$
- $x_2 = [0,1,0,0]$  $a_{4,2} = 0.2$
- $x_3 = [0,0,2,0]$  $a_{4,3} = 0.6$
- $x_4 = [0,0,0,1]$  $a_{4,4} = 0.1$

And $W_v = I$. Then we can compute $x_4'$.

Now suppose we swap the embeddings $x_2$ and $x_3$ such that
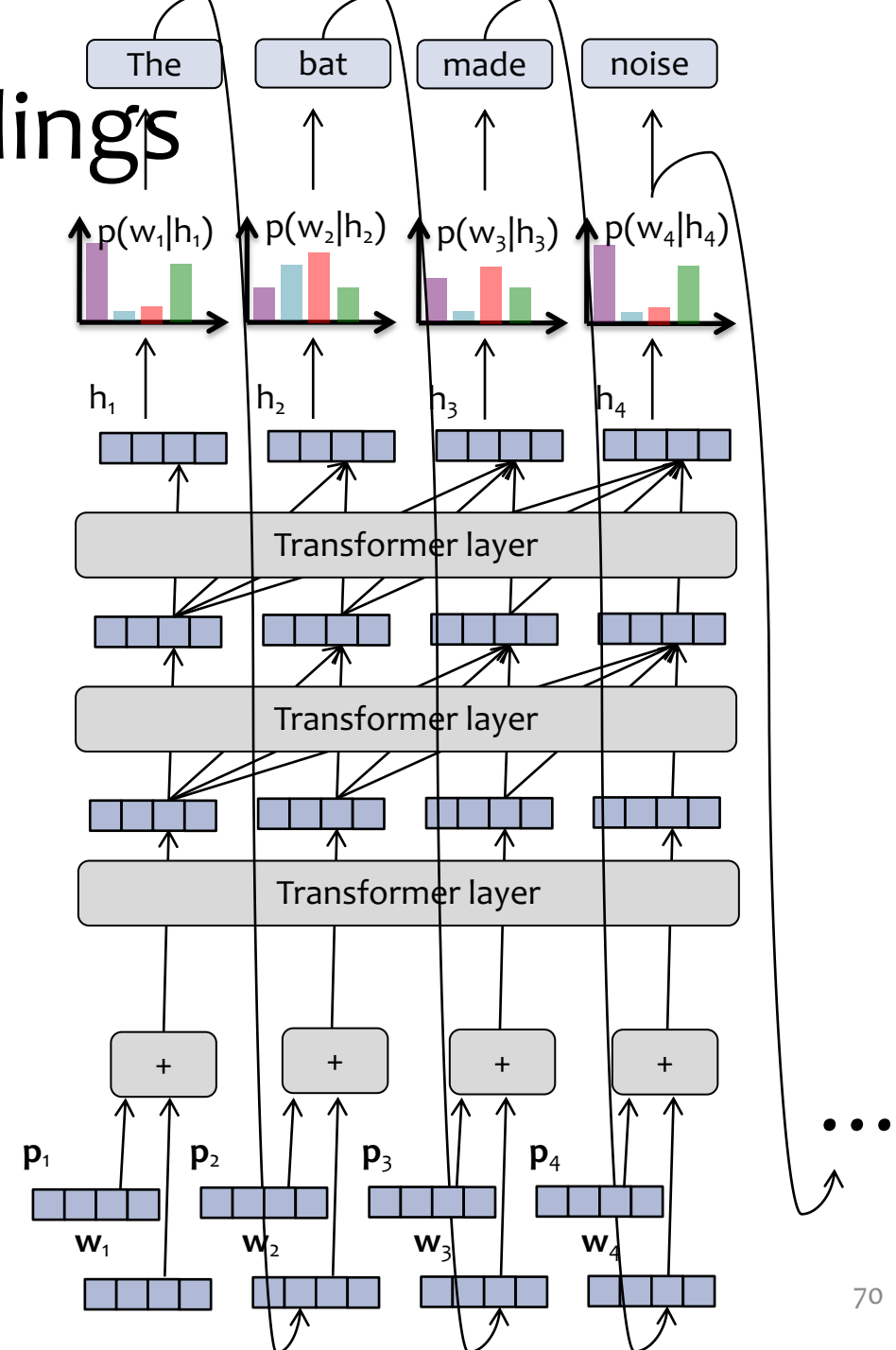
- $x_2 = [0,0,2,0]$
- $x_3 = [0,1,0,0]$

What is the new value of $x_4'$?

$\mathbf{a}_4 = \mathrm{softmax}(\mathbf{s}_4)$ attention weights

$s_{4,j} = \mathbf{k}_j^T \mathbf{q}_4 / \sqrt{d_k}$ scores

$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j$ queries

$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j$ keys

$\mathbf{v}_j = \mathbf{W}_v^T \mathbf{x}_j$ values

softmax

**Answer:**

# Position Embeddings

- **The Problem:** Because attention is position invariant, we **need** a way to learn about positions
- **The Solution:** Use (or learn) a collection of position specific embeddings: $p_t$ represents what it means to be in position t. And add this to the word embedding $w_t$.
  The **key idea** is that every word that appears in position t uses the same position embedding $p_t$
- There are a number of varieties of position embeddings:
  - Some are fixed (based on sine and cosine), whereas others are learned (like word embeddings)
  - Some are absolute (as described above) but we can also use relative position embeddings (i.e. relative to the position of the query vector)
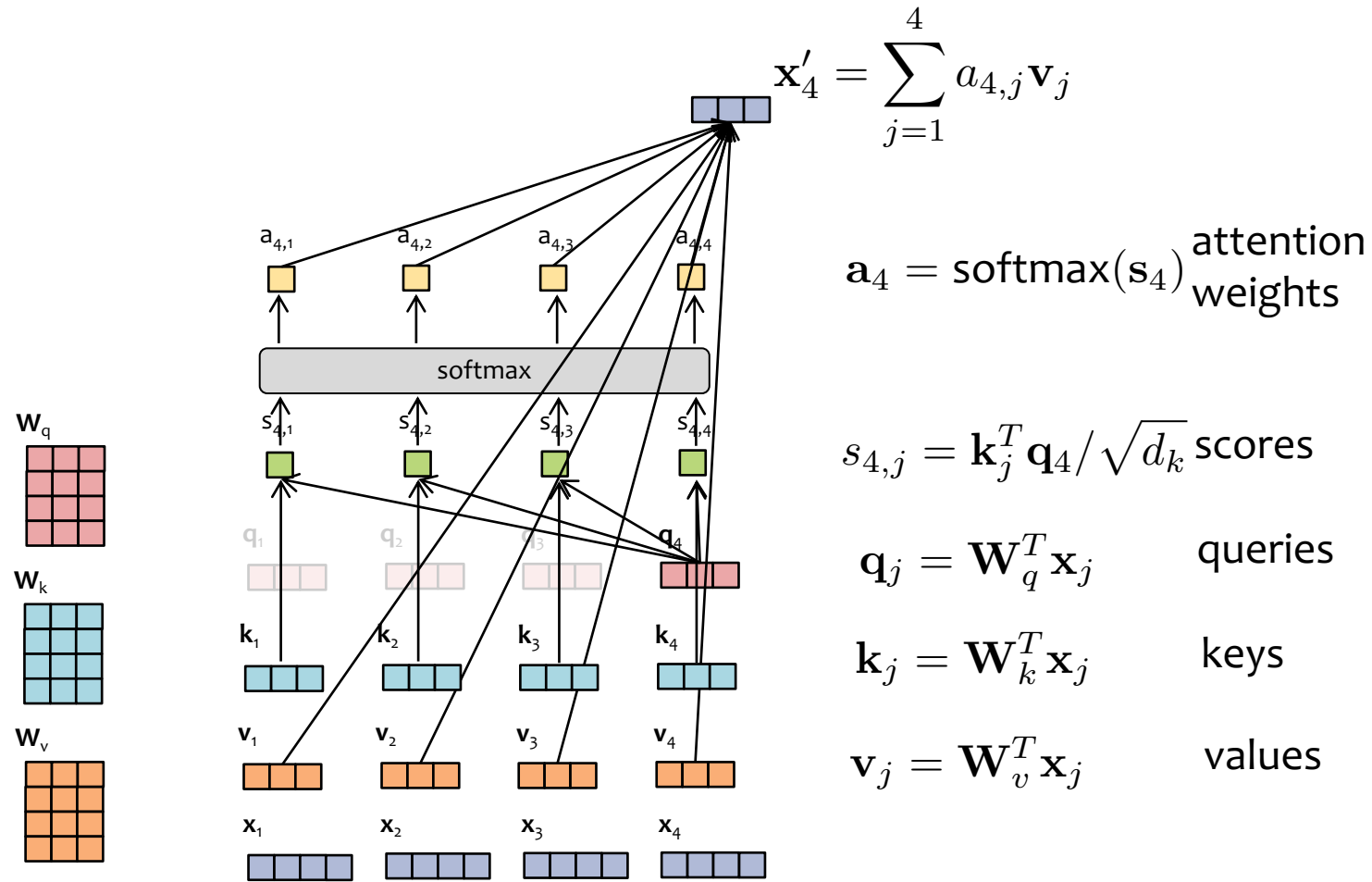
# GPT-3

- GPT stands for Generative Pre-trained Transformer
- GPT is just a Transformer LM, but with a huge number of parameters

| Model | # layers | dimension of states | dimension of inner states | # attention heads | # params |
|-------|----------|---------------------|---------------------------|-------------------|----------|
| GPT (2018) | 12 | 768 | 3072 | 12 | 117M |
| GPT-2 (2019) | 48 | 1600 | -- | -- | 1542M |
| GPT-3 (2020) | 96 | 12288 | 4*12288 | 96 | 175000M |

# Matrix Version of Scaled Dot-Product Attention



$$\mathbf{x}'_4 = \sum_{j=1}^{4} a_{4,j} \mathbf{v}_j$$

$$\mathbf{a}_4 = \text{softmax}(\mathbf{s}_4) \quad \begin{array}{l}\text{attention} \\ \text{weights}\end{array}$$

$$s_{4,j} = \mathbf{k}_j^T \mathbf{q}_4 / \sqrt{d_k} \quad \text{scores}$$

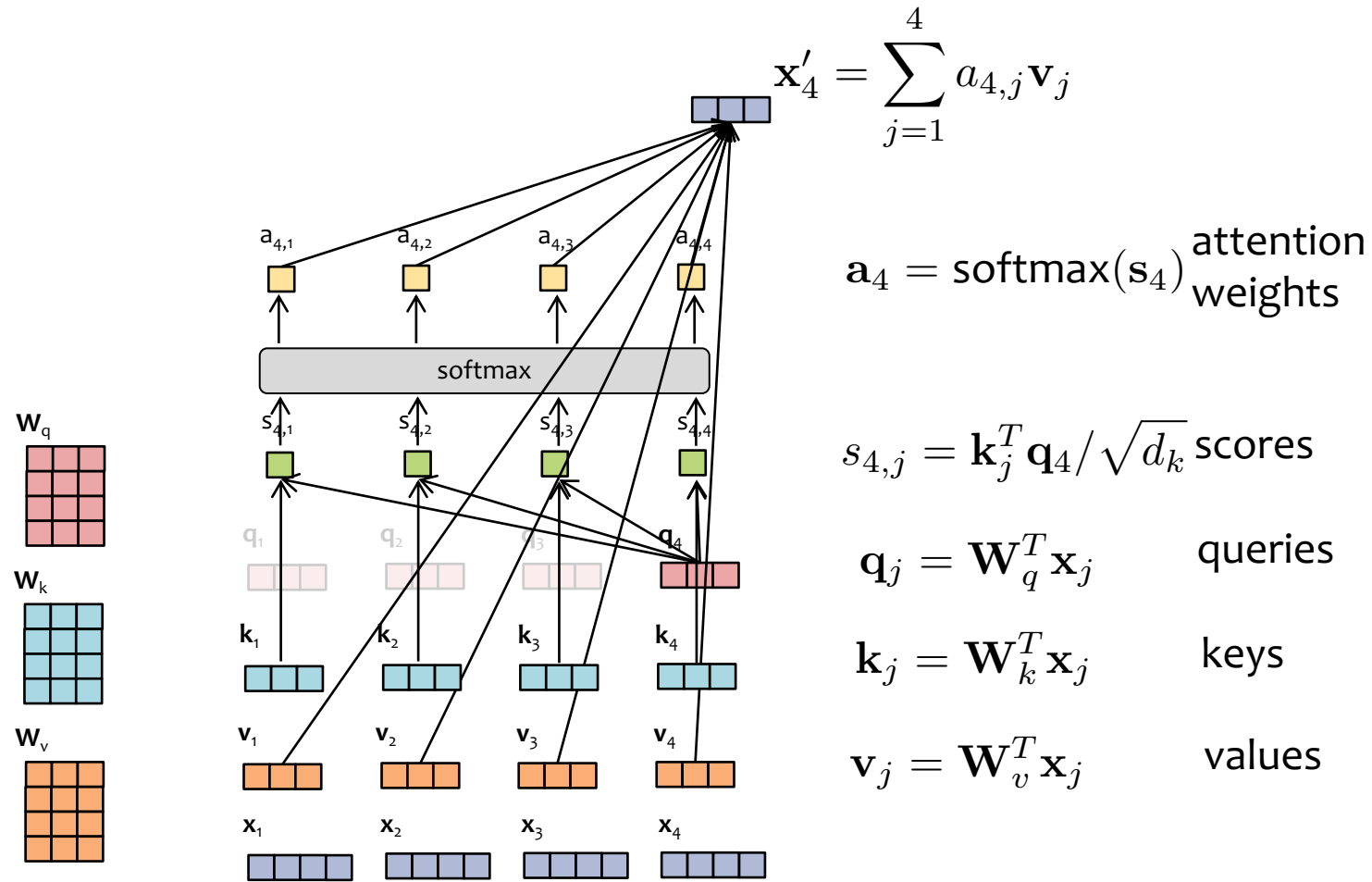$$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j \qquad \text{queries}$$

$$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j \qquad \text{keys}$$

$$\mathbf{v}_j = \mathbf{W}_v^T \mathbf{x}_j \qquad \text{values}$$

- For speed, we compute all the queries at once using matrix operations
- First we pack the queries, keys, values into matrices:
  - $Q = [q_1, \ldots, q_N]^T$
  - $K = [k_1, \ldots, k_N]^T$
  - $V = [v_1, \ldots, v_N]^T$
- Then we compute all the queries at once:

$$\text{Attn}(\mathbf{x}_{1:N}) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}V\right)$$

74

# Matrix Version of Scaled Dot-Product Attention



$$\mathbf{x}'_4 = \sum_{j=1}^{4} a_{4,j} \mathbf{v}_j$$

$$\mathbf{a}_4 = \text{softmax}(\mathbf{s}_4) \quad \begin{array}{l}\text{attention}\\ \text{weights}\end{array}$$

$$s_{4,j} = \mathbf{k}_j^T \mathbf{q}_4 / \sqrt{d_k} \quad \text{scores}$$

$$\mathbf{q}_j = \mathbf{W}_q^T \mathbf{x}_j \quad \text{queries}$$

$$\mathbf{k}_j = \mathbf{W}_k^T \mathbf{x}_j \quad \text{keys}$$

$$\mathbf{v}_j = \mathbf{W}_v^T \mathbf{x}_j \quad \text{values}$$

- For speed, we compute all the queries at once using matrix operations
- First we pack the queries, keys, values into matrices:
  - $Q = [q_1, \ldots, q_N]^T$
  - $K = [k_1, \ldots, k_N]^T$
  - $V = [v_1, \ldots, v_N]^T$
- Then we compute all the queries at once:

$$\text{Attn}(\mathbf{x}_{1:N}) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}} V\right)$$

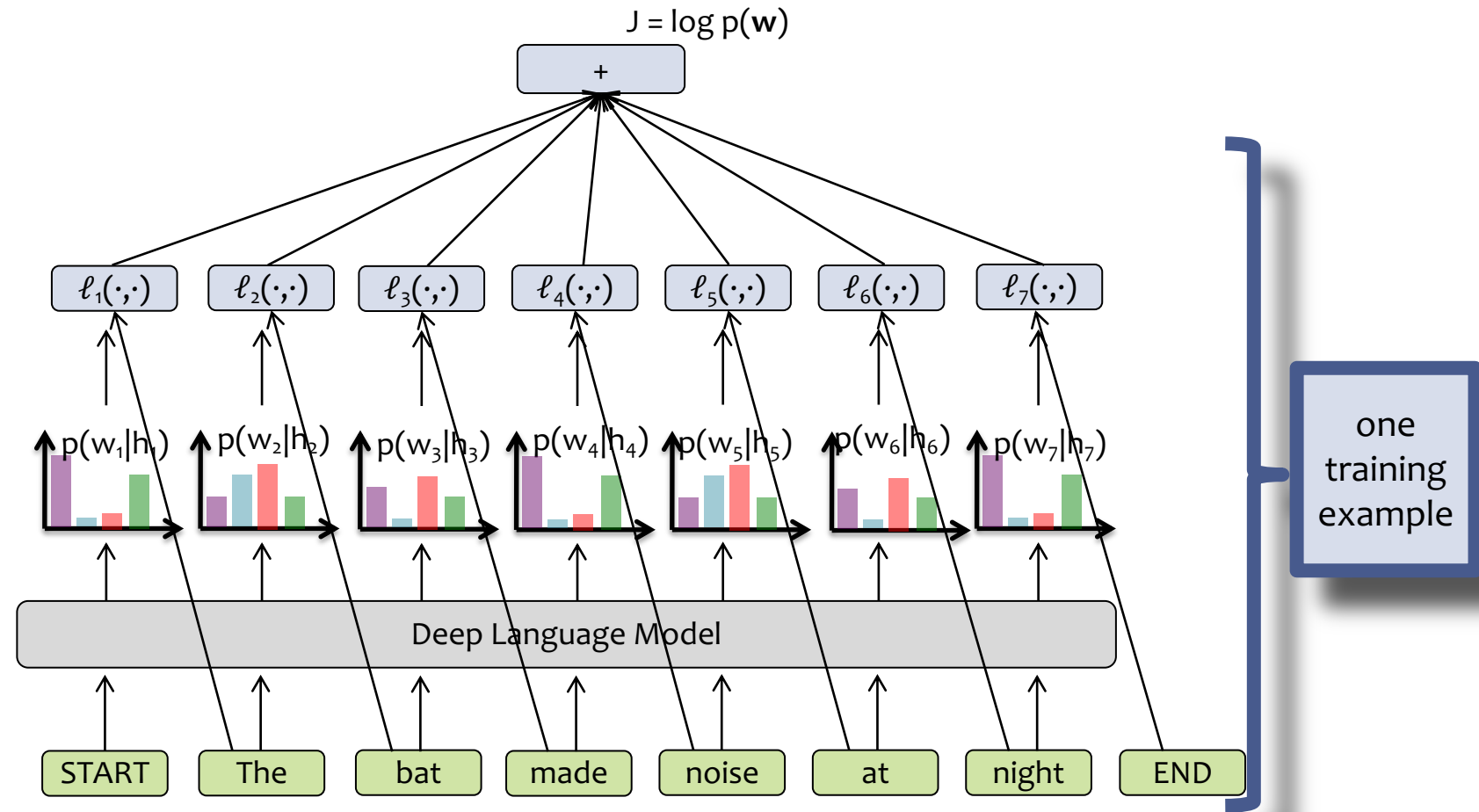In practice, the attention weights are computed for all time steps T, then we mask out (by setting to –inf) all the inputs to the softmax that are for the timesteps to the right of the query.

# LEARNING A TRANSFORMER LM

# Learning a Deep Language Model

- Each training example is a sequence (e.g. sentence), so we have training data $D = \{\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \dots, \mathbf{w}^{(N)}\}$

- The objective function for a Deep LM (e.g. RNN-LM or Tranformer-LM) is typically the log-likelihood of the training examples:

$$J(\boldsymbol{\theta}) = \Sigma_i \log p_{\boldsymbol{\theta}}(\mathbf{w}^{(i)})$$

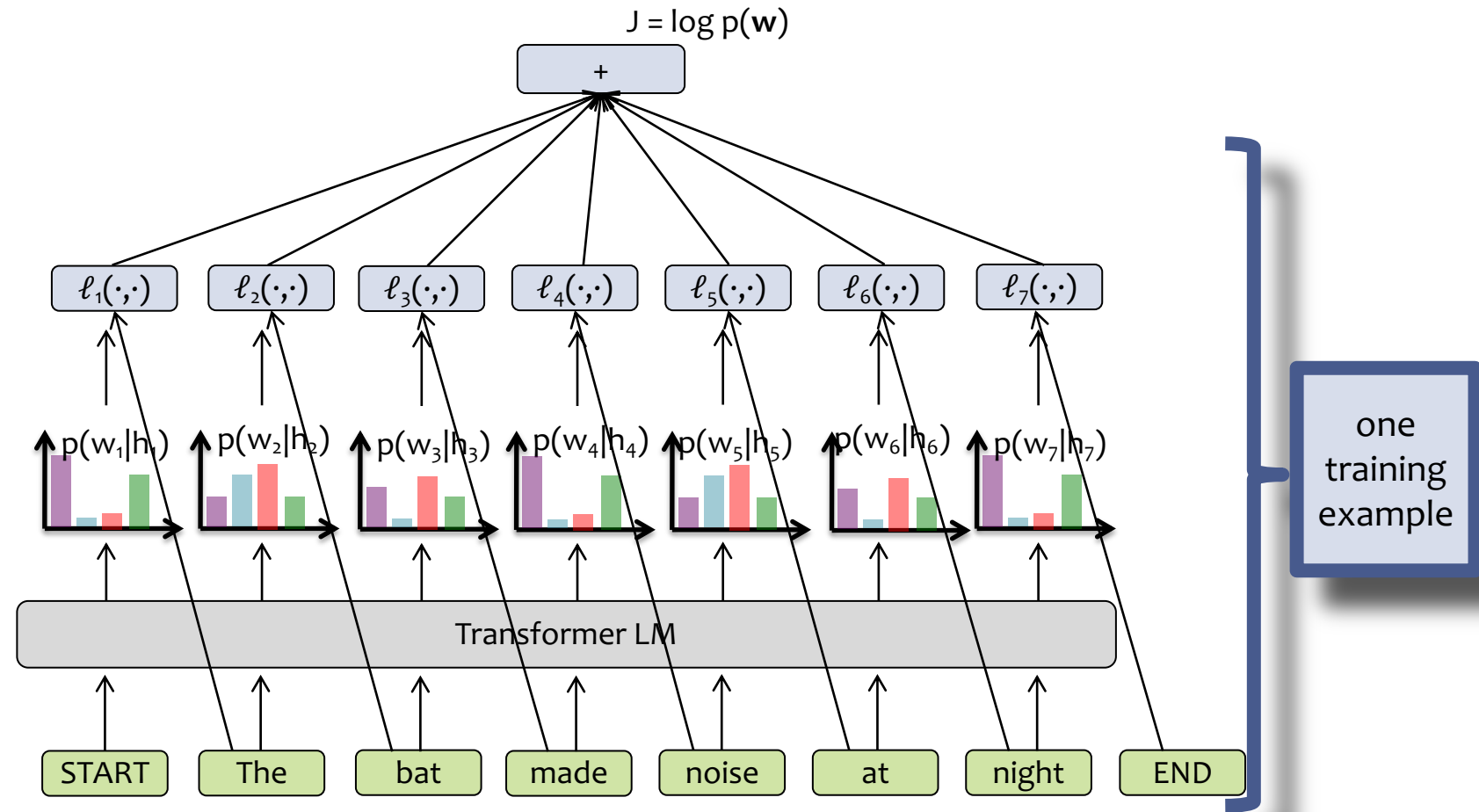- We train by mini-batch SGD (or your favorite flavor of mini-batch SGD)

$$\log p(\mathbf{w}) = \log p(w_1, w_2, w_3, \dots, w_T)$$
$$= \log p(w_1 \mid h_1) + \log p(w_2 \mid h_2) + \dots + \log p(w_2 \mid h_T)$$

# Learning a Transformer Language Model

- Each training example is a sequence (e.g. sentence), so we have training data $D = \{\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \dots, \mathbf{w}^{(N)}\}$

- The objective function for a Deep LM (e.g. RNN-LM or Tranformer-LM) is typically the log-likelihood of the training examples:

$$J(\boldsymbol{\theta}) = \Sigma_i \log p_{\boldsymbol{\theta}}(\mathbf{w}^{(i)})$$

- We train by mini-batch SGD (or your favorite flavor of mini-batch SGD)



$\log p(\mathbf{w}) = \log p(w_1, w_2, w_3, \dots, w_T)$
$\qquad = \log p(w_1 \mid h_1) + \log p(w_2 \mid h_2) + \dots + \log p(w_2 \mid h_T)$

$J = \log p(\mathbf{w})$

# Language Modeling

**An aside:**

- State-of-the-art language models currently tend to rely on **transformer networks** (e.g. GPT-2)
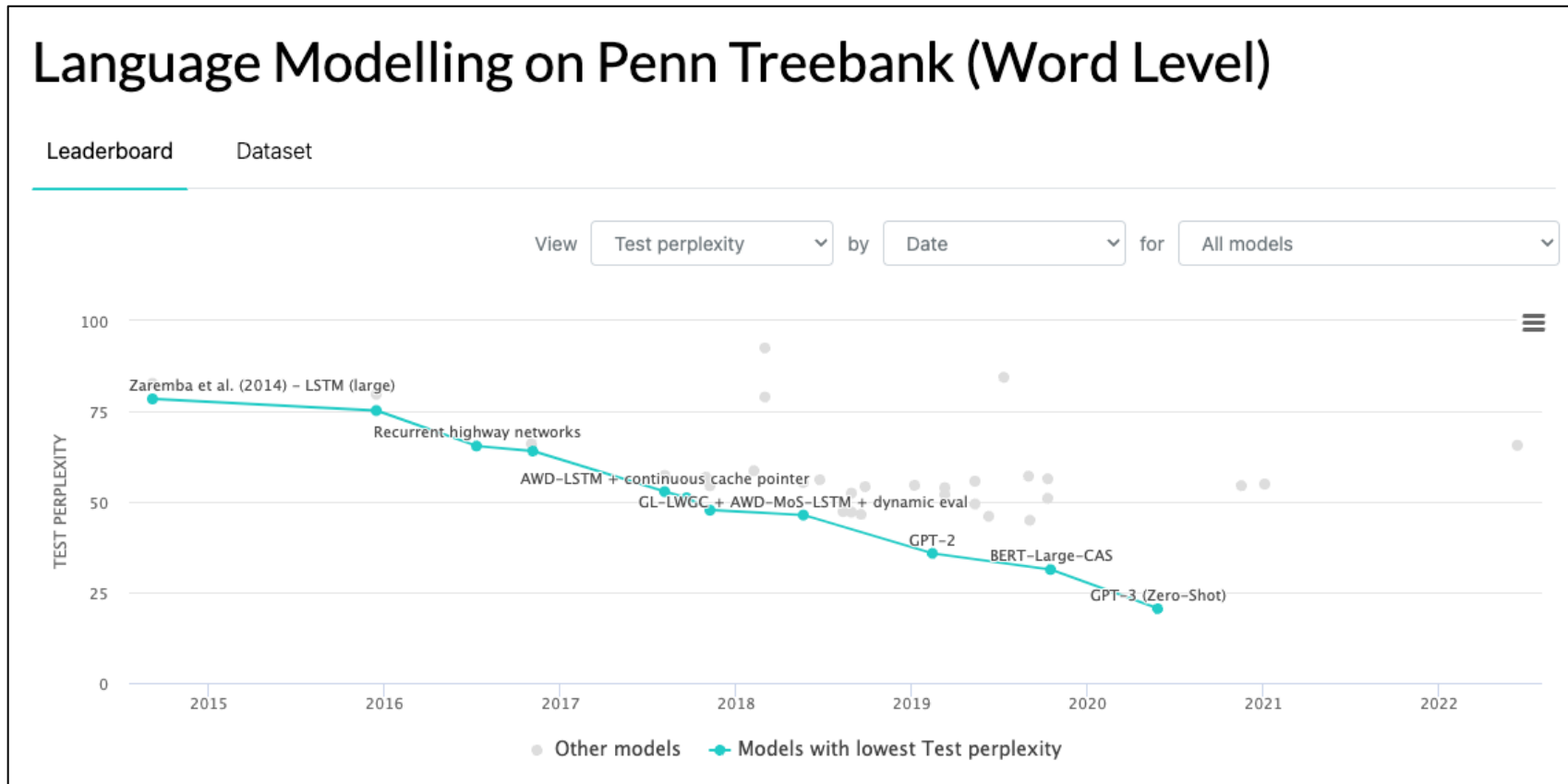- RNN-LMs comprised most of the early neural LMs that **led to** current SOTA architectures

Figure from https://paperswithcode.com/sota/language-modelling-on-penn-treebank-word

# Why does efficiency matter?

**Case Study: GPT-3**

- # of training tokens = 500 billion

- # of parameters = 175 billion

- # of cycles = 50 petaflop/s-days (each of which are 8.64e+19 flops)

| Dataset | Quantity (tokens) | Weight in training mix | Epochs elapsed when training for 300B tokens |
|---|---|---|---|
| Common Crawl (filtered) | 410 billion | 60% | 0.44 |
| WebText2 | 19 billion | 22% | 2.9 |
| Books1 | 12 billion | 8% | 1.9 |
| Books2 | 55 billion | 8% | 0.43 |
| Wikipedia | 3 billion | 3% | 3.4 |

**Table 2.2: Datasets used to train GPT-3**. "Weight in training mix" refers to the fraction of examples during training that are drawn from a given dataset, which we intentionally do not make proportional to the size of the dataset. As a result, when we train for 300 billion tokens, some datasets are seen up to 3.4 times during training while other datasets are seen less than once.

| Model Name | $n_{params}$ | $n_{layers}$ | $d_{model}$ | $n_{heads}$ | $d_{head}$ | Batch Size | Learning Rate |
|---|---|---|---|---|---|---|---|
| GPT-3 Small | 125M | 12 | 768 | 12 | 64 | 0.5M | $6.0 \times 10^{-4}$ |
| GPT-3 Medium | 350M | 24 | 1024 | 16 | 64 | 0.5M | $3.0 \times 10^{-4}$ |
| GPT-3 Large | 760M | 24 | 1536 | 16 | 96 | 0.5M | $2.5 \times 10^{-4}$ |
| GPT-3 XL | 1.3B | 24 | 2048 | 24 | 128 | 1M | $2.0 \times 10^{-4}$ |
| GPT-3 2.7B | 2.7B | 32 | 2560 | 32 | 80 | 1M | $1.6 \times 10^{-4}$ |
| GPT-3 6.7B | 6.7B | 32 | 4096 | 32 | 128 | 2M | $1.2 \times 10^{-4}$ |
| GPT-3 13B | 13.0B | 40 | 5140 | 40 | 128 | 2M | $1.0 \times 10^{-4}$ |
| GPT-3 175B or "GPT-3" | 175.0B | 96 | 12288 | 96 | 128 | 3.2M | $0.6 \times 10^{-4}$ |

**Table 2.1:** Sizes, architectures, and learning hyper-parameters (batch size in tokens and learning rate) of the models which we trained. All models were trained for a total of 300 billion tokens.
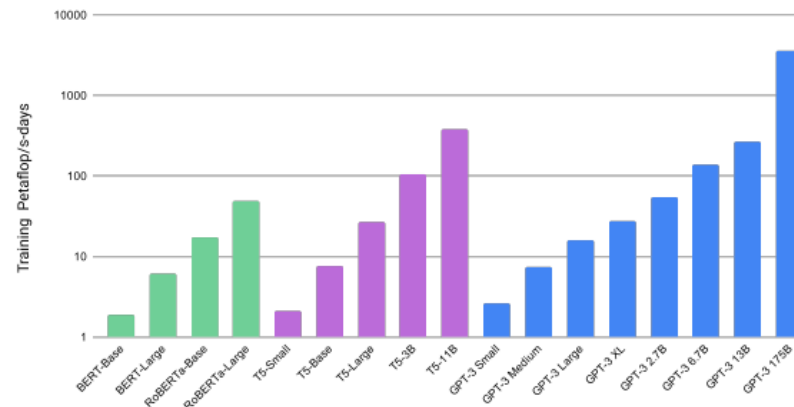


**Figure 2.2: Total compute used during training**. Based on the analysis in Scaling Laws For Neural Language Models [KMH+20] we train much larger models on many fewer tokens than is typical. As a consequence, although GPT-3 3B is almost 10x larger than RoBERTa-Large (355M params), both models took roughly 50 petaflop/s-days of compute during pre-training. Methodology for these calculations can be found in Appendix D.

Figure from https://arxiv.org/pdf/2005.14165

# Summary

- Task: Language Modeling
  - noisy channel models (speech / MT)
  - (historical) Large LMs (n-gram models)
- Model: GPT
  - Attention (computation graph)
  - Transformer-LM (cf. RNN-LM)
- Learning for LLMs
  - Pre-training (unsupervised learning)
  - Reinforcement Learning with Human Feedback (deep RL)
- Optimization for LLMs
  - Adam (cf. SGD)
  - Distributed training
- Societal Impacts of LLMs