

CS11-711 Advanced NLP

# Recurrent Neural Networks

Graham Neubig



**Carnegie Mellon University**

Language Technologies Institute

Site

<https://phontron.com/class/anlp2021/>

# NLP and Sequential Data

- NLP is full of sequential data
  - Words in sentences
  - Characters in words
  - Sentences in discourse
  - ....

# Long-distance Dependencies in Language

- Agreement in number, gender, etc.

**He** does not have very much confidence in **himself**.  
**She** does not have very much confidence in **herself**.

- Selectional preference

The **reign** has lasted as long as the life of the **queen**.  
The **rain** has lasted as long as the life of the **clouds**.

# Can be Complicated!

- What is the referent of “it”?

The trophy would not fit in the brown suitcase because it was too **big**.

Trophy

The trophy would not fit in the brown suitcase because it was too **small**.

Suitcase

(from Winograd Schema Challenge:

<http://commonsensereasoning.org/winograd.html>)

An Aside:

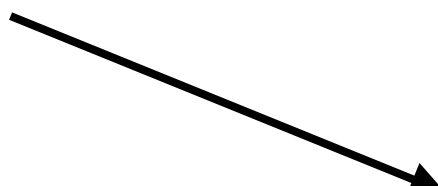
Types of Prediction Problems

# Types of Prediction: Binary, Multi-class, Structured


- Two classes (**binary classification**)

I hate this movie  positive  
negative

- Multiple classes (**multi-class classification**)

I hate this movie  very good  
good  
neutral  
bad  
very bad

- Exponential/infinite labels (**structured prediction**)

I hate this movie  PRP VBP DT NN

I hate this movie  *kono eiga ga kirai*

# Types of Prediction: Unconditioned vs. Conditioned

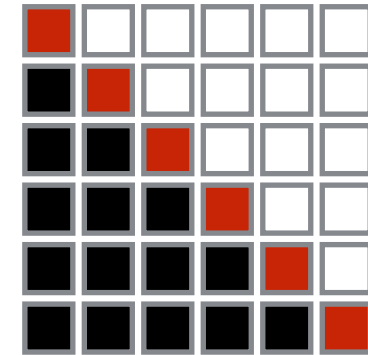
- **Unconditioned Prediction:** Predict the probability of a single variable  $P(X)$
- **Conditioned Prediction:** Predict the probability of an output variable given an input  $P(Y|X)$

# Types of Unconditioned Prediction

*Left-to-right Autoregressive Prediction*

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

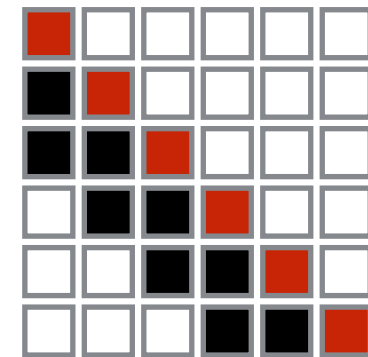
(e.g. RNN LM)



*Left-to-right Markov Chain (order n-1)*

$$P(X) = \prod_{i=1}^{|X|} P(x_i | x_{i-n+1}, \dots, x_{i-1})$$

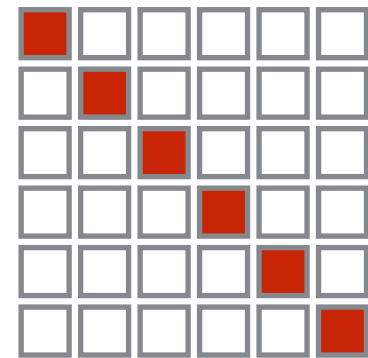
(e.g. n-gram LM, feed-forward LM)



*Independent Prediction*

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

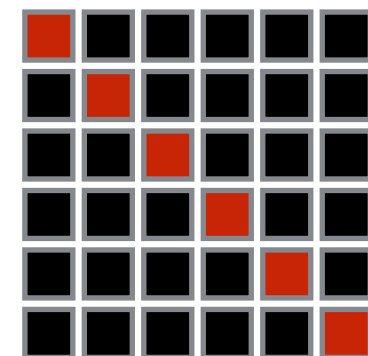
(e.g. unigram model)



*Bidirectional Prediction*

$$P(X) \neq \prod_{i=1}^{|X|} P(x_i | x_{\neq i})$$

(e.g. masked language model)





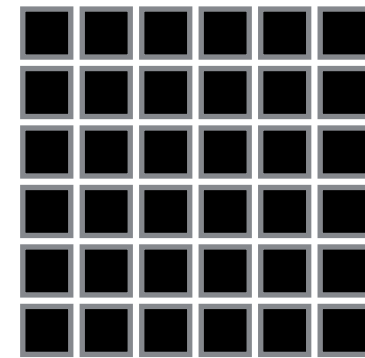
# Types of Conditioned Prediction

*Autoregressive Conditioned Prediction*

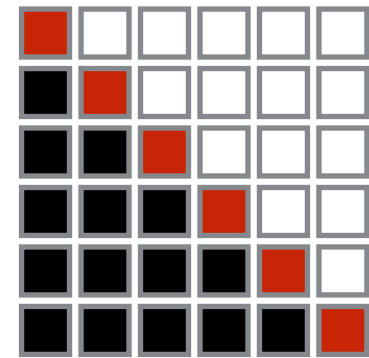
$$P(Y|X) = \prod_{i=1}^{|Y|} P(y_i|X, y_1, \dots, y_{i-1})$$

*(e.g. seq2seq model)*

*Source X*



*Target Y*

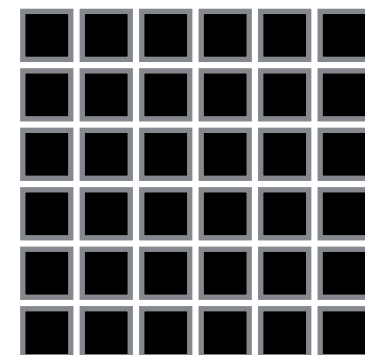


*Non-autoregressive Conditioned Prediction*

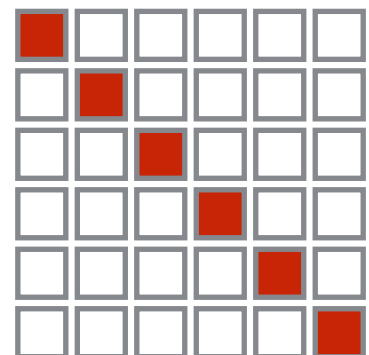
$$P(Y|X) = \prod_{i=1}^{|Y|} P(y_i|X)$$

*(e.g. sequence labeling, non-autoregressive MT)*

*Source X*



*Target Y*



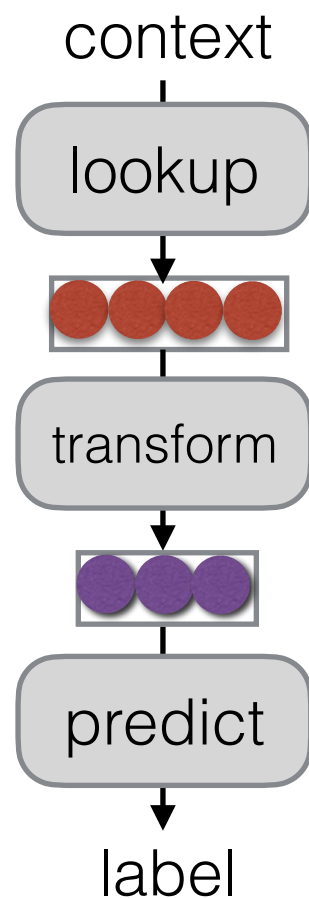
# Recurrent Neural Networks

# Recurrent Neural Networks

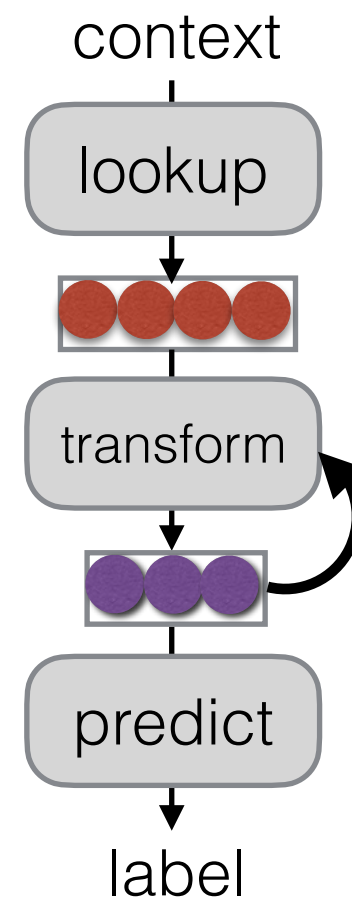
(Elman 1990)

- Tools to “remember” information

Feed-forward NN

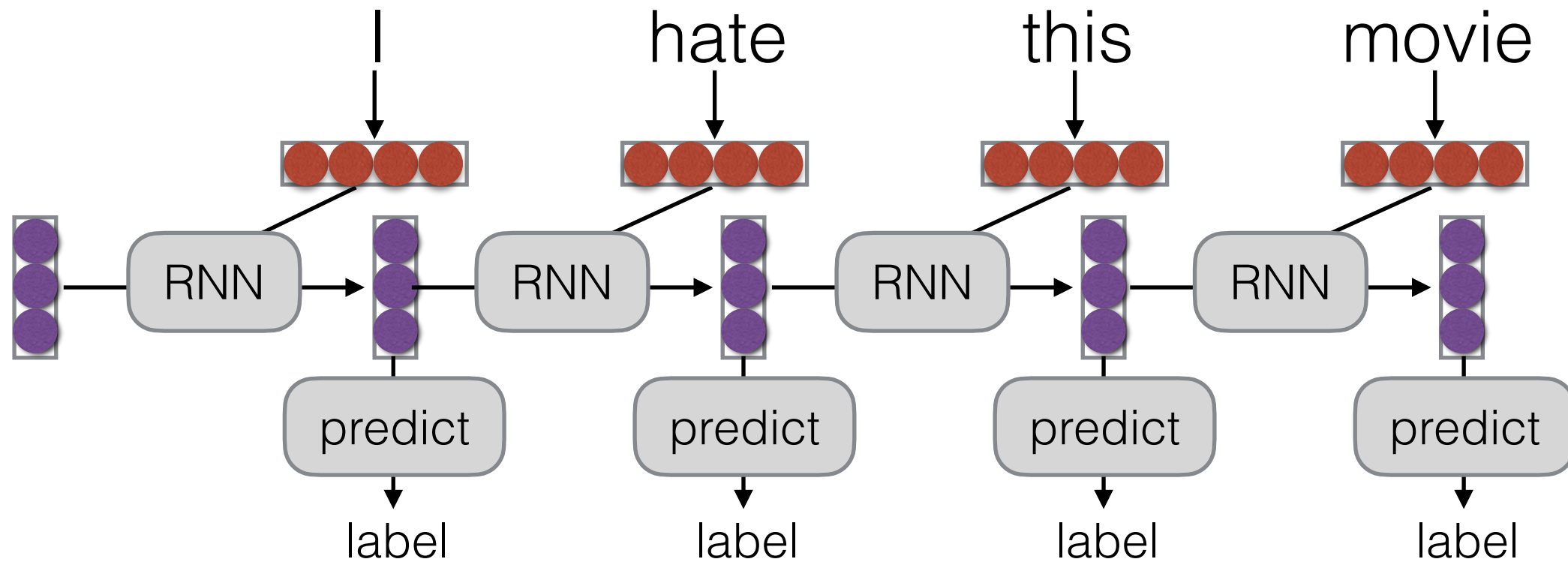


Recurrent NN

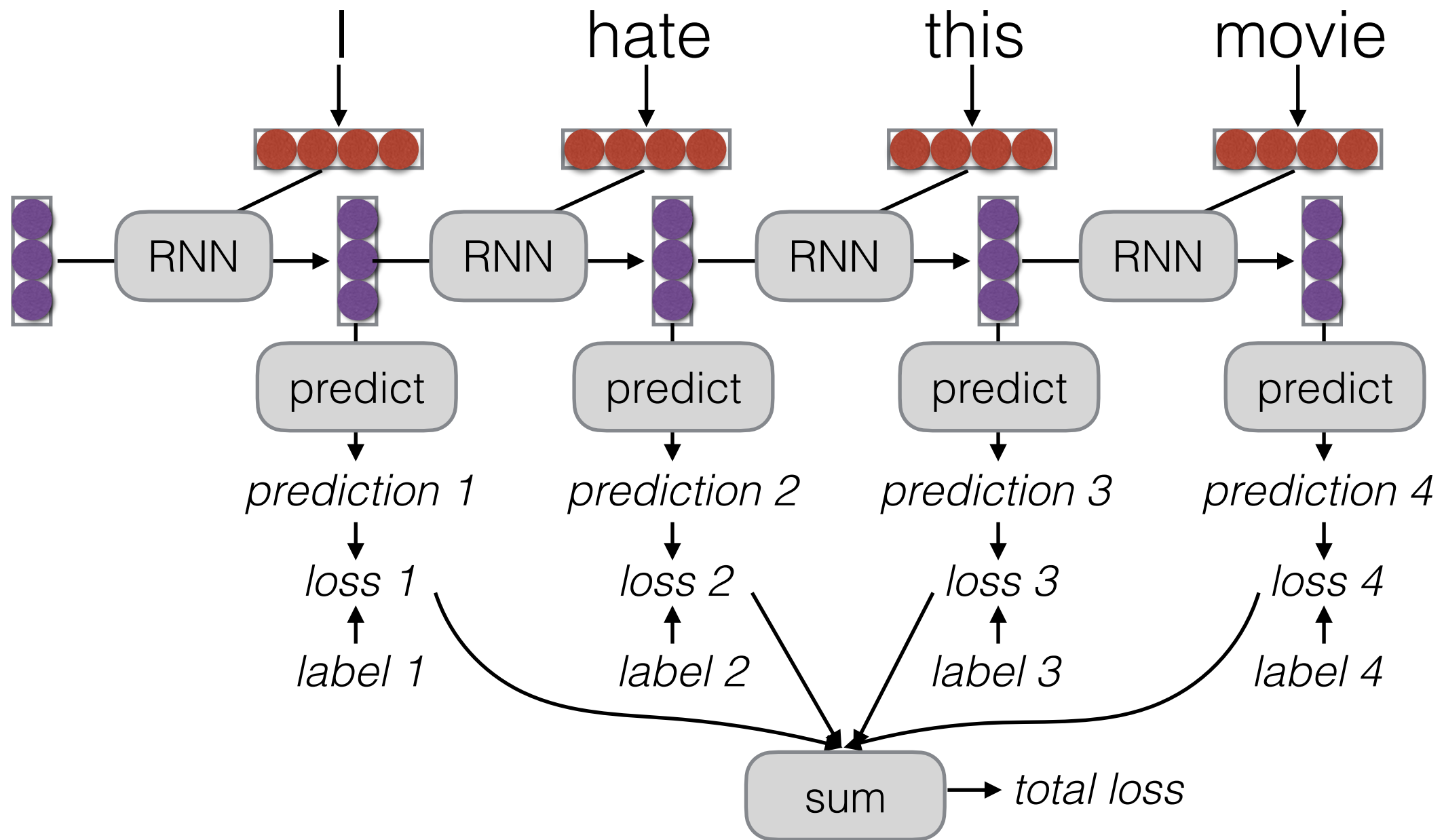


# Unrolling in Time

- What does processing a sequence look like?

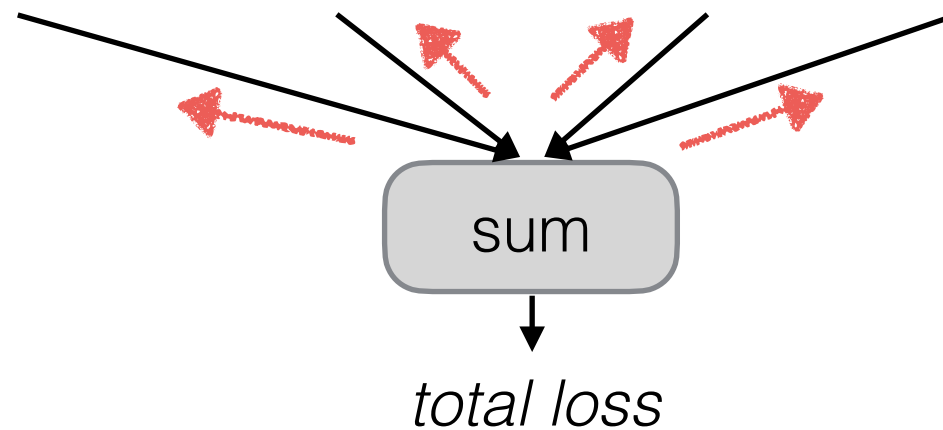


# Training RNNs



# RNN Training

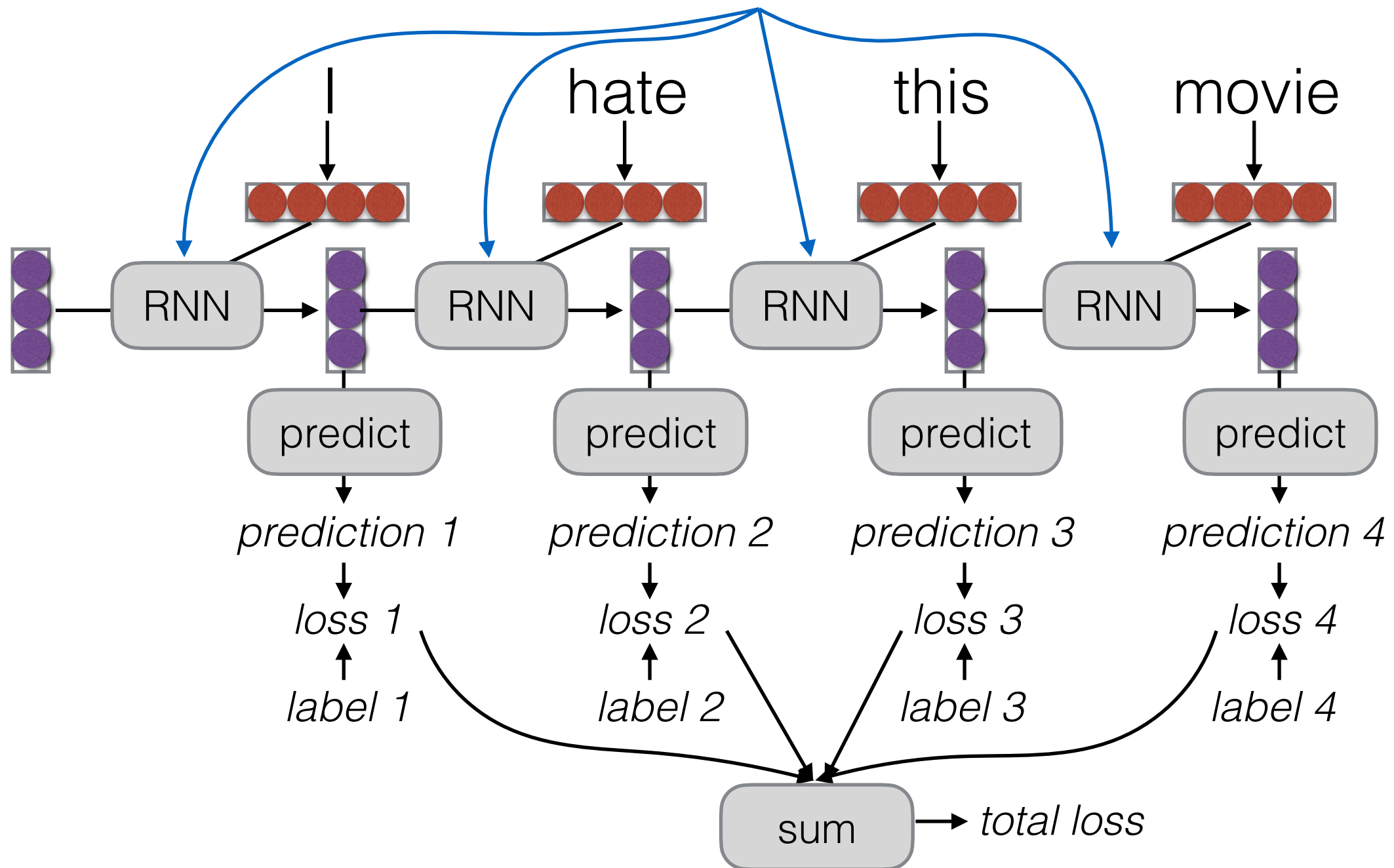
- The unrolled graph is a well-formed (DAG) computation graph—we can run backprop



- Parameters are tied across time, derivatives are aggregated across all time steps
- This is historically called “backpropagation through time” (BPTT)

# Parameter Tying

Parameters are shared! Derivatives are accumulated.



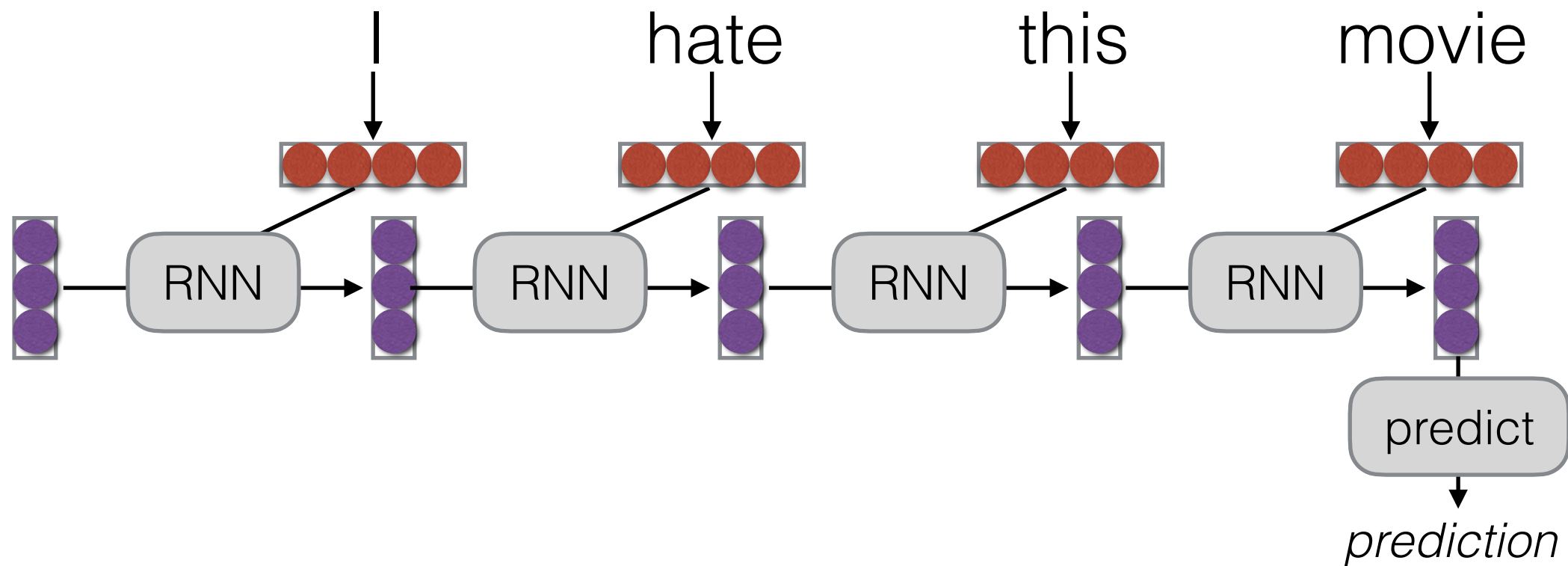
# Applications of RNNs



# What Can RNNs Do?

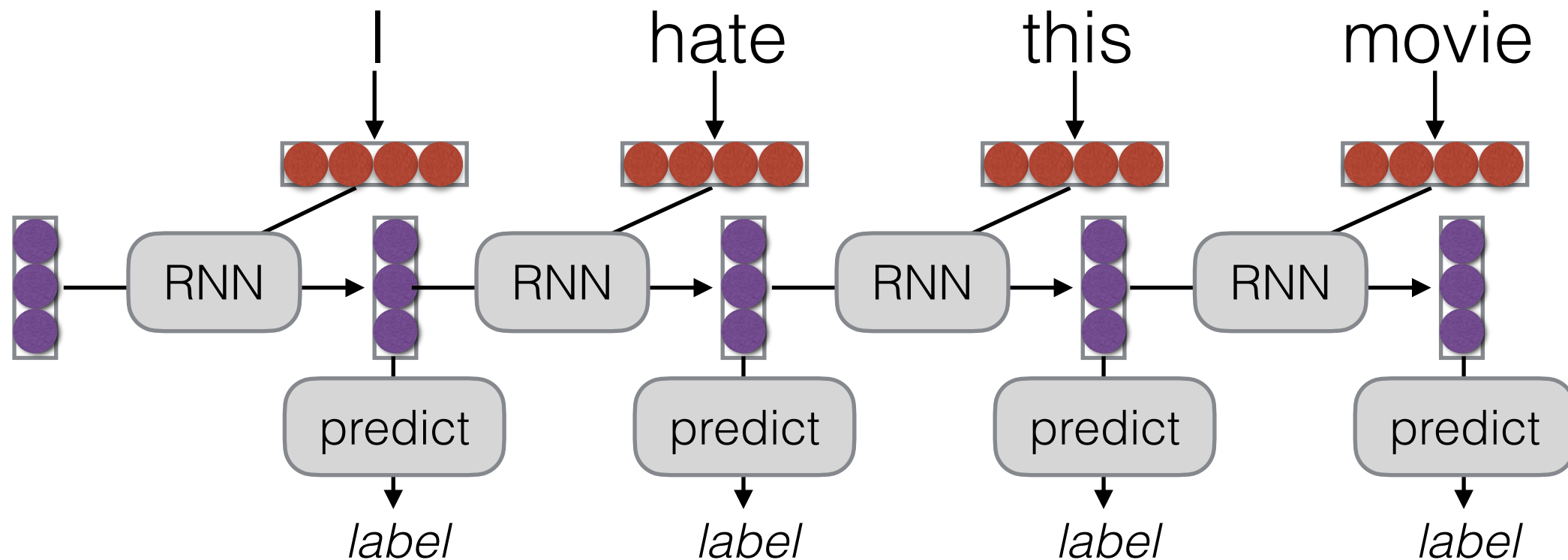
- Represent a sentence
  - Read whole sentence, make a prediction
- Represent a context within a sentence
  - Read context up until that point

# Encoding Sentences



- Binary or multi-class prediction
- Sentence representation for retrieval, sentence comparison, etc.

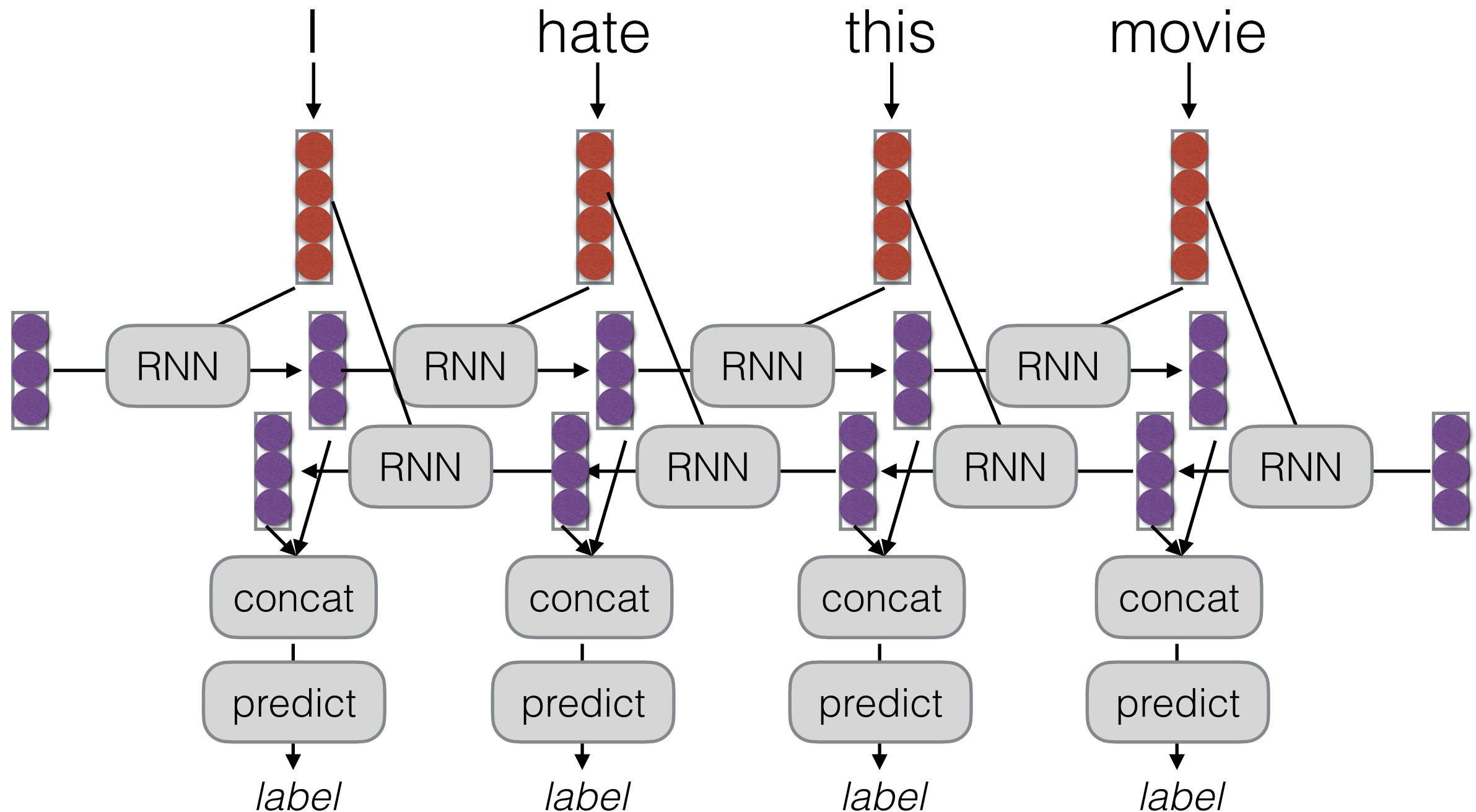
# Representing Contexts



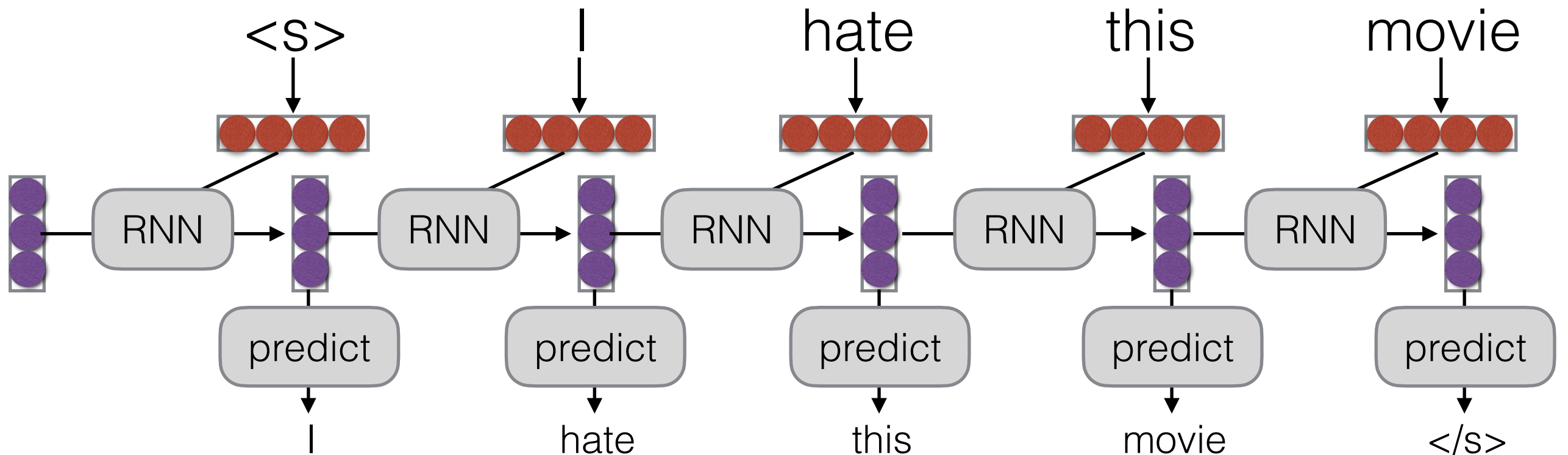
- Sequence labeling
- Calculating representations for parsing, etc.

# Bi-RNNs

- A simple extension, run the RNN in both directions



# e.g. Language Modeling



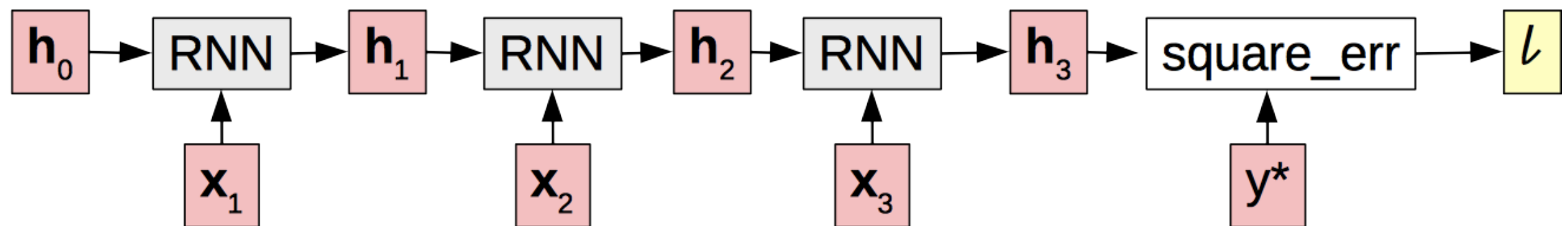
- Language modeling is like a tagging task, where each tag is the next word!
- Note: this is an autoregressive model

# Vanishing Gradients

# Vanishing Gradient

- Gradients decrease as they get pushed back

$$\frac{dl}{d_{h_0}} = \text{tiny} \quad \frac{dl}{d_{h_1}} = \text{small} \quad \frac{dl}{d_{h_2}} = \text{med.} \quad \frac{dl}{d_{h_3}} = \text{large}$$



- Why? “Squashed” by non-linearities or small weights in matrices.

# A Solution:

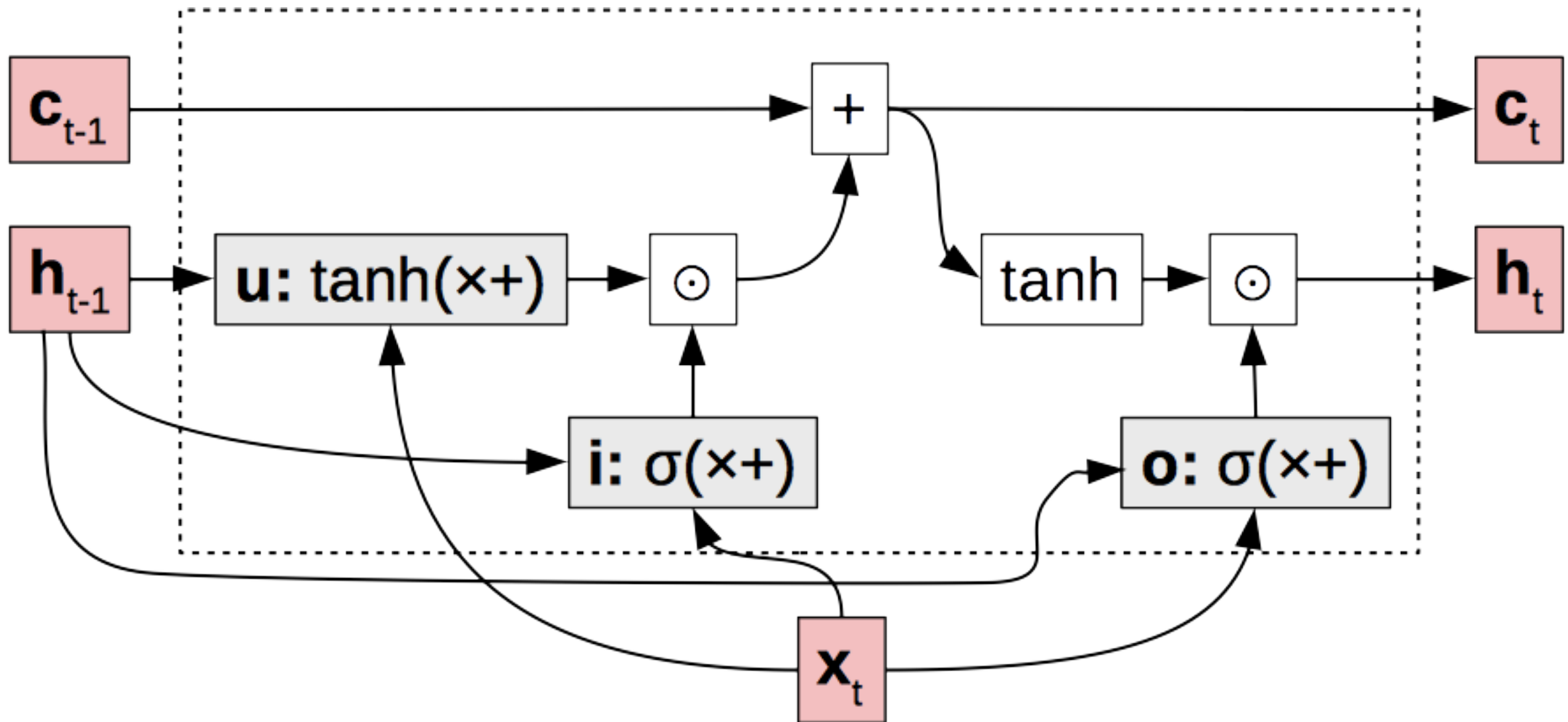
## Long Short-term Memory

(Hochreiter and Schmidhuber 1997)

- **Basic idea:** make additive connections between time steps
- Addition does not modify the gradient, no vanishing
- Gates to control the information flow



# LSTM Structure



update  $u$ : what value do we try to add to the memory cell?  
input  $i$ : how much of the update do we allow to go through?  
output  $o$ : how much of the cell do we reflect in the next state?

# Understanding RNNs

# What can LSTMs Learn? (1)

(Karpathy et al. 2015)

- Additive connections make single nodes surprisingly interpretable

Cell sensitive to position in line:

```
The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.
```

Cell that turns on inside quotes:

```
"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.
```

```
Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."
```

Cell that robustly activates inside if statements:

```
static int __dequeue_signal(struct sigpending *pending, sigset_t *mask, siginfo_t *info)
{
    int sig = next_signal(pending, mask);
    if (sig) {
        if (current->notifier) {
            if (sigismember(current->notifier_mask, sig)) {
                if (!(current->notifier)(current->notifier_data)) {
                    clear_thread_flag(TIF_SIGPENDING);
                    return 0;
                }
            }
        }
        collect_signal(sig, pending, info);
    }
    return sig;
}
```

A large portion of cells are not easily interpretable. Here is a typical example:

```
/* Unpack a filter field's string representation from user-space
 * buffer. */
char *audit_unpack_string(void **bufp, size_t *remain, size_t len)
{
    char *str;
    if (!*bufp || (len == 0) || (len > *remain))
        return ERR_PTR(-EINVAL);
    /* Of the currently implemented string fields, PATH_MAX
     * defines the longest valid length.
     */
}
```

Cell that turns on inside comments and quotes:

```
/* duplicate LSM field information. The lsm_rule is opaque, so
 * re-initialized. */
static inline int audit_dupe_lsm_field(struct audit_field *df,
                                     struct audit_field *sf)
{
    int ret = 0;
    char *lsm_str;
    /* our own copy of lsm_str */
    lsm_str = kstrdup(sf->lsm_str, GFP_KERNEL);
    if (unlikely(!lsm_str))
        return -ENOMEM;
    df->lsm_str = lsm_str;
    /* our own (refreshed) copy of lsm_rule */
    ret = security_audit_rule_init(df->type, df->op, df->lsm_str,
                                   (void **)&df->lsm_rule);
    /* Keep currently invalid fields around in case they
     * become valid after a policy reload. */
    if (ret == -EINVAL) {
        pr_warn("audit rule for LSM '%s' is invalid\n",
              df->lsm_str);
        ret = 0;
    }
    return ret;
}
```

Cell that is sensitive to the depth of an expression:

```
#ifdef CONFIG_AUDITSYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```

Cell that might be helpful in predicting a new line. Note that it only turns on for some "):

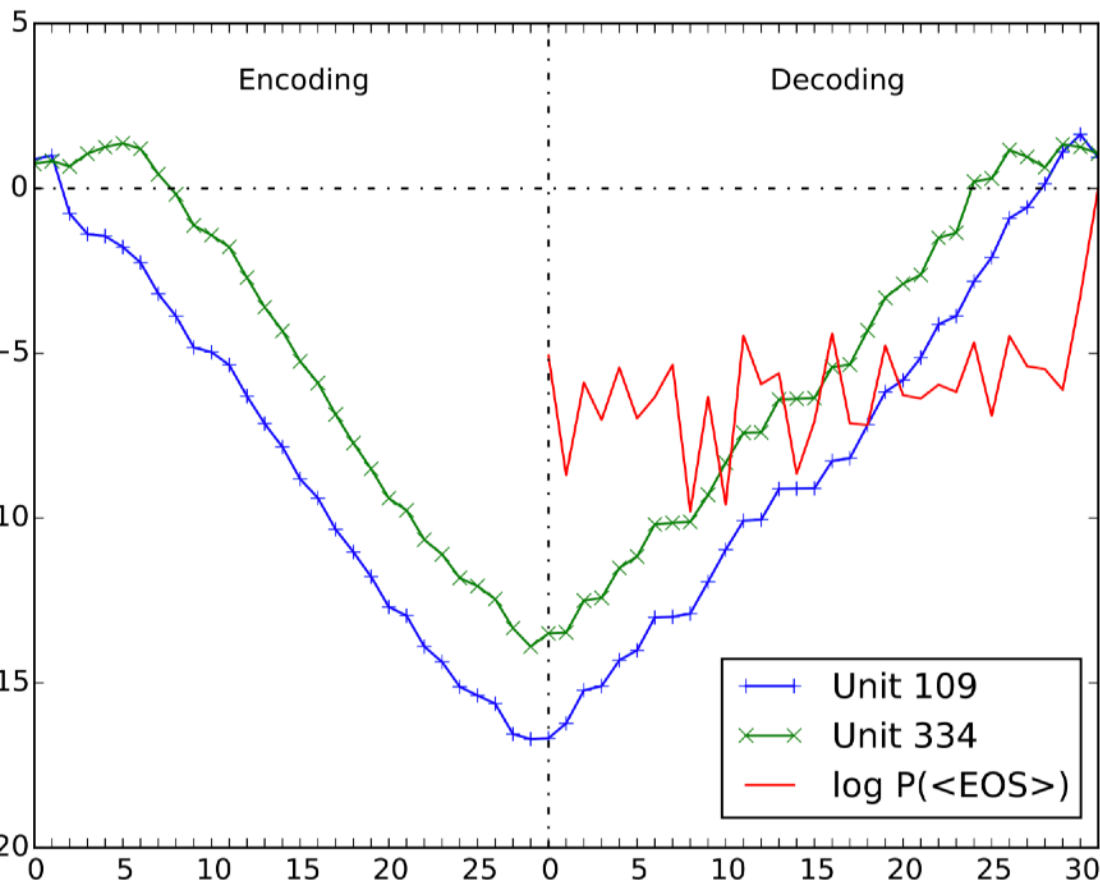
```
char *audit_unpack_string(void **bufp, size_t *remain, si
{
    char *str;
    if (!*bufp || (len == 0) || (len > *remain))
        return ERR_PTR(-EINVAL);
    /* of the currently implemented string fields, PATH_MAX
     * defines the longest valid length.
     */
    if (len > PATH_MAX)
        return ERR_PTR(-ENAMETOOLONG);
    str = kmalloc(len + 1, GFP_KERNEL);
    if (unlikely(!str))
        return ERR_PTR(-ENOMEM);
    memcpy(str, *bufp, len);
    str[len] = 0;
    *bufp += len;
    *remain -= len;
    return str;
}
```

# What can LSTMs Learn? (2)

(Shi et al. 2016, Radford et al. 2017)

## Count length of sentence

## Sentiment



25 August 2003 League of Extraordinary Gentlemen: Sean Connery is one of the all time greats and I have been a fan of his since the 1950's. I went to this movie because Sean Connery was the main actor. I had not read reviews or had any prior knowledge of the movie. The movie surprised me quite a bit. The scenery and sights were spectacular, but the plot was unreal to the point of being ridiculous. In my mind this was not one of his better movies it could be the worst. Why he chose to be in this movie is a mystery. For me, going to this movie was a waste of my time. I will continue to go to his movies and add his movies to my video collection. But I can't see wasting money to put this movie in my collection

I found this to be a charming adaptation, very lively and full of fun. With the exception of a couple of major errors, the cast is wonderful. I have to echo some of the earlier comments -- Chynna Phillips is horribly miscast as a teenager. At 27, she's just too old (and, yes, it DOES show), and lacks the singing "chops" for Broadway-style music. Vanessa Williams is a decent-enough singer and, for a non-dancer, she's adequate. However, she is NOT Latina, and her character definitely is. She's also very STRIDENT throughout, which gets tiresome. The girls of Sweet Apple's Conrad Birdie fan club really sparkle -- with special kudos to Brigitta Dau and Chiara Zanni. I also enjoyed Tyne Daly's performance, though I'm not generally a fan of her work. Finally, the dancing Shriners are a riot, especially the dorky three in the bar. The movie is suitable for the whole family, and I highly recommend it.

Notably, difficult for GRUs!



# RNNs as Turing Machines

(Siegelmann and Sontag 1992)

- Real-valued recurrent nets can calculate arbitrary functions

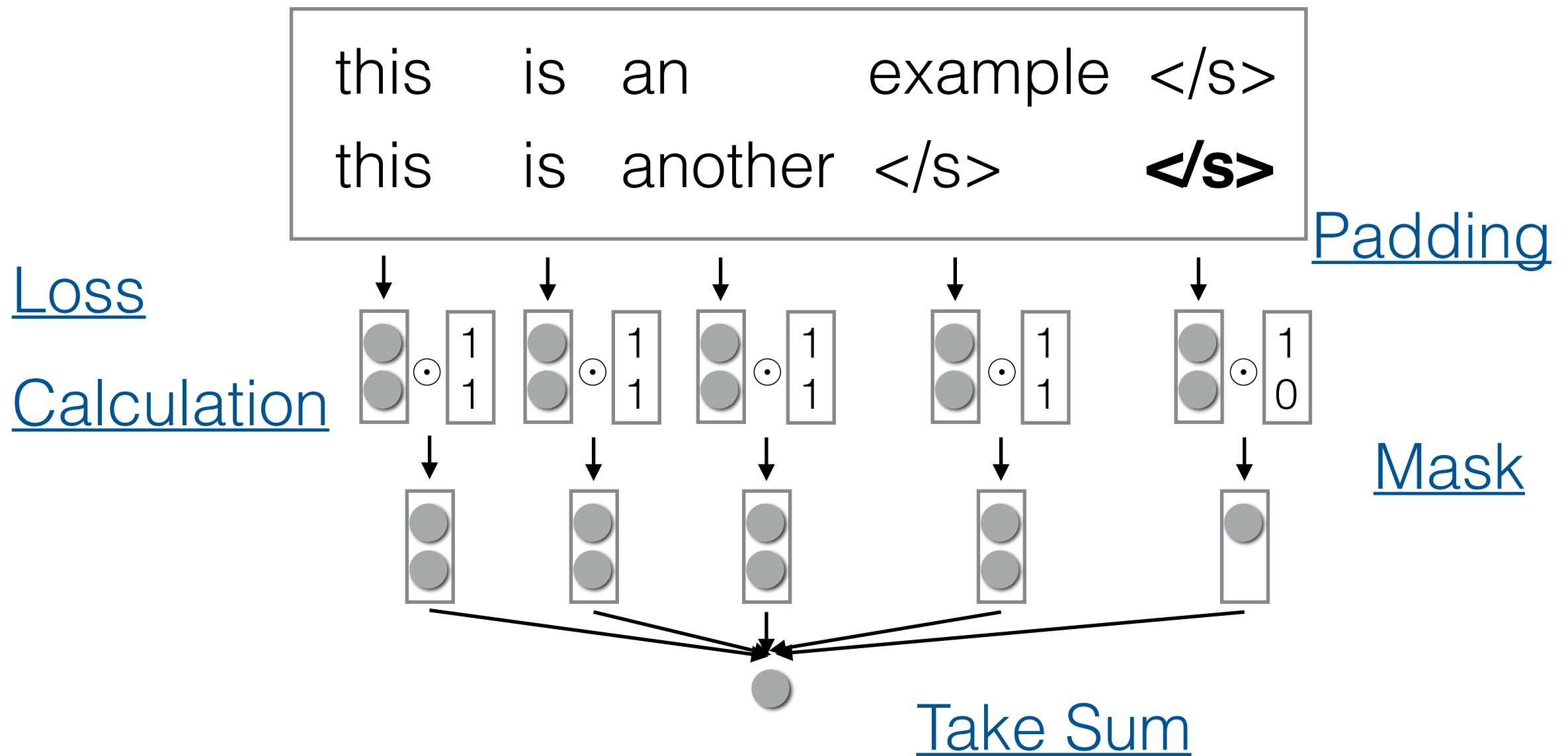
<b>Weights</b>	<b>Recognition Power</b>
integer	regular
rational	recursive
real	arbitrary

# Efficiency Tricks

# Handling Mini-batching

- Mini-batching makes things much faster!
- But mini-batching in RNNs is harder than in feed-forward networks
  - Each word depends on the previous word
  - Sequences are of various length

# Mini-batching Method



(Or use DyNet automatic mini-batching, much easier but a bit slower)



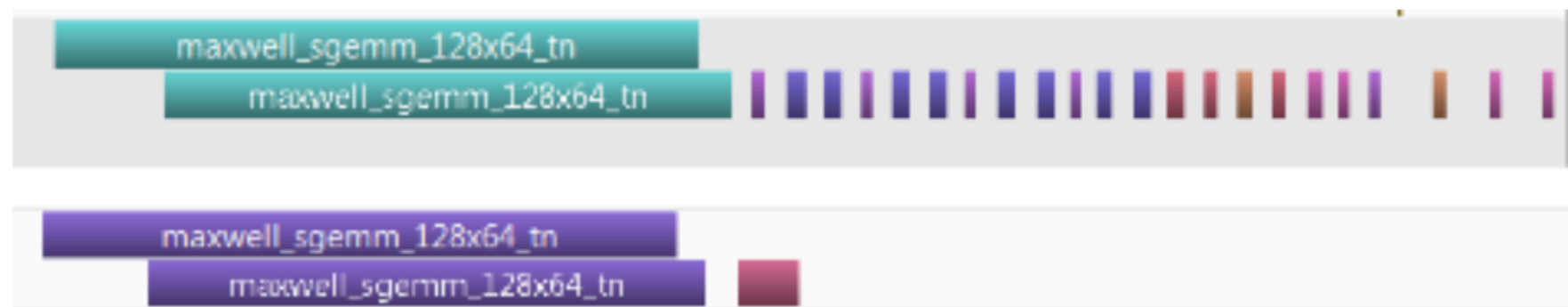
# Bucketing/Sorting

- If we use sentences of different lengths, too much padding and sorting can **result in decreased performance**
- To remedy this: **sort sentences** so similarly-lengthed sentences are in the same batch

# Optimized Implementations of LSTMs

(Appleyard 2015)

- In simple implementation, still need one GPU call for each time step
- For some RNN variants (e.g. LSTM) efficient full-sequence computation supported by CuDNN



- **Basic process:** combine inputs into tensor, single GPU call combine inputs into tensor, single GPU call
- **Downside:** significant loss of flexibility

# Handling Long Sequences

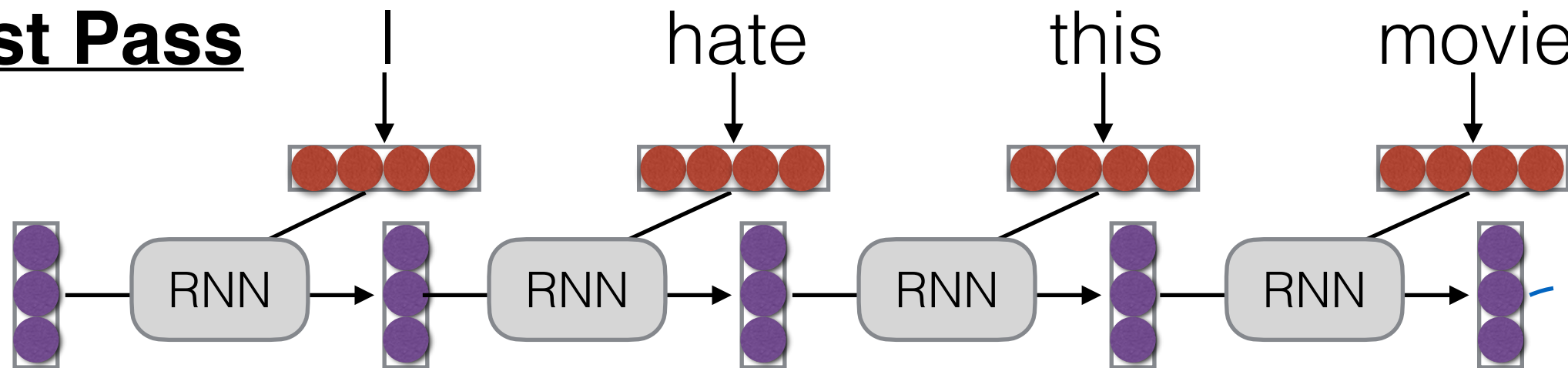
# Handling Long Sequences

- Sometimes we would like to capture long-term dependencies over long sequences
- e.g. words in full documents
- However, this may not fit on (GPU) memory

# Truncated BPTT

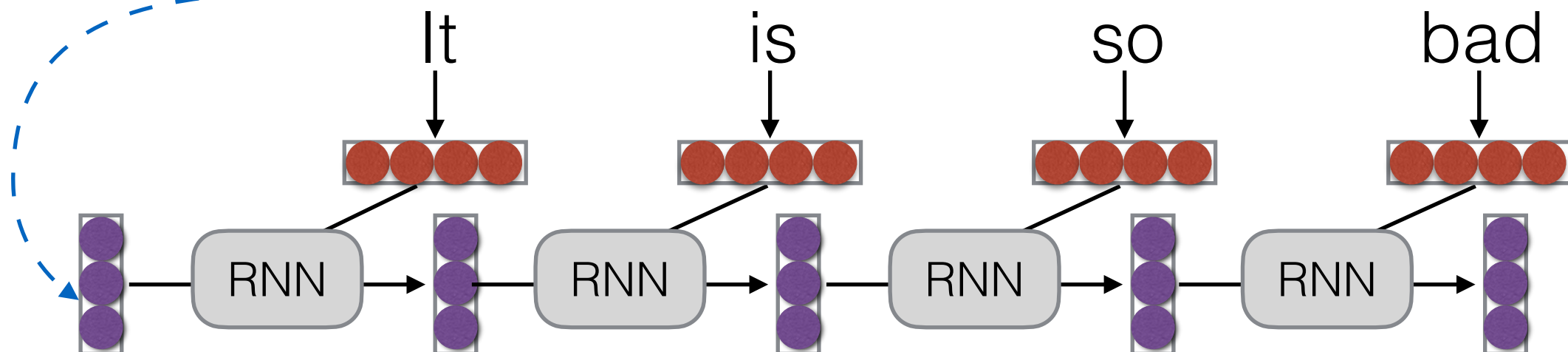
- Backprop over shorter segments, initialize w/ the state from the previous segment

## 1st Pass



## 2nd Pass

state only, no backprop



# RNN Variants

# Gated Recurrent Units

(Cho et al. 2014)

- A simpler version that preserves the additive connections

$$z_t = \sigma_g(W_z x_t + U_z h_{t-1} + b_z)$$

$$r_t = \sigma_g(W_r x_t + U_r h_{t-1} + b_r)$$

$$h_t = \boxed{(1 - z_t)} \circ h_{t-1} + \boxed{z_t} \circ \sigma_h(W_h x_t + U_h (r_t \circ h_{t-1}) + b_h)$$

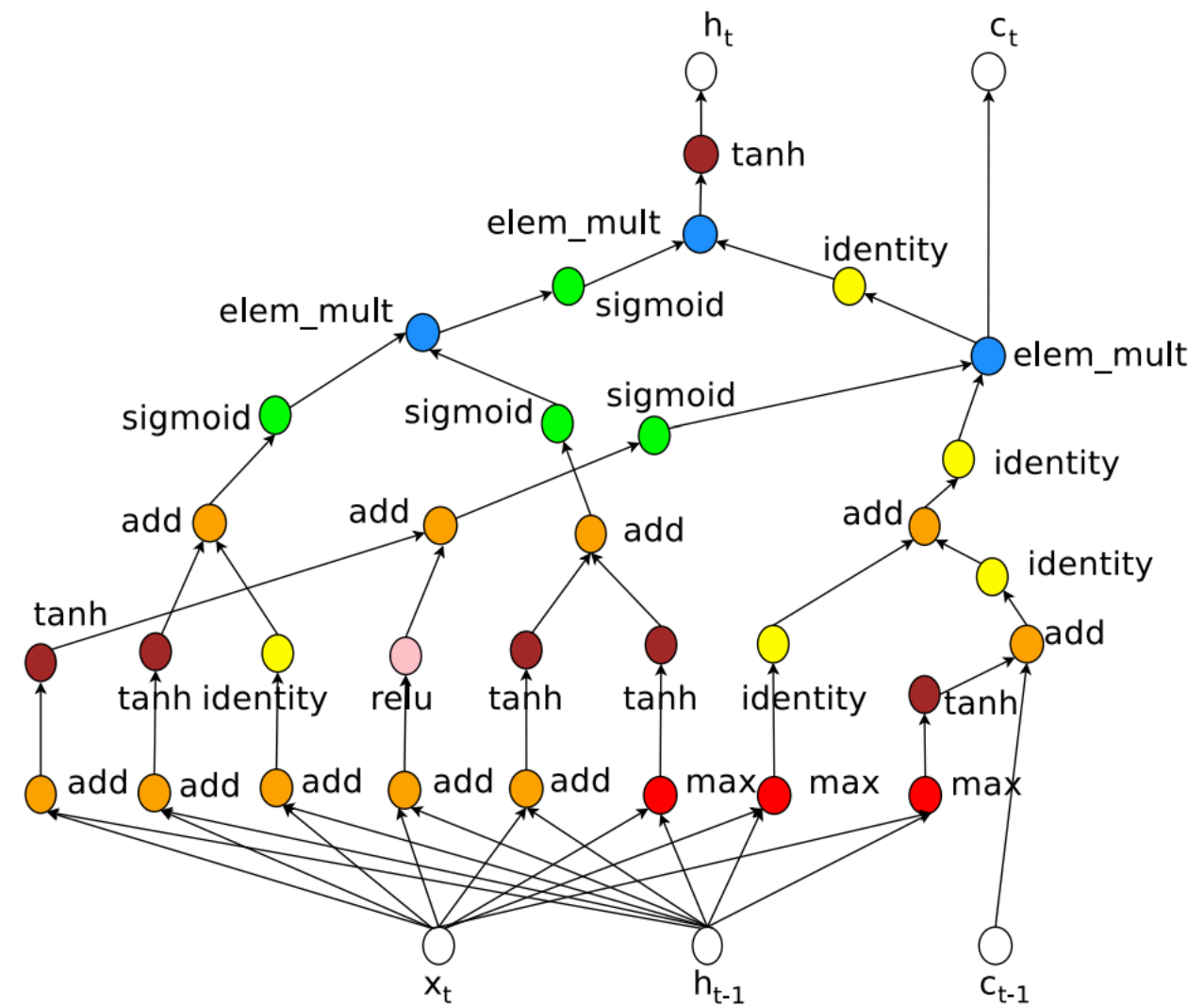
Additive or Non-linear

- **Note:** GRUs cannot do things like simply count

# Extensive Architecture Search for LSTMs

(Greffen et al. 2015, Zoph and Le 2017)

- Many different types of architectures tested for LSTMs (manually or automatically)
- **Conclusion:** basic LSTM quite good, other variants (e.g. coupled input/forget gates) reasonable

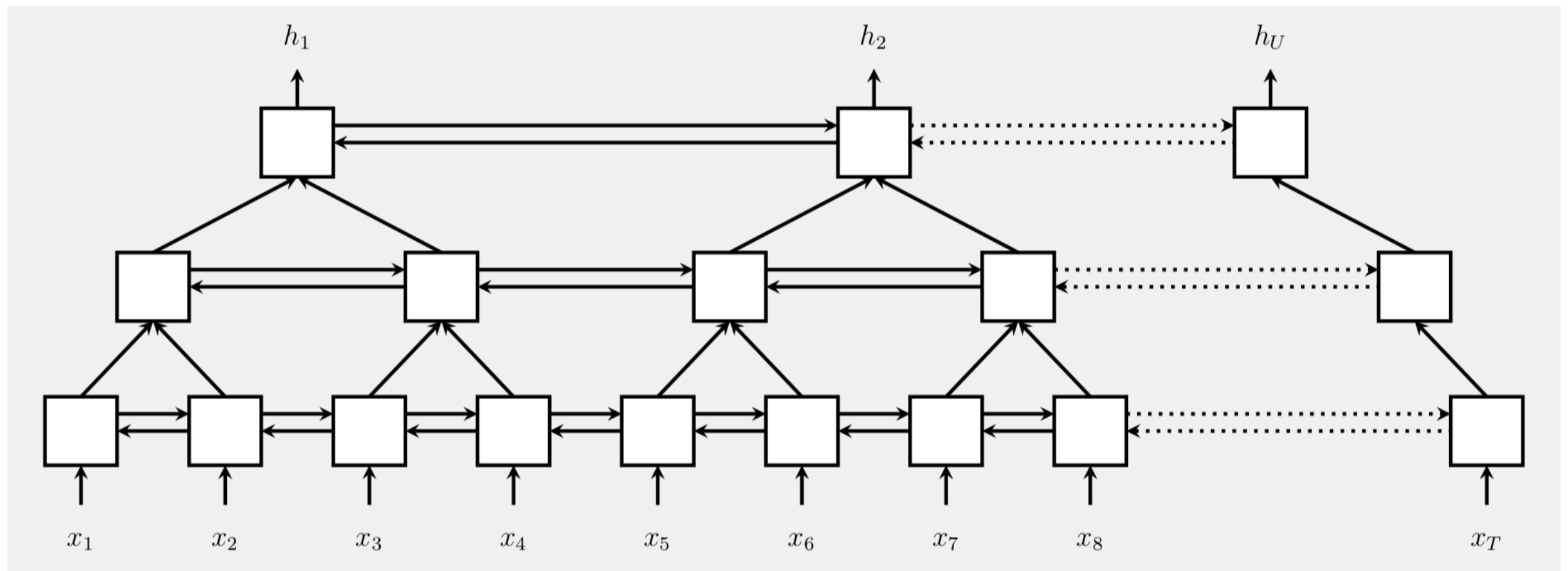




# Multi-scale/Pyramidal RNN

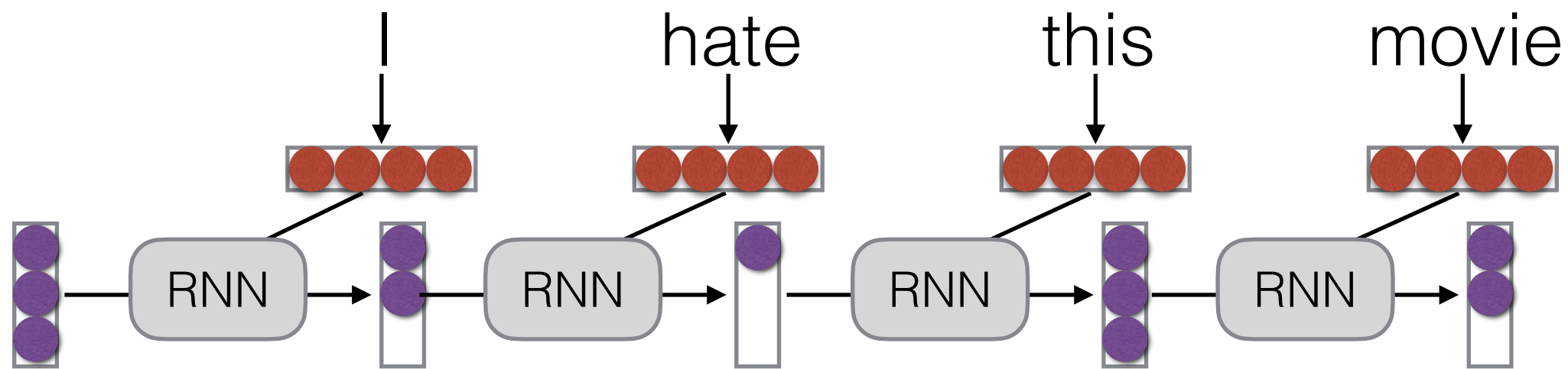
(e.g. Chan et al. 2015)

- Have different RNNs of different scales



# Soft Hierarchical Structure

- e.g. "Ordered Neurons" Shen et al. 2019



$$\tilde{f}_t = \text{cumax}(W_{\tilde{f}}x_t + U_{\tilde{f}}h_{t-1} + b_{\tilde{f}})$$

$$\tilde{i}_t = 1 - \text{cumax}(W_{\tilde{i}}x_t + U_{\tilde{i}}h_{t-1} + b_{\tilde{i}})$$

- Makes it possible to discover hierarchies

Types of Prediction  
Recurrent Neural Nets  
RNN Applications  
Vanishing Gradients

Understanding RNNs  
Efficiency Tricks  
Handling Long Sequences  
RNN Variants

Questions?