# Topics in AI (CPSC 532S):
## Multimodal Learning with Vision, Language and Sound

**Lecture 8: Word2Vec, Language Models and RNNs**

# Course **Logistics**

— **Assignment 3**

— **Final project group** Goolge form will be out **tomorrow**

# Representing a **Word:** One Hot Encoding

| **Vocabulary** | | **one-hot** encodings |
|---|---|---|
| dog | 1 | [ 1, 0, 0, 0, 0, 0, 0, 0, 0, 0 ] |
| cat | 2 | [ 0, 1, 0, 0, 0, 0, 0, 0, 0, 0 ] |
| person | 3 | [ 0, 0, 1, 0, 0, 0, 0, 0, 0, 0 ] |
| holding | 4 | [ 0, 0, 0, 1, 0, 0, 0, 0, 0, 0 ] |
| tree | 5 | [ 0, 0, 0, 0, 1, 0, 0, 0, 0, 0 ] |
| computer | 6 | [ 0, 0, 0, 0, 0, 1, 0, 0, 0, 0 ] |
| using | 7 | [ 0, 0, 0, 0, 0, 0, 1, 0, 0, 0 ] |

# Representing **Phrases**: Bag-of-Words

**bag-of-words** representation

person holding dog {3, 4, 1} [ 1, 0, 1, 1, 0, 0, 0, 0, 0, 0 ]

person holding cat {3, 4, 2} [ 1, 1, 0, 1, 0, 0, 0, 0, 0, 0 ]

person using computer {3, 7, 6} [ 0, 0, 0, 1, 0, 1, 1, 0, 0, 0 ]

dog cat person holding tree computer using

person using computer
person holding cat {3, 3, 7, 6, 2} [ 0, 1, 2, 1, 0, 1, 1, 0, 0, 0 ]

*slide from V. Ordonex

# **Distributional** Hypothesis    <span>[ Lenci, 2008 ]</span>

— At least certain aspects of the meaning of lexical expressions depend on their distributional properties in the linguistic contexts

— The degree of semantic similarity between two linguistic expressions is a function of the similarity of the two linguistic contexts in which they can appear

# What is the meaning of "**bardiwac**"?

— He handed her glass of **bardiwac**.

— Beef dishes are made to complement the **bardiwacs**.

— Nigel staggered to his feet, face flushed from too much **bardiwac**.

— Malbec, one of the lesser-known **bardiwac** grapes, responds well to Australia's sunshine.

— I dined off bread and cheese and this excellent **bardiwac**.

—The drinks were delicious: blood-red **bardiwac** as well as light, sweet Rhenish.

**bardic** is an alcoholic beverage made from grapes

# The **Use Theory** of Meaning

"If you can understand and predict in which context a word will appear in, then you understood the meaning of the word"  [Paul Horwich]

# **Geometric Interpretation**: Co-occurrence as feature

— Row vector describes usage of word in a corpus of text

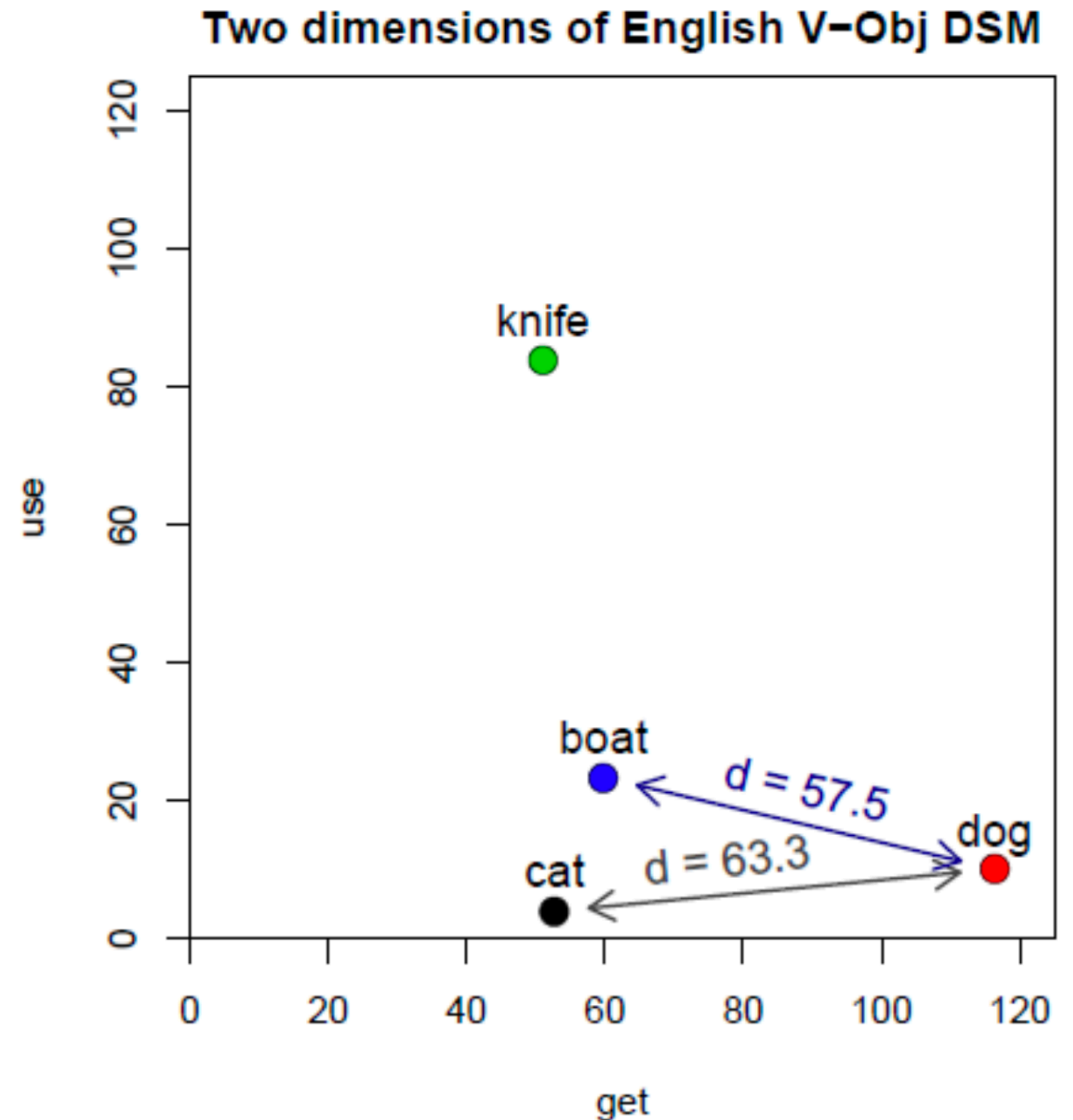— Can be seen as coordinates o the point in an n-dimensional Euclidian space

|        | get | see | use | hear | eat | kill |
|--------|-----|-----|-----|------|-----|------|
| knife  | 51  | 20  | 84  | 0    | 3   | 0    |
| cat    | 52  | 58  | 4   | 4    | 6   | 26   |
| dog    | 115 | 83  | 10  | 42   | 33  | 17   |
| boat   | 59  | 39  | 23  | 4    | 0   | 0    |
| cup    | 98  | 14  | 6   | 2    | 1   | 0    |
| pig    | 12  | 17  | 3   | 2    | 9   | 27   |
| banana | 11  | 2   | 2   | 0    | 18  | 0    |

**Co-occurrence** Matrix

# **Distance** and Similarity

— Illustrated in two dimensions

— Similarity = spatial proximity (Euclidian distance)

— Location depends on frequency of noun (dog is 27 times as frequent as ca)



* Slides from Louis-Philippe Morency

# **Angle** and Similarity

— direction is more important than location

— normalize length of vectors

— or use angle as a distance measure



Two dimensions of English V-Obj DSM

# **Geometric Interpretation**: Co-occurrence as feature

— Row vector describes usage of word in a corpus of text

— Can be seen as coordinates of the point in an n-dimensional Euclidian space

**Way too high dimensional!**

|        | get | see | use | hear | eat | kill |
|--------|-----|-----|-----|------|-----|------|
| knife  | 51  | 20  | 84  | 0    | 3   | 0    |
| cat    | 52  | 58  | 4   | 4    | 6   | 26   |
| dog    | 115 | 83  | 10  | 42   | 33  | 17   |
| boat   | 59  | 39  | 23  | 4    | 0   | 0    |
| cup    | 98  | 14  | 6   | 2    | 1   | 0    |
| pig    | 12  | 17  | 3   | 2    | 9   | 27   |
| banana | 11  | 2   | 2   | 0    | 18  | 0    |

**Co-occurrence** Matrix

# **SVD** for Dimensionality Reduction

# **Learned** Word Vector Visualization

We can also use other methods, like LLE here:



Nonlinear dimensionality reduction by locally linear embedding. Sam Roweis & Lawrence Saul. Science, v.290,2000

[ Roweis and Saul, 2000 ]

# Issues with **SVD**

**Computational** cost for a $d \times n$ matrix is $\mathcal{O}(dn^2)$, where $d < n$

— Makes it not possible for large number of word vocabularies or documents

It is hard to incorporate out of sample (**new**) words or documents

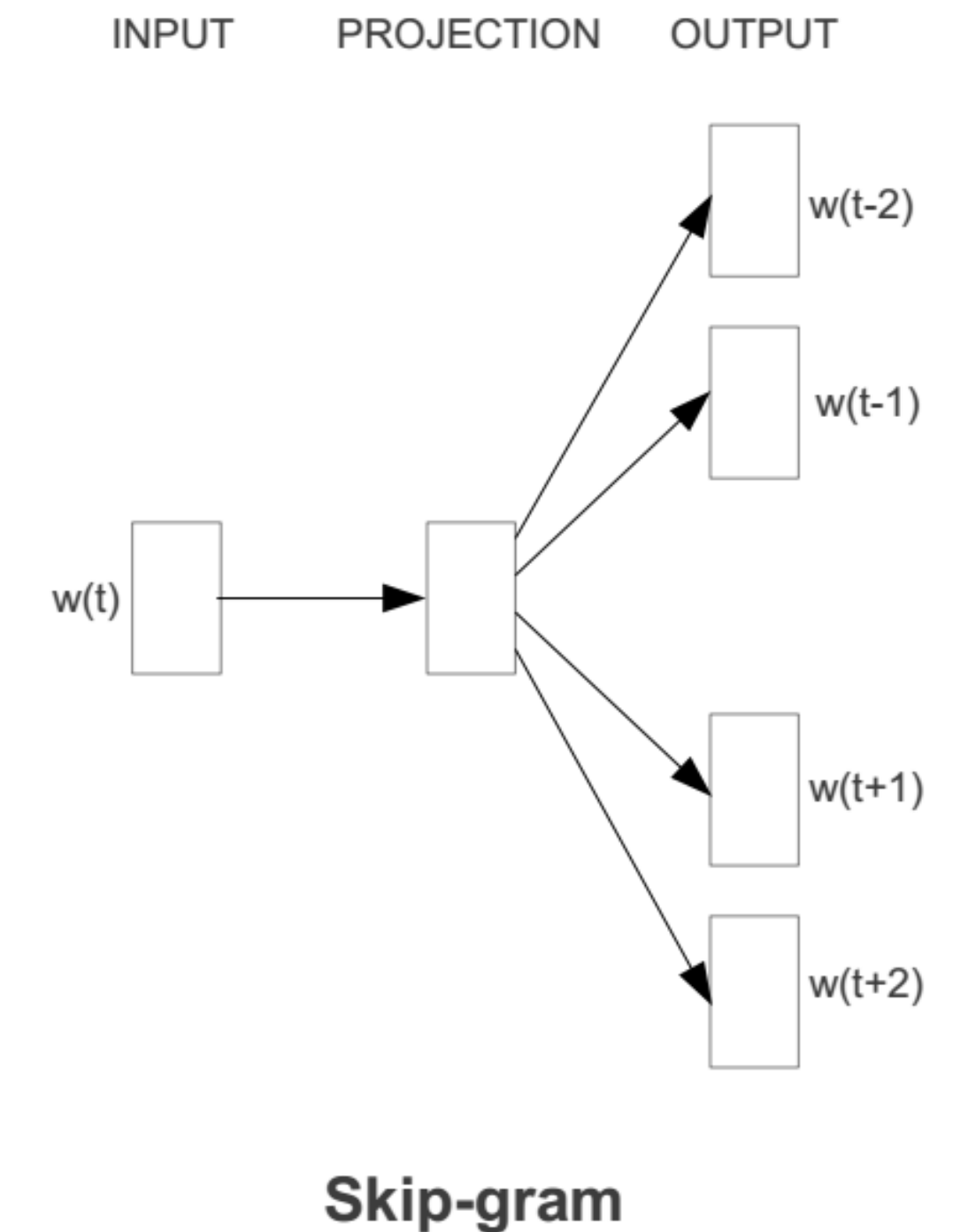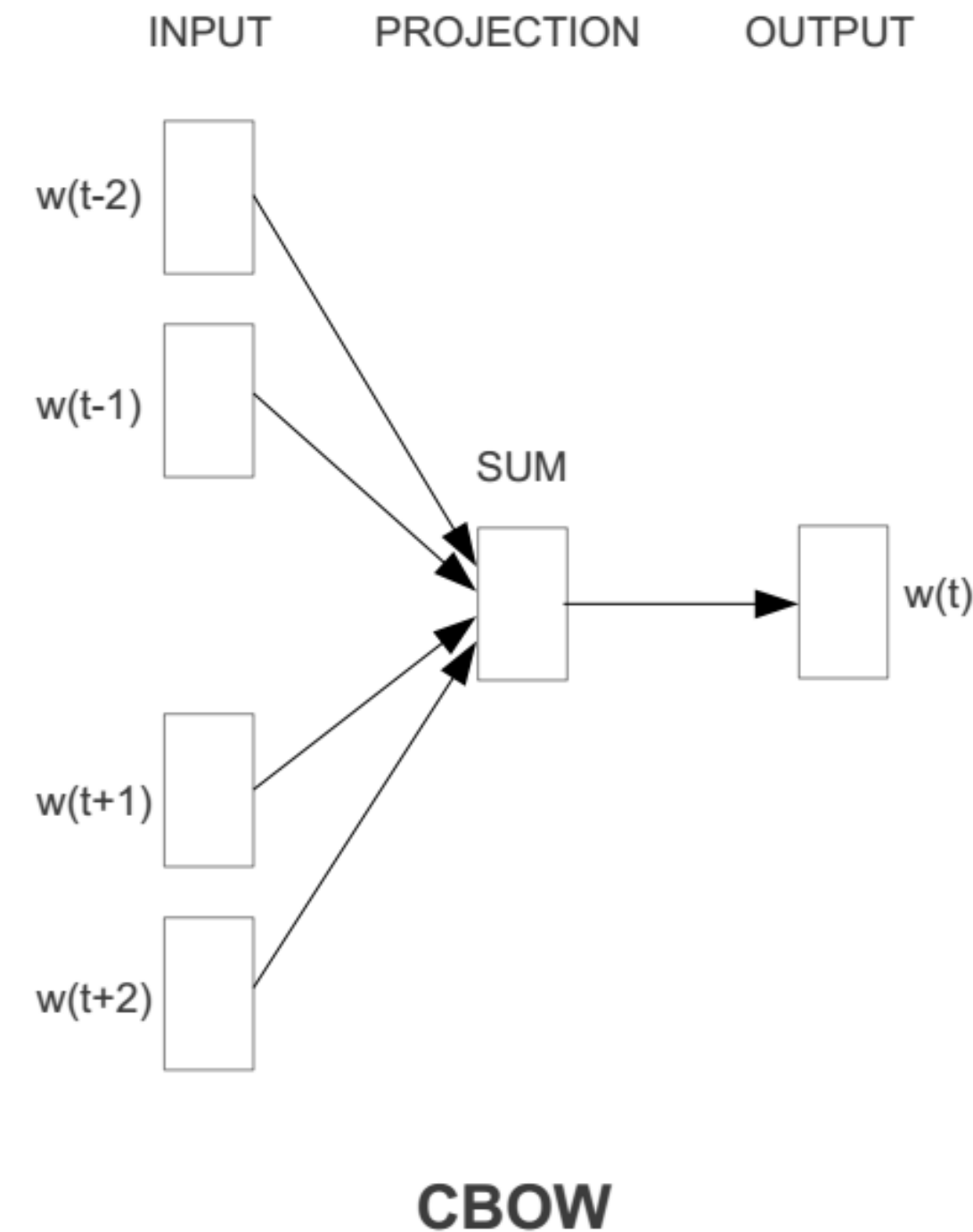# **word2vec**: Representing the Meaning of Words  [ Mikolov et al., 2013 ]

**Key idea:** Predict surrounding words
of every word

**Benefits:** Faster and easier to
incorporate new document, words, etc.

# word2vec: Representing the Meaning of Words   [ Mikolov et al., 2013 ]

**Key idea:** Predict surrounding words
of every word

**Benefits:** Faster and easier to
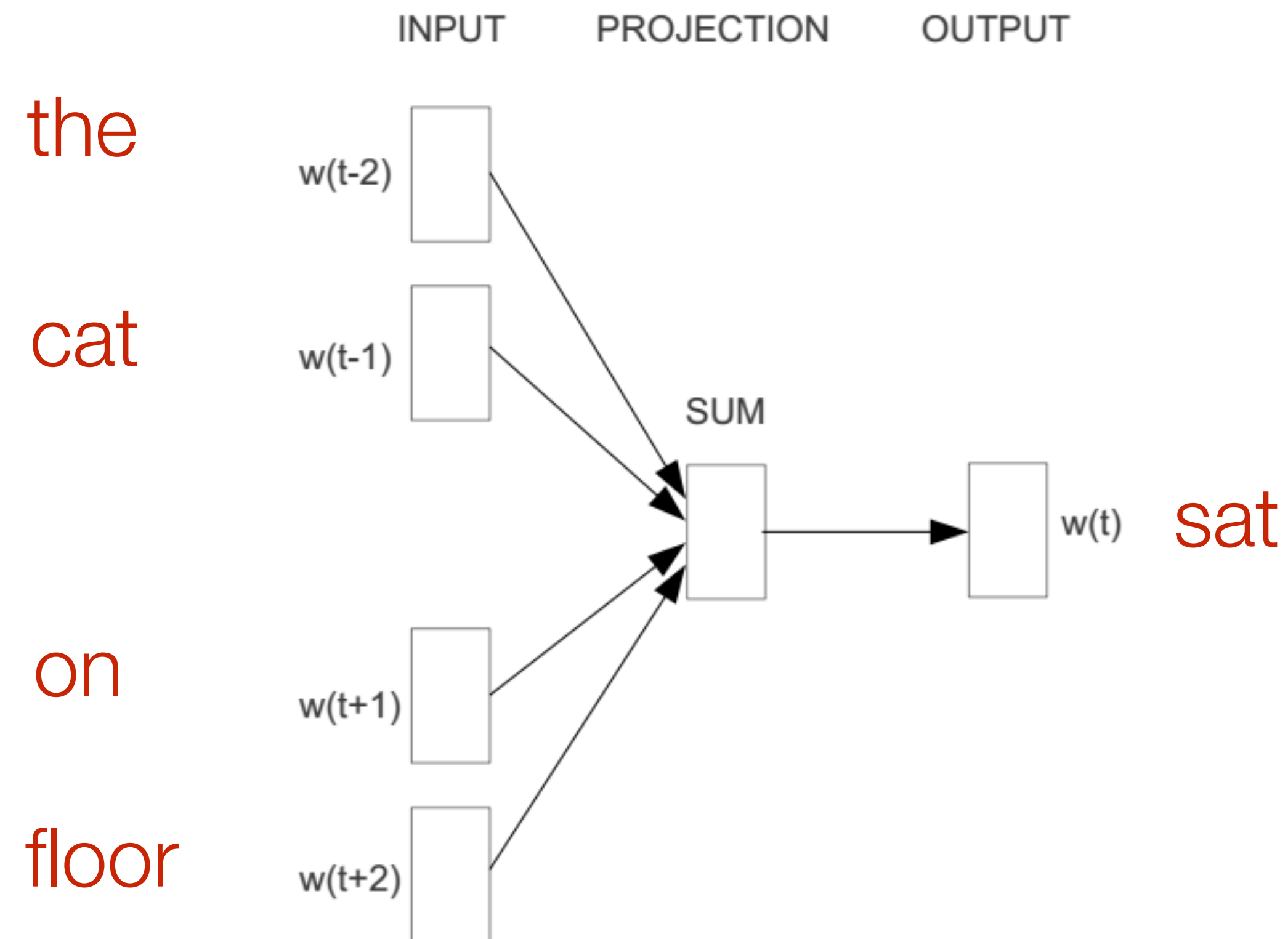incorporate new document, words, etc.



Continuous Bag of Words (**CBOW**): use context words in a window to predict
middle word

**Skip-gram:** use the middle word to predict surrounding ones in a window

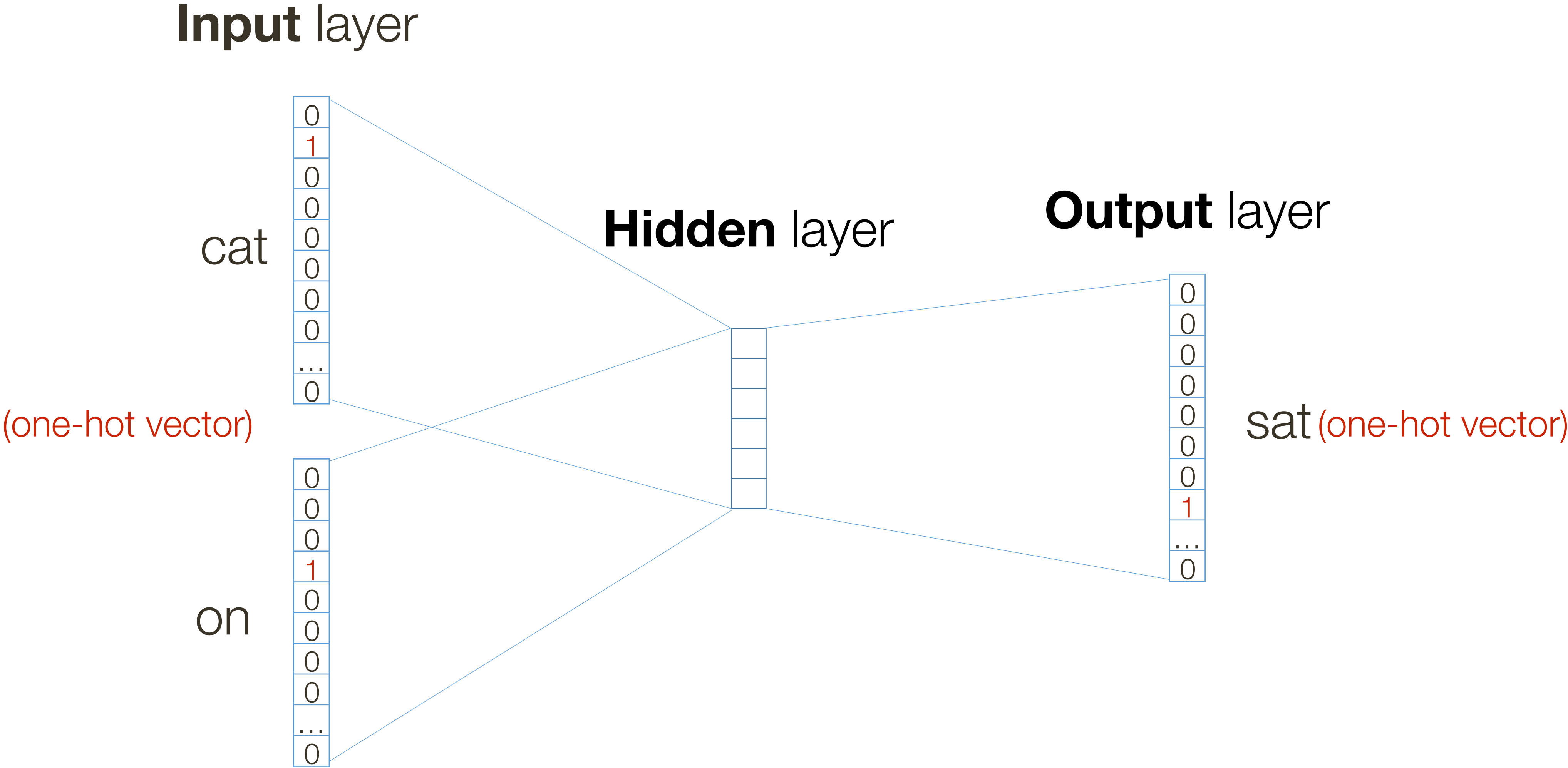*slide from Vagelis Hristidis

# **CBOW**: Continuous Bag of Words

**Example:** "The cat sat on floor" (window size 2)



the

cat

on

floor

sat

# **CBOW**: Continuous Bag of Words

[ Mikolov et al., 2013 ]



*slide from Vagelis Hristidis
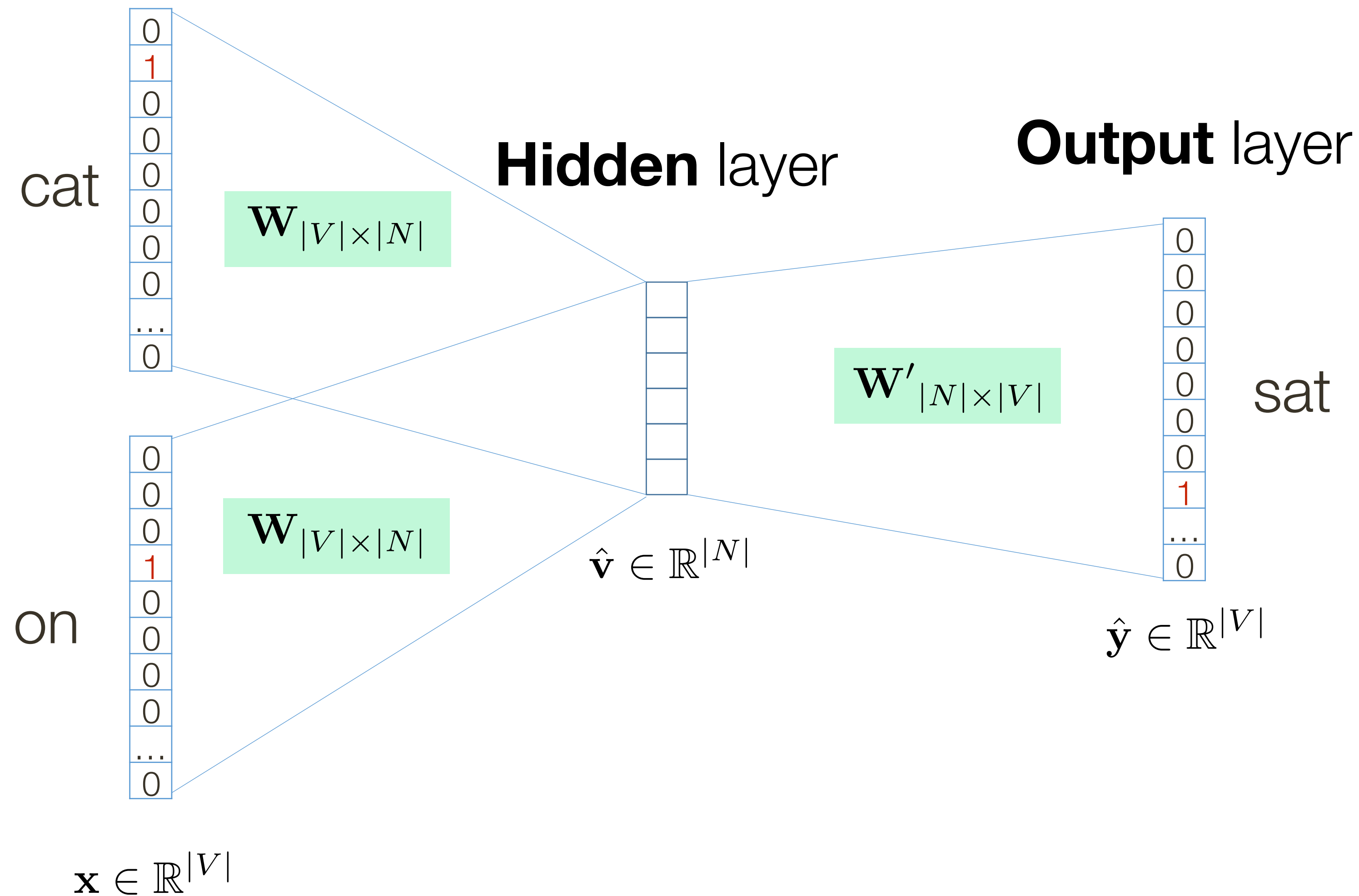
# CBOW: Continuous Bag of Words

**Input** layer

cat

on

$\mathbf{x} \in \mathbb{R}^{|V|}$

$\mathbf{W}_{|V| \times |N|}$

**Hidden** layer

$\hat{\mathbf{v}} \in \mathbb{R}^{|N|}$

**Output** layer

$\mathbf{W'}_{|N| \times |V|}$

sat

$\hat{\mathbf{y}} \in \mathbb{R}^{|V|}$

*slide from Vagelis Hristidis

# **CBOW**: Continuous Bag of Words

**Input** layer

**Parameters** to be learned

cat

$\mathbf{W}_{|V| \times |N|}$

**Hidden** layer

**Output** layer

$\mathbf{W}'_{|N| \times |V|}$

sat

$\hat{\mathbf{v}} \in \mathbb{R}^{|N|}$

on

$\mathbf{W}_{|V| \times |N|}$

$\hat{\mathbf{y}} \in \mathbb{R}^{|V|}$

$\mathbf{x} \in \mathbb{R}^{|V|}$

*slide from Vagelis Hristidis

# **CBOW**: Continuous Bag of Words  [ Mikolov et al., 2013 ]

**Input** layer

**Parameters** to be learned

cat

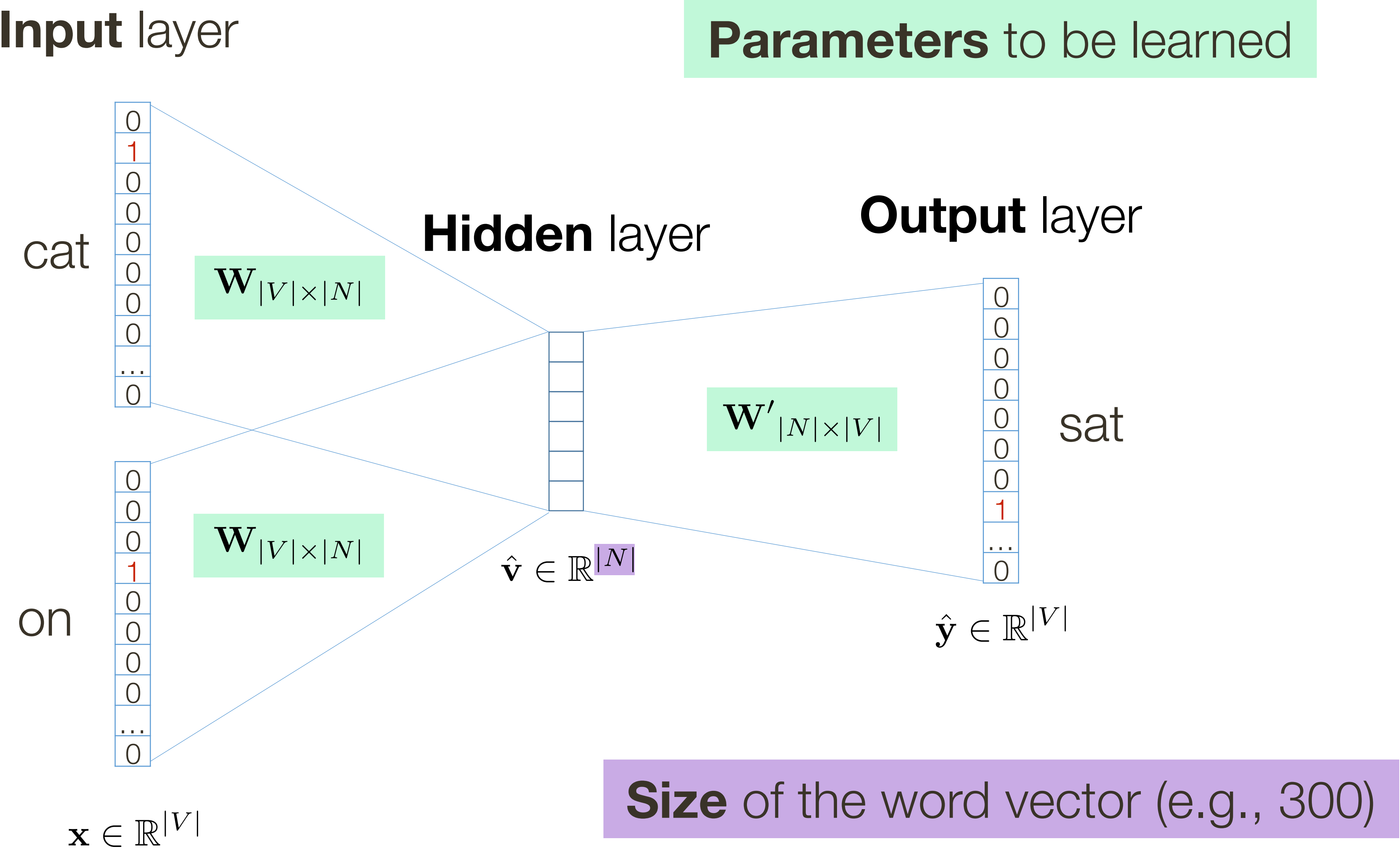$\mathbf{W}_{|V| \times |N|}$

**Hidden** layer

**Output** layer

$\mathbf{W}'_{|N| \times |V|}$

sat

on

$\mathbf{W}_{|V| \times |N|}$

$\hat{\mathbf{v}} \in \mathbb{R}^{|N|}$

$\hat{\mathbf{y}} \in \mathbb{R}^{|V|}$

$\mathbf{x} \in \mathbb{R}^{|V|}$

**Size** of the word vector (e.g., 300)

*slide from Vagelis Hristidis

# **CBOW**: Continuous Bag of Words

[ Mikolov et al., 2013 ]

**Input** layer



**Hidden** layer

**Output** layer

$\mathbf{x}_{cat}$

$\mathbf{W}_{|N|\times|V|} \times \mathbf{x}_{cat} = \mathbf{v}_{cat}$

$\mathbf{x}_{on}$

$\mathbf{W}_{|N|\times|V|} \times \mathbf{x}_{on} = \mathbf{v}_{on}$

$\hat{\mathbf{v}} \in \mathbb{R}^{|N|}$

sat

$\hat{\mathbf{y}} \in \mathbb{R}^{|V|}$

$\mathbf{x} \in \mathbb{R}^{|V|}$

*slide from Vagelis Hristidis

# CBOW: Continuous Bag of Words

**Input** layer



$$\mathbf{W}^T_{|V|\times|N|} \quad \times \quad \mathbf{x}_{cat} \quad = \quad \mathbf{v}_{cat}$$

$\mathbf{W}_{|N|\times|V|} \times \mathbf{x}_{cat}$

$\mathbf{x}_{cat}$

$\mathbf{x}_{on}$

$\mathbf{W}_{|N|\times|V|} \times \mathbf{x}$

$\mathbf{x} \in \mathbb{R}^{|V|}$

| 0.1 | 2.4 | 1.6 | 1.8 | 0.5 | 0.9 | ... | ... | ... | 3.2 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 0.5 | 2.6 | 1.4 | 2.9 | 1.5 | 3.6 | ... | ... | ... | 6.1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 0.6 | 1.8 | 2.7 | 1.9 | 2.4 | 2.0 | ... | ... | ... | 1.2 |

*slide from Vagelis Hristidis

# CBOW: Continuous Bag of Words

**Input** layer



$$\mathbf{x}_{cat}$$

$$W_{|N|\times|V|} \times \mathbf{x}_{cat}$$

$$\mathbf{x}_{on}$$

$$W_{|N|\times|V|} \times \mathbf{x}$$

$$\mathbf{x} \in \mathbb{R}^{|V|}$$

$$\mathbf{W}^T_{|V|\times|N|} \quad \times \quad \mathbf{x}_{on} \quad = \quad \mathbf{v}_{on}$$

| 0.1 | 2.4 | 1.6 | 1.8 | 0.5 | 0.9 | ... | ... | ... | 3.2 |
| 0.5 | 2.6 | 1.4 | 2.9 | 1.5 | 3.6 | ... | ... | ... | 6.1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 0.6 | 1.8 | 2.7 | 1.9 | 2.4 | 2.0 | ... | ... | ... | 1.2 |

$$\times$$

| 0 |
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |

$$=$$

| **1.8** |
| **2.9** |
| **...** |
| **...** |
| **1.9** |

*slide from Vagelis Hristidis

# CBOW: Continuous Bag of Words

**Input** layer

**Hidden** layer

**Output** layer

$\mathbf{x}_{cat}$

$0$
$1$
$0$
$0$
$0$
$0$
$0$
$0$
$0$
$...$
$0$

$W_{|N| \times |V|} \times \mathbf{x}_{cat} = \mathbf{v}_{cat}$

$\mathbf{x}_{on}$

$0$
$0$
$0$
$1$
$0$
$0$
$0$
$0$
$...$
$0$

$\mathbf{W}_{|N| \times |V|} \times \mathbf{x}_{on} = \mathbf{v}_{on}$

$\hat{\mathbf{v}} = \dfrac{\mathbf{v}_{cat} + \mathbf{v}_{on}}{2}$

$\hat{\mathbf{v}} \in \mathbb{R}^{|N|}$

$0$
$0$
$0$
$0$
$0$
$0$
$0$
$0$
$1$
$...$
$0$

sat

$\hat{\mathbf{y}} \in \mathbb{R}^{|V|}$

$\mathbf{x} \in \mathbb{R}^{|V|}$

*slide from Vagelis Hristidis

# **CBOW**: Continuous Bag of Words

**Input** layer

$\mathbf{x}_{cat}$

$$W_{|N|\times|V|} \times \mathbf{x}_{cat} = \mathbf{v}_{cat}$$

**Hidden** layer

**Output** layer

$$\hat{\mathbf{y}} = \mathbf{softmax}(\mathbf{z})$$

$$\mathbf{W'}_{|V|\times|N|} \times \hat{\mathbf{v}} = \mathbf{z}$$

$$\hat{\mathbf{y}}_{sat}$$

$$\hat{\mathbf{v}} \in \mathbb{R}^{|N|}$$

$\mathbf{x}_{on}$

$$W_{|N|\times|V|} \times \mathbf{x}_{on} = \mathbf{v}_{on}$$

$$\hat{\mathbf{y}} \in \mathbb{R}^{|V|}$$

$$\mathbf{x} \in \mathbb{R}^{|V|}$$

*slide from Vagelis Hristidis

# **CBOW**: Continuous Bag of Words

[ Mikolov et al., 2013 ]

**Input** layer

$\mathbf{x}_{cat}$

$\mathbf{W}^T_{|V|\times|N|}$

| 0.1 | 2.4 | 1.6 | 1.8 | 0.5 | 0.9 | ... | ... | ... | 3.2 |
| 0.5 | 2.6 | 1.4 | 2.9 | 1.5 | 3.6 | ... | ... | ... | 6.1 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 0.6 | 1.8 | 2.7 | 1.9 | 2.4 | 2.0 | ... | ... | ... | 1.2 |

**Word** vectors

**Output** layer

$\hat{\mathbf{y}} = \mathbf{softmax}(\mathbf{z})$

$\hat{\mathbf{y}}_{sat}$

$= \mathbf{z}$

$\mathbf{x}_{on}$

$\mathbf{W}_{|N|\times|V|} \times \mathbf{x}_{on} = \mathbf{v}_{on}$

$\hat{\mathbf{v}} \in \mathbb{R}^{|N|}$

$\hat{\mathbf{y}} \in \mathbb{R}^{|V|}$

$\mathbf{x} \in \mathbb{R}^{|V|}$

*slide from Vagelis Hristidis

# CBOW: Interesting Observation

**Another way to look at it**: Maximize similarity between context word representation and the word representation itself

$$p(w|c) = \frac{\exp\left[\left(\sum_c \mathbf{W}\mathbf{x}_c\right)^T (\mathbf{W}\mathbf{x}_w)\right]}{\sum_i^{|V|} \exp\left[(\mathbf{W}\mathbf{x}_i)^T (\mathbf{W}\mathbf{x}_w)\right]}$$

# CBOW: Interesting Observation

**Another way to look at it**: Maximize similarity between context word representation and the word representation itself

$$J(\mathbf{W}) = -\frac{1}{T} \sum_{t=1}^{T} \sum_{-m \le j \le m; j \ne 0} \log p(w_{t+j}|w_t)$$

$$p(w_{t+j}|w_t) = \frac{\exp(\mathbf{w}_{t+j}^T \mathbf{w}_t)}{\sum_{i=1}^{|V|} \exp(\mathbf{w}_i^T \mathbf{w}_t)}$$

# Comparison

— **CBOW** is not great for rare words and typically needs less data to train

— **Skip-gram** better for rate words and needs more data to train the model

| Model | Vector Dimensionality | Training words | Accuracy [%] | | |
|---|---|---|---|---|---|
| | | | Semantic | Syntactic | Total |
| Collobert-Weston NNLM | 50 | 660M | 9.3 | 12.3 | 11.0 |
| Turian NNLM | 50 | 37M | 1.4 | 2.6 | 2.1 |
| Turian NNLM | 200 | 37M | 1.4 | 2.2 | 1.8 |
| Mnih NNLM | 50 | 37M | 1.8 | 9.1 | 5.8 |
| Mnih NNLM | 100 | 37M | 3.3 | 13.2 | 8.8 |
| Mikolov RNNLM | 80 | 320M | 4.9 | 18.4 | 12.7 |
| Mikolov RNNLM | 640 | 320M | 8.6 | 36.5 | 24.6 |
| Huang NNLM | 50 | 990M | 13.3 | 11.6 | 12.3 |
| Our NNLM | 20 | 6B | 12.9 | 26.4 | 20.3 |
| Our NNLM | 50 | 6B | 27.9 | 55.8 | 43.2 |
| Our NNLM | 100 | 6B | 34.2 | **64.5** | 50.8 |
| CBOW | 300 | 783M | 15.5 | 53.1 | 36.1 |
| Skip-gram | 300 | 783M | **50.0** | 55.9 | **53.3** |

# Interesting Results: **Word Analogies**

Test for linear relationships, examined by Mikolov et al. (2014)

$$a:b :: c:?$$

$$d = \arg\max_{x} \frac{(w_b - w_a + w_c)^T w_x}{\|w_b - w_a + w_c\|}$$

man:woman :: king:?

| | | |
|---|---|---|
| + | king | [ 0.30 0.70 ] |
| - | man | [ 0.20 0.20 ] |
| + | woman | [ 0.60 0.30 ] |
| | queen | [ 0.70 0.80 ] |

# **Language** Models

Model the **probability of a sentence**; ideally be able to sample plausible sentences

Why is this useful?

$$\underset{wordsequence}{\arg\max} \; P(wordsequence \,|\, acoustics) =$$

$$\underset{wordsequence}{\arg\max} \; \frac{P(acoustics \,|\, wordsequence) \times P(wordsequence)}{P(acoustics)}$$

$$\underset{wordsequence}{\arg\max} \; P(acoustics \,|\, wordsequence) \times P(wordsequence)$$

# Simple **Language Models**: N-Grams

Given a word sequence: $w_{1:n} = [w_1, w_2, ..., w_n]$

We want to estimate $p(w_{1:n})$

Using **Chain Rule** of probabilities:

$$p(w_{1:n}) = p(w_1)p(w_2|w_1)p(w_3|w_1, w_2) \cdots p(w_n|w_{1:n-1})$$

**Bi-gram** Approximation:

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

**N-gram** Approximation:

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

# Estimating **Probabilities**

N-gram conditional probabilities can be estimated based on raw concurrence counts in the observed sequences

**Bi-gram**:

$$p(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$$

**N-gram**:

$$p(w_n|w_{n-N-1:n-1}) = \frac{C(w_{n-N-1:n-1}w_n)}{C(w_{n-N-1:n-1})}$$

# **Neural-based** Unigram Language Mode



**Problem:** Does not model sequential information (too local)

**We need sequence modeling!**

# **Sequence** Modeling

# Why Model **Sequences**?



$$a_1=2 \quad a_2=0 \quad a_3=1 \quad a_4=3 \quad a_5=4 \quad a_6=2 \quad a_7=5$$

$x$ = bringen sie bitte das auto zurück .

$y$ = please return the car .

* slide from Dhruv Batra

# **Multi-modal** tasks

[ Vinyals *et al.*, 2015 ]



Vision Deep CNN → Language Generating RNN → A group of people shopping at an outdoor market. There are many vegetables at the fruit stand.

# **Sequences** where you don't expect them …

Classify images by taking a
series of "glimpses"



[ Gregor et al., ICML 2015 ]

[ Mnih et al., ICLR 2015 ]

# **Sequences** in Inputs or Outputs?

one to one

**Input:** No sequence

**Output:** No seq.

**Example:**
"standard"
classification /
regression problems

# **Sequences** in Inputs or Outputs?



**one to one**

**one to many**

**Input:** No sequence

**Output:** No seq.

**Example:** "standard" classification / regression problems

**Input:** No sequence

**Output:** Sequence

**Example:** Im2Caption

# **Sequences** in Inputs or Outputs?



**one to one**

**Input:** No sequence
**Output:** No seq.

**Example:** "standard" classification / regression problems

**one to many**

**Input:** No sequence
**Output:** Sequence
**Example:** Im2Caption

**many to one**

**Input:** Sequence
**Output:** No seq.
**Example:** sentence classification, multiple-choice question answering

# **Sequences** in Inputs or Outputs?



**one to one**

**Input:** No sequence
**Output:** No seq.
**Example:** "standard" classification / regression problems

**one to many**

**Input:** No sequence
**Output:** Sequence
**Example:** Im2Caption

**many to one**

**Input:** Sequence
**Output:** No seq.
**Example:** sentence classification, multiple-choice question answering

**many to many**

**Input:** Sequence
**Output:** Sequence
**Example:** machine translation, video captioning, open-ended question answering, video question answering

**many to many**

# **Key Conceptual** Ideas

**Parameter Sharing**

— in computational graphs = adding gradients


"**Unrolling**"

— in computational graphs with parameter sharing


Parameter Sharing + "Unrolling"

— Allows modeling **arbitrary length sequences**!

— Keeps number of parameters in check

# **Recurrent** Neural Network

# **Recurrent** Neural Network

y

**RNN**

x

usually want to predict a
vector at some time steps

# **Recurrent** Neural Network

We can process a sequence of vectors **x** by applying a **recurrence formula** at every time step:

new **state**        old **state**

$$h_t = f_W(h_{t-1}, x_t)$$

some **function**
with parameters W

**input** vector at
some time step

**y**

**RNN**

**x**

# **Recurrent** Neural Network

We can process a sequence of vectors **x** by applying a **recurrence formula** at every time step:

$$h_t = f_W(h_{t-1}, x_t)$$

**Note:** the same function and the same set of parameters are used at every time step

**y**

**RNN**

**x**

# (Vanilla) **Recurrent** Neural Network

$$y_t = W_{hy} h_t + b_y$$

$$h_t = f_W(h_{t-1}, x_t)$$

$$\downarrow$$

$$h_t = \tanh(W_{hh} h_{t-1} + W_{xh} x_t + b_h)$$

**y**

**RNN**

**x**

# RNN **Computational Graph**

# RNN **Computational Graph**

# RNN **Computational Graph**

# RNN **Computational Graph**

Re-use the same weight matrix at every time-step

# RNN **Computational Graph**: Many to Many

# RNN **Computational Graph**: Many to Many

# RNN **Computational Graph**: Many to Many

# RNN **Computational Graph**: Many to One

# RNN **Computational Graph**: One to Many

# Sequence to Sequence: Many to One + One to Many

**Many to one:** Encode input sequence in a single vector

**One to many:** Produce output sequence from single input vector

# **Example**: Character-level Language Model

**Vocabulary:**

['h', 'e', 'l', 'o']

Example training sequence:

"hello"

# **Example**: Character-level Language Model

**Vocabulary:**
['h', 'e', 'l', 'o']

Example training sequence:
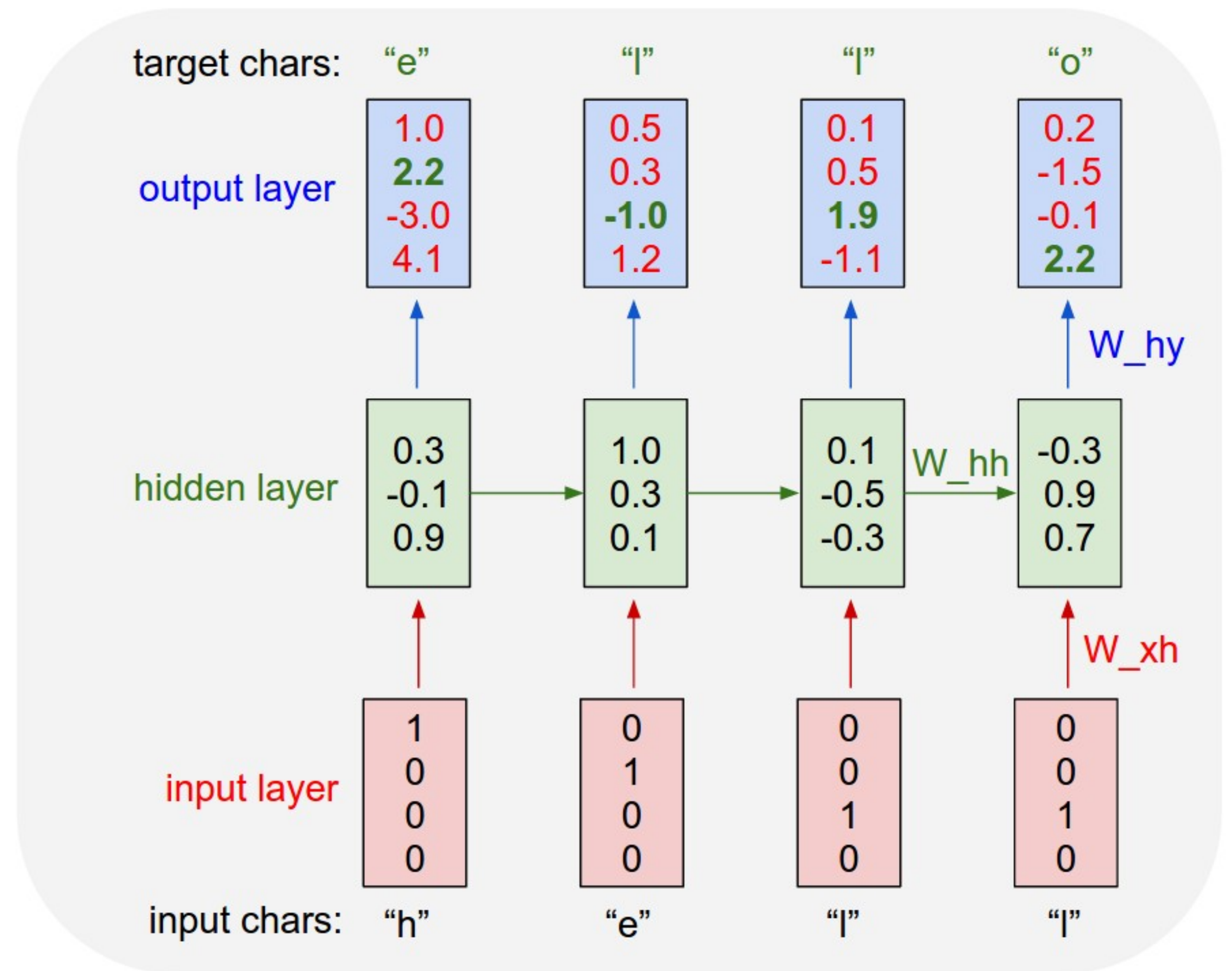"hello"

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$$

# **Example**: Character-level Language Model

**Vocabulary:**

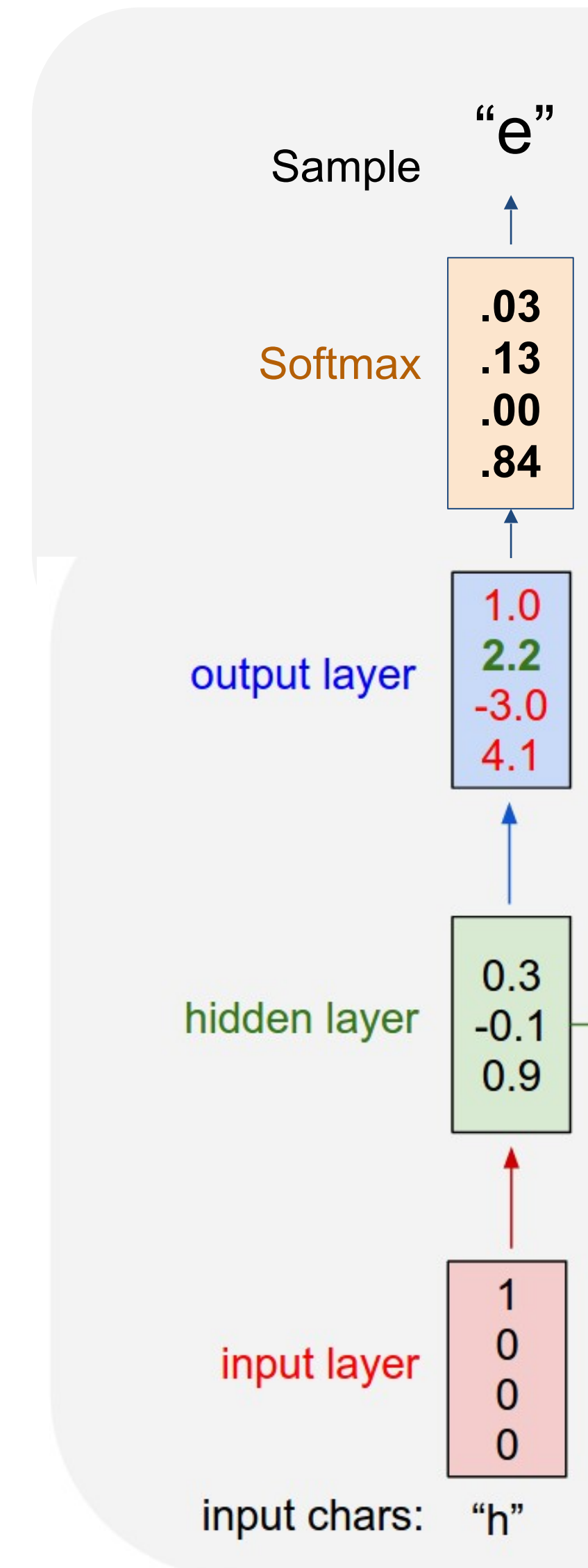['h', 'e', 'l', 'o']

Example training sequence:
"hello"

# **Example**: Character-level Language Model (**Sampling**)

**Vocabulary:**

['h', 'e', 'l', 'o']

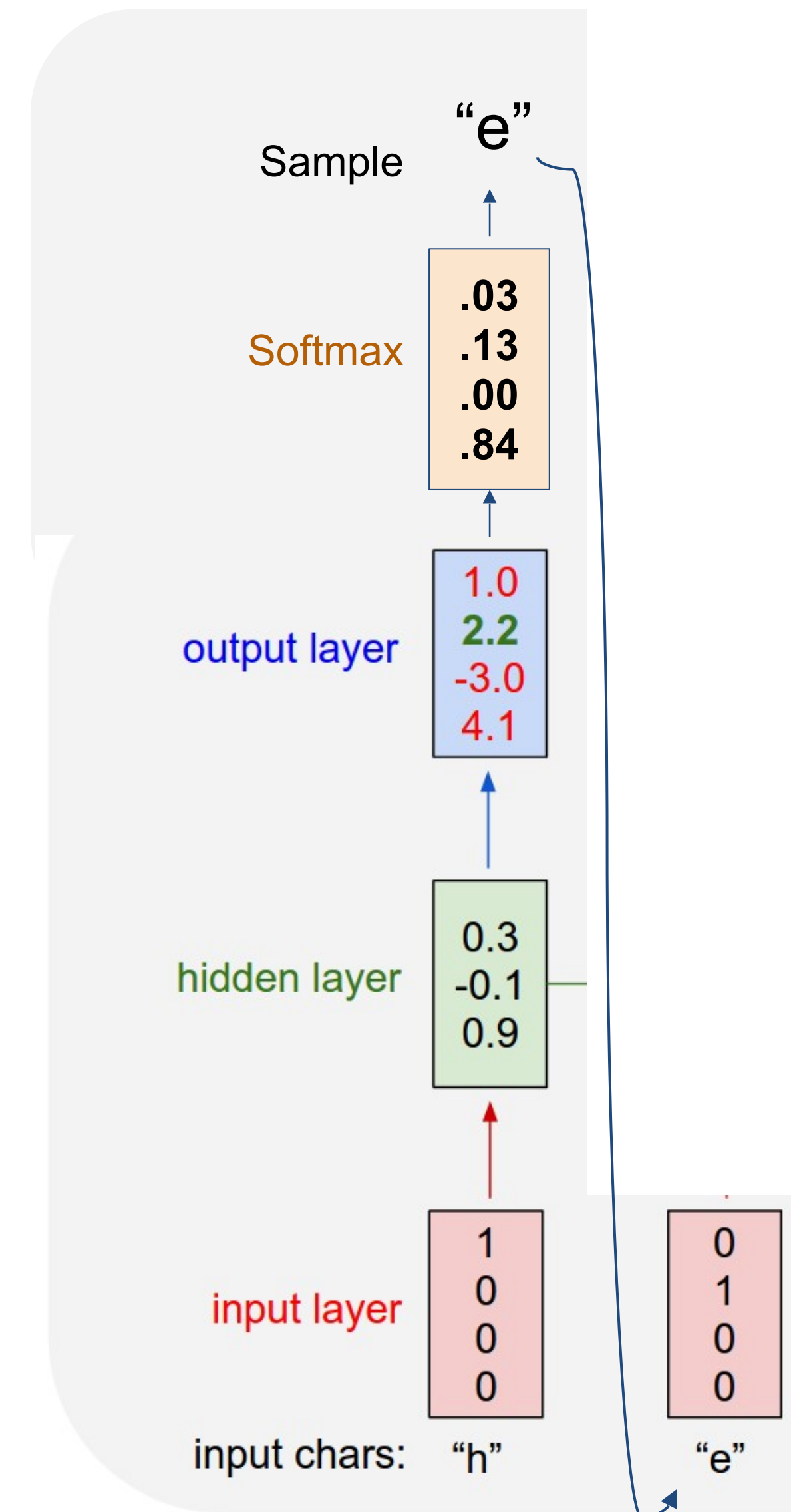At test time sample one character at a time and feed back to the model

# **Example**: Character-level Language Model (**Sampling**)

**Vocabulary:**

['h', 'e', 'l', 'o']

At test time sample one character at a time and feed back to the model
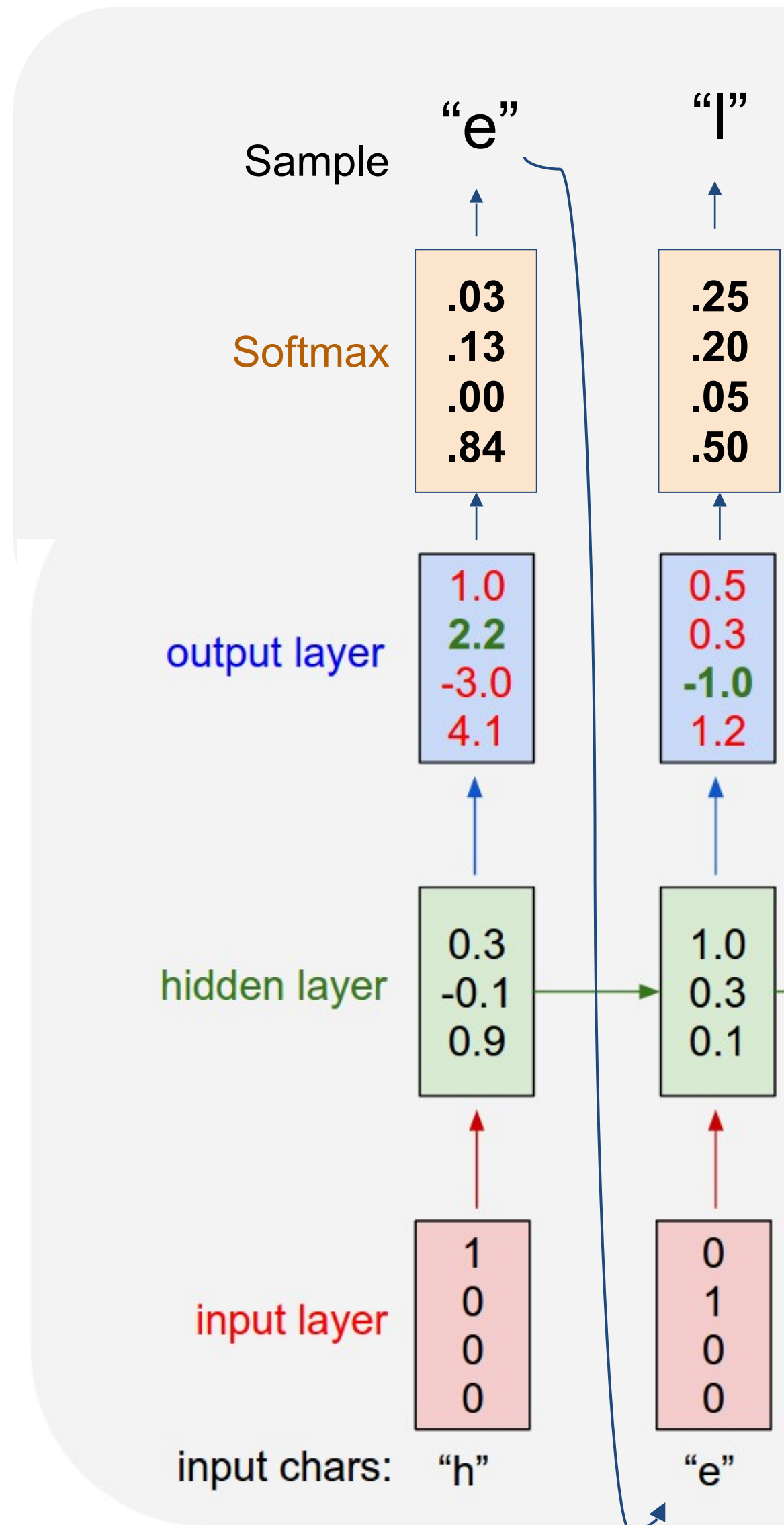
# **Example**: Character-level Language Model (**Sampling**)

**Vocabulary:**

['h', 'e', 'l', 'o']

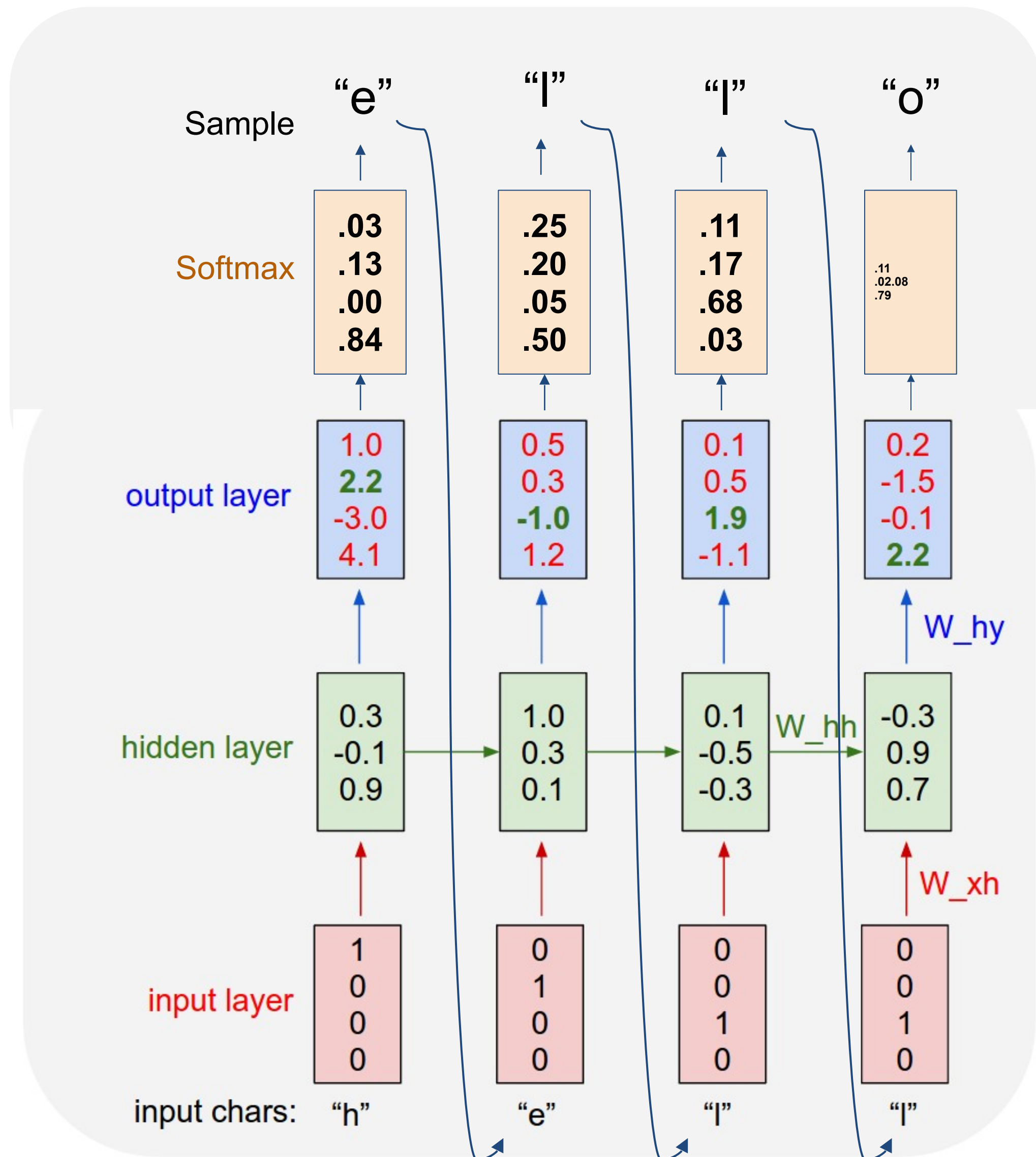At test time sample one character at a time and feed back to the model
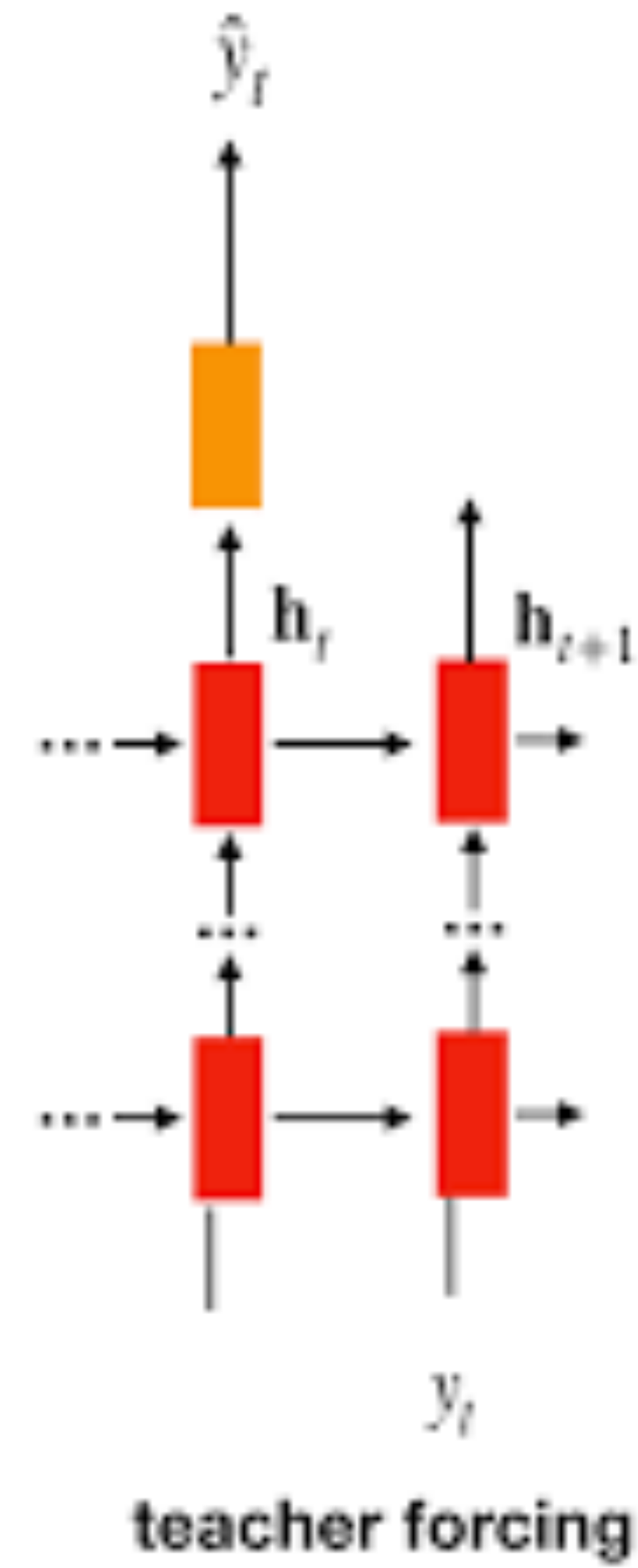
# Example: Character-level Language Model (Sampling)
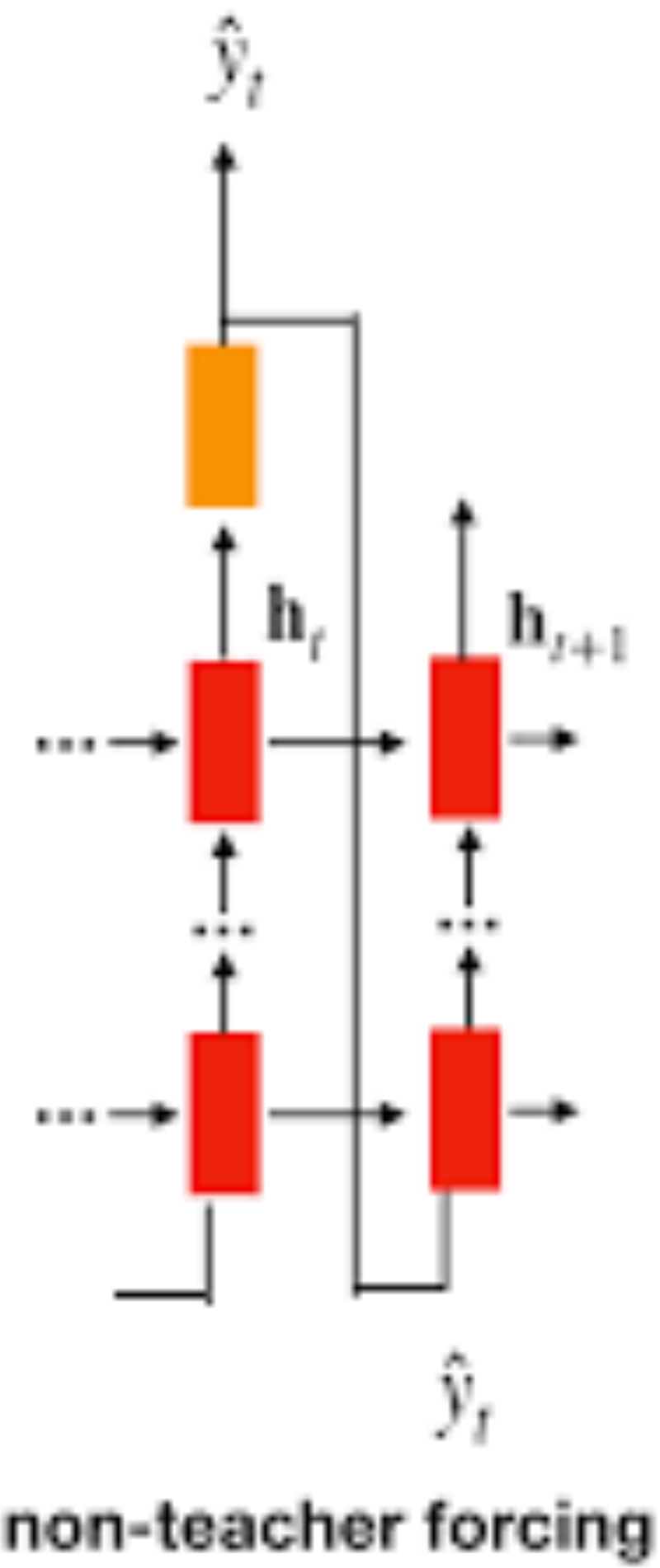
**Vocabulary:**

['h', 'e', 'l', 'o']

At test time sample one character at a time and feed back to the model

# **Teacher** Forcing



non-teacher forcing      teacher forcing

# **Sampling** vs. **ArgMax**

**Sampling**: allows to generate diverse outputs

**ArgMax**: could be more stable in practice