

AN

Can you please write a sonnet for my upcoming Advanced Natural Language Processing class? i will be teaching about neural networks and NLP, transformers, syntax, morphology, and translation.



Here is a sonnet about NLP and neural networks:



In the realm of language, deep neural nets
Unlock the secrets of syntax and word forms,
With Transformers leading the way, it sets
A new standard for language norms.

Morphology and translation, once so tough,
Are now made simple by these powerful tools.
We train them on mountains of data, enough
To make language learning a breeze for fools.

But let us not forget the human touch,
For even machines can make mistakes.
In this class, we'll learn how to make them such
That language understanding they'll undertake.

So join me as we delve into this quest,
To unlock the secrets of language at its best.

CS678 - Advanced NLP

Intro:

Neural Language Models

Antonis Anastasopoulos

antonis@gmu.edu



Our goal today

What is **natural language processing**?

Specifically today:

- Class Logistics
- Neural Networks for NLP and Language Modeling

Hello, everyone!

Research Interests: NLP and AI

NLP for Low-Resource Languages

Machine Translation

Multilinguality and Cross-Lingual Learning

Fairness in NLP

Collaborations:

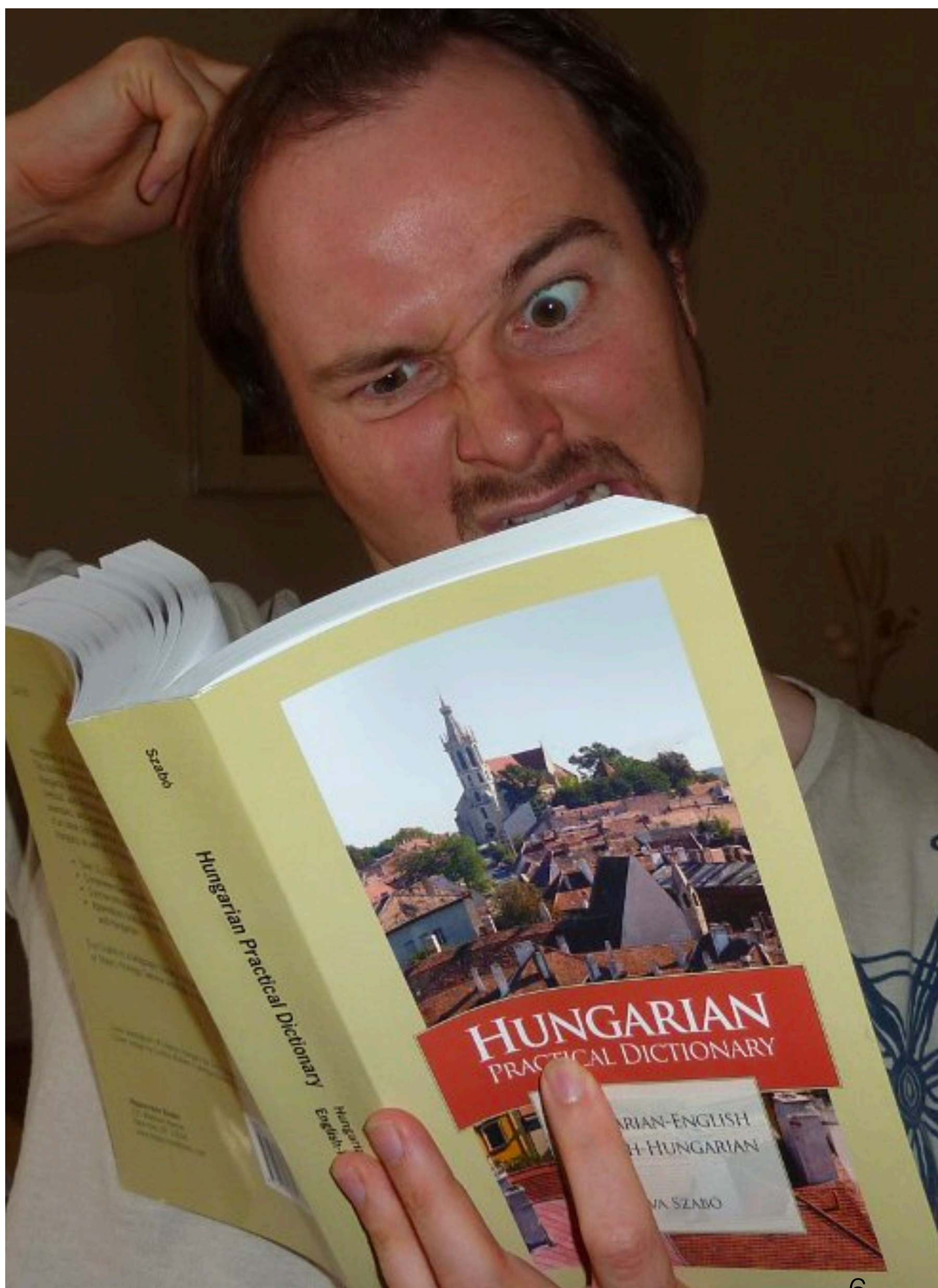
Carnegie Mellon University, U. Washington, Google, Amazon (AWS), Meta, U. of Notre Dame, Microsoft, Karya.



Antonis Anastasopoulos,
Asst. Prof.@GMU CS
Researcher @Archimedes AI



A bit about you



**Language
is Hard!**

Are These Sentences OK?

Are These Sentences OK?

Jane went to the store.

Are These Sentences OK?

Jane went to the store.

store to Jane went the.

Are These Sentences OK?

Jane went to the store.

store to Jane went the.

Jane went store.

Are These Sentences OK?

Jane went to the store.

store to Jane went the.

Jane went store.

Jane goed to the store.

Are These Sentences OK?

Jane went to the store.

store to Jane went the.

Jane went store.

Jane goed to the store.

The store went to Jane.

Are These Sentences OK?

Jane went to the store.

store to Jane went the.

Jane went store.

Jane goed to the store.

The store went to Jane.

The food truck went to Jane.

Engineering Solutions

Jane went to the store.

store to Jane went the.

Jane went store.

Jane goed to the store.

The store went to Jane.

The food truck went to Jane.

Engineering Solutions

Jane went to the store.

store to Jane went the.

Jane went store.

Jane goed to the store.

The store went to Jane.

The food truck went to Jane.

} Create a grammar of
the language

Engineering Solutions

Jane went to the store.

store to Jane went the.

Jane went store.

Jane goed to the store.

The store went to Jane.

The food truck went to Jane.

} Create a grammar of
the language

} Consider
morphology and exceptions

Engineering Solutions

Jane went to the store.

store to Jane went the.

Jane went store.

Jane goed to the store.

The store went to Jane.

The food truck went to Jane.

} Create a grammar of
the language

} Consider
morphology and exceptions

} Semantic categories,
preferences

Engineering Solutions

Jane went to the store.

store to Jane went the.

Jane went store.

Jane goed to the store.

The store went to Jane.

The food truck went to Jane.

} Create a grammar of
the language

} Consider
morphology and exceptions

} Semantic categories,
preferences

} And their exceptions

Are These Sentences OK?

ジェインは店へ行った。
は店行ったジェインは。
ジェインは店へ行た。
店はジェインへ行った。
屋台はジェインのところへ行った。

Μπορείτε να διαβάσετε αυτήν
την πρόταση;

Potete leggere questa frase?

इस वाक्य क्या आप को पढ़ सकते हैं?

mungawerenge chiganizo ichi?

Phenomena to Handle

Morphology

Syntax

Semantics/World Knowledge

Discourse

Pragmatics

Multilinguality

Neural Nets for NLP

Neural Nets for NLP

Neural nets are a tool to do hard things!

Neural Nets for NLP

Neural nets are a tool to do hard things!

Combined with lots of data and compute, they can approximate any function!



**What do you think of
when you think of NLP?**

What will you learn?

What will you learn?

1. The generics of how large pre-trained neural language models are trained and operate

What will you learn?

1. The generics of how large pre-trained neural language models are trained and operate
2. Understand the limitations of current technologies

What will you learn?

1. The generics of how large pre-trained neural language models are trained and operate
2. Understand the limitations of current technologies
3. Learn about the harms associated with their use and ways to mitigate them

How will you learn?
(syllabus highlights)

Syllabus Highlights

Syllabus Highlights

Mostly lectures, but

I'll try to have some sort of interactive element every day

Syllabus Highlights

Mostly lectures, but

I'll try to have some sort of interactive element every day

Syllabus Highlights

Mostly lectures, but

I'll try to have some sort of interactive element every day

- I will provide additional readings for anyone interested

Lectures

Lectures

You should ask lots of questions

Lectures

You should ask lots of questions

- interrupting (by raising a hand) to ask your question is *strongly* encouraged

Lectures

You should ask lots of questions

- interrupting (by raising a hand) to ask your question is *strongly* encouraged
- Asking questions later (or in real time)

Lectures

You should ask lots of questions

- interrupting (by raising a hand) to ask your question is *strongly* encouraged
- Asking questions later (or in real time)

Lectures

You should ask lots of questions

- interrupting (by raising a hand) to ask your question is *strongly* encouraged
- Asking questions later (or in real time)

Interaction improves learning!

Logistics

We will use Discord for

Announcements

Distributing course materials before/after class

Language Models

Calculating the Probability of a Sentence


Jane went to the store .

$$P(X) = \prod_{i=1}^n P(x_i)$$

Calculating the Probability of a Sentence

Jane went to the store .

$$P(X) = \prod_{i=1}^n P(x_i)$$

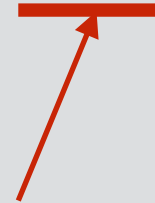

Unigram

Calculating the Probability of a Sentence

Jane went to the store .

$$P(X) = \prod_{i=1}^n P(x_i)$$

Unigram



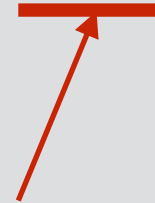
$$P(\text{Jane went to the store}) = P(\text{Jane}) \times P(\text{went}) \times P(\text{to}) \times P(\text{the}) \times P(\text{store}) \times P(.).$$

Calculating the Probability of a Sentence

Jane went to the store .

$$P(X) = \prod_{i=1}^n P(x_i)$$

Unigram



$$P(\text{Jane went to the store}) = P(\text{Jane}) \times P(\text{went}) \times P(\text{to}) \times P(\text{the}) \times P(\text{store}) \times P(.).$$

But word order and context matters!

Calculating the Probability of a Sentence

$$P(X) = \prod_{i=1}^I P(x_i \mid x_1, \dots, x_{i-1})$$

Next Word Context

Calculating the Probability of a Sentence

$$P(X) = \prod_{i=1}^I P(x_i \mid x_1, \dots, x_{i-1})$$

The diagram illustrates the components of the probability formula. A red horizontal line is positioned under the variable x_i in the numerator of the product term. A red arrow points from the text "Next Word" below to this red line. A blue horizontal line is positioned under the entire denominator x_1, \dots, x_{i-1} . A blue arrow points from the text "Context" below to this blue line.

$$\begin{aligned} P(\text{Jane went to the store}) &= P(\text{Jane} \mid \langle s \rangle) \times P(\text{went} \mid \text{Jane}) \times \\ &\quad P(\text{to} \mid \text{went}) \times P(\text{the} \mid \text{to}) \times \\ &\quad P(\text{store} \mid \text{the}) \times P(. \mid \text{store}) \\ &\quad P(\langle /s \rangle \mid .) \end{aligned}$$

Calculating the Probability of a Sentence

$$P(X) = \prod_{i=1}^I P(x_i \mid x_1, \dots, x_{i-1})$$

Next Word Context

The big problem: How do we predict

$$P(x_i \mid x_1, \dots, x_{i-1})$$

?!?!

Count-based Language Models

Count-based Language Models

- Count up the frequency and divide:

$$P_{ML}(x_i \mid x_{i-n+1}, \dots, x_{i-1}) := \frac{c(x_{i-n+1}, \dots, x_i)}{c(x_{i-n+1}, \dots, x_{i-1})}$$

Count-based Language Models

- Count up the frequency and divide:

$$P_{ML}(x_i \mid x_{i-n+1}, \dots, x_{i-1}) := \frac{c(x_{i-n+1}, \dots, x_i)}{c(x_{i-n+1}, \dots, x_{i-1})}$$

Corpus:

The cat sat on the mat . A mouse ate some cheese .
A dog chased the cat . The mouse ran under a mat .

Count-based Language Models

- Count up the frequency and divide:

$$P_{ML}(x_i | x_{i-n+1}, \dots, x_{i-1}) := \frac{c(x_{i-n+1}, \dots, x_i)}{c(x_{i-n+1}, \dots, x_{i-1})}$$

Corpus:

The cat sat on the mat . A mouse ate some cheese .
A dog chased the cat . The mouse ran under a mat .

$$p(chased | dog) = ? \quad p(cat | the) = ? \quad p(the | \langle s \rangle) = ?$$

Count-based Language Models

- Count up the frequency and divide:

$$P_{ML}(x_i | x_{i-n+1}, \dots, x_{i-1}) := \frac{c(x_{i-n+1}, \dots, x_i)}{c(x_{i-n+1}, \dots, x_{i-1})}$$

Corpus:

The cat sat on the mat . A mouse ate some cheese .

A dog chased the cat . The mouse ran under a mat .

$$p(\textit{chased} | \textit{dog}) = \frac{1}{1} = 1 \quad p(\textit{cat} | \textit{the}) = \frac{1}{4} = 0.25 \quad p(\textit{the} | \langle s \rangle) = 0.5$$

Count-based Language Models

- Count up the frequency and divide:

$$P_{ML}(x_i | x_{i-n+1}, \dots, x_{i-1}) := \frac{c(x_{i-n+1}, \dots, x_i)}{c(x_{i-n+1}, \dots, x_{i-1})}$$

Corpus:

The cat sat on the mat . A mouse ate some cheese .
A dog chased the cat . The mouse ran under a mat .

$p(\text{A cat chased the mouse .}) = ?$

Count-based Language Models

Corpus:

The cat sat on the mat . A mouse ate some cheese .
A dog chased the cat . The mouse ran under a mat .

$$\begin{aligned} p(\text{A cat chased the mouse .}) = & \\ & p(\langle s \rangle | A) \times \\ & p(\text{cat} | a) \times \\ & p(\text{chased} | \text{cat}) \times \\ & p(\text{the} | \text{chased}) \times \\ & p(\text{mouse} | \text{the}) \times \\ & p(. | \text{mouse}) \end{aligned}$$

Count-based Language Models

Corpus:

The cat sat on the mat . A mouse ate some cheese .
A dog chased the cat . The mouse ran under a mat .

$p(\text{A cat chased the mouse .}) =$

$p(\langle s \rangle | A) \times$



$p(\text{cat} | a) \times$

$p(\text{chased} | \text{cat}) \times$

$p(\text{the} | \text{chased}) \times$



$p(\text{mouse} | \text{the}) \times$



$p(. | \text{mouse})$

Count-based Language Models

- Count up the frequency and divide:

$$P_{ML}(x_i | x_{i-n+1}, \dots, x_{i-1}) := \frac{c(x_{i-n+1}, \dots, x_i)}{c(x_{i-n+1}, \dots, x_{i-1})}$$

- Add smoothing to deal with zero counts:

$$p(x_i | x_{i-n+1:i-1}) = \frac{c(x_{i-n+1:i}) + \alpha}{c(x_{i-n+1:i-1}) + \alpha |V|}$$

Count-based Language Models

Corpus:

The cat sat on the mat . A mouse ate some cheese .
A dog chased the cat . The mouse ran under a mat .

$$|V| = |\{the, a, cat, sat, \dots\}| = 15 \quad \alpha = 1$$

Count-based Language Models

Corpus:

The cat sat on the mat . A mouse ate some cheese .
A dog chased the cat . The mouse ran under a mat .

$$|V| = |\{the, a, cat, sat, \dots\}| = 15 \quad \alpha = 1$$

$$p(\text{A cat chased the mouse .}) =$$

$$\begin{aligned} & p(\langle s \rangle \| A) \times \checkmark \\ & p(cat | a) \times \checkmark \\ & p(chased | cat) \times \checkmark \\ & p(the | chased) \times \checkmark \\ & p(mouse | the) \times \checkmark \\ & p(. | mouse) \checkmark \end{aligned}$$

An Alternative:
Featurized Log-Linear Models

An Alternative: Featurized Models

An Alternative: Featurized Models

- Calculate features of the context

An Alternative: Featurized Models

- Calculate features of the context
- Based on the features, calculate probabilities

An Alternative: Featurized Models

- Calculate features of the context
- Based on the features, calculate probabilities
- Optimize feature weights using gradient descent, etc.

Example:

Previous words: "giving a"

Example:

Previous words: “giving a”

a
the
talk
gift
hat
...

Words we're
predicting

Example:

Previous words: "giving a"

a	3.0
the	2.5
talk	-0.2
gift	0.1
hat	1.2
...	...

Words we're
predicting

How likely
are they?

Example:

Previous words: "giving a"

a
the
talk
gift
hat
...

$$b = \begin{pmatrix} 3.0 \\ 2.5 \\ -0.2 \\ 0.1 \\ 1.2 \\ \dots \end{pmatrix}$$

$$w_{1,a} = \begin{pmatrix} -6.0 \\ -5.1 \\ 0.2 \\ 0.1 \\ 0.5 \\ \dots \end{pmatrix}$$

Words we're
predicting

How likely
are they?

How likely
are they
given prev.
word is "a"?

Example:

Previous words: "giving a"

a
the
talk
gift
hat
...

$$b = \begin{pmatrix} 3.0 \\ 2.5 \\ -0.2 \\ 0.1 \\ 1.2 \\ \dots \end{pmatrix}$$

$$W_{1,a} = \begin{pmatrix} -6.0 \\ -5.1 \\ 0.2 \\ 0.1 \\ 0.5 \\ \dots \end{pmatrix}$$

$$W_{2,giving} = \begin{pmatrix} -0.2 \\ -0.3 \\ 1.0 \\ 2.0 \\ -1.2 \\ \dots \end{pmatrix}$$

Words we're predicting

How likely are they?

How likely are they given prev. word is "a"?

How likely are they given 2nd prev. word is "giving"?

Example:

Previous words: "giving a"

a
the
talk
gift
hat
...

$$b = \begin{pmatrix} 3.0 \\ 2.5 \\ -0.2 \\ 0.1 \\ 1.2 \\ \dots \end{pmatrix}$$

$$w_{1,a} = \begin{pmatrix} -6.0 \\ -5.1 \\ 0.2 \\ 0.1 \\ 0.5 \\ \dots \end{pmatrix}$$

$$w_{2,giving} = \begin{pmatrix} -0.2 \\ -0.3 \\ 1.0 \\ 2.0 \\ -1.2 \\ \dots \end{pmatrix}$$

$$s = \begin{pmatrix} -3.2 \\ -2.9 \\ 1.0 \\ 2.2 \\ 0.6 \\ \dots \end{pmatrix}$$

Words we're predicting

How likely are they?

How likely are they given prev. word is "a"?

How likely are they given 2nd prev. word is "giving"?

Total score

Softmax

- Convert scores into probabilities by taking the exponent and normalizing (softmax)

Softmax

- Convert scores into probabilities by taking the exponent and normalizing (softmax)

$$P(x_i \mid x_{i-n+1}^{i-1}) = \frac{e^{s(x_i \mid x_{i-n+1}^{i-1})}}{\sum_{\tilde{x}_i} e^{s(\tilde{x}_i \mid x_{i-n+1}^{i-1})}}$$

Softmax

- Convert scores into probabilities by taking the exponent and normalizing (softmax)

$$P(x_i | x_{i-n+1}^{i-1}) = \frac{e^{s(x_i | x_{i-n+1}^{i-1})}}{\sum_{\tilde{x}_i} e^{s(\tilde{x}_i | x_{i-n+1}^{i-1})}}$$

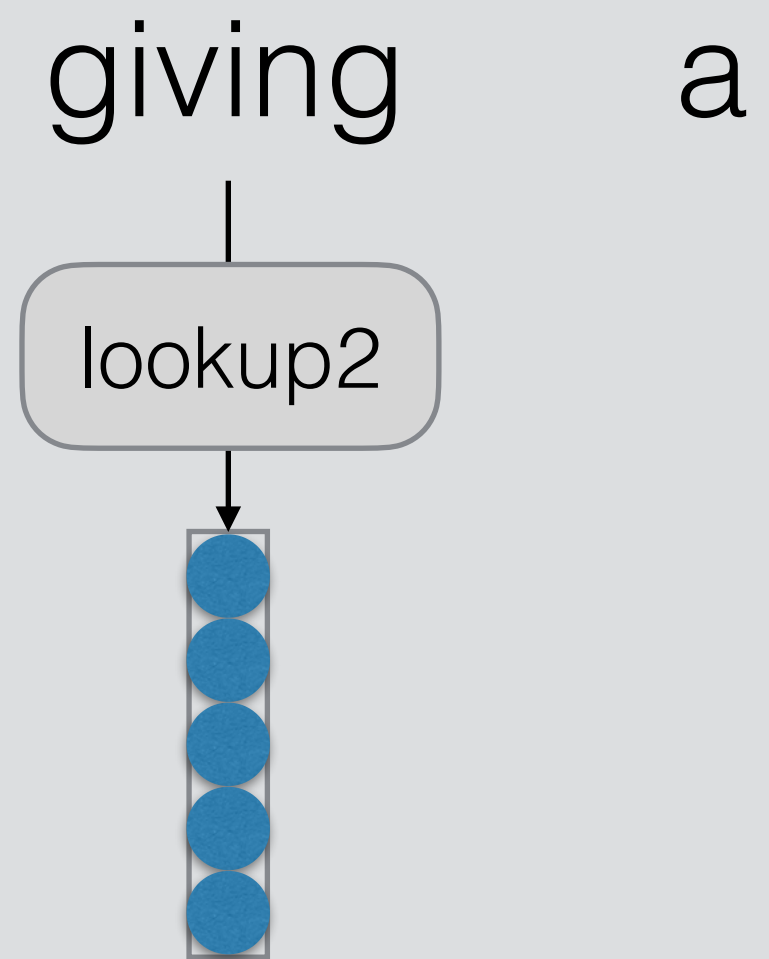
$$s = \begin{pmatrix} -3.2 \\ -2.9 \\ 1.0 \\ 2.2 \\ 0.6 \\ \dots \end{pmatrix} \longrightarrow p = \begin{pmatrix} 0.002 \\ 0.003 \\ 0.329 \\ 0.444 \\ 0.090 \\ \dots \end{pmatrix}$$

A Computation Graph View

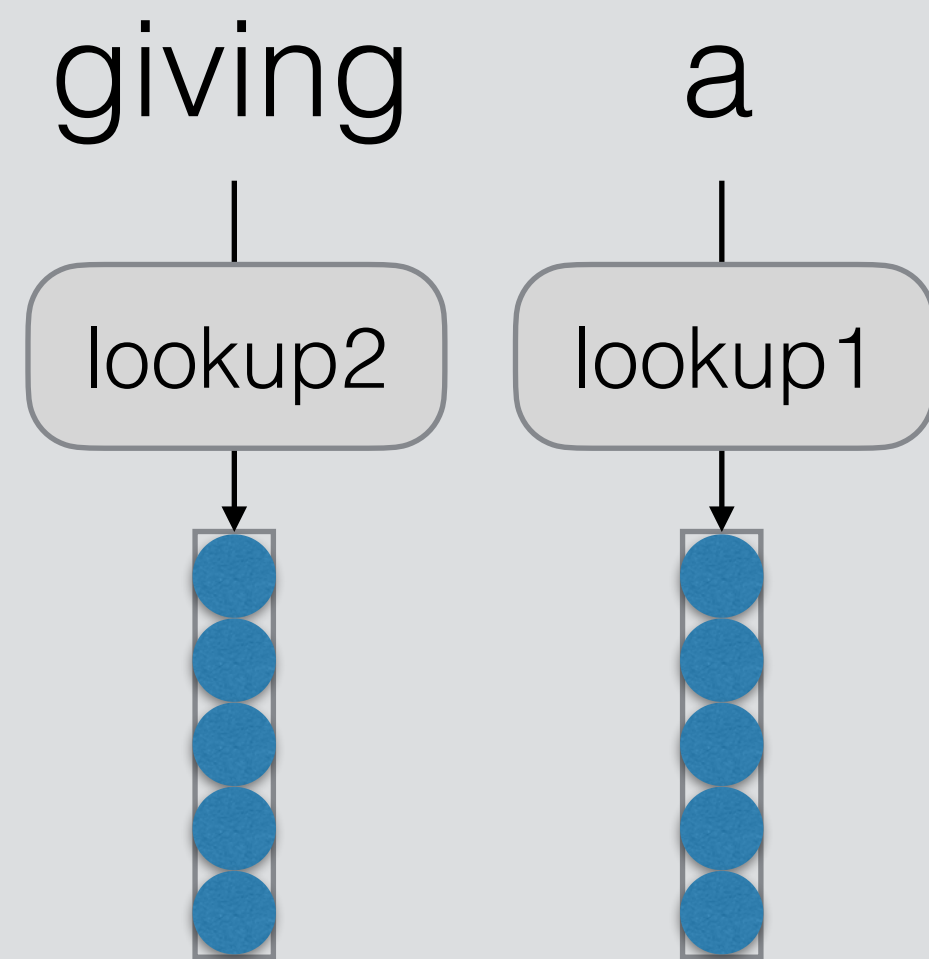
giving a

Each vector is size of output vocabulary

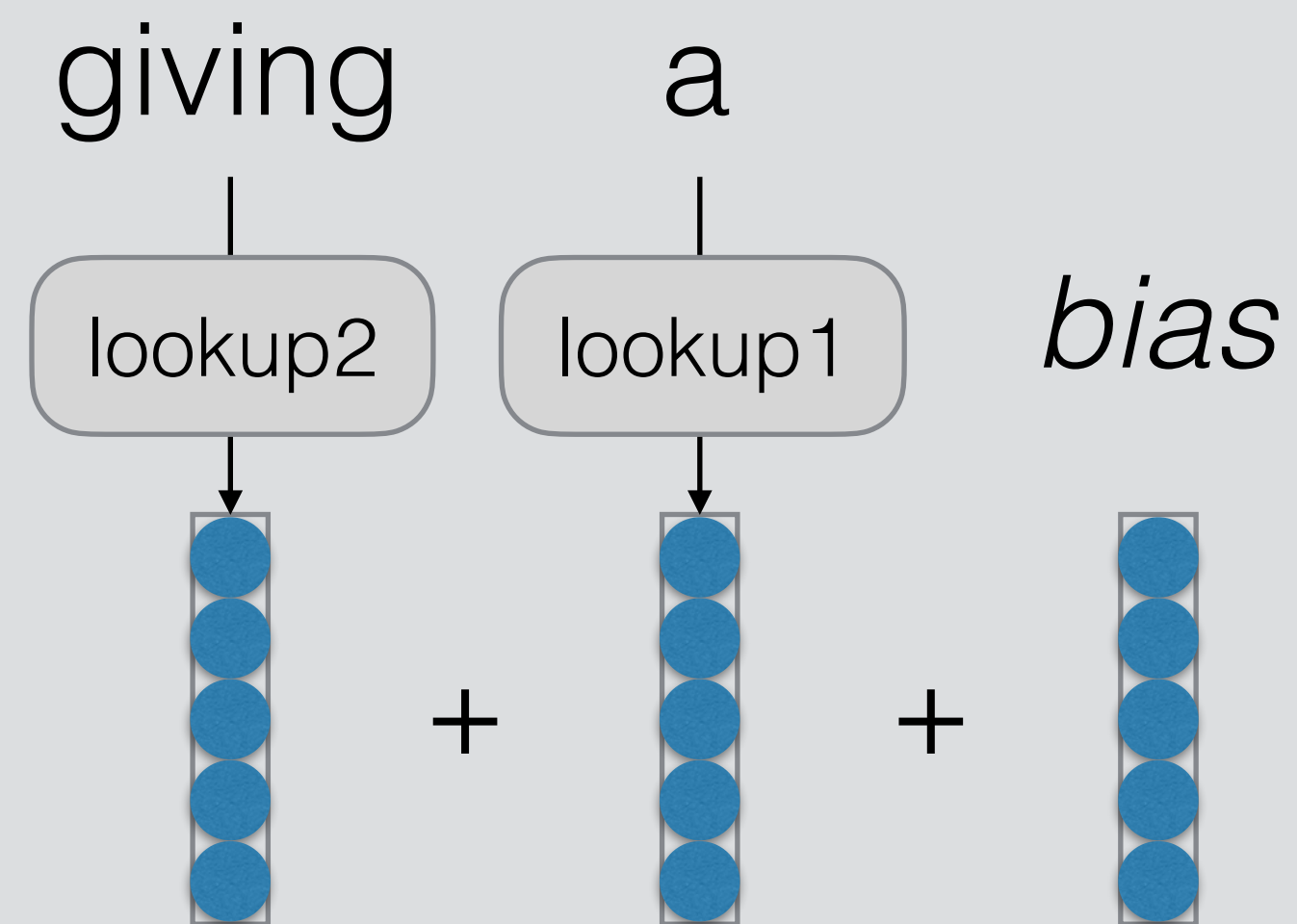
A Computation Graph View



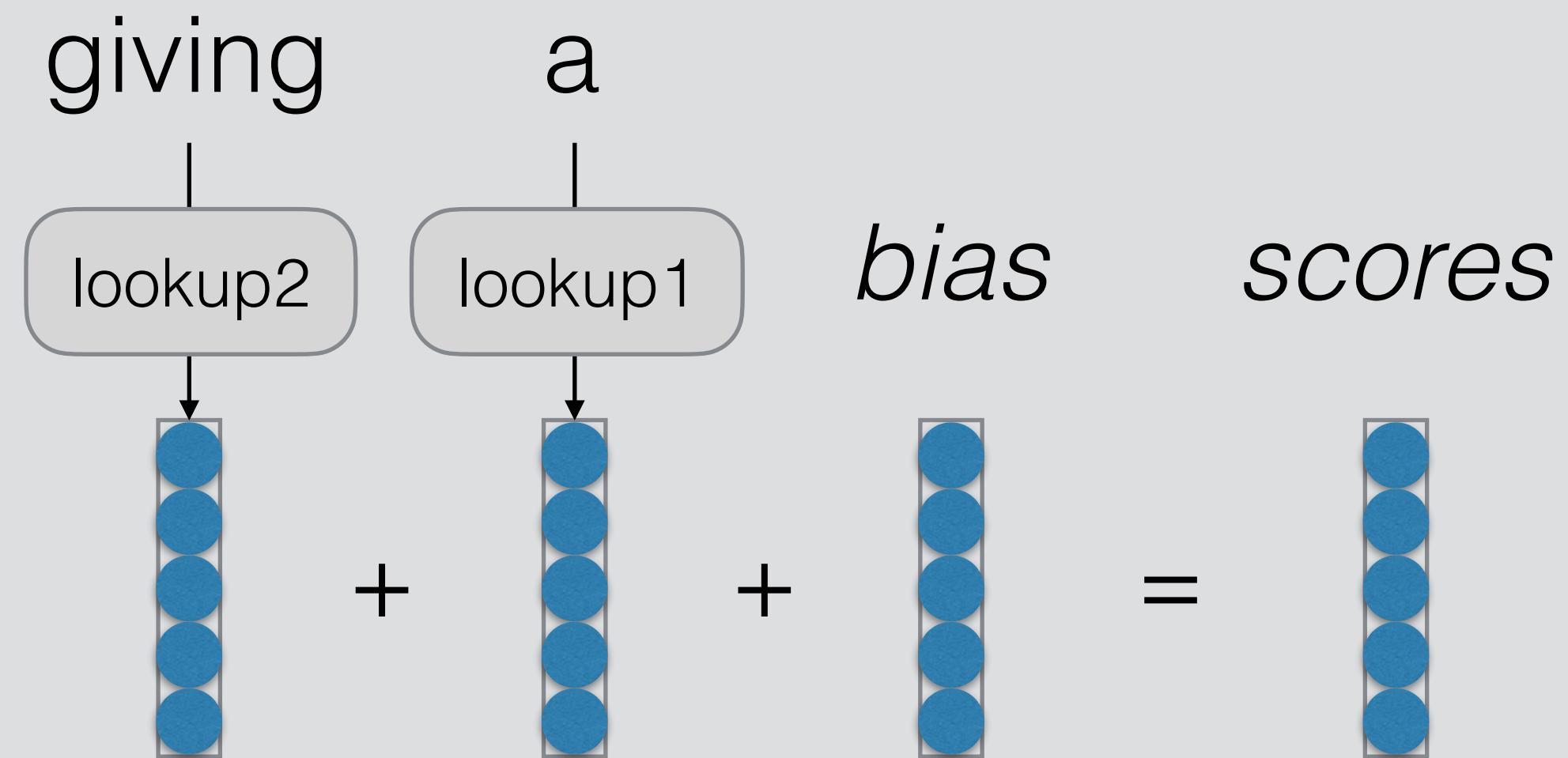
A Computation Graph View



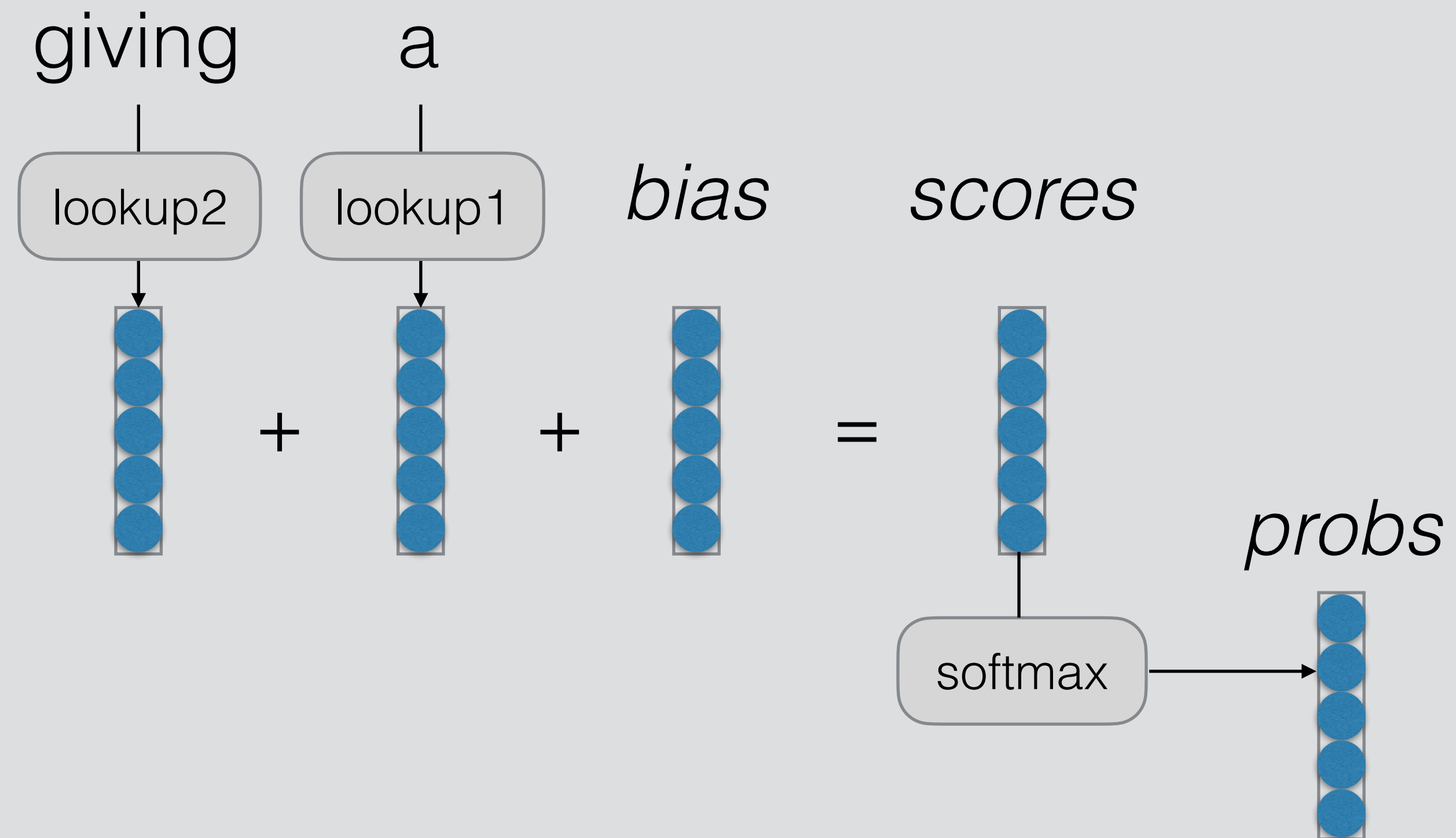
A Computation Graph View



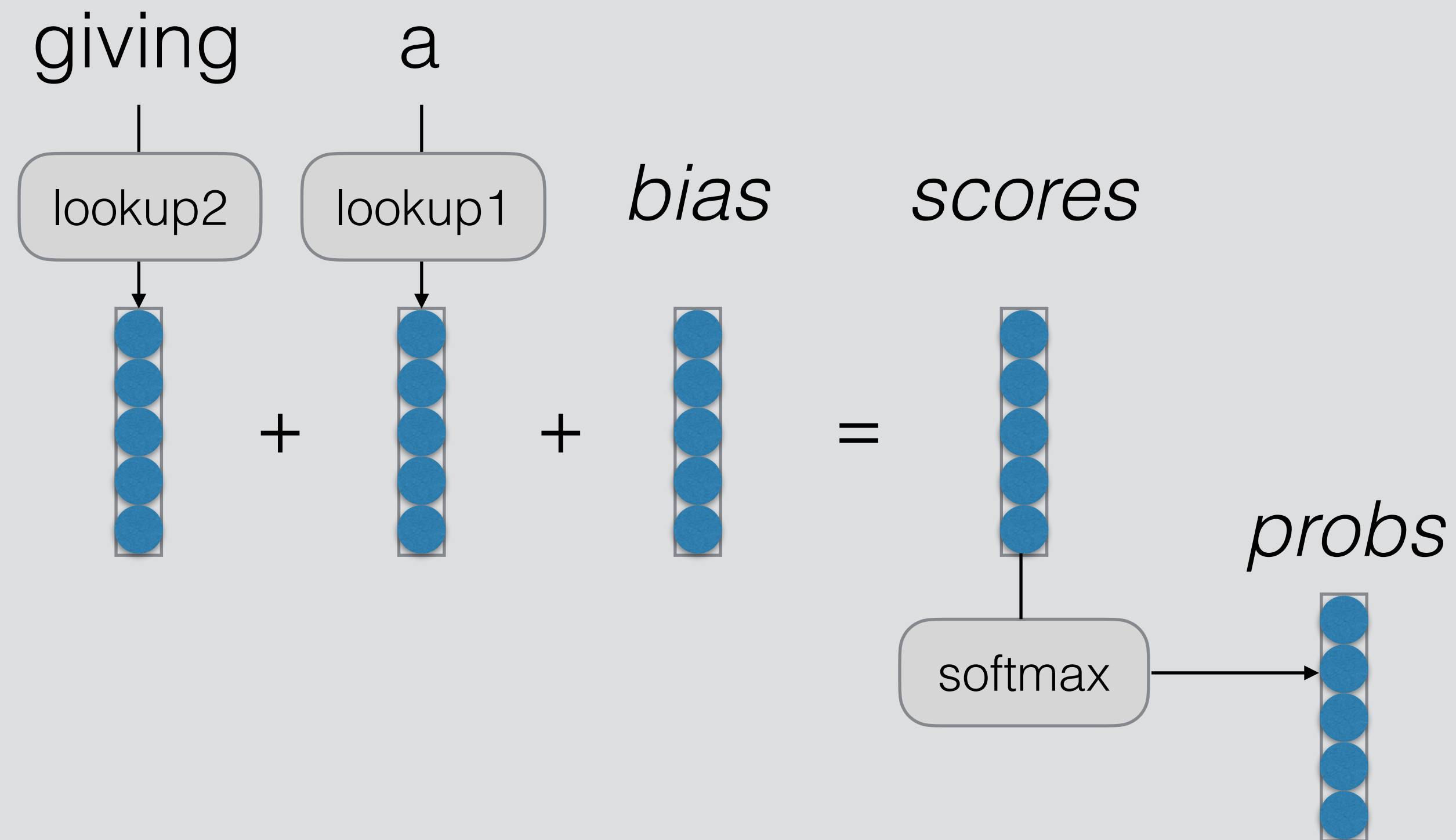
A Computation Graph View



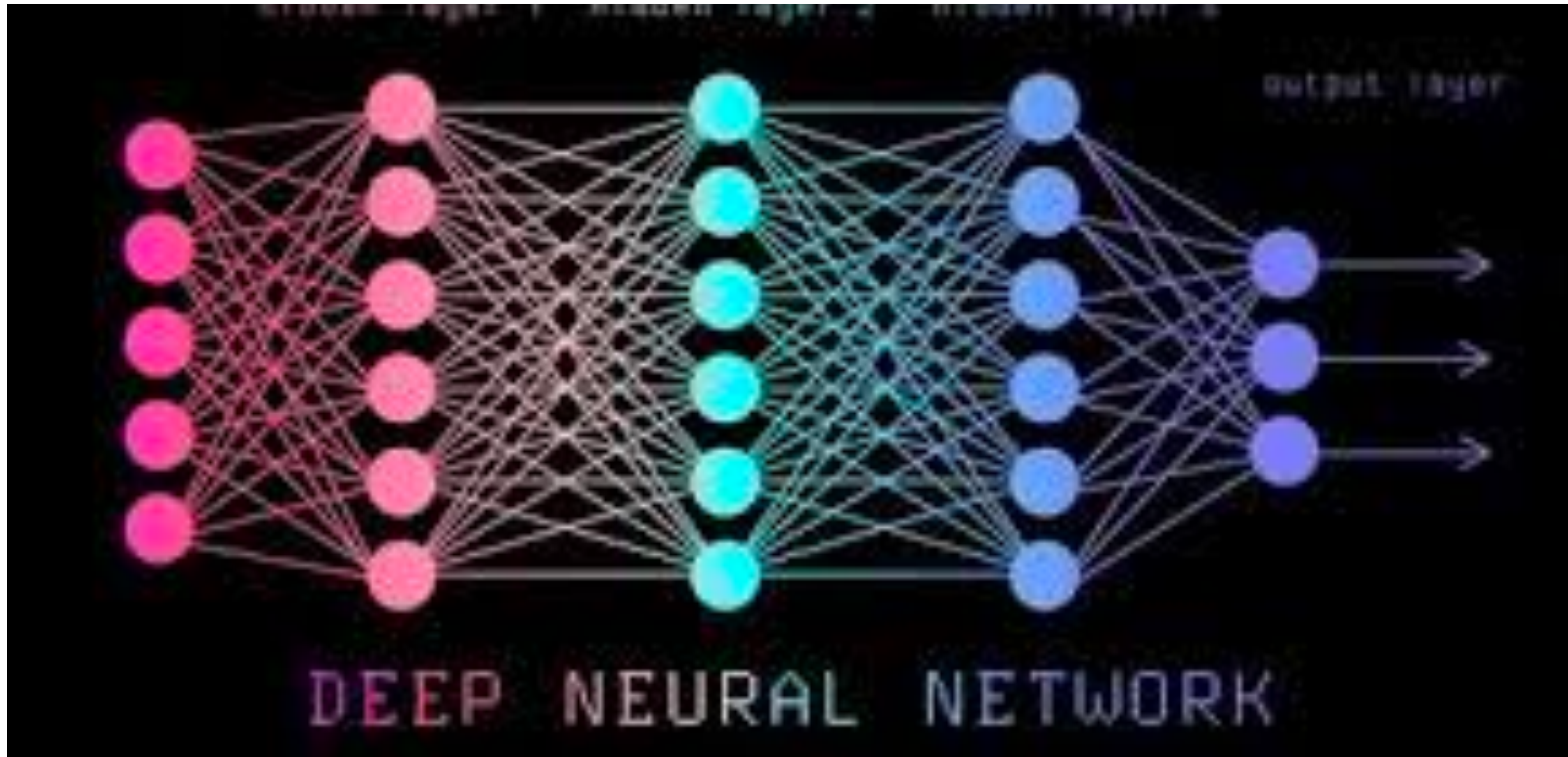
A Computation Graph View



A Computation Graph View



Each vector is size of output vocabulary



Neural Networks: A Tool for Doing Hard Things

An Example Prediction Problem: Sentence Classification

An Example Prediction Problem: Sentence Classification

I hate this movie

An Example Prediction Problem: Sentence Classification

I hate this movie

I love this movie

An Example Prediction Problem: Sentence Classification

I hate this movie

very good
good
neutral
bad
very bad

I love this movie

An Example Prediction Problem: Sentence Classification

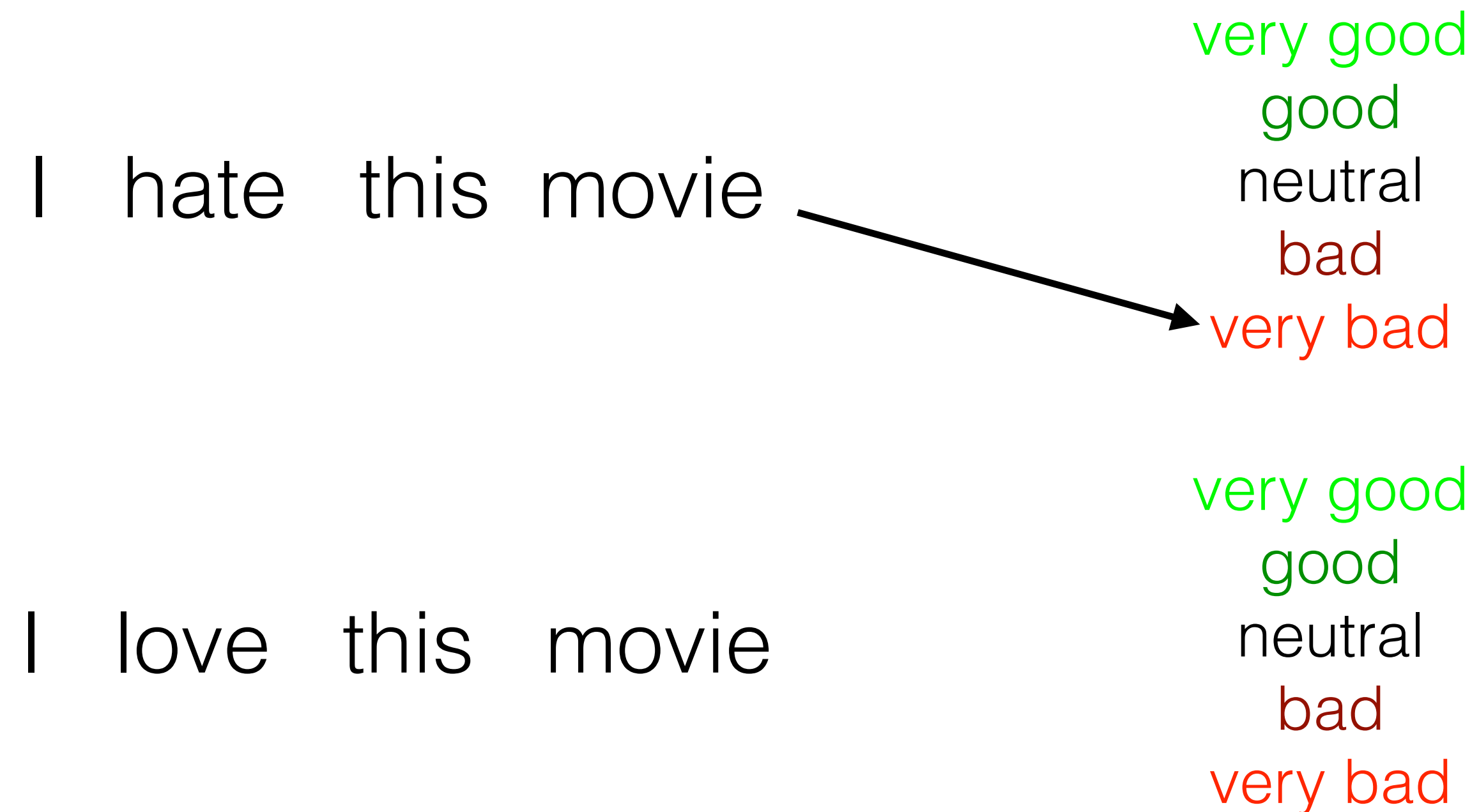
I hate this movie

very good
good
neutral
bad
very bad

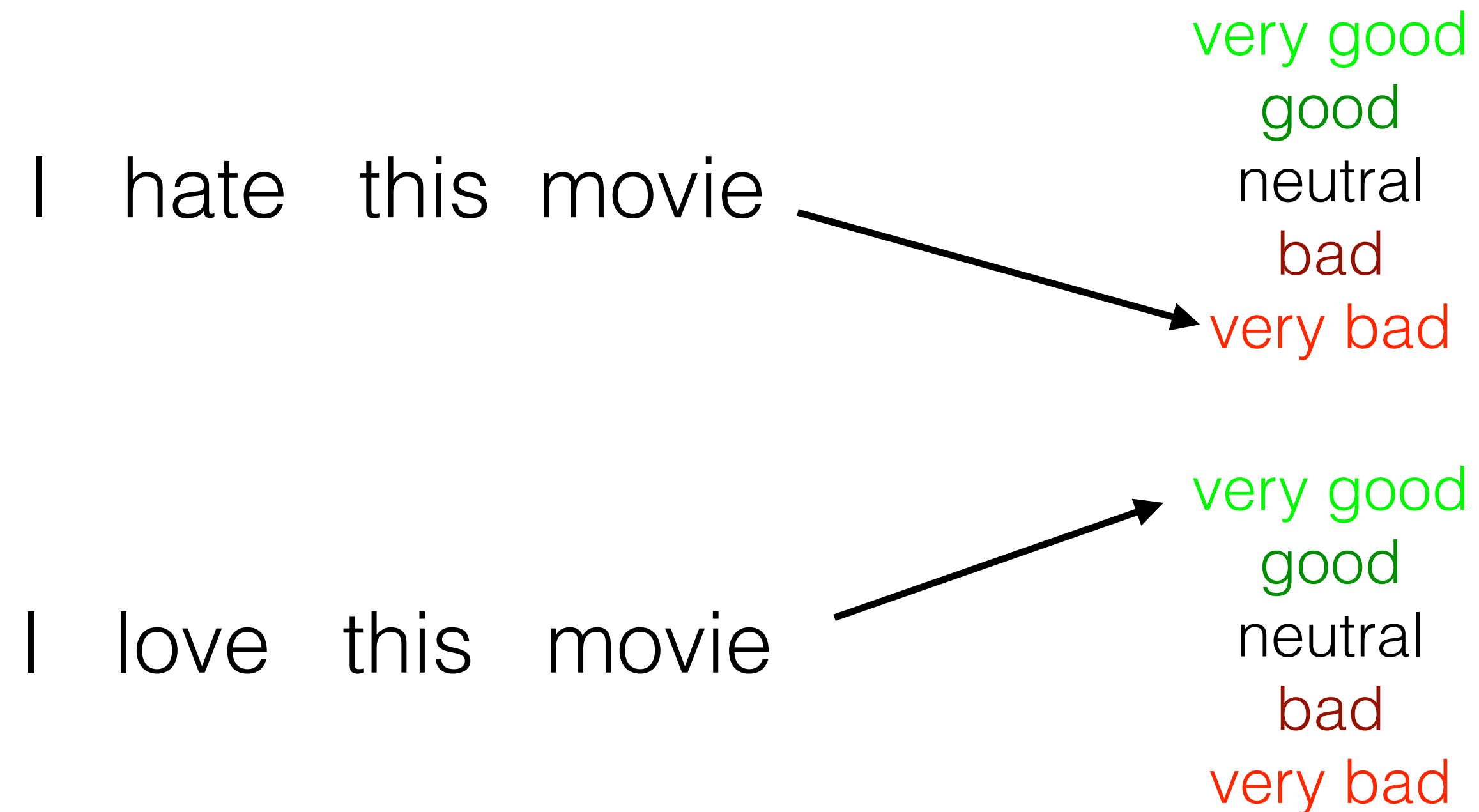
I love this movie

very good
good
neutral
bad
very bad

An Example Prediction Problem: Sentence Classification



An Example Prediction Problem: Sentence Classification



A First Try:

Bag of Words (BOW)

Each word has a vector of weights for each tag

I hate this movie

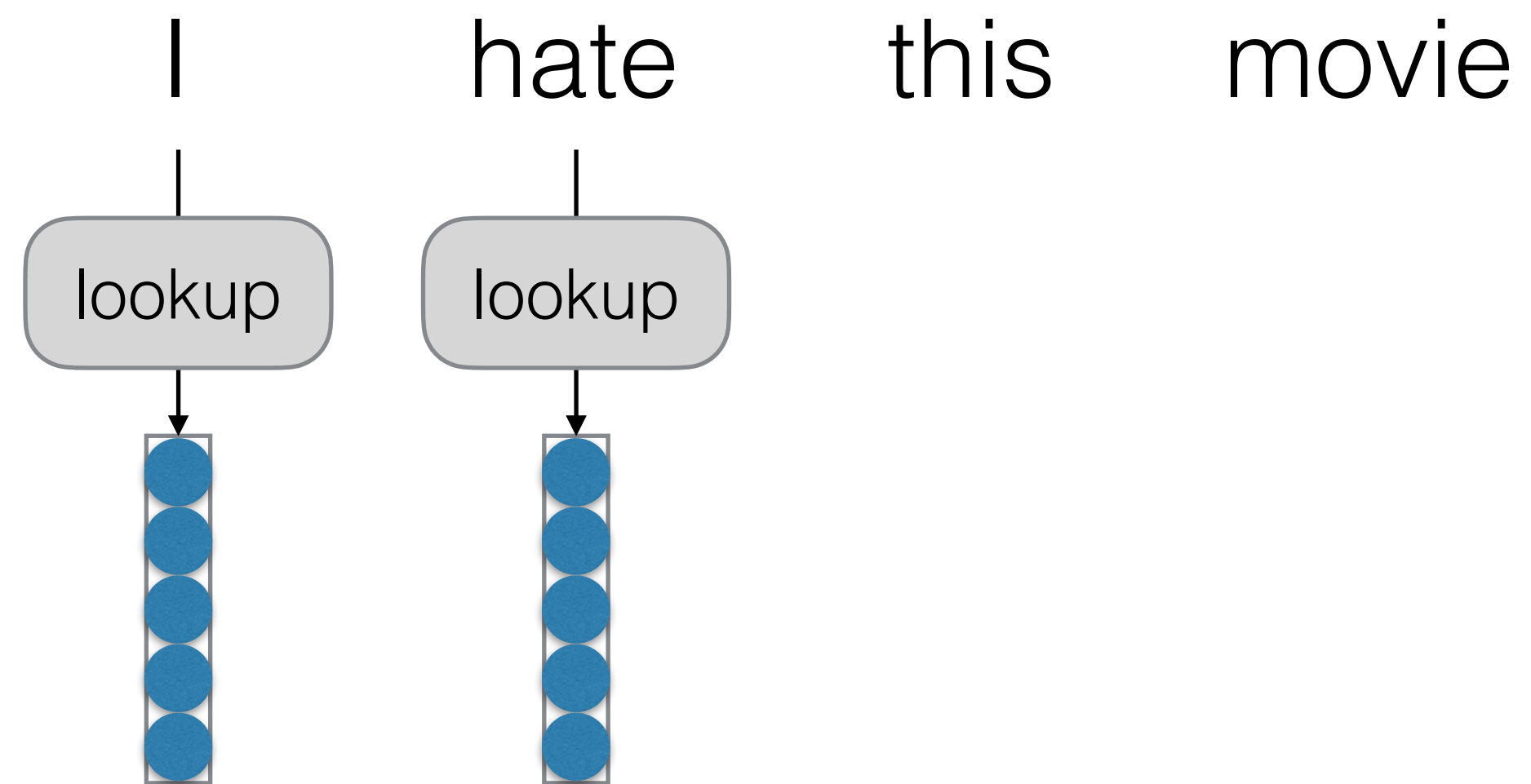
A First Try: Bag of Words (BOW)

Each word has a vector of weights for each tag



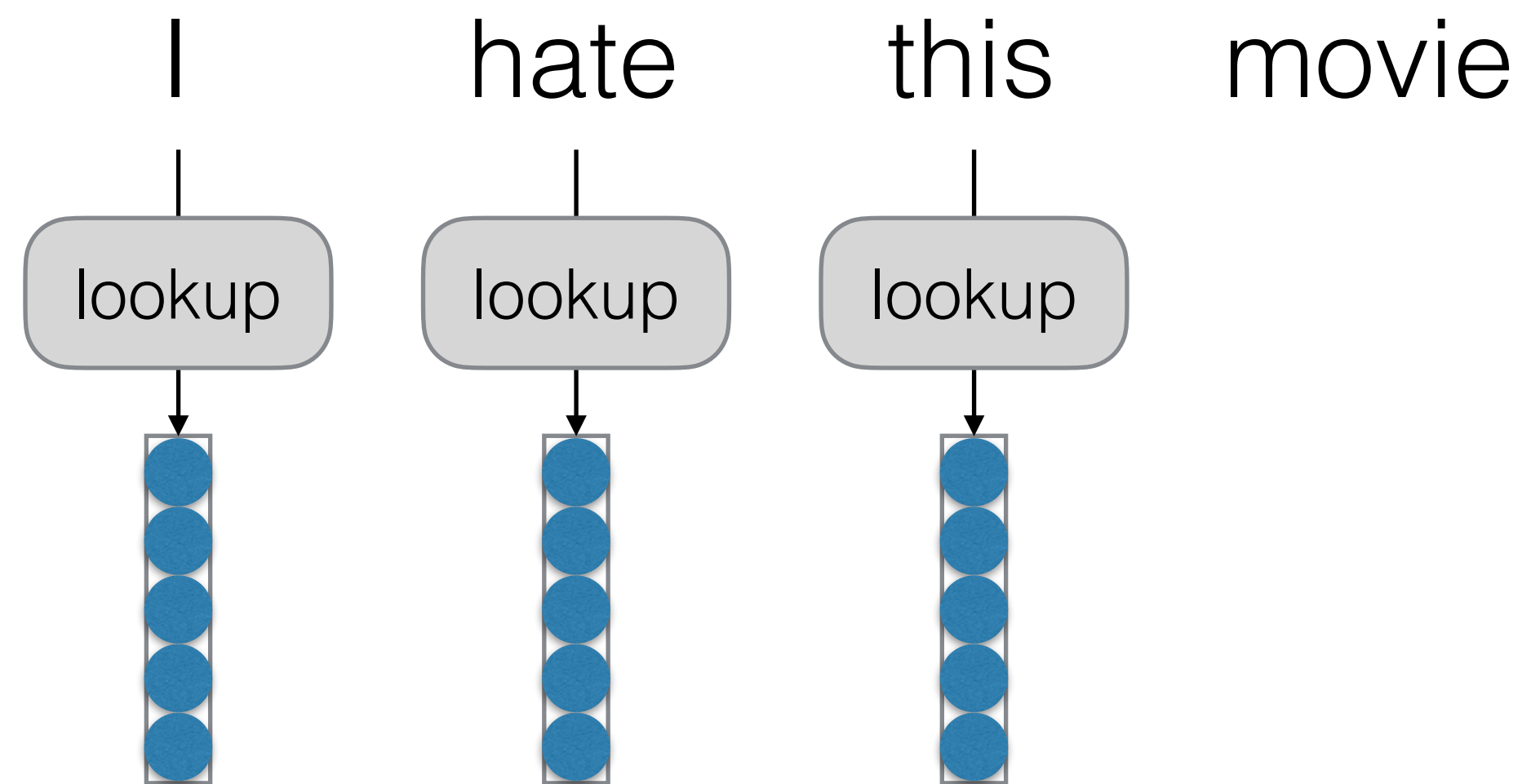
A First Try: Bag of Words (BOW)

Each word has a vector of weights for each tag



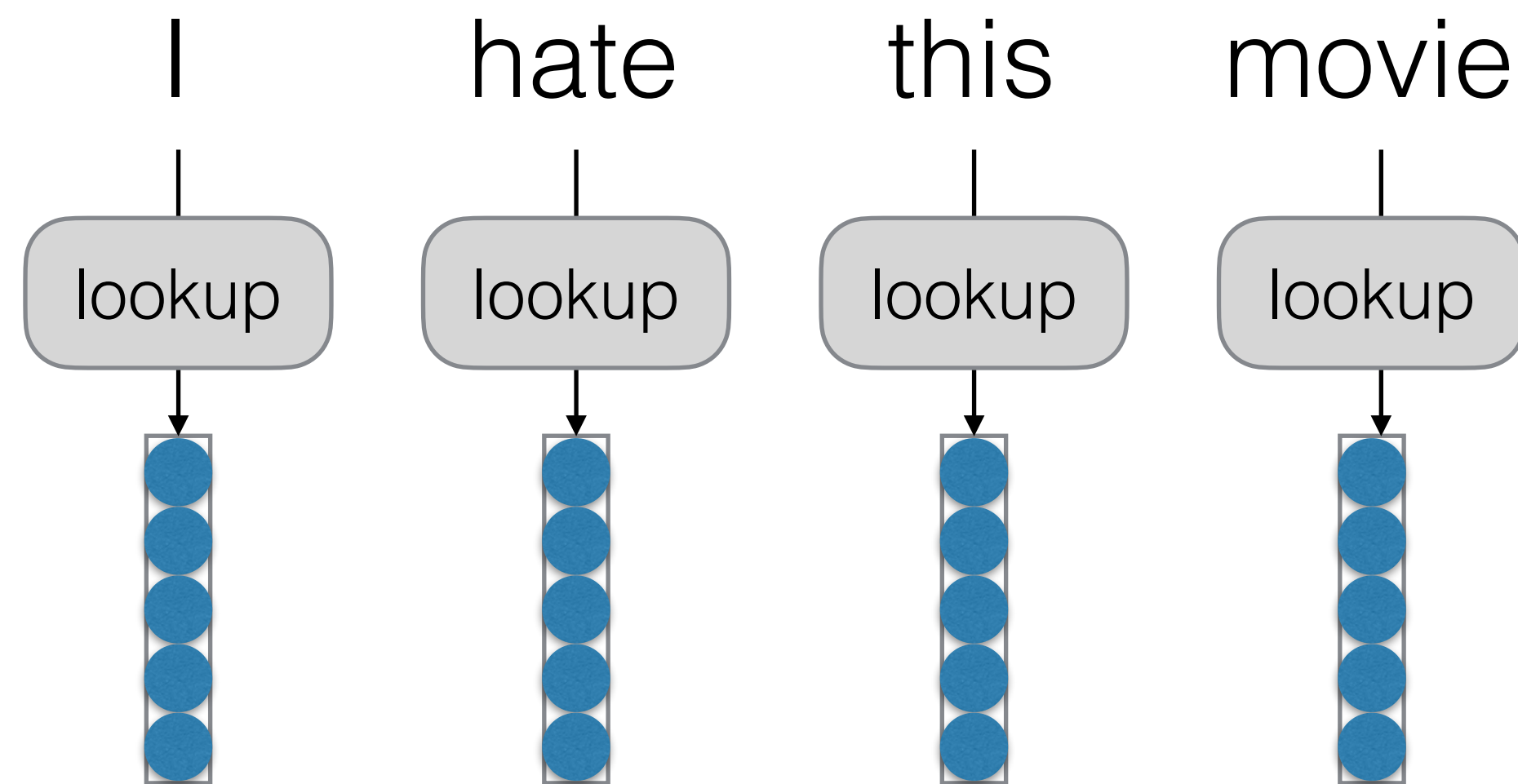
A First Try: Bag of Words (BOW)

Each word has a vector of weights for each tag



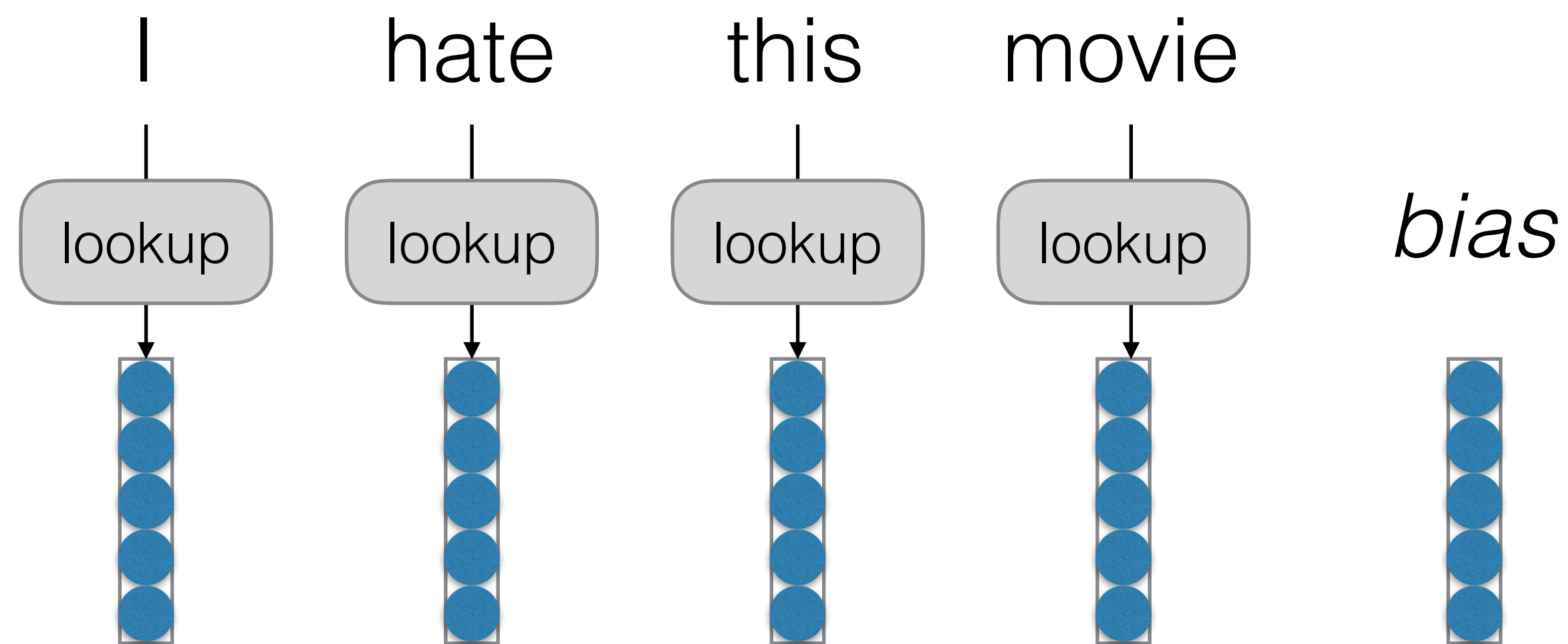
A First Try: Bag of Words (BOW)

Each word has a vector of weights for each tag



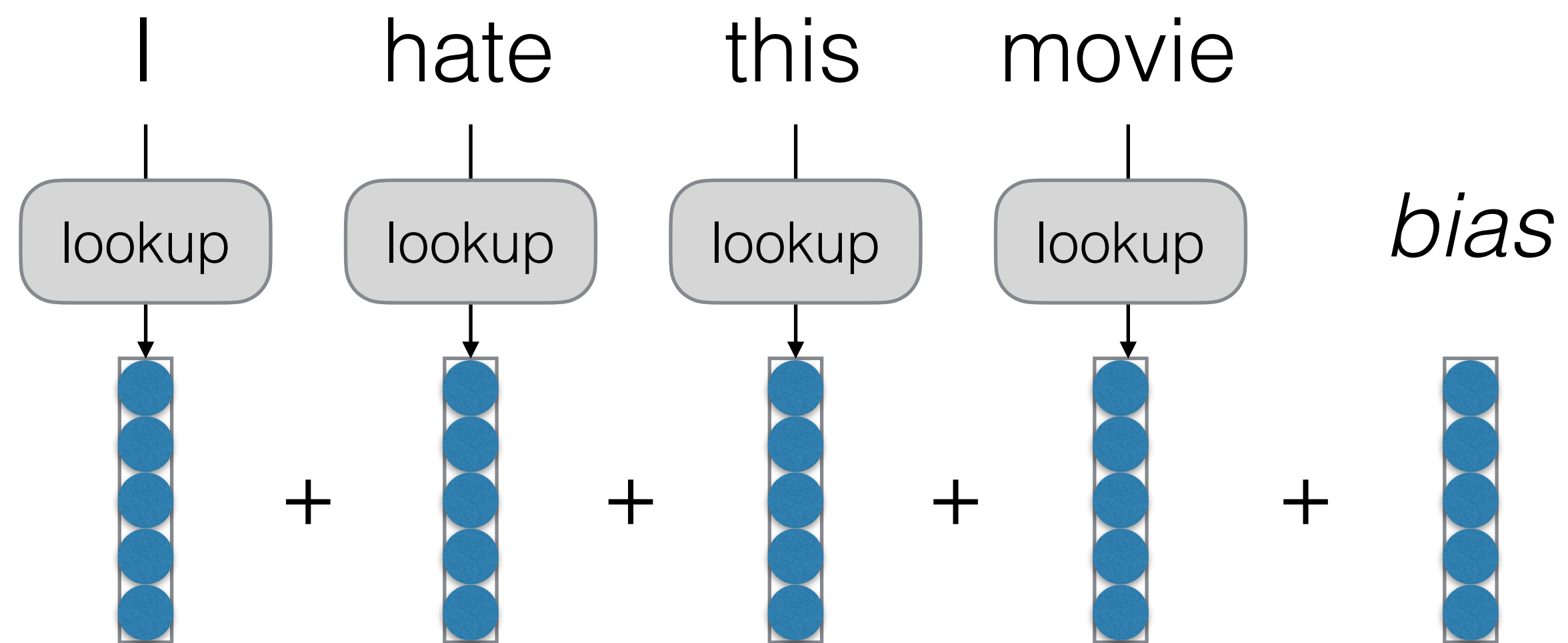
A First Try: Bag of Words (BOW)

Each word has a vector of weights for each tag



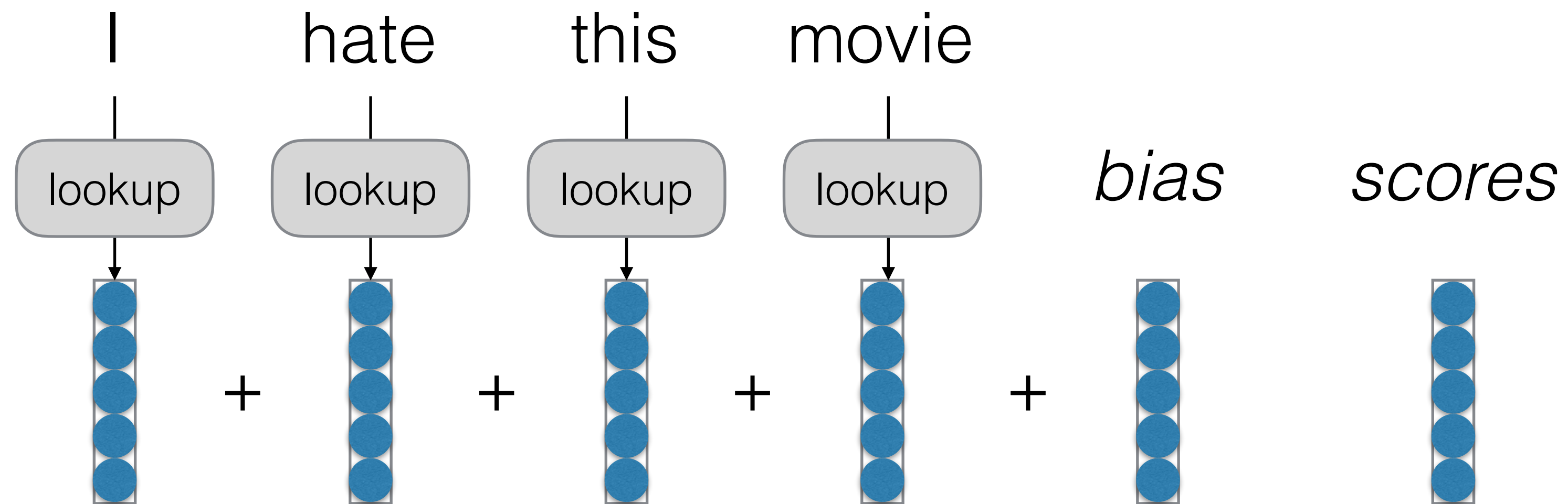
A First Try: Bag of Words (BOW)

Each word has a vector of weights for each tag



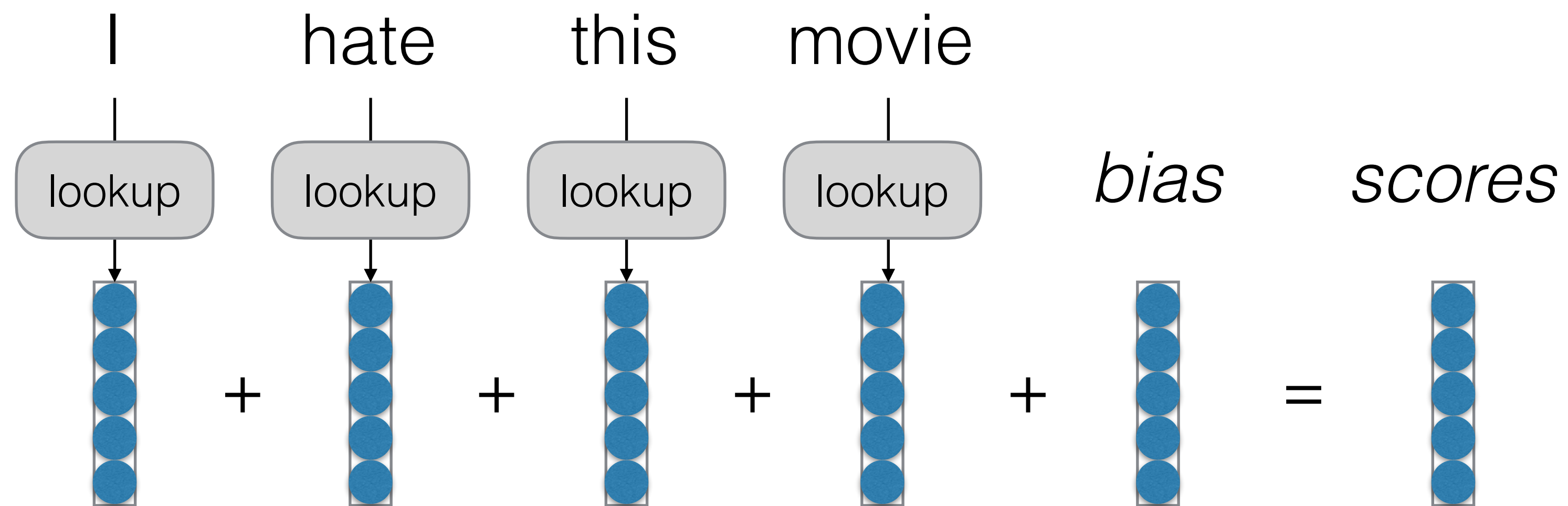
A First Try: Bag of Words (BOW)

Each word has a vector of weights for each tag



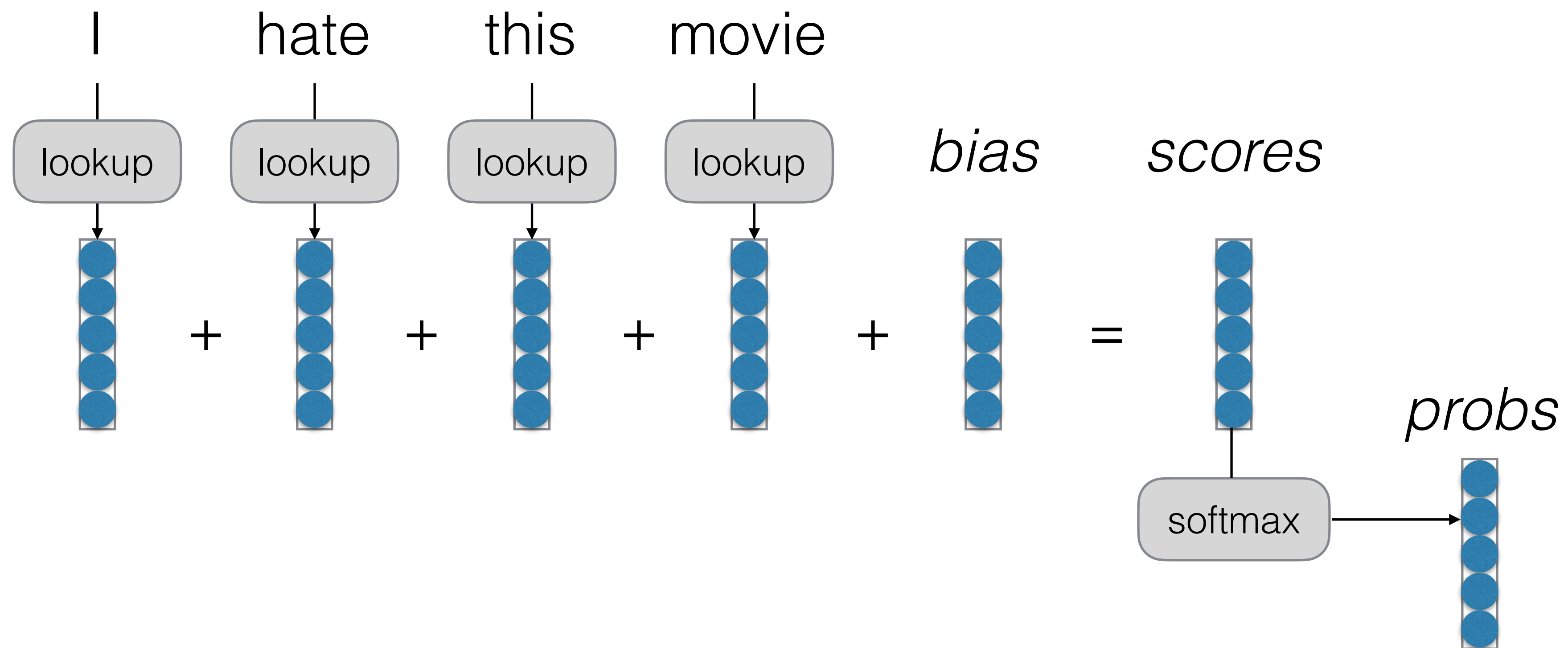
A First Try: Bag of Words (BOW)

Each word has a vector of weights for each tag



A First Try: Bag of Words (BOW)

Each word has a vector of weights for each tag



What do Our Vectors Represent?

What do Our Vectors Represent?

Each word has its own 5 elements corresponding to [very good, good, neutral, bad, very bad]

What do Our Vectors Represent?

Each word has its own 5 elements corresponding to [very good, good, neutral, bad, very bad]

“hate” will have a high value for “very bad”, etc.

Build It, Break It

Build It, Break It

I don't love this movie

Build It, Break It

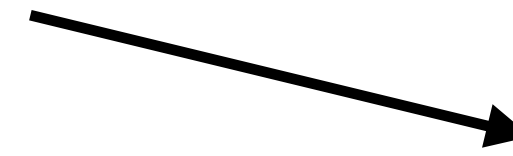
I don't love this movie

very good
good
neutral
bad
very bad

Build It, Break It

I don't love this movie

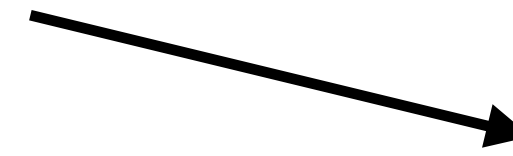
very good
good
neutral
bad
very bad



Build It, Break It

I don't love this movie

very good
good
neutral
bad
very bad

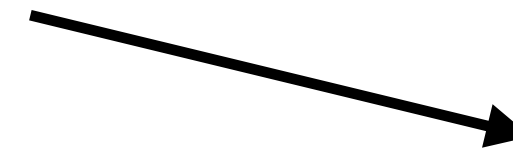


There's nothing I don't
love about this movie

Build It, Break It

I don't love this movie

very good
good
neutral
bad
very bad



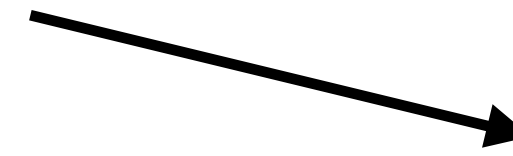
There's nothing I don't love about this movie

very good
good
neutral
bad
very bad

Build It, Break It

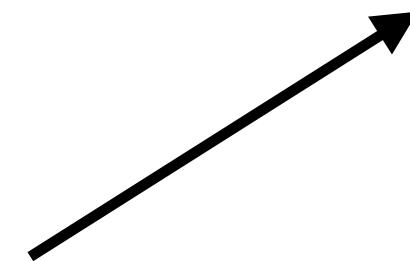
I don't love this movie

very good
good
neutral
bad
very bad



There's nothing I don't
love about this movie

very good
good
neutral
bad
very bad



Combination Features

Combination Features

Does it contain "don't" and "love"?

Combination Features

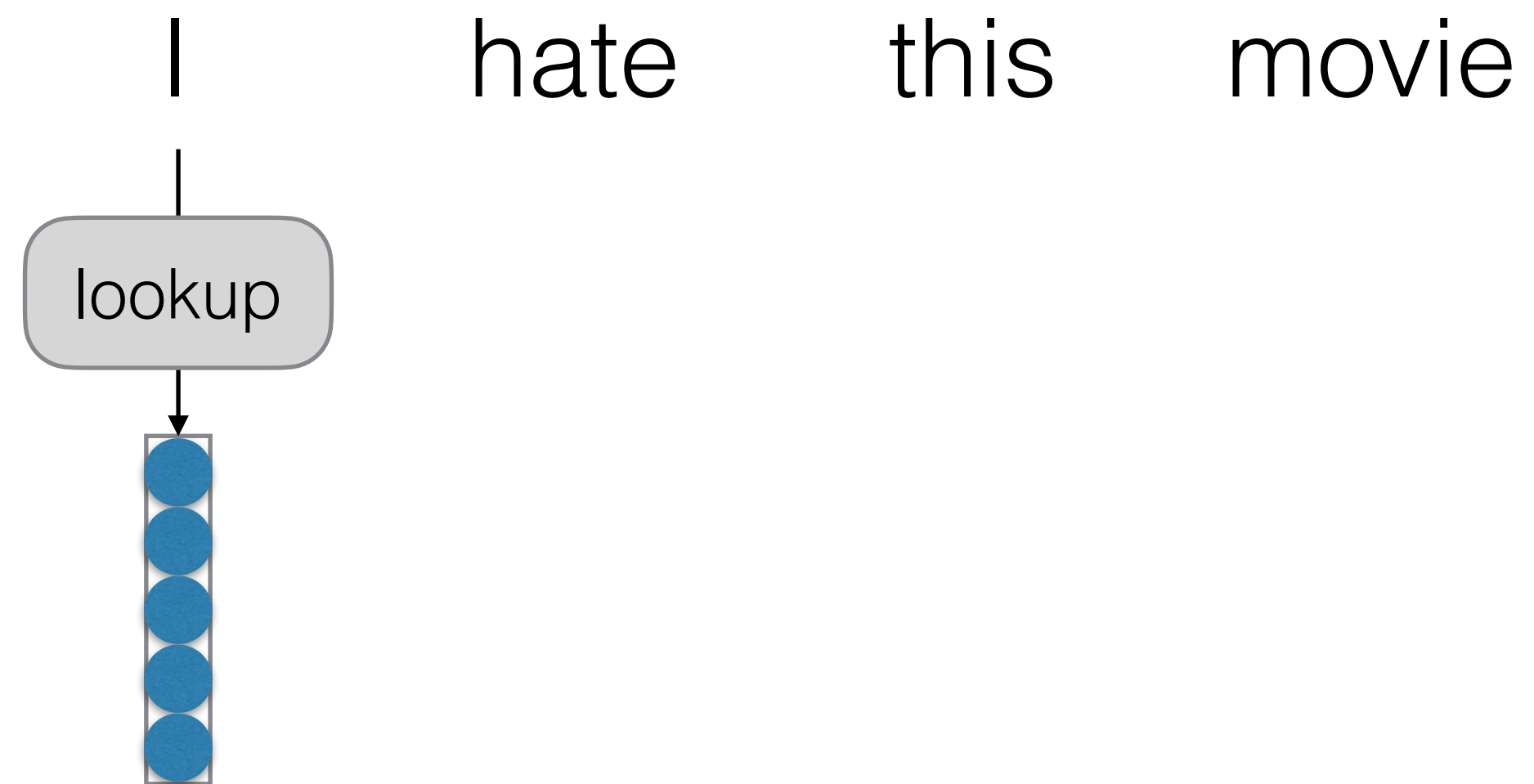
Does it contain "don't" and "love"?

Does it contain "don't", "i", "love", and "nothing"?

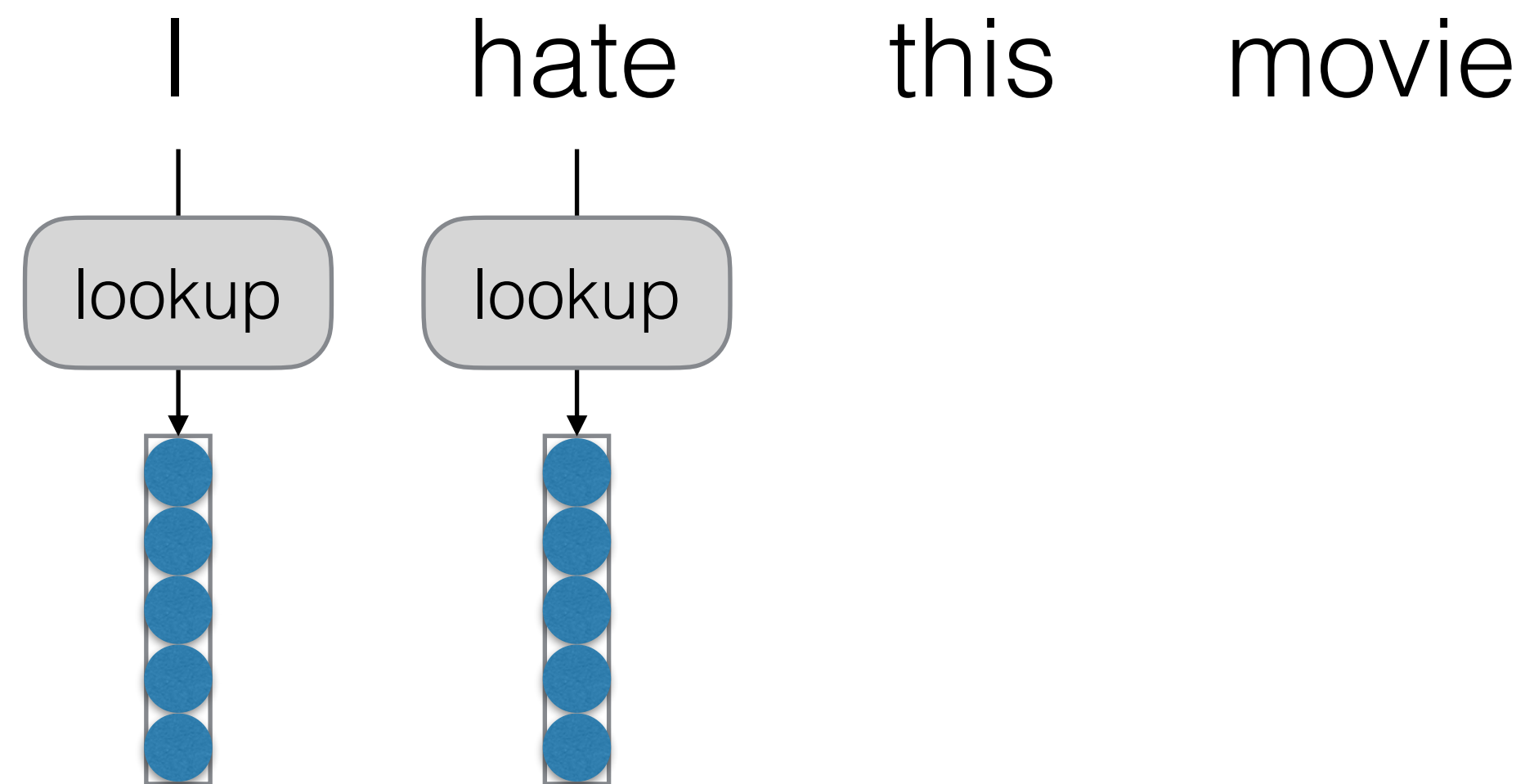
Basic Idea of Neural Networks (for NLP Prediction Tasks)

I hate this movie

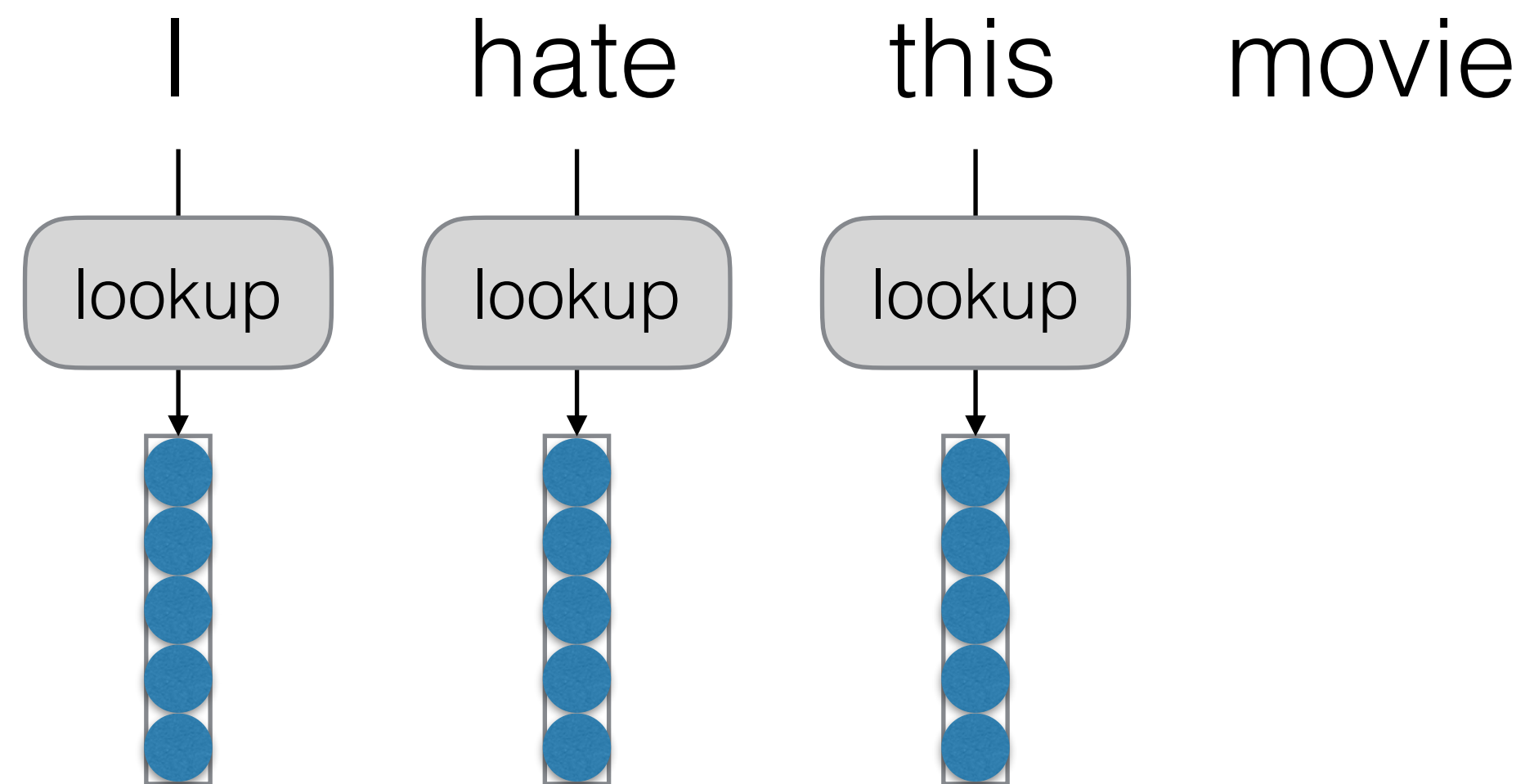
Basic Idea of Neural Networks (for NLP Prediction Tasks)



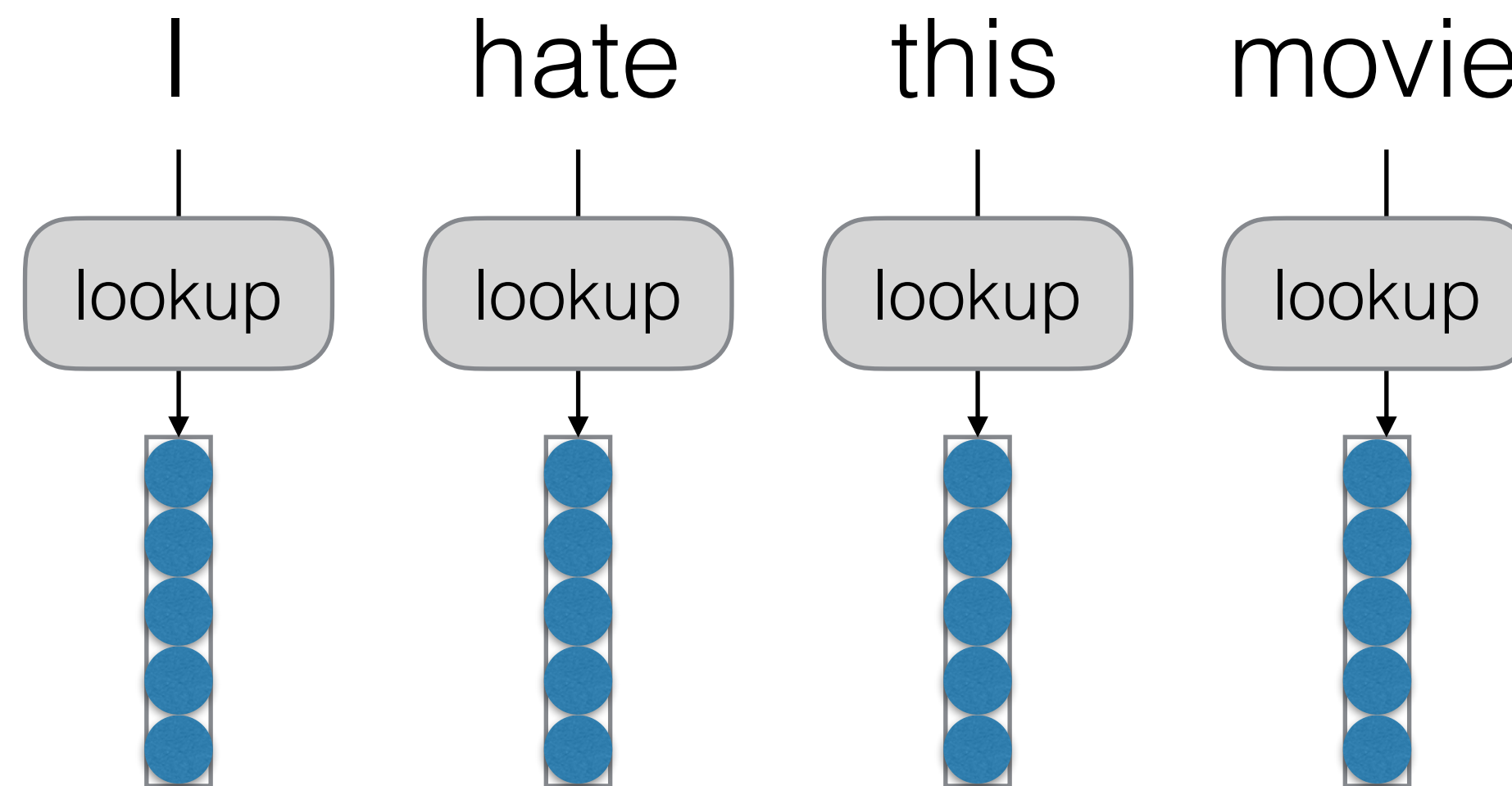
Basic Idea of Neural Networks (for NLP Prediction Tasks)



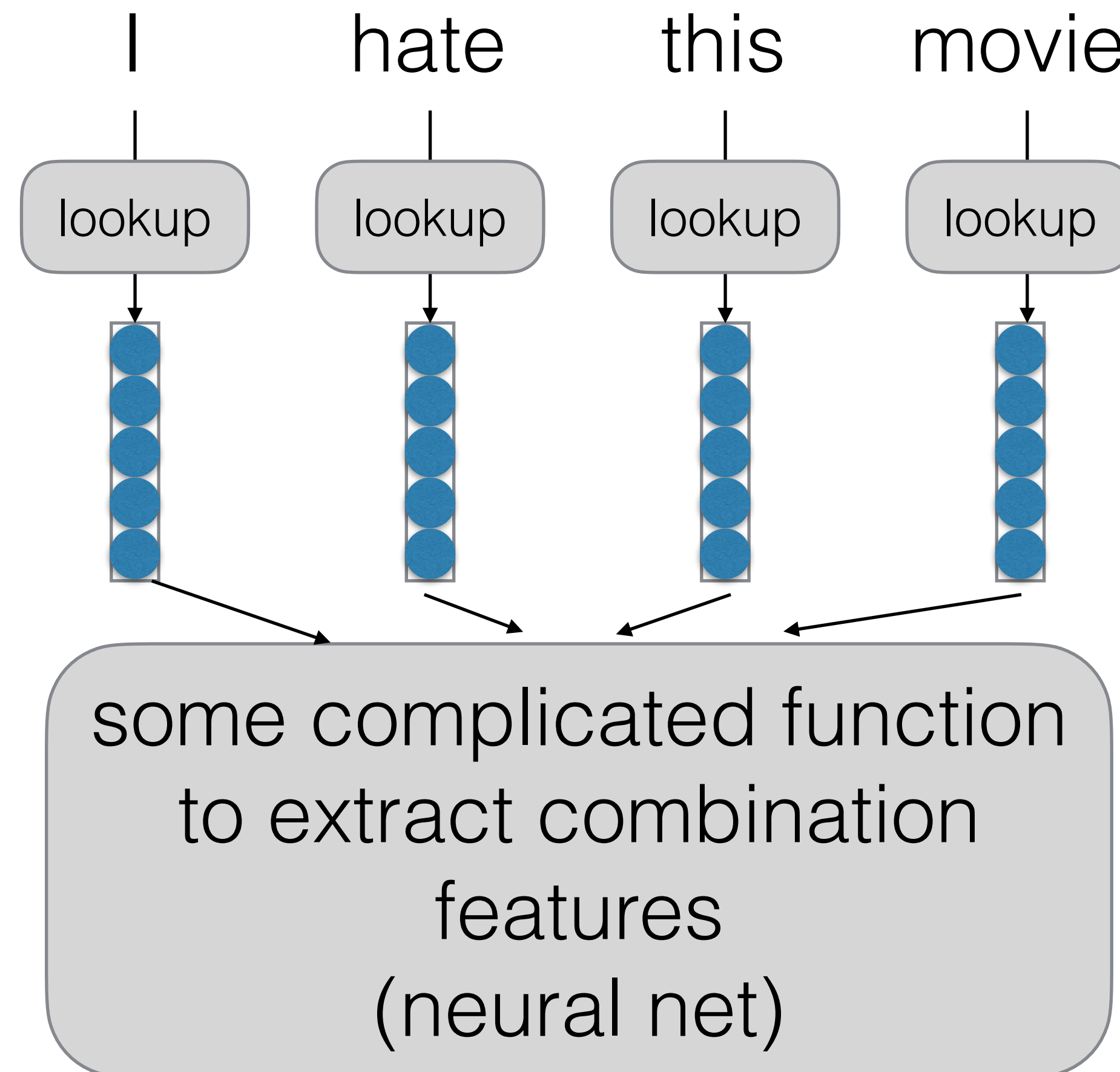
Basic Idea of Neural Networks (for NLP Prediction Tasks)



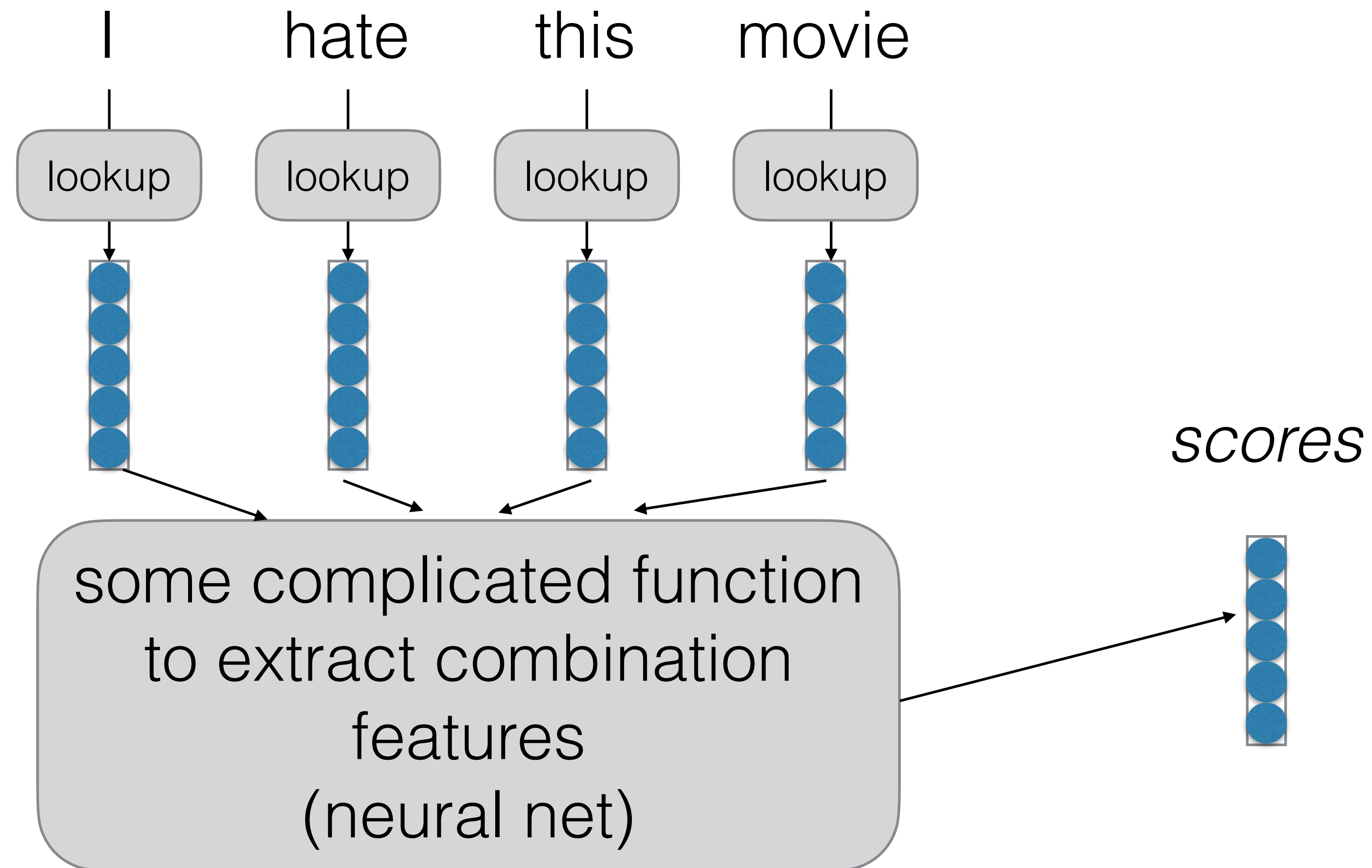
Basic Idea of Neural Networks (for NLP Prediction Tasks)



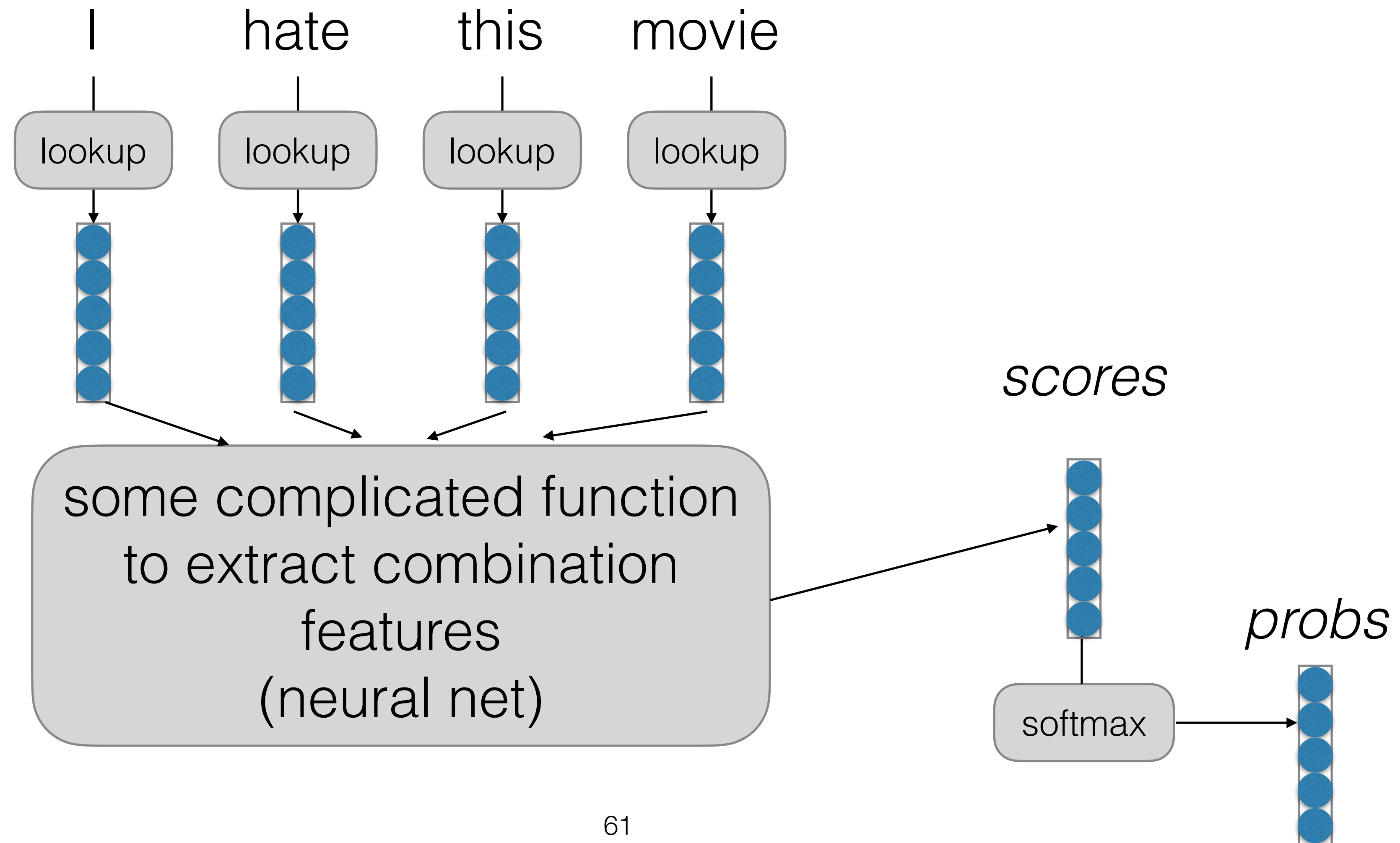
Basic Idea of Neural Networks (for NLP Prediction Tasks)



Basic Idea of Neural Networks (for NLP Prediction Tasks)



Basic Idea of Neural Networks (for NLP Prediction Tasks)



Continuous Bag of Words (CBOW)

Each word has a feature vector, each feature has weights

I hate this movie

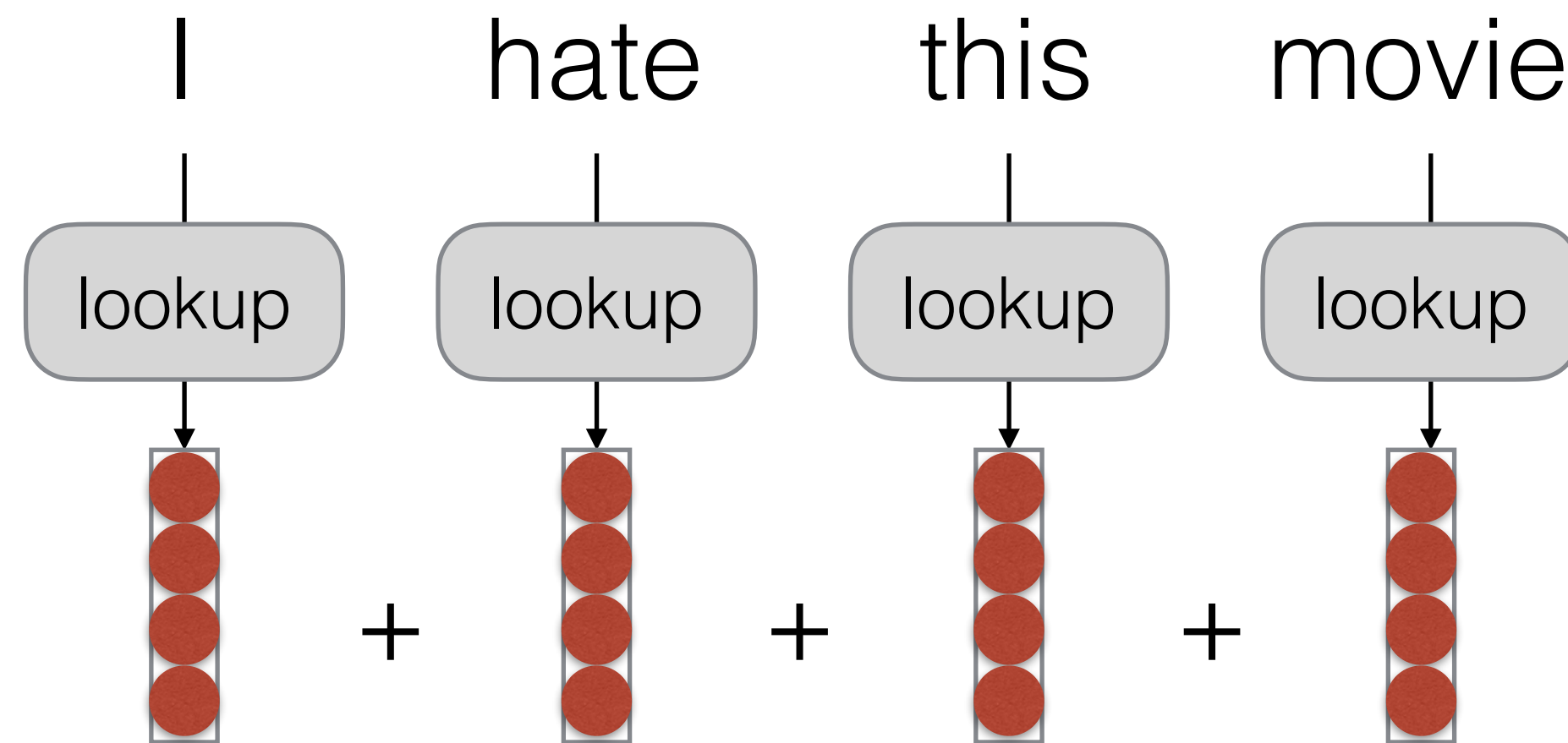
Continuous Bag of Words (CBOW)

Each word has a feature vector, each feature has weights



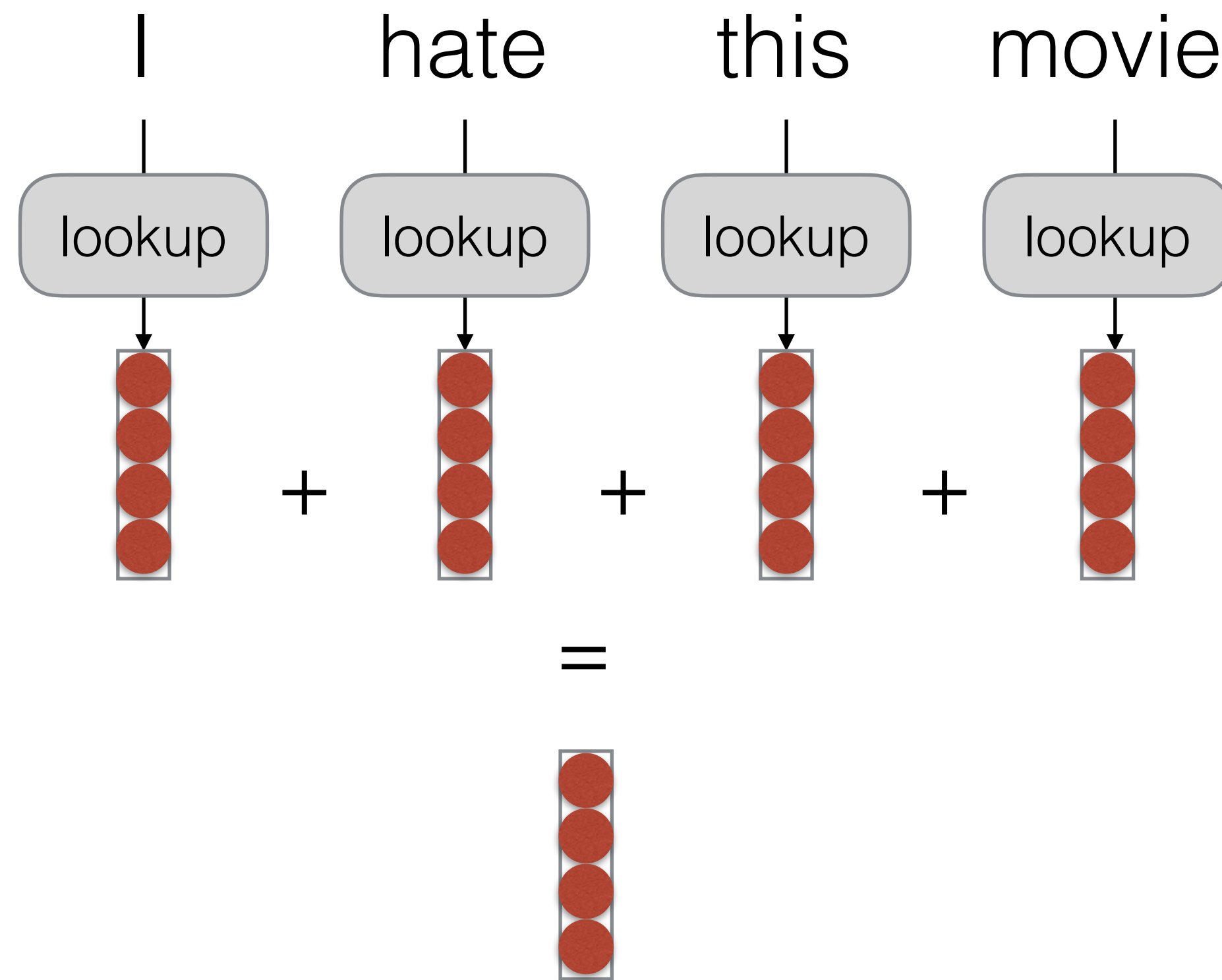
Continuous Bag of Words (CBOW)

Each word has a feature vector, each feature has weights



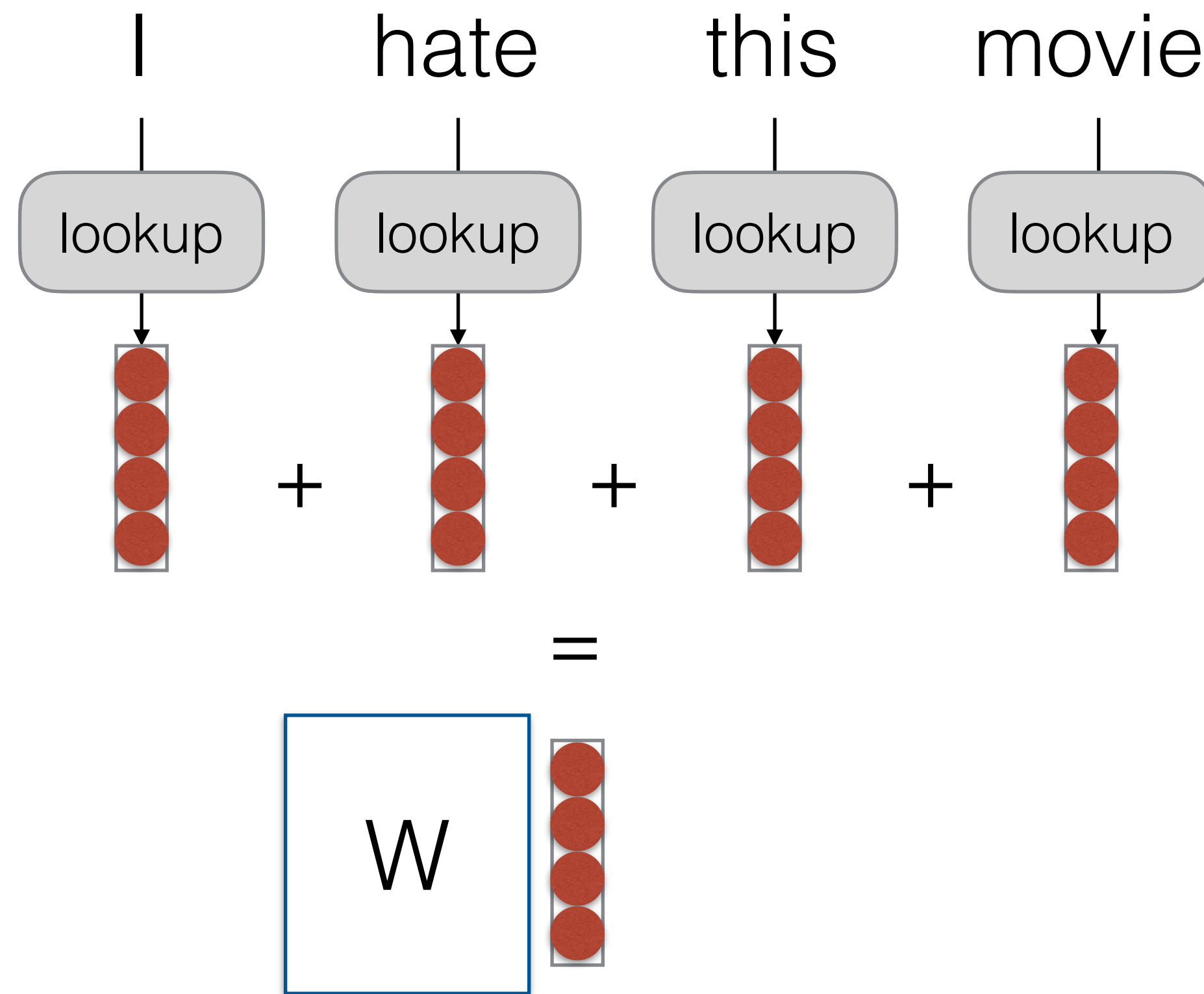
Continuous Bag of Words (CBOW)

Each word has a feature vector, each feature has weights



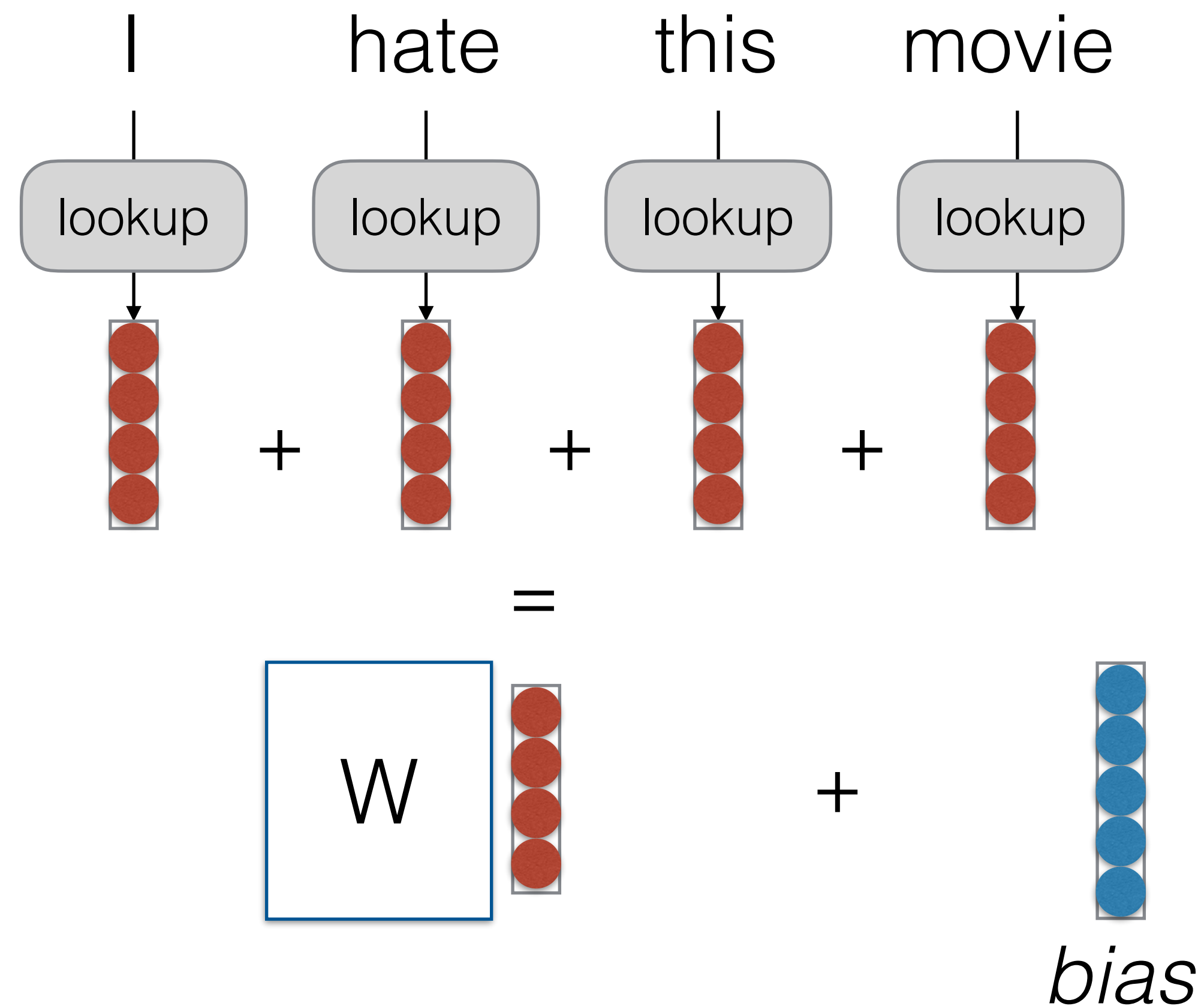
Continuous Bag of Words (CBOW)

Each word has a feature vector, each feature has weights



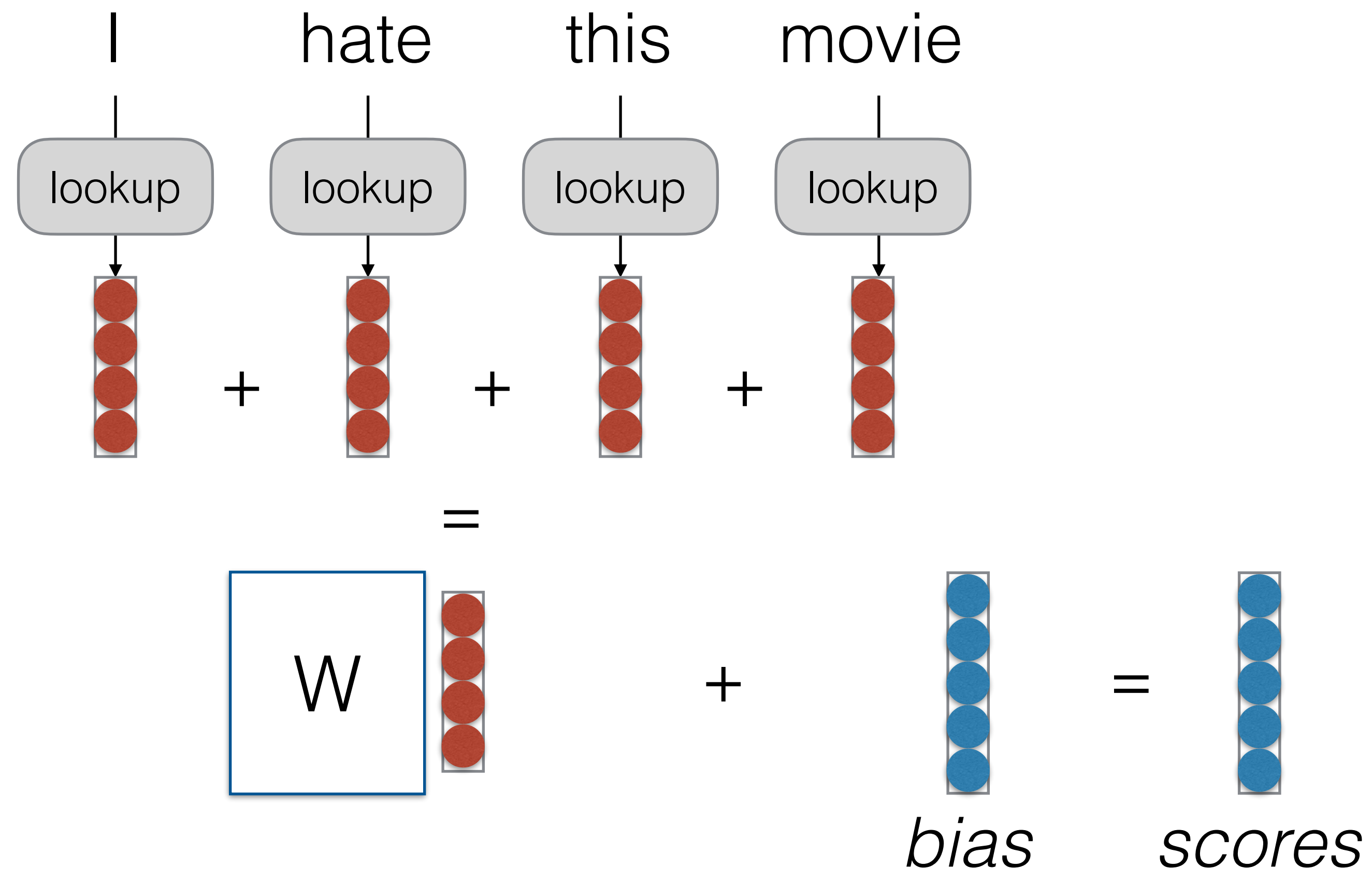
Continuous Bag of Words (CBOW)

Each word has a feature vector, each feature has weights



Continuous Bag of Words (CBOW)

Each word has a feature vector, each feature has weights



What do Our Vectors Represent?

What do Our Vectors Represent?

Each vector has “features” (e.g. is this an animate object? is this a positive word, etc.)

What do Our Vectors Represent?

Each vector has "features" (e.g. is this an animate object? is this a positive word, etc.)

We sum these features, then use these to make predictions

What do Our Vectors Represent?

Each vector has "features" (e.g. is this an animate object? is this a positive word, etc.)

We sum these features, then use these to make predictions

Still no combination features: only the expressive power of a linear model, but dimension reduced

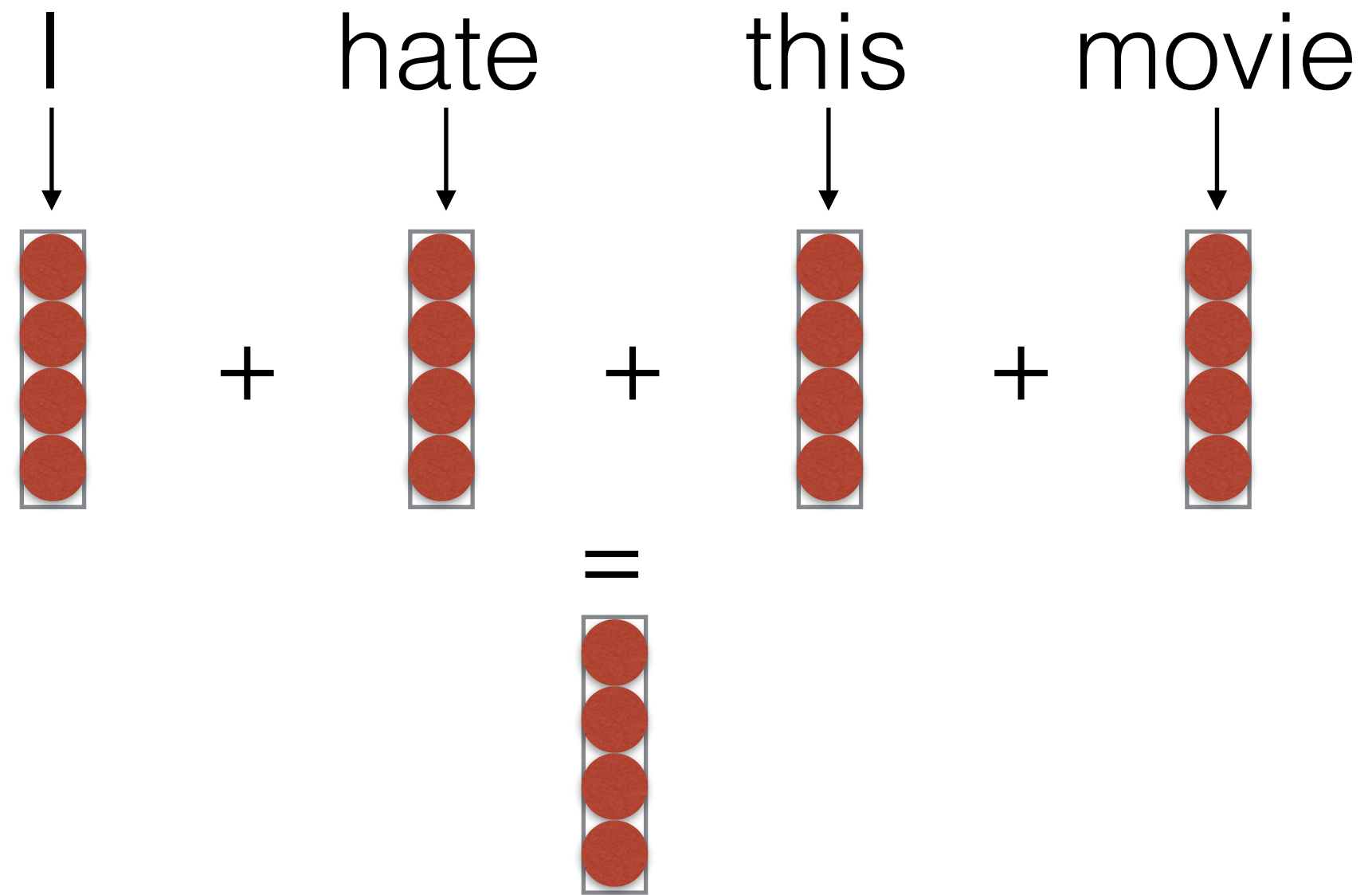
Deep CBOW

Add several feature transforms

I hate this movie

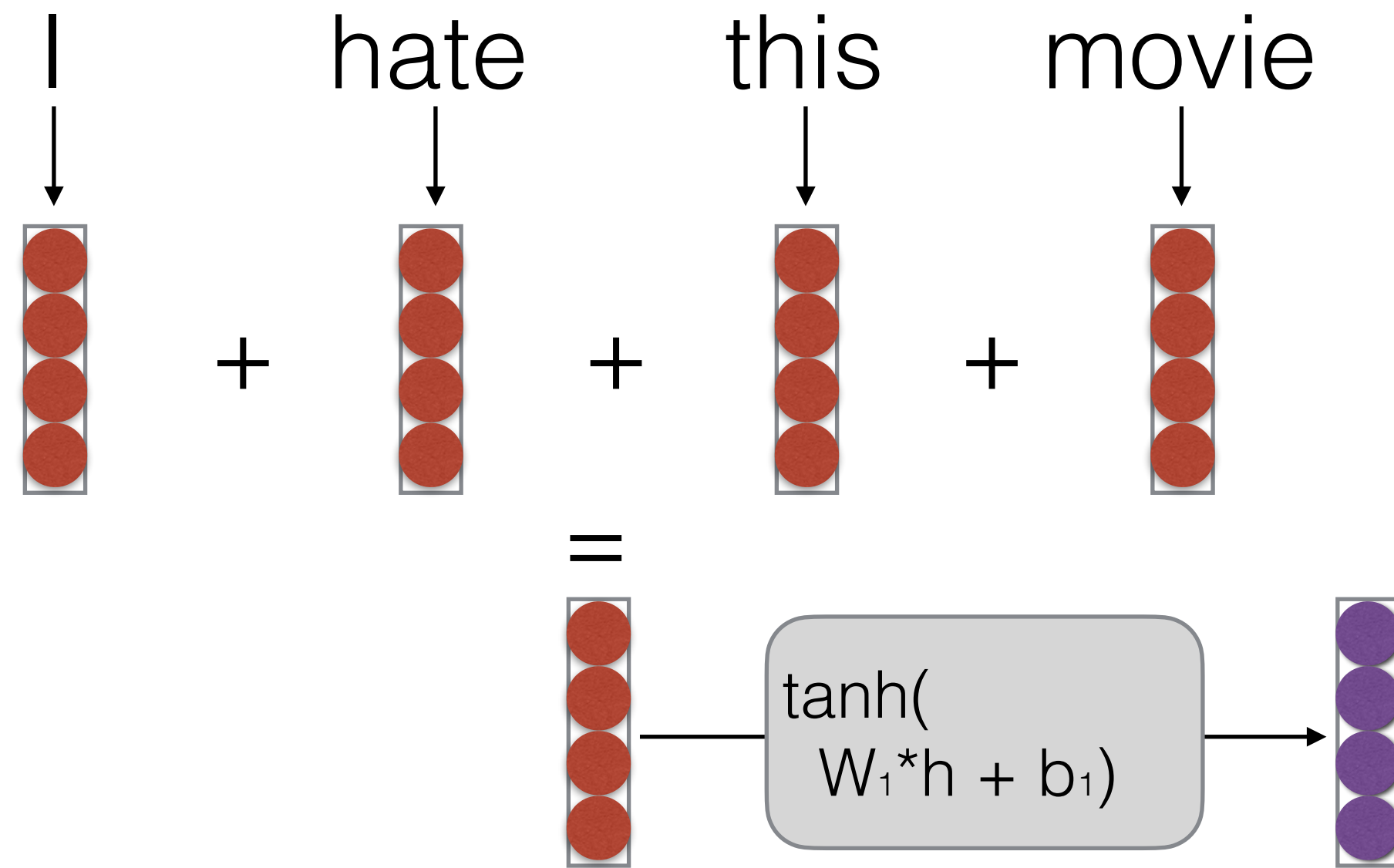
Deep CBOW

Add several feature transforms



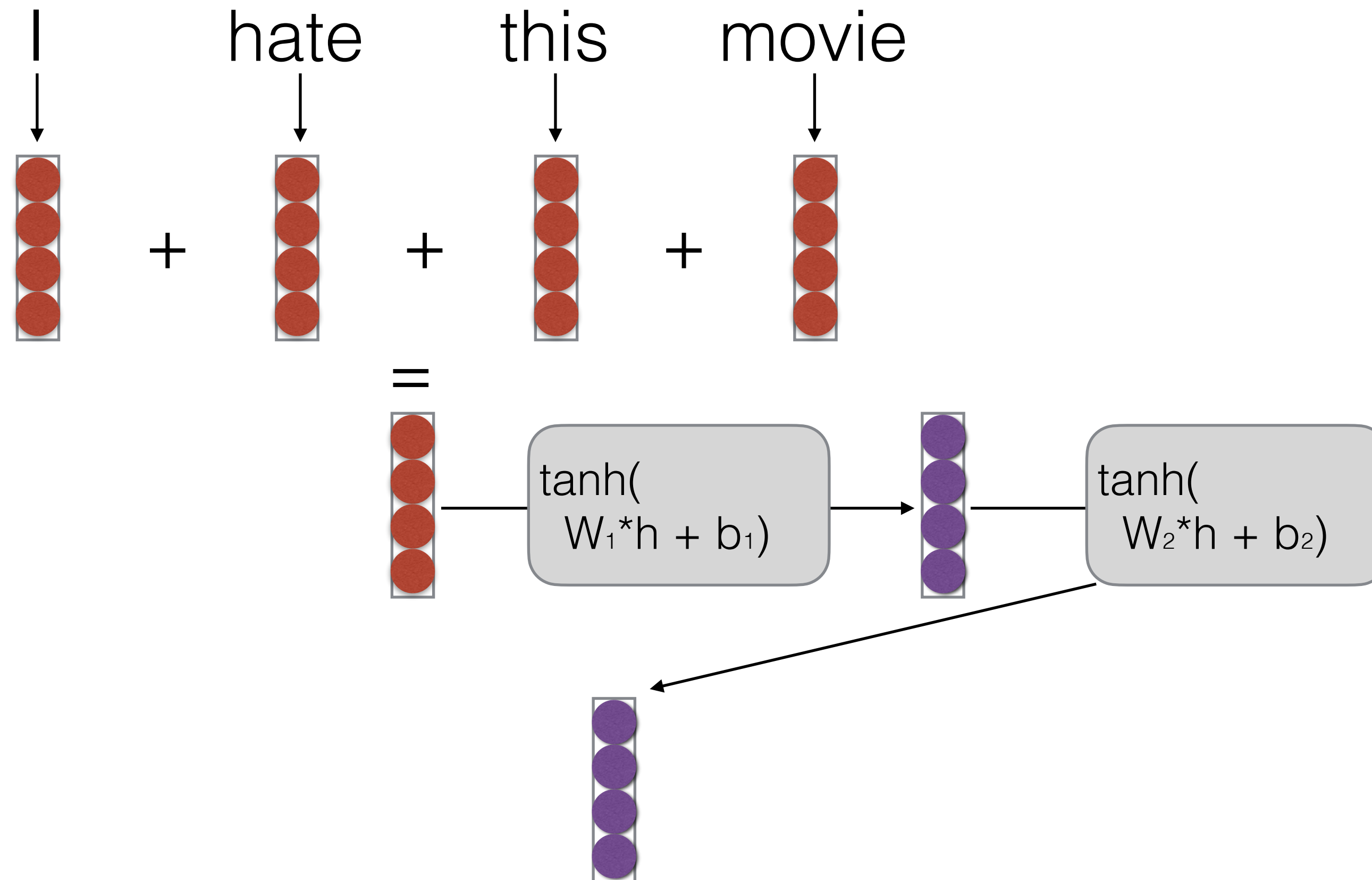
Deep CBOW

Add several feature transforms



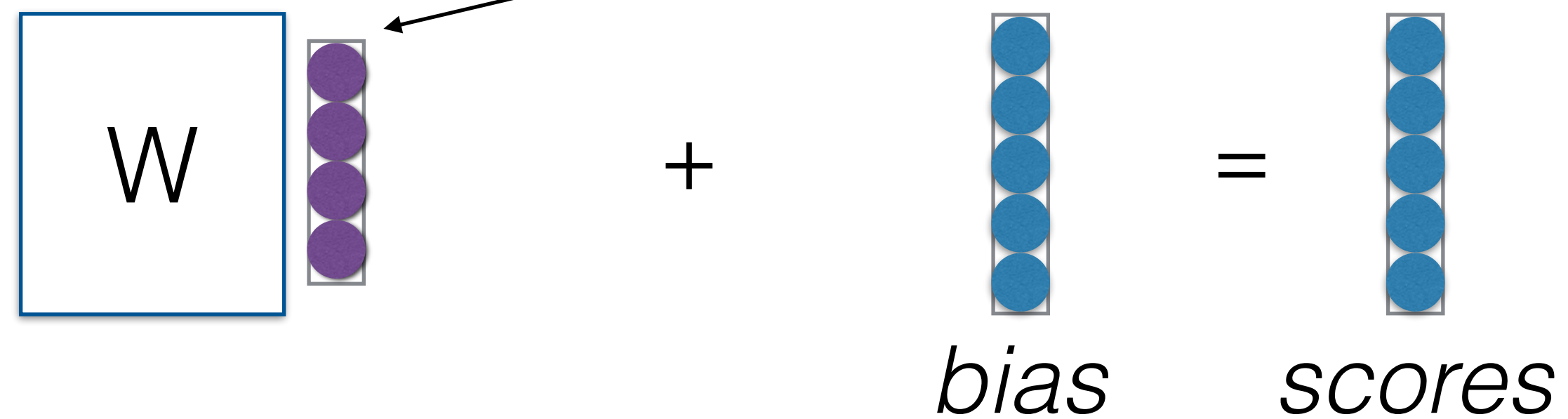
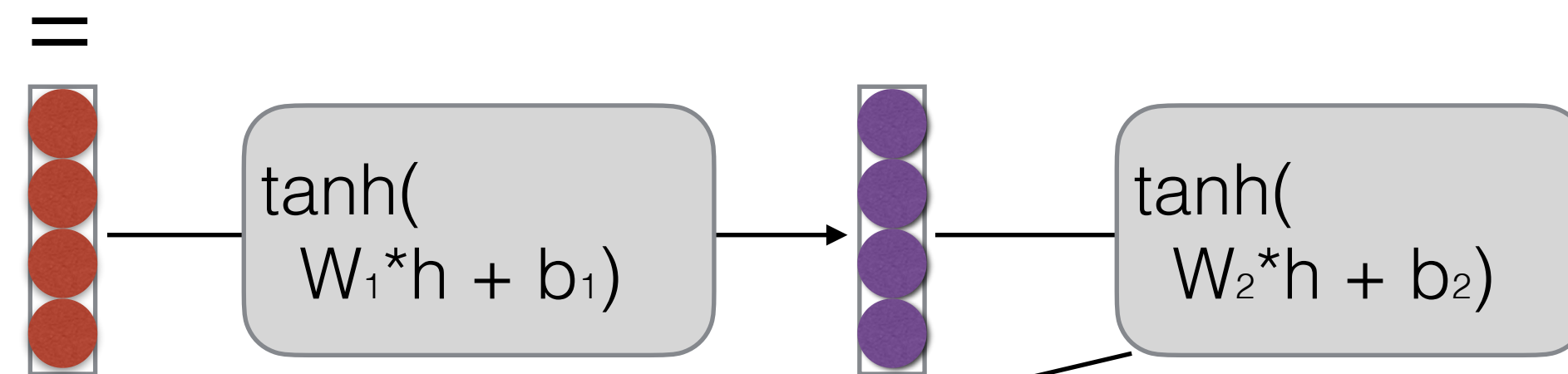
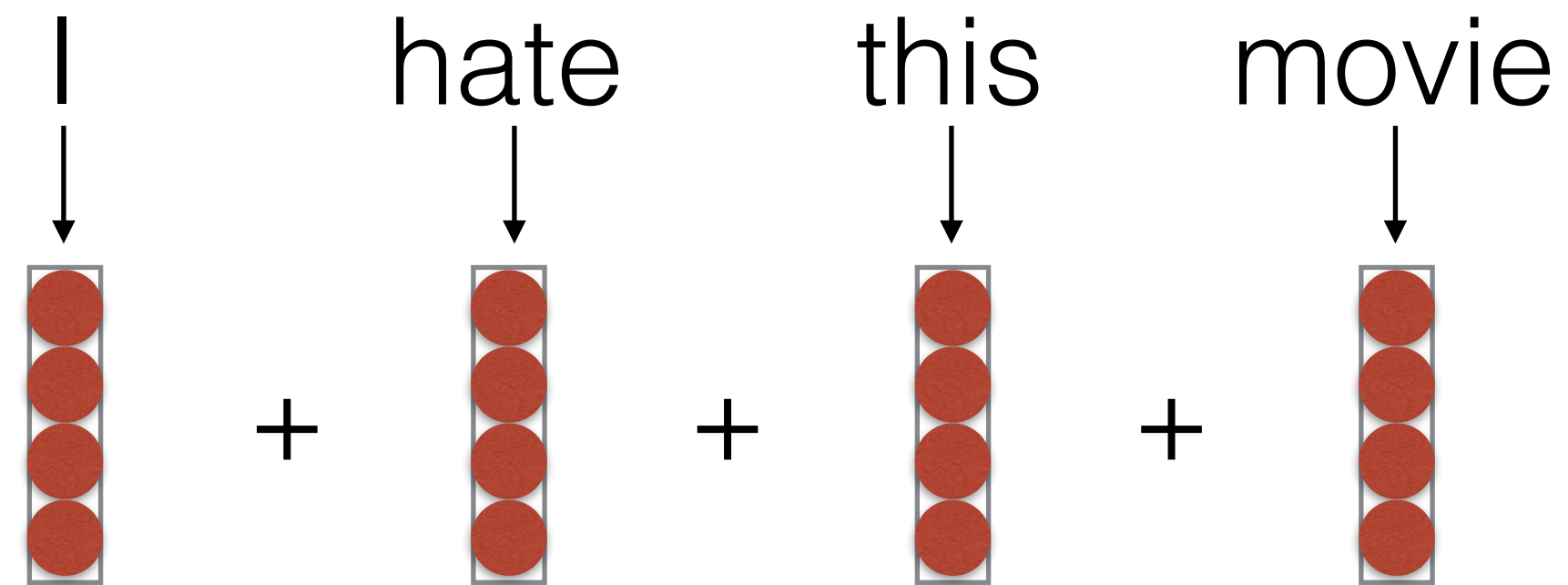
Deep CBOW

Add several feature transforms



Deep CBOW

Add several feature transforms



What do Our Vectors Represent?

What do Our Vectors Represent?

Now things are more interesting!

What do Our Vectors Represent?

Now things are more interesting!

We can learn feature combinations (a node in the second layer might be “feature 1 AND feature 5 are active”)

What do Our Vectors Represent?

Now things are more interesting!

We can learn feature combinations (a node in the second layer might be "feature 1 AND feature 5 are active")

e.g. capture things such as "not" AND "hate"

What is a Neural Net?

Computation Graphs

“Neural” Nets

Original Motivation: The Brain

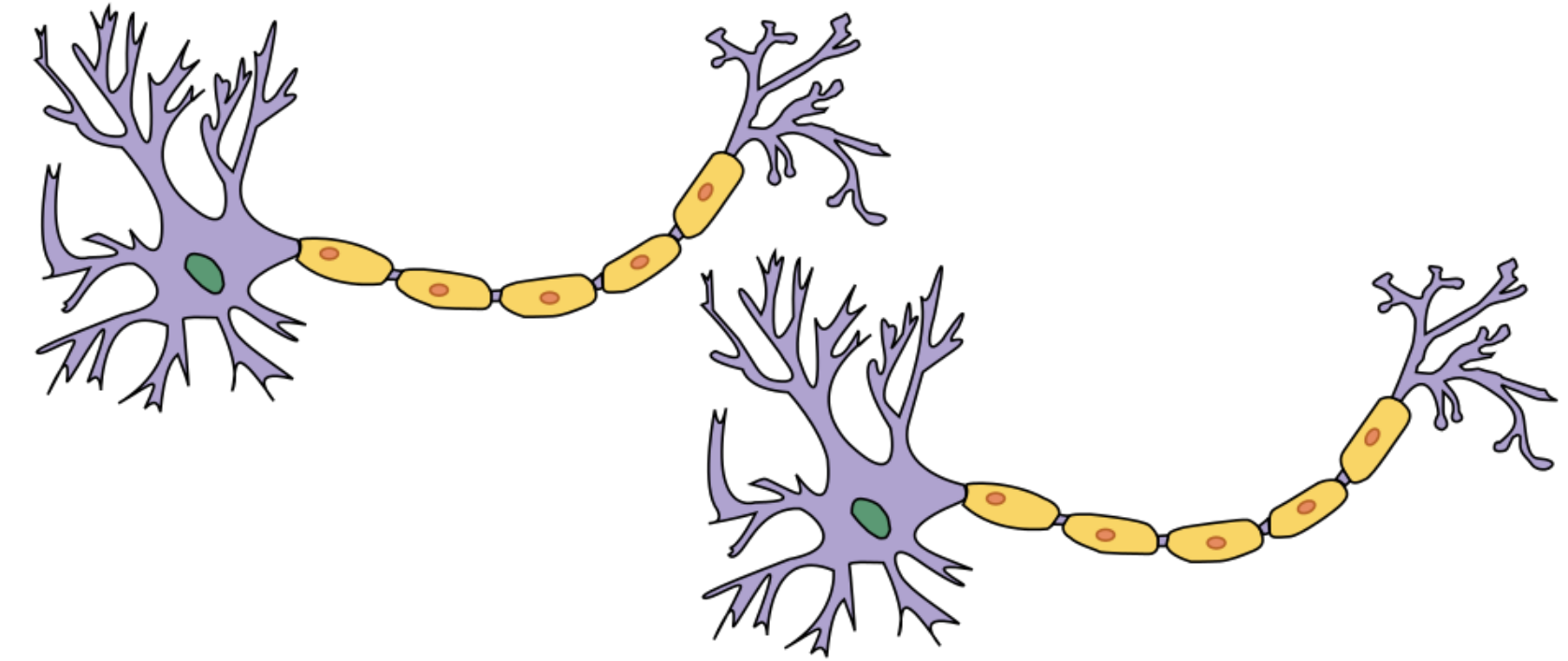
Current Implementation

“Neural” Nets

Original Motivation: The Brain

Original Motivation: Neurons in the Brain

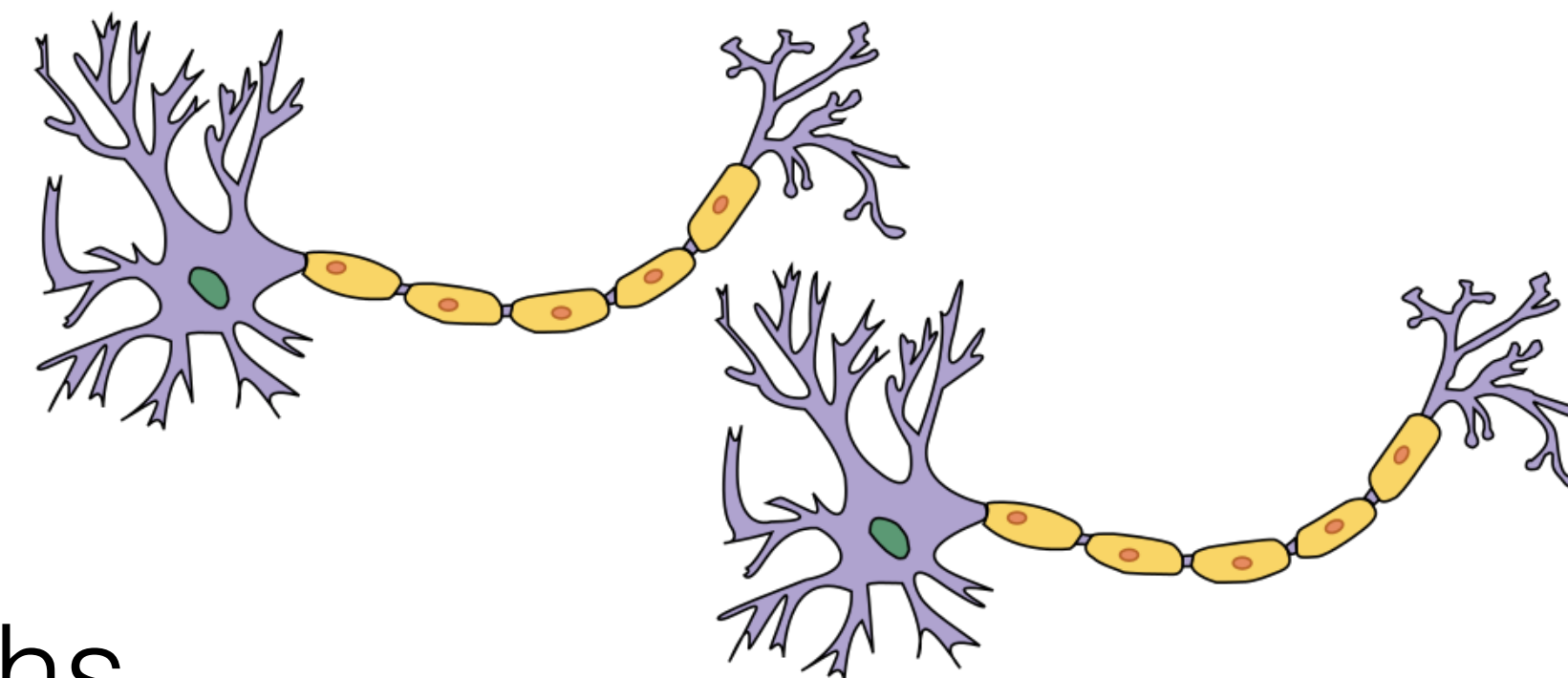
Current Implementation



"Neural" Nets

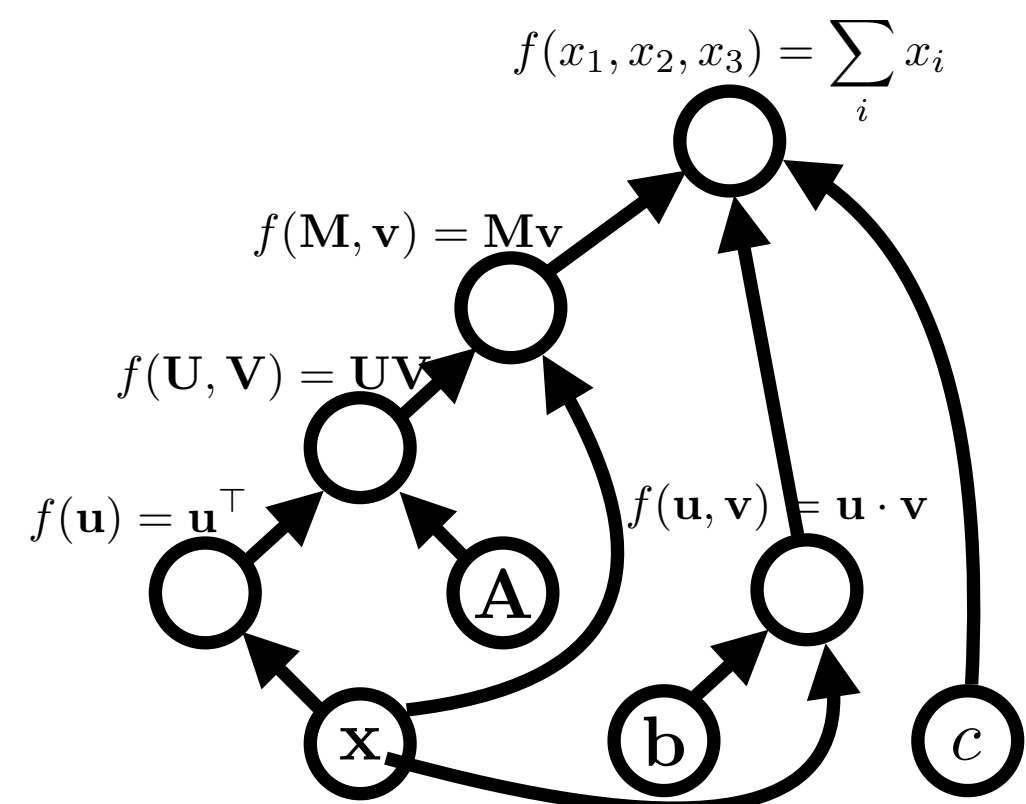
Original Motivation: The Brain

Original Motivation: Neurons in the Brain



Current Implementation

Current Conception: Computation Graphs



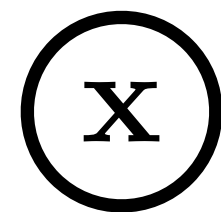
expression:

x

expression:

x

graph:



expression:

x

graph:

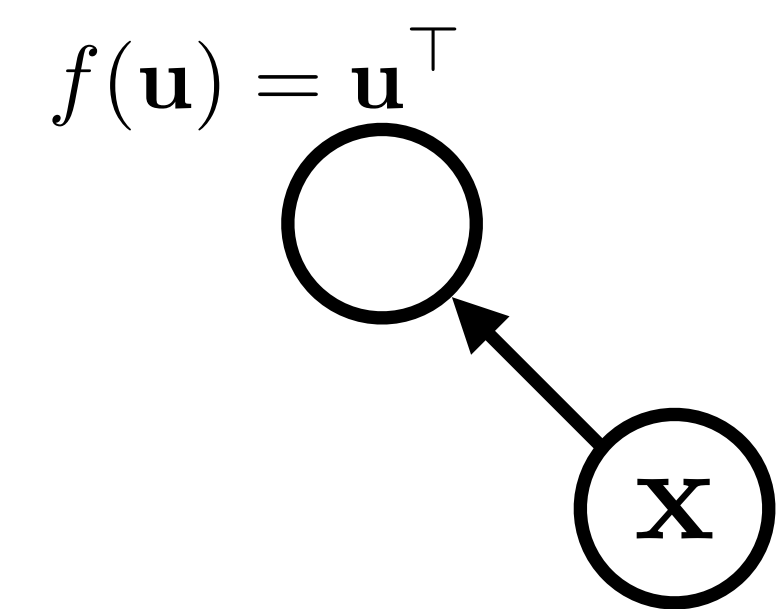
A **node** is a {tensor, matrix, vector, scalar} value

x

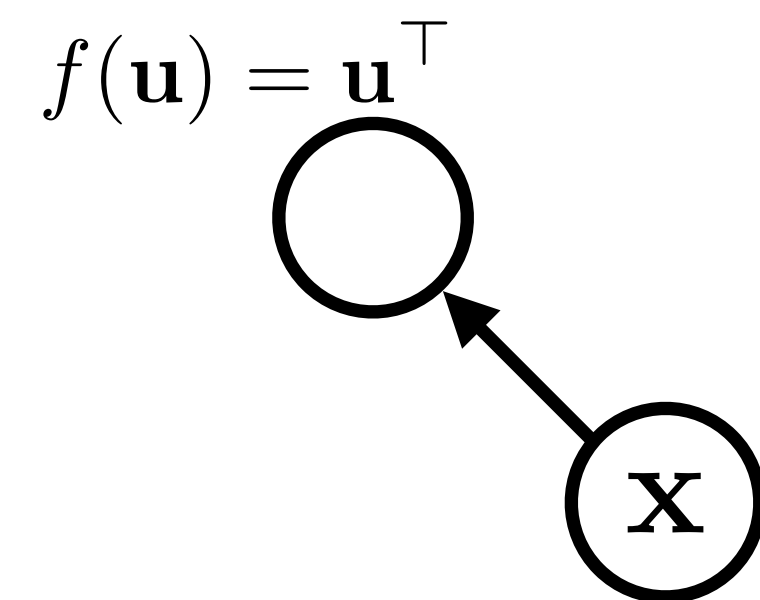
expression:

$$\mathbf{x}^\top$$

graph:

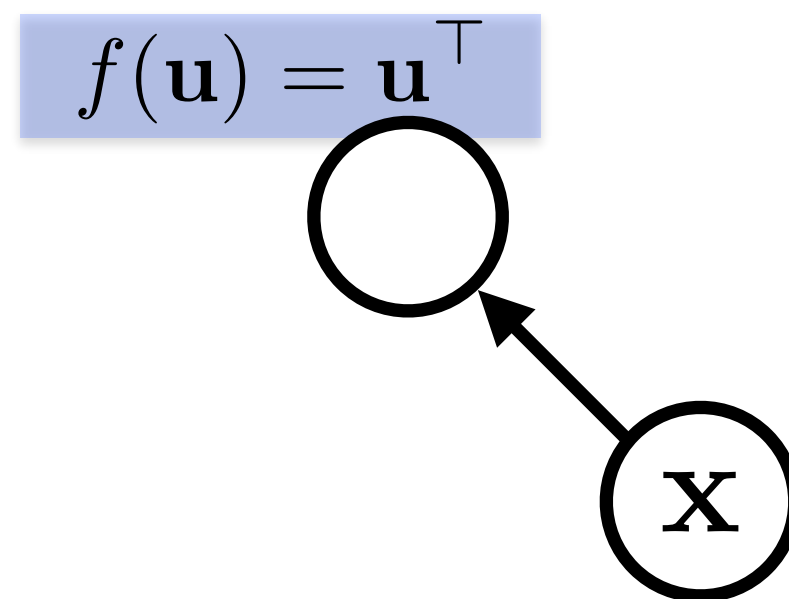


An **edge** represents a function argument (and also an data dependency). They are just pointers to nodes.



An **edge** represents a function argument (and also an data dependency). They are just pointers to nodes.

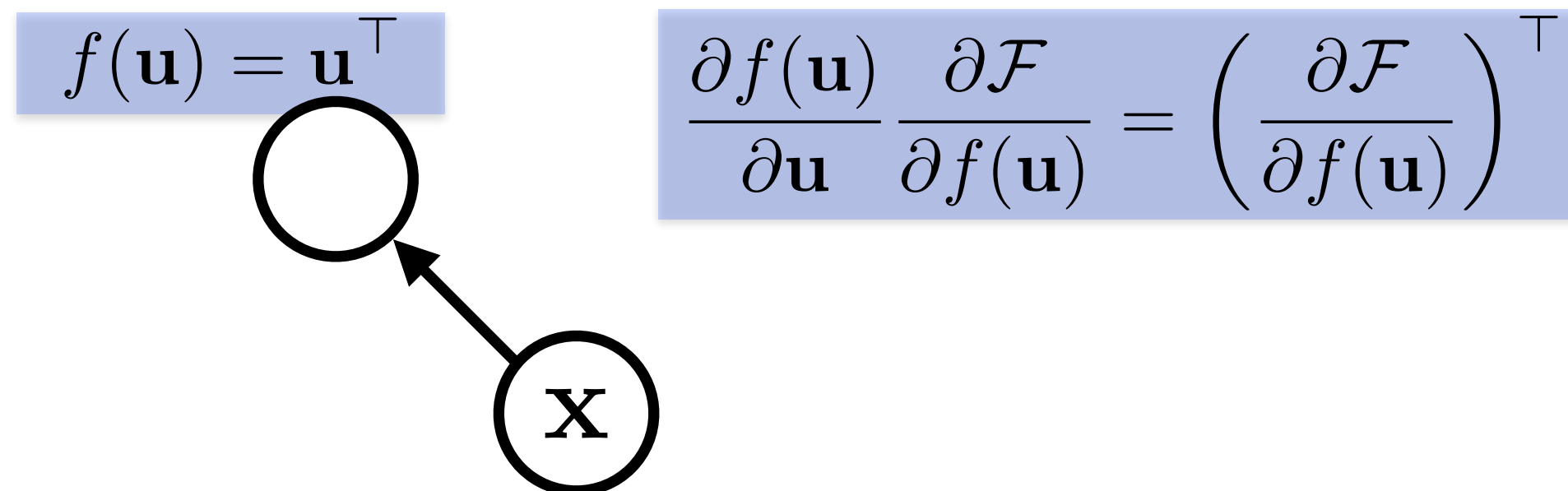
A **node** with an incoming **edge** is a **function** of that edge's tail node.



An **edge** represents a function argument (and also an data dependency). They are just pointers to nodes.

A **node** with an incoming **edge** is a **function** of that edge's tail node.

A **node** knows how to compute its value and the *value of its derivative w.r.t each argument (edge) times a derivative of an arbitrary input* $\frac{\partial \mathcal{F}}{\partial f(\mathbf{u})}$.

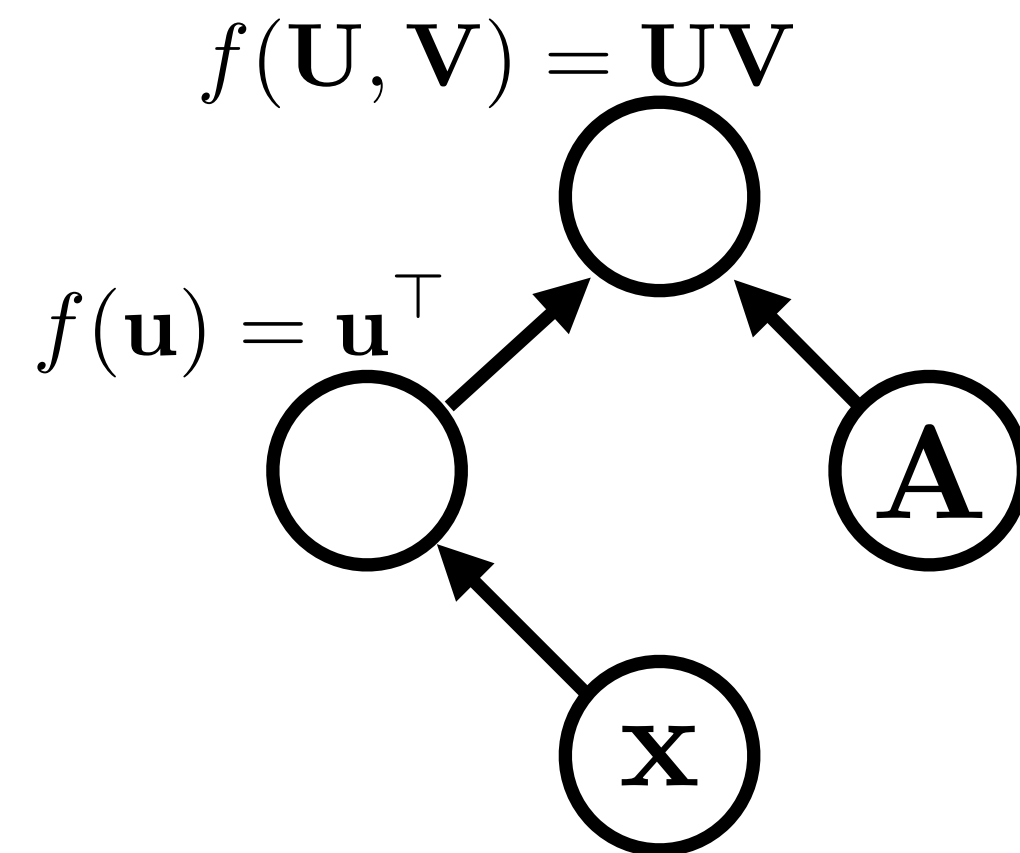


expression:

$$\mathbf{x}^\top \mathbf{A}$$

graph:

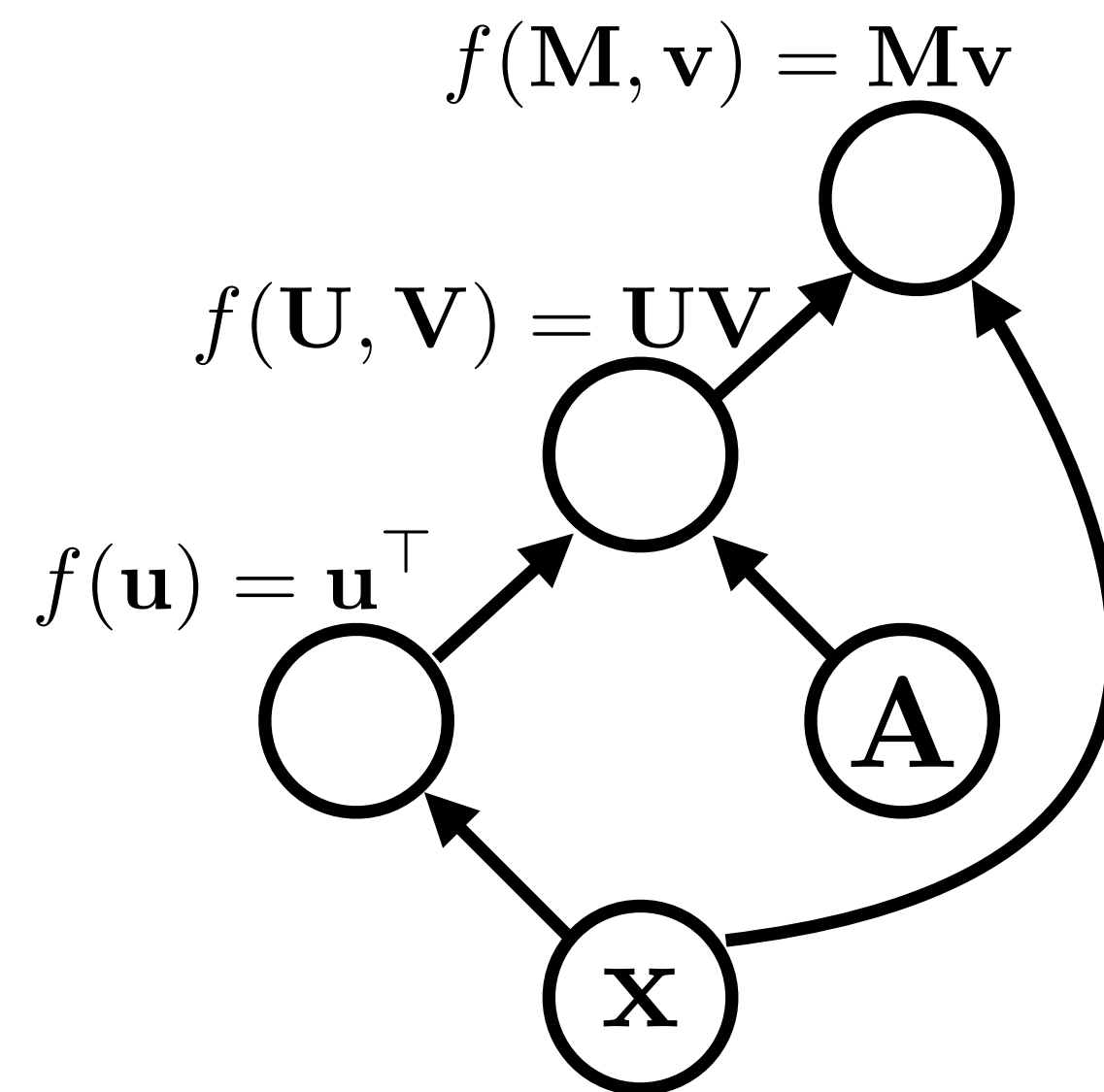
Functions can be nullary, unary,
binary, ... n -ary. Often they are unary or binary.



expression:

$$\mathbf{x}^\top \mathbf{A} \mathbf{x}$$

graph:

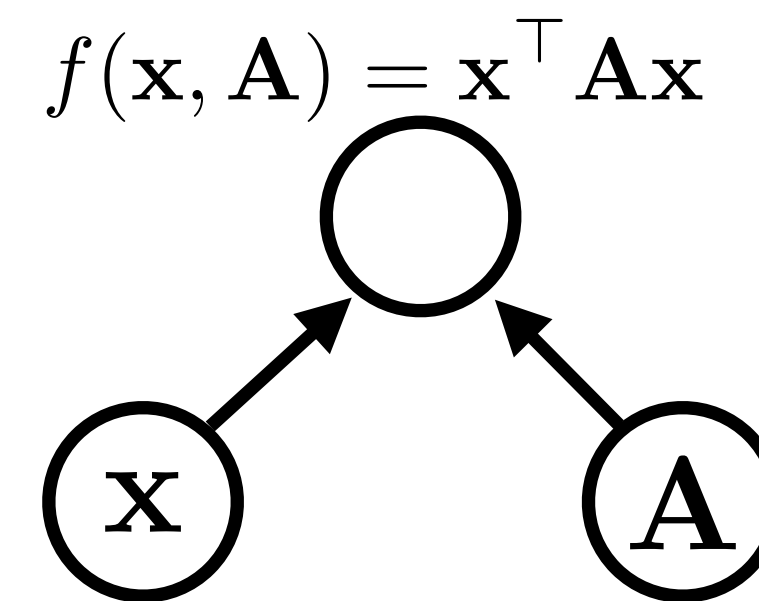
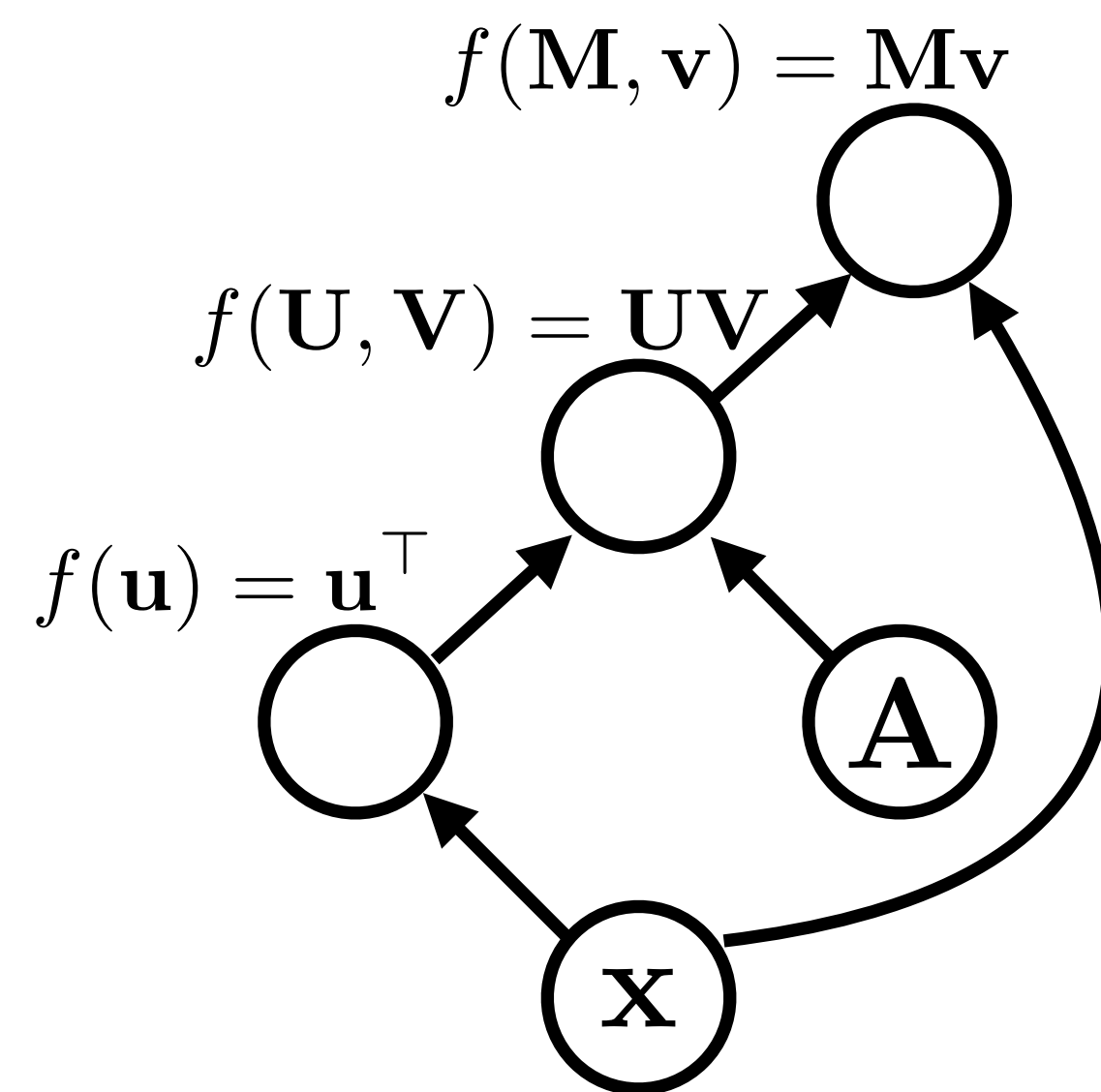


Computation graphs are directed and acyclic

expression:

$$\mathbf{x}^\top \mathbf{A} \mathbf{x}$$

graph:

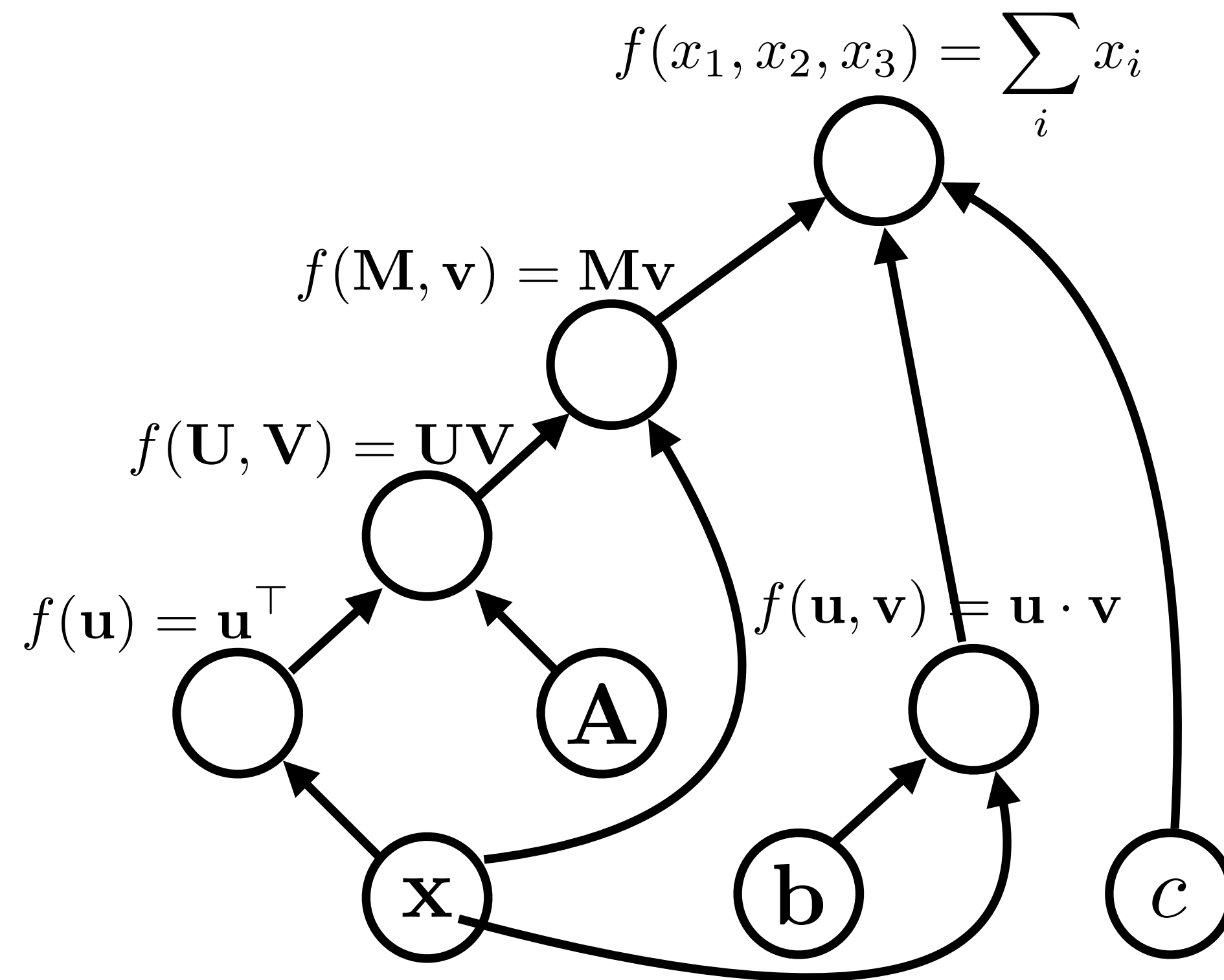


$$\frac{\partial f(\mathbf{x}, \mathbf{A})}{\partial \mathbf{x}} = (\mathbf{A}^\top + \mathbf{A})\mathbf{x}$$
$$\frac{\partial f(\mathbf{x}, \mathbf{A})}{\partial \mathbf{A}} = \mathbf{x}\mathbf{x}^\top$$

expression:

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} + \mathbf{b} \cdot \mathbf{x} + c$$

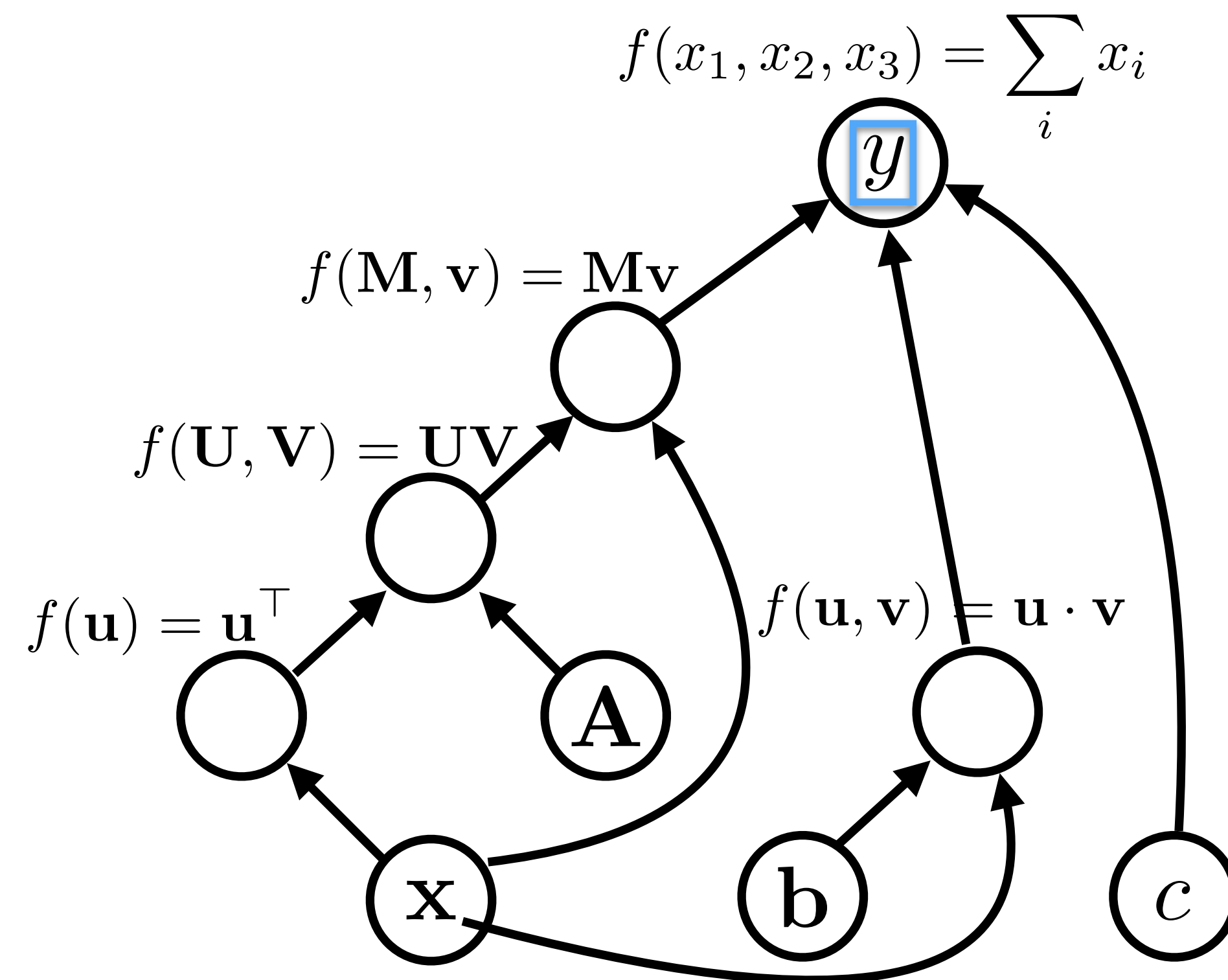
graph:



expression:

$$y = \mathbf{x}^\top \mathbf{A} \mathbf{x} + \mathbf{b} \cdot \mathbf{x} + c$$

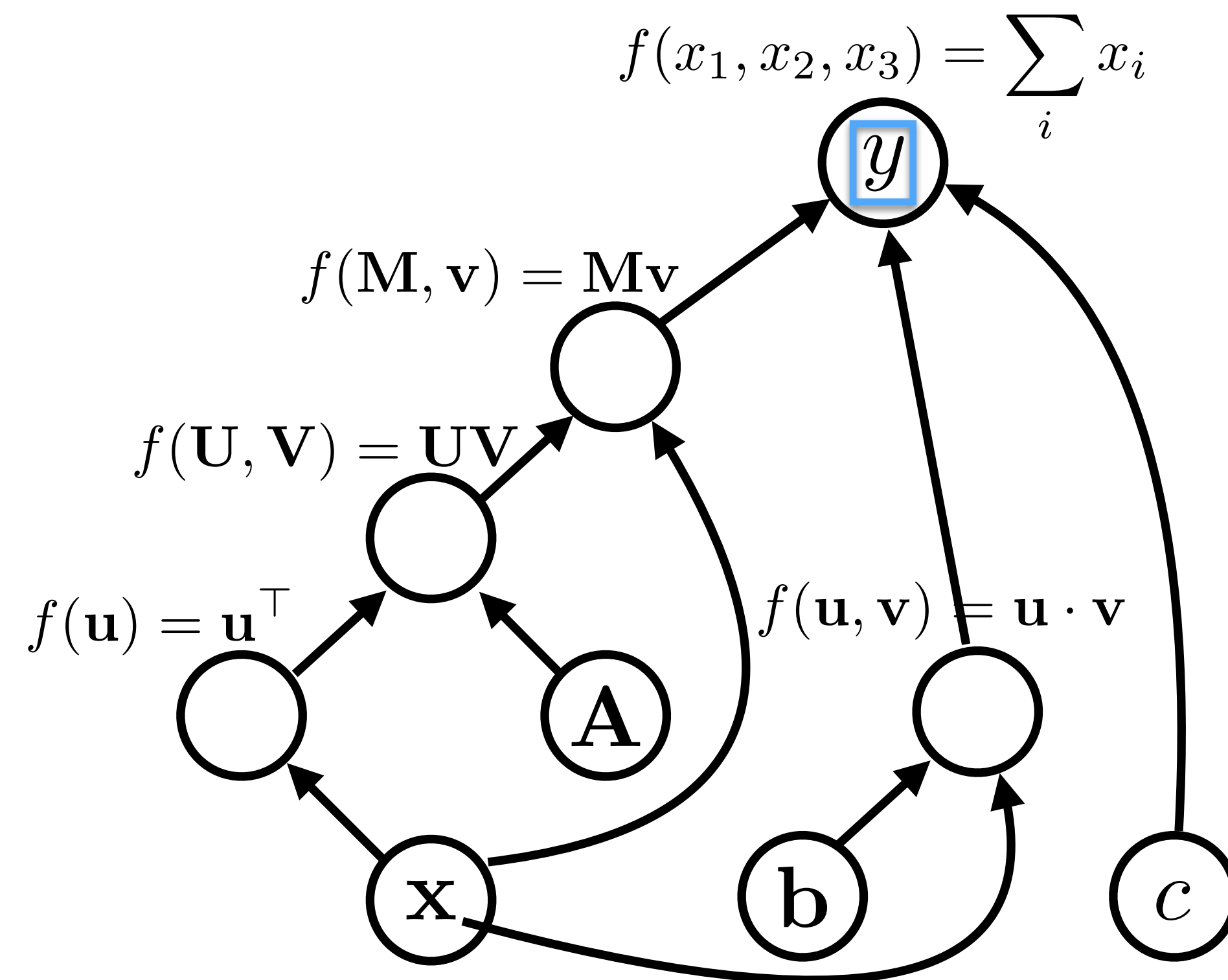
graph:



expression:

$$y = \mathbf{x}^\top \mathbf{A} \mathbf{x} + \mathbf{b} \cdot \mathbf{x} + c$$

graph:



variable names are just labelings of nodes.

Algorithms (1)

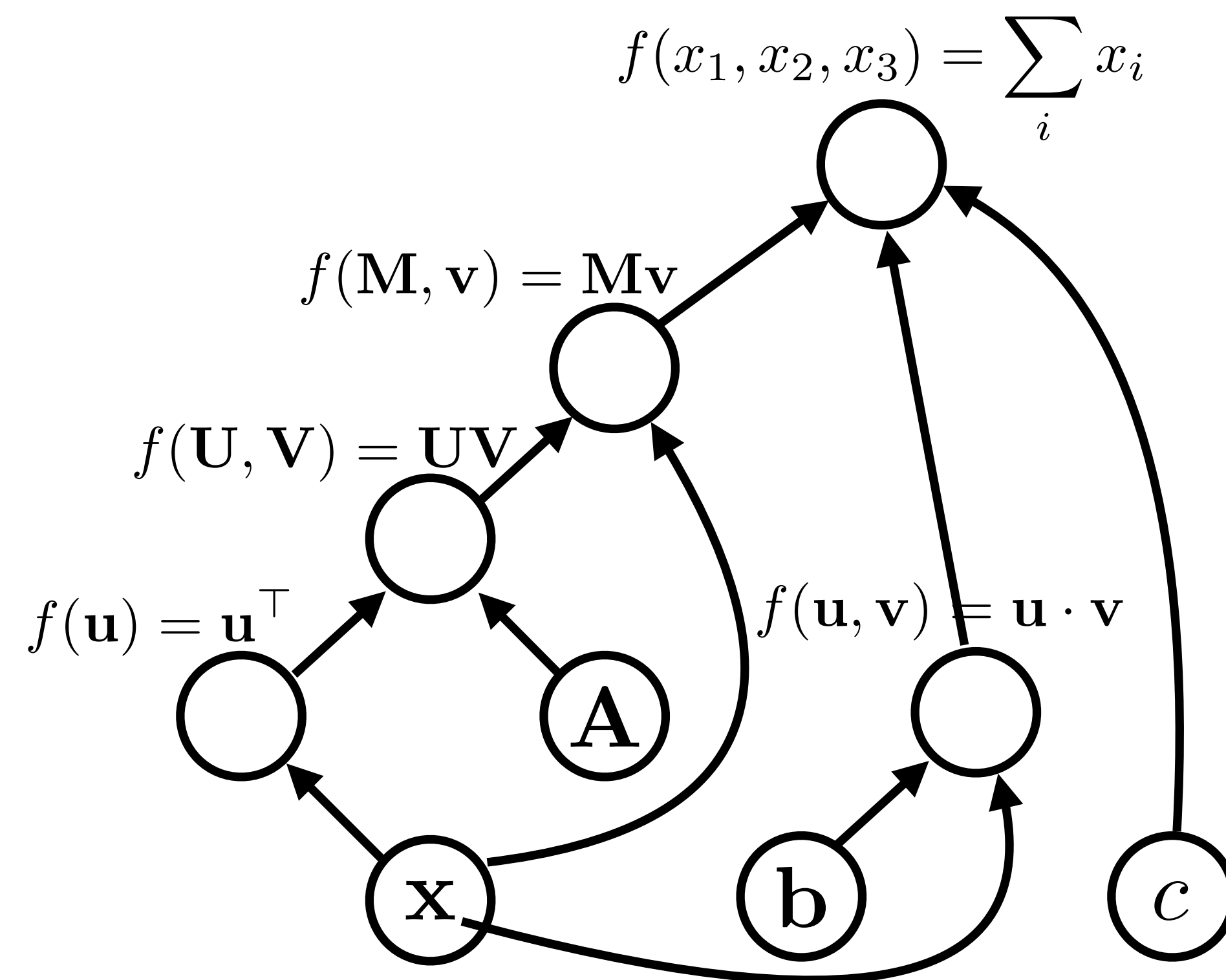
Graph construction

Forward propagation

In topological order, compute the **value** of the node given its inputs

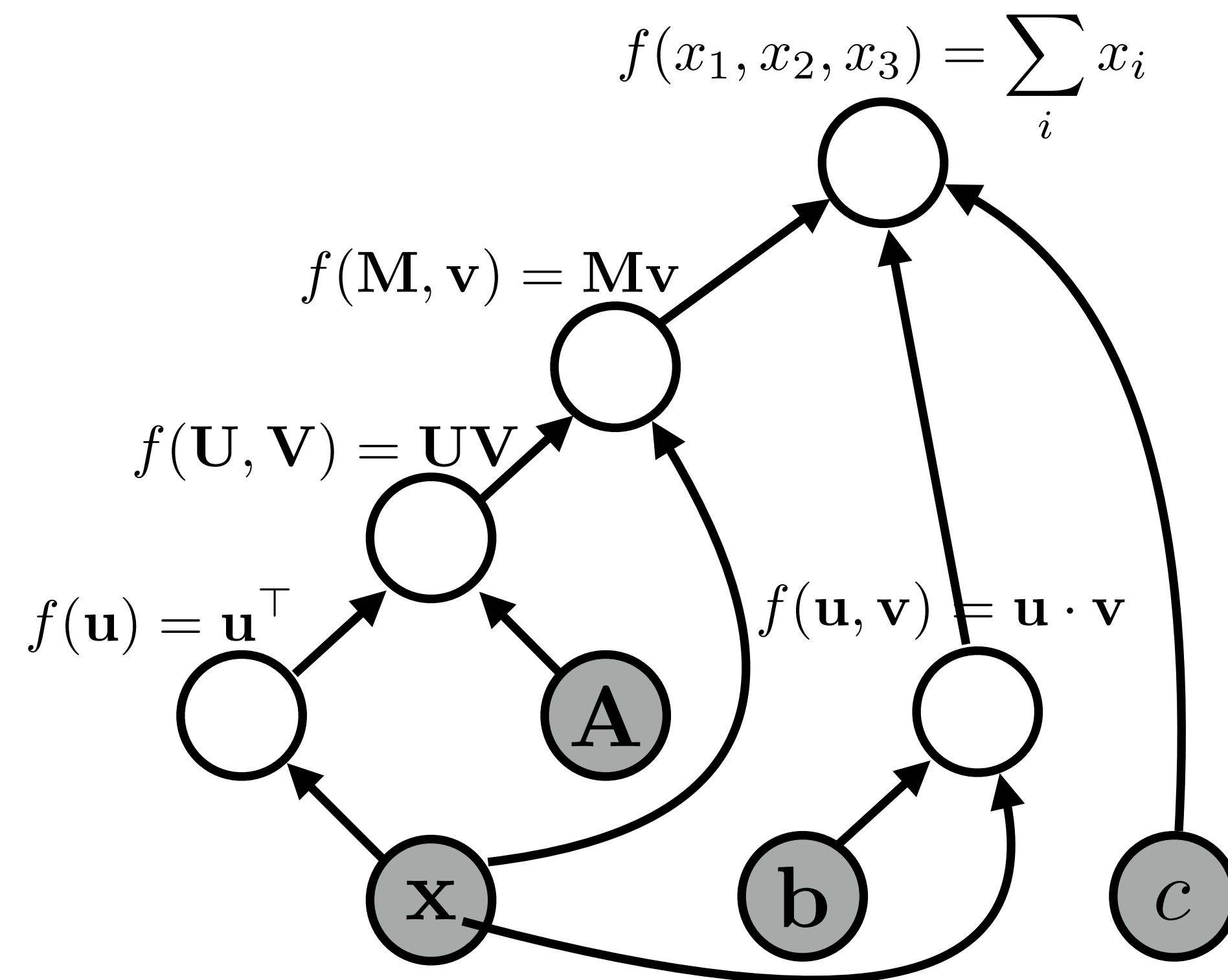
Forward Propagation

graph:



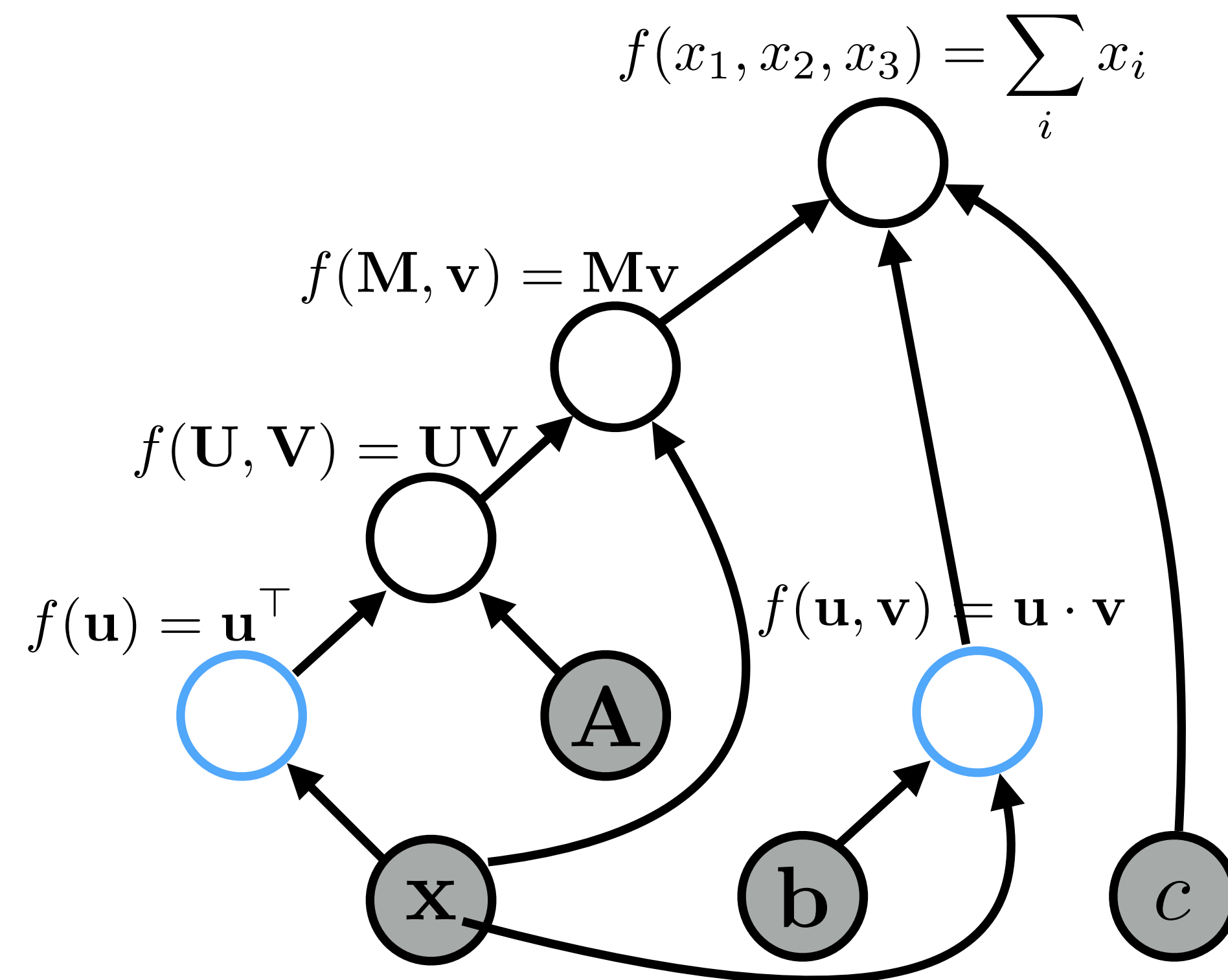
Forward Propagation

graph:



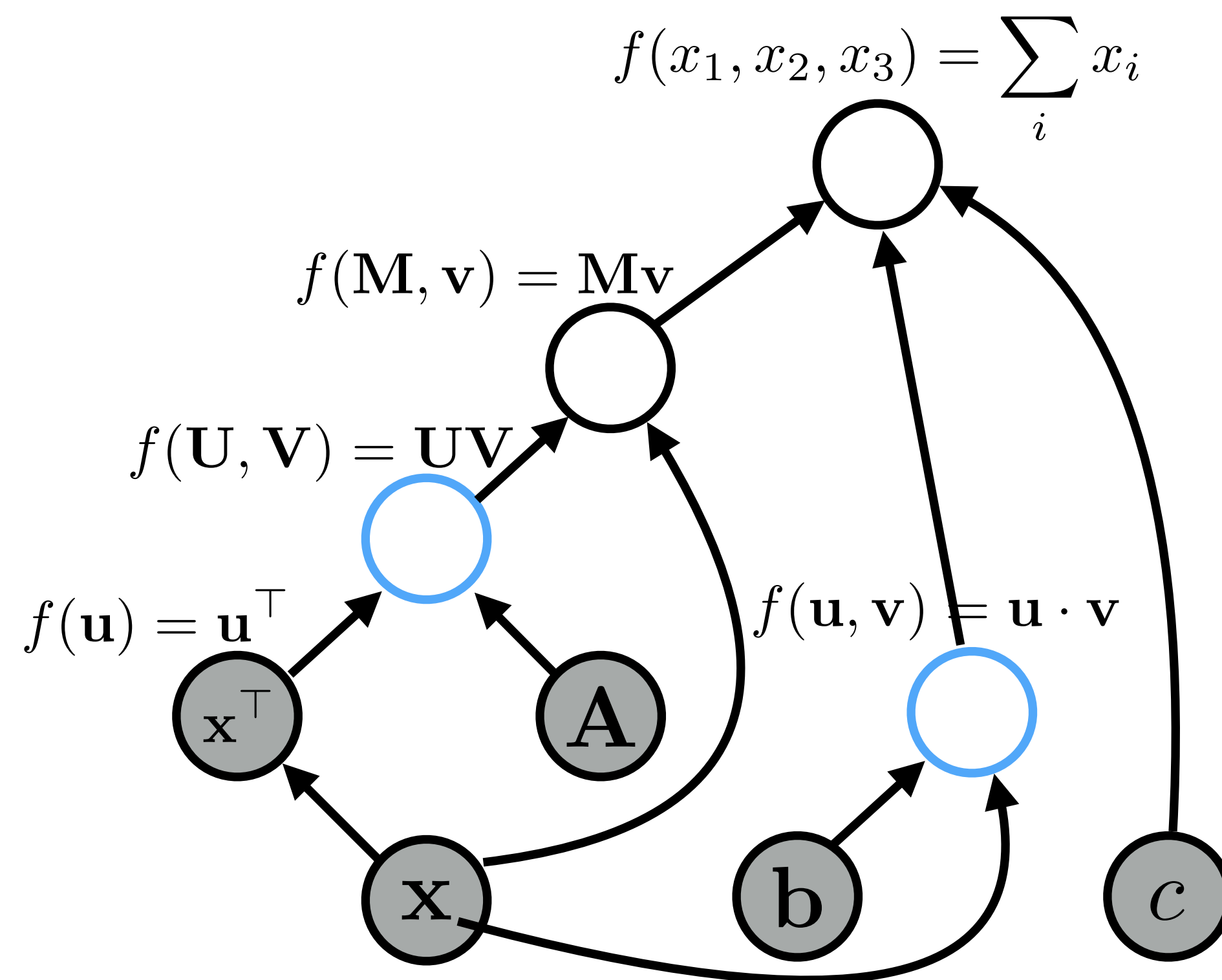
Forward Propagation

graph:



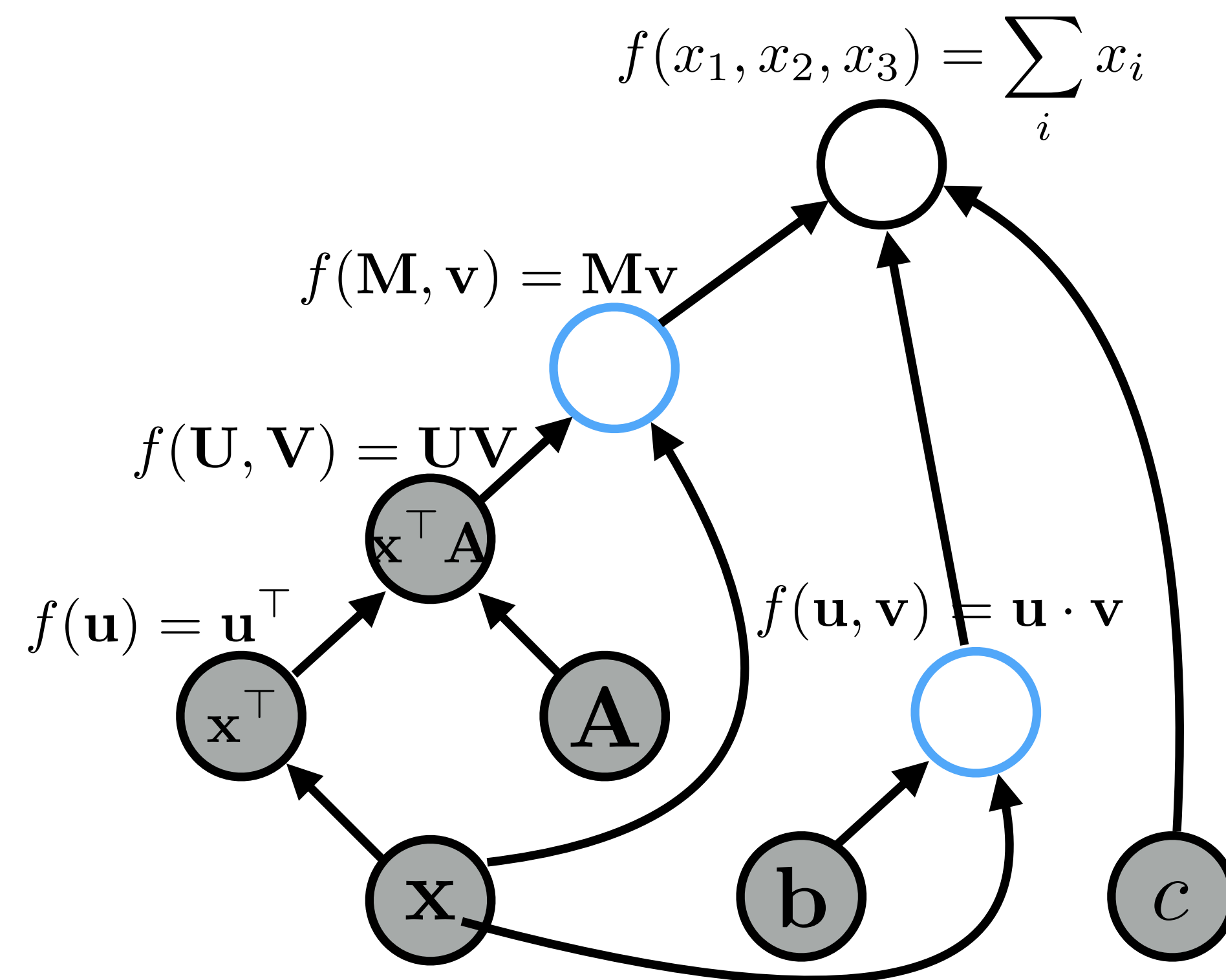
Forward Propagation

graph:



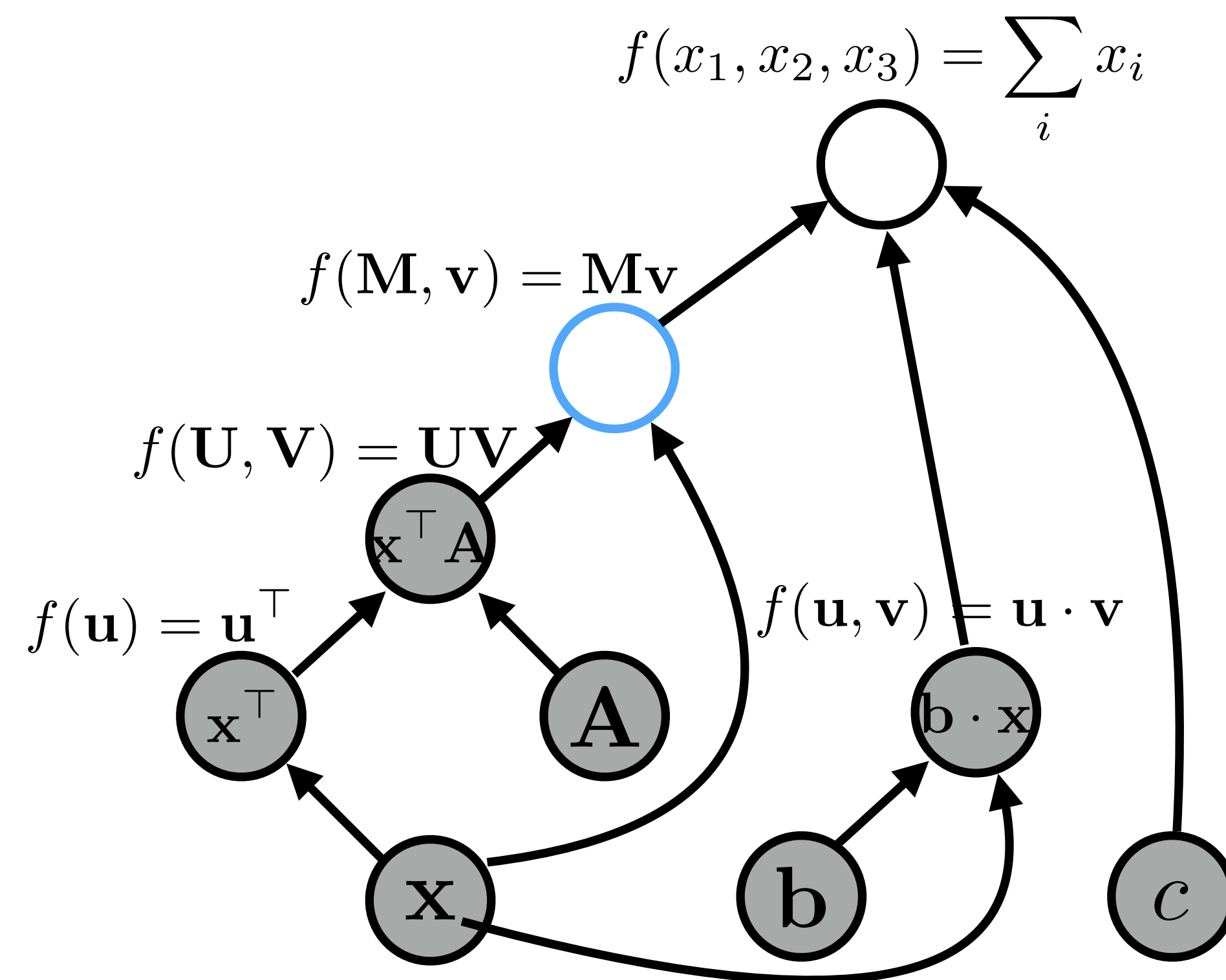
Forward Propagation

graph:



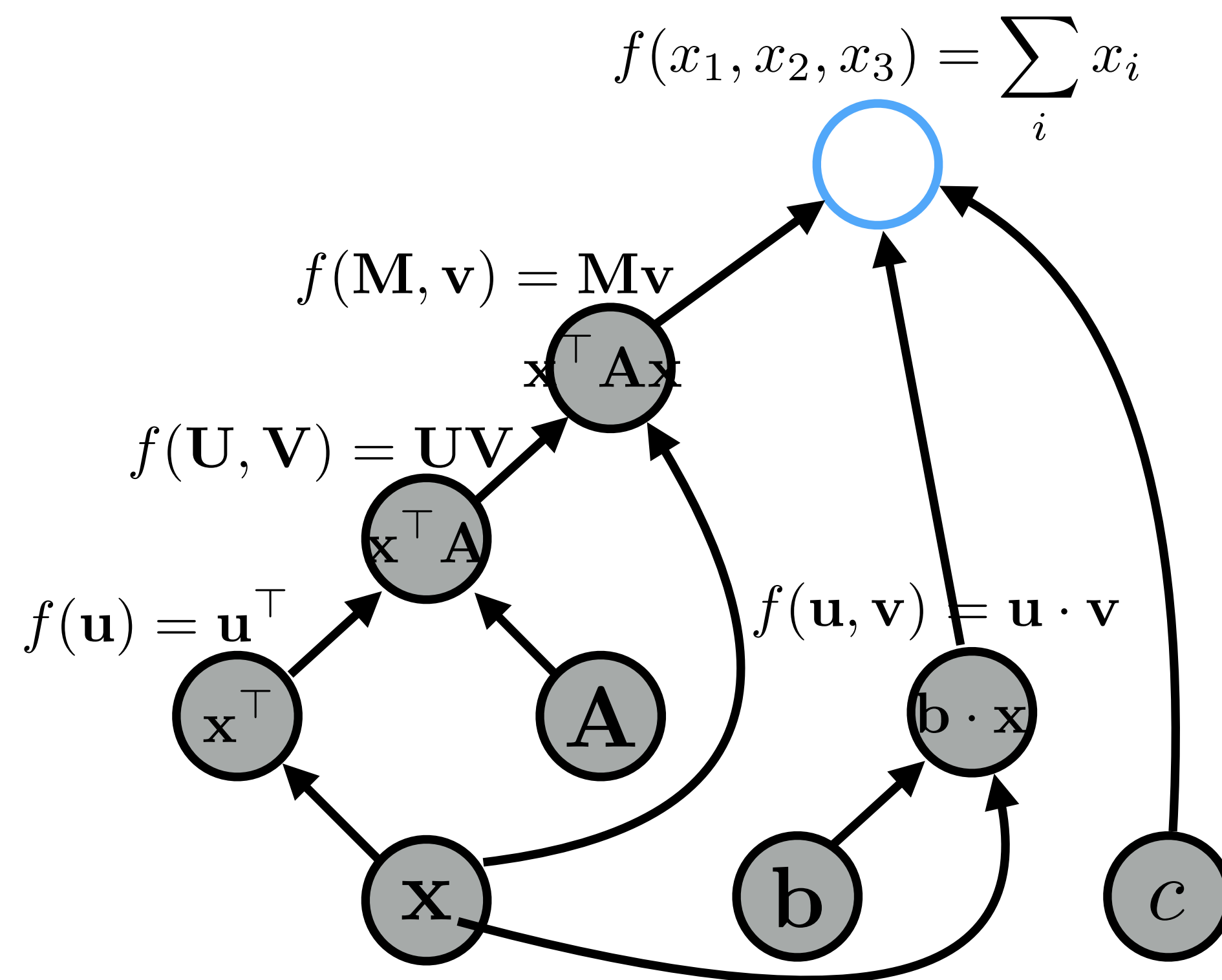
Forward Propagation

graph:



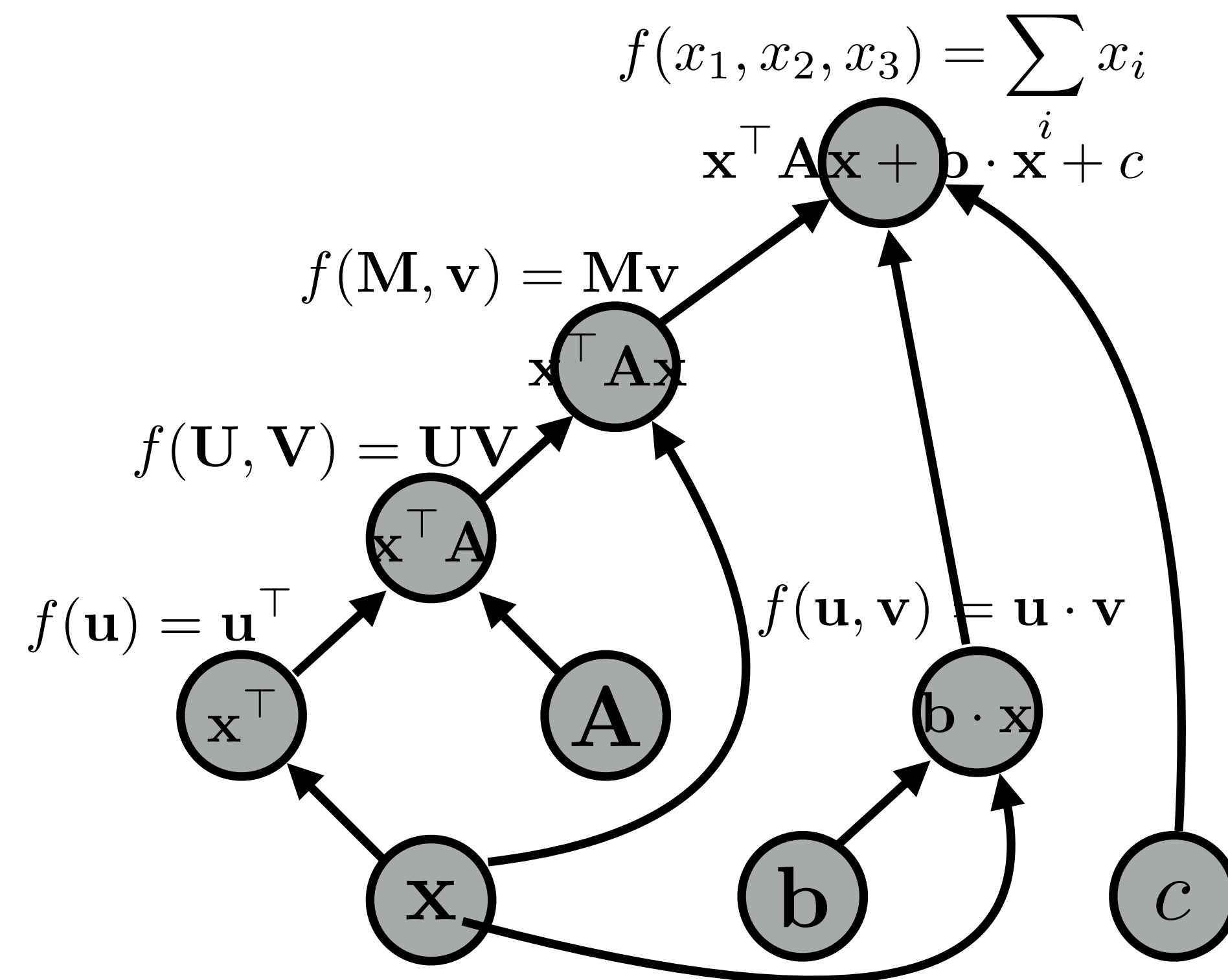
Forward Propagation

graph:



Forward Propagation

graph:



Algorithms (2)

Algorithms (2)

Back-propagation:

Process examples in reverse topological order

Calculate the derivatives of the parameters with respect to the final value
(This is usually a “loss function”, a value we want to minimize)

Algorithms (2)

Back-propagation:

Process examples in reverse topological order

Calculate the derivatives of the parameters with respect to the final value
(This is usually a “loss function”, a value we want to minimize)

Parameter update:

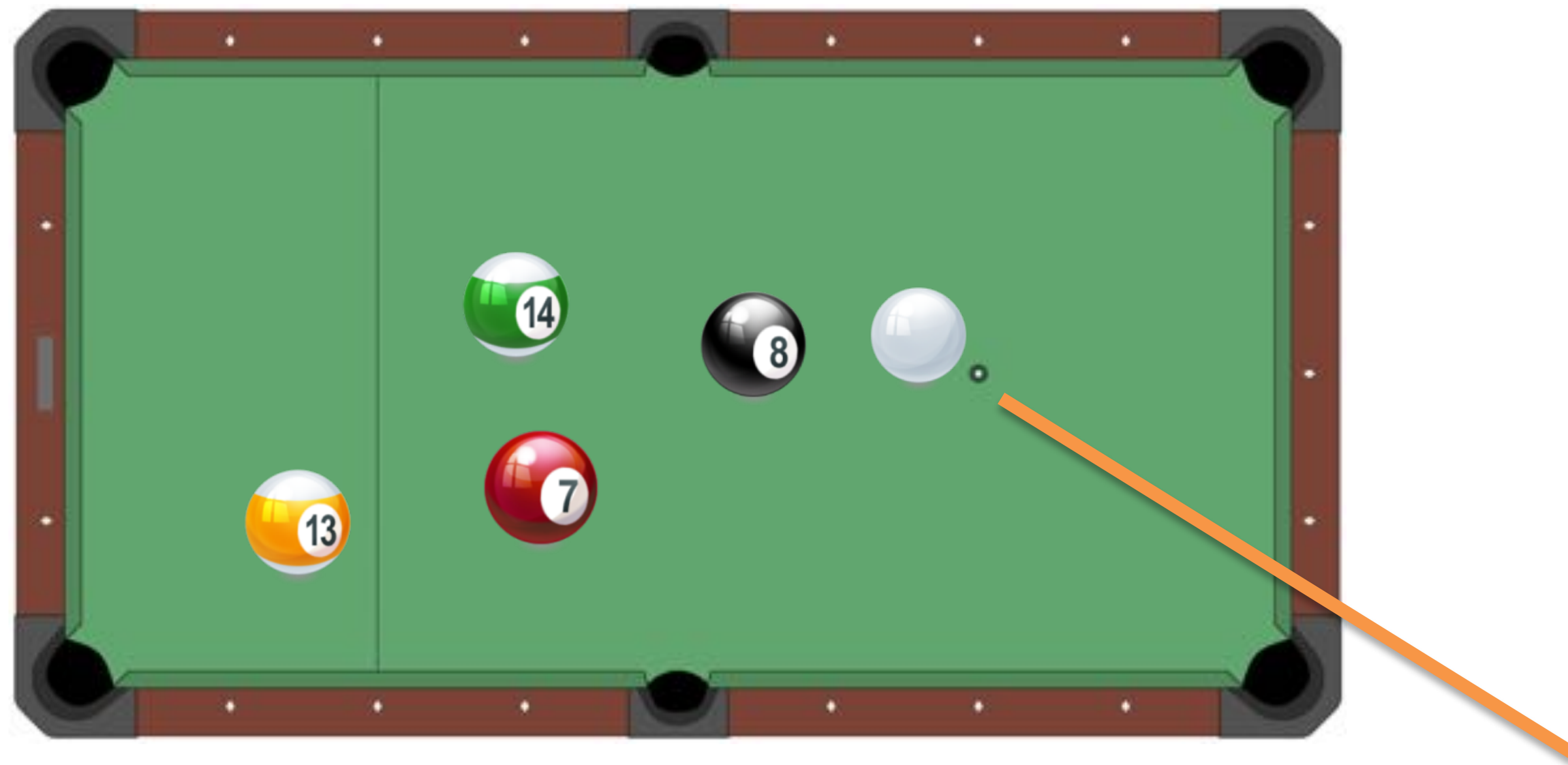
Move the parameters in the direction of this derivative

$$W -= \alpha * dl/dW$$

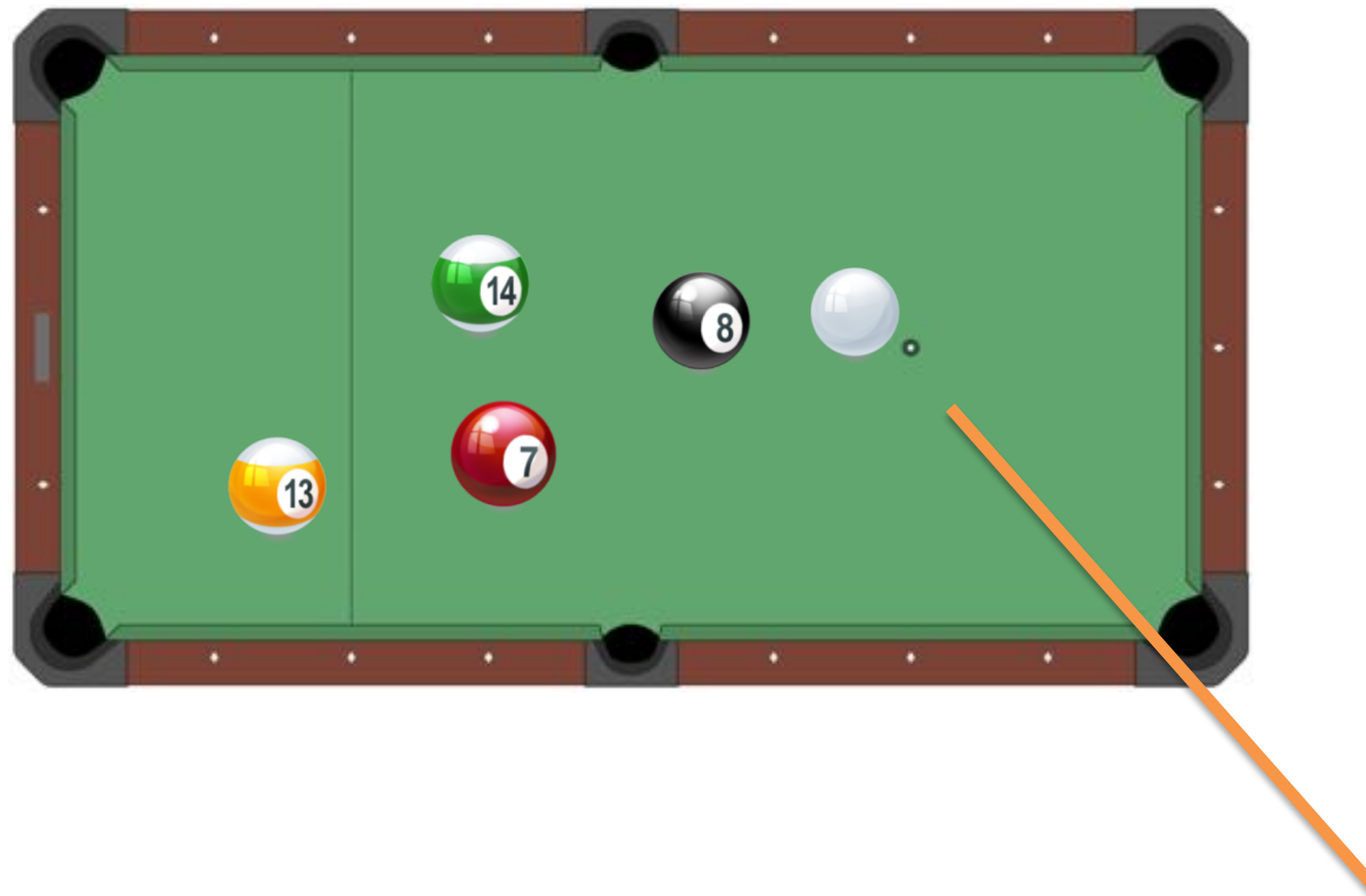
Intuitions

BACKPROPAGATION OF ERRORS

Error Back-Propagation

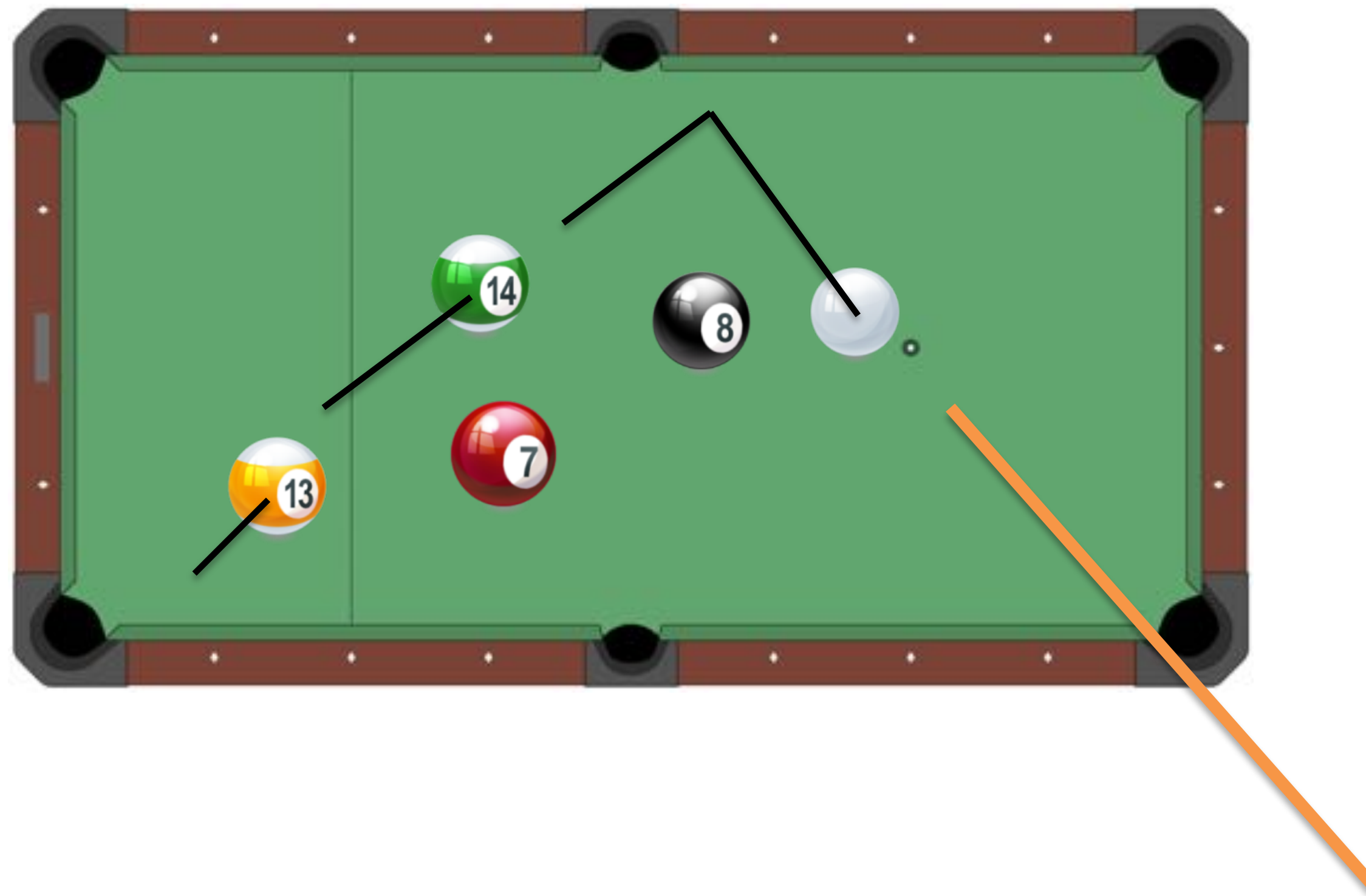


Error Back-Propagation

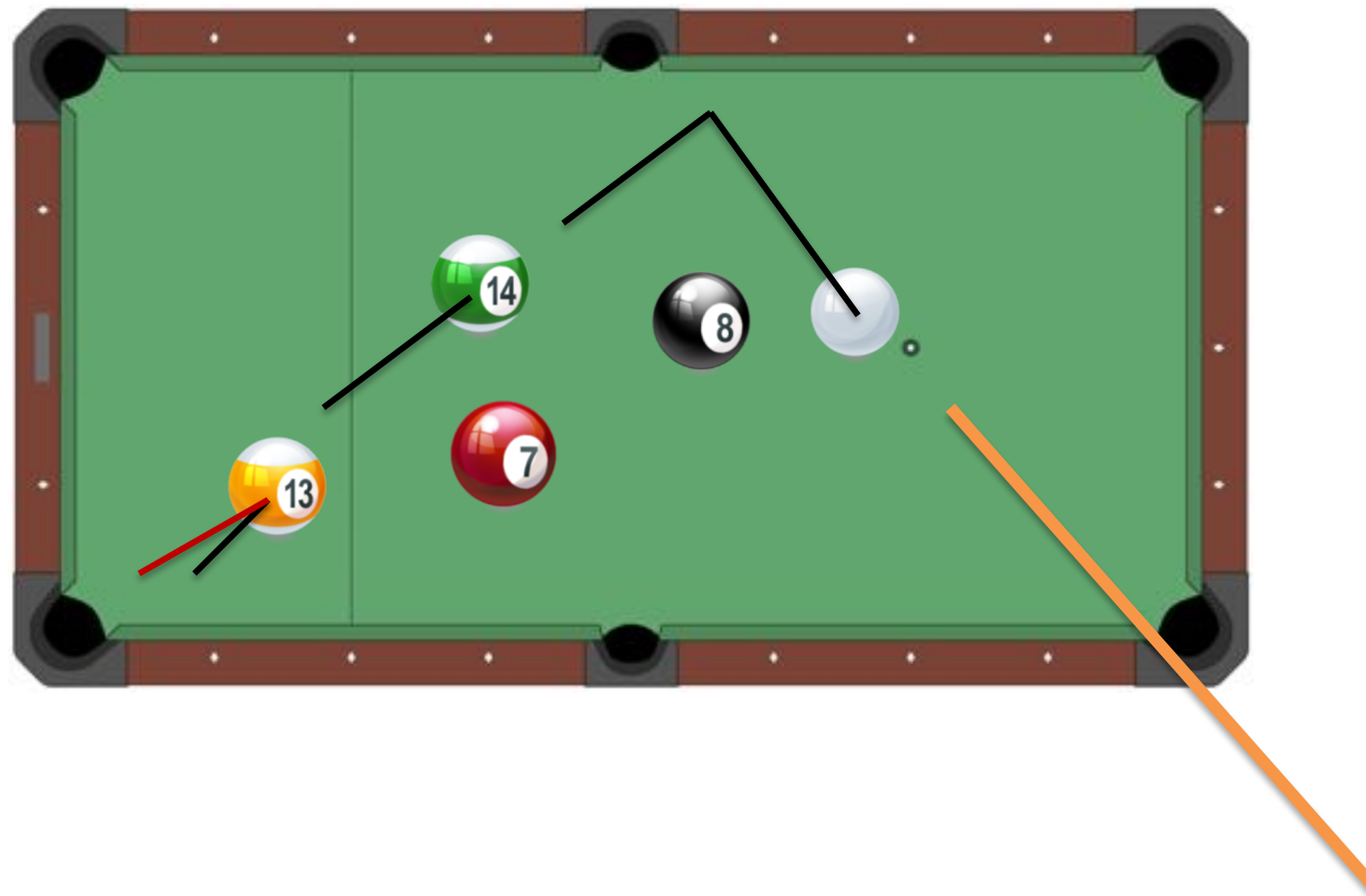


Slide from (Stoyanov & Eisner, 2012)

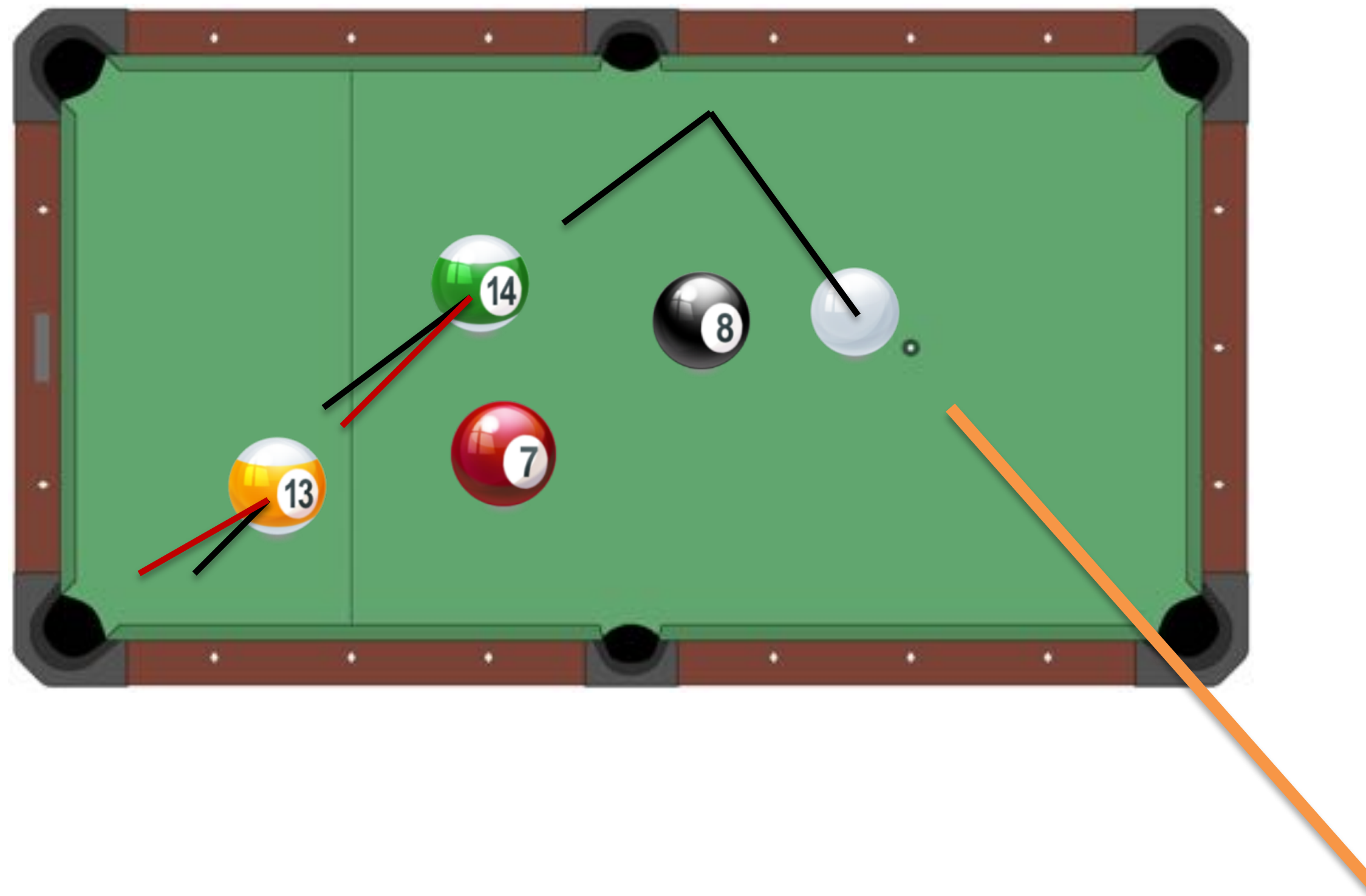
Error Back-Propagation



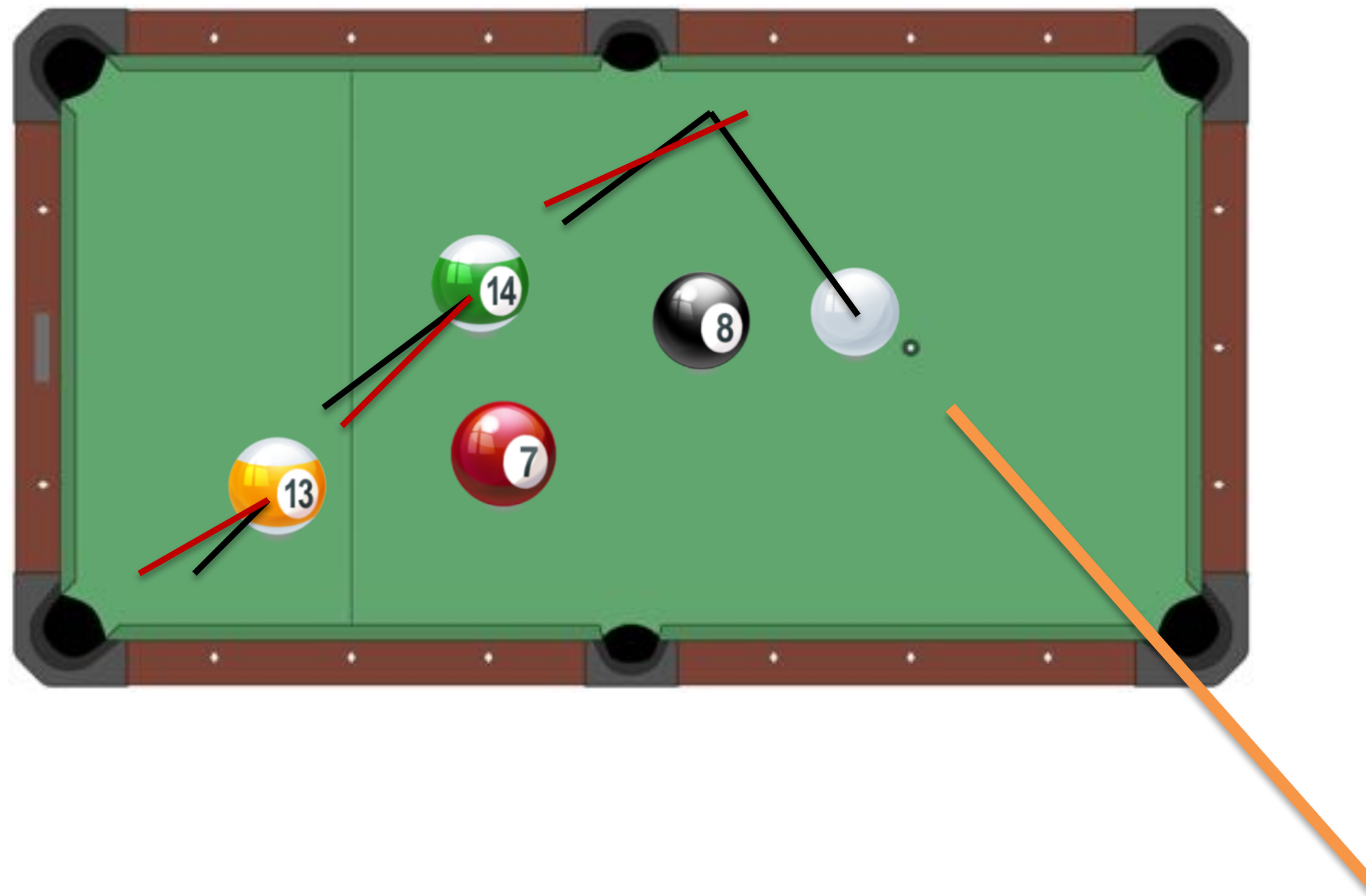
Error Back-Propagation



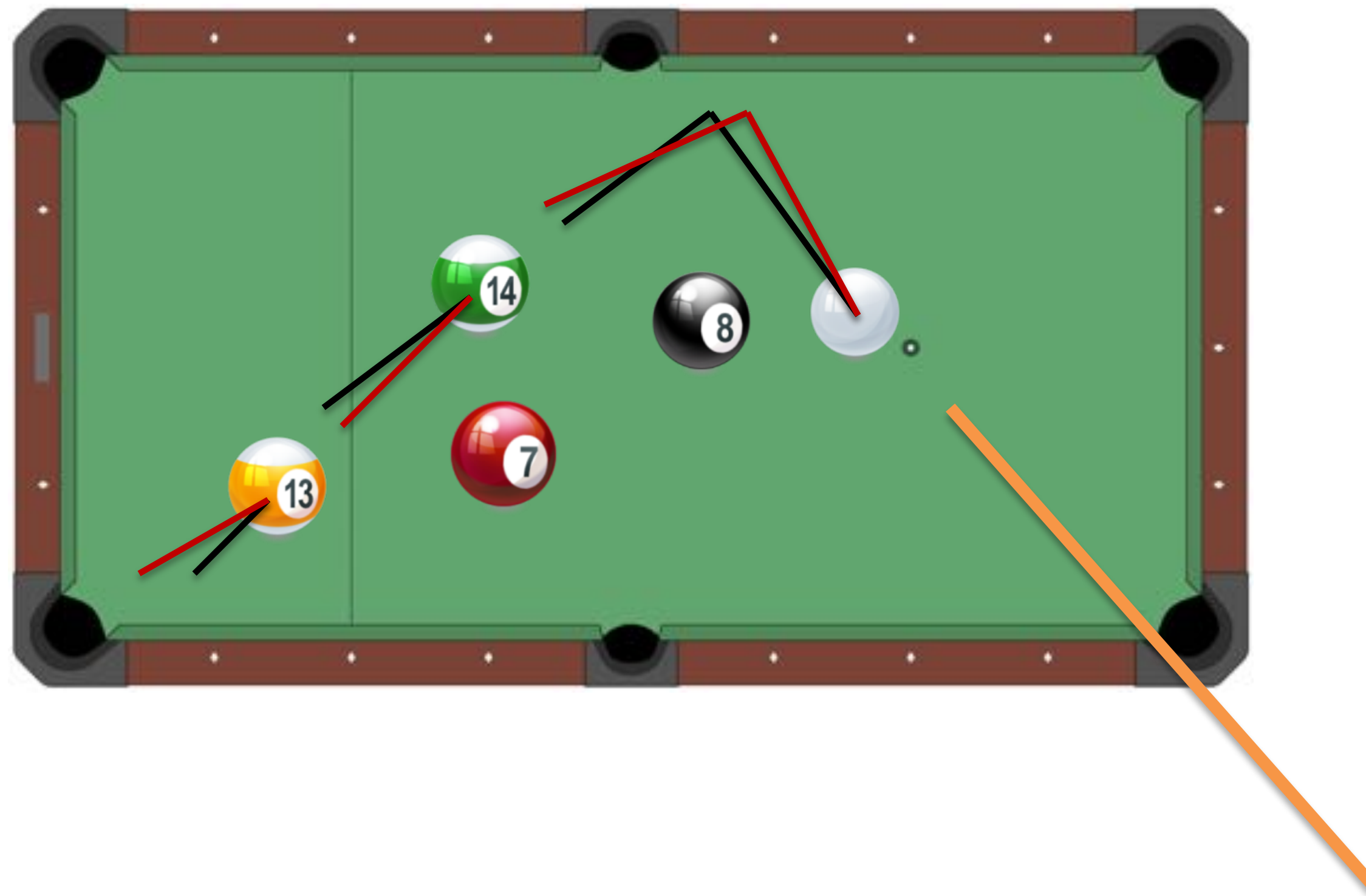
Error Back-Propagation



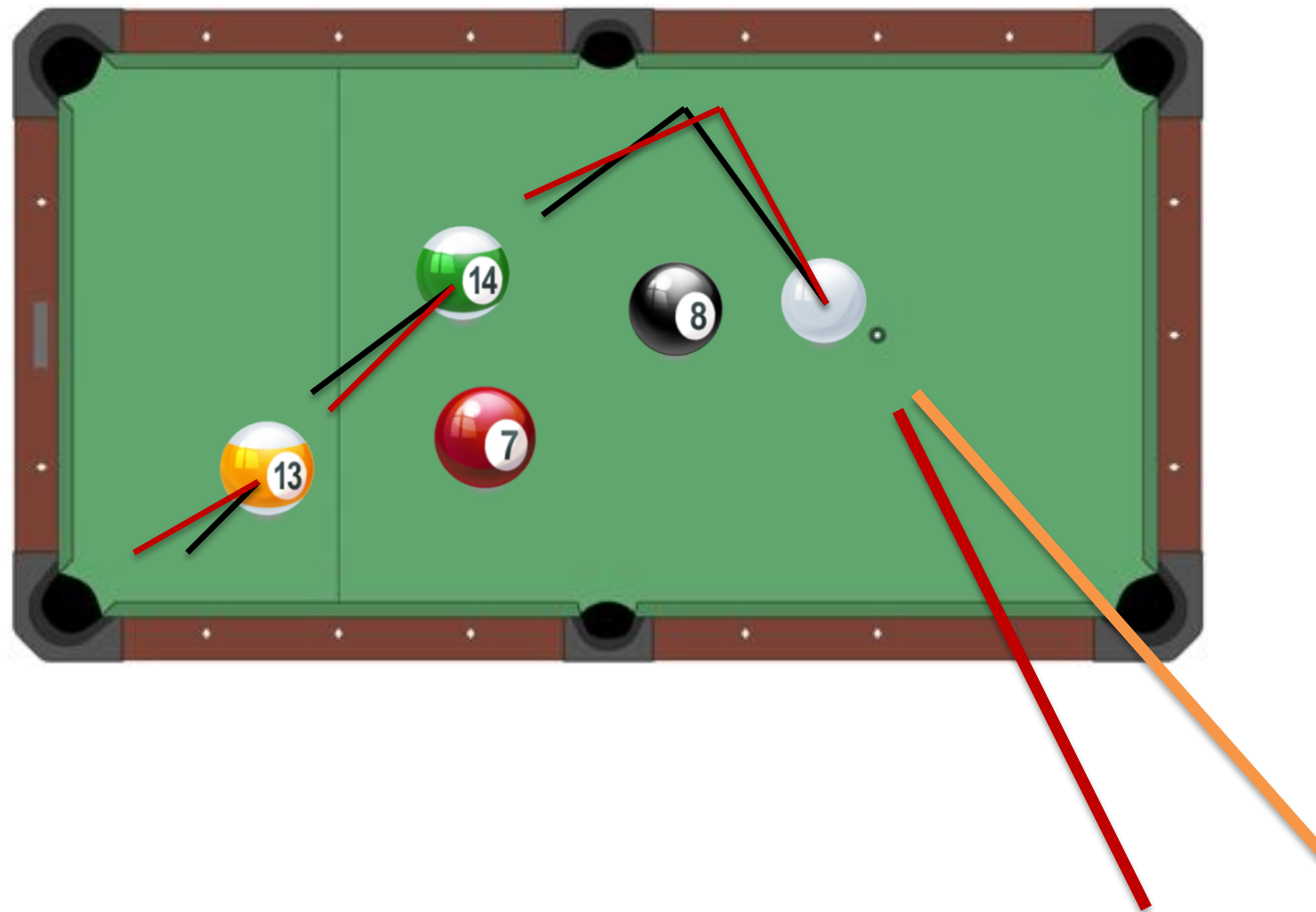
Error Back-Propagation



Error Back-Propagation

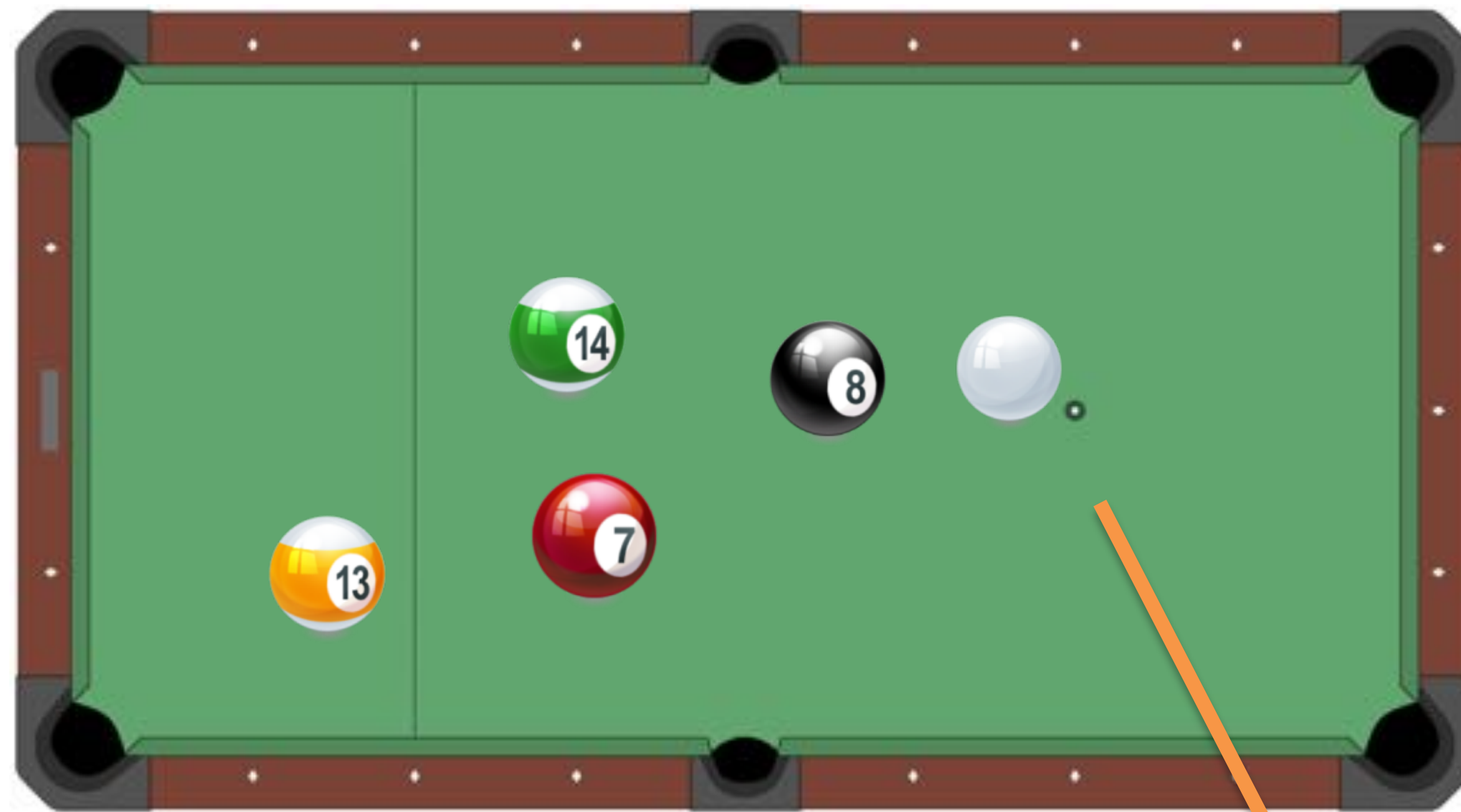


Error Back-Propagation



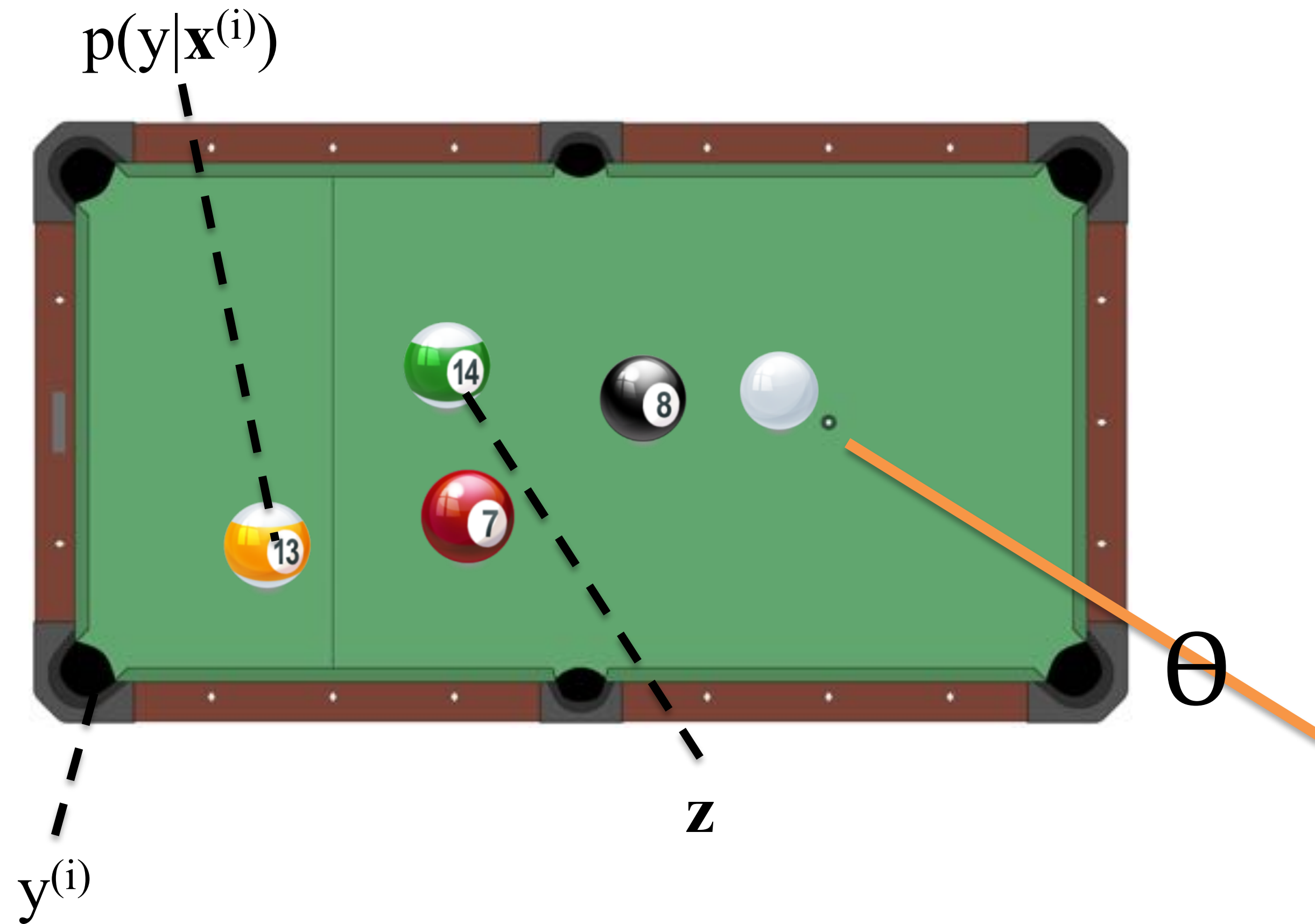
Slide from (Stoyanov & Eisner, 2012)

Error Back-Propagation



Slide from (Stoyanov & Eisner, 2012)

Error Back-Propagation



Language Models

- Language models are generative models of text

Language Models

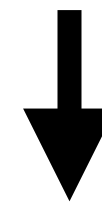
- Language models are generative models of text

$$s \sim P(x)$$

Language Models

- Language models are generative models of text

$$s \sim P(x)$$



“The Malfoys!” said Hermione.

Harry was watching him. He looked like Madame Maxime. When she strode up the wrong staircase to visit himself.

“I’m afraid I’ve definitely been suspended from power, no chance—indeed?” said Snape. He put his head back behind them and read groups as they crossed a corner and fluttered down onto their ink lamp, and picked up his spoon. The doorbell rang. It was a lot cleaner down in London.

The Generation Problem

The Generation Problem

- We have a model of $P(Y|X)$, how do we use it to generate a sentence?

The Generation Problem

- We have a model of $P(Y|X)$, how do we use it to generate a sentence?
- Two methods:

The Generation Problem

- We have a model of $P(Y|X)$, how do we use it to generate a sentence?
- Two methods:
 - **Sampling:** Try to generate a *random* sentence according to the probability distribution.

The Generation Problem

- We have a model of $P(Y|X)$, how do we use it to generate a sentence?
- Two methods:
 - **Sampling:** Try to generate a *random* sentence according to the probability distribution.
 - **Argmax:** Try to generate the sentence with the *highest* probability.

Ancestral Sampling

Ancestral Sampling

- **Randomly generate** words one-by-one.

Ancestral Sampling

- **Randomly generate** words one-by-one.

```
while  $y_{j-1} \neq \text{"</s>"}$ :  
   $y_j \sim P(y_j \mid X, y_1, \dots, y_{j-1})$ 
```

Ancestral Sampling

- **Randomly generate** words one-by-one.

```
while  $y_{j-1} \neq \text{"</s>"}$ :  
   $y_j \sim P(y_j \mid X, y_1, \dots, y_{j-1})$ 
```

- An **exact method** for sampling from $P(X)$, no further work needed.

Greedy Search

Greedy Search

- One by one, pick the single highest-probability word

Greedy Search

- One by one, pick the single highest-probability word

```
while  $y_{j-1} \neq \text{"</s>"}$ :  
   $y_j = \operatorname{argmax} P(y_j \mid X, y_1, \dots, y_{j-1})$ 
```

Greedy Search

- One by one, pick the single highest-probability word

```
while  $y_{j-1} \neq \text{"</s>"}$ :  
   $y_j = \operatorname{argmax} P(y_j \mid X, y_1, \dots, y_{j-1})$ 
```

- **Not exact, real problems:**

Greedy Search

- One by one, pick the single highest-probability word

```
while  $y_{j-1} \neq \text{"</s>"}$ :  
   $y_j = \operatorname{argmax} P(y_j \mid X, y_1, \dots, y_{j-1})$ 
```

- **Not exact, real problems:**
 - Will often generate the “easy” words first

Greedy Search

- One by one, pick the single highest-probability word

```
while  $y_{j-1} \neq \text{"</s>"}$ :  
   $y_j = \operatorname{argmax} P(y_j \mid X, y_1, \dots, y_{j-1})$ 
```

- **Not exact, real problems:**
 - Will often generate the “easy” words first
 - Will prefer multiple common words to one rare word