# 자연어 처리 DAY 2

Jaegul Choo

Associate Professor, Graduate School of AI, KAIST
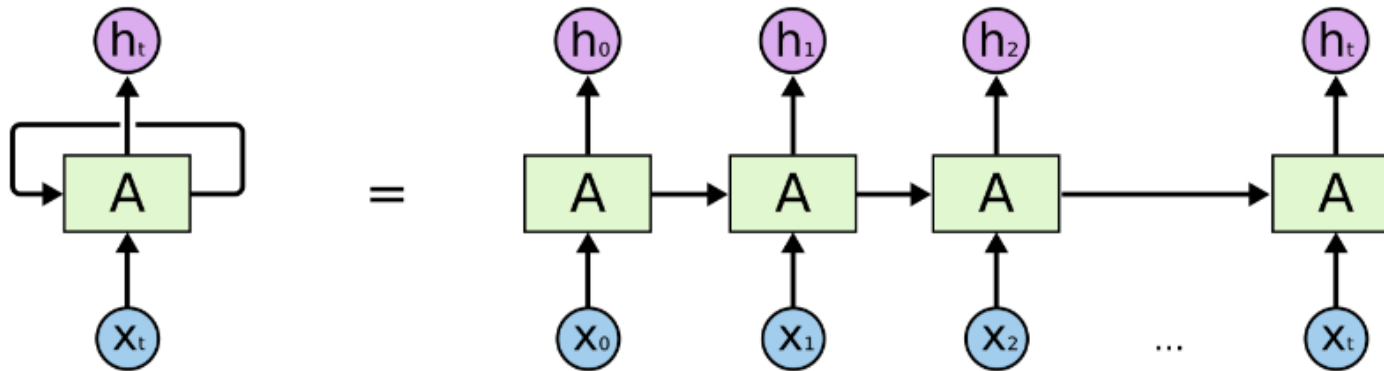
# 1.
# Basics of Recurrent Neural Networks (RNNs)

Basic problem settings
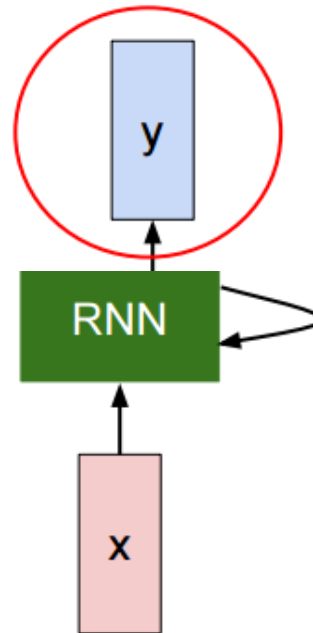
Model architecture and how it works

# Types of RNNs

- ## Basic structure



An unrolled recurrent neural network.

© NAVER Connect Foundation

3

- ## Inputs and outputs of RNNs (rolled version)

  - We usually want to predict a vector at some time steps

- ## How to calculate the hidden state of RNNs

  - We can process a sequence of vectors by applying a recurrence formula at every time step

  - $h_{t-1}$: old hidden-state vector

  - $x_t$: input vector at some time step

  - $h_t$: new hidden-state vector

  - $f_W$: RNN function with parameters $W$
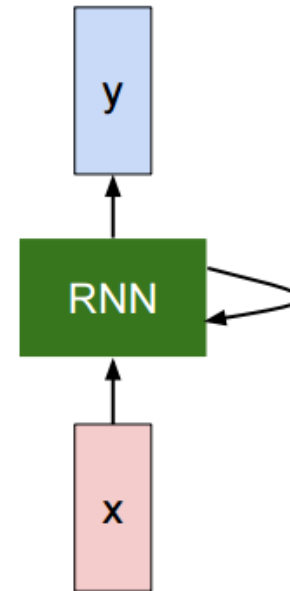
  - $y_t$: output vector at time step $t$

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

http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture10.pdf
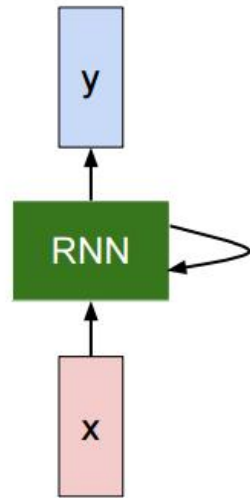
boostcamp AI Tech

- ## How to calculate the hidden state of RNNs

  - Notice: The same function and the same set of parameters are used at every time step

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

y

RNN

x

- ## How to calculate the hidden state of RNNs

  - The state consists of a single "hidden" vector <u>h</u>



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

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

$$y_t = W_{hy}h_t$$
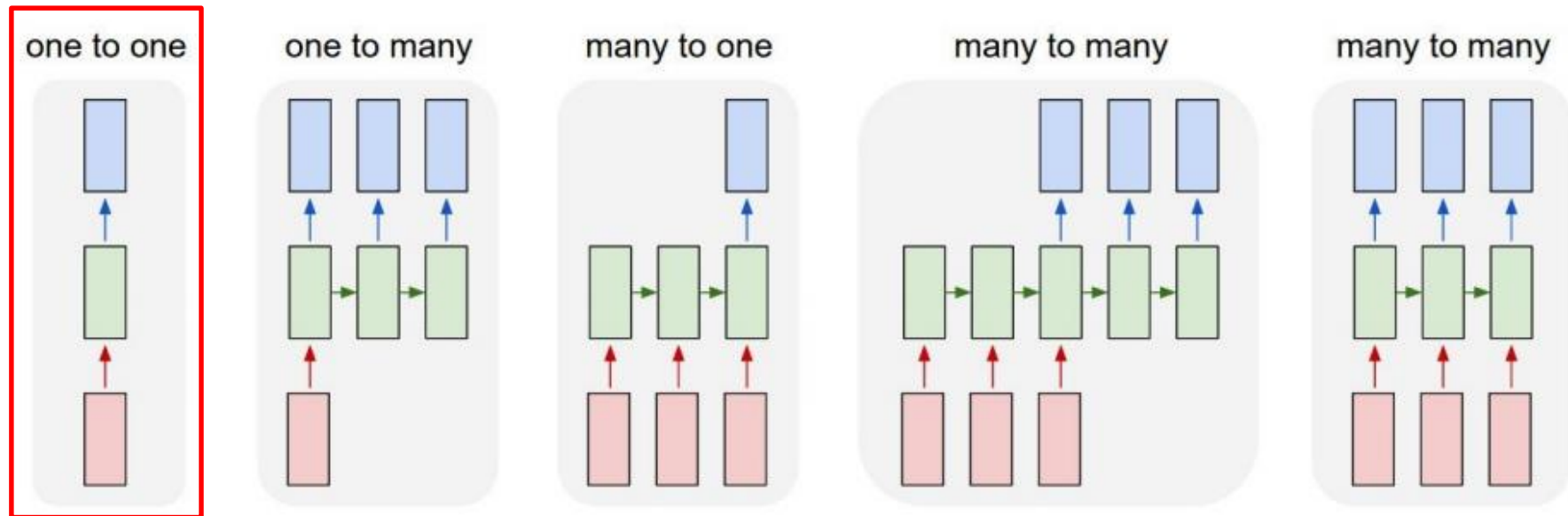
# End of Document
# Thank You.
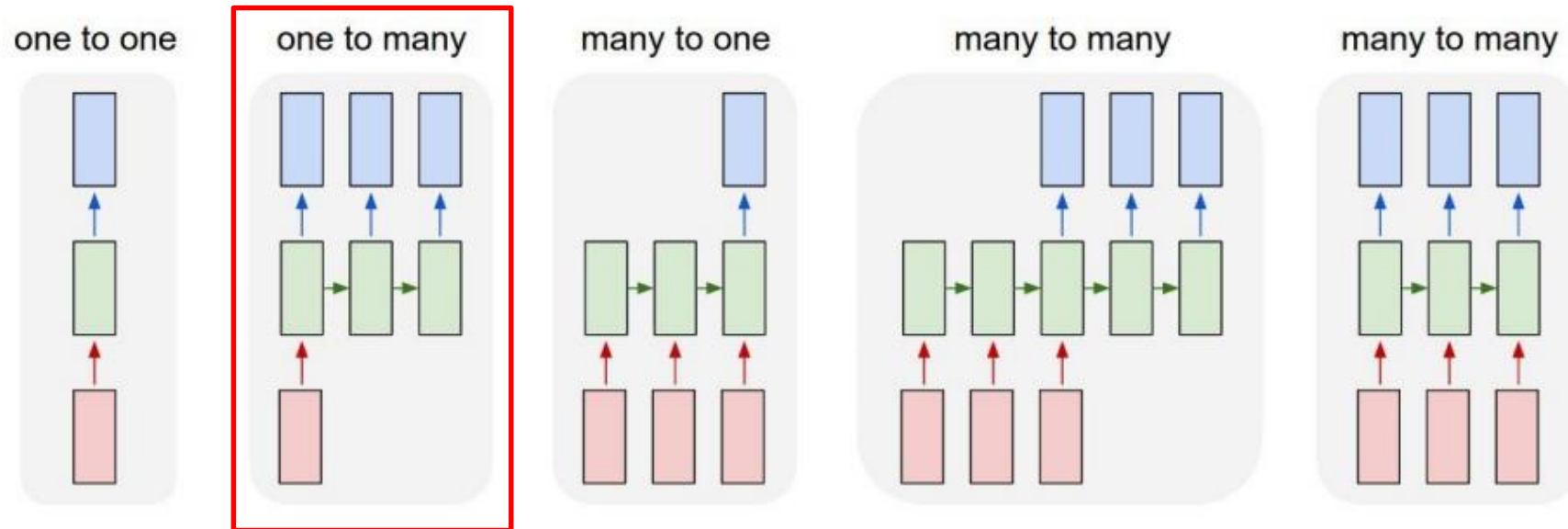
# 2.
# Types of RNNs

one-to-one

one-to-many

many-to-one

many-to-many

# Types of RNNs

- ## One-to-one

  - Standard Neural Networks

# Types of RNNs

- ## One-to-many

  - Image Captioning



one to one    one to many    many to one    many to many    many to many

# Types of RNNs

- ## Many-to-one

  - Sentiment Classification

# Types of RNNs

- ## Sequence-to-sequence

  - Machine Translation



http://karpathy.github.io/2015/05/21/rnn-effectiveness/

# Types of RNNs

- ## Sequence-to-sequence

  - Video classification on frame level

© NAVER Connect Foundation

# End of Document
# Thank You.

boostcamp AI Tech

# 3.
# Character-level Language Model

# Character-level Language Model

- ## Example of training sequence "hello"

  - Vocabulary: [h, e, l, o]

  - Example training sequence: "hello"



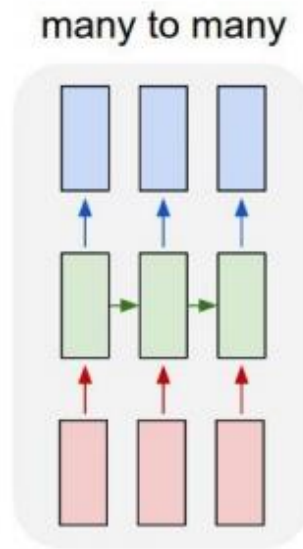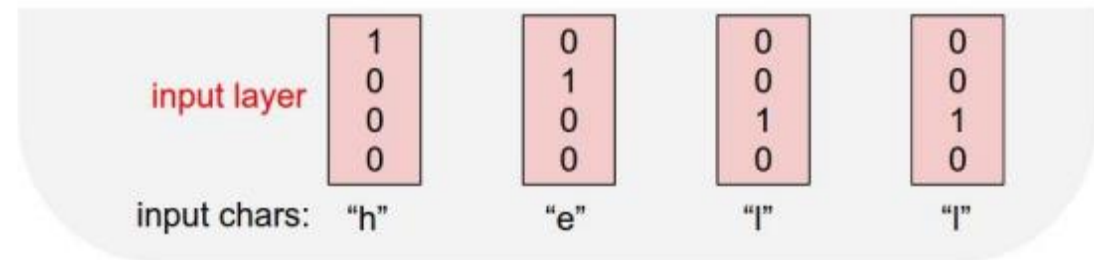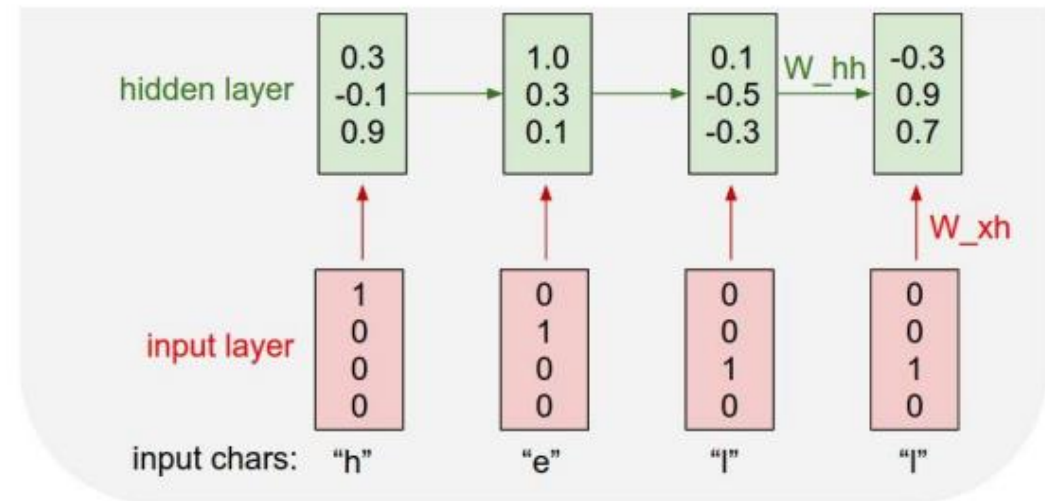many to many

# Character-level Language Model

- ## Example of training sequence "hello"

    - Vocabulary: [h, e, l, o]
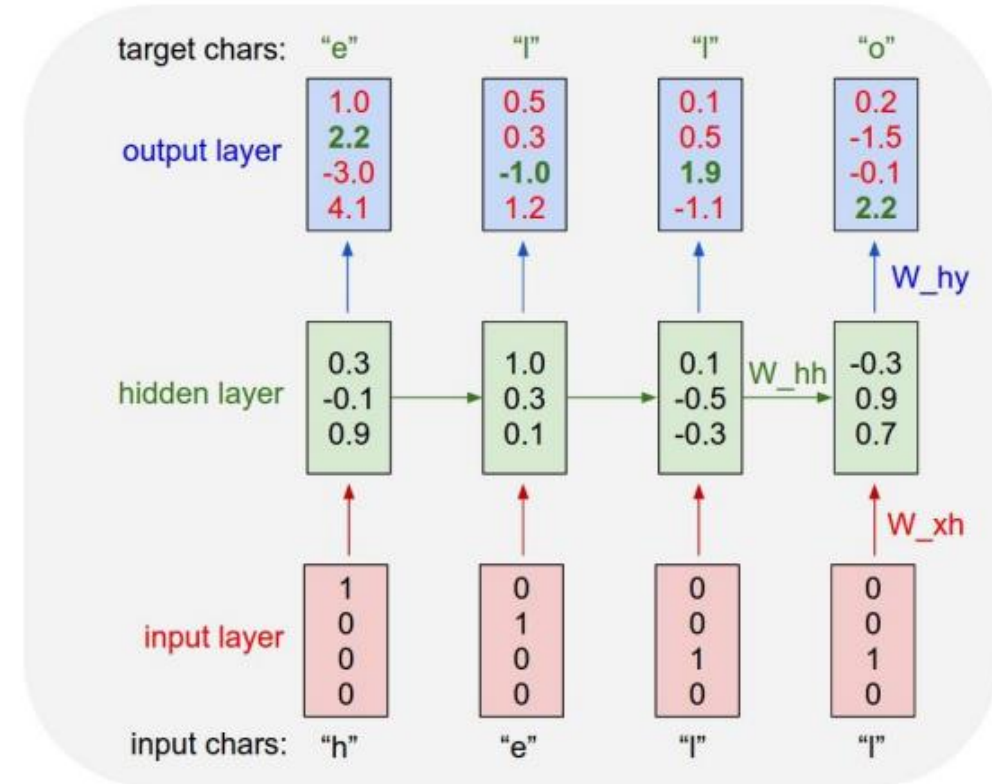
    - Example training sequence: "hello"

- ## Example of training sequence "hello"

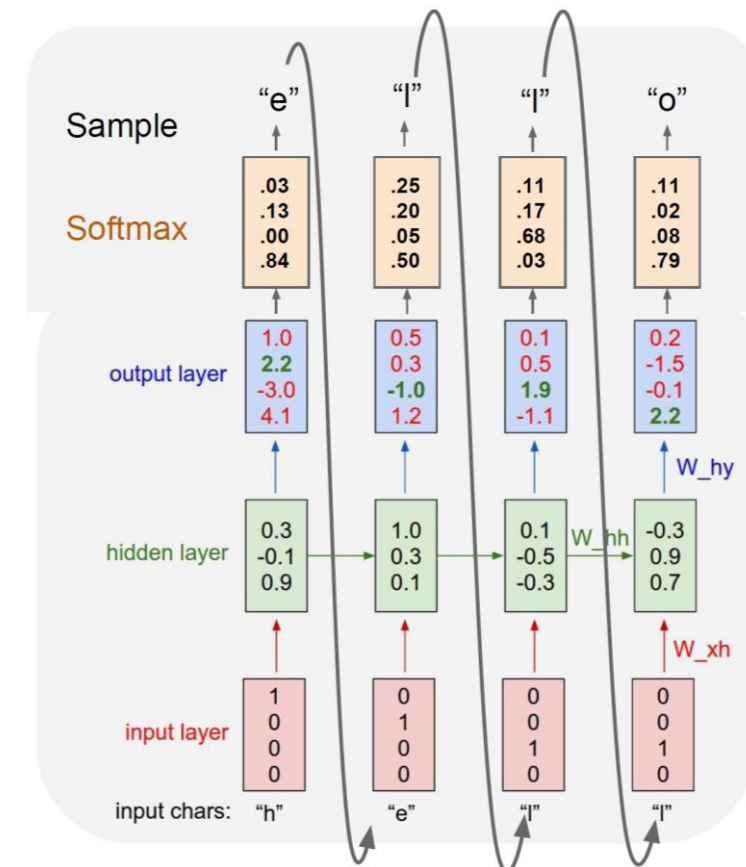  - $h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b)$

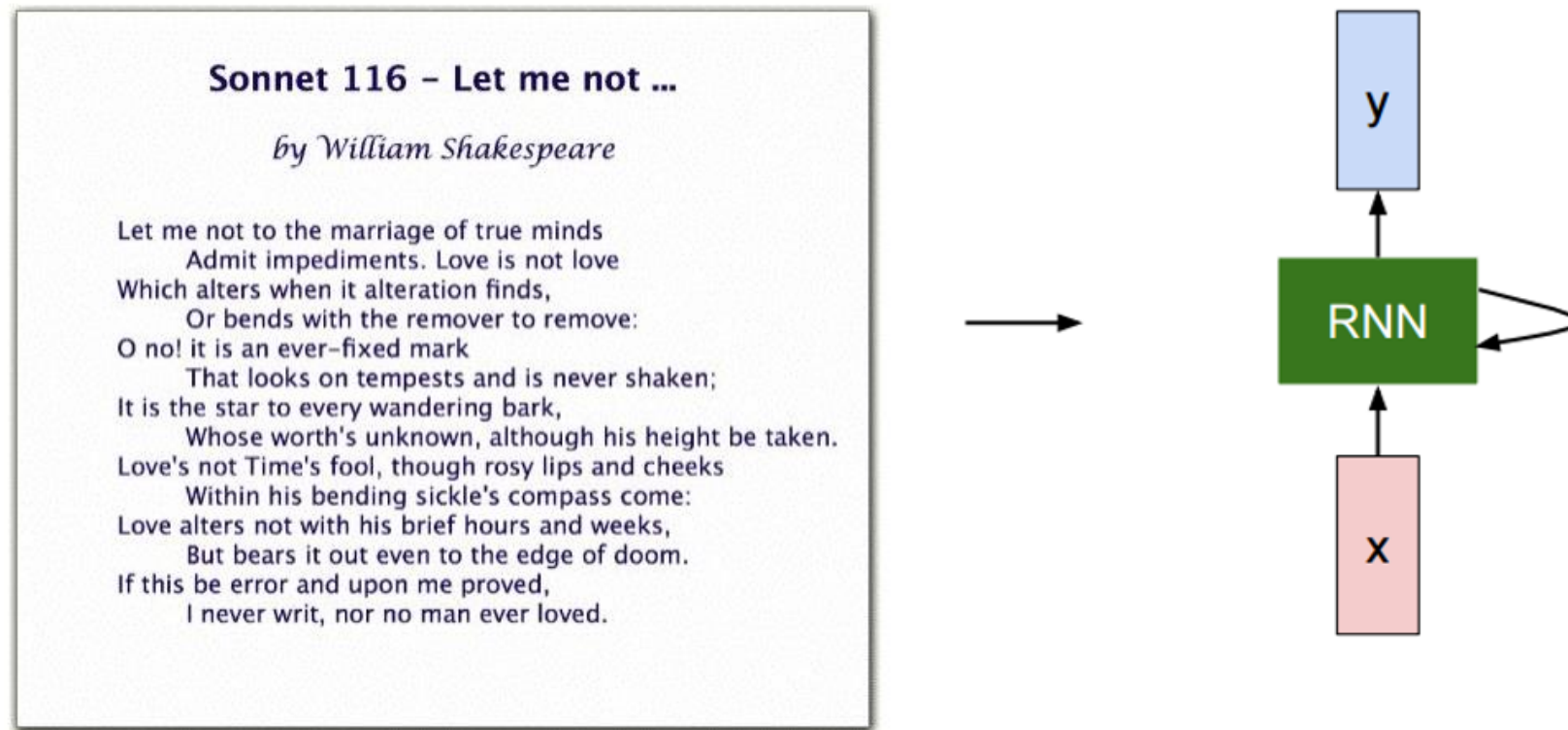- ## Example of training sequence "hello"

  - Logit $= W_{hy} h_t + b$



http://karpathy.github.io/2015/05/21/rnn-effectiveness/

# Character-level Language Model

Character-level Language Model

- ## Example of training sequence "hello"

  - At test-time, sample characters one at a time, feed back to model



http://karpathy.github.io/2015/05/21/rnn-effectiveness/

© NAVER Connect Foundation

21

- **Training a RNN on Shakespeare's plays**

# Character-level Language Model

- **Training process of RNN**

at first:
> tyntd-iafhatawiaoihrdemot  lytdws  e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e
> plia tklrgd t o idoe ns,smtt    h ne etie h,hregtrs nigtike,aoaenns lng

↓ train more

> "Tmont thithey" fomesscerliund
> Keushey. Thom here
> sheulke, anmerenith ol sivh I lalterthend Bleipile shuwy fil on aseterlome
> coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

↓ train more

> Aftair fall unsuch that the hall for Prince Velzonski's that me of
> her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort
> how, and Gogition is so overelical and ofter.

↓ train more

> "Why do what that day," replied Natasha, and wishing to himself the fact the
> princess, Princess Mary was easier, fed in had oftened him.
> Pierre aking his soul came to the packs and drove up his father-in-law women.

http://karpathy.github.io/2015/05/21/rnn-effectiveness/

- ## Results of trained RNN

PANDARUS:
Alas, I think he shall be come approached and the day
When little srain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:
They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

DUKE VINCENTIO:
Well, your wit is in the care of side and that.

Second Lord:
They would be ruled after this chamber, and
my fair nues begun out of the fact, to be conveyed,
Whose noble souls I'll have the heart of the wars.

Clown:
Come, sir, I will make did behold your worship.

VIOLA:
I'll drink it.

VIOLA:
Why, Salisbury must find his flesh and thought
That which I am not aps, not a man and in fire,
To show the reining of the raven and the wars
To grace my hand reproach within, and not a fair are hand,
That Caesar and my goodly father's world;
When I was heaven of presence and our fleets,
We spare with hours, but cut thy council I am great,
Murdered and by thy master's ready there
My power to give thee but so much as hell:
Some service in the noble bondman here,
Would show him to her wine.

KING LEAR:
O, if you were a feeble sight, the courtesy of your law,
Your sight and several breath, will wear the gods
With his heads, and my hands are wonder'd at the deeds,
So drop upon your lordship's head, and your opinion
Shall be against your honour.

http://karpathy.github.io/2015/05/21/rnn-effectiveness/

# Character-level Language Model

- A paper written by RNN

# Character-level Language Model

- ## C code generated by RNN
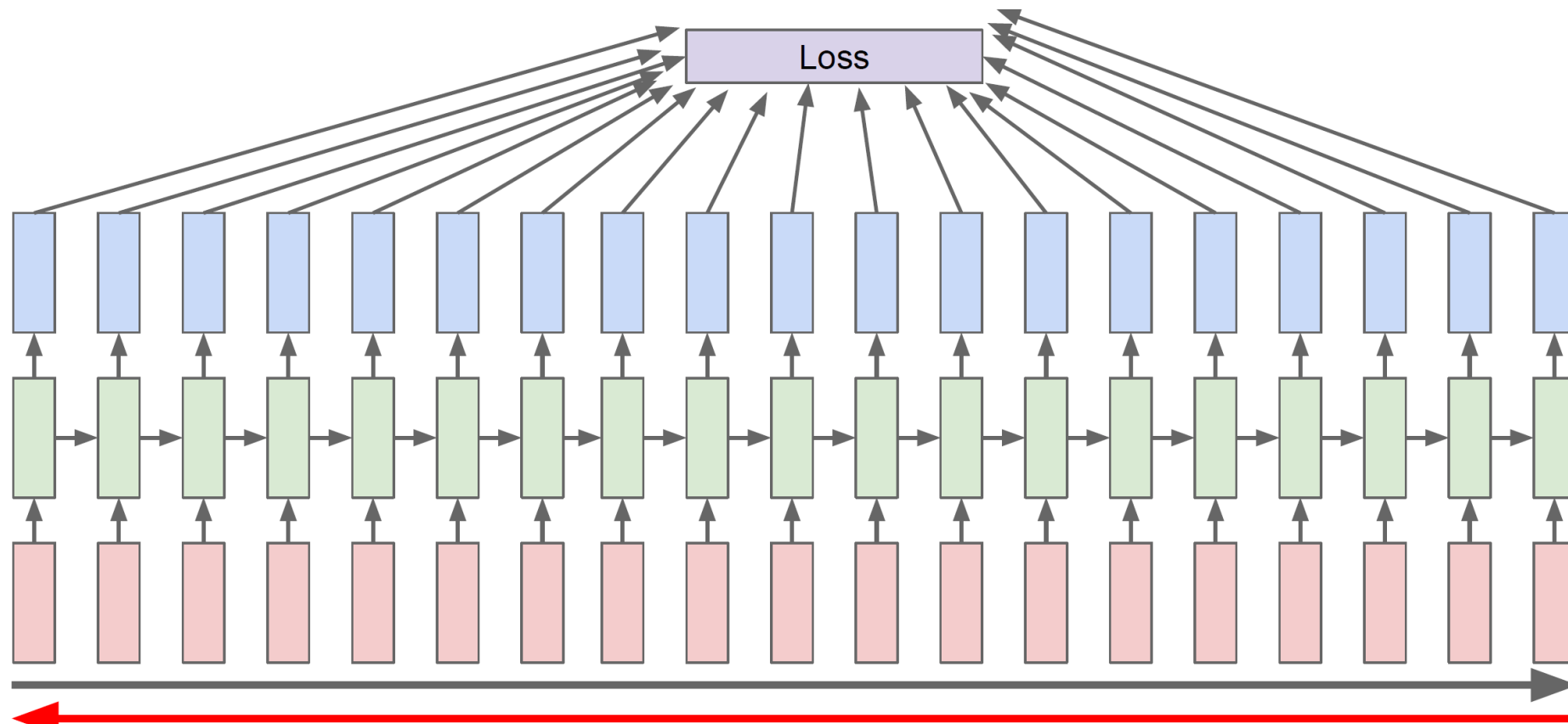
```c
static void do_command(struct seq_file *m, void *v)
{
  int column = 32 << (cmd[2] & 0x80);
  if (state)
    cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
  else
    seq = 1;
  for (i = 0; i < 16; i++) {
    if (k & (1 << 1))
      pipe = (in_use & UMXTHREAD_UNCCA) +
        ((count & 0x00000000ffffffff8) & 0x000000f) << 8;
    if (count == 0)
      sub(pid, ppc_md.kexec_handle, 0x20000000);
    pipe_set_bytes(i, 0);
  }
  /* Free our user pages pointer to place camera if all dash */
  subsystem_info = &of_changes[PAGE_SIZE];
  rek_controls(offset, idx, &soffset);
  /* Now we want to deliberately put it to device */
  control_check_polarity(&context, val, 0);
  for (i = 0; i < COUNTER; i++)
    seq_puts(s, "policy ");
}
```

http://karpathy.github.io/2015/05/21/rnn-effectiveness/
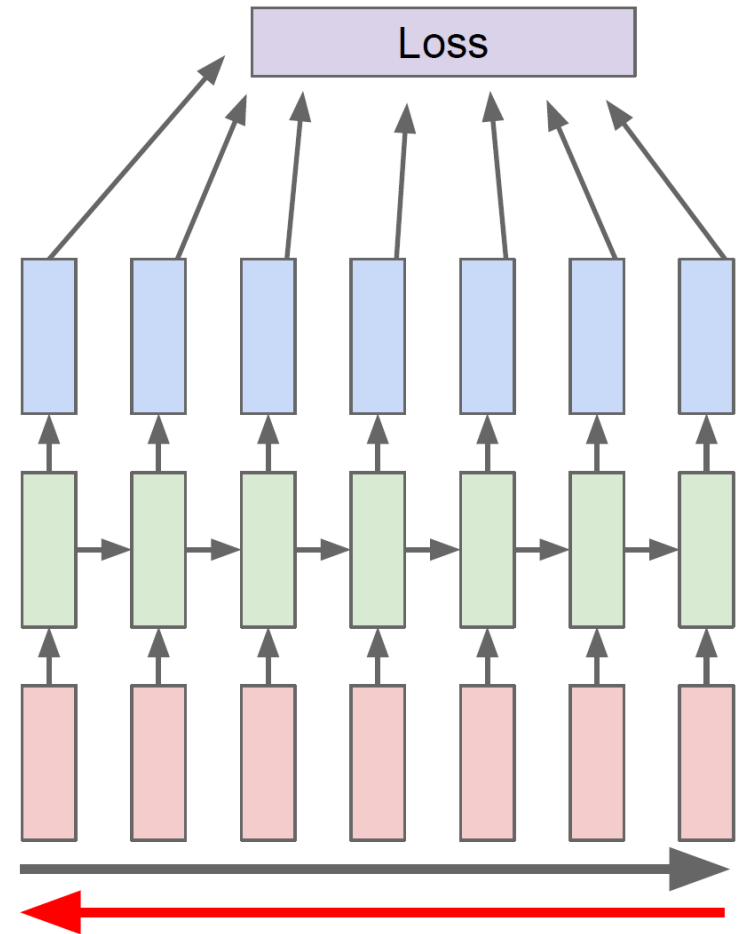
# Backpropagation through time (BPTT)

- Forward through entire sequence to compute loss, then backward through entire sequence to comput gradient
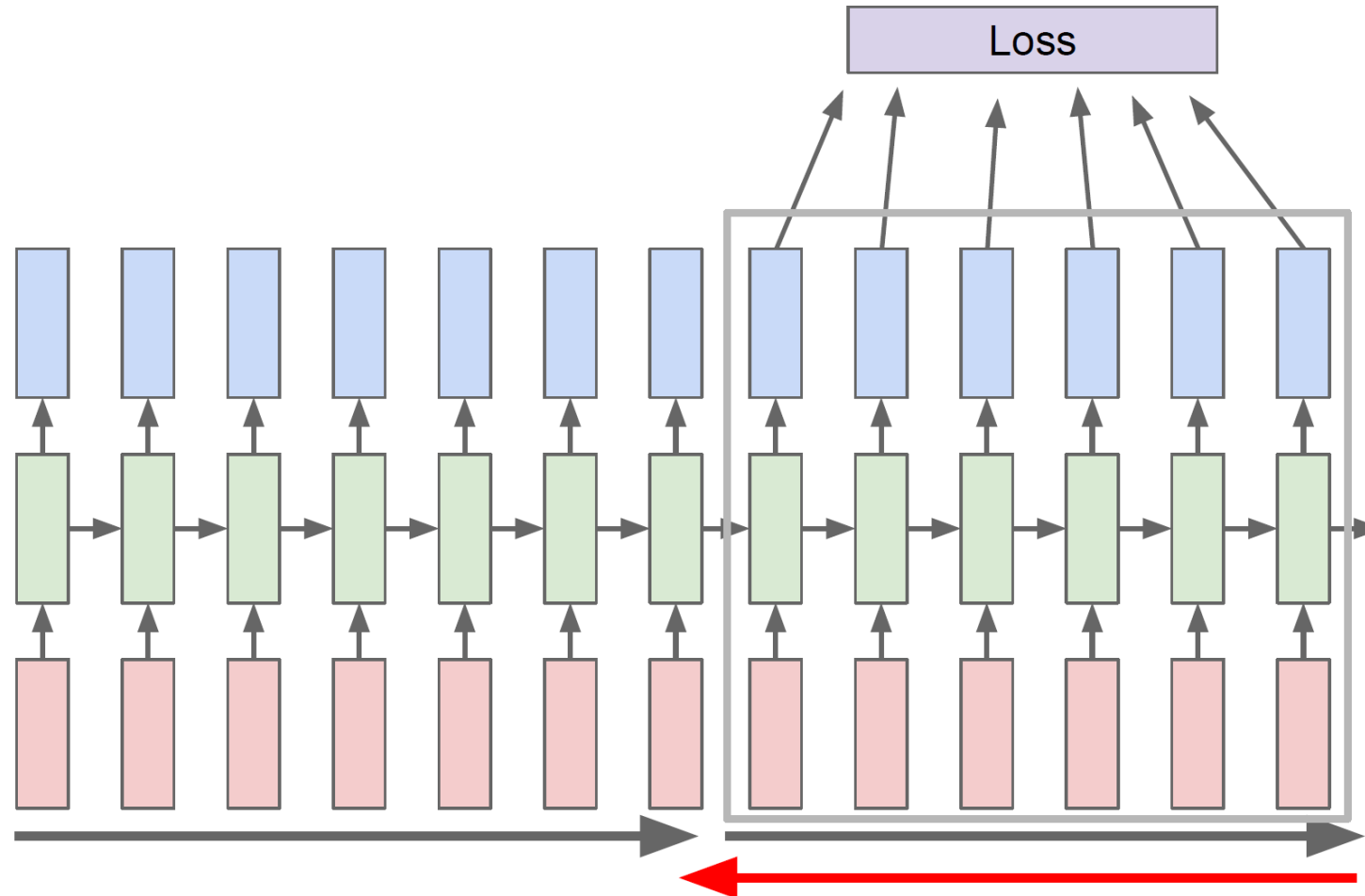
# Backpropagation through time (BPTT)

- Run forward and backward through chunks of the sequence instead of whole sequence

boostcamp AI Tech

© NAVER Connect Foundation

# Backpropagation through time (BPTT)
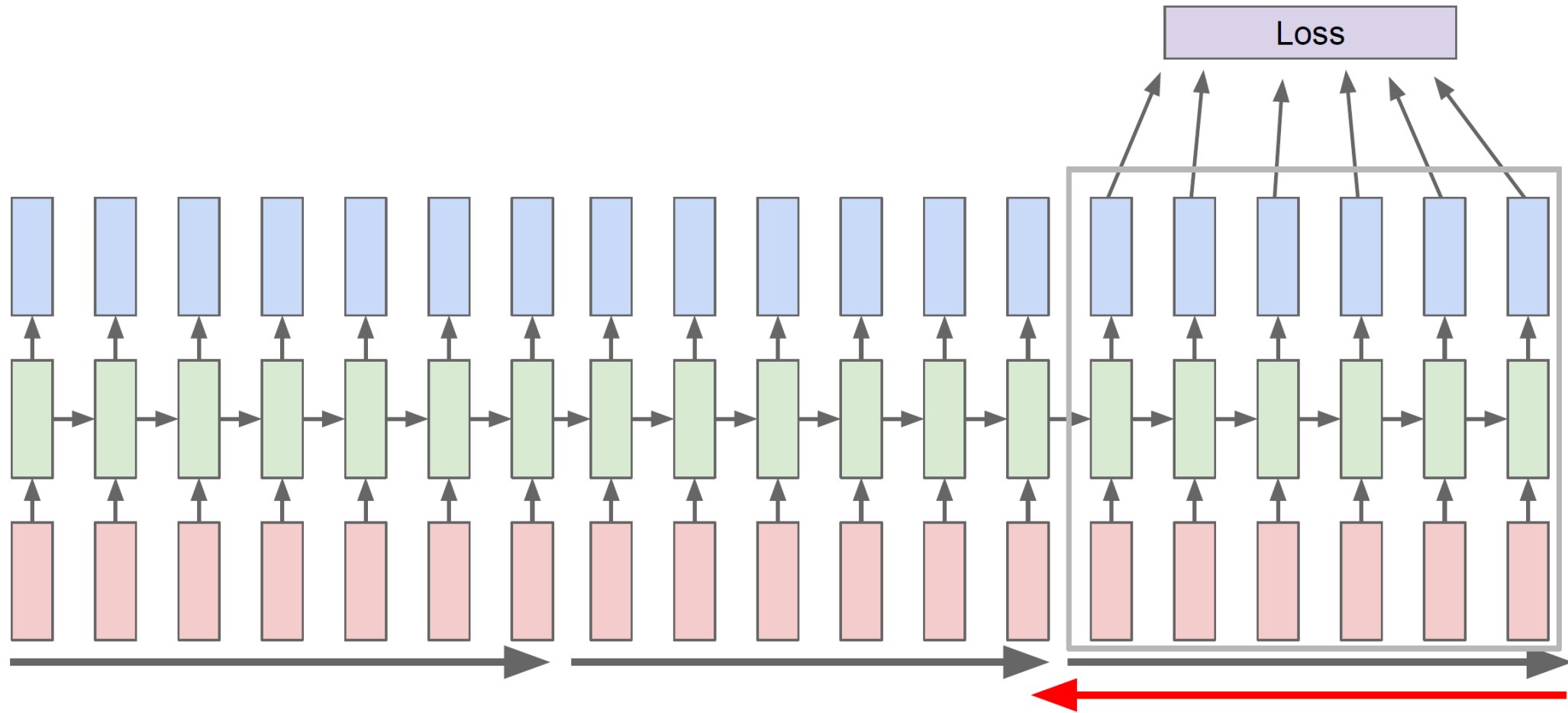
- Carry hidden states forward in time forever, but only backpropagate for some smaller number of steps
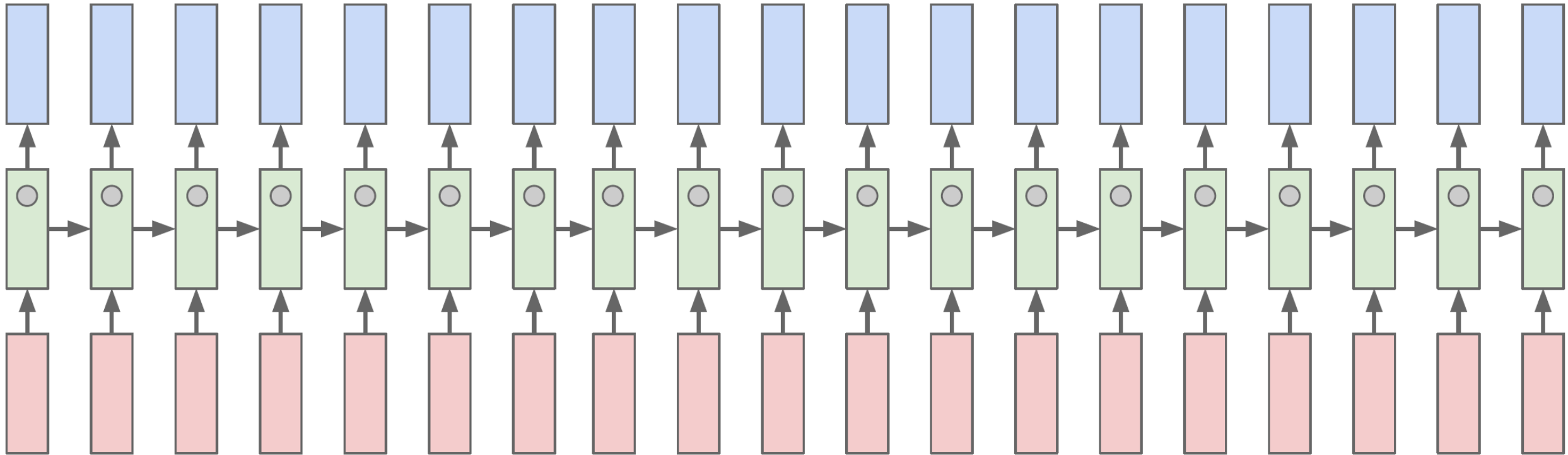
# Backpropagation through time (BPTT)

Loss

http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture10.pdf

# Searching for Interpretable Cells

http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture10.pdf

© NAVER Connect Foundation

# Character-level Language Model

- ## How RNN works

# Character-level Language Model

- ## How RNN works

  - Quote detection cell



http://karpathy.github.io/2015/05/21/rnn-effectiveness/

- ## How RNN works

  - ### If statement cell



http://karpathy.github.io/2015/05/21/rnn-effectiveness/

- ## RNN is excellent, but...

  - Multiplying the same matrix at each time step during backpropagation causes gradient vanishing or exploding

- ## Toy Example

  - $h_t = \tanh(w_{xh}x_t + w_{hh}h_{t-1} + b), t = 1,2,3$

  - For $w_{hh} = 3, w_{xh} = 2, b = 1$

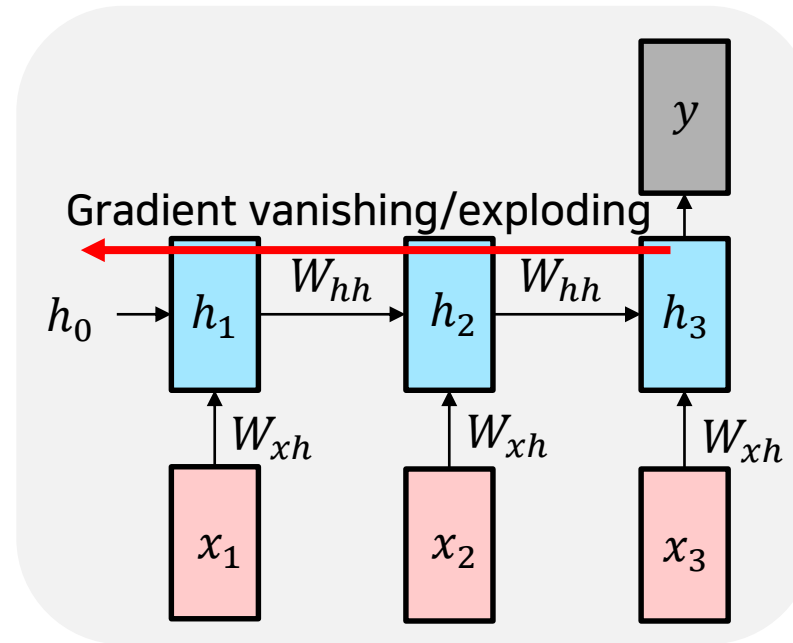$h_3 = \tanh(2x_3 + 3h_2 + 1)$

$h_2 = \tanh(2x_2 + 3h_1 + 1))$

$h_1 = \tanh(2x_1 + 3h_0 + 1))$

...

$h_3 = \tanh(2x_3 + 3\tanh(2x_2 + 3h_1 + 1) + 1)$

Gradient vanishing/exploding

- The reason why the vanishing gradient problem is important



https://imgur.com/gallery/vaNahKE

# End of Document
# Thank You.

# 4.
# Long Short-Term Memory (LSTM)
# Gated Recurrent Unit (GRU)

- ## Core Idea: pass cell state information straightly without any transformation

  - Solving long-term dependency problem



The repeating module in an LSTM contains four interacting layers.

http://colah.github.io/posts/2015-08-Understanding-LSTMs/

# Long Short-Term Memory (LSTM)

- What is LSTM (Long Short-Term Memory)?



The repeating module in an LSTM contains four interacting layers.

http://colah.github.io/posts/2015-08-Understanding-LSTMs/

- ## Long short-term memory

    - i: Input gate, Whether to write to cell
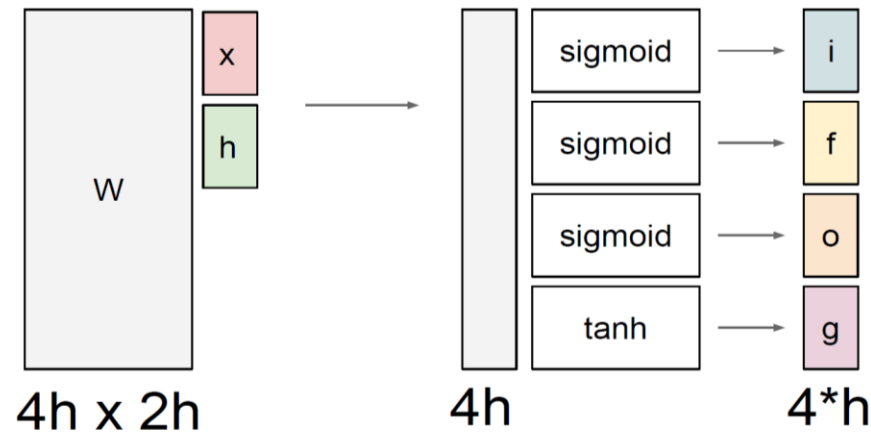
    - f: Forget gate, Whether to erase cell

    - o: Output gate, How much to reveal cell

    - g: Gate gate, How much to write to cell



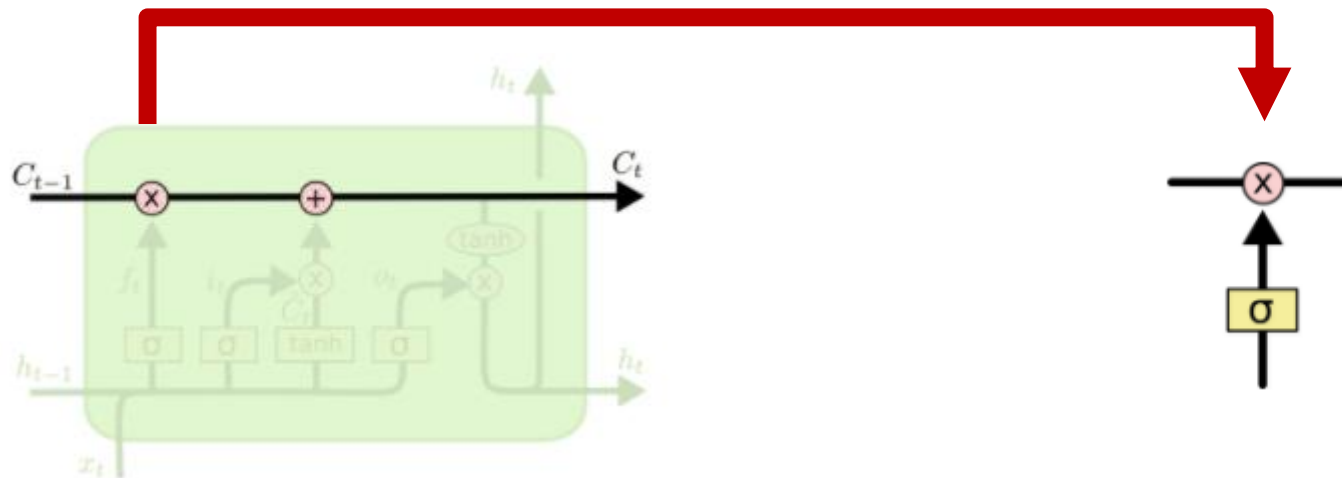$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

4h x 2h      4h      4*h

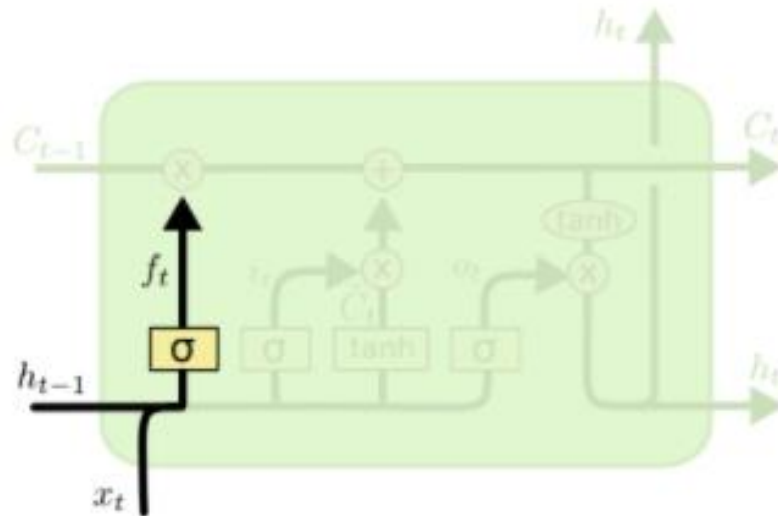Long short-term memory, Neural computation'97

# Long Short-Term Memory (LSTM)

- A gate exists for controlling how much information could flow from cell state

# Long Short-Term Memory (LSTM)

- ## Forget gate

  - $f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$

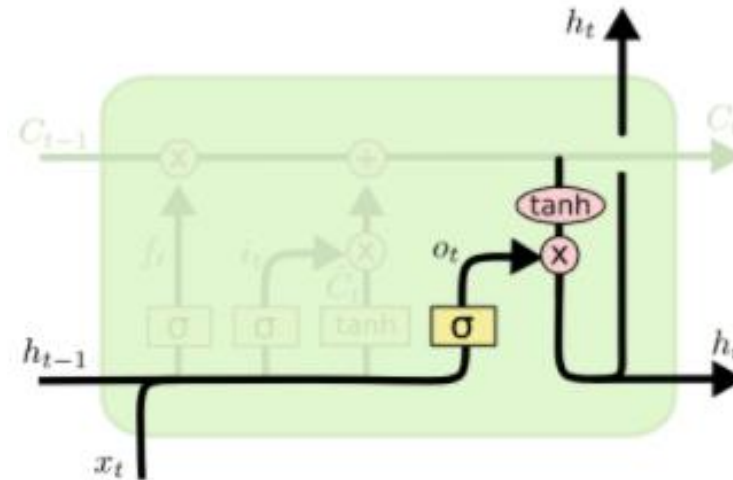# Long Short-Term Memory (LSTM)

- ## Generate information to be added and cut it by input gate

  - $i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$

  - $\widetilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$

- ## Generate new cell state by adding current information to previous cell state

  - $C_t = f_t \cdot C_{t-1} + i_t \cdot \widetilde{C}_t$



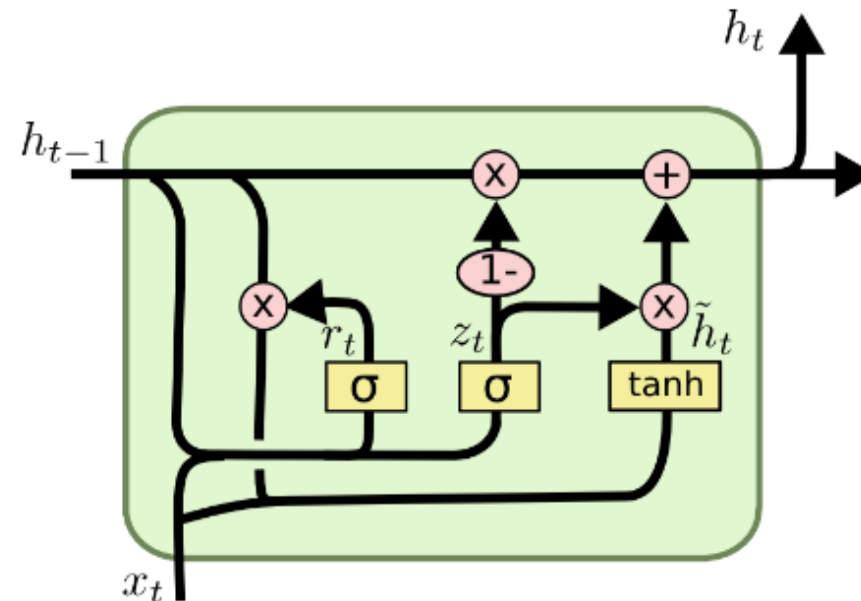/colah.github.io/posts/2015-08-Understanding-LSTMs/

- Generate hidden state by passing cell state to tanh and output gate

- Pass this hidden state to next time step, and output or next layer if needed

  - $o_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$

  - $h_t = o_t \cdot \tanh(C_t)$

# Gated Recurrent Unit (GRU)

- ## What is GRU?

  - $z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$

  - $r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$

  - $\widetilde{h}_t = \tanh(W \cdot [r_t \cdot h_{t-1}, x_t])$

  - $h_t = (1 - z_t) \cdot h_{t-1} + z_t \cdot \widetilde{h}_t$

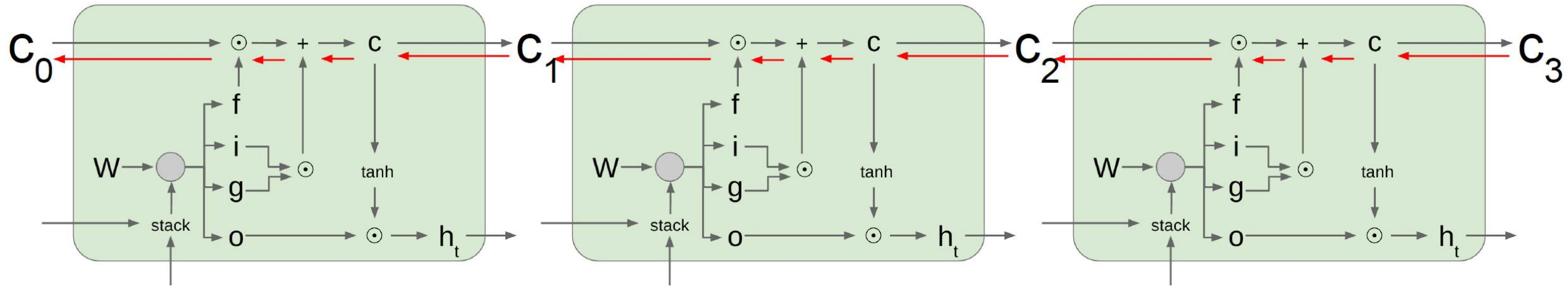  - c.f) $C_t = f_t \cdot C_{t-1} + i_t \cdot \widetilde{C}_t$

    in LSTM



http://colah.github.io/posts/2015-08-Understanding-LSTMs/

- Uninterrupted gradient flow !

# Summary on RNN/LSTM/GRU

- RNNs allow a lot of flexibility in architecture design

- Vanilla RNNs are simple but don't work very well

- Backward flow of gradients in RNN can explode or vanish

- Common to use LSTM or GRU: their additive interactions improve gradient flow

# References

- https://imgur.com/gallery/vaNahKE
- http://karpathy.github.io/2015/05/21/rnn-effectiveness/
- http://colah.github.io/posts/2015-08-Understanding-LSTMs/
- http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture10.pdf

# End of Document
# Thank You.

boostcamp AI Tech