



# 学习会

START



01

# 赤潮预测项目



## 模型参数

	原模型		当前模型
隐层单元数	5	➡	32
优化器	Nadam	➡	Adam
BatchSize	100	➡	128
Stateful	True	➡	False

```
In [176]: class Model(nn.Module):
            def __init__(self, input_size, hidden_size, output_dim):
                super(Model, self).__init__()
                self.lstm = nn.LSTM(input_size=input_size, hidden_size=hidden_size, batch_first=True)
            #         self.hidden = (torch.zeros(1, batch_size, hidden_size), torch.zeros(1, batch_size, hidden_size))
                self.linear = nn.Linear(hidden_size, output_dim)
            def forward(self, x):
                x, _ = self.lstm(x, None)
                last_output = x[:, -1, :]
                x = self.linear(last_output)
                out = F.relu(x, inplace=True)
                return out
```

```
In [154]: model = Model(1, 32, 3).to(device)
```

```
In [183]: criterion = nn.MSELoss().to(device)
            ops = optim.Adam(model.parameters(), lr=1e-5)
```

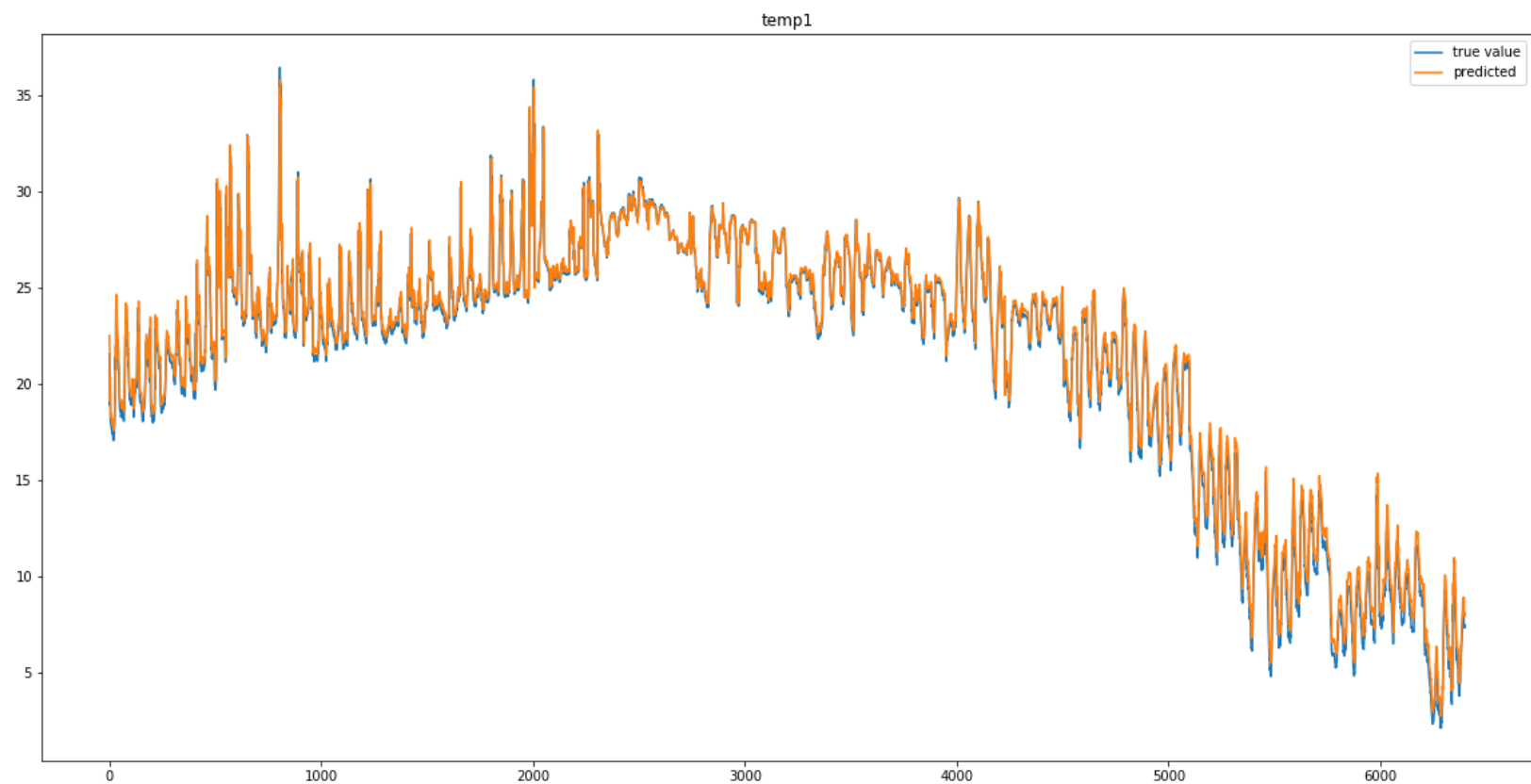


01

## 模型编写

```
for epoch in range(300):
    total_loss = 0.0
    for batch in range(num_batch):
        model.train()
        x = x_train[batch*batch_size:(batch+1)*batch_size]
        y = y_train[batch*batch_size:(batch+1)*batch_size]
        logits = model(x)
        loss = criterion(logits, y)
        ops.zero_grad()
        loss.backward()
        ops.step()
        total_loss += loss
    model.eval()
    y_pre = model(x_test)
    test_loss = criterion(y_pre, y_test)
    print('epoch:%d train_loss:%e, val_loss:%e' % (epoch, total_loss.item()/(len(x_train)/batch_size), test_loss.
```

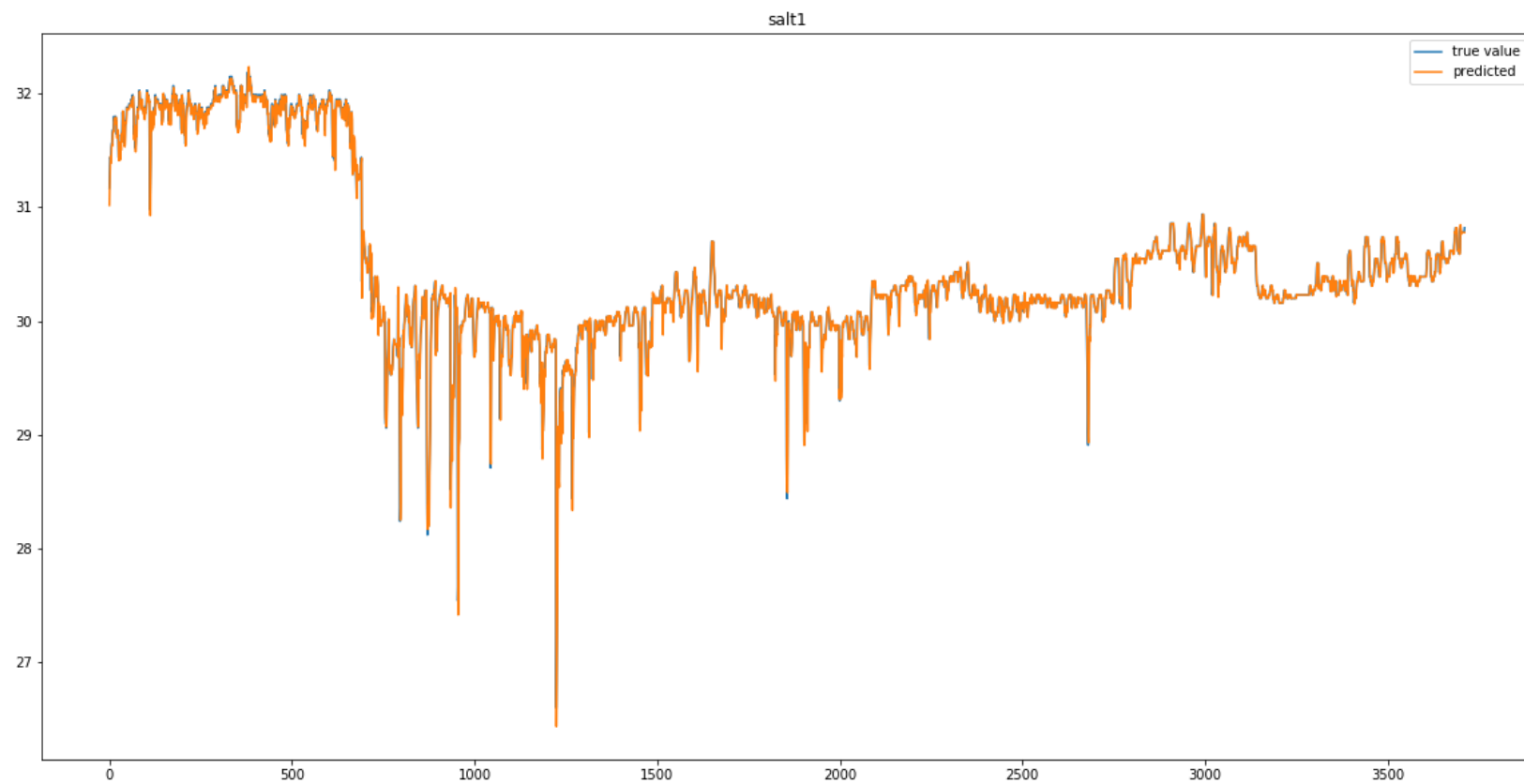




## 气温

1号13-17年训练18年预测

蓝色为真实值，橙色为预测值

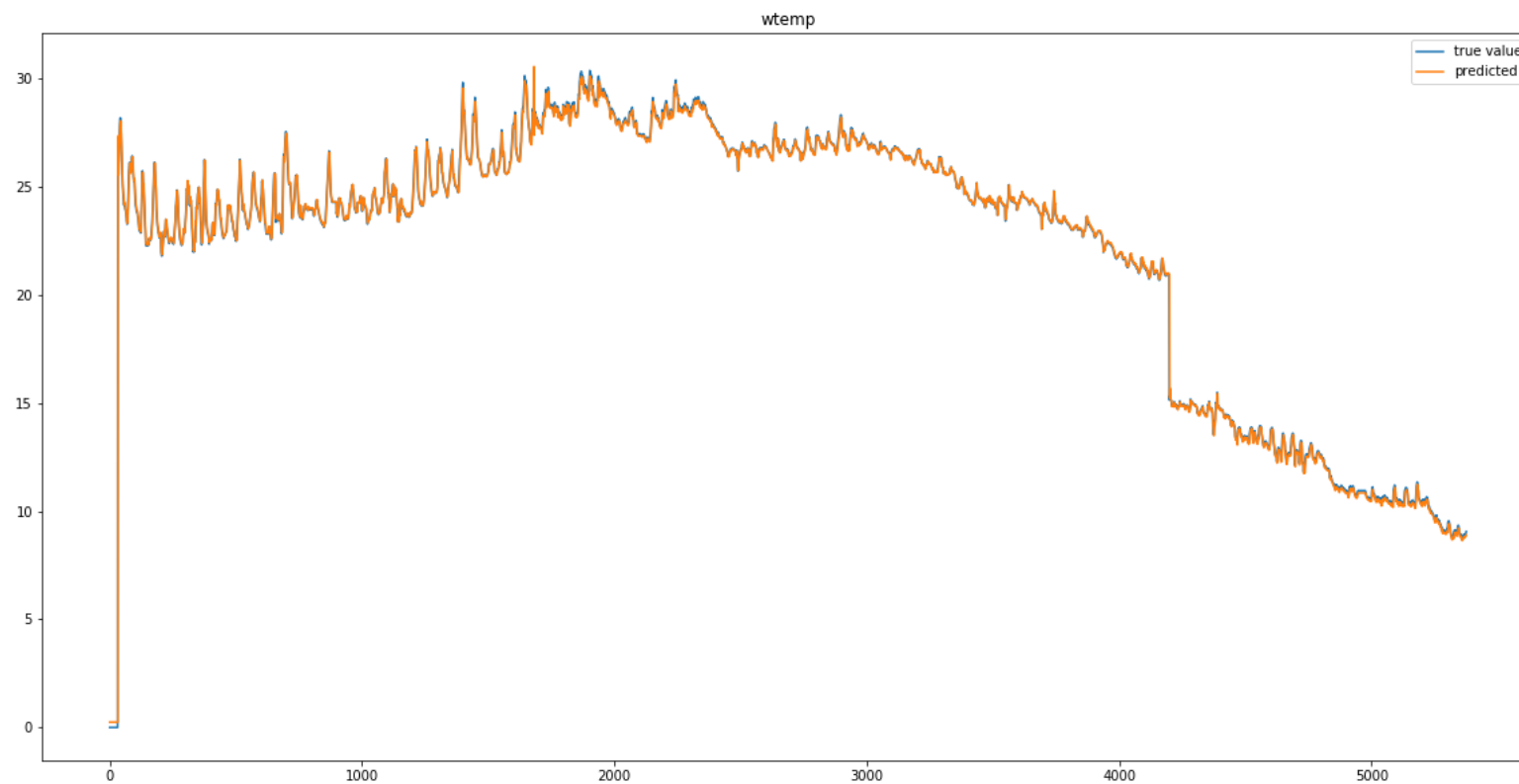


盐度

1号13-17训练18预测

蓝色为真实值，橙色为预测值



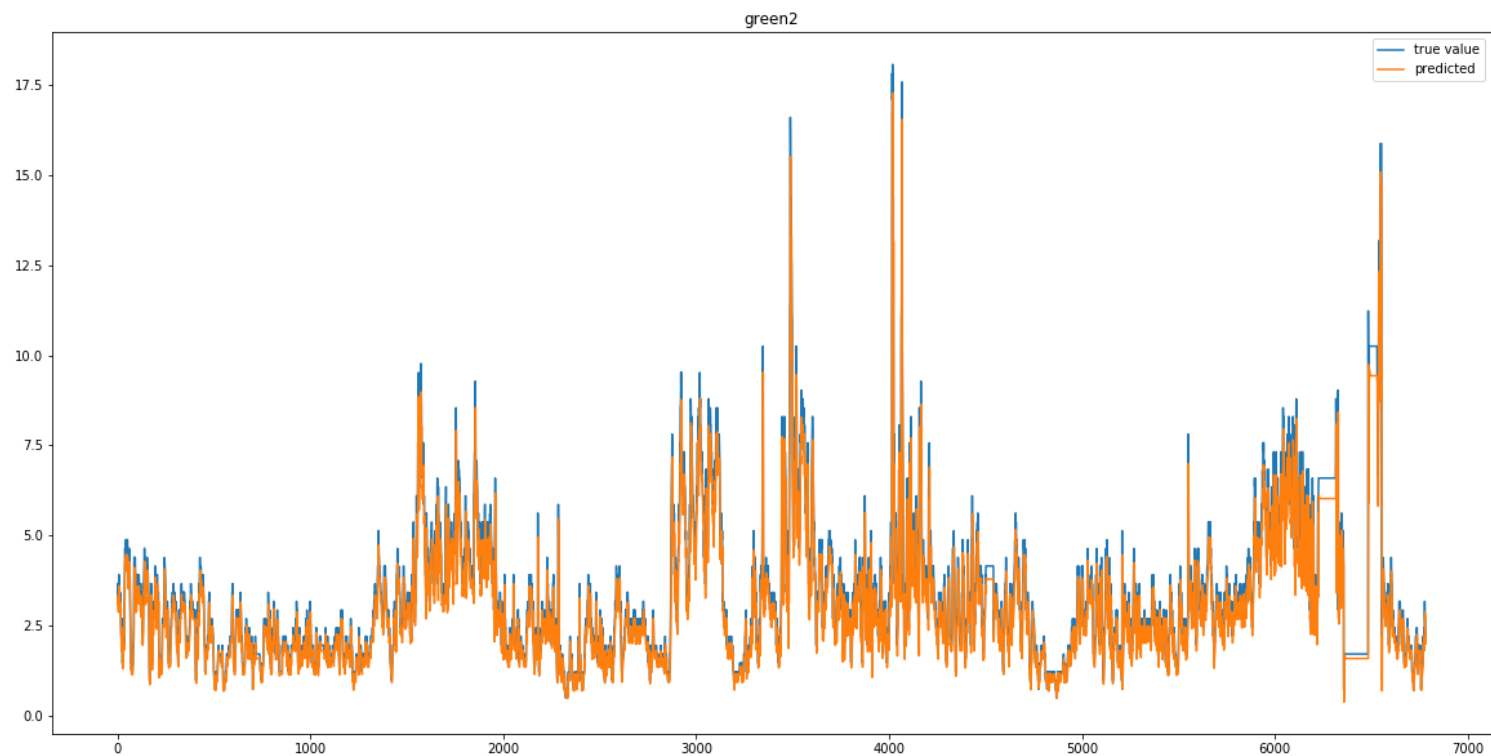


## 水温

1号13-17年训练18年预测

蓝色为真实值，橙色为预测值

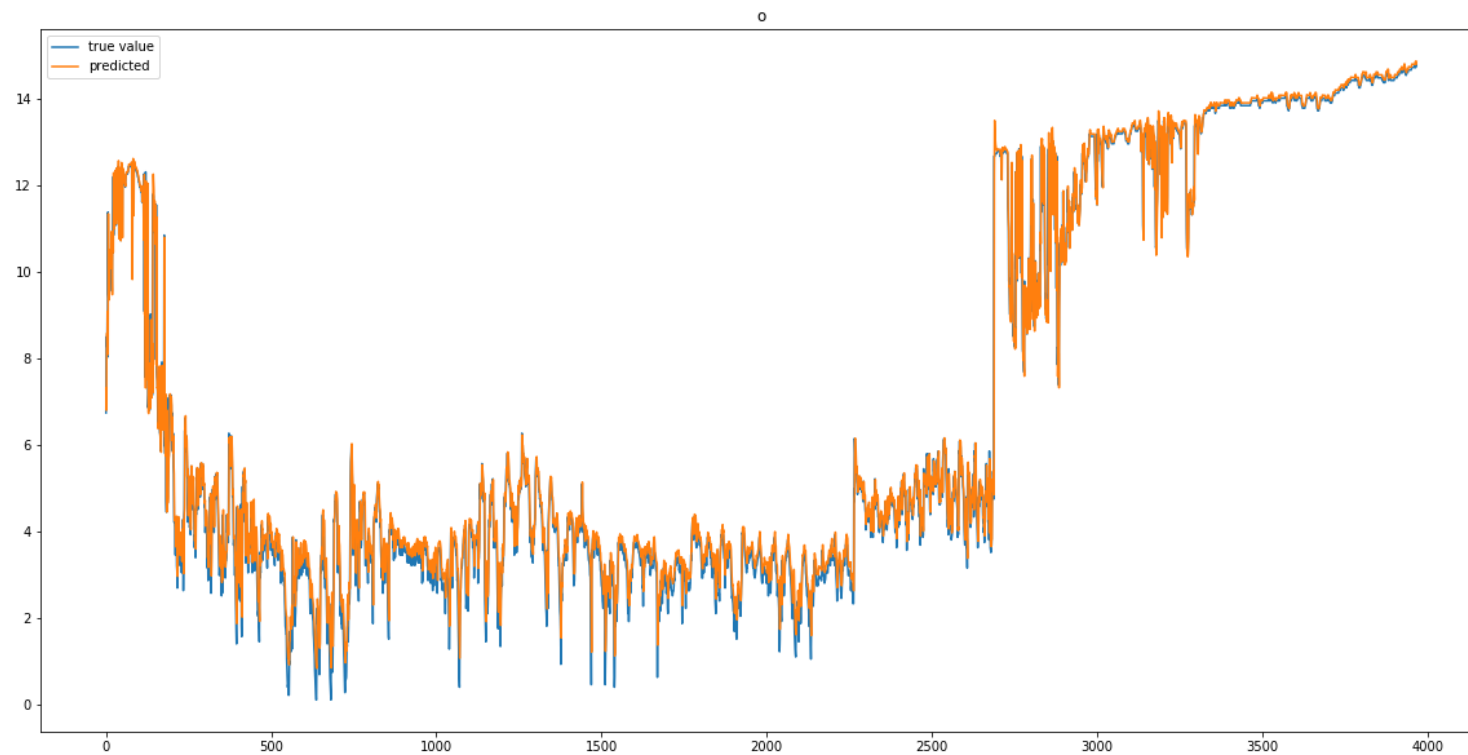




# 叶绿素

1号13-17训练18预测

蓝色为真实值，橙色为预测值



## 溶解氧

1号13-17训练18预测

蓝色为真实值，橙色为预测值



02

# 学习与分享



## **1.The Graph Neural Network Model**

## **2. Convolutional LSTM Network**



## The Graph Neural Network Model

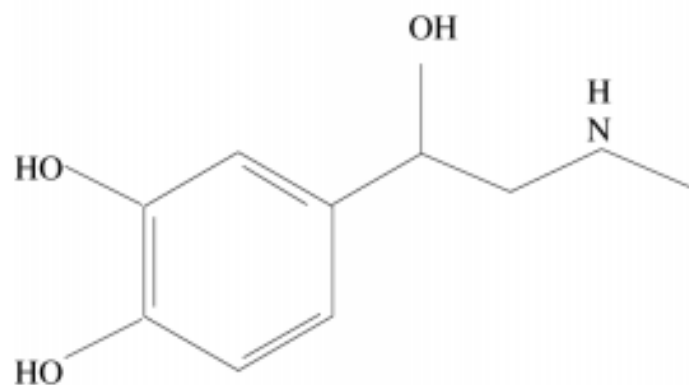
### I. INTRODUCTION

**D**ATA can be naturally represented by graph structures in several application areas, including proteomics [1], image analysis [2], scene description [3], [4], software engineering [5], [6], and natural language processing [7]. The simplest kinds of graph structures include single nodes and sequences. But in several applications, the information is organized in more complex graph structures such as trees, acyclic graphs, or cyclic graphs. Traditionally, data relationships exploitation has been the subject of many studies in the community of inductive logic programming and, recently, this research theme has been evolving in different directions [8], also because of the applications of relevant concepts in statistics and neural networks to such areas (see, for example, the recent workshops [9]–[12]).

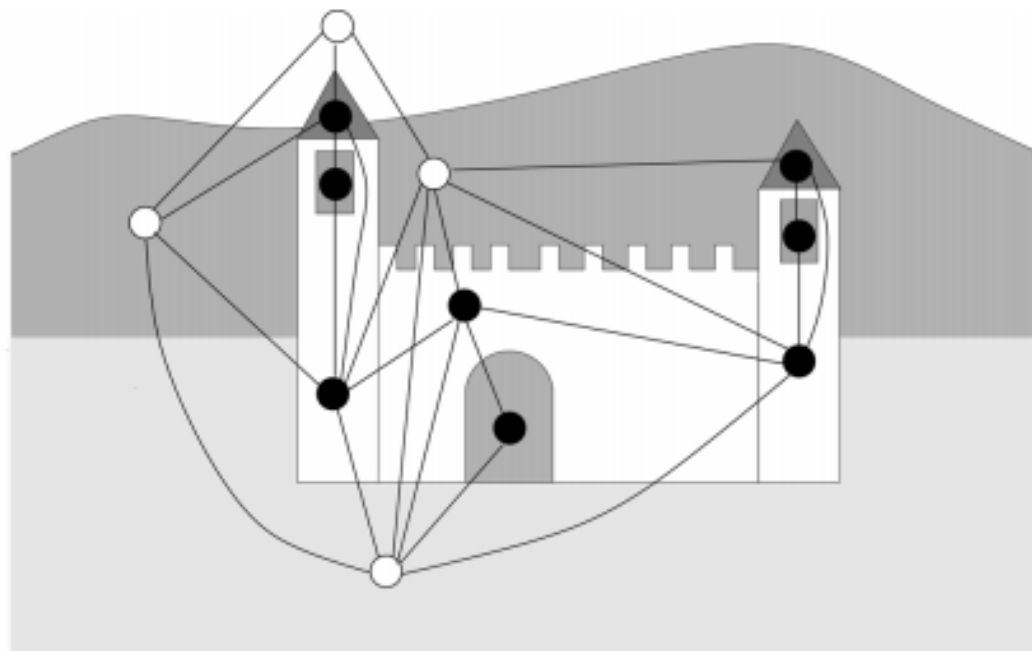




## The Graph Neural Network Model



(a)



(b)







## The Graph Neural Network Model

$$\begin{aligned}\mathbf{x}_n &= f_{\mathbf{w}}(\mathbf{l}_n, \mathbf{l}_{\text{co}[n]}, \mathbf{x}_{\text{ne}[n]}, \mathbf{l}_{\text{ne}[n]}) \\ \mathbf{o}_n &= g_{\mathbf{w}}(\mathbf{x}_n, \mathbf{l}_n)\end{aligned}\tag{1}$$

where  $\mathbf{l}_n$ ,  $\mathbf{l}_{\text{co}[n]}$ ,  $\mathbf{x}_{\text{ne}[n]}$ , and  $\mathbf{l}_{\text{ne}[n]}$  are the label of  $n$ , the labels of its edges, the states, and the labels of the nodes in the neighborhood of  $n$ , respectively.

Let  $\mathbf{x}$ ,  $\mathbf{o}$ ,  $\mathbf{l}$ , and  $\mathbf{l}_N$  be the vectors constructed by stacking all the states, all the outputs, all the labels, and all the node labels, respectively. Then, (1) can be rewritten in a compact form as

$$\begin{aligned}\mathbf{x} &= F_{\mathbf{w}}(\mathbf{x}, \mathbf{l}) \\ \mathbf{o} &= G_{\mathbf{w}}(\mathbf{x}, \mathbf{l}_N)\end{aligned}\tag{2}$$







## The Graph Neural Network Model

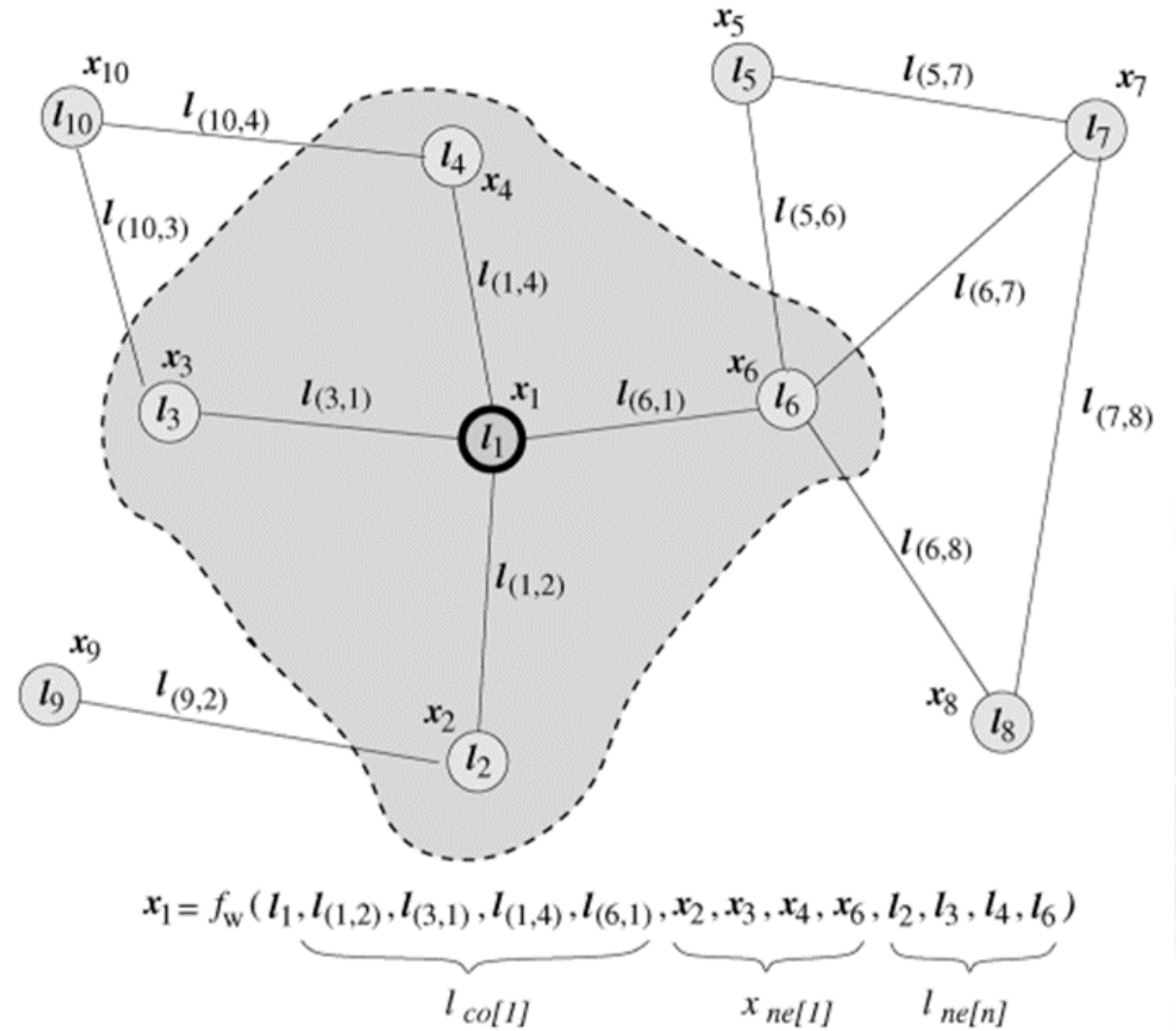
Note that (1) makes it possible to process both positional and nonpositional graphs. For positional graphs,  $f_{\mathbf{w}}$  must receive the positions of the neighbors as additional inputs. In practice, this can be easily achieved provided that information contained in  $\mathbf{x}_{\text{ne}[n]}$ ,  $\mathbf{l}_{\text{co}[n]}$ , and  $\mathbf{l}_{\text{ne}[n]}$  is sorted according to neighbors' positions and is properly padded with special null values in positions corresponding to nonexisting neighbors. For example,  $\mathbf{x}_{\text{ne}[n]} = [\mathbf{y}_1, \dots, \mathbf{y}_M]$ , where  $M = \max_{n,u} \nu_n(u)$  is the maximal number of neighbors of a node;  $\mathbf{y}_i = \mathbf{x}_u$  holds, if  $u$  is the  $i$ th neighbor of  $n$  ( $\nu_n(u) = i$ ); and  $\mathbf{y}_i = \mathbf{x}_0$ , for some predefined null state  $\mathbf{x}_0$ , if there is no  $i$ th neighbor.

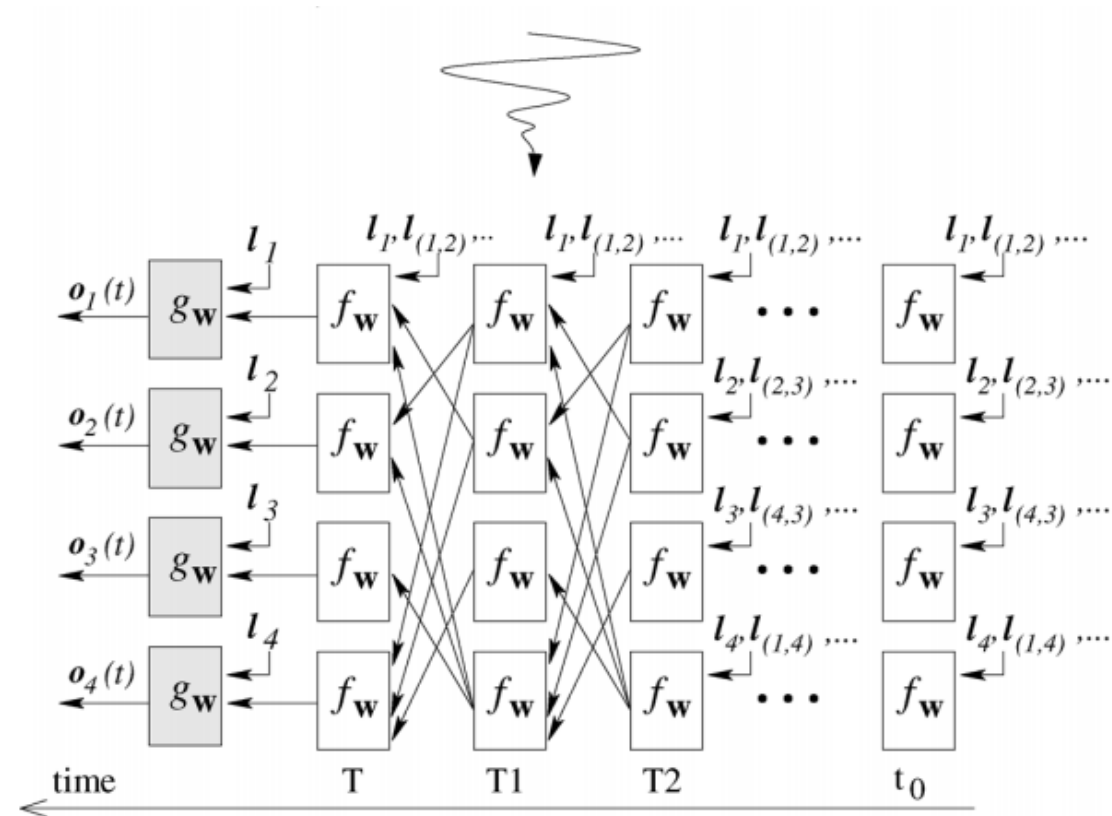
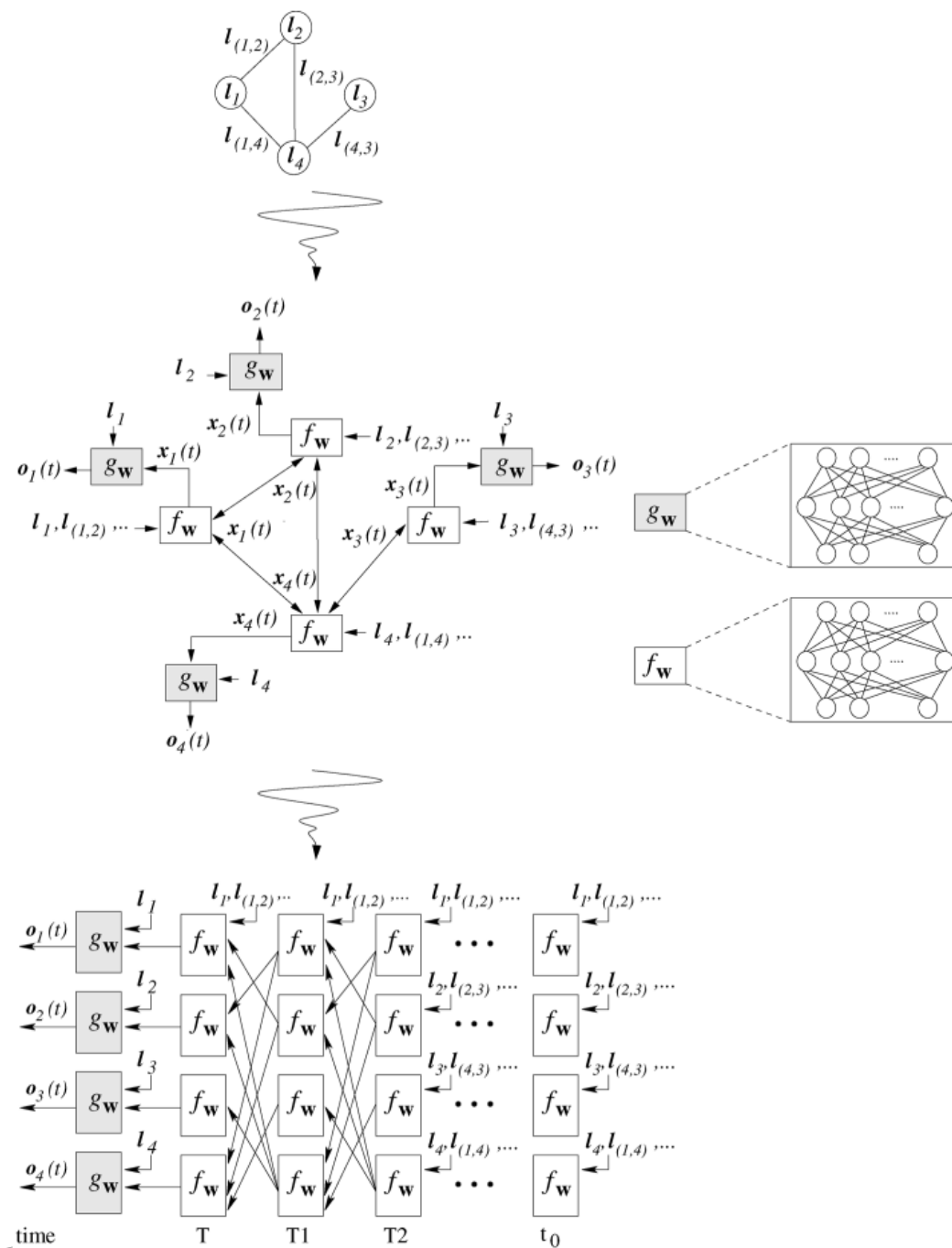
However, for nonpositional graphs, it is useful to replace function  $f_{\mathbf{w}}$  of (1) with

$$\mathbf{x}_n = \sum_{u \in \text{ne}[n]} h_{\mathbf{w}}(\mathbf{l}_n, \mathbf{l}_{(n,u)}, \mathbf{x}_u, \mathbf{l}_u), \quad n \in \mathbf{N} \quad (3)$$



# The Graph Neural Network Model







## The Graph Neural Network Model

### B. Computation of the State

Banach's fixed point theorem [53] does not only ensure the existence and the uniqueness of the solution of (1) but it also suggests the following classic iterative scheme for computing the state:

$$\mathbf{x}(t+1) = F_{\mathbf{w}}(\mathbf{x}(t), \mathbf{l}) \quad (4)$$

where  $\mathbf{x}(t)$  denotes the  $t$ th iteration of  $\mathbf{x}$ . The dynamical system (4) converges exponentially fast to the solution of (2) for any initial value  $\mathbf{x}(0)$ . We can, therefore, think of  $\mathbf{x}(t)$  as the state that is updated by the transition function  $F_{\mathbf{w}}$ . In fact, (4) implements the Jacobi iterative method for solving nonlinear equations [55]. Thus, the outputs and the states can be computed by iterating

$$\begin{aligned} \mathbf{x}_n(t+1) &= f_{\mathbf{w}}(\mathbf{l}_n, \mathbf{l}_{\text{co}[n]}, \mathbf{x}_{\text{ne}[n]}(t), \mathbf{l}_{\text{ne}[n]}) \\ \mathbf{o}_n(t) &= g_{\mathbf{w}}(\mathbf{x}_n(t), \mathbf{l}_n), \quad n \in \mathbf{N}. \end{aligned} \quad (5)$$





1) *Linear (nonpositional) GNN*. Equation (3) can naturally be implemented by

$$h_{\mathbf{w}}(\mathbf{l}_n, \mathbf{l}_{(n,u)}, \mathbf{x}_u, \mathbf{l}_u) = \mathbf{A}_{n,u} \mathbf{x}_u + \mathbf{b}_n \quad (12)$$

where the vector  $\mathbf{b}_n \in \mathbb{R}^s$  and the matrix  $\mathbf{A}_{n,u} \in \mathbb{R}^{s \times s}$  are defined by the output of two feedforward neural networks (FNNs), whose parameters correspond to the parameters of the GNN. More precisely, let us call *transition network* an FNN that has to generate  $\mathbf{A}_{n,u}$  and *forcing network* another FNN that has to generate  $\mathbf{b}_n$ . Moreover, let  $\phi_{\mathbf{w}} : \mathbb{R}^{2l_N + l_E} \rightarrow \mathbb{R}^{s^2}$  and  $\rho_{\mathbf{w}} : \mathbb{R}^{l_N} \rightarrow \mathbb{R}^s$  be the functions implemented by the transition and the forcing network, respectively. Then, we define

$$\mathbf{A}_{n,u} = \frac{\mu}{s|\text{ne}[u]|} \cdot \Xi \quad (13)$$

$$\mathbf{b}_n = \rho_{\mathbf{w}}(\mathbf{l}_n) \quad (14)$$

where  $\mu \in (0, 1)$  and  $\Xi = \text{resize}(\phi_{\mathbf{w}}(\mathbf{l}_n, \mathbf{l}_{(n,u)}, \mathbf{l}_u))$  hold, and  $\text{resize}(\cdot)$  denotes the operator that allocates the elements of a  $s^2$ -dimensional vector into as  $s \times s$  matrix. Thus,  $\mathbf{A}_{n,u}$  is obtained by arranging the outputs of the transition network into the square matrix  $\Xi$  and by multiplication with the factor  $\mu/s|\text{ne}[u]|$ . On the other hand,  $\mathbf{b}_n$  is just a vector that contains the outputs of the forcing network. Here, it is further assumed that  $\|\phi_{\mathbf{w}}(\mathbf{l}_n, \mathbf{l}_{(n,u)}, \mathbf{l}_u)\|_1 \leq s$  holds<sup>7</sup>; this can be straightforwardly verified if the output neurons of the transition network use an appropriately

bounded activation function, e.g., a hyperbolic tangent. Note that in this case  $F_{\mathbf{w}}(\mathbf{x}, \mathbf{l}) = \mathbf{A}\mathbf{x} + \mathbf{b}$ , where  $\mathbf{b}$  is the vector constructed by stacking all the  $\mathbf{b}_n$ , and  $\mathbf{A}$  is a block matrix  $\{\bar{\mathbf{A}}_{n,u}\}$ , with  $\bar{\mathbf{A}}_{n,u} = \mathbf{A}_{n,u}$  if  $u$  is a neighbor of  $n$  and  $\bar{\mathbf{A}}_{n,u} = \mathbf{0}$  otherwise. Moreover, vectors  $\mathbf{b}_n$  and matrices  $\mathbf{A}_{n,u}$  do not depend on the state  $\mathbf{x}$ , but only on node and edge labels. Thus,  $\partial F_{\mathbf{w}}/\partial \mathbf{x} = \mathbf{A}$ , and, by simple algebra

$$\begin{aligned} \left\| \frac{\partial F_{\mathbf{w}}}{\partial \mathbf{x}} \right\|_1 &= \|\mathbf{A}\|_1 \leq \max_{u \in \mathbf{N}} \left( \sum_{n \in \text{ne}[u]} \|\mathbf{A}_{n,u}\|_1 \right) \\ &\leq \max_{u \in \mathbf{N}} \left( \frac{\mu}{s|\text{ne}[u]|} \cdot \sum_{n \in \text{ne}[u]} \|\Xi\|_1 \right) \leq \mu \end{aligned}$$

which implies that  $F_{\mathbf{w}}$  is a contraction map (w.r.t.  $\|\cdot\|_1$ ) for any set of parameters  $\mathbf{w}$ .

2) *Nonlinear (nonpositional) GNN*. In this case,  $h_{\mathbf{w}}$  is realized by a multilayered FNN. Since three-layered neural networks are universal approximators [67],  $h_{\mathbf{w}}$  can approximate any desired function. However, not all the parameters  $\mathbf{w}$  can be used, because it must be ensured that the corresponding transition function  $F_{\mathbf{w}}$  is a contraction map. This can be achieved by adding a penalty term to (6), i.e.,

$$e_{\mathbf{w}} = \sum_{i=1}^p \sum_{j=1}^{q_i} (t_{i,j} - \varphi_{\mathbf{w}}(\mathbf{G}_i, n_{i,j}))^2 + \beta L \left( \left\| \frac{\partial F_{\mathbf{w}}}{\partial \mathbf{x}} \right\| \right)$$

where the penalty term  $L(y)$  is  $(y - \mu)^2$  if  $y > \mu$  and 0 otherwise, and the parameter  $\mu \in (0, 1)$  defines the desired contraction constant of  $F_{\mathbf{w}}$ . More generally, the penalty term can be any expression, differentiable w.r.t.  $\mathbf{w}$ , that is monotone increasing w.r.t. the norm of the Jacobian. For example, in our experiments, we use the penalty term  $p_{\mathbf{w}} = \sum_{i=1}^s L(\|\mathbf{A}^i\|_1)$ , where  $\mathbf{A}^i$  is the  $i$ th column of  $\partial F_{\mathbf{w}}/\partial \mathbf{x}$ . In fact, such an expression is an approximation of  $L(\|\partial F_{\mathbf{w}}/\partial \mathbf{x}\|_1) = L(\max_i \|\mathbf{A}^i\|_1)$ .



MAIN

```

initialize  $w$ ;
 $x = \text{Forward}(w)$ ;
repeat
     $\frac{\partial e_w}{\partial w} = \text{BACKWARD}(x, w)$ ;
     $w = w - \lambda \cdot \frac{\partial e_w}{\partial w}$ ;
     $x = \text{FORWARD}(w)$ ;
until (a stopping criterion);
return  $w$ ;

```

end

FORWARD( $w$ )

```

initialize  $x(0)$ ,  $t = 0$ ;
repeat
     $x(t+1) = F_w(x(t), l)$ ;
     $t = t + 1$ ;
until  $\|x(t) - x(t-1)\| \leq \varepsilon_f$ 
return  $x(t)$ ;

```

end

BACKWARD( $x, w$ )

```

 $o = G_w(x, l_N)$ ;
 $A = \frac{\partial F_w}{\partial x}(x, l)$ ;
 $b = \frac{\partial e_w}{\partial o} \cdot \frac{\partial G_w}{\partial x}(x, l_N)$ ;
initialize  $z(0)$ ,  $t=0$ ;
repeat
     $z(t) = z(t+1) \cdot A + b$ ;
     $t = t - 1$ ;
until  $\|z(t-1) - z(t)\| \leq \varepsilon_b$ ;
 $c = \frac{\partial e_w}{\partial o} \cdot \frac{\partial G_w}{\partial w}(x, l_N)$ ;
 $d = z(t) \cdot \frac{\partial F_w}{\partial w}(x, l)$ ;
 $\frac{\partial e_w}{\partial w} = c + d$ ;
return  $\frac{\partial e_w}{\partial w}$ ;

```

end

Learning in GNNs consists of estimating the parameter  $w$  such that  $\varphi_w$  approximates the data in the learning data set

$$\mathcal{L} = \{(\mathbf{G}_i, n_{i,j}, \mathbf{t}_{i,j}) \mid \mathbf{G}_i = (\mathbf{N}_i, \mathbf{E}_i) \in \mathcal{G}; \\ n_{i,j} \in \mathbf{N}_i; \mathbf{t}_{i,j} \in \mathbb{R}^m, 1 \leq i \leq p, 1 \leq j \leq q_i\}$$

where  $q_i$  is the number of supervised nodes in  $\mathbf{G}_i$ . For graph-focused tasks, one special node is used for the target ( $q_i = 1$  holds), whereas for node-focused tasks, in principle, the supervision can be performed on every node. The learning task can be posed as the minimization of a quadratic cost function

$$e_w = \sum_{i=1}^p \sum_{j=1}^{q_i} (\mathbf{t}_{i,j} - \varphi_w(\mathbf{G}_i, n_{i,j}))^2. \quad (6)$$

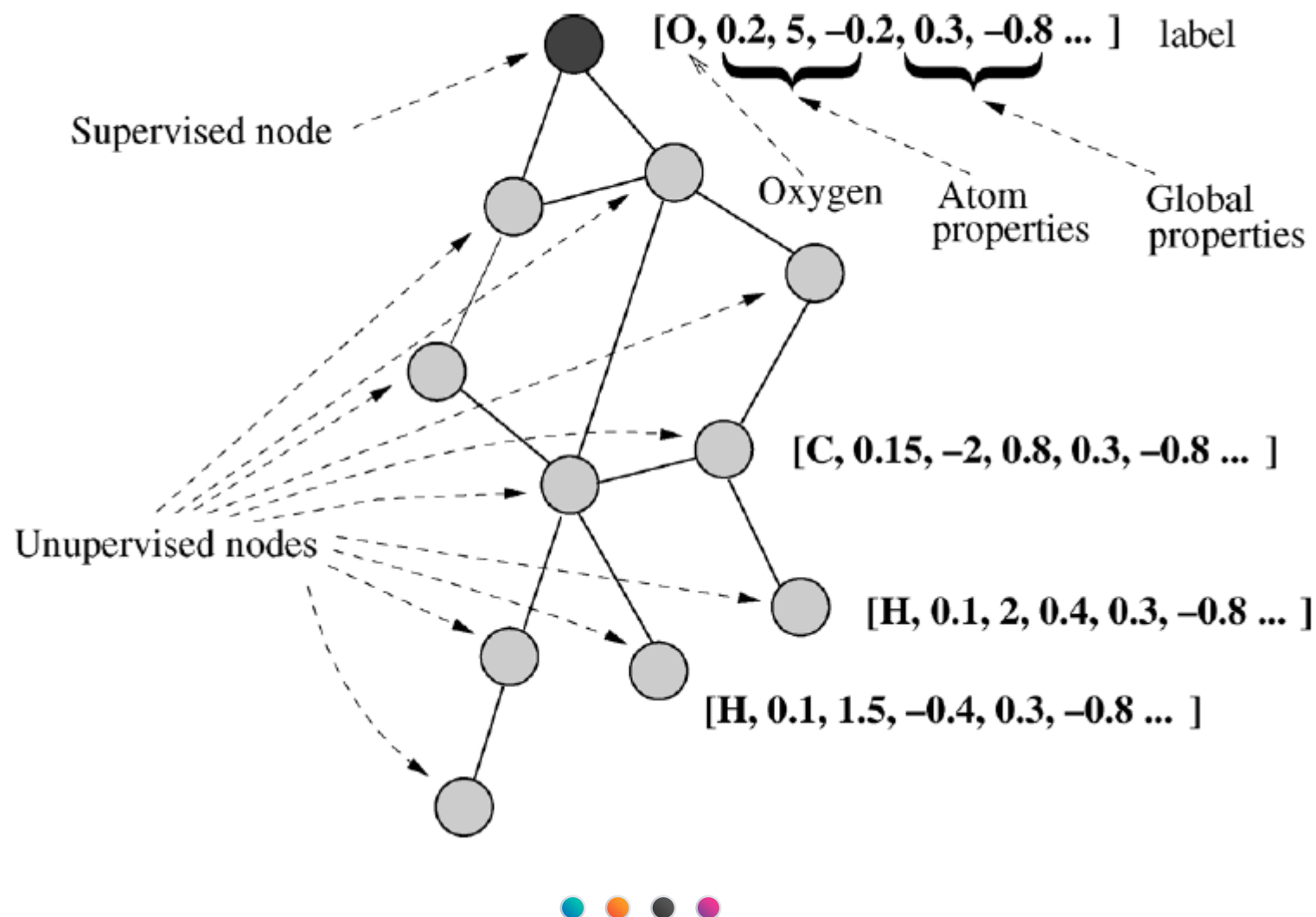
The learning algorithm is based on a gradient-descent strategy and is composed of the following steps.

- The states  $\mathbf{x}_n(t)$  are iteratively updated by (5) until at time  $T$  they approach the fixed point solution of (2):  $\mathbf{x}(T) \approx \mathbf{x}$ .
- The gradient  $\partial e_w(T)/\partial w$  is computed.
- The weights  $w$  are updated according to the gradient computed in step b).





## The Graph Neural Network Model







## Convolutional LSTM Network

$$\begin{aligned}i_t &= \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i) \\f_t &= \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f) \\c_t &= f_t \circ c_{t-1} + i_t \circ \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \\o_t &= \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co} \circ c_t + b_o) \\h_t &= o_t \circ \tanh(c_t)\end{aligned}\tag{2}$$



# Convolutional LSTM Network

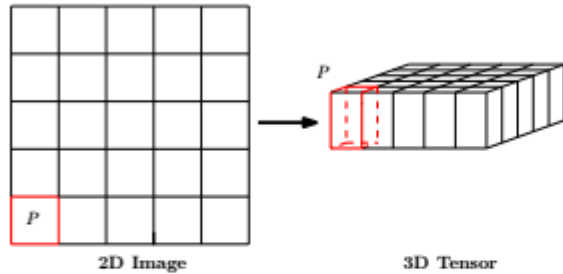


Figure 1: Transforming 2D image into 3D tensor

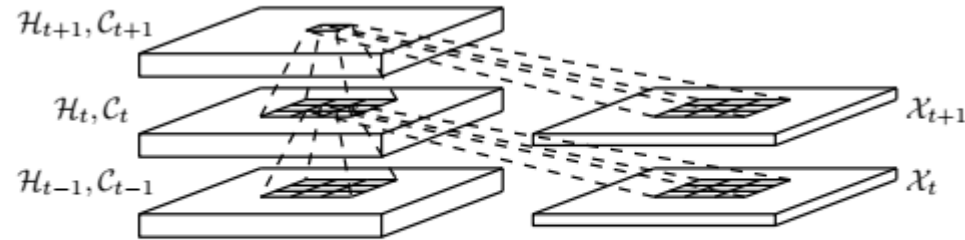


Figure 2: Inner structure of ConvLSTM

$$i_t = \sigma(W_{xi} * \mathcal{X}_t + W_{hi} * \mathcal{H}_{t-1} + W_{ci} \circ \mathcal{C}_{t-1} + b_i)$$

$$f_t = \sigma(W_{xf} * \mathcal{X}_t + W_{hf} * \mathcal{H}_{t-1} + W_{cf} \circ \mathcal{C}_{t-1} + b_f)$$

$$\mathcal{C}_t = f_t \circ \mathcal{C}_{t-1} + i_t \circ \tanh(W_{xc} * \mathcal{X}_t + W_{hc} * \mathcal{H}_{t-1} + b_c)$$

$$o_t = \sigma(W_{xo} * \mathcal{X}_t + W_{ho} * \mathcal{H}_{t-1} + W_{co} \circ \mathcal{C}_t + b_o)$$

$$\mathcal{H}_t = o_t \circ \tanh(\mathcal{C}_t)$$



## Convolutional LSTM Network

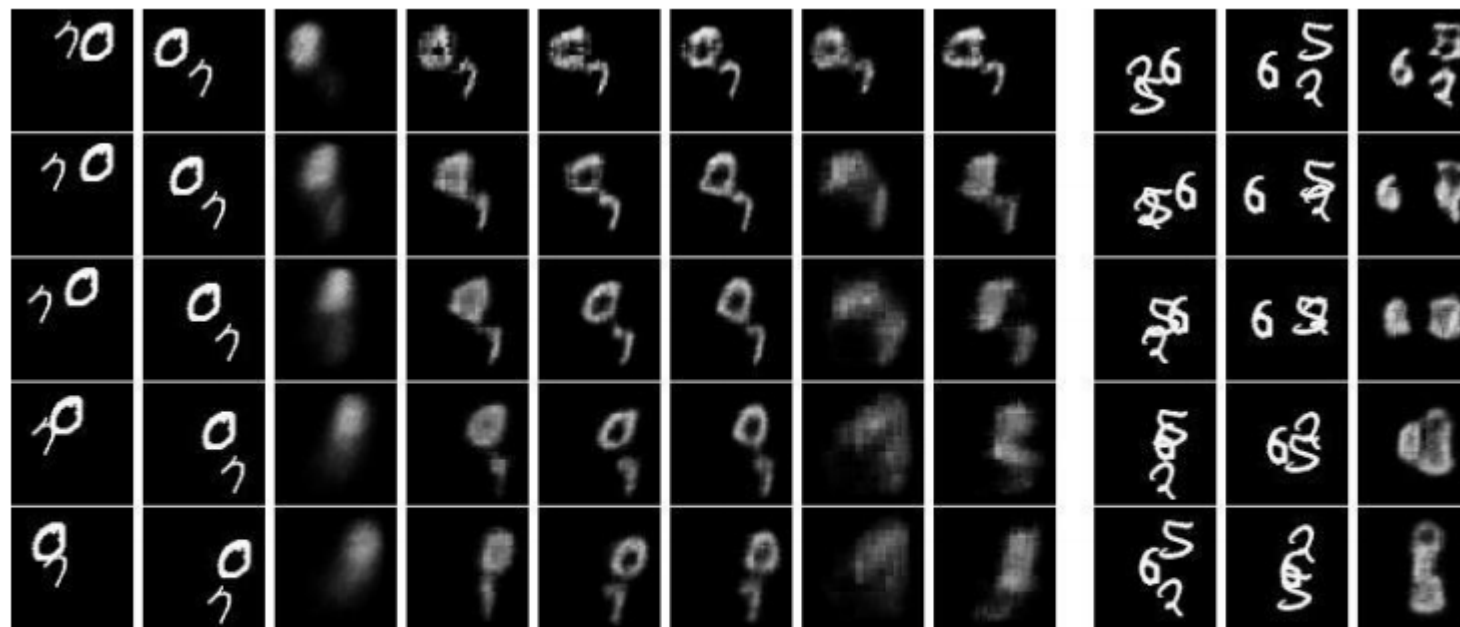


Figure 4: **Left part** is an example showing the prediction results of different models. From left to right: input frames; ground truth; FC-LSTM; 1-layer; 2-layer; 3-layer; 2-layer without convolution; 3-layer without convolution. **Right part** is an example showing an “out-of-domain” run. From left to right: input frames; ground truth; 3-layer. All sequences are sampled with an interval of 2.



## Convolutional LSTM Network

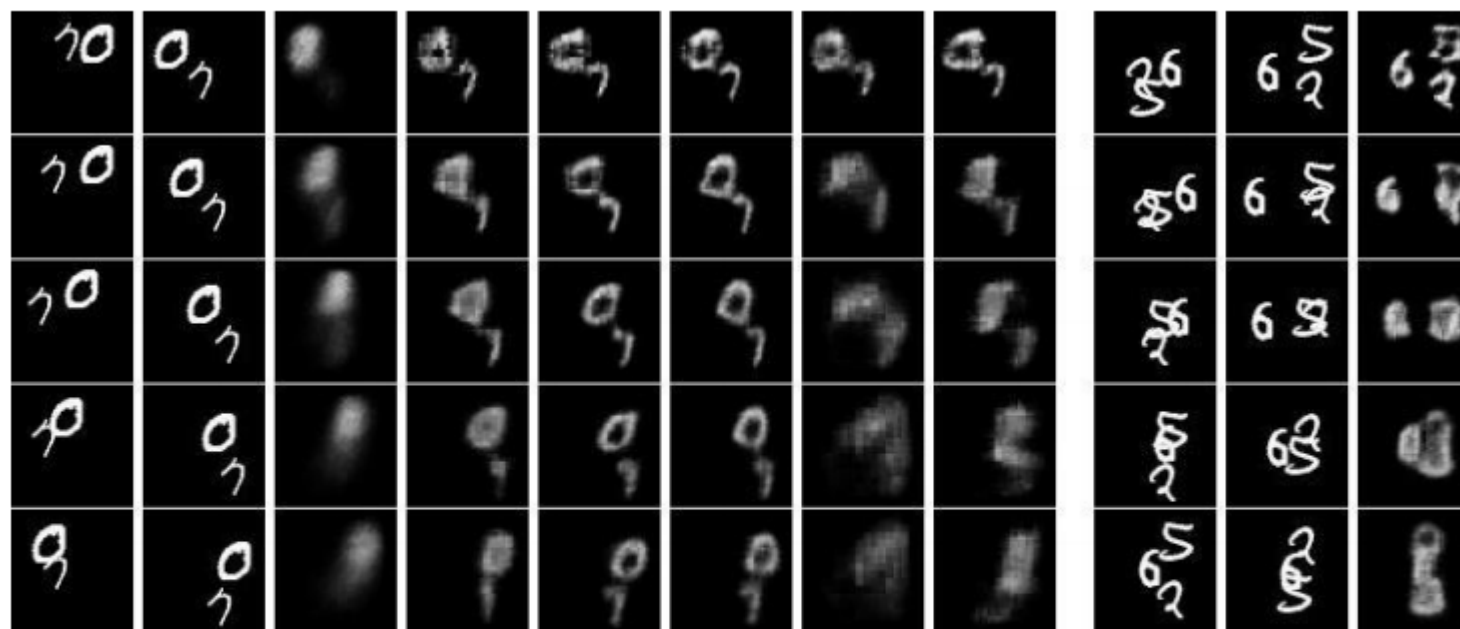
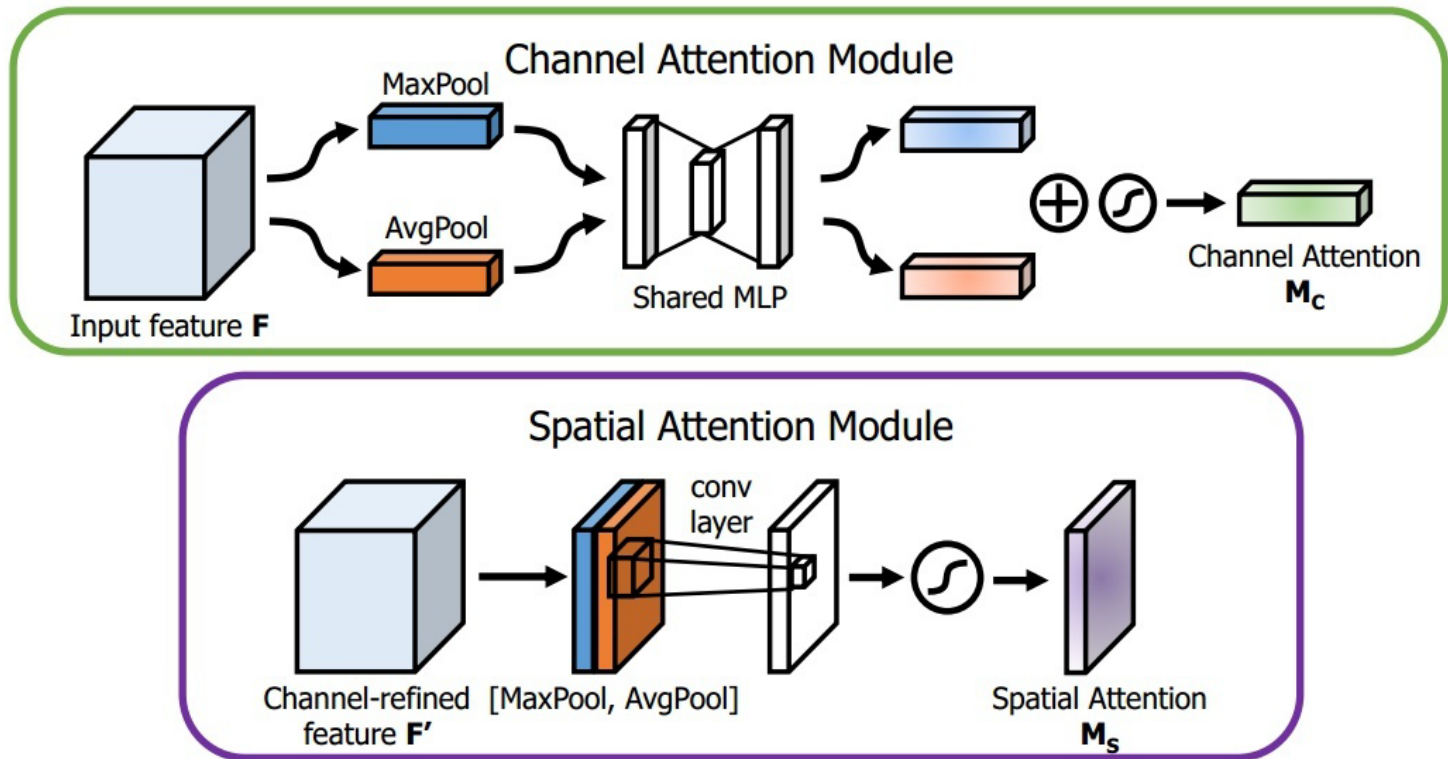


Figure 4: **Left part** is an example showing the prediction results of different models. From left to right: input frames; ground truth; FC-LSTM; 1-layer; 2-layer; 3-layer; 2-layer without convolution; 3-layer without convolution. **Right part** is an example showing an “out-of-domain” run. From left to right: input frames; ground truth; 3-layer. All sequences are sampled with an interval of 2.



## Stand-Alone Self-Attention in Vision Models

卷积捕获长距离交互能力比较差，这是因为卷积的感受野大时缩放特性弱。此前，针对这个问题的处理方法就是用注意力机制。在卷积中运用的注意力机制主要有两种。一个是，基于通道的注意力机制；一个是，基于空间的注意力机制。这些方法有个特点，就是用全局注意力层作为卷积的附加模块；还有一个限制，因为关注输入的所有位置，要求输入比较小，否则计算成本太大。





## Stand-Alone Self-Attention in Vision Models

The count of parameters in CNN

$$K_W * K_H * N_F * N_I$$

The count of parameters in SASA

$$1 * 1 * N_F * N_I * 3$$

The computational cost of attention also grows slower with spatial extent compared to convolution with typical values of  $d_{in}$  and  $d_{out}$ . For example, if  $d_{in} = d_{out} = 128$ , a convolution layer with  $k = 3$  has the same computational cost as an attention layer with  $k = 19$ .

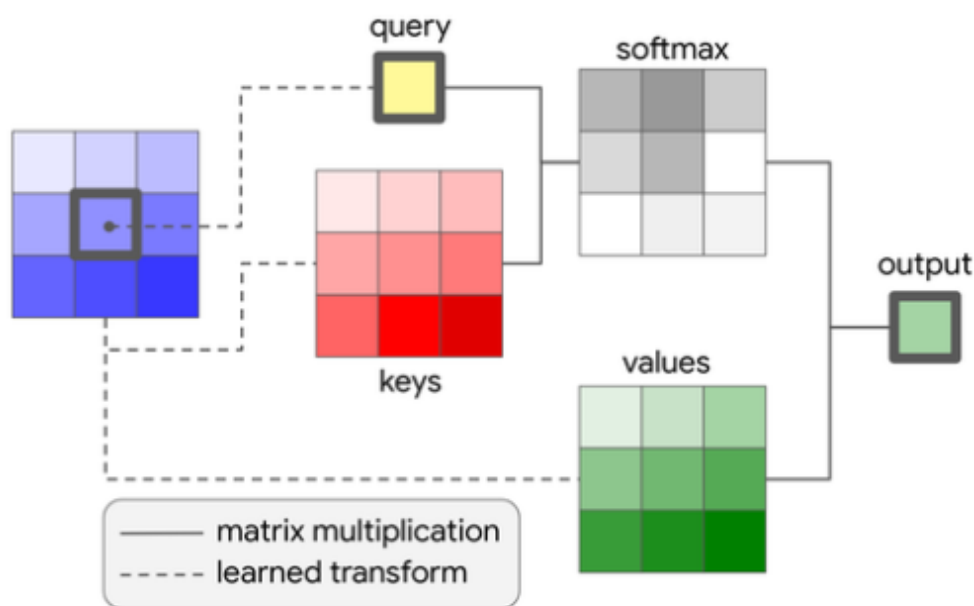


Figure 3: An example of a local attention layer over spatial extent of  $k = 3$ .

-1, -1	-1, 0	-1, 1	-1, 2
0, -1	0, 0	0, 1	0, 2
1, -1	1, 0	1, 1	1, 2
2, -1	2, 0	2, 1	2, 2

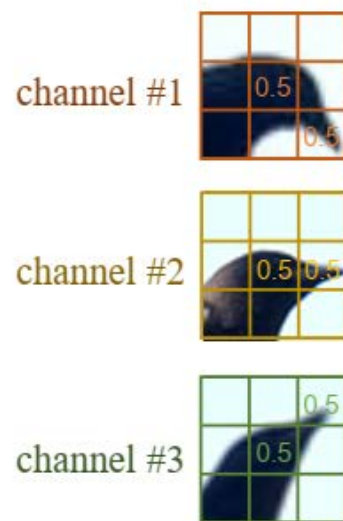
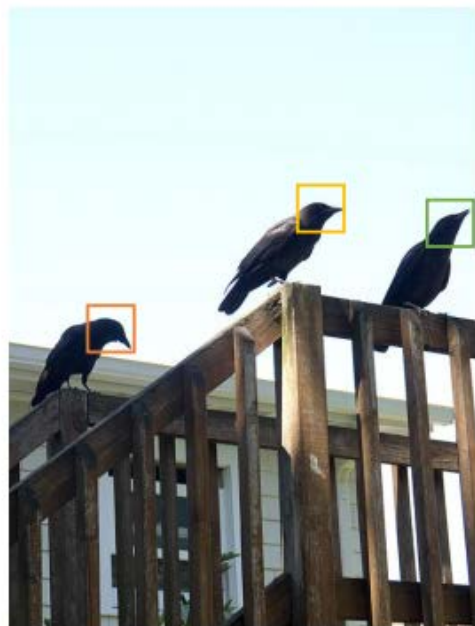
Figure 4: An example of relative distance computation. The relative distances are computed with respect to the position of the highlighted pixel. The format of distances is *row offset*, *column offset*.



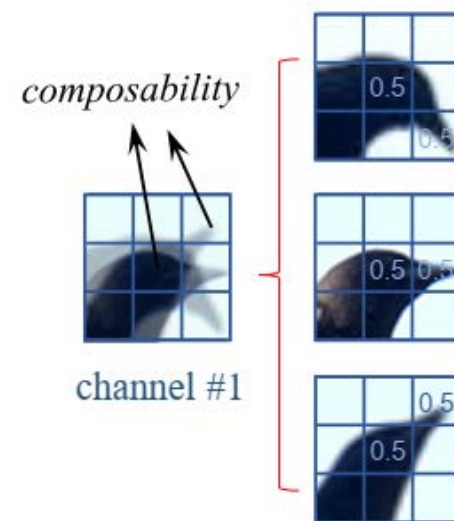




# Local Relation Networks for Image Recognition



convolution

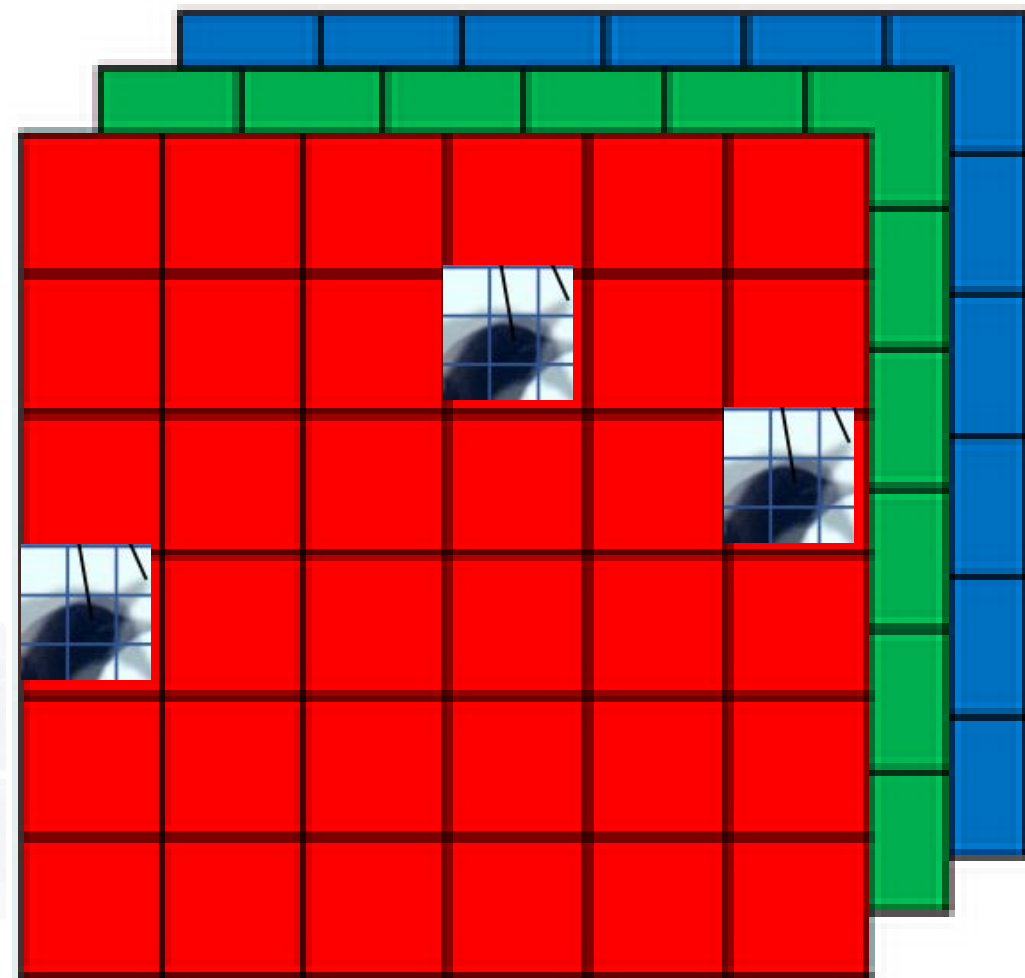
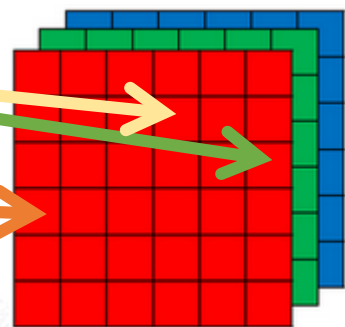
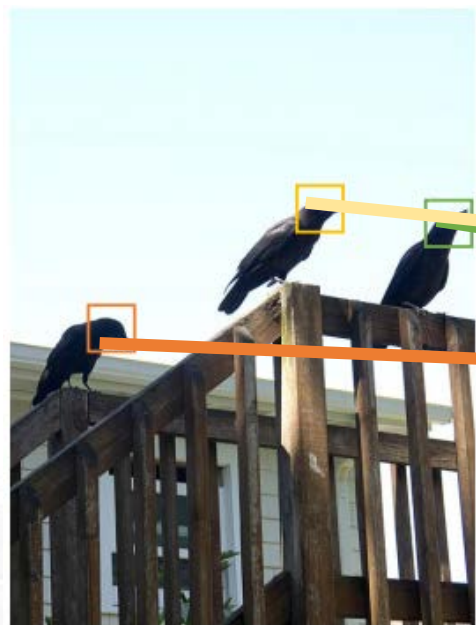


local relation



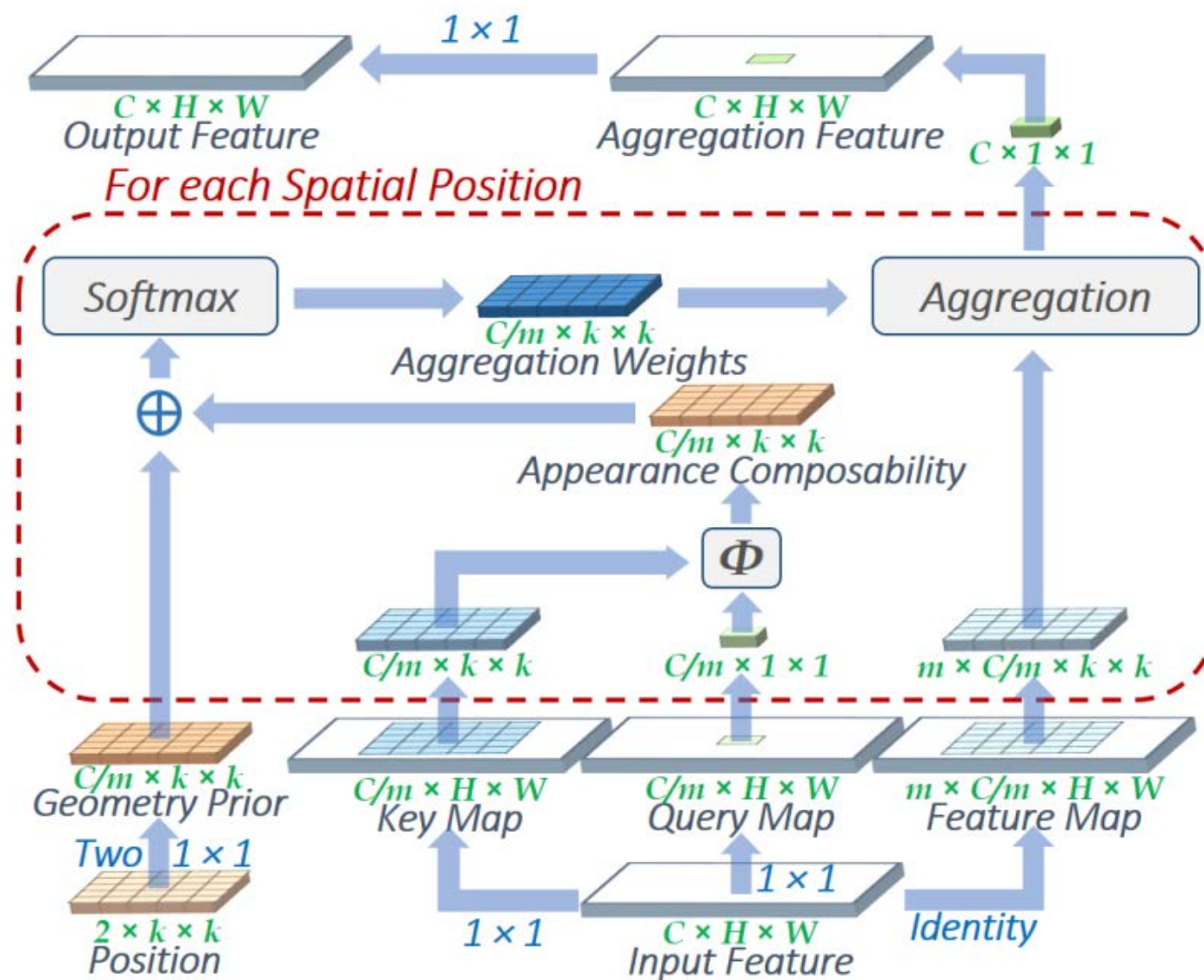


## Local Relation Networks for Image Recognition





# Local Relation Networks for Image Recognition





感谢您的观看

END