

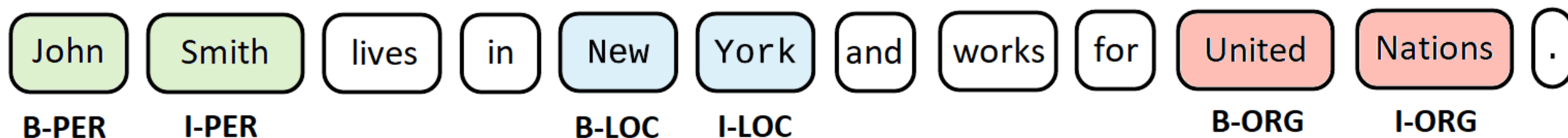
深度学习下的序列标注

王瑞

上海交大计算机系

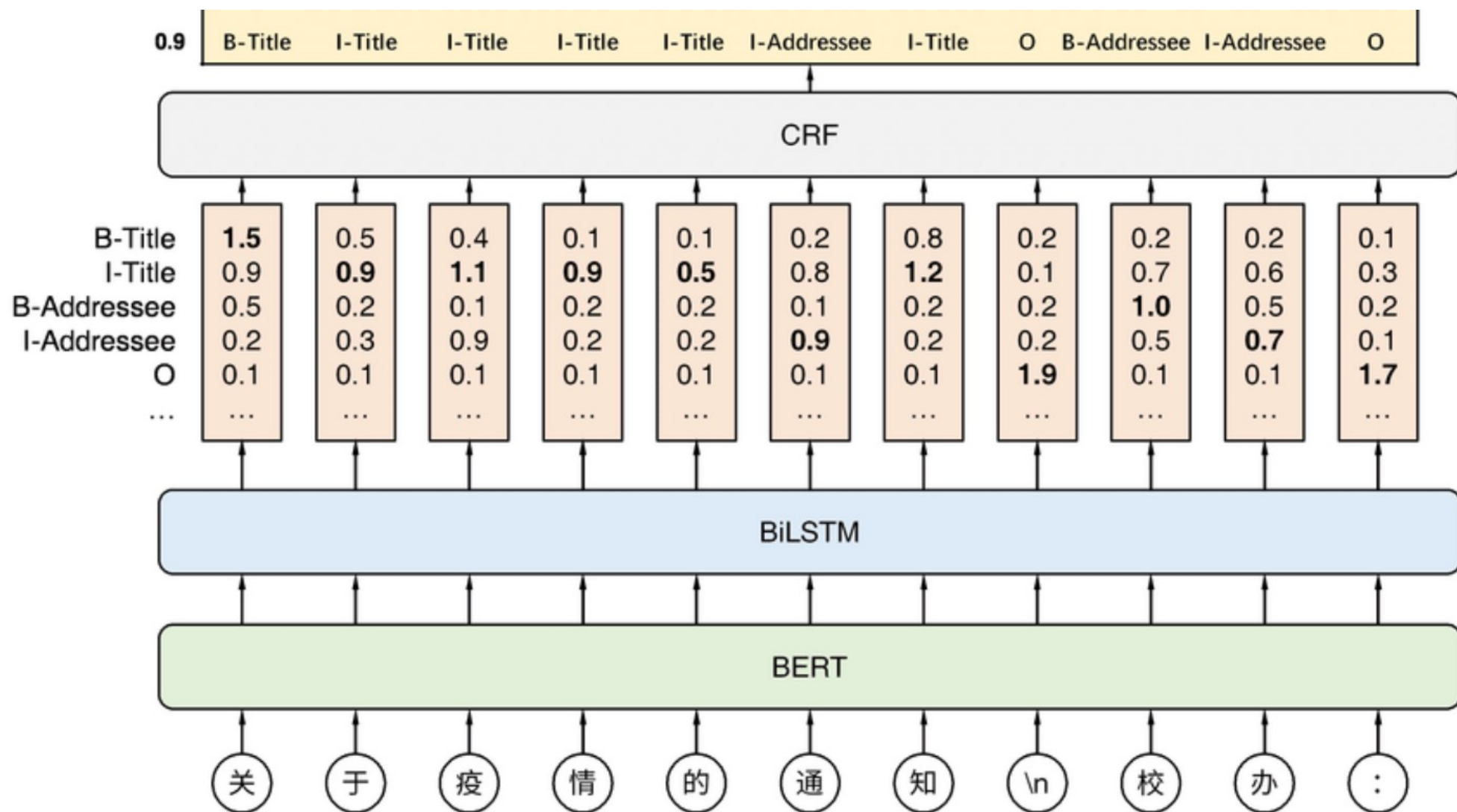
wangrui12@sjtu.edu.cn

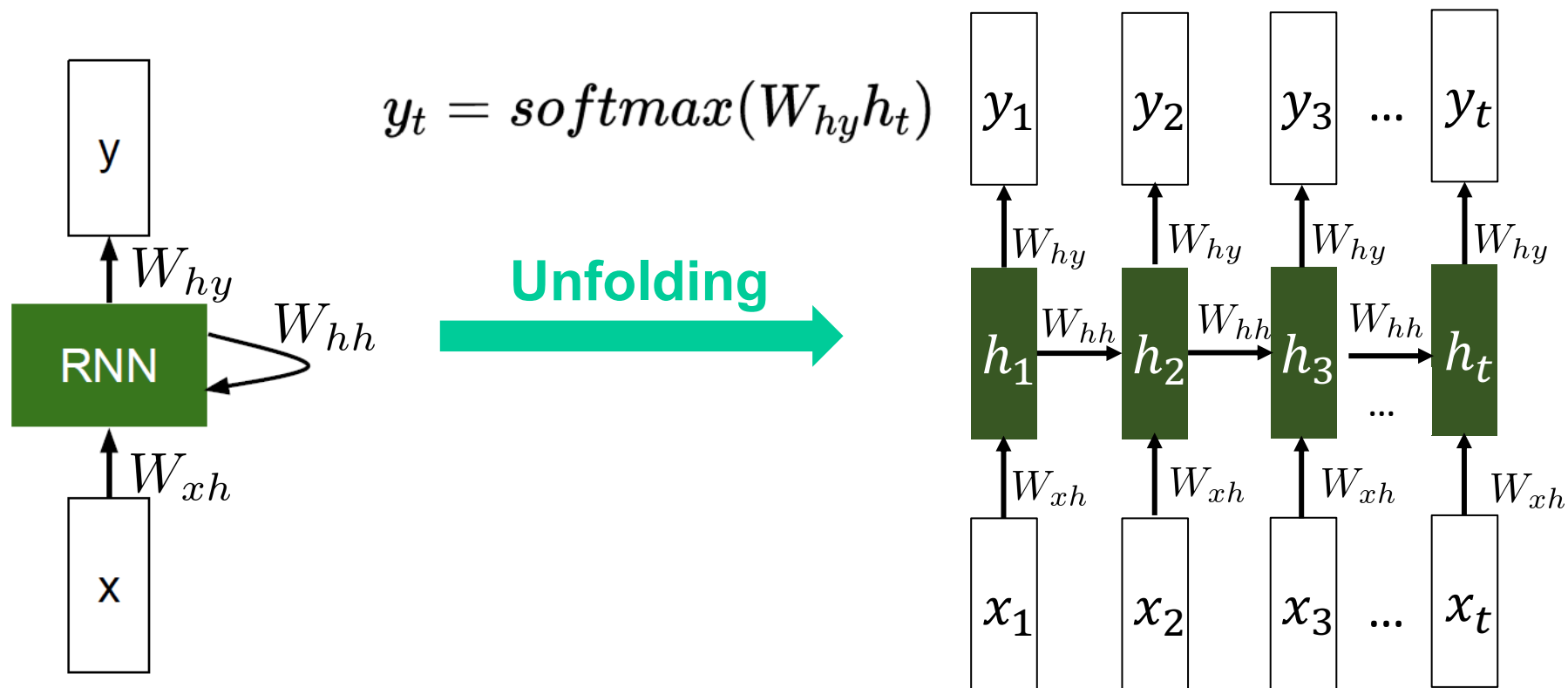
- ▶ **序列标注 (Sequence Tagging/Labeling)** 是NLP中最基础的任务，应用十分广泛，如分词、词性标注 (POS tagging)、命名实体识别 (Named Entity Recognition, NER)、关键词抽取、语义角色标注 (Semantic Role Labeling)、槽位抽取 (Slot Filling) 等实质上都属于序列标注的范畴



- ▶ 输入序列 $\mathbf{X}=\{x_1, x_2, x_3, \dots, x_n\}$, 输出标注序列 $\mathbf{Y}=\{y_1, y_2, y_3, \dots, y_n\}$, 其中 $|y|=m$ (y 的种类数一共有 m 类)
- ▶ Y一般标注法(以NER为例)为:
 - ▶ B (begin): signifies beginning of an NE
 - ▶ I (inside): signifies that the word is inside an NE
 - ▶ E (end): signifies that the word is the end of an NE
 - ▶ S (singleton): signifies that the single word is an NE
 - ▶ O (outside): signifies that the word is just a regular word outside of an NE

序列标注：深度学习框架





LSTM回顾

- ▶ 存储细胞的内容可以在很长时间内保持不变（在输入门和遗忘门都关闭时）
- ▶ 允许信息跨多个时间步保留
- ▶ 允许梯度跨多个时间步流动
- ▶ 克服了梯度消失问题

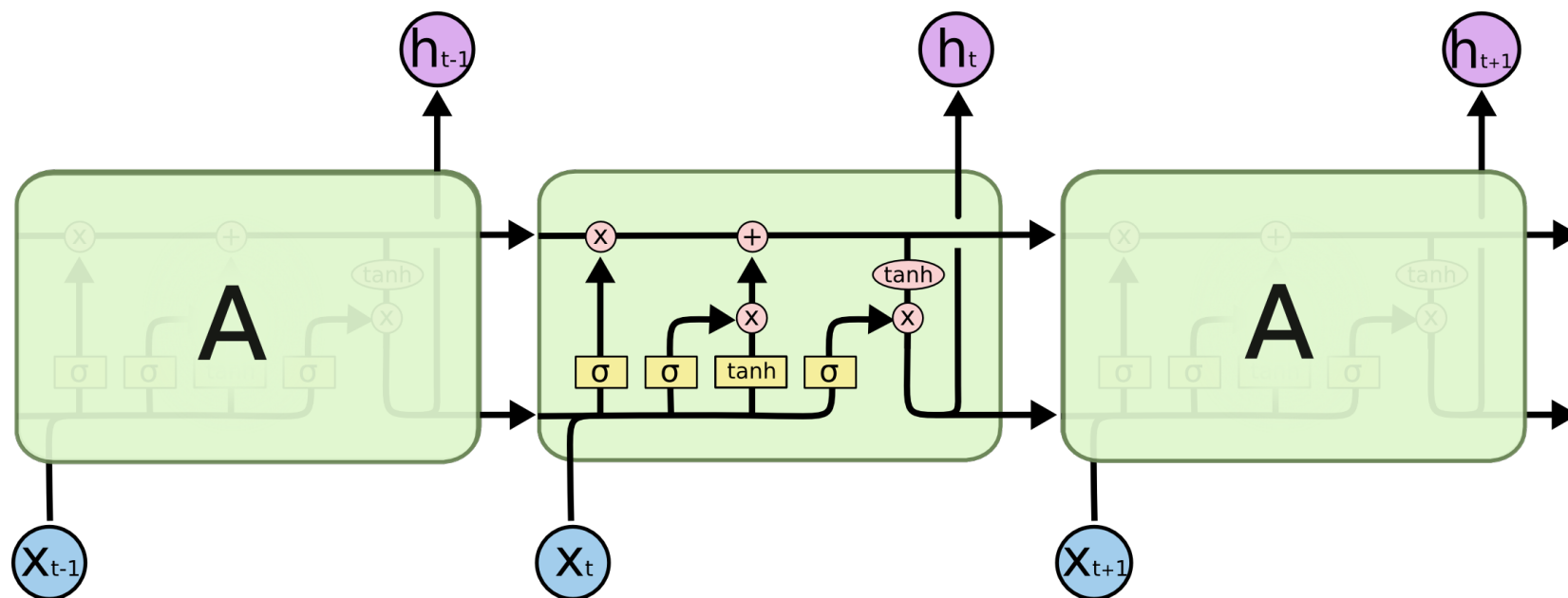


Table 2: Comparison of tagging performance on POS, chunking and NER tasks for various models.

		POS	CoNLL2000	CoNLL2003
Random	Conv-CRF (Collobert et al., 2011)	96.37	90.33	81.47
	LSTM	97.10	92.88	79.82
	BI-LSTM	97.30	93.64	81.11
	CRF	97.30	93.69	83.02
	LSTM-CRF	97.45	93.80	84.10
	BI-LSTM-CRF	97.43	94.13	84.26
Senna	Conv-CRF (Collobert et al., 2011)	97.29	94.32	88.67 (89.59)
	LSTM	97.29	92.99	83.74
	BI-LSTM	97.40	93.92	85.17
	CRF	97.45	93.83	86.13
	LSTM-CRF	97.54	94.27	88.36
	BI-LSTM-CRF	97.55	94.46	88.83 (90.10)

[Huang et al., 2015]

- ▶ 对标签之间的转移(transition)概率进行建模
- ▶ 特征函数:
 - ▶ 输入序列, X ;
 - ▶ 当前输入的位置信息 i ;
 - ▶ 前一个输入的标签 l_{i-1} ;
 - ▶ 当前输入的标签 l_i ;

$$f(X, i, l_{i-1}, l_i)$$

► 从特征函数到概率

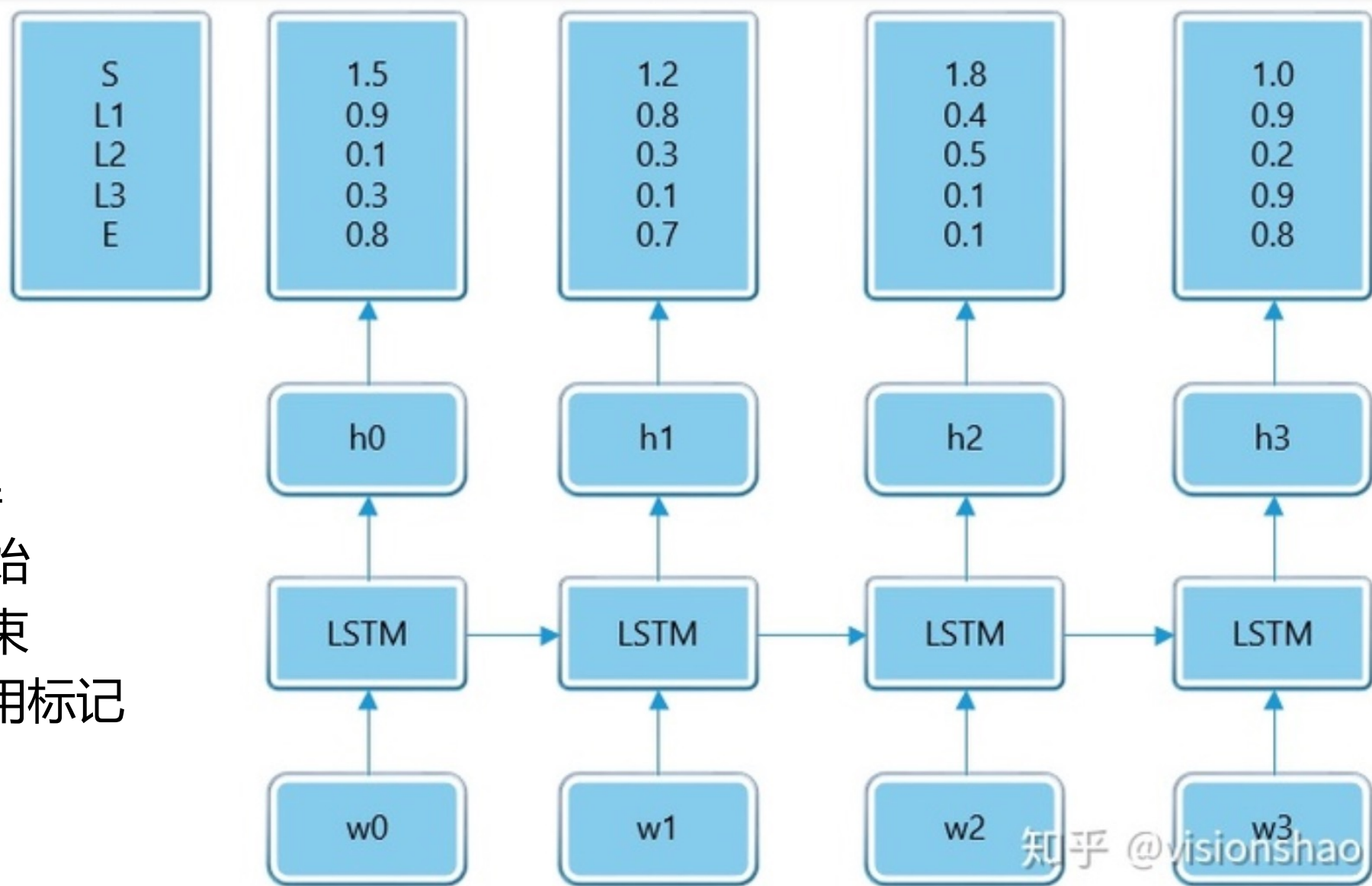
$$score(l|s) = \sum_{j=1}^m \sum_{i=1}^n \lambda_j f_j(s, i, l_i, l_{i-1})$$

► 归一化

$$p(l|s) = \frac{\exp[score(l|s)]}{\sum_{l'} \exp[score(l'|s)]} = \frac{\exp[\sum_{j=1}^m \sum_{i=1}^n \lambda_j f_j(s, i, l_i, l_{i-1})]}{\sum_{l'} \exp[\sum_{j=1}^m \sum_{i=1}^n \lambda_j f_j(s, i, l'_i, l'_{i-1})]}$$

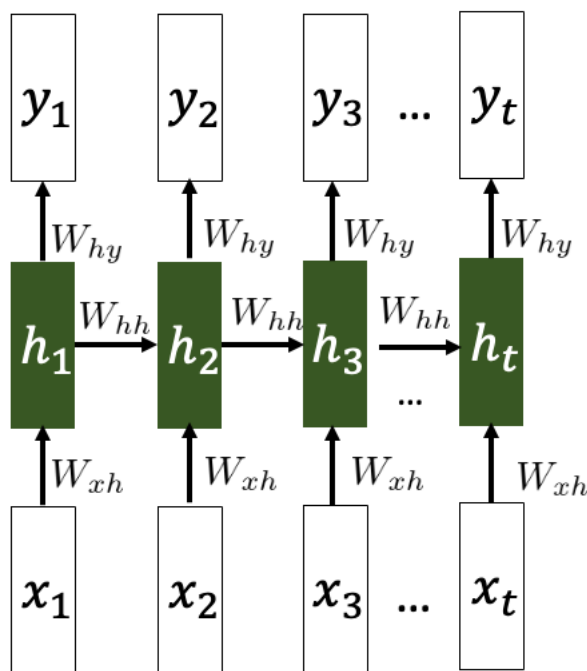
序列标注：深度学习框架 (LSTM+CRF)

- ▶ 便于理解
 - ▶ S: 起始
 - ▶ E: 结束
 - ▶ L: 通用标记



- ▶ 训练：给定训练数据(\mathbf{X}, \mathbf{Y})，如何对LSTM+CRF模型进行训练？
- ▶ 推理：对于训练好的LSTM+CRF模型，给定输入 \mathbf{X} ，如何求得最可能的结果 \hat{Y} ($\hat{}$ 表示预测值)。

- ▶ 发射(*emission*)概率 e : 时刻 i 处于状态 h_i 的条件下生成观测 y_i 的概率, 在LSTM-CRF中, h_i 是 x_i ——对应的隐层状态, 所以这个发射概率也可以理解成从序列 x 到标签 y 的概率。



$$e_i = p(y_i | x_i) = \frac{e^{\text{score}(y_i)}}{\sum_y e^{\text{score}(y)}}$$

- ▶ **为什么要取对数?**
 - ▶ 幂运算可能会导致数值很大，浮点数溢出
 - ▶ 相对于乘法，减法看起来简单一些

$$\log p(y|x) = \text{score}(y) - \log(\sum_y e^{\text{score}(y)})$$

- ▶ 每个词 x_i 输入LSTM后，得到对应的发射概率 e_i
- ▶ 整个序列 $\mathbf{X}=\{x_1, x_2, x_3, \dots, x_n\}$ 输入，得到对应的发射矩阵 E ，其中

$$e_i = E[i] \in R^m$$

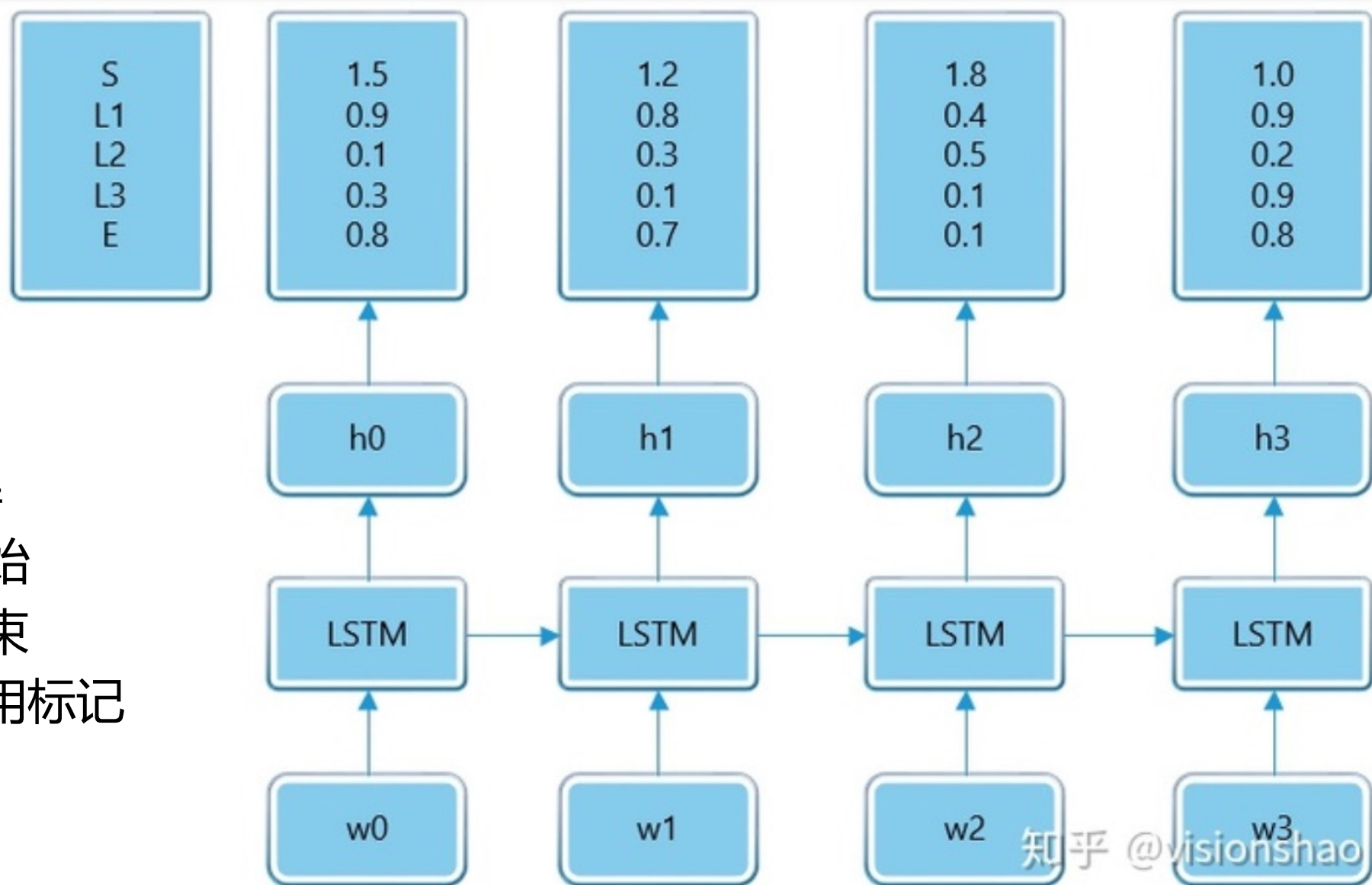
是一个 m 维的向量（ y 的种类数一共有 m 个）

- ▶ 其中第 k 类标签 y_k 的标签为：

$$e_i(y_k)$$

词的标签分布概率

- ▶ 便于理解
 - ▶ S: 起始
 - ▶ E: 结束
 - ▶ L: 通用标记



知乎 @visionshao

- ▶ $Y = \{y_1, y_2, y_3, \dots, y_n\}$ 是一个序列 (sequence), 彼此间有依赖关系
 - ▶ 动词后面一般不再出现动词, 代词后面动词概率大
 - ▶ 两个命名实体一般不连续出现
 - ▶ 中文中时间一般出现在地点前面
 - ▶ 中文人名一般不超过4个字
 - ▶ . . .
- ▶ 如何能构建这个转移概率?

标签转移矩阵

- ▶ 我们将标签 y_{i-1} 转移到 y_i 的概率定义为标签转移概率 $\pi[i-1, i]$
- ▶ 这样就构建了一个标签转移矩阵 T 。

	S	L1	L2	L3	E
S	0	0.7	0.8	0.2	0.3
L1	0	0.2	0.1	0.5	1.3
L2	0	0.4	1.4	1.5	0.1
L3	0	0.7	0.3	0.4	0.5
E	0	0	0	0	0

知乎 @visionshao

- ▶ 对于一个标签序列Y综合打分为：发射概率+转移概率（log前提下）

$$\textit{score}(Y) = \sum_{i=1}^n e_i + \sum_{i=2}^n (T[i-1, i])$$

- ▶ $\log p(y|x) = \text{score}(y) - \log(\sum_y e^{\text{score}(y)})$
 - ▶ 我们已经知道了 $\text{score}(Y)$
 - ▶ 还需要知道全部可能概率 $\log(\sum_y e^{\text{score}(y)})$
- ▶ 但是，这种情况下共有 m^n 个 (m 是标签种类, n 是序列长度) 可能的标记结果，如果列举所有的路径再求总分，那么计算量就很大。

► 动态规划

- 大致上，若要解一个给定问题，我们需要解其不同部分（即子问题），再根据子问题的解以得出原问题的解。
- 通常许多子问题非常相似，为此动态规划法试图仅仅解决每个子问题一次，从而减少计算量：一旦某个给定子问题的解已经算出，则将其记忆化存储，以便下次需要同一个子问题解之时直接查表。
- 这种做法在重复子问题的数目关于输入的规模呈指数增长时特别有用。

- ▶ 我们知道对于 t 时刻的输出 y_t 来说, 其有 m 种可能, 我们记录下从开始时刻到 t 时刻时, y_t 取某一个值 i 时经过的所有路径的分数的指数和的对数:

$$\log Z(y_t = i) = \log \sum_{y_t} e^{\text{score}(\text{start} \rightarrow y_t)}$$

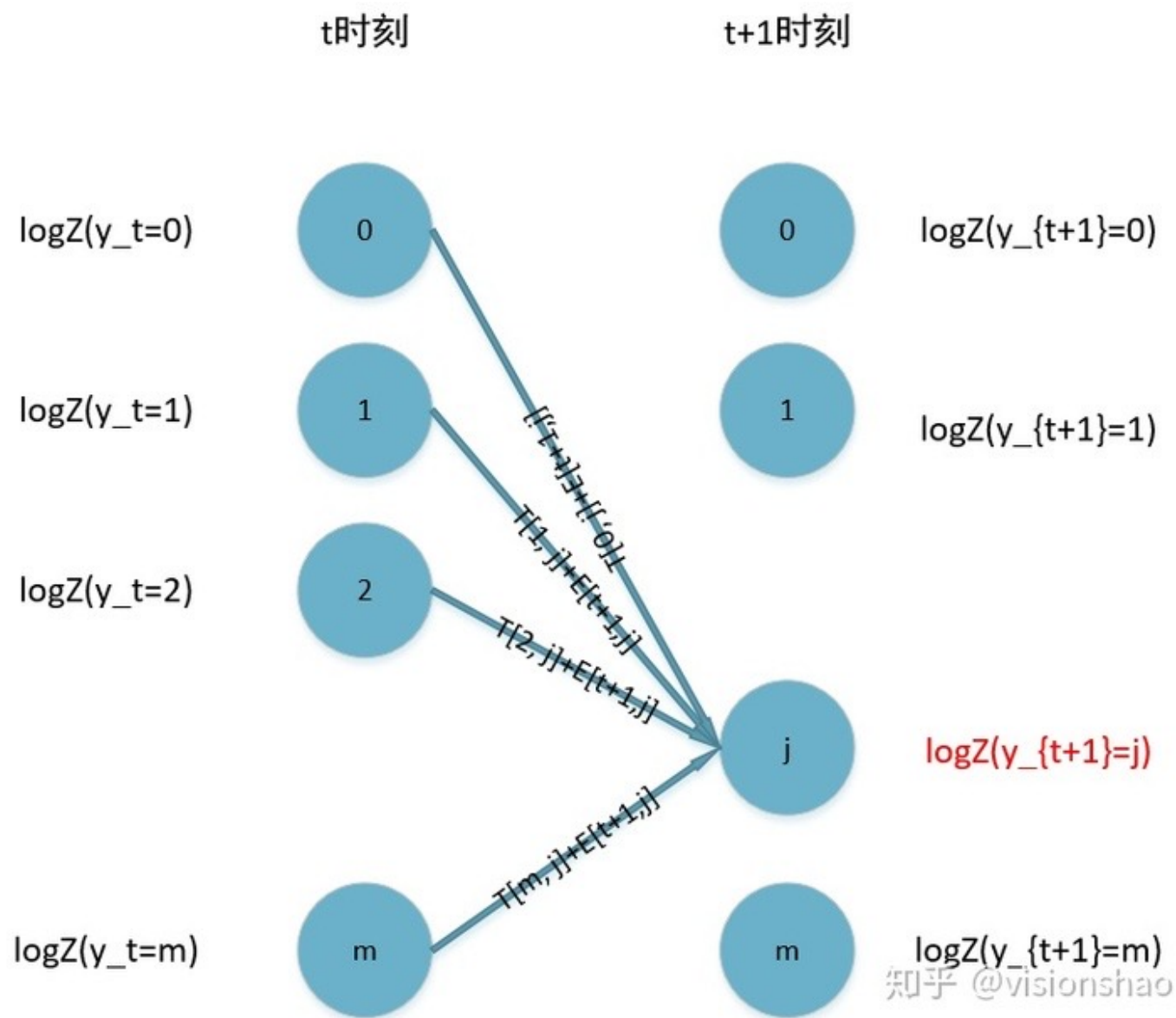
下一个子问题

- ▶ 我们就可以利用此信息获得 $t+1$ 时刻, $y_{t+1}=j$ 的所有路径的得分的指数和的对数.
- ▶ 我们知道对于 $y_t=i$ 到 $y_{t+1}=j$, 路径得分变化是在原来的基础上加上

$$T[i, j] + E[t+1, j]$$

- ▶ 所以:

$$\log Z(y_{t+1} = j) = \log\left(\sum_{y_t=0}^{m-1} e^{\text{score}(\text{start} \rightarrow y_t)}\right) + (T[i, j] + E[t+1, j])$$



- ▶ 这样就需要我们维护一个 m 维向量 α_t , $\alpha_t[i]$ 表示时刻 t , $y_t = i$ 时所有路径指数和的对数 $\log Z(y_t = i)$, 即:

$$\alpha_t = [\log Z(y_t = 0), \log Z(y_t = 1), \dots, \log Z(y_t = m)]$$

- ▶ 如果我们想根据 α_t 得到 α_{t+1} ，则按照前面的公式来说，

$$\alpha_{t+1} = [\sum(\alpha_t + T[:, 0] + E[t + 1, 1]), \sum(\alpha_t + T[:, 0] + E[t + 1, 1]), \dots, \sum(\alpha_t + T[:, m] + E[t + 1, m])]$$

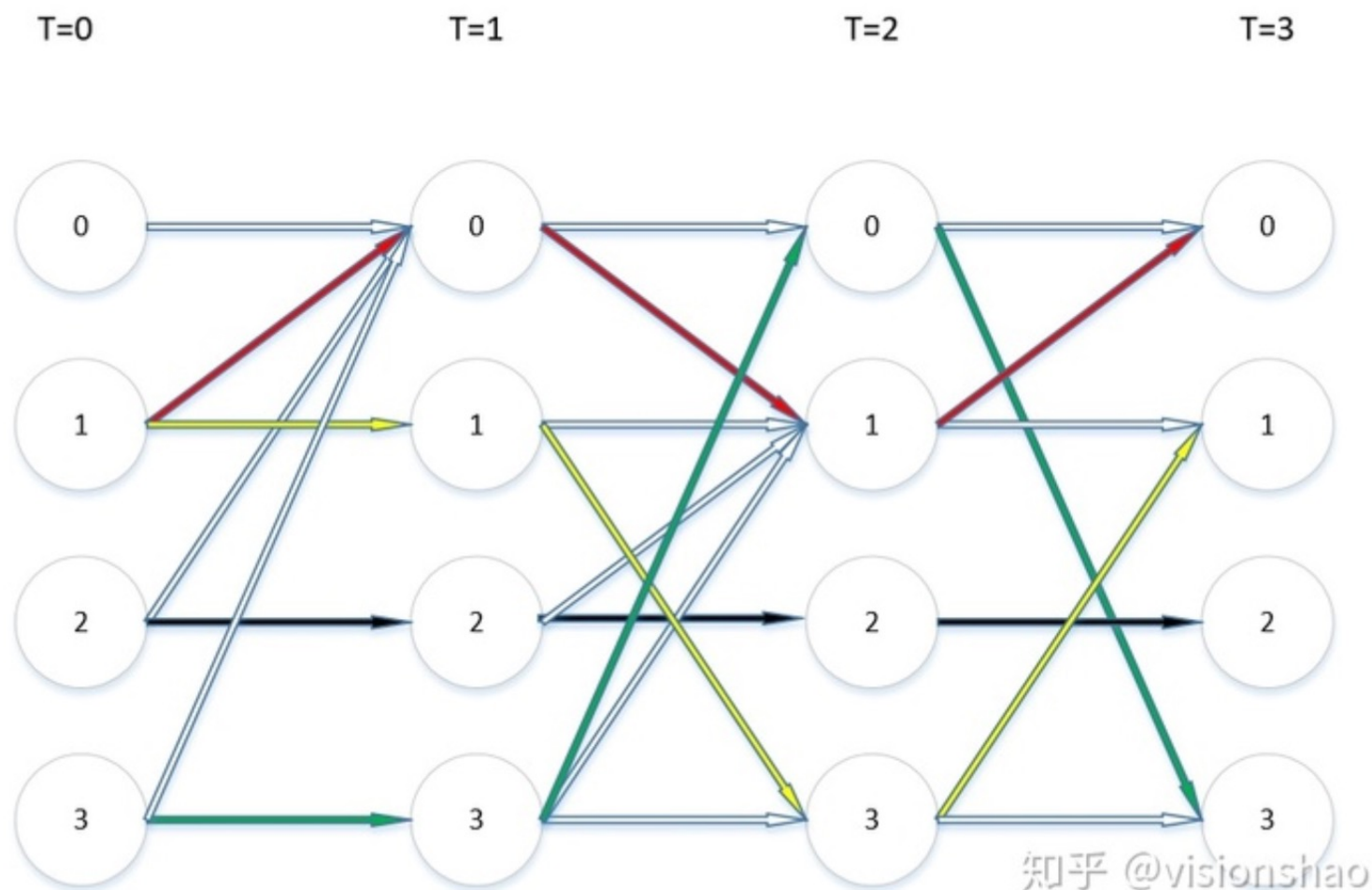
当得到最后一个时刻 n 的 α 时，我们只需要用 $\text{logsumexp}(\alpha_n)$ 就可以得到最终所有路径的分数的指数和的对数，即 $\log Z = \log(\sum_y e^{\text{score}(y)})$ 。

这样，我们就能最终得到 $\log p(y|x) = \text{score}(y) - \log(\sum_y e^{\text{score}(y)})$ 。

- ▶ 训练：给定训练数据 (\mathbf{X}, \mathbf{Y}) ，如何对LSTM+CRF模型进行训练？
- ▶ 推理：对于训练好的LSTM+CRF模型，给定输入 \mathbf{X} ，如何求得最可能的结果 $\hat{\mathbf{Y}}$ ($\hat{\cdot}$ 表示预测值)。

得分最高的路径

带有颜色的箭头表示的路径是指向该节点的得分最高的路径

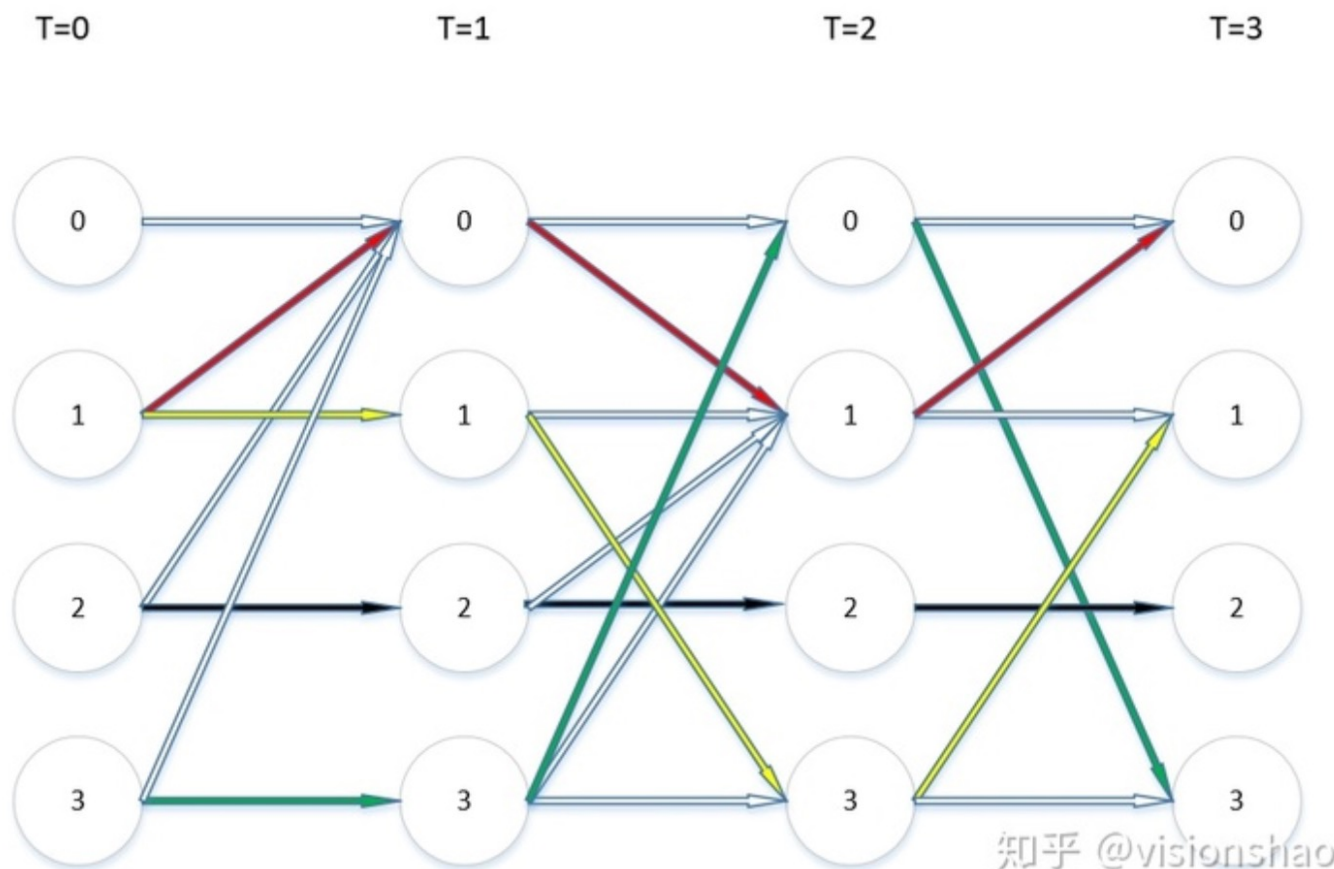


得分最高的路径

我们只需要找出得分最高的路径（某种颜色表示的路径）指向的节点，然后不断回溯即可找到整体的最优路径。

n 个词

m 类标签



得分最高的路径

- ▶ 要想实现这种方法，我们需要保存每个节点对应的得分最高的路径及其分数。
- ▶ 经过LSTM后得到的发射矩阵为 $E \in R^{(n,m)}$
- ▶ 标签转移矩阵为 $T \in R^{(m,m)}$
- ▶ 其中R表示矩阵的维度

- ▶ 首先, $t=0$ 时刻, 我们获取各个节点对应的最大路径的得分,

$$\beta_0 \in R^m$$

- ▶ 即节点的发射得分

$$E[0, :]$$

- ▶ t_0 时刻节点 i 的最大得分路径到节点 j 的路径得分:

$$M_1[i, j] = \beta \circ T$$

- ▶ \circ 表示element-wise乘积 (如果取log就是相加)

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \circ \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \end{bmatrix} = \begin{bmatrix} a_{11} b_{11} & a_{12} b_{12} & a_{13} b_{13} \\ a_{21} b_{21} & a_{22} b_{22} & a_{23} b_{23} \\ a_{31} b_{31} & a_{32} b_{32} & a_{33} b_{33} \end{bmatrix}.$$

- ▶ 如果我们想求得 $t = 1$ 时刻到节点 j 的最大得分路径，只需要对 $M_1[:, j]$ 求最大值。
- ▶ 同理， $t = 1$ 时刻，我们获取各个节点对应的最大路径的得分，

$$\beta_1 \in R^m$$

- ▶ 也得到了 $t = 1$ 时刻各个节点对应的最大路径在 $t = 0$ 时刻对应的节点向量

$$P_1 \in R^m$$

- ▶ 其中 $/$ 表示 index

$$I_0 = P_1[I_1]$$

- ▶ 不断重复上面的步骤，直到得到 $t = n - 1$ 时刻的 β_{n-1} 和 P_{n-1} 。这时 β_{n-1} 对应的是最后时刻的各个节点对应的最大路径得分。我们求得 β_{n-1} 中最大值对应的索引 I_{n-1}
- ▶ 这个索引表示最优路径的最后一个节点。然后将 I_{n-1} 带入 P_{n-1} ，得到 $I_{n-2} = P_{n-1}[I_{n-1}]$ ，表示最优路径倒数第二个节点，然后依次类推，知道得到 $I_0 = P_1[I_1]$
- ▶ 这时，最优路径为 $l = [l_0, l_1, \dots, l_{n-1}]$

序列标注：深度学习框架

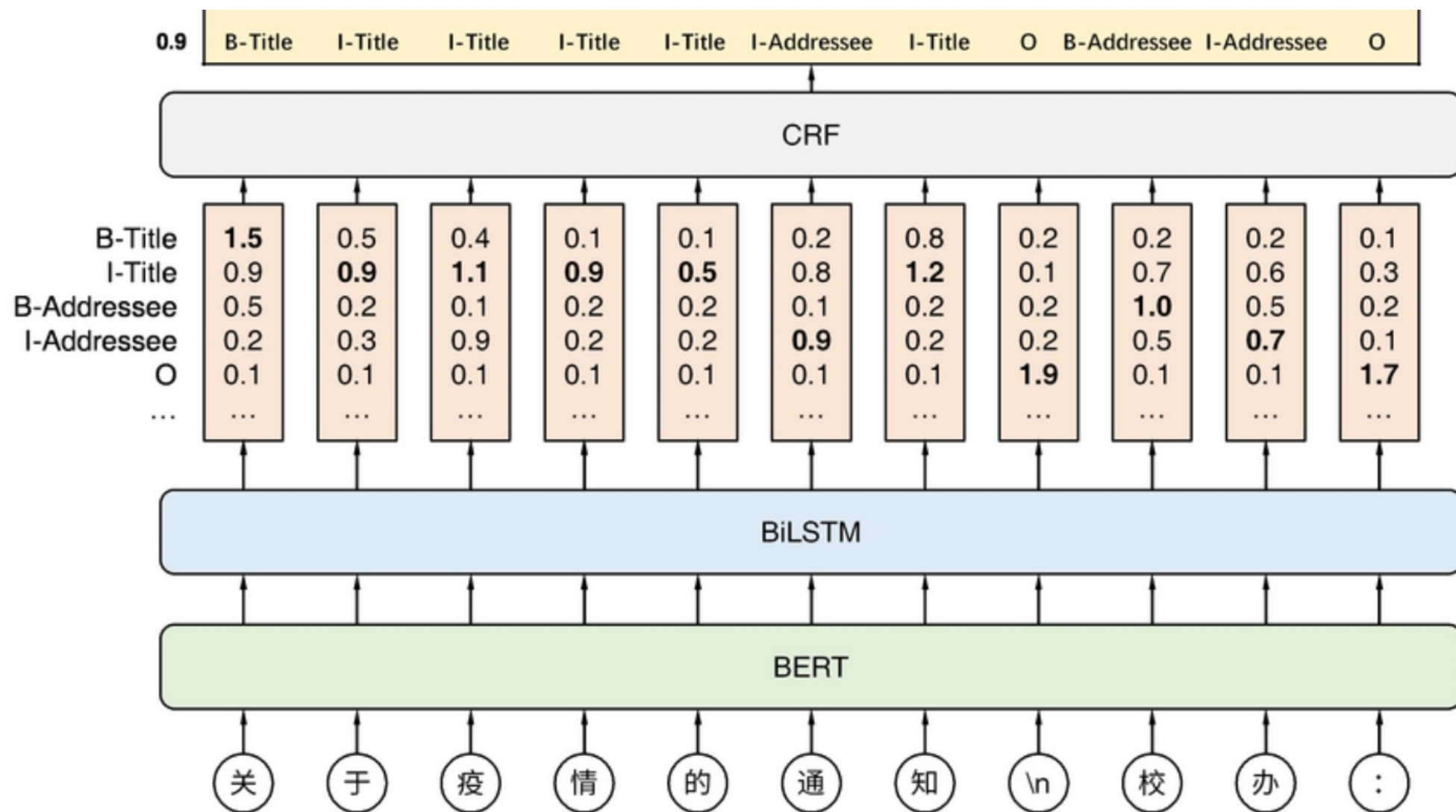


Table 2: Comparison of tagging performance on POS, chunking and NER tasks for various models.

		POS	CoNLL2000	CoNLL2003
Random	Conv-CRF (Collobert et al., 2011)	96.37	90.33	81.47
	LSTM	97.10	92.88	79.82
	BI-LSTM	97.30	93.64	81.11
	CRF	97.30	93.69	83.02
	LSTM-CRF	97.45	93.80	84.10
	BI-LSTM-CRF	97.43	94.13	84.26
Senna	Conv-CRF (Collobert et al., 2011)	97.29	94.32	88.67 (89.59)
	LSTM	97.29	92.99	83.74
	BI-LSTM	97.40	93.92	85.17
	CRF	97.45	93.83	86.13
	LSTM-CRF	97.54	94.27	88.36
	BI-LSTM-CRF	97.55	94.46	88.83 (90.10)