



Deep Learning Explained

Module 5: Recurrence (RNN) and Long-Short Term Memory (LSTM)

Sayan D. Pathak, Ph.D., Principal ML Scientist, Microsoft

Roland Fernandez, Senior Researcher, Microsoft

Module outline

Application:

Time series forecasting with IOT data

Model:

Recurrence

Long-short term memory cell

Concepts:

Recurrence

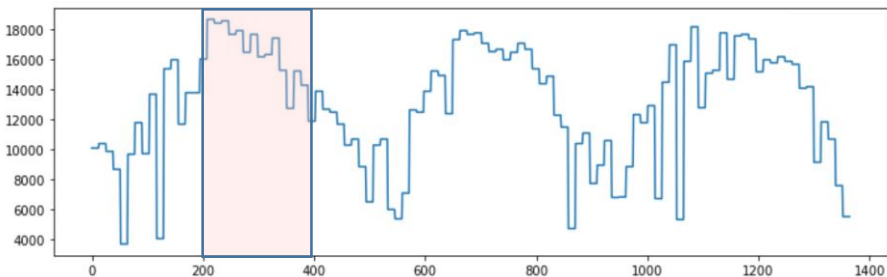
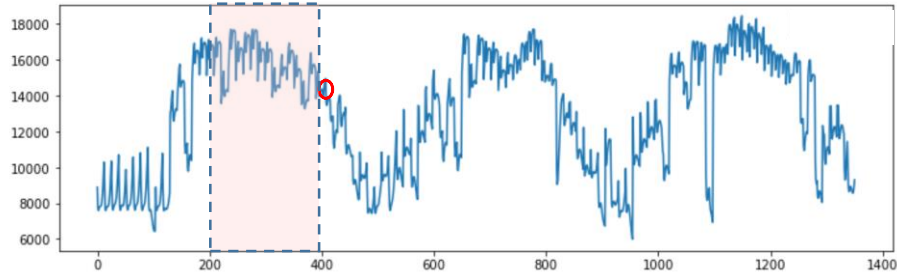
LSTM

Dropout

Train-Test-Predict workflow

Sequences (many to one)

Problem: Time series prediction with IOT data

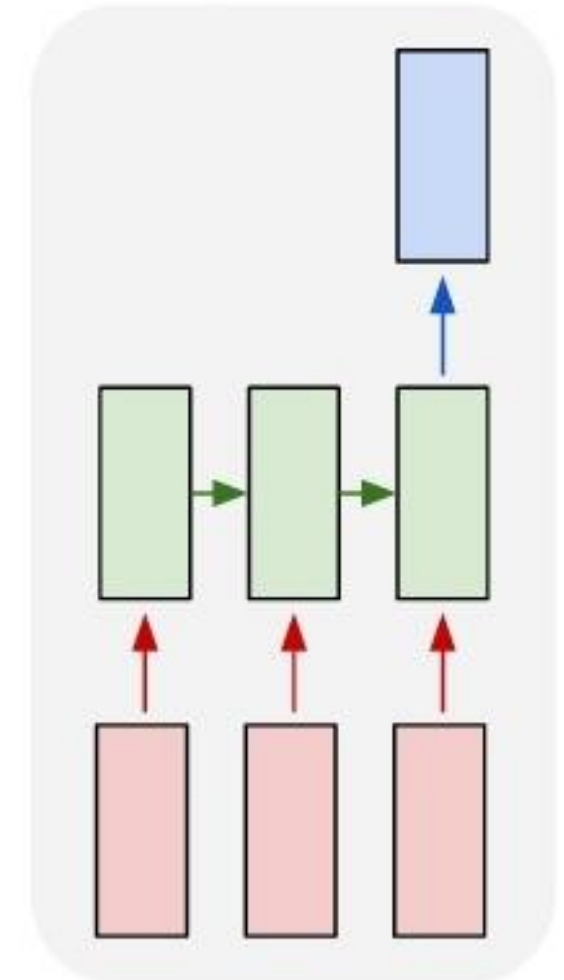


Output
(Y : $n \times$ future prediction)

Model
Rec = Recurrence

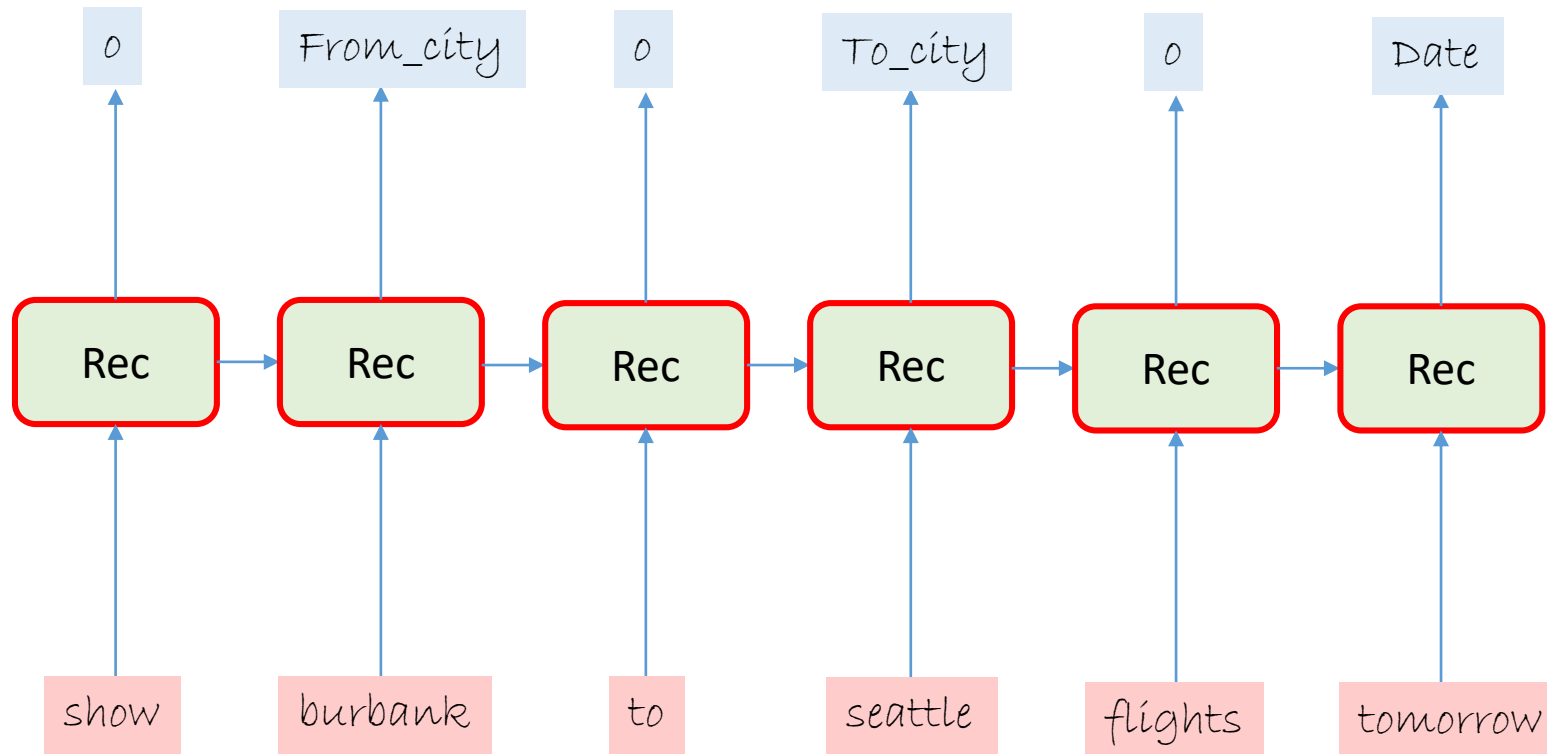
Input feature
(X : $n \times 14$ data pts)

many to one

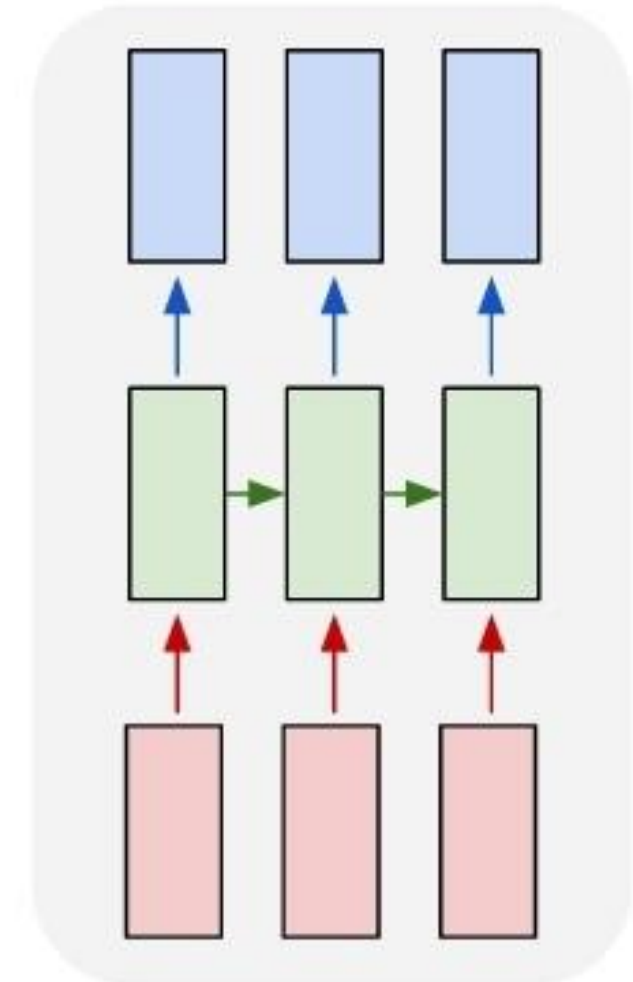


Sequences (many to many + 1:1)

Problem: Tagging entities in Air Traffic Controller (ATIS) data



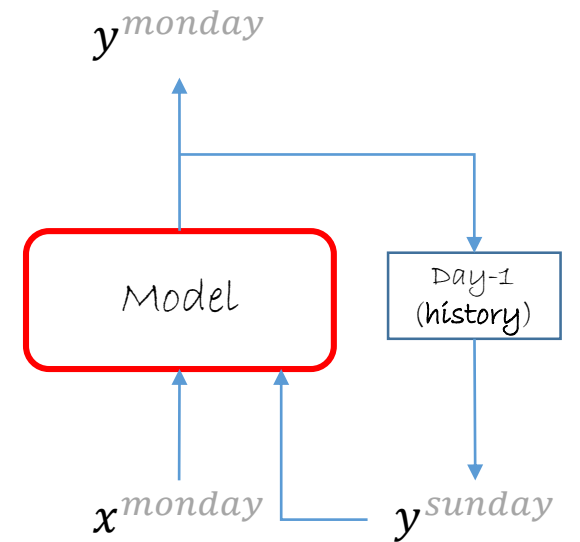
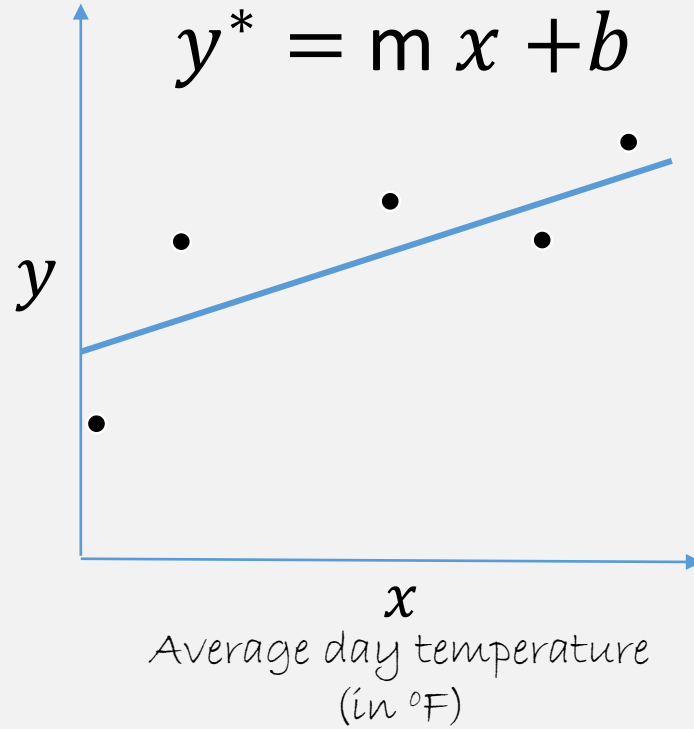
many to many



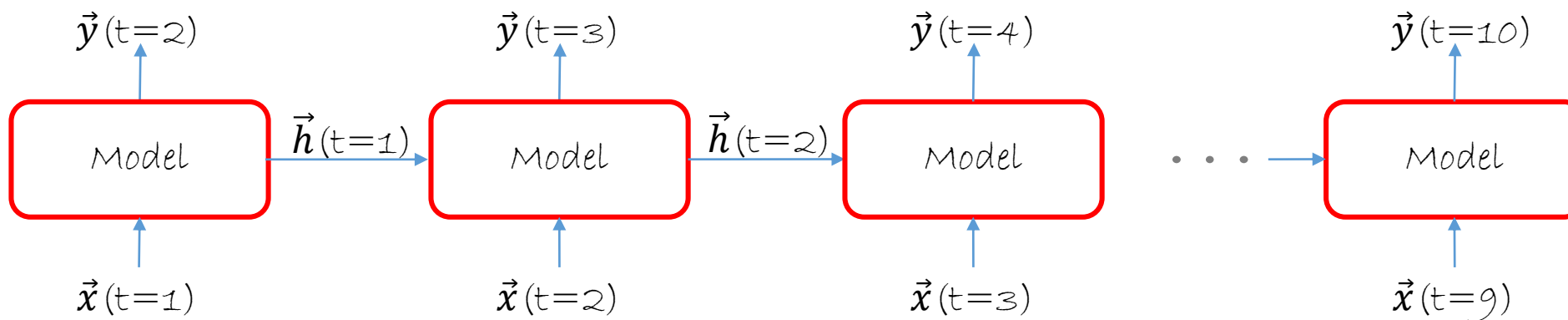
Forecasting



Solar panel Output
(in W)



Recurrence



$\vec{x}(t)$: Input (n -dimensional array) at time t

$\vec{y}(t)$: Output (c -dimensional array) at time t

$\vec{h}(t)$: Internal State [m -dimensional array] at time t a.k.a history

Input:

For numeric:

Array of numeric values coming from different sensor

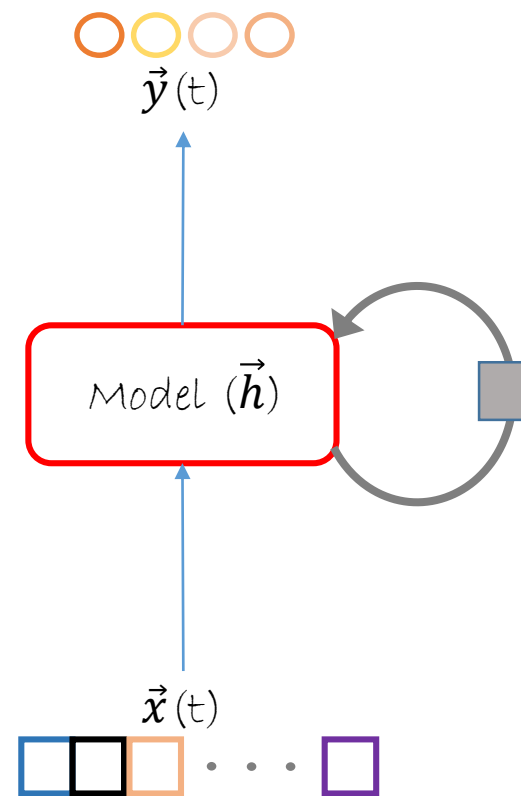
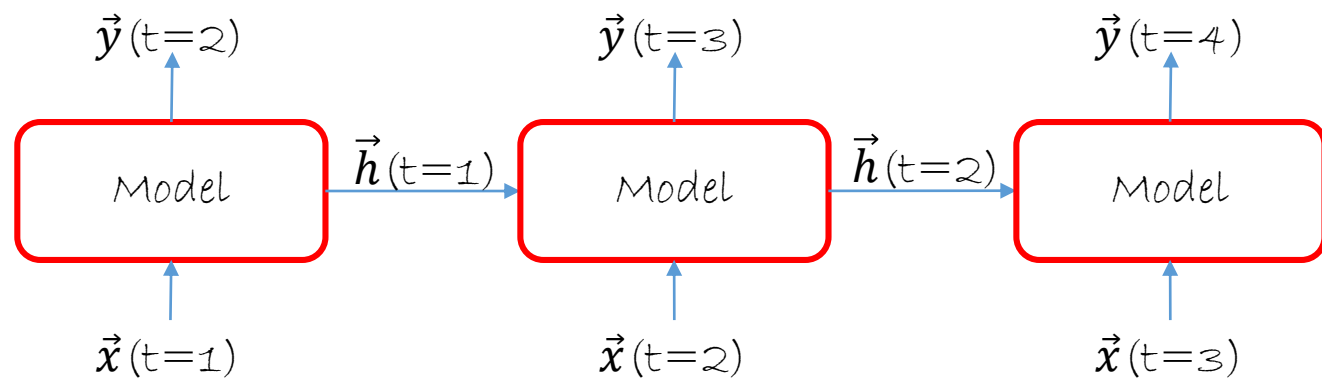
For an image:

Pixels in an array, Map the image pixels to a compact representation (say n values)

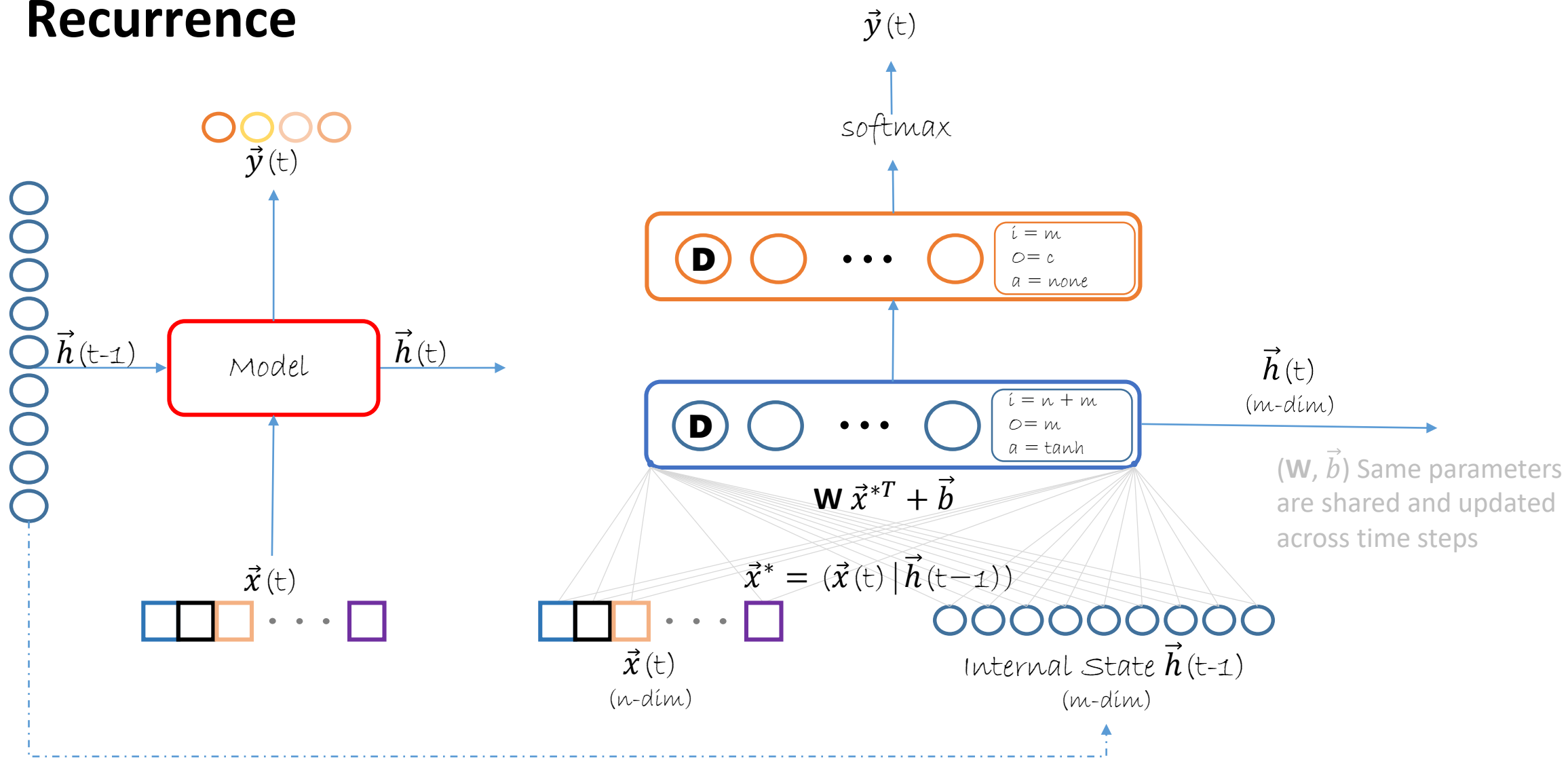
For word in text:

Represent words as a numeric vector using embeddings (word2vec or GloVe)

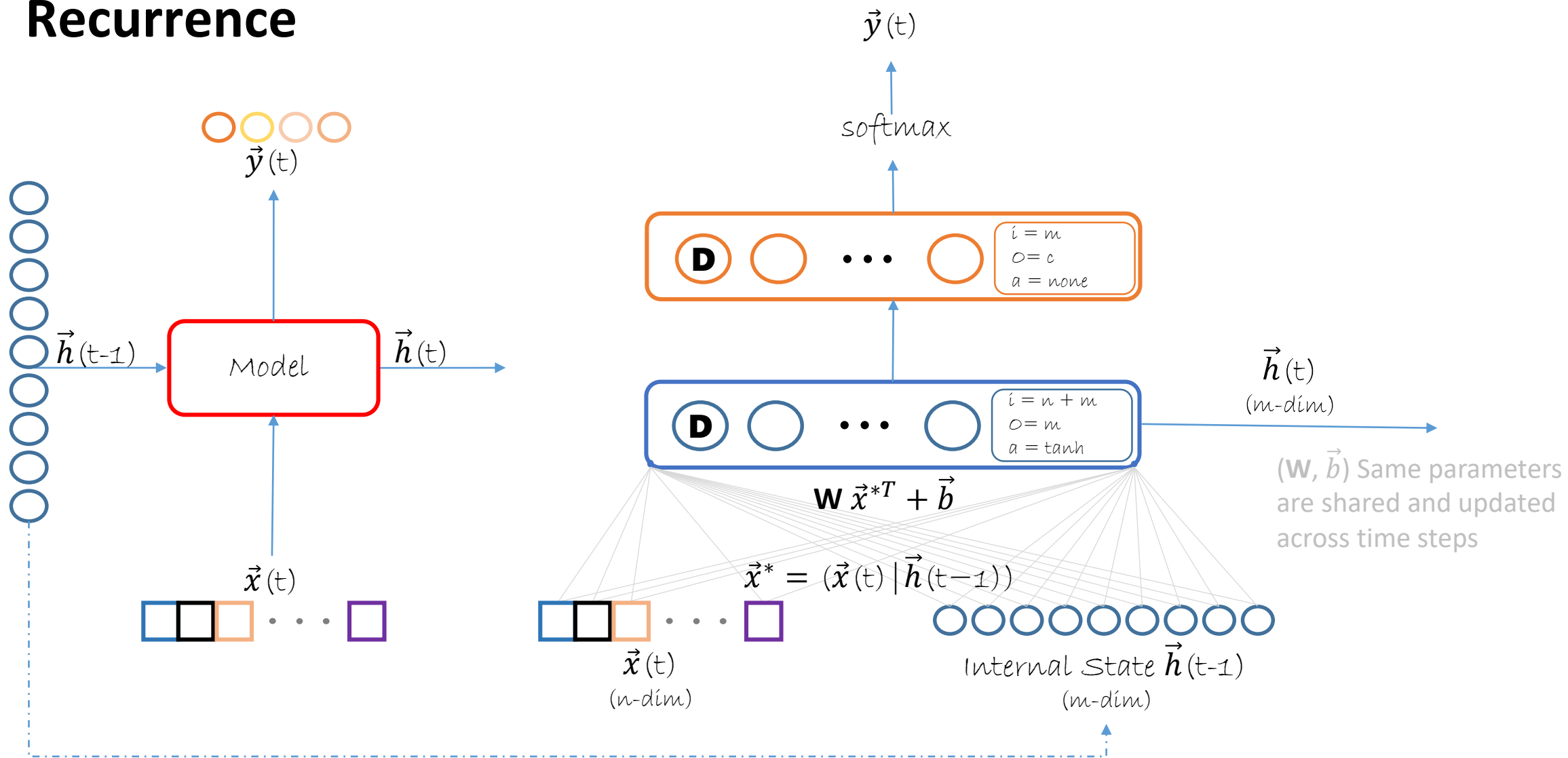
Recurrence



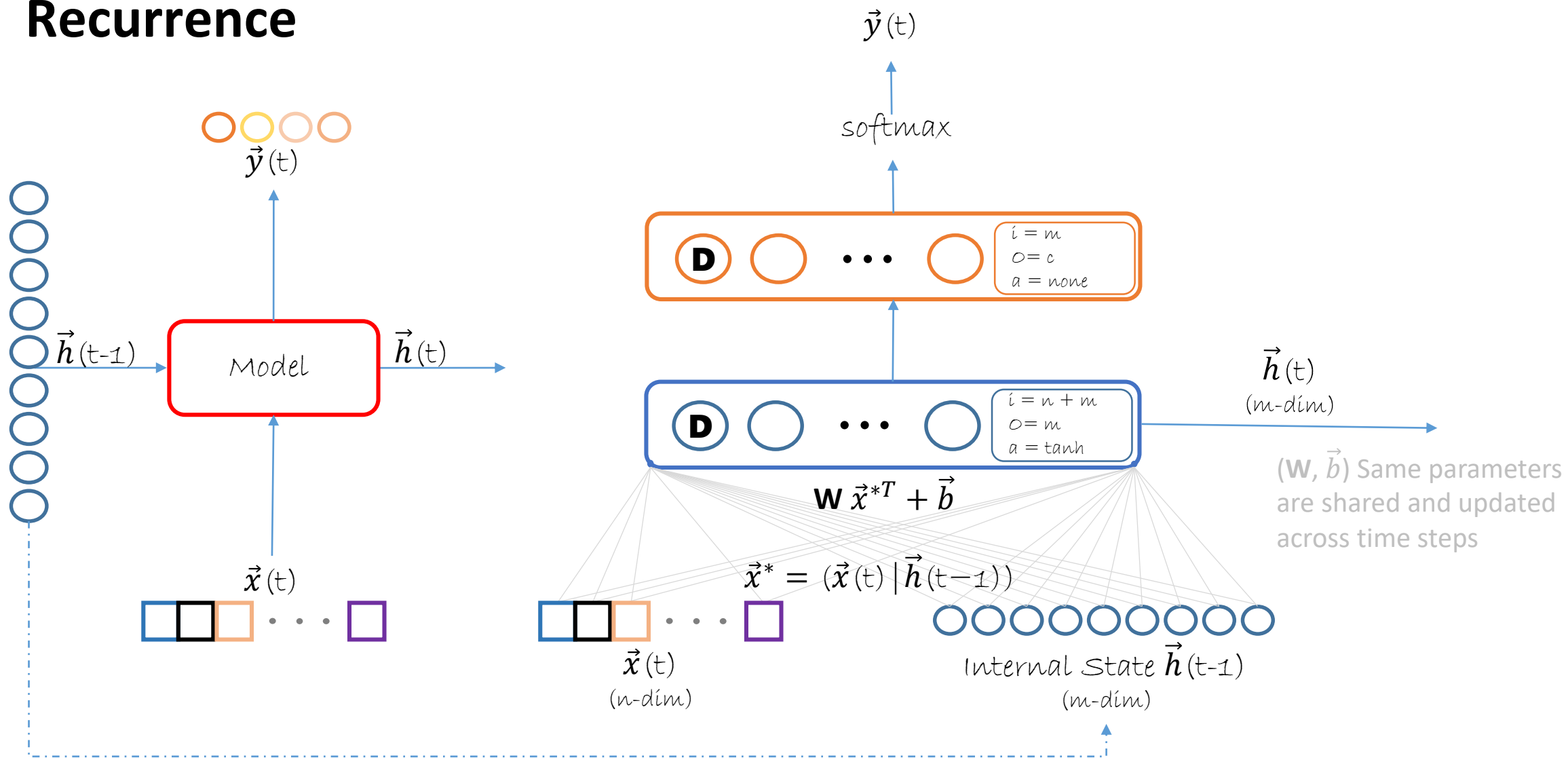
Recurrence



Recurrence

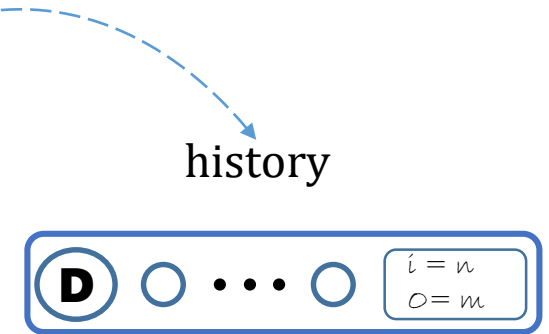
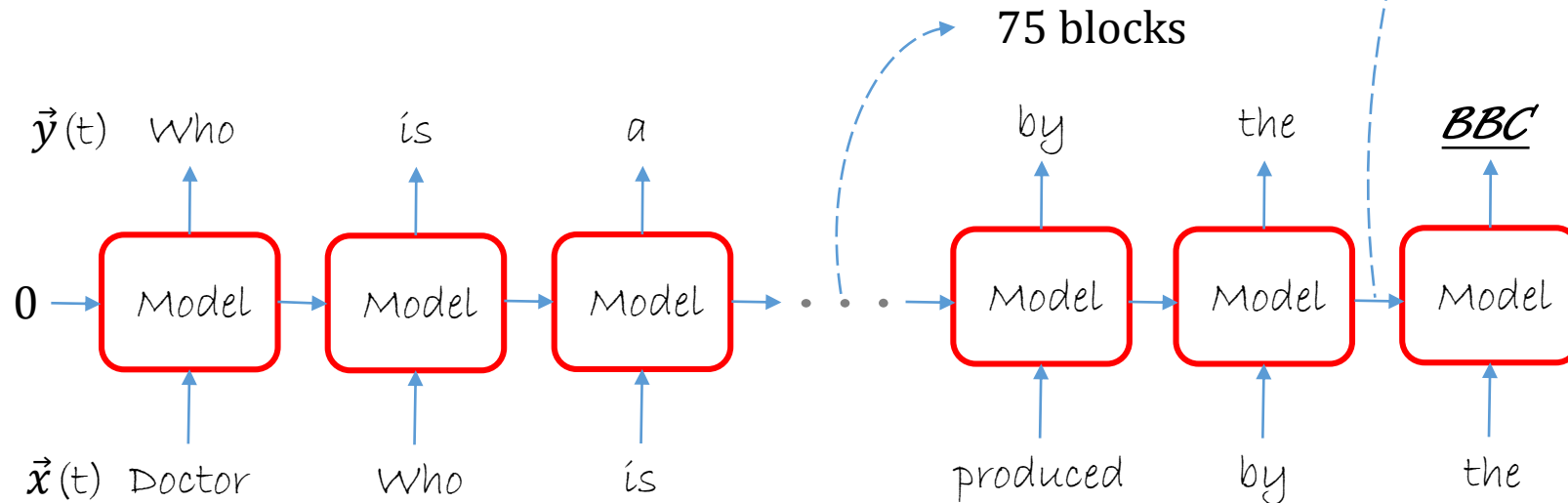


Recurrence



Recurrence (Vanishing Gradients)

Doctor Who is a British science-fiction television programme produced by the BBC since 1963. The programme depicts the adventures of the Doctor, a Time Lord—a space and time-travelling humanoid alien. He explores the universe in his TARDIS, a sentient time-travelling space ship. Accompanied by companions, the Doctor combats a variety of foes, while working to save civilizations and help people in need. This television series produced by the (...)

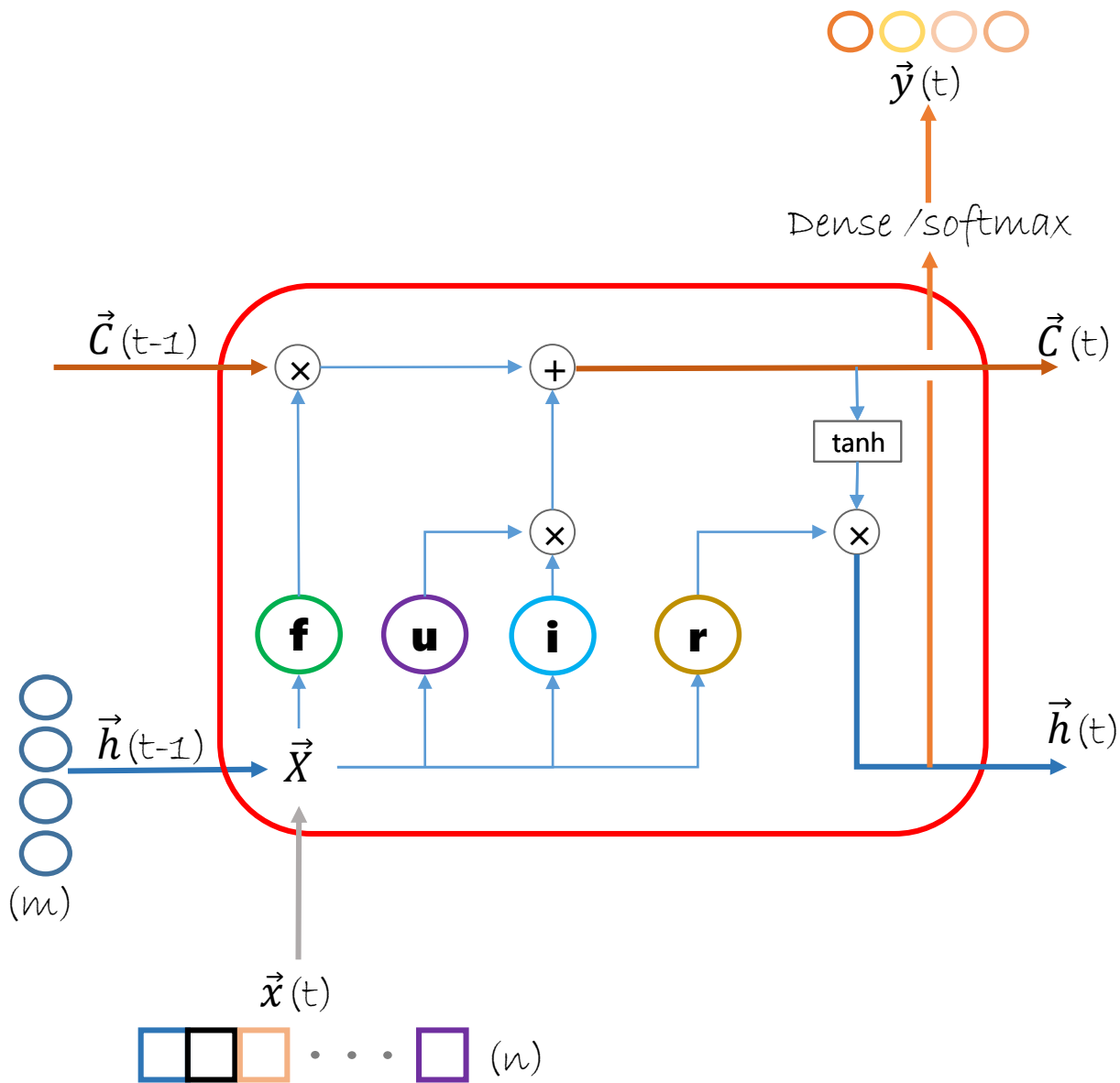


$$\vec{h} = \mathbf{W} \vec{x}^T + \vec{b}$$

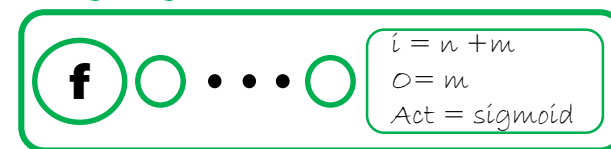


A single set of (\mathbf{W}, \vec{b})
has
limited memory

Long-Short Term Memory (LSTM)

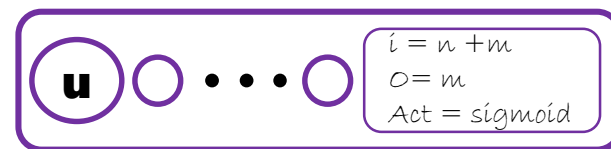


Forget gate



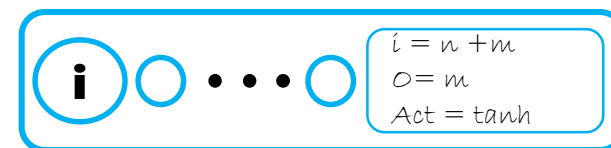
$$\vec{f} = \text{sigmoid}(\mathbf{W}_f \vec{X}^T + \vec{b}_f)$$

update gate



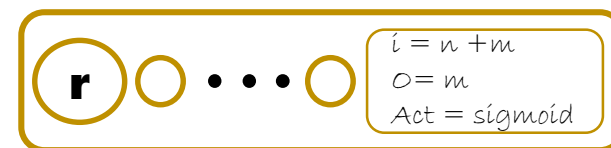
$$\vec{u} = \text{sigmoid}(\mathbf{W}_u \vec{X}^T + \vec{b}_u)$$

Input



$$\vec{X}^* = \tanh(\mathbf{W}_i \vec{X}^T + \vec{b}_i)$$

Result gate



$$\vec{r} = \text{sigmoid}(\mathbf{W}_r \vec{X}^T + \vec{b}_r)$$

New cell memory

$$\vec{C}(t) = \vec{C}(t-1) \times \mathbf{f} + \mathbf{i} \times \mathbf{u}$$

New history

$$\vec{h}(t) = \tanh(\vec{C}(t)) \times \mathbf{r}$$

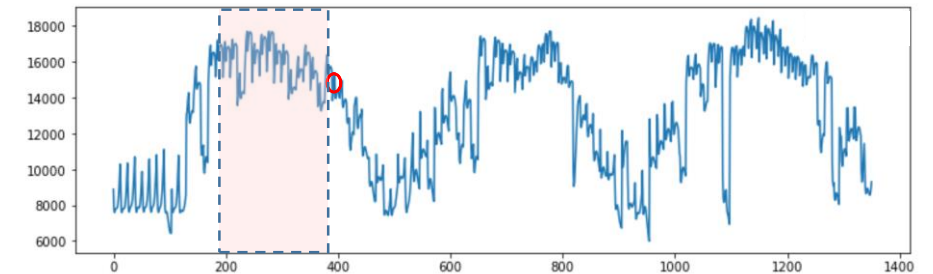
Time-series forecasting

Problem: Time series prediction with IOT data

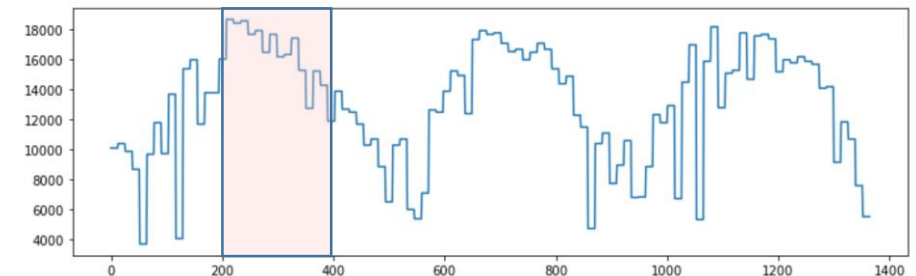
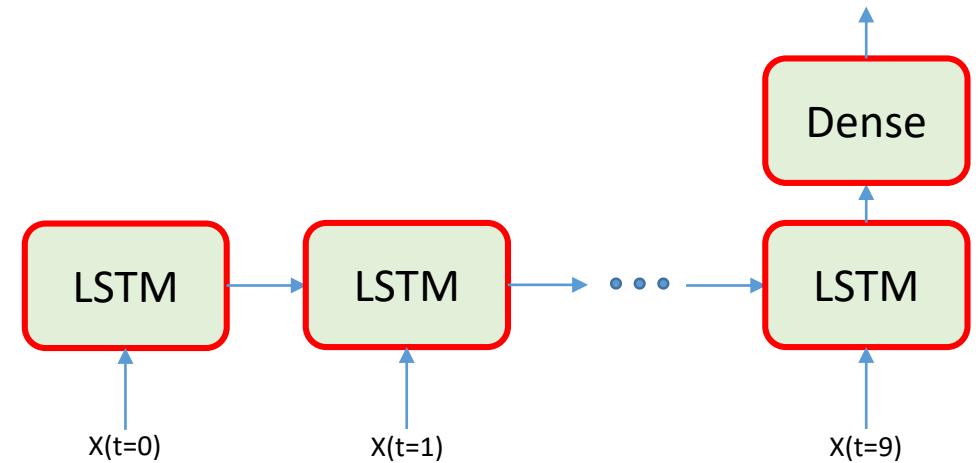
Output
(Y : $n \times$ future prediction)

```
z = create_model(x):  
    m = C.layers.Recurrence(C.layers.LSTM(TIMESTEPS))(x)  
    m = C.sequence.last(m)  
    m = C.layers.Dense(1)(m)  
    return m
```

Input feature
(X : $n \times 14$ data pnts)



Predict (Y^*)



Dropout

Problem:

Overfitting

Model works great with training data

With new data (unseen during training): high prediction error

Classical Approach:

L1 / L2 regularization

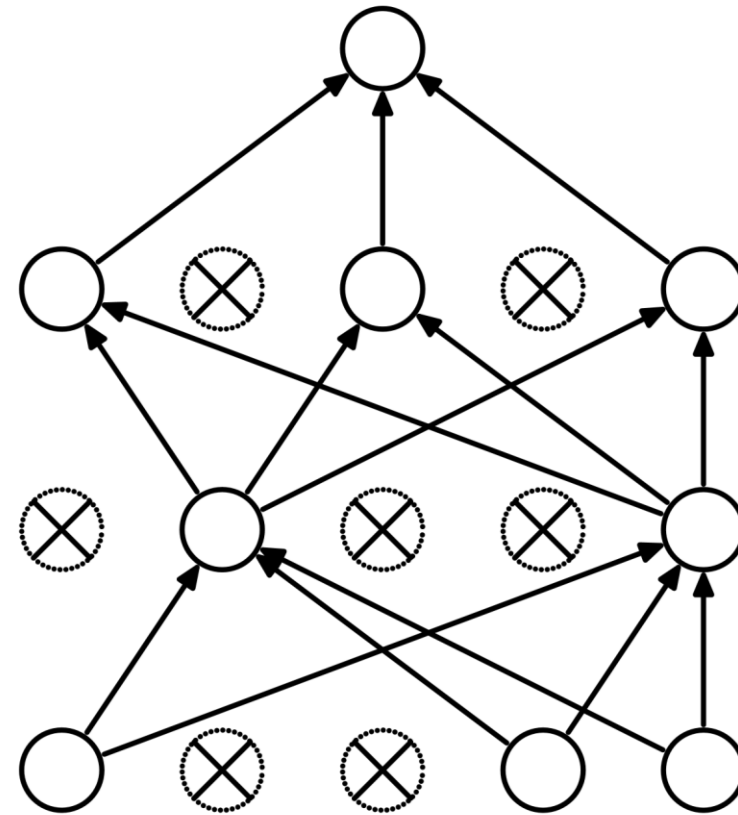
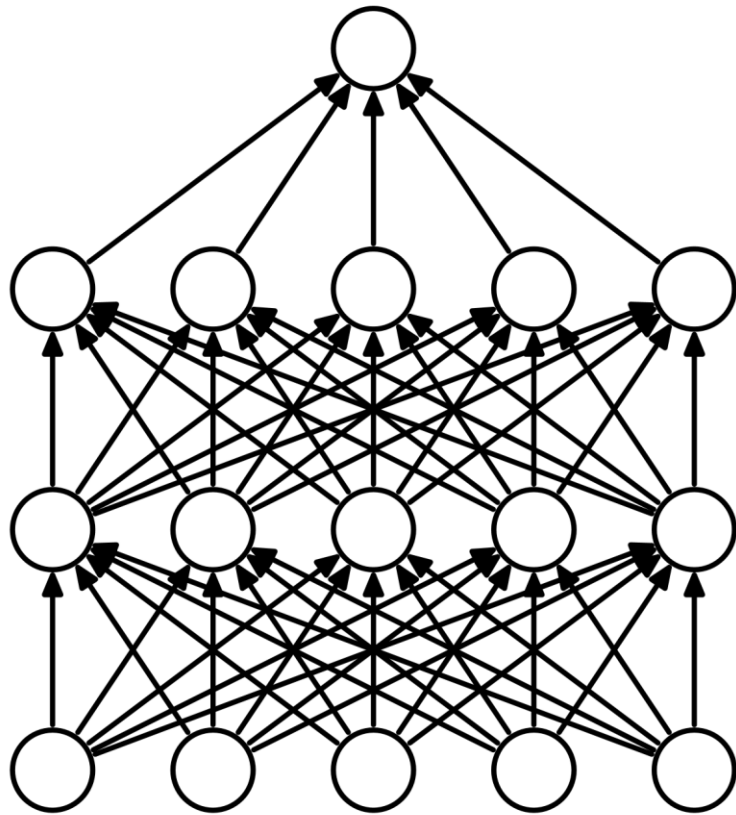
Data augmentation / train with noise added

Early stopping

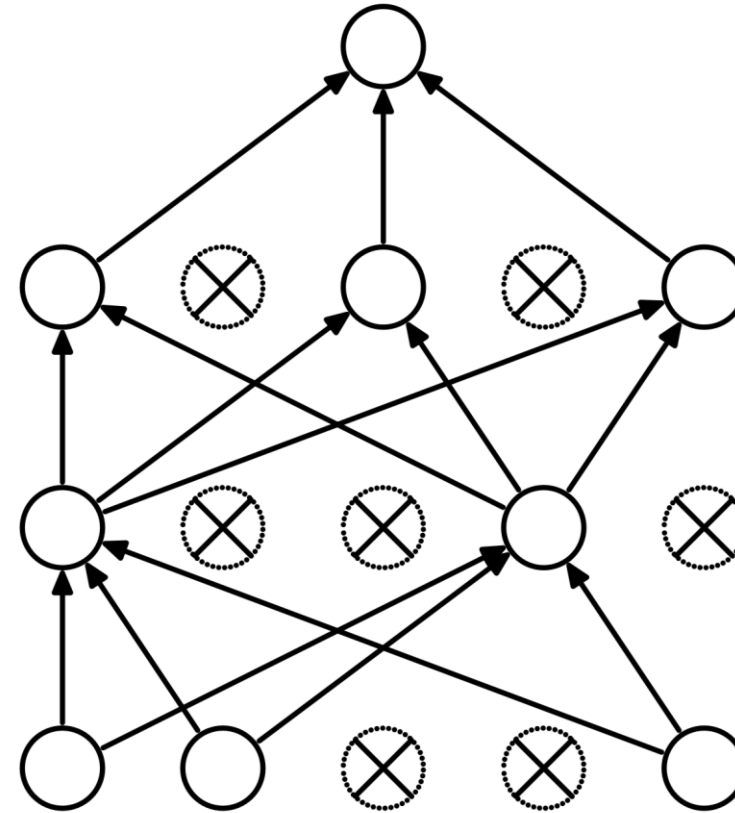
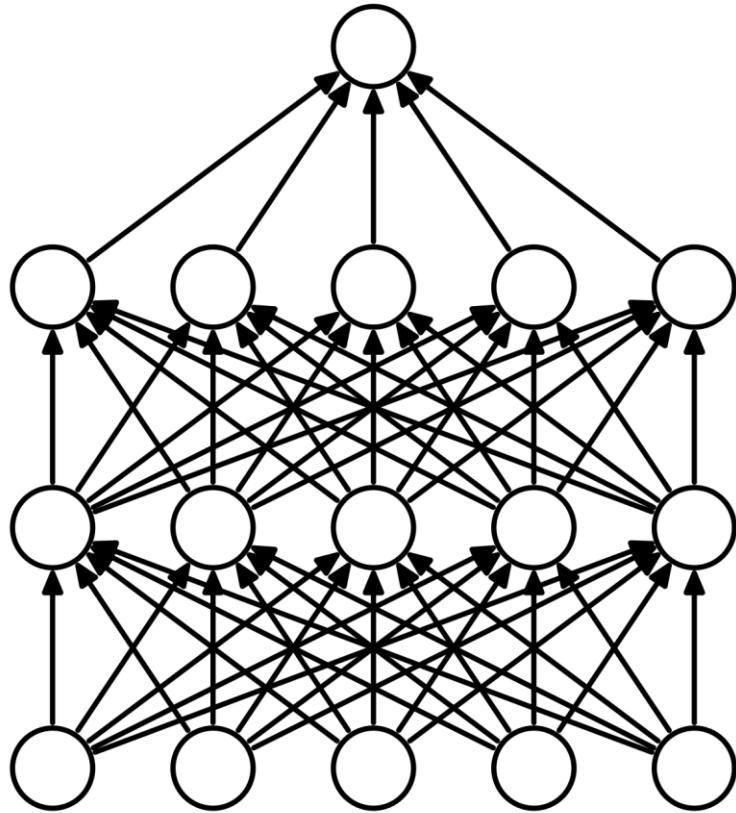
Dropout

Extremely effective technique to tackle overfitting in neural networks

Dropout



Dropout



Time-series forecasting

IOT data:

- ✓ Output of a solar panel, measurements are recorded at every 30 min interval:
 - solar.current: Current production in Watts
 - solar.total: Total production for the day so far in Watt/hour

Data Summary:

- ✓ Starting at a time in the day, two values are recorded

time, solar.current, solar.total

7am, 6.3, 1.7

7:30am, 44.3, 11.4

- ✓ 3 years of data

...

- ✓ The input data is not cleansed i.e., errors (panel failed to report) is included

Data pre-processing

Goal:

- ✓ Compose sequence such that each training instance would be:
 - $X = [\text{solar.current} @ t = 1 - t = 14]$ ($t=1 - 14$: corresponds to 1 day)
 - $Y = \text{Predicted total production for a future day}$

Pre-processing:

✓ Steps:

- read raw data into a pandas dataframe,
- normalize the data,
- group by day,
- append the columns "solar.current.max" and "solar.total.max", and
- generate the sequences for each day.

✓ Data filtering:

- If X has less than 8 data points - we skip
- If X has more than 14 data points - we truncate

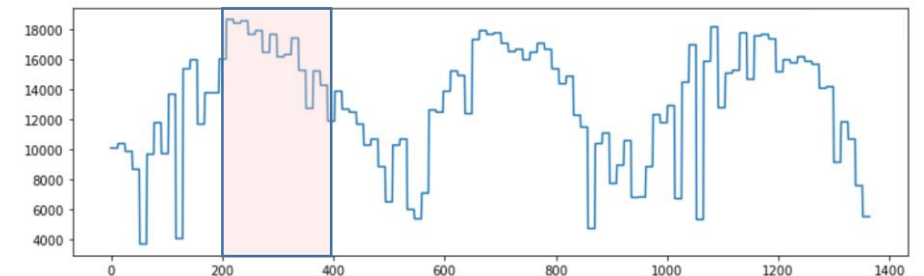
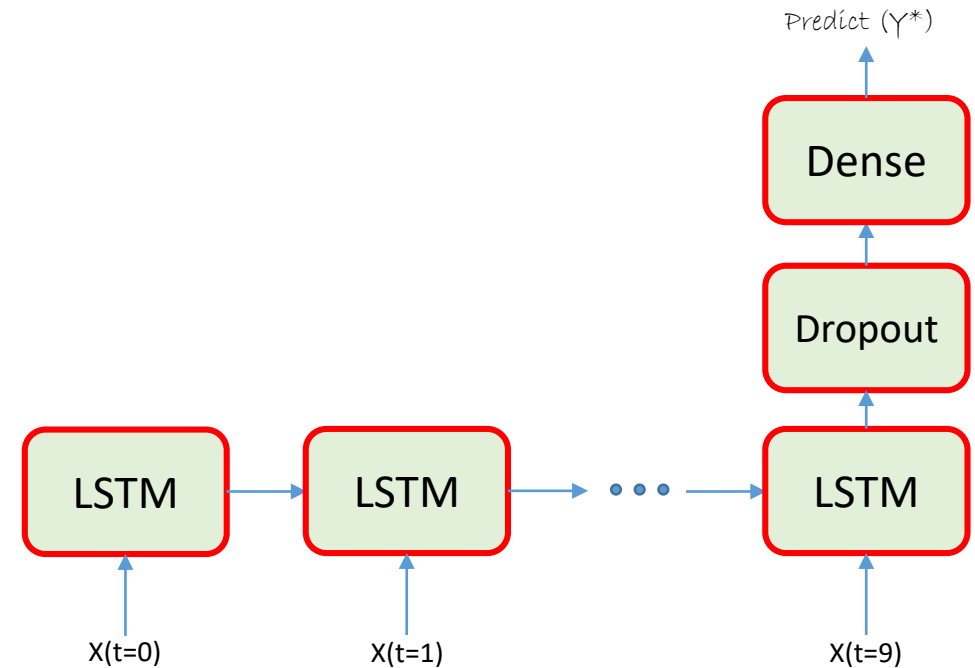
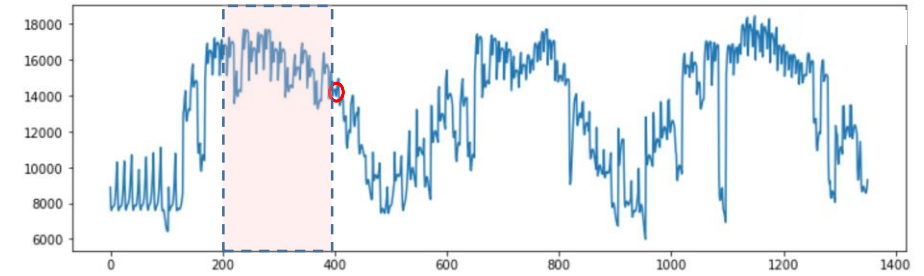
Time-series forecasting

Problem: Time series prediction with IOT data

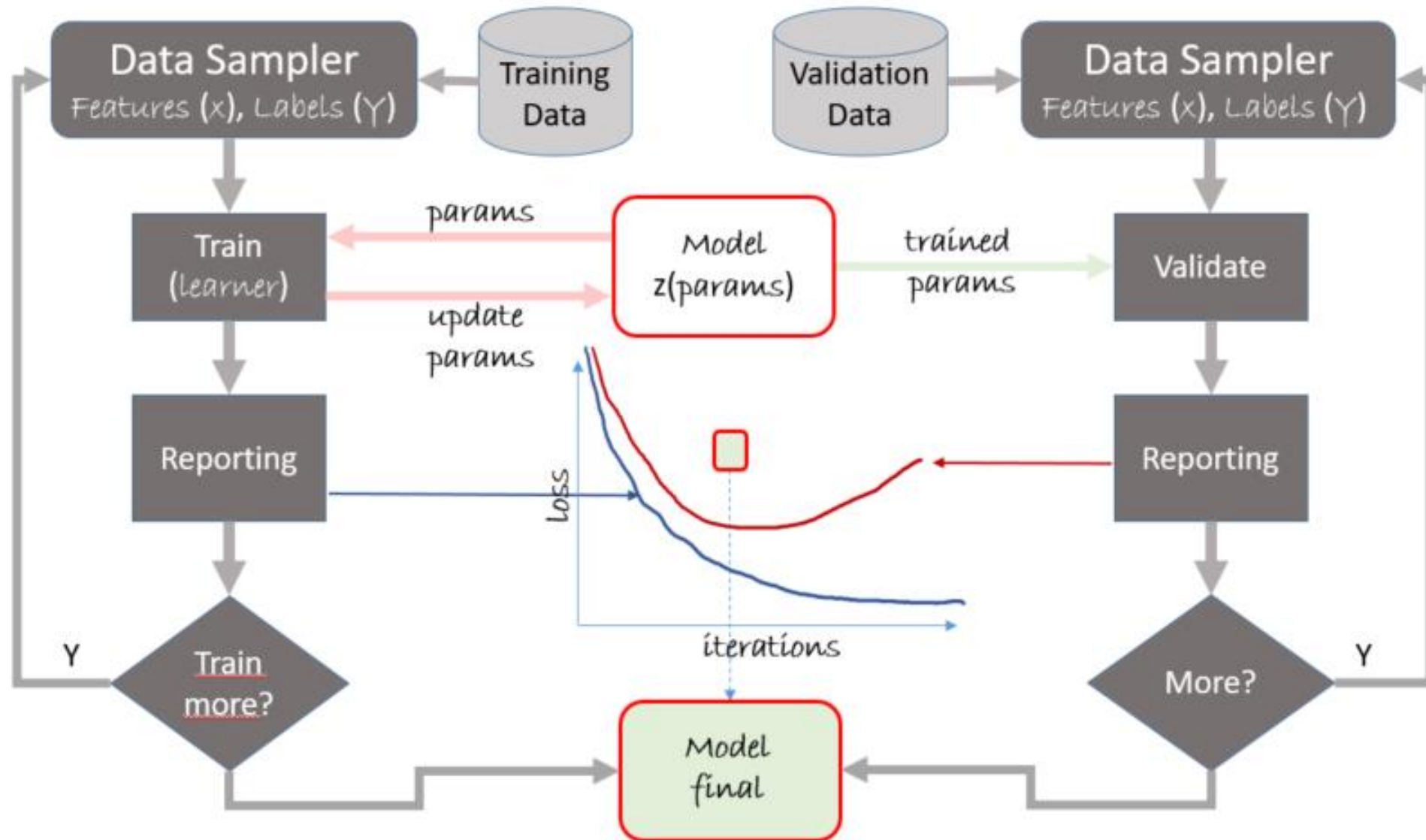
Output
(Y : $n \times$ future prediction)

```
z = create_model(x):  
    m = C.layers.Recurrence(C.layers.LSTM(TIMESTEPS))(x)  
    m = C.sequence.last(m)  
    m = C.layers.Dropout(0.2)(m)  
    m = C.layers.Dense(1)(m)  
    return m
```

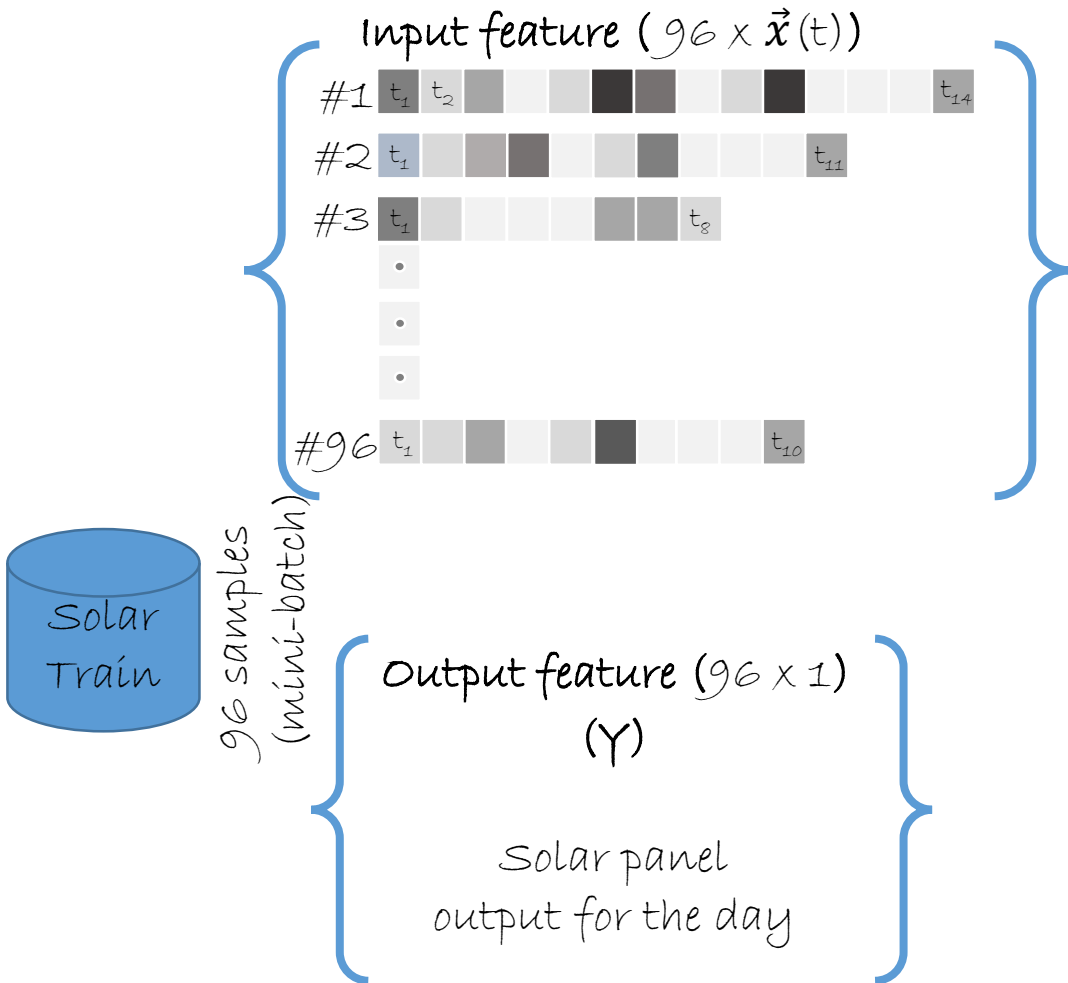
Input feature
(X : $n \times 14$ data pts)



Train / Validation Workflow



Train workflow



```
z = create_model(x):  
    m = C.layers.Recurrence(C.layers.LSTM(H_DIMS))(x)  
    m = C.sequence.last(m)  
    m = C.layers.Dropout(0.2)(m)  
    m = C.layers.Dense(1)(m)  
    return m
```



Loss
Error

squared_error(z, Y)

squared_error(z, Y)

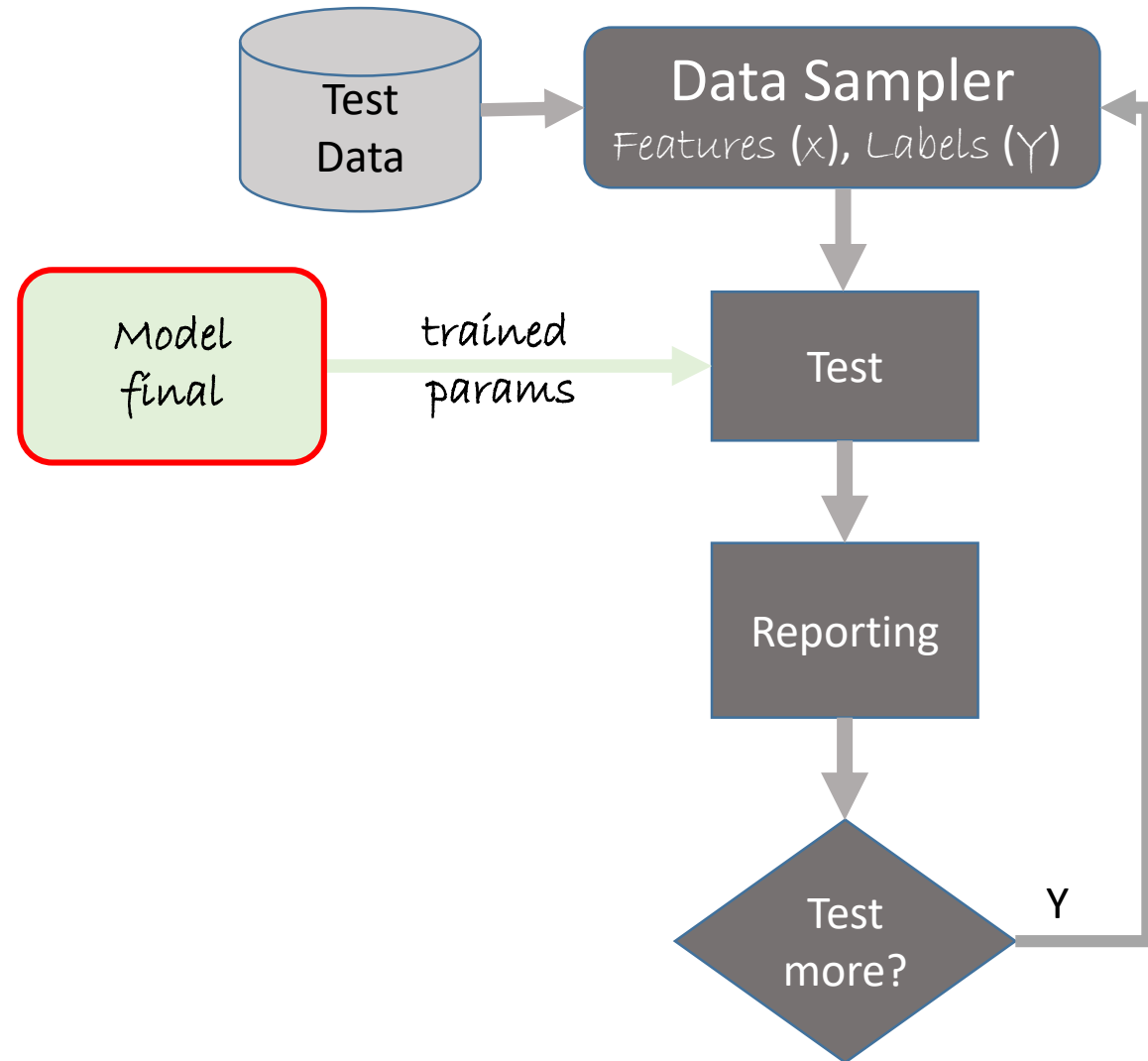
Trainer(model, (loss, error), learner)

Trainer.train_minibatch({X, Y})

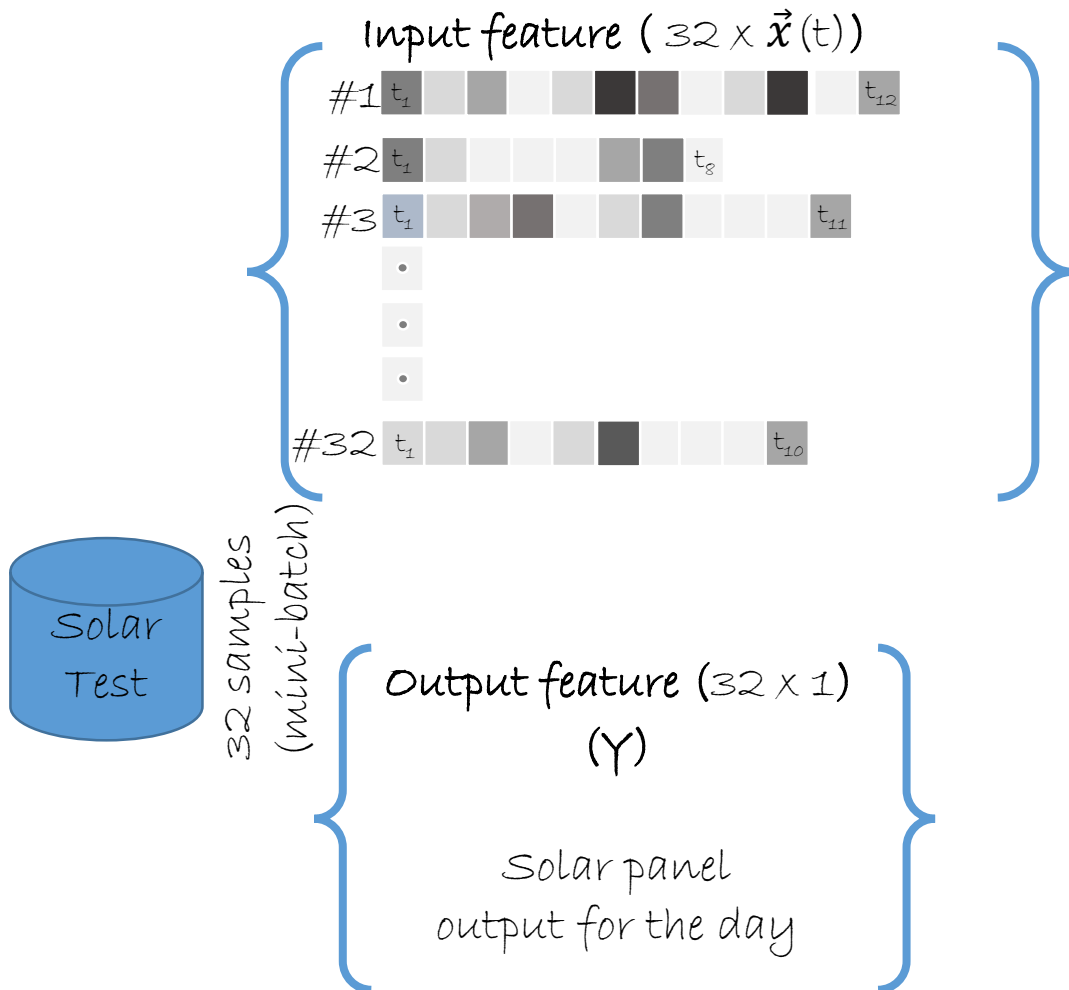
Learner

sgd, adagrad etc, are solvers to estimate

Test workflow



Test workflow



```
z = create_model(x):  
    m = C.layers.Recurrence(C.layers.LSTM(H_DIMS))(x)  
    m = C.sequence.last(m)  
    m = C.layers.Dropout(0.2)(m)  
    m = C.layers.Dense(1)(m)  
    return m
```

Trainer.**test**_minibatch($\{X, Y\}$)

Returns the squared error between the observed and predicted output from the solar panel

Prediction workflow

