# Max-margin learning

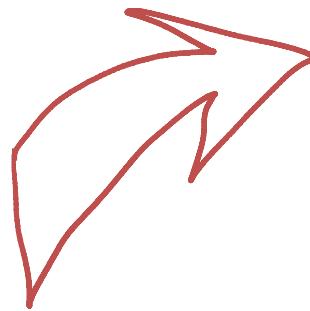Nando de Freitas

UNIVERSITY OF
OXFORD

# Outline of the lecture

Max margin learning is an extremely powerful idea for learning features with auxiliary tasks, and then use these features to solve tasks with few data. The goal of this lecture is for you to learn

- ❑ Transfer, multi-task, and multi-instance learning
- ❑ Harnessing auxiliary tasks to learn features
    - ❑ Matchings
    - ❑ Preferences
    - ❑ Corruption
- ❑ Formulations of multi-task learning
- ❑ Applications:
    - ❑ Cross-lingual embeddings
    - ❑ Relation learning
    - ❑ Question answering
    - ❑ Memory networks

# Embedding discrete objects in metric spaces

Code
Words
Formulae
Logical expressions
Symbols
DNA sequences

$w_i \in \mathbb{R}^2$

marmars
decembre
dezember
avril
april
august
août
juju ni
januar
janvier
oktober
octobre
novembre
novembre
septembre
september
juillet
februar
fevrier
mai
mai

$w_i \in \mathbb{R}^d$

$$\begin{matrix}1\\d\end{matrix}\begin{bmatrix} w_3 \end{bmatrix} = \begin{bmatrix} w_1 & w_2 & w_3 & \dots & w_K \end{bmatrix} \begin{matrix}1\\ \\ \\ \\ \\ \\ \\ K\end{matrix}\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

d-dimensional encoding of word.

$d \times K$

← vector of zeros with a 1 indicating a word token in the vocabulary.

**Document Embedding**

Vectorisation

Stackable Layers

Pooled representation

K-max pooling

Feature map

Wide convolution

**Sentence Embedding - Document Matrix**

Vectorisation

Stackable Layers

Pooled representation

K-max pooling

Feature map

Wide convolution

**Word Embedding - Sentence Matrix**

The cat sat on the mat .

They found it really really funny .

*Handwritten annotations:* sentence 1, sentence 2, 2 features, pooling, $d=4$, $w_i$, sentence1 sentence2, sentence1, sentence 2

# **Idea**: learn embeddings (features) in one task and transfer these to solve new tasks
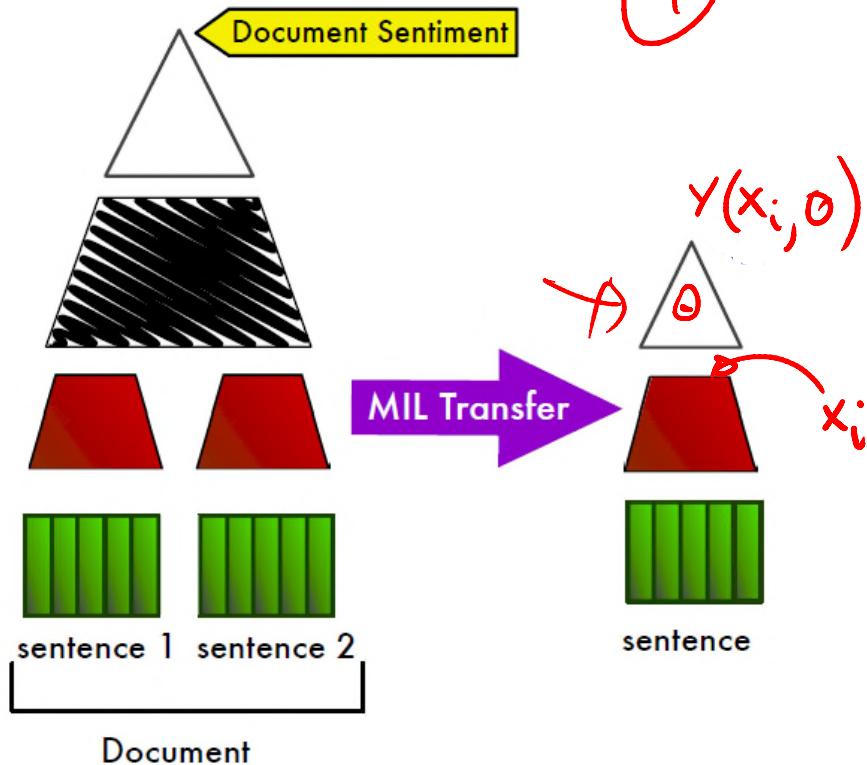


[Kotzias, Denil & NdF, 2014]

# Deep Multi-Instance Learning

$$W(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|_2^2)$$

← Similarity

← sentence embeddings

$$J(\boldsymbol{\theta}) = \sum_{i,j \in I} W(\mathbf{x}_i, \mathbf{x}_j)\,(y(\mathbf{x}_i, \boldsymbol{\theta}) - y(\mathbf{x}_j, \boldsymbol{\theta}))^2 + \lambda \sum_{g \in G} \left( \frac{1}{|g|} \sum_{i \in g} y(\mathbf{x}_i, \boldsymbol{\theta}) - s_g \right)^2$$

①   ②



Document Sentiment

MIL Transfer

$y(x_i, \theta)$

$\theta$

$x_i$

sentence 1   sentence 2

Document

sentence

# Deep Multi-Instance Learning

Paul Bettany did a great role as the tortured father whose favorite little girl dies tragically of disease.
For that, he deserves all the credit.
However, the movie was mostly about exactly that, keeping the adventures of Darwin as he gathered data for his theories as incomplete stories told to children and skipping completely the disputes regarding his ideas.
Two things bothered me terribly: the soundtrack, with its whiny sound, practically shoving sadness down the throat of the viewer, and the movie trailer, showing some beautiful sceneries, the theological musings of him and his wife and the enthusiasm of his best friends as they prepare for a battle against blind faith, thus misrepresenting the movie completely.
To put it bluntly, if one were to remove the scenes of the movie trailer from the movie, the result would be a non descript family drama about a little child dying and the hardships of her parents as a result.
Clearly, not what I expected from a movie about Darwin, albeit the movie was beautifully interpreted.

# Auxiliary tasks to learn features that can be transferred to learn tasks with few labels
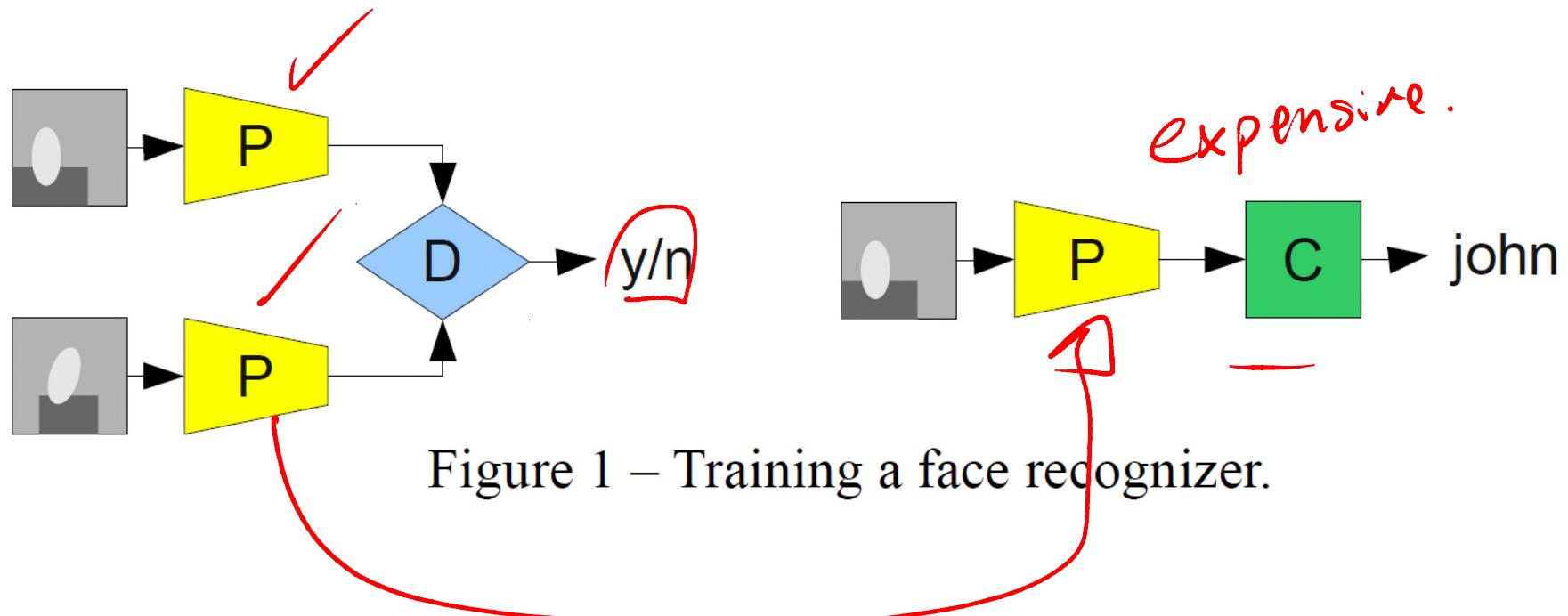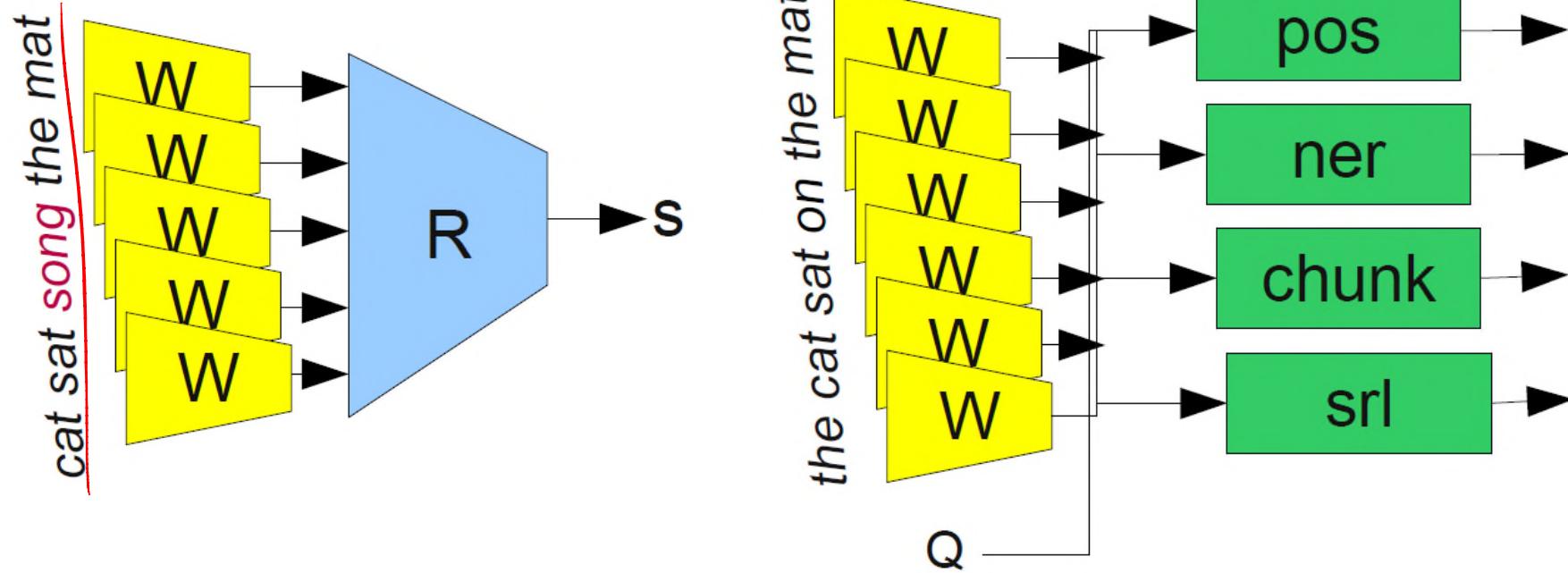
## 1. Matchings



Figure 1 – Training a face recognizer.

[*From machine learning to machine reasoning*, Leon Bottou]

# Auxiliary tasks to learn features that can be transferred to learn tasks with few labels

## 2. Corruption



[*NLP almost from scratch*, Ronan Collobert and colleagues]

# Auxiliary tasks to learn features that can be transferred to learn tasks with few labels
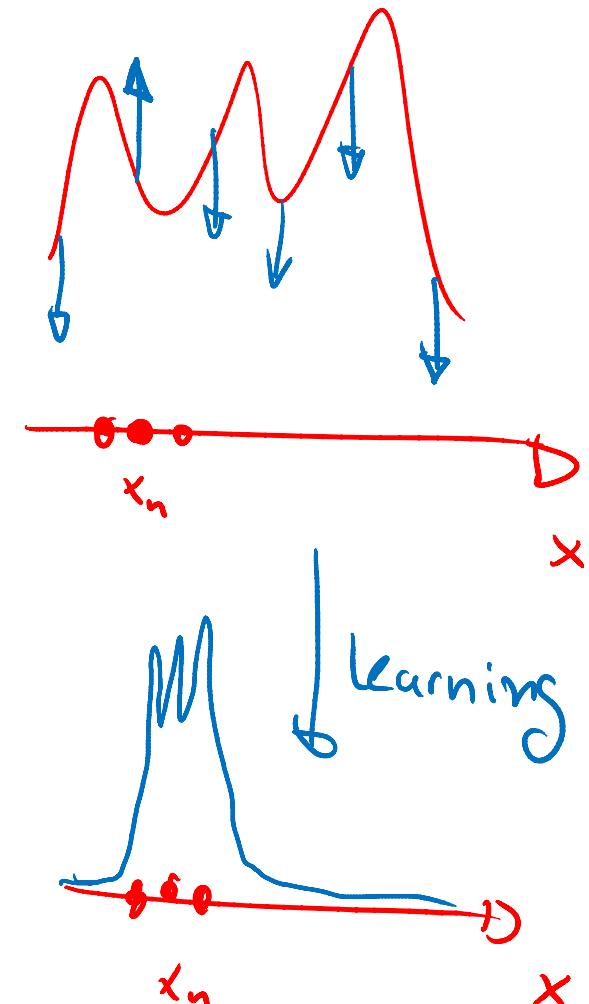
## 1. Preferences

# Max-margin formulations

data instance $\mathbf{x}_n \in \mathcal{X}$

$$p(\mathbf{x}_n|\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \exp[-E(\mathbf{x}_n, \boldsymbol{\theta})]$$

$$Z(\boldsymbol{\theta}) = \int_{\mathcal{X}} \exp[-E(\mathbf{x}_n, \boldsymbol{\theta})] d\mathbf{x}_n$$

$$\boldsymbol{\theta}^\star = \underset{\boldsymbol{\theta}: \, p(\mathbf{x}_n|\boldsymbol{\theta}) \geq p(\mathbf{x}|\boldsymbol{\theta}) + m_p}{\arg\min} \|\boldsymbol{\theta}\|_2^2$$

$$\forall n, \ \forall \mathbf{x} : \|\mathbf{x} - \mathbf{x}_n\| > \epsilon$$

# Max-margin formulations

$$\theta^{\star} = \underset{\theta : E(\mathbf{x}, \theta) \geq E(\mathbf{x}_n, \theta) + m}{\arg\min} \|\theta\|_2^2$$
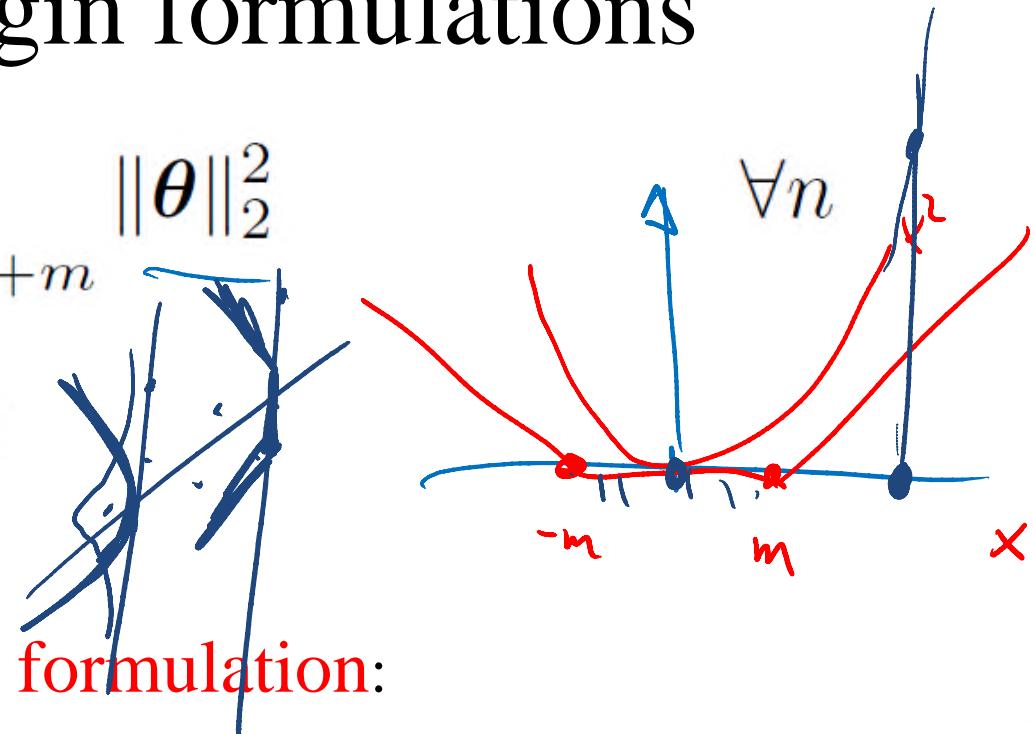
$$\forall n, \ \forall \mathbf{x} : \|\mathbf{x} - \mathbf{x}_n\| > \epsilon$$

*[A tutorial on energy based learning, Yann LeCun et al]*

# Max-margin formulations

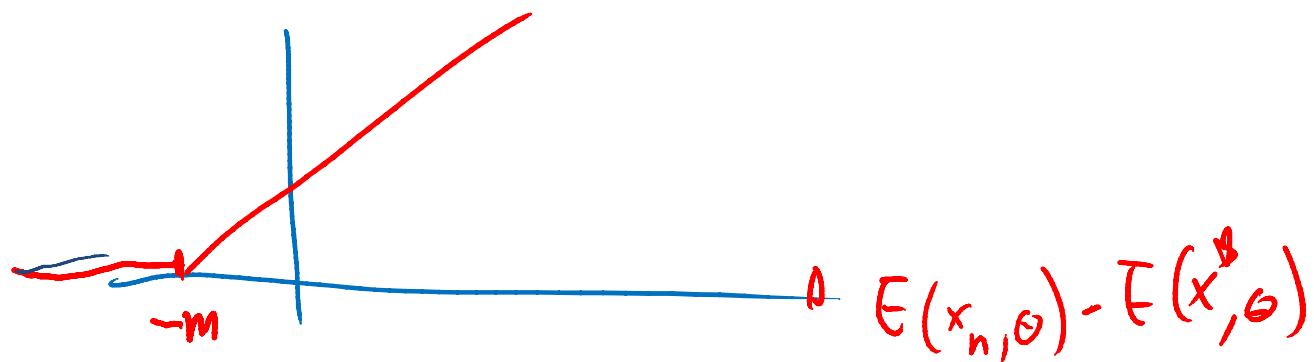$$\boldsymbol{\theta}^{\star} = \underset{\boldsymbol{\theta}: E(\mathbf{x}^{\star}, \boldsymbol{\theta}) \geq E(\mathbf{x}_n, \boldsymbol{\theta}) + m}{\arg\min} \|\boldsymbol{\theta}\|_2^2 \qquad \forall n$$

$$\mathbf{x}^{\star} = \underset{\mathbf{x}: \|\mathbf{x} - \mathbf{x}_n\| > \epsilon}{\arg\min} E(\mathbf{x}, \boldsymbol{\theta})$$

<span style="color:red">Hinge Loss unconstrained formulation</span>:

$$\boldsymbol{\theta}^{\star} = \underset{\boldsymbol{\theta}}{\arg\min} \left\{ \sum_{n=1}^{N} \max\left[0, m + E(\mathbf{x}_n, \boldsymbol{\theta}) - E(\mathbf{x}^{\star}, \boldsymbol{\theta})\right] + \lambda \|\boldsymbol{\theta}\|_2^2 \right\}$$

# Max-margin formulations

introduce corrupted data $\widetilde{\mathbf{x}}_n$

correct

incorrect

$$\boldsymbol{\theta}^{\star} = \arg \min_{\boldsymbol{\theta}} \left\{ \sum_{n=1}^{N} \max \left[ 0, m_e + E(\mathbf{x}_n, \boldsymbol{\theta}) - E(\widetilde{\mathbf{x}}_n, \boldsymbol{\theta}) \right] + \lambda \|\boldsymbol{\theta}\|_2^2 \right\}$$

# Hinge loss layer

+1    -1

$$z = E(x_1, x_2, m) = \max\left[0, m + x_2 - x_1\right]$$

scalar loss $E = \left(\mathbf{z}^{L+1}\right)$  $\left(\frac{\partial E}{\partial \mathbf{z}^{L+1}} = 1\right)$
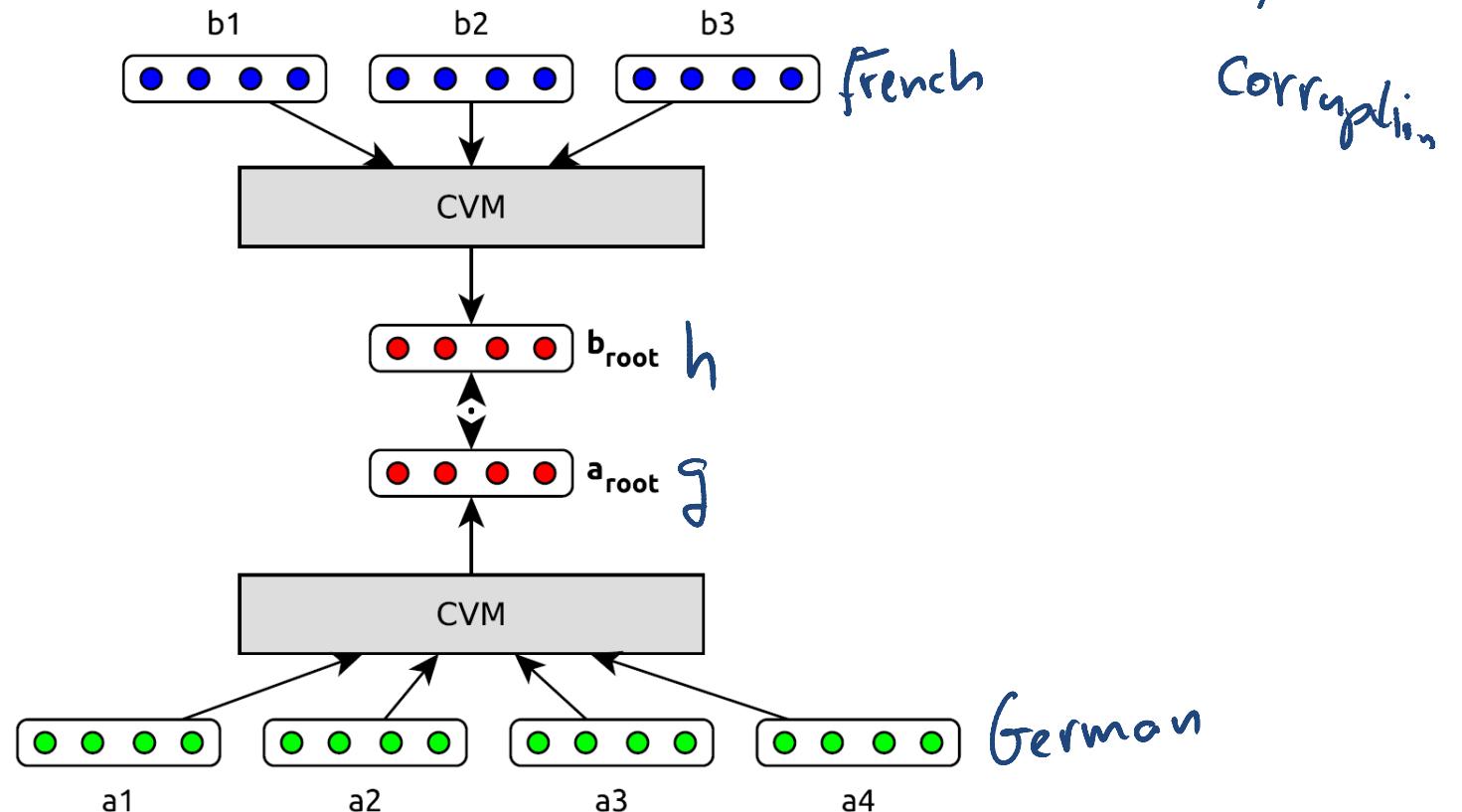
$$\frac{\partial E}{\partial x_i} = \frac{\partial E}{\partial z}\frac{\partial z}{\partial x_i} = (-1)^{i}\,\mathbb{1}_{(m+x_2-x_1>0)}$$

loss (layer $L$)

prediction $\mathbf{z}^L$  $\quad \frac{\partial E}{\partial \mathbf{z}^L}$

layer $L-1$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad \frac{\partial E}{\partial \boldsymbol{\theta}}$

$\boldsymbol{\theta}^l \longrightarrow$  layer $l$  $\quad \dashrightarrow \frac{\partial E}{\partial \boldsymbol{\theta}^l}$

$\mathbf{z}^l$  $\quad \frac{\partial E}{\partial \mathbf{z}^l}$  $\quad\quad\quad\quad\quad\quad \longleftarrow \boldsymbol{\theta}$

$\boldsymbol{\theta}^{l-1} \longrightarrow$  layer $l-1$  $\quad \dashrightarrow \frac{\partial E}{\partial \boldsymbol{\theta}^{l-1}}$

layer 1

input $\mathbf{x} = \left(\mathbf{z}^1\right)$  $\left(\frac{\partial E}{\partial \mathbf{z}^1}\right) = \frac{\partial E}{\partial \mathbf{x}}$

# Example: Bi-lingual word embeddings

$$E_{bi}(\mathbf{a}, \mathbf{b}, \boldsymbol{\theta}) = \|\mathbf{g}(\mathbf{a}, \boldsymbol{\theta}_g) - \mathbf{h}(\mathbf{b}, \boldsymbol{\theta}_h)\|^2$$

$$L_{hinge}(\mathbf{a}, \mathbf{b}, \mathbf{n}, \boldsymbol{\theta}) = \max\left[0, m + E_{bi}(\mathbf{a}, \mathbf{b}, \boldsymbol{\theta}) - E_{bi}(\mathbf{a}, \mathbf{n}, \boldsymbol{\theta})\right]$$

Corruption



[Karl Hermann and Phil Blunsom, 2014]

# Siamese networks (Yann LeCun)

$$L_{siamese}(\boldsymbol{\theta}, S_{ij}) = \begin{cases} ||\mathbf{h}(\mathbf{x}_i; \boldsymbol{\theta}) - \mathbf{h}(\mathbf{x}_j; \boldsymbol{\theta})||_2^2 & \text{if } S_{ij} = 1 \\ \max\left[0, m_h - ||\mathbf{h}(\mathbf{x}_i; \boldsymbol{\theta}) - \mathbf{h}(\mathbf{x}_j; \boldsymbol{\theta})||_2^2\right] & \text{if } S_{ij} = 0 \end{cases}$$

Semi-supervised deep learning (Jason Weston et al)

$$\sum_{n \in \mathcal{L}} \text{NLL}(\mathbf{h}(\mathbf{x}_n; \boldsymbol{\theta}), \mathbf{y}_n) + \lambda \sum_{i,j \in \mathcal{U}} L_{siamese}(\boldsymbol{\theta}, S_{ij})$$
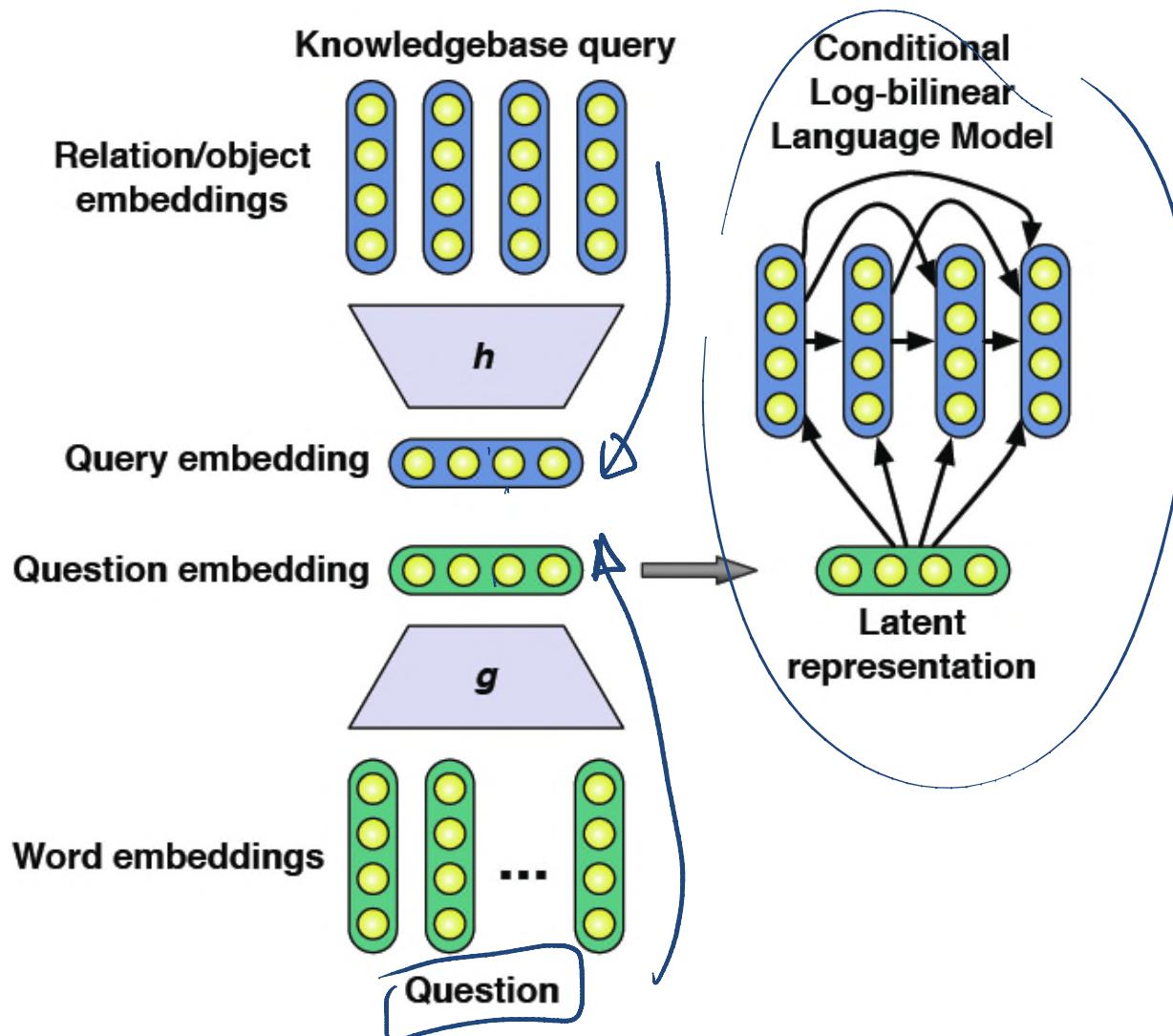
labels

no labels.

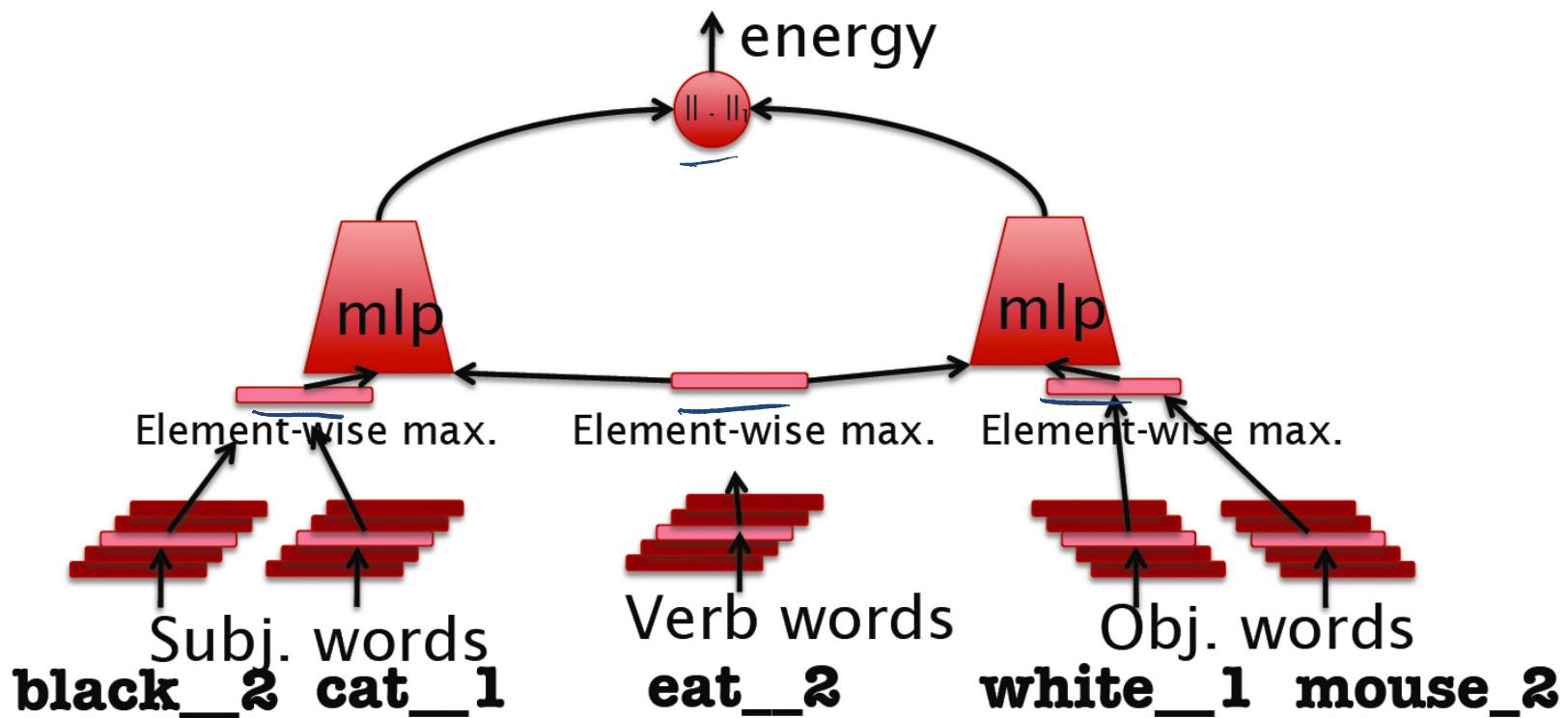# Applications: Paraphrase detection



[Phil Blunsom et al]

# Transfer: Question answering



[Grefenstette, Blunsom, NdF, Hermann, 2014]

# Relation learning



*Learning Structured Embeddings of Knowledge Bases*, Bordes, Weston, Collobert & Bengio, AAAI 2011

# Relation learning

- Intuition: if an entity of a triplet was missing, we would like our model to predict it correctly i.e. to give it the lowest energy. For example, this would allow us to answer questions like "what is part of a car?"

- Hence, for any training triplet $x_i = (lhs_i, rel_i, rhs_i)$ we would like:

    (1)     $E(lhs_i, rel_i, rhs_i) < E(lhs_j, rel_i, rhs_i)$,

    (2)     $E(lhs_i, rel_i, rhs_i) < E(lhs_i, rel_j, rhs_i)$,

    (3)     $E(lhs_i, rel_i, rhs_i) < E(lhs_i, rel_i, rhs_j)$,

That is, the energy function E is trained to rank training samples below all other triplets.

[Yoshua Bengio et al]

# Relation learning

Train by stochastic gradient descent:
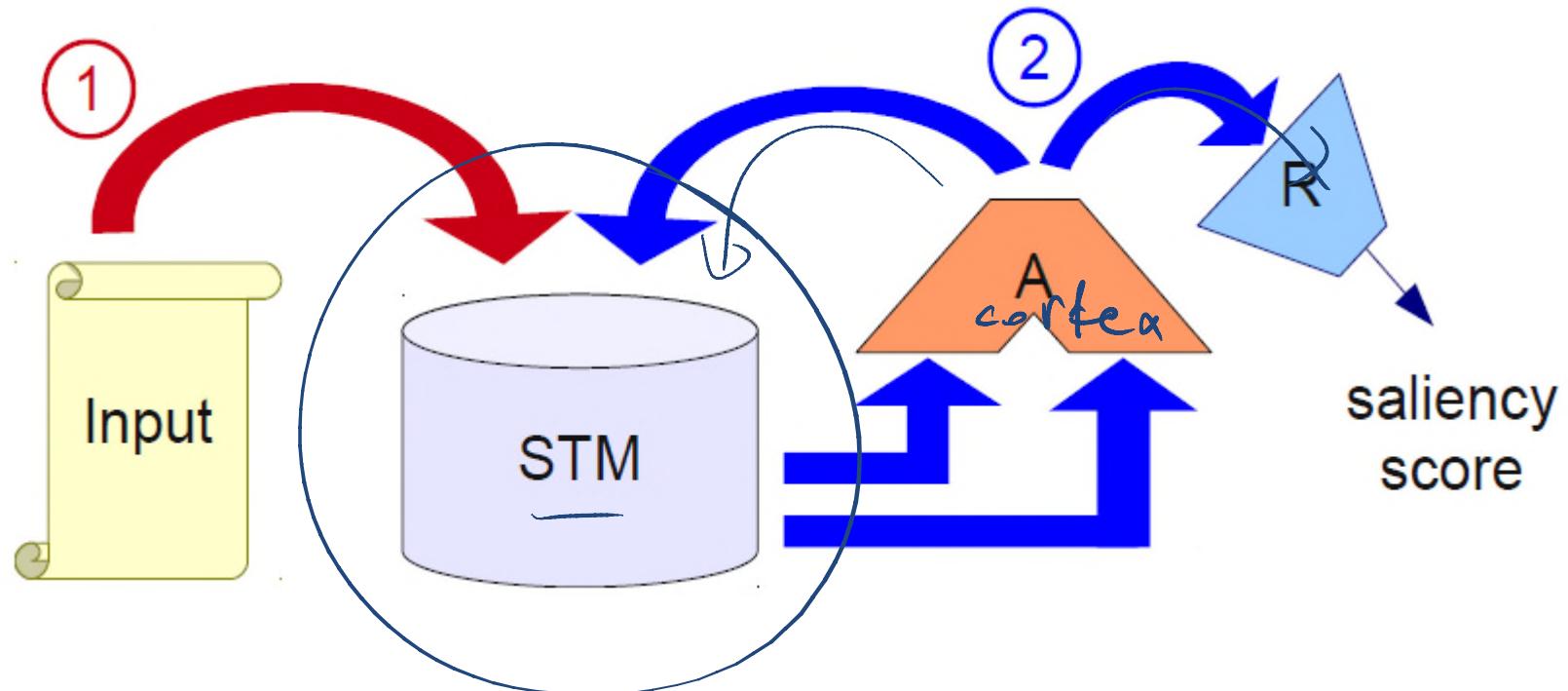
1. Randomly select a positive training triplet $x_i = (\text{lhs}_i, \text{rel}_i, \text{rhs}_i)$.

2. Randomly select constraint (1), (2) or (3) and an entity $\tilde{e}$:

   - If constraint (1), construct negative triplet $\tilde{x} = (\tilde{e}, \text{rel}_i, \text{rhs}_i)$.

   - Else if constraint (2), construct $\tilde{x} = (\text{lhs}_i, \tilde{e}, \text{rhs}_i)$.

   - Else, construct $\tilde{x} = (\text{lhs}_i, \text{rel}_i, \tilde{e})$.

3. If $E(x_i) > E(\tilde{x}) - 1$ make a gradient step to minimize:

   $$\max(0, 1 - E(\tilde{x}) + E(x_i)).$$

4. Constraint embedding vectors to norm 1

[Yoshua Bengio et al]

# Relation learning

| | |
|---|---|
| *lhs* | **_army_NN_1** |
| *rel* | **_attack_VB_1** |
| top ranked *rhs* | _troop_NN_4<br>_armed_service_NN_1<br>_ship_NN_1<br>_territory_NN_1<br>_military_unit_NN_1 |

| | |
|---|---|
| top ranked *lhs* | _business_firm_NN_1<br>_person_NN_1<br>_family_NN_1<br>_payoff_NN_3<br>_card_game_NN_1 |
| *rel* | **_earn_VB_1** |
| *rhs* | **_money_NN_1** |

[Yoshua Bengio et al]

# Short term & long term memory



[*From machine learning to machine reasoning*, Leon Bottou]

# Memory networks

1. Convert $x$ to an internal feature representation $I(x)$.
2. Update memories $\mathbf{m}_i$ given the new input: $\mathbf{m}_i = G(\mathbf{m}_i, I(x), \mathbf{m})$, $\forall i$.
3. Compute output features $o$ given the new input and the memory: $o = O(I(x), \mathbf{m})$.
4. Finally, decode output features $o$ to give the final response: $r = R(o)$.

"slot" in the memory: $\quad \mathbf{m}_{S(x)} = I(x)$

next empty memory slot $N$: $\mathbf{m}_N = x$, $N = N + 1$

[Jason Weston, Sumit Chopra & Antoine Bordes, 2015]

# Memory networks

1. Convert $x$ to an internal feature representation $I(x)$.

2. Update memories $\mathbf{m}_i$ given the new input: $\mathbf{m}_i = G(\mathbf{m}_i, I(x), \mathbf{m})$, $\forall i$.

3. Compute output features $o$ given the new input and the memory: $o = O(I(x), \mathbf{m})$.

4. Finally, decode output features $o$ to give the final response: $r = R(o)$.

$$o_1 = O_1(x, \mathbf{m}) = \arg\max_{i=1,\ldots,N} s_O(x, \mathbf{m}_i)$$

$$o_2 = O_2(x, \mathbf{m}) = \arg\max_{i=1,\ldots,N} s_O([x, \mathbf{m}_{o_1}], \mathbf{m}_i)$$

$$r = \mathrm{argmax}_{w \in W}\ s_R([x, \mathbf{m}_{o_1}, \mathbf{m}_{o_2}], w)$$

[Jason Weston, Sumit Chopra & Antoine Bordes, 2015]

# Memory networks

1. Convert $x$ to an internal feature representation $I(x)$.
2. Update memories $\mathbf{m}_i$ given the new input: $\mathbf{m}_i = G(\mathbf{m}_i, I(x), \mathbf{m})$, $\forall i$.
3. Compute output features $o$ given the new input and the memory: $o = O(I(x), \mathbf{m})$.
4. Finally, decode output features $o$ to give the final response: $r = R(o)$.

$$\sum_{\bar{f} \neq f_1} \max(0, \gamma - s_O(x, f_1^{\overset{good}{}}) + s_O(x, \underset{bad}{\bar{f}})) +$$

$$\sum_{\bar{f}' \neq f_2} \max(0, \gamma - s_O([x, \mathbf{m}_{o_1}], f_2)) + s_O([x, \mathbf{m}_{o_1}], \bar{f}'])) +$$

$$\sum_{\bar{r} \neq r} \max(0, \gamma - s_R([x, \mathbf{m}_{o_1}, \mathbf{m}_{o_2}], r) + s_R([x, \mathbf{m}_{o_1}, \mathbf{m}_{o_2}], \bar{r}]))$$

[Jason Weston, Sumit Chopra & Antoine Bordes, 2015]

# Memory networks

$$\sum_{\bar{f} \neq f_1} \max(0, \gamma - s_O(x, f_1) + s_O(x, \bar{f})) +$$

$$\sum_{\bar{f}' \neq f_2} \max(0, \gamma - s_O([x, \mathbf{m}_{o_1}], f_2) + s_O([x, \mathbf{m}_{o_1}], \bar{f}')) +$$

$$\sum_{\bar{r} \neq r} \max(0, \gamma - s_R([x, \mathbf{m}_{o_1}, \mathbf{m}_{o_2}], r) + s_R([x, \mathbf{m}_{o_1}, \mathbf{m}_{o_2}], \bar{r}))$$

[Jason Weston, Sumit Chopra & Antoine Bordes, 2015]

# Memory networks

Joe went to the garden then Fred picked up the milk; Joe moved to the bathroom and Fred dropped the milk, and then Dan moved to the living_room.
Where is Dan? A: living room I believe
Where is Joe? A: the bathroom
Fred moved to the bedroom and Joe went to the kitchen then Joe took the milk there and Dan journeyed to the bedroom; Joe discarded the milk.
Where is the milk now? A: the milk is in the kitchen
Where is Dan now? A: I think he is in the bedroom
Joe took the milk there, after that Mike travelled to the office, then Joe went to the living_room, next Dan went back to the kitchen and Joe travelled to the office.
Where is Joe now? A: I think Joe is in the office

[Jason Weston, Sumit Chopra & Antoine Bordes, 2015]

# Next lecture

In the next lecture, we will look models for sequences of data.