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

**Lecture 11: Word Vector Representations**

# Logistics

**Assignment 3** … was due **last night**

— This is the most difficult assignment in the course

**Assignment 4** … will be out **today**

— Do not wait

**Assignment 5** … will be delayed to enable project proposals

# Logistics

**Paper readings** coming up

Project **groups** and topics

**Invited talks**

# Fun Example: **Code Deobfuscating** with DOBF

**Obfuscated Code**

```python
class CLASS_0(nn.Module):

    def __init__(VAR_0, VAR_1, VAR_2, VAR_3):
        super(CLASS_0, VAR_0).__init__()
        VAR_0.VAR_1 = VAR_1
        VAR_0.VAR_2 = VAR_2
        VAR_0.VAR_4 = nn.Linear(VAR_1, (4 * VAR_2), bias=VAR_3)
        VAR_0.VAR_5 = nn.Linear(VAR_2, (4 * VAR_2), bias=VAR_3)
        VAR_0.FUNC_0()

    def FUNC_0(VAR_6):
        VAR_7 = (1.0 / math.sqrt(VAR_6.VAR_8))
        for VAR_9 in VAR_6.VAR_10():
            VAR_9.data.uniform_((- VAR_7), VAR_7)

    def FUNC_1(VAR_11, VAR_12, VAR_13):
        (VAR_14, VAR_15) = VAR_13
        VAR_14 = VAR_14.view(VAR_14.size(1), (- 1))
        VAR_15 = VAR_15.view(VAR_15.size(1), (- 1))
        VAR_12 = VAR_12.view(VAR_12.size(1), (- 1))
        VAR_16 = (VAR_11.VAR_4(VAR_12) + VAR_11.VAR_5(VAR_14))
        VAR_17 = VAR_16[:, :(3 * VAR_11.VAR_8)].sigmoid()
        VAR_18 = VAR_16[:, (3 * VAR_11.VAR_8):].tanh()
        VAR_19 = VAR_17[:, :VAR_11.VAR_8]
        VAR_20 = VAR_17[:, VAR_11.VAR_8:(2 * VAR_11.VAR_8)]
        VAR_21 = VAR_17[:, (- VAR_11.VAR_8):]
        VAR_22 = (th.mul(VAR_15, VAR_20) + th.mul(VAR_19, VAR_18))
        VAR_23 = th.mul(VAR_21, VAR_22.tanh())
        VAR_23 = VAR_23.view(1, VAR_23.size(0), (- 1))
        VAR_22 = VAR_22.view(1, VAR_22.size(0), (- 1))
        return (VAR_23, (VAR_23, VAR_22))
```

**Code Deobfuscated using DOBF**

```python
class LSTM(nn.Module):

    def __init__(self, input_size, hidden_size, bias):
        super(LSTM, self).__init__()
        self.input_size = input_size
        self.hidden_size = hidden_size
        self.h1 = nn.Linear(input_size, (4 * hidden_size), bias=bias)
        self.h2 = nn.Linear(hidden_size, (4 * hidden_size), bias=bias)
        self.init_weights()

    def init_weights(self):
        stdv = (1.0 / math.sqrt(self.hidden_size))
        for m in self.modules():
            m.data.uniform_((- stdv), stdv)

    def forward(self, x, prev_state):
        (prev_h, prev_c) = prev_state
        prev_h = prev_h.view(prev_h.size(1), (- 1))
        prev_c = prev_c.view(prev_c.size(1), (- 1))
        x = x.view(x.size(1), (- 1))
        h = (self.h1(x) + self.h2(prev_h))
        s = h[:, :(3 * self.hidden_size)].sigmoid()
        c = h[:, (3 * self.hidden_size):].tanh()
        r = s[:, :self.hidden_size]
        g = s[:, self.hidden_size:(2 * self.hidden_size)]
        o = s[:, (- self.hidden_size):]
        c = (th.mul(prev_c, g) + th.mul(r, c))
        h = th.mul(o, c.tanh())
        h = h.view(1, h.size(0), (- 1))
        c = c.view(1, c.size(0), (- 1))
        return (h, (h, c))
```

[ Roziere et al., ArXiv, 2021 ]

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

**Vocabulary**

dog

cat

person

holding

tree

computer

using

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

**Vocabulary**

dog         1

cat         2

person      3

holding     4

tree        5

computer    6

using       7

# 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

dog cat person holding tree computer using

*slide from V. Ordonex

# 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 ]

dog  cat  person  holding  tree  computer  using

# 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 ]

dog   cat   person   holding   tree   computer   using

*slide from V. Ordonex

# 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

| Vocabulary | |
| --- | --- |
| dog | 1 |
| cat | 2 |
| person | 3 |
| holding | 4 |
| tree | 5 |
| computer | 6 |
| using | 7 |

# 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

# **Word** Representations

1. **One-hot encodings** — only non-zero at the index of the word

    e.g., [ 0, 1, 0, 0, 0, ...., 0, 0, 0 ]

    **Good:** simple

    **Bad:** not compact, distance between words always same (e.g., synonyms vs. antonyms)

2. **Word feature representations** — manually define "good" features

    e.g., [ 1, 1, 0, 30, 0, ...., 0, 0, 0 ] -> 300-dimensional irrespective of dictionary

    e.g., word ends on -ing

3. **Learned word representations** — vector should approximate "meaning" of the word

    e.g., [ 1, 1, 0, 30, 0, ...., 0, 0, 0 ] -> 300-dimensional irrespective of dictionary

    **Good:** compact, distance between words is semantic

# **Distributional** Hypothesis  [ Lenci, 2008 ]

— 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.

# 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]

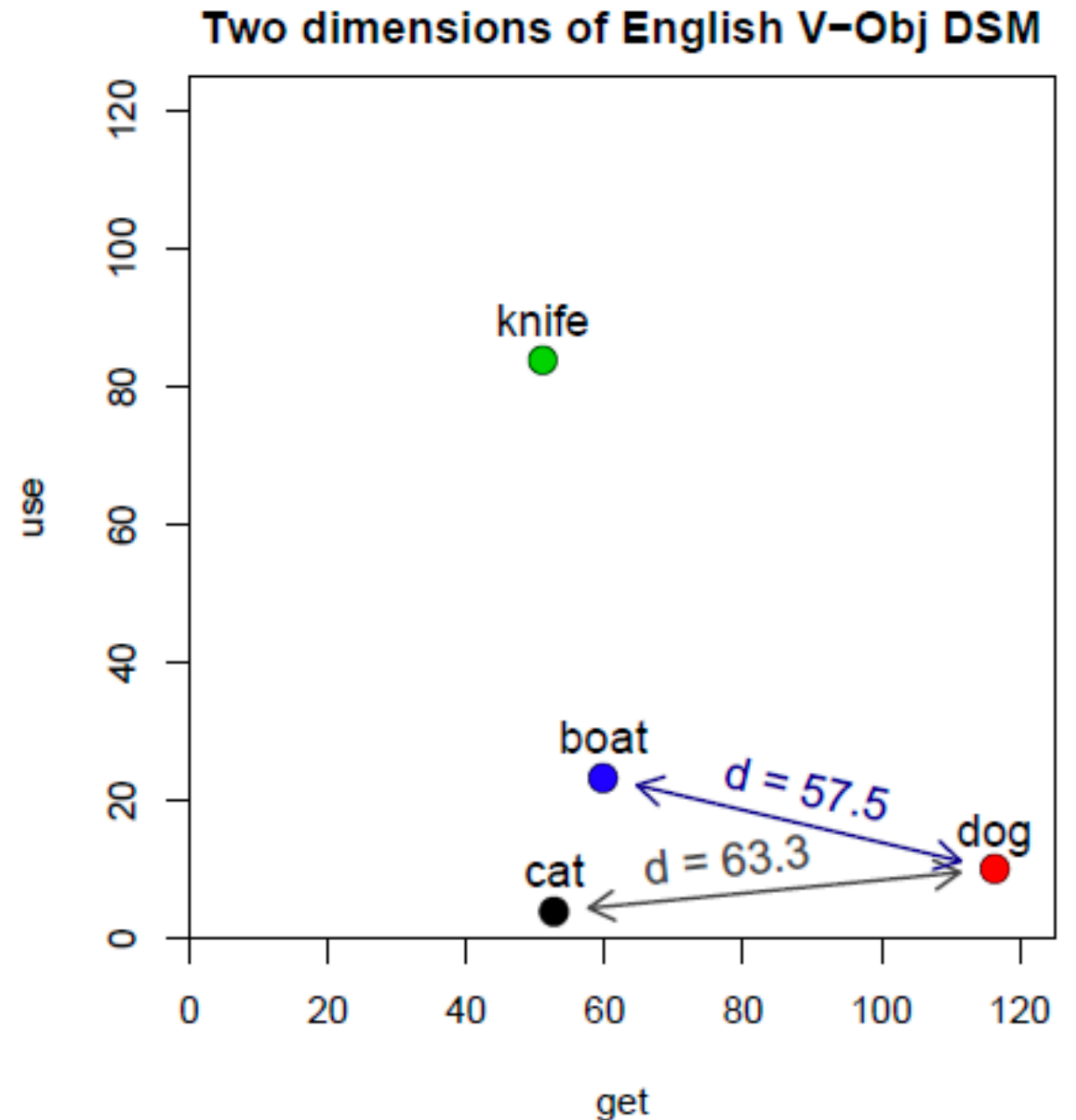# **Geometric Interpretation**: Co-occurrence as feature

— 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

# **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

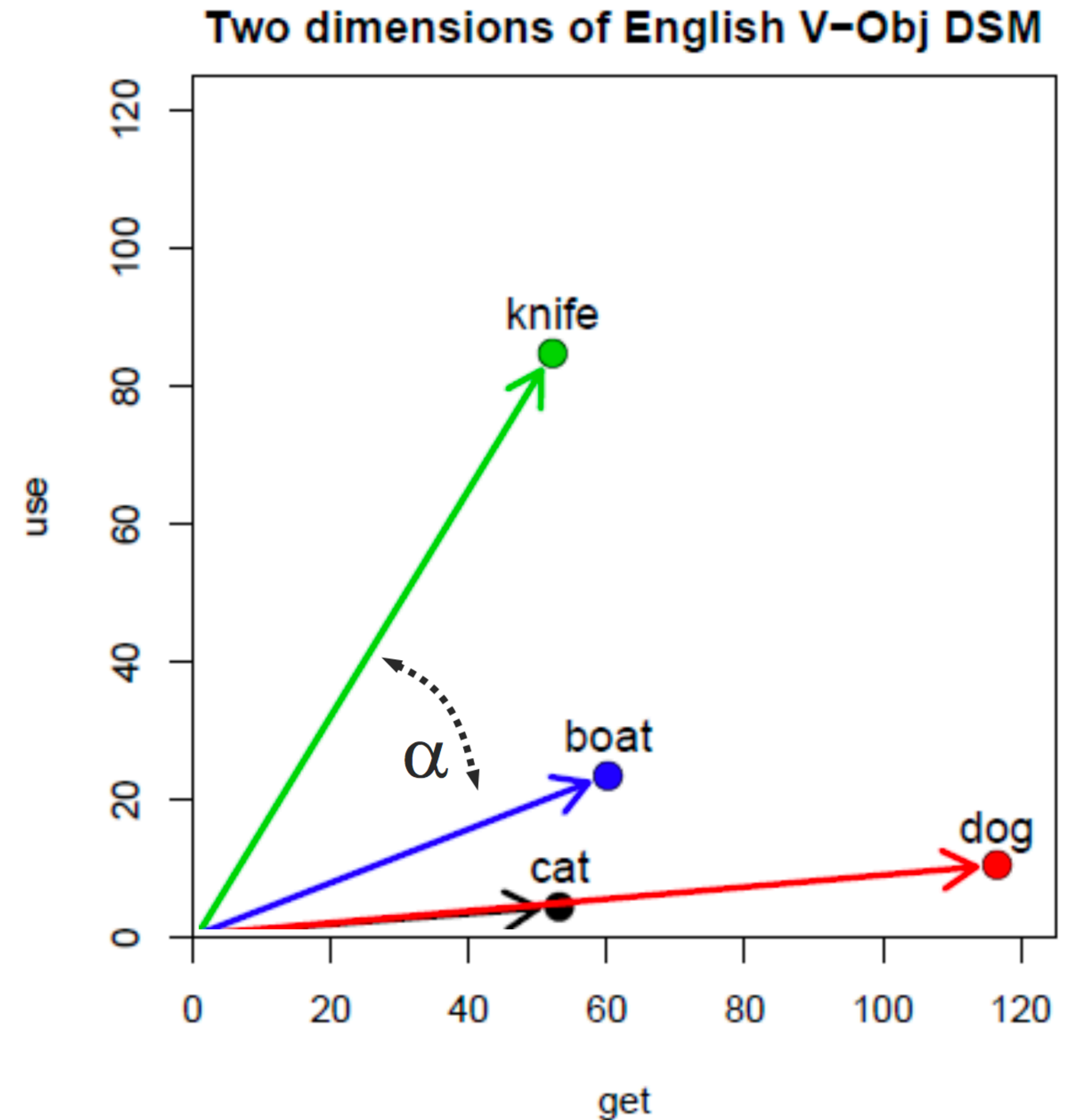|        | get | see | use | hear | eat | kill |
|--------|-----|-----|-----|------|-----|------|
| knife  | 51  | 20  | 84  | 0    | 3   | 0    |
| cat    | 52  | 58  | 4   | 4    | 6   | 26   |
| dog    | 115 | 83  | 10  | 42   | 33  | 17   |
| boat   | 59  | 39  | 23  | 4    | 0   | 0    |
| cup    | 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 cat)



Two dimensions of English V-Obj DSM

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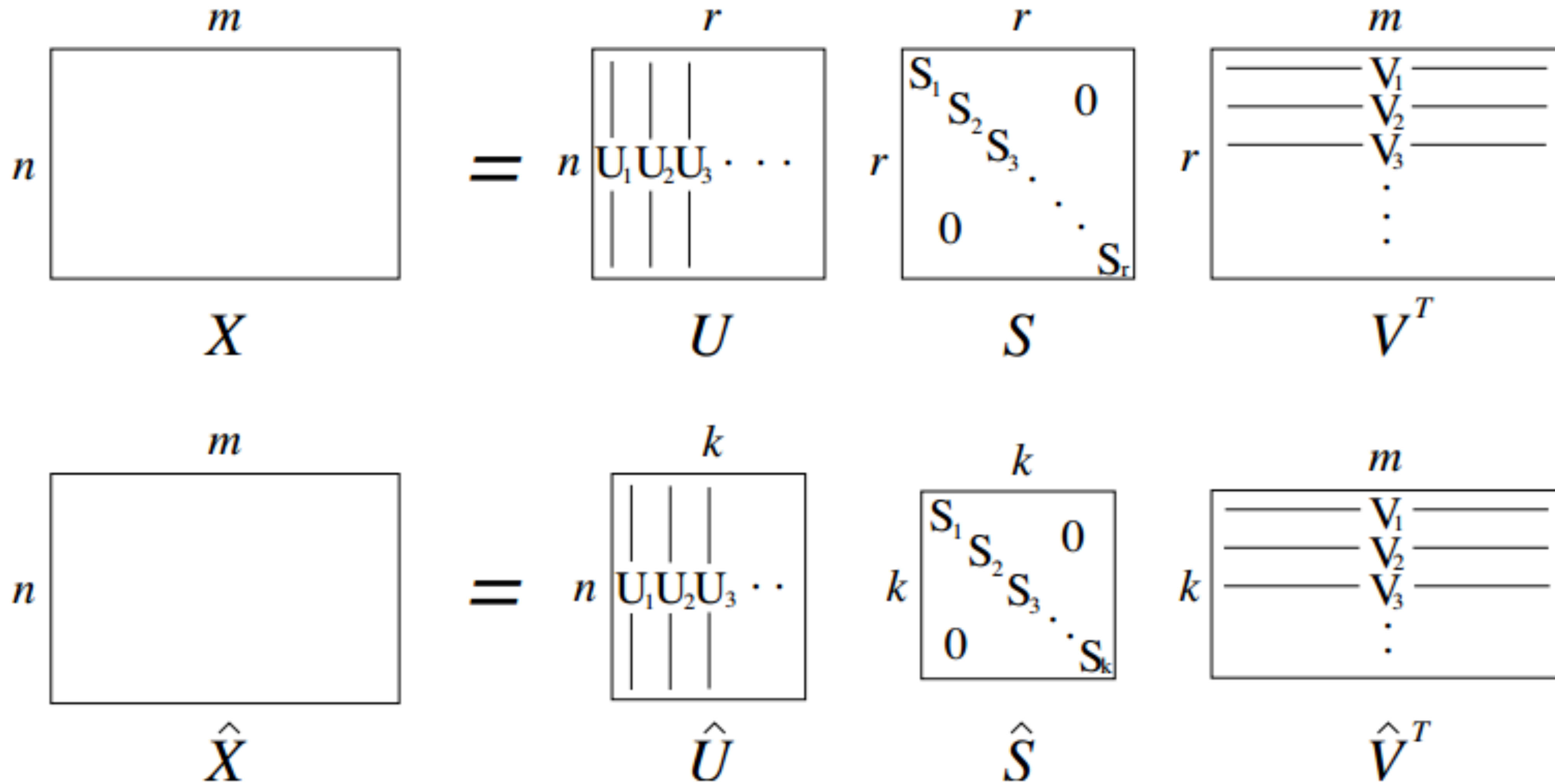
# Geometric Interpretation: Co-occurrence as feature

— 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

| | 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

# **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!**

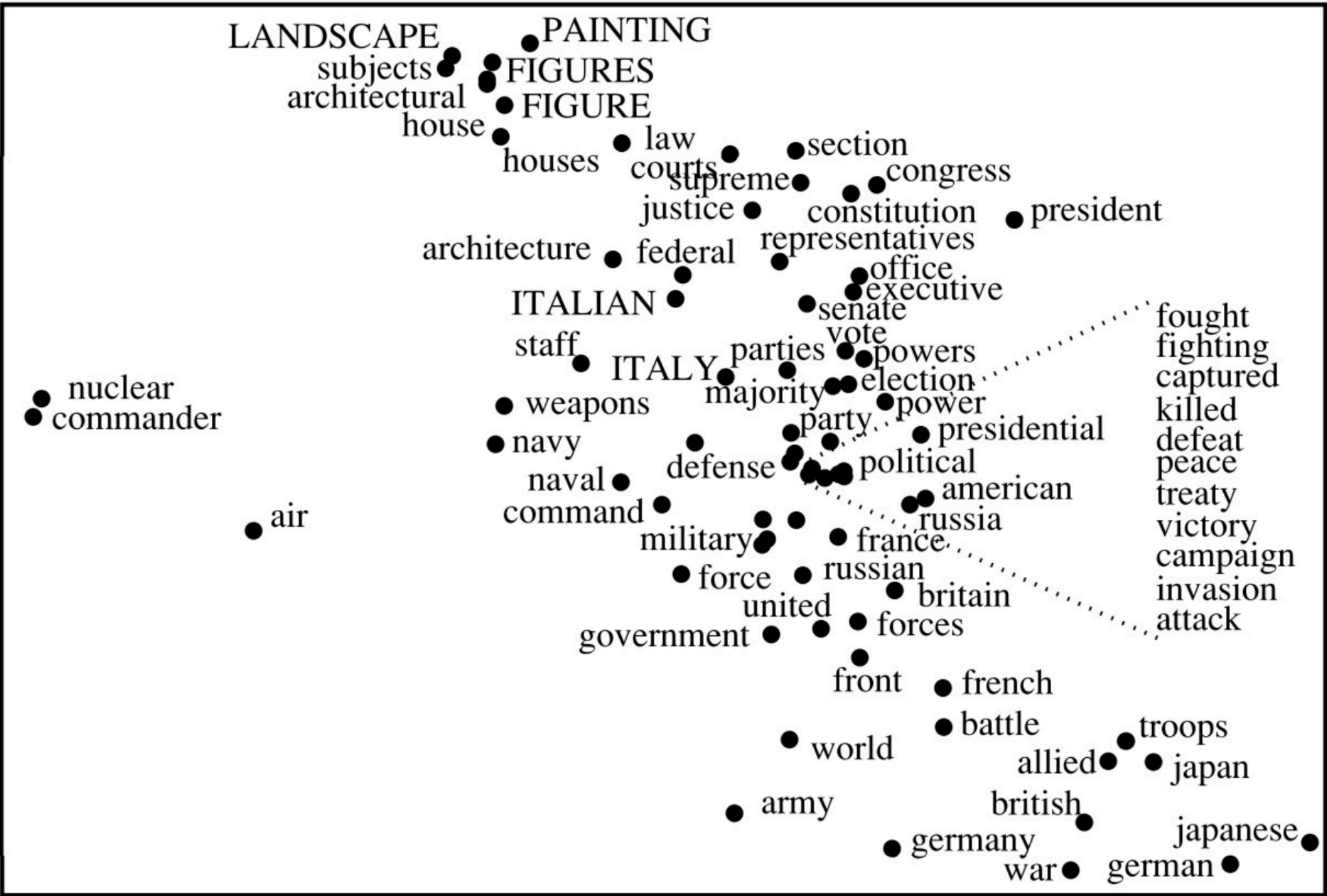|        | get | see | use | hear | eat | kill |
|--------|-----|-----|-----|------|-----|------|
| knife  | 51  | 20  | 84  | 0    | 3   | 0    |
| cat    | 52  | 58  | 4   | 4    | 6   | 26   |
| dog    | 115 | 83  | 10  | 42   | 33  | 17   |
| boat   | 59  | 39  | 23  | 4    | 0   | 0    |
| cup    | 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

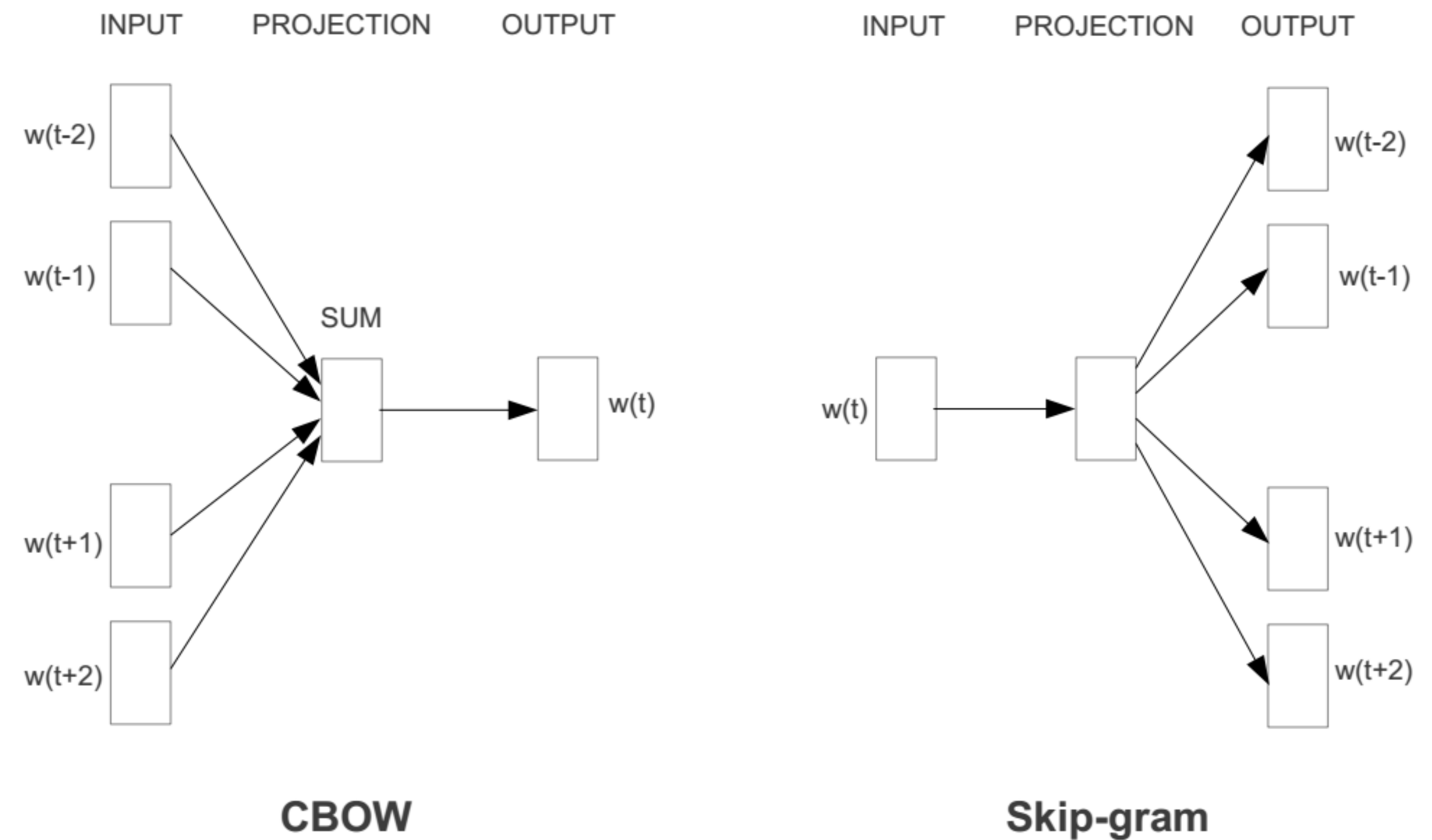# **word2vec**: Representing the Meaning of Words   <span style="color:red">[ Mikolov et al., 2013 ]</span>

**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.
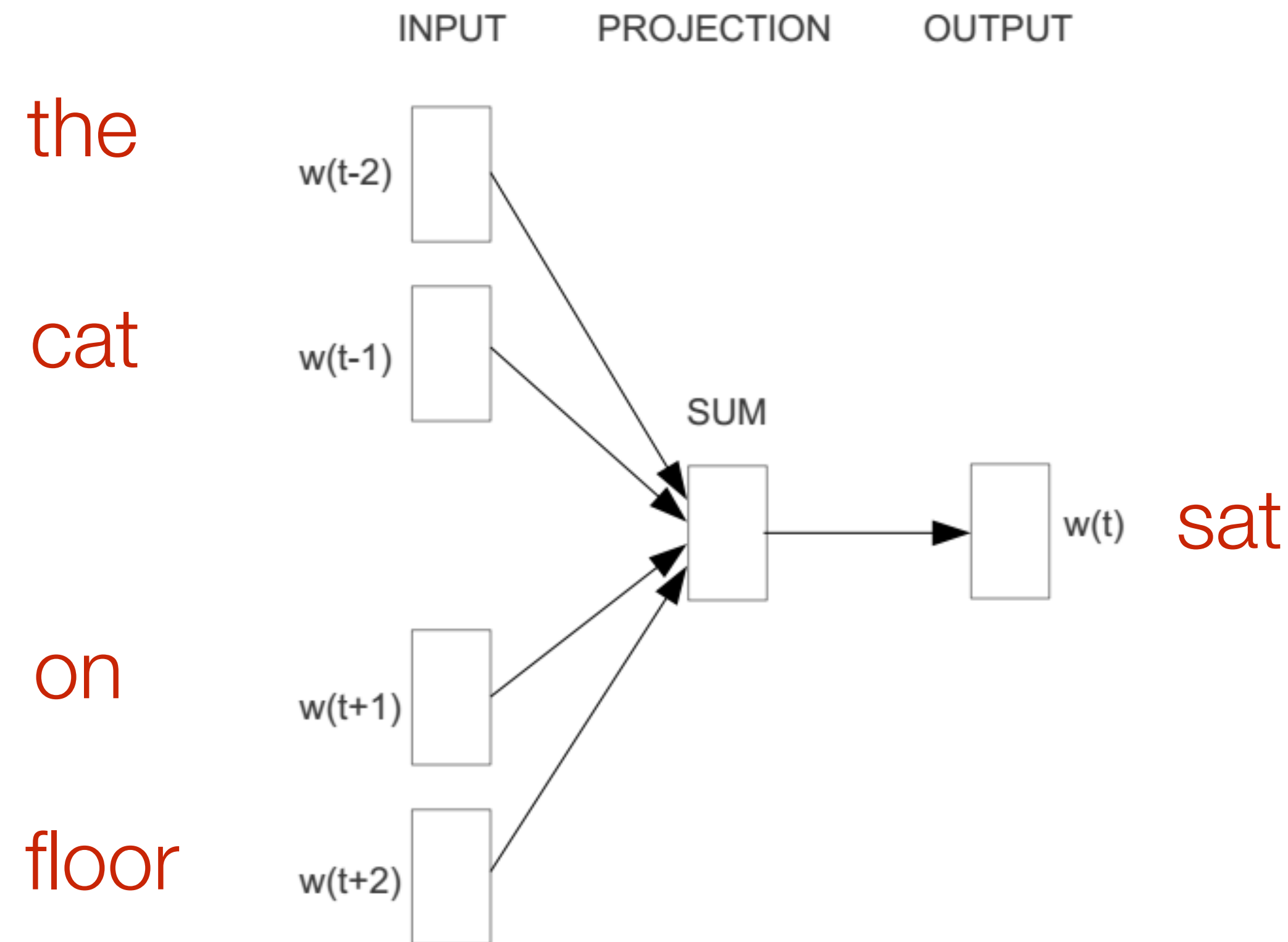


Continuous Bag of Words (**CBOW**): use context words in a window to predict middle word

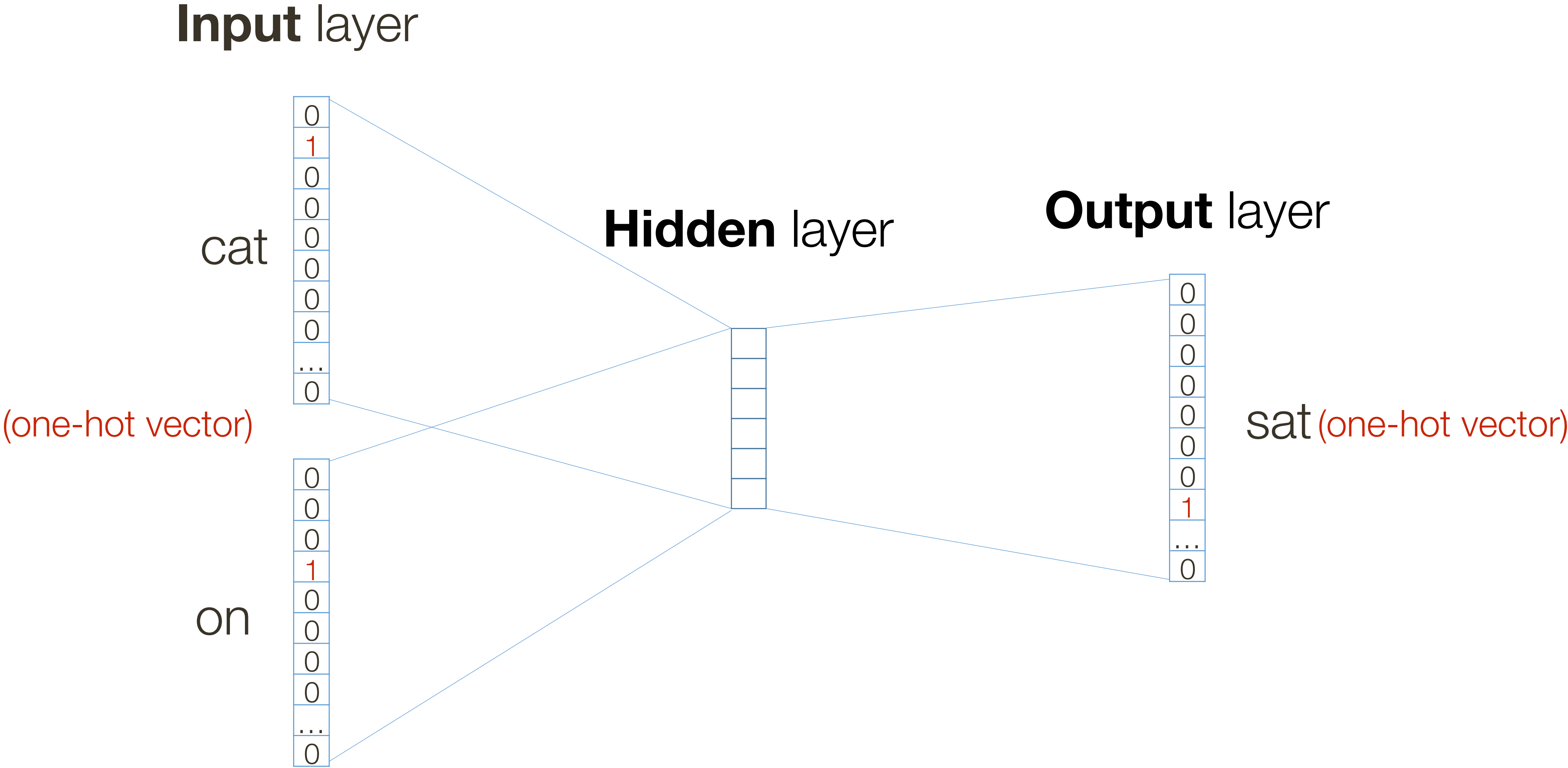**Skip-gram:** use the middle word to predict surrounding ones in a window

*slide from Vagelis Hristidis

# **CBOW**: Continuous Bag of Words

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

# **CBOW**: Continuous Bag of Words

**Input** layer

**Hidden** layer

**Output** layer

cat

0
1
0
0
0
0
0
0
...
0

(one-hot vector)

on

0
0
0
1
0
0
0
0
...
0

sat (one-hot vector)

0
0
0
0
0
0
0
1
...
0

# **CBOW**: Continuous Bag of Words

**Input** layer

**Hidden** layer

**Output** layer

cat

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

on

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

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

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

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

**Parameters** to be learned

**Hidden** layer

**Output** layer

cat

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

on

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

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

sat

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

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

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

*slide from Vagelis Hristidis

# **CBOW**: Continuous Bag of Words

**Input** layer

**Parameters** to be learned

**Hidden** layer

**Output** layer

cat

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

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

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

**Input** layer



$\mathbf{x}_{cat}$

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

**Hidden** layer

**Output** layer

$\mathbf{x}_{on}$

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

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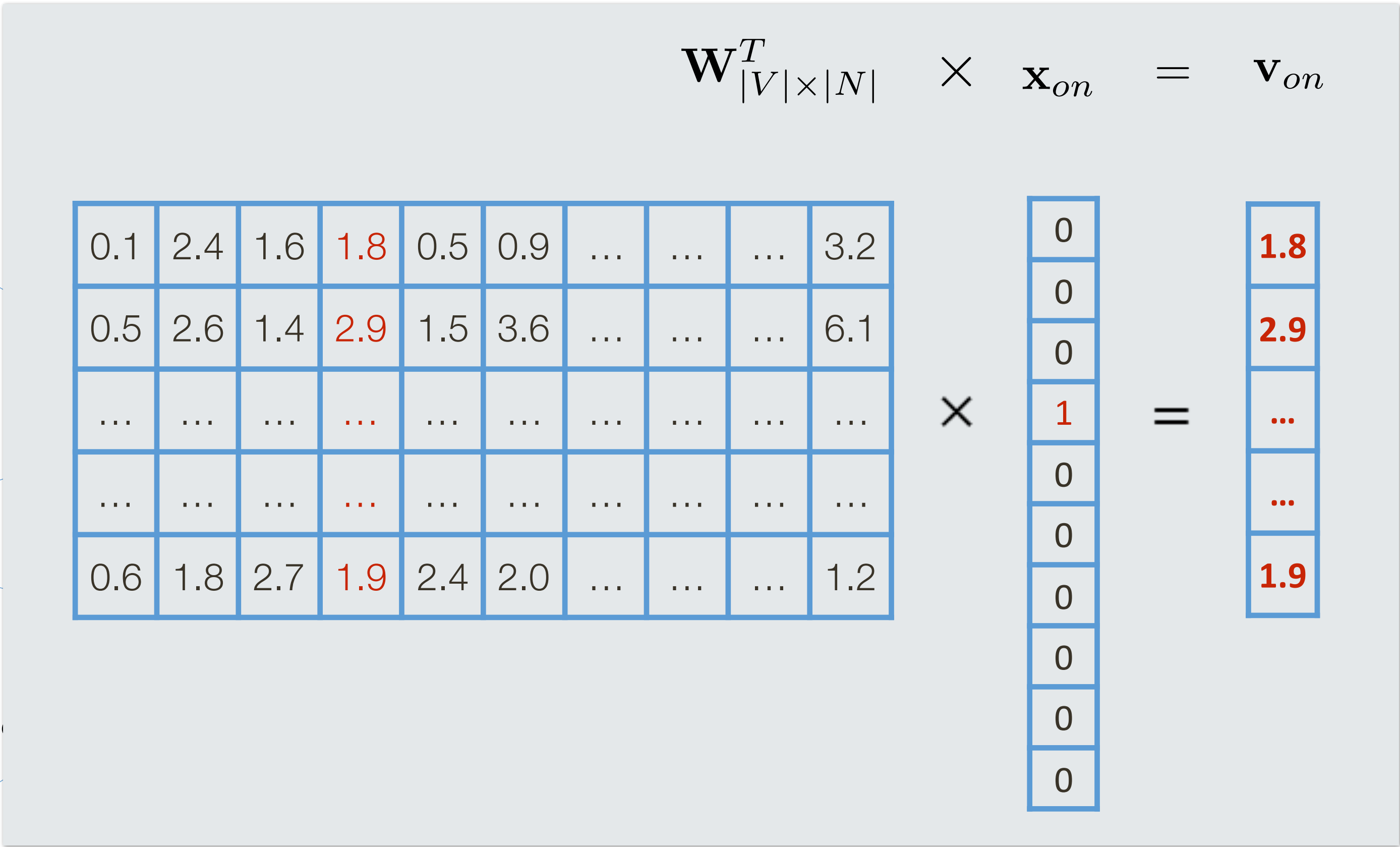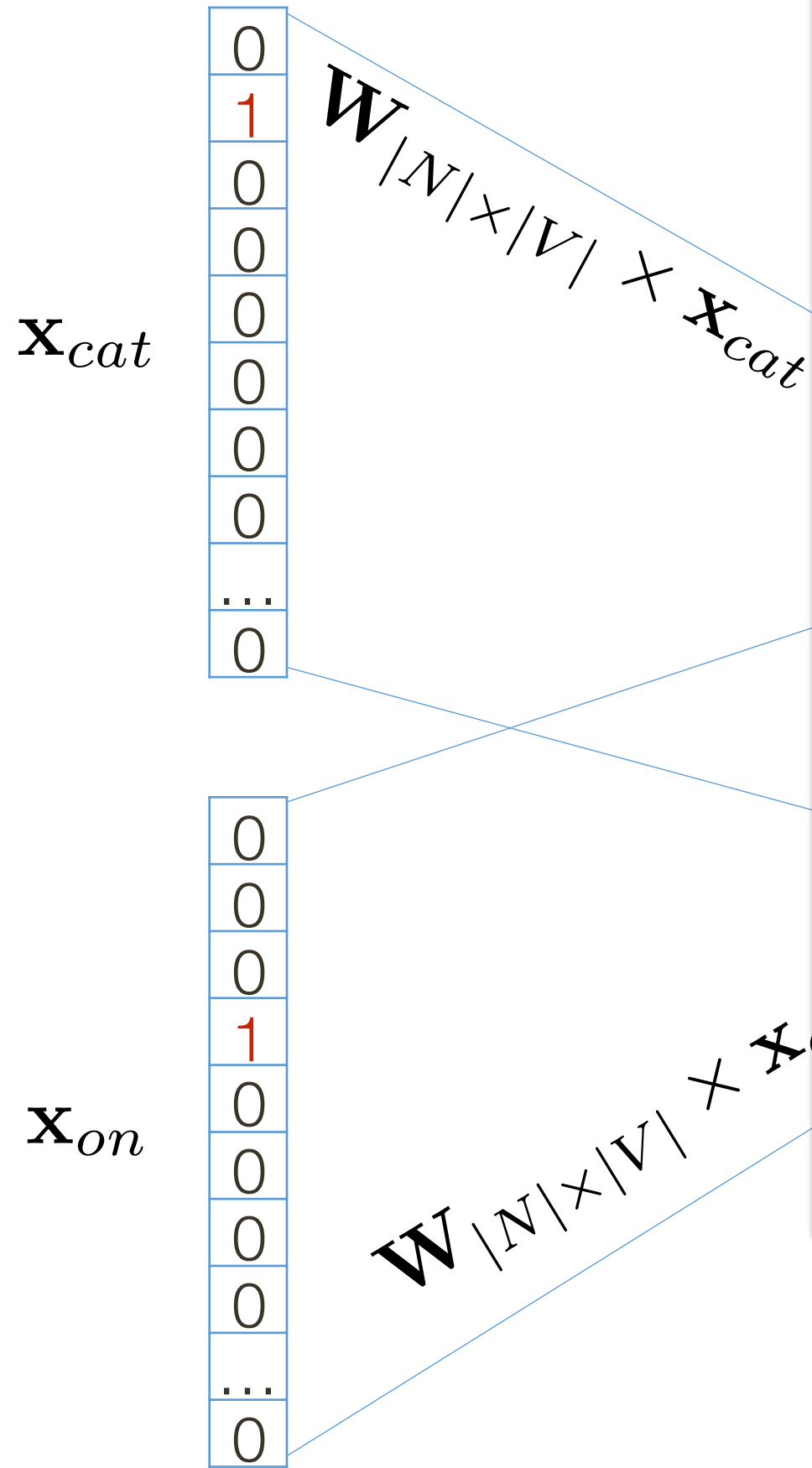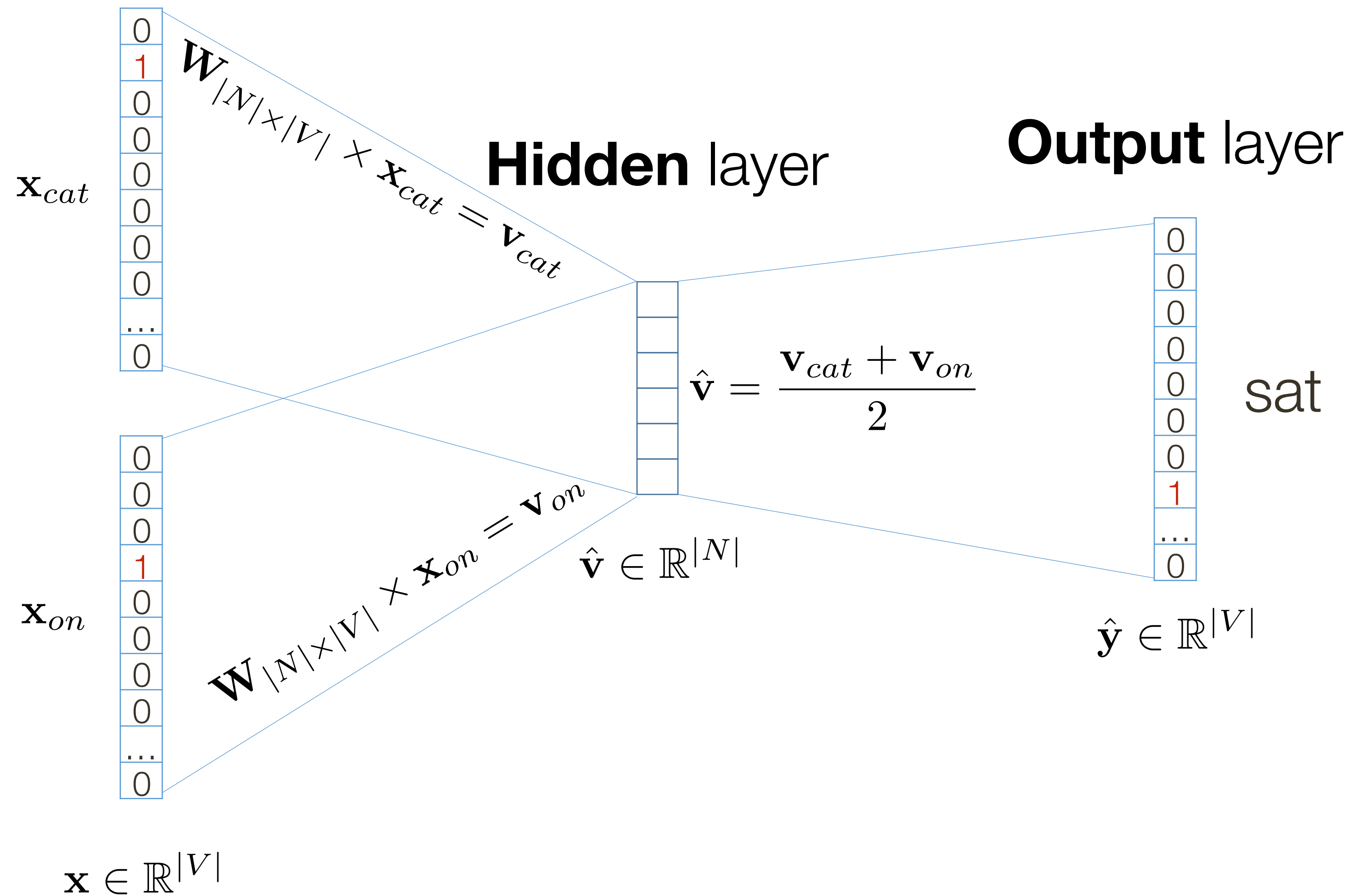
sat

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

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

*slide from Vagelis Hristidis

# CBOW: Continuous Bag of Words

[ Mikolov et al., 2013 ]

**Input** layer

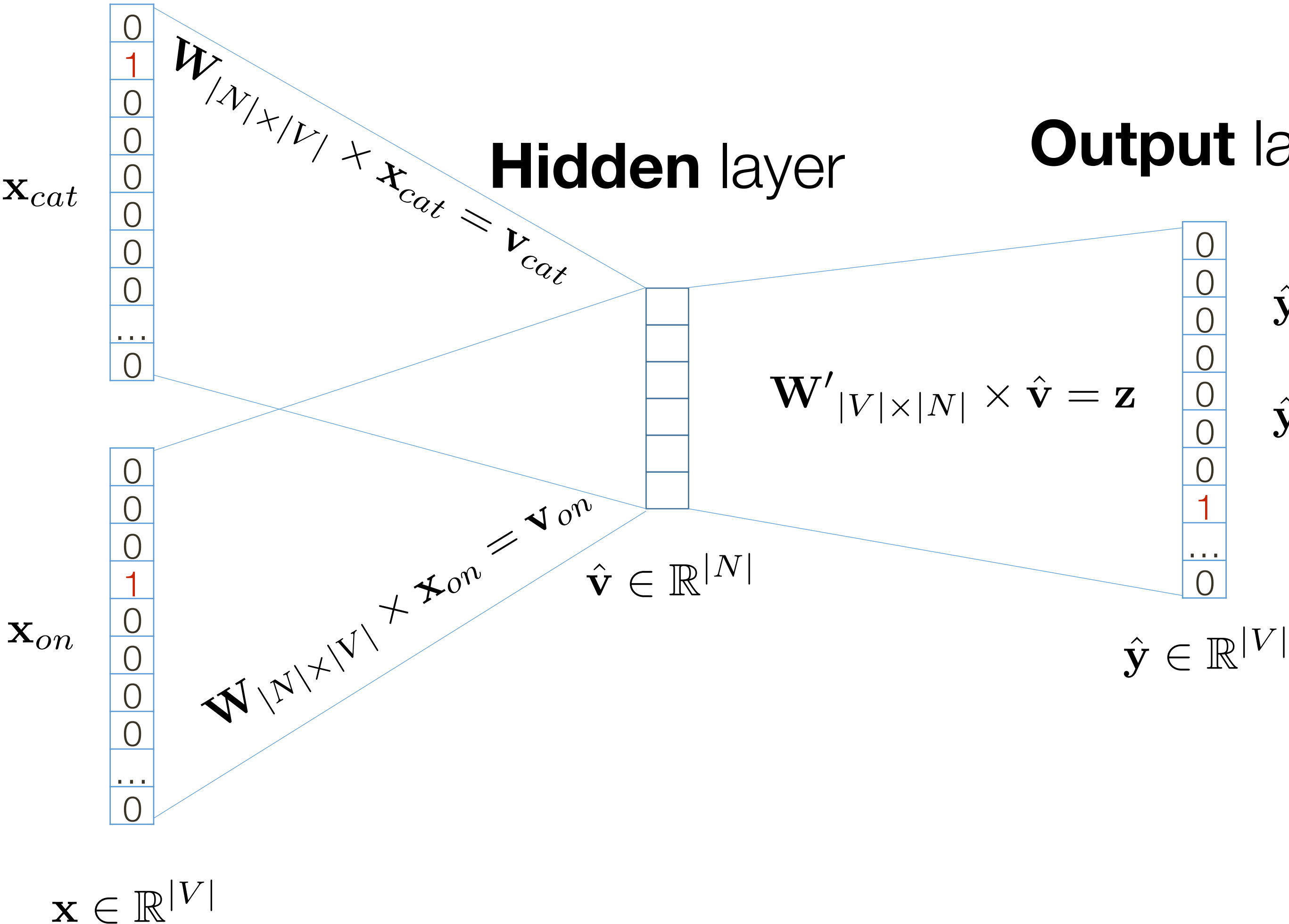$$\mathbf{W}^T_{|V| \times |N|} \quad \times \quad \mathbf{x}_{cat} \quad = \quad \mathbf{v}_{cat}$$

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

$\mathbf{x}_{cat}$

| 0 |
|---|
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| ... |
| 0 |

$\mathbf{x}_{on}$

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

| 0 |
|---|
| 0 |
| 0 |
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| ... |
| 0 |

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

| 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 |

$\times$

| 0 |
|---|
| 1 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |
| 0 |

$=$

| 2.4 |
|-----|
| 2.6 |
| ... |
| ... |
| 1.8 |

*slide from Vagelis Hristidis

# CBOW: Continuous Bag of Words

**Input** layer



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

| 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 |

$\mathbf{x}_{cat}$

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

$\mathbf{x}_{on}$

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

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

*slide from Vagelis Hristidis

# **CBOW**: Continuous Bag of Words

**Input** layer

**Hidden** layer

**Output** layer

$\mathbf{x}_{cat}$

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

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

sat

$\mathbf{x}_{on}$

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

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

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

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

*slide from Vagelis Hristidis

# **CBOW**: Continuous Bag of Words

**Input** layer



$\mathbf{x}_{cat}$

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

**Hidden** layer

**Output** layer

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

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

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

$\mathbf{x}_{on}$

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

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

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

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

*slide from Vagelis Hristidis

# **CBOW**: Continuous Bag of Words

**Input** layer



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

**Hidden** layer

$\mathbf{x}_{cat}$

$\mathbf{x}_{on}$

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

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

**Output** layer

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

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

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

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

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

Optimize to get close to 1-hot encoding

*slide from Vagelis Hristidis

# CBOW: Continuous Bag of Words

**Input** lay

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

| 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

$\mathbf{x}_{cat}$

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

$= \mathbf{z}$

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

$\mathbf{x}_{on}$

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

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

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

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

*slide from Vagelis Hristidis

# CBOW: Interesting Observation

**Input** layer      There are two representations for same word!



$\mathbf{x}_{cat}$

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

**Hidden** layer

**Output** layer

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

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

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

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

$\mathbf{x}_{on}$

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

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

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

*slide from Vagelis Hristidis

# CBOW: 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 [ Mikolov et al., 2013 ]

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

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

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

# **Skip-Gram** Model

# 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 ]

# **Dynamic** Word Embeddings

# **Dynamic** Word Embeddings

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

**Lecture 11: RNN Applications**

Let us look at some actual practical uses of RNNs

# Applications: Skip-thought Vectors

word2vec but for sentences, where each sentence is processed by an LSTM



[ Kiros et al., 2015 ]

# **Applications:** Google Language Translation

One model to translate from **any language** to any other language



[ Johnson et al., 2017 ]

# **Applications:** Google Language Translation

One model to translate from **any language** to any other language



Token designating
**target** language

[ Johnson et al., 2017 ]

# Applications: Google Language Translation

One model to translate from **any language** to any other language



**Flipped** order encoding

Token designating **target** language

[ Johnson et al., 2017 ]

# Applications: Google Language Translation

One model to translate from **any language** to any other language



**Flipped** order encoding

**Why?**

Token designating
**target** language

[ Johnson et al., 2017 ]

# Applications: Google Language Translation

One model to translate from **any language** to any other language



**Flipped** order encoding

Token designating
**target** language

[ Johnson et al., 2017 ]

8! layer LSTM decoder and encoder

# **Applications:** Google Language Translation

One model to translate from **any language** to any other language



**Residual** at
other layers
(ResNet style)

**Bi-directional**
at lower layers

**Flipped** order encoding

Token designating
**target** language

[ Johnson et al., 2017 ]

8! layer LSTM decoder and encoder

# **Applications:** Google Language Translation

One model to translate from **any language** to any other language



**Residual** at other layers (ResNet style)

**Bi-directional** at lower layers

**Flipped** order encoding

Token designating **target** language

[ Johnson et al., 2017 ]

8! layer LSTM decoder and encoder

# Applications: BERT and SoTA

To learn relationships between sentences, predict whether Sentence B is actual sentence that **proceeds** Sentence A, or a random sentence

**Sentence A** = The man went to the store.
**Sentence B** = He bought a gallon of milk.
**Label** = IsNextSentence

**Sentence A** = The man went to the store.
**Sentence B** = Penguins are flightless
**Label** = NotNextSentence

# Applications: BERT and SoTA



Use 30,000 WordPiece **vocabulary**

Each token is a **sum of three** embeddings

# **Applications:** BERT and SoTA

Multi-headed self **attention**

— Models context

Feed-forward layers

— Computes non-linear hierarchical features

Layer norm and **residuals**

— Makes training deep neural network (e.g., 12 layers possible)

**Positional Embeddings**

— Allows model to learn relative positioning

# **Applications:** BERT and SoTA



Pre-training

# **Applications:** BERT and SoTA



Pre-training

Fine-Tuning

# **Applications:** BERT and SoTA

# Applications: BERT and SoTA



Transfer: ASNQ Dataset

Adapt: Target Dataset

# **Applications:** Neural Image Captioning

# Applications: Neural Image Captioning

Image Embedding (VGGNet)



4096-dim

Convolution Layer + Non-Linearity   Pooling Layer   Convolution Layer + Non-Linearity   Pooling Layer   Fully-Connected MLP

# **Applications:** Neural Image Captioning

Image Embedding (VGGNet)



4096-dim

Convolution Layer
+ Non-Linearity

Pooling Layer

Convolution Layer
+ Non-Linearity

Pooling Layer

Fully-Connected MLP

# Applications: Neural Image Captioning

# Applications: Neural Image Captioning



* slide from Dhruv Batra

# **Applications:** Neural Image Captioning

**Good** results



A cat sitting on a
suitcase on the floor

A cat is sitting on a tree
branch

A dog is running in the
grass with a frisbee

A white teddy bear sitting in
the grass

Two people walking on
the beach with surfboards

A tennis player in action
on the court

Two giraffes standing in a
grassy field

A man riding a dirt bike on
a dirt track

# **Applications:** Neural Image Captioning

## **Failure** cases



*A woman is holding a cat in her hand*



*A person holding a computer mouse on a desk*



*A woman standing on a beach holding a surfboard*



*A bird is perched on a tree branch*



*A man in a baseball uniform throwing a ball*

# **Applications:** Image Captioning with Attention

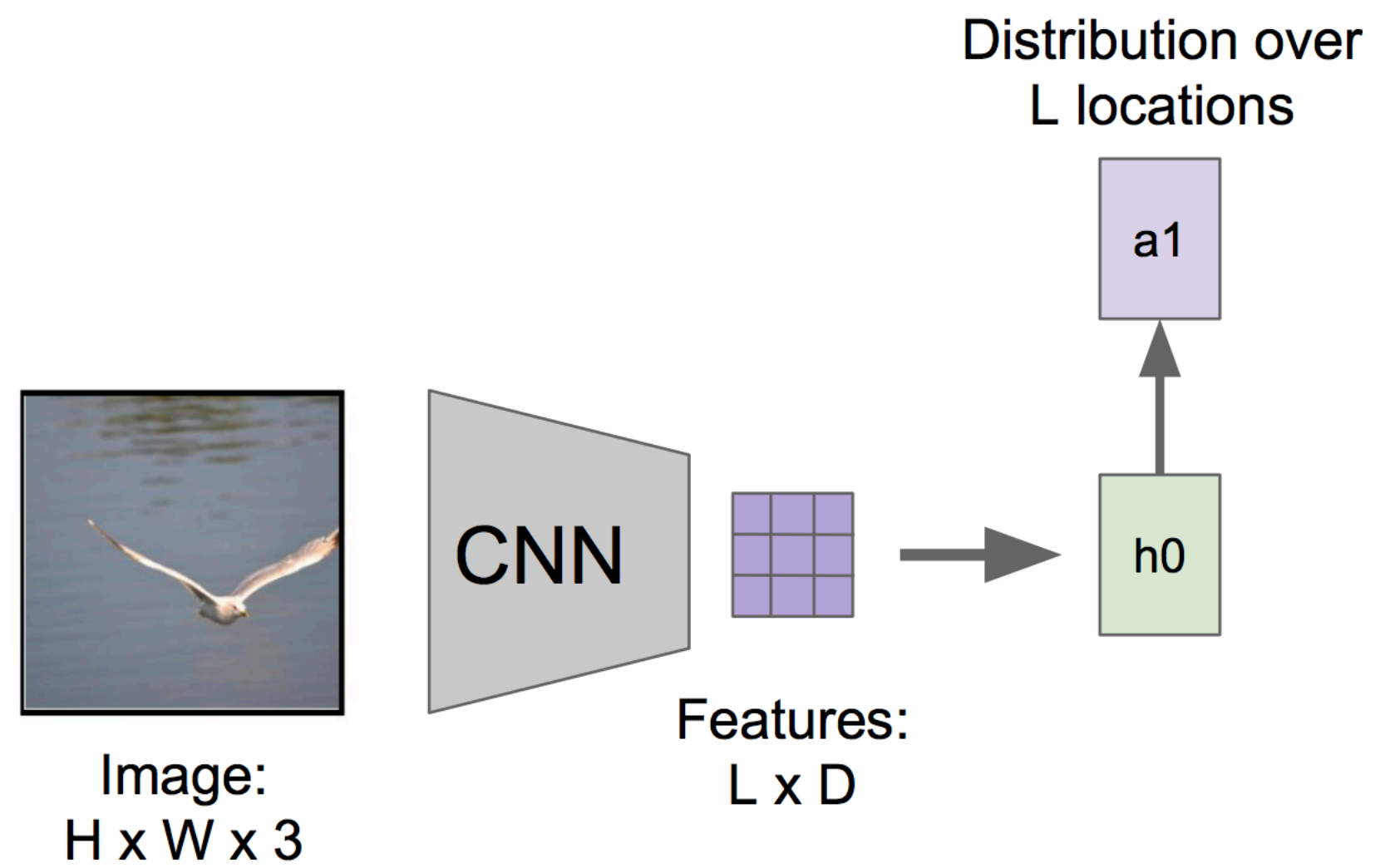RNN focuses its attention at a different spatial location
when generating each word



14x14 Feature Map

LSTM

A bird flying over a body of water

1. Input Image

2. Convolutional Feature Extraction

3. RNN with attention over the image

4. Word by word generation

# Applications: Image Captioning with Attention

Image:
H x W x 3

CNN

Features:
L x D

h0

# Applications: Image Captioning with Attention

Distribution over
L locations

a1

CNN

h0

Features:
L x D

Image:
H x W x 3

# Applications: Image Captioning with Attention

Distribution over L locations

a1

h0

CNN

Image:
H x W x 3

Features:
L x D

Weighted combination of features

Weighted features: D

z1

$$z = \sum_{i=1}^{L} p_i v_i$$

# **Applications:** Image Captioning with Attention

Distribution over
L locations

a1

h0 → h1

Features:
L x D

Weighted
features: D

z1    y1

Image:
H x W x 3

Weighted
combination
of features

First word

# **Applications:** Image Captioning with Attention

# **Applications:** Image Captioning with Attention

# Applications: Image Captioning with Attention

# **Applications:** Image Captioning with Attention

Soft attention

Hard attention

A bird flying over a body of water .

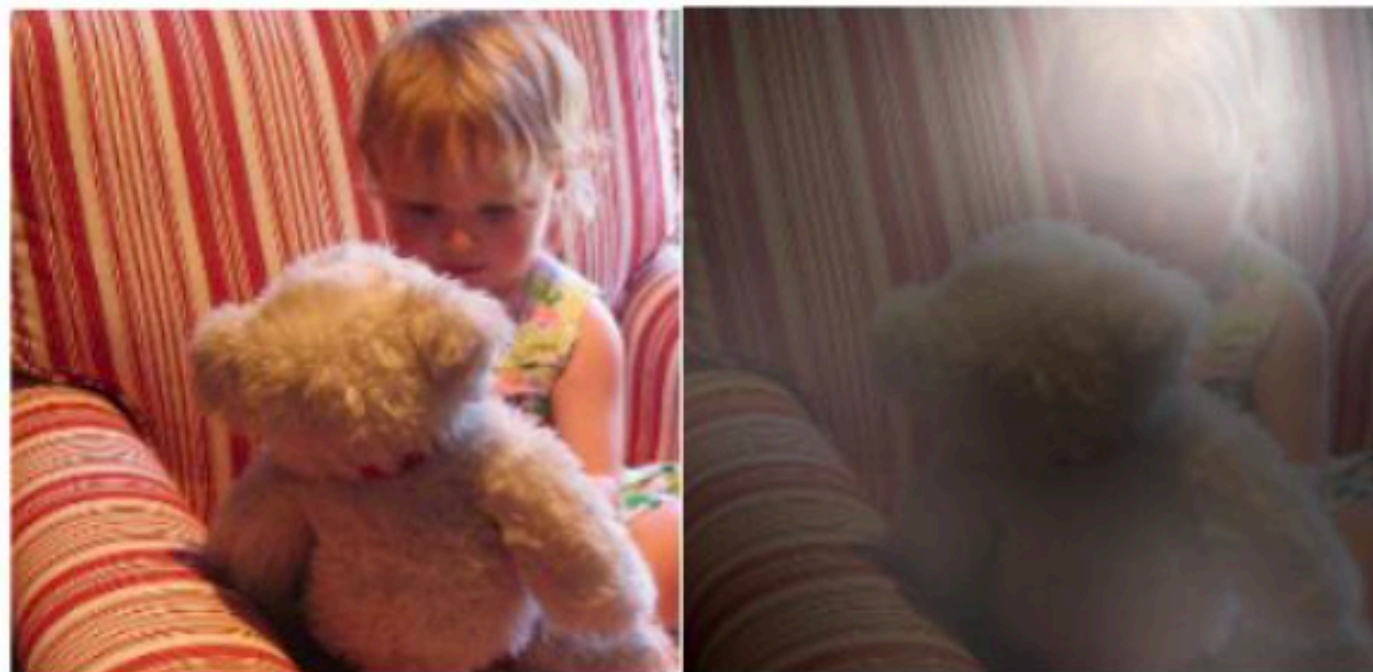**Good** results



A woman is throwing a <u>frisbee</u> in a park.

A <u>dog</u> is standing on a hardwood floor.

A <u>stop</u> sign is on a road with a mountain in the background.

A little <u>girl</u> sitting on a bed with a teddy bear.

A group of <u>people</u> sitting on a boat in the water.

A giraffe standing in a forest with <u>trees</u> in the background.

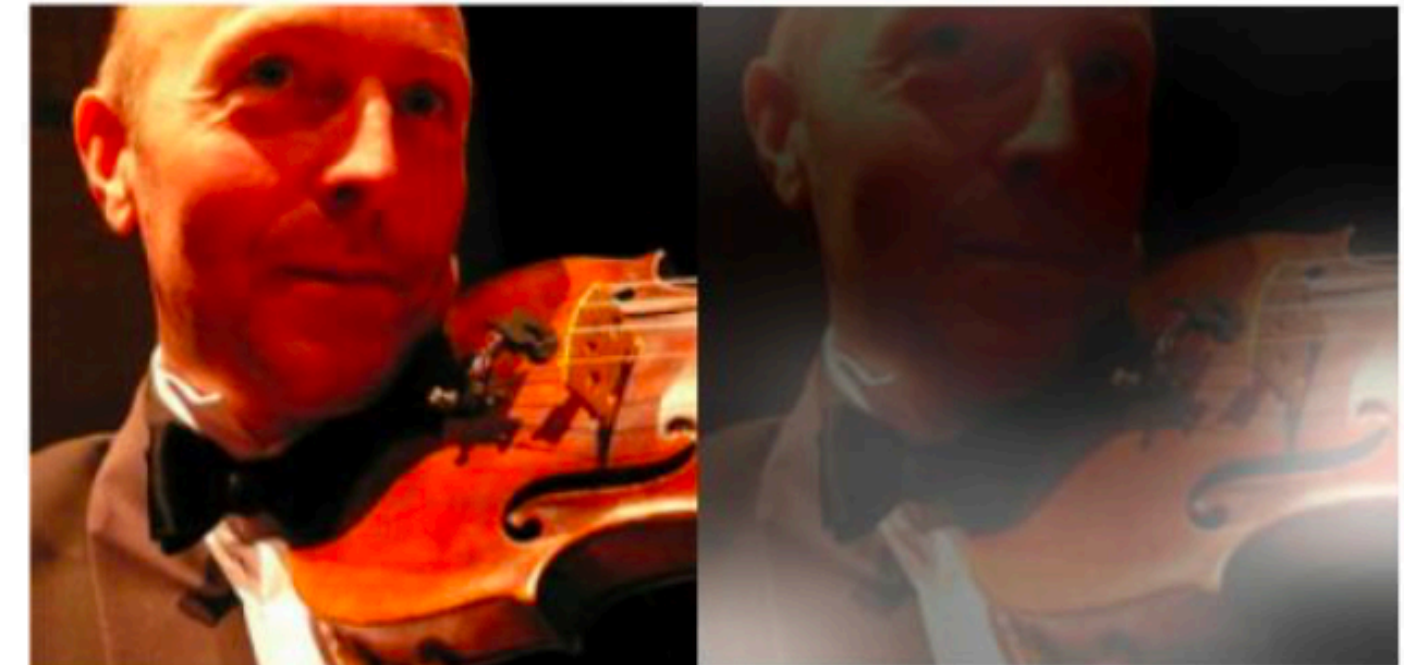# Applications: Image Captioning with Attention

[ Xu et al., ICML 2015 ]

**Failure** results



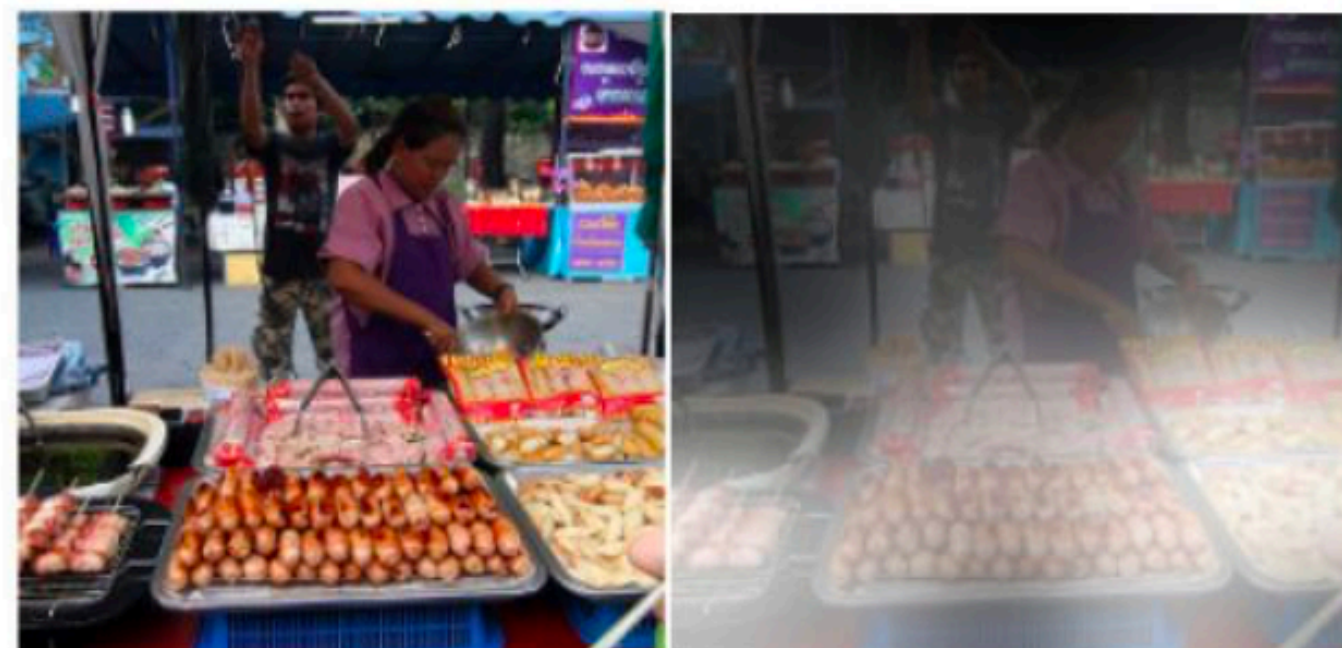A large white <u>bird</u> standing in a forest.

A woman holding a <u>clock</u> in her hand.

A man wearing a hat and a hat on a <u>skateboard</u>.

A person is standing on a beach with a <u>surfboard</u>.

A woman is sitting at a table with a large <u>pizza</u>.

A man is talking on his cell <u>phone</u> while another man watches.

* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**