# Sequences in Caffe

Jeff Donahue
CVPR Caffe Tutorial
June 6, 2015
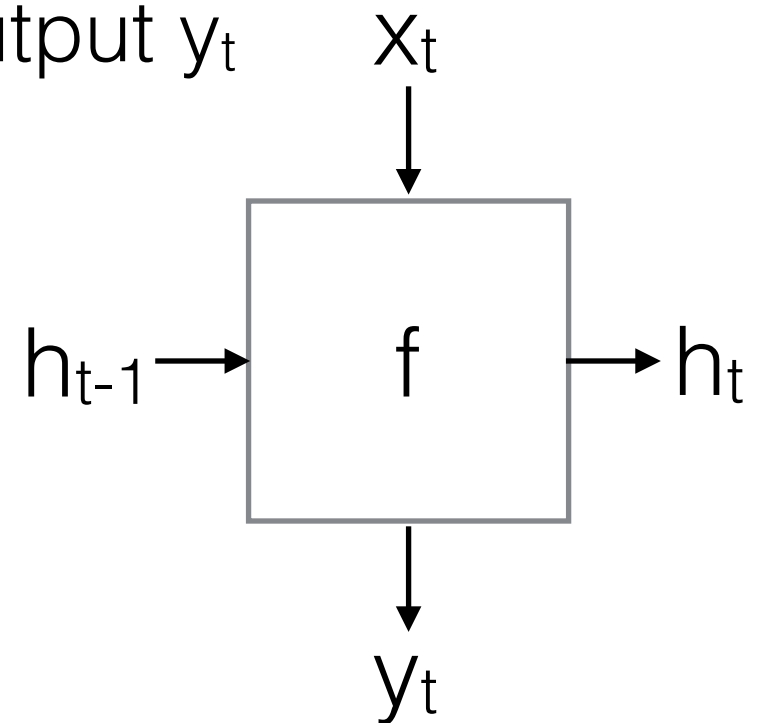
# Sequence Learning

- Instances of the form $\mathbf{x} = <x_1, x_2, x_3, \ldots, x_T>$

- Variable sequence length T

- Learn a transition function f with parameters W:

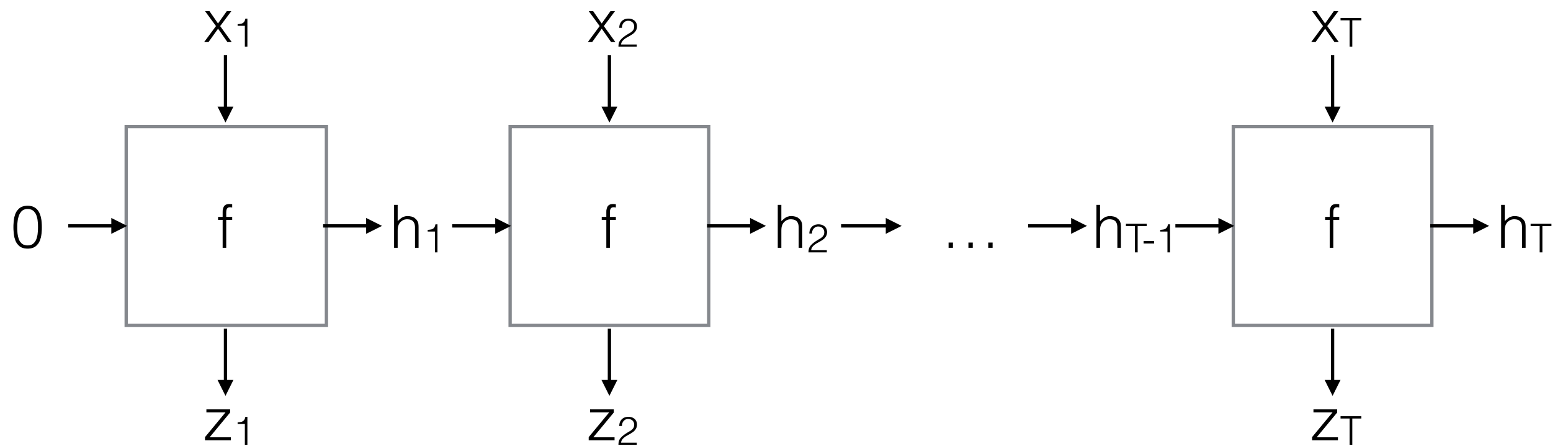- f should update hidden state $h_t$ and output $y_t$

$$h_0 := 0$$

for t = 1, 2, 3, …, T:

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

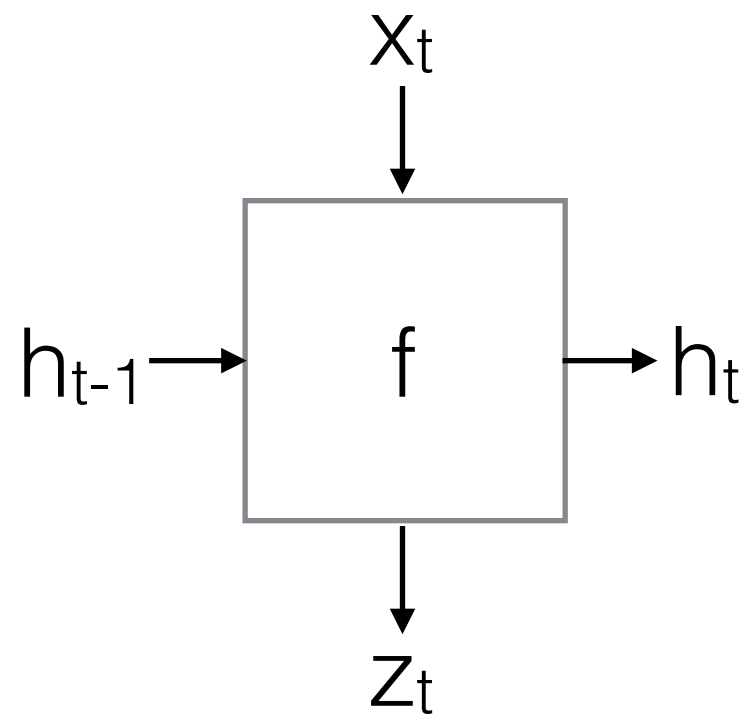# Sequence Learning

Equivalent to a T-layer deep network, unrolled in time

# Sequence Learning
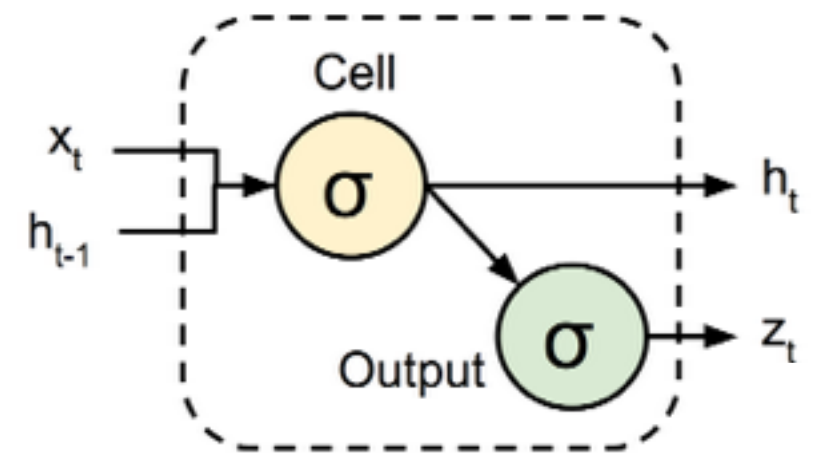
- What should the transition function f be?

$x_t$

$h_{t-1} \rightarrow \boxed{f} \rightarrow h_t$

$z_t$

- At a minimum, we want something **non-linear** and **differentiable**

# Sequence Learning
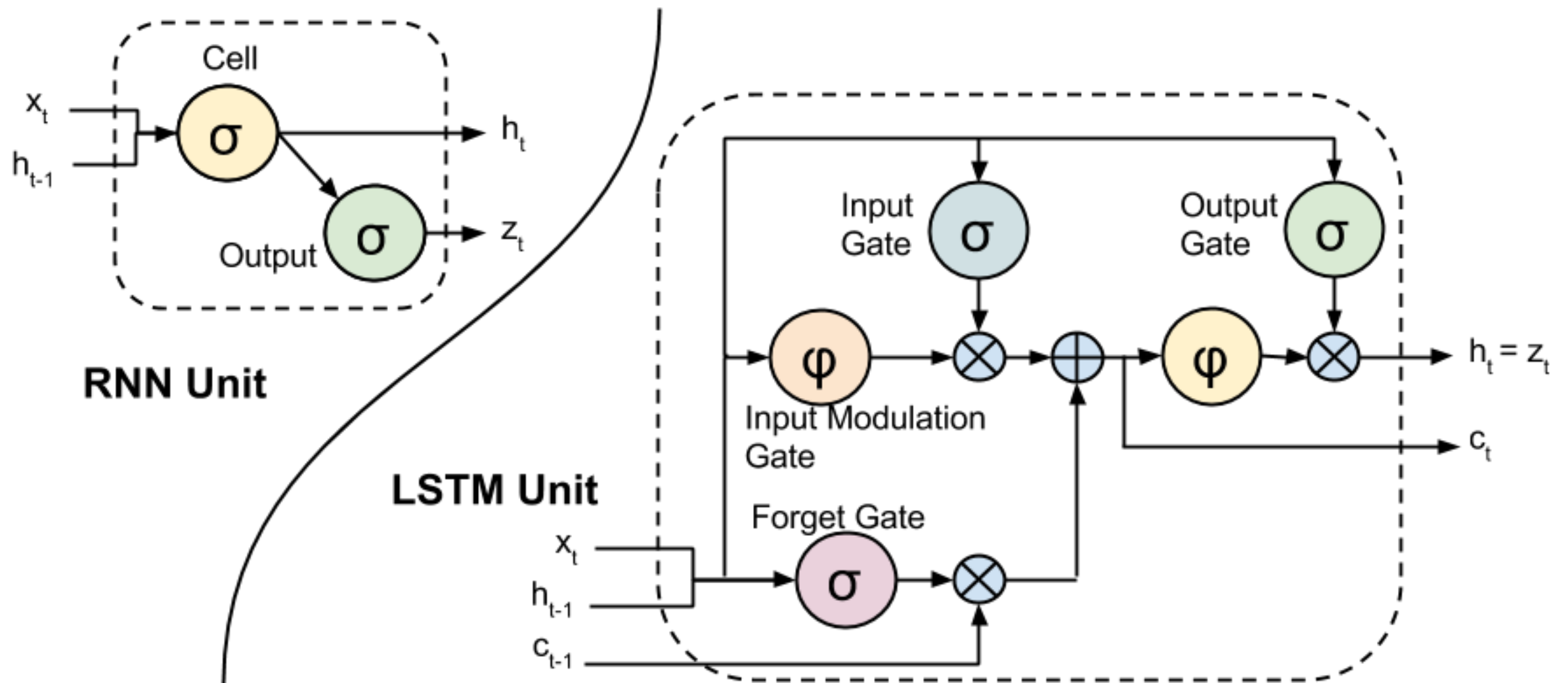
- A "vanilla" RNN:

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

  $$z_t = \sigma(W_{hz}h_t + b_z)$$



- Problems

  - Difficult to train — vanishing/exploding gradients

  - Unable to "select" inputs, hidden state, outputs
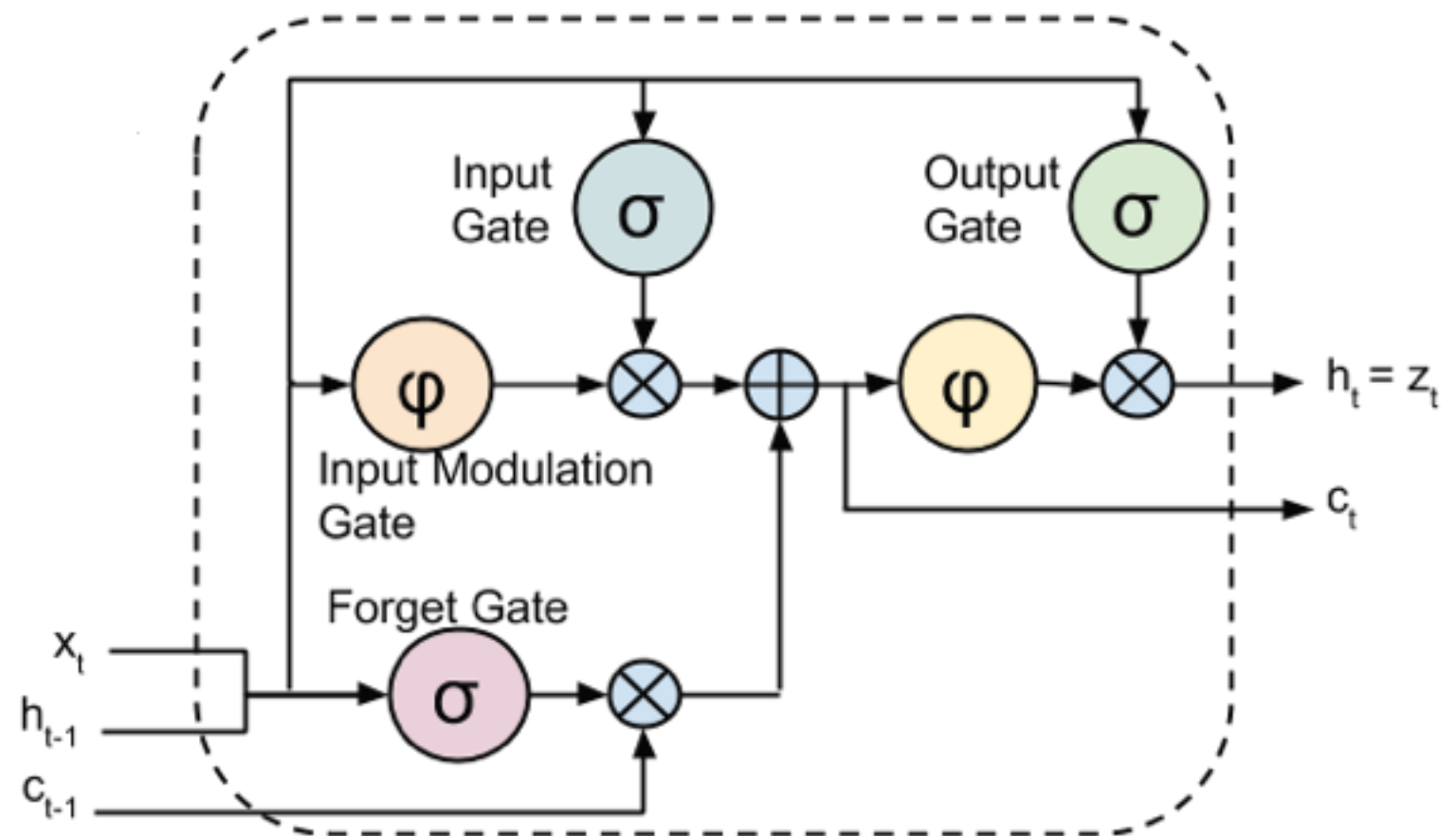
# Sequence Learning



Long Short-Term Memory (LSTM)
Proposed by Hochreiter and Schmidhuber, 1997

# Sequence Learning

- Allows long-term dependencies to be learned
- Effective for
  - speech recognition
  - handwriting recognition
  - translation
  - parsing



LSTM
(Hochreiter &
Schmidhuber, 1997)

# Sequence Learning

Exactly remember previous cell state — discard input



LSTM
(Hochreiter &
Schmidhuber, 1997)

# Sequence Learning



Forget previous cell state — only remember input

LSTM
(Hochreiter & Schmidhuber, 1997)

# Activity Recognition



sequential input

# Image Description

# Video Description



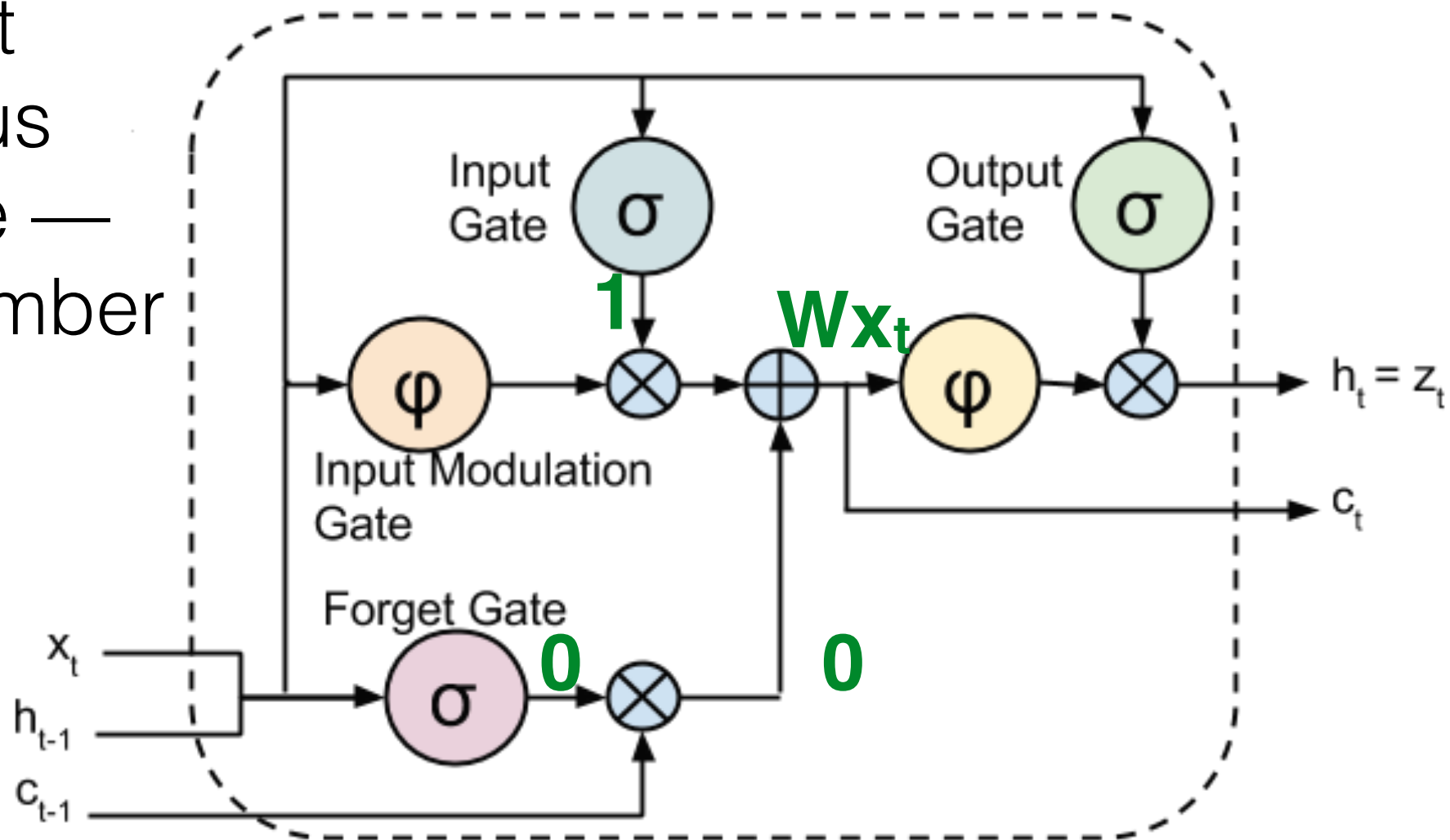sequential input & output

Sequence learning features now available in Caffe.
Check out PR #2033
"Unrolled recurrent layers (RNN, LSTM)"

# Training Sequence Models

- At training time, want the model to predict the next time step given all previous time steps: $p(w_{t+1} \mid w_{1:t})$
- Example: *A bee buzzes.*

|   | input | output |
|---|-------|--------|
| **0** | <BOS> | a |
| **1** | a | bee |
| **2** | bee | buzzes |
| **3** | buzzes | <EOS> |

# Sequence Input Format

- First input: "cont" (continuation) indicators (T x N)
- Second input: data (T x N x D)

N = 2, T = 6

| batch 1 | | | | | | | batch 2… | | |
|---|---|---|---|---|---|---|---|---|---|
| a | dog | fetches | <EOS> | the | bee | | buzzes | <EOS> | a |
| 0 | 1 | 1 | 1 | 0 | 1 | | 1 | 1 | 0 |
| cat | in | a | hat | <EOS> | a | | tree | falls | <EOS> |
| 0 | 1 | 1 | 1 | 1 | 0 | | 1 | 1 | 1 |

# Sequence Input Format

- Inference is exact over infinite batches
- Backpropagation approximate — truncated at batch boundaries

$N = 2, T = 6$

| batch 1 | | | | | | | batch 2... | | |
|---------|-----|---------|-------|-----|-----|---|--------|-------|-----|
| a | dog | fetches | <EOS> | the | bee | | buzzes | <EOS> | a |
| 0 | 1 | 1 | 1 | 0 | 1 | | 1 | 1 | 0 |
| cat | in | a | hat | <EOS> | a | | tree | falls | <EOS> |
| 0 | 1 | 1 | 1 | 1 | 0 | | 1 | 1 | 1 |

# Sequence Input Format

- Words are usually represented as one-hot vectors

| | | | |
|---|---|---|---|
| a | 0 | | 0 |
| the | 1 | *8000D vector;* | 0 |
| am | 2 | *all 0s except* | 0 |
| . | . | *index of* | . |
| . | . | *word* | . |
| cat | 3999 | | 0 |
| bee | 4000 | "bee" | 1 |
| dog | 4001 | | 0 |
| . | . | | . |
| . | . | | . |
| run | 7999 | | 0 |

*vocabulary with 8000 words*

# Sequence Input Format

- EmbedLayer projects one-hot vector

*8000D vector; all 0s except index of word*

"bee"

```
0
0
0
.
.
0
1
0
.
.
0
```

4000
(index of the word "bee" in our vocabulary)

```
layer {
  type: "Embed"
  embed_param {
    input_dim: 8000
    num_output: 500
  }
}
```

500D embedding of "bee"

```
2.7e-2
1.5e-1
3.1e+2
...
```
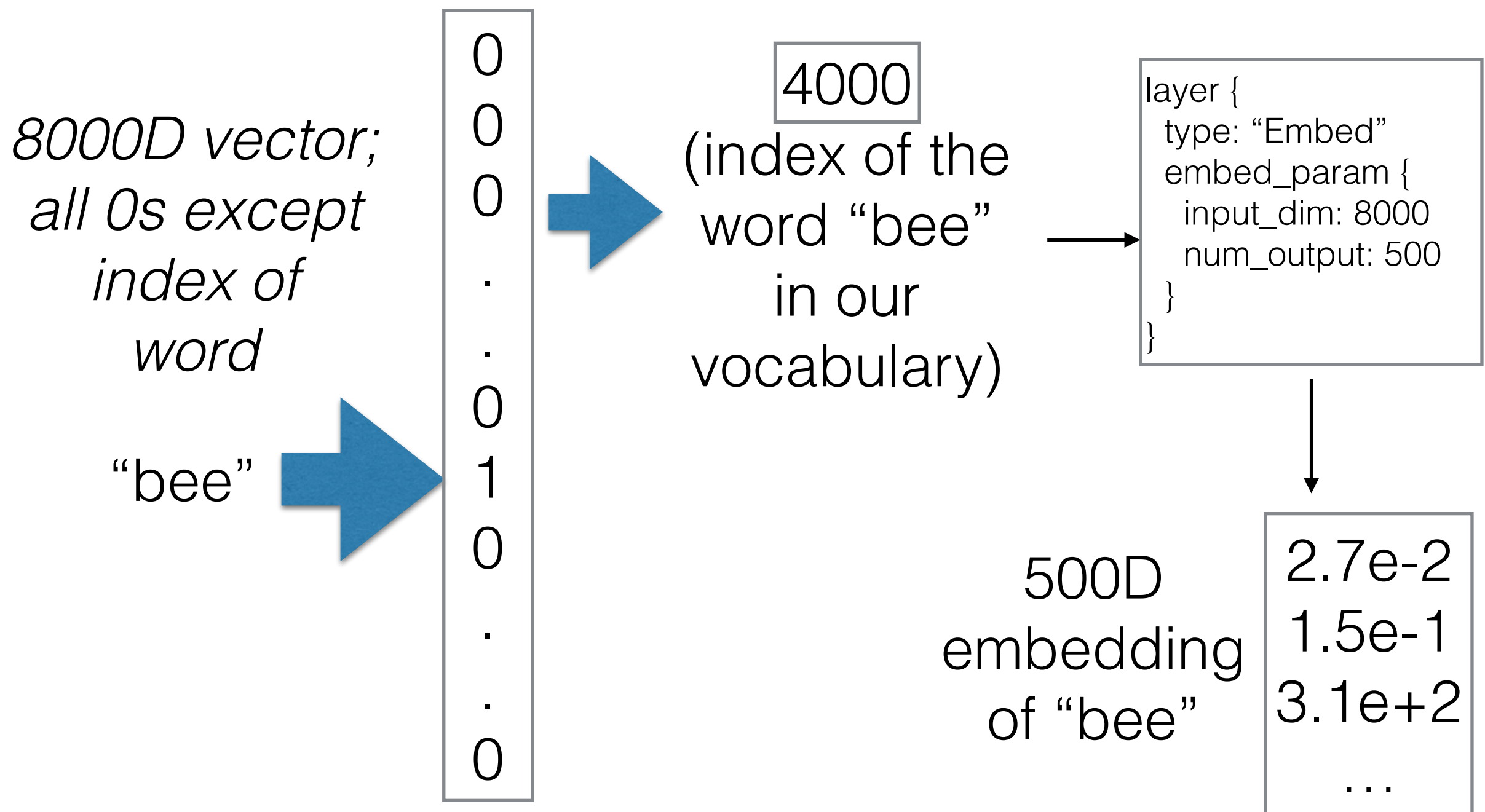
# Image Description



A female tennis player in action on the court.



A group of young men playing a game of soccer.



A man riding a wave on top of a surfboard.

# Image Description

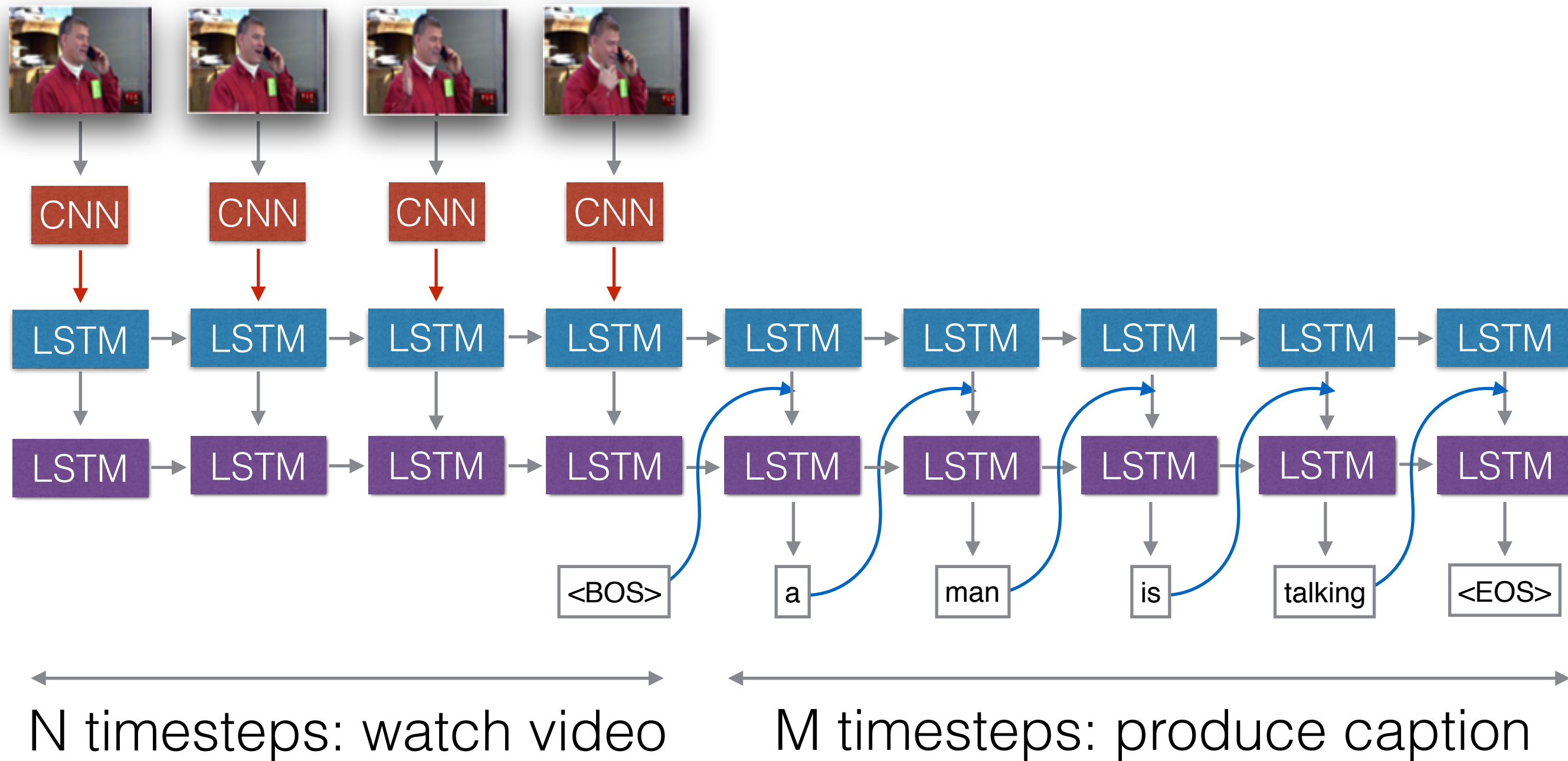

A black and white cat is sitting on a chair.



A large clock mounted to the side of a building.



A bunch of fruit that are sitting on a table.

# Video Description



N timesteps: watch video          M timesteps: produce caption

Venugopalan et al., "Sequence to Sequence -- Video to Text," 2015.
http://arxiv.org/abs/1505.00487