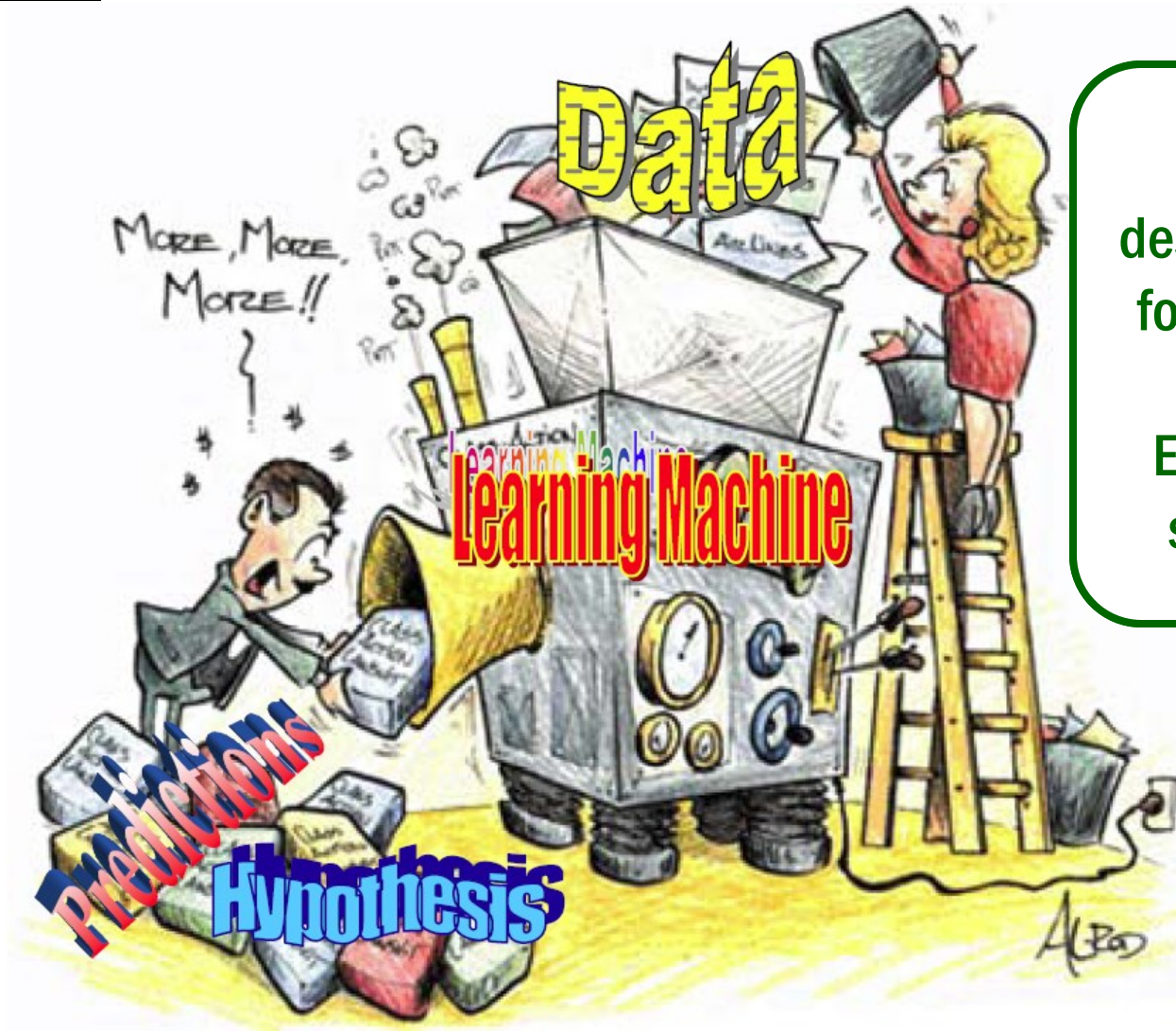# Embedding as a Tool for Algorithm Design

## Le Song

College of Computing

Center for Machine Learning

Georgia Institute of Technology

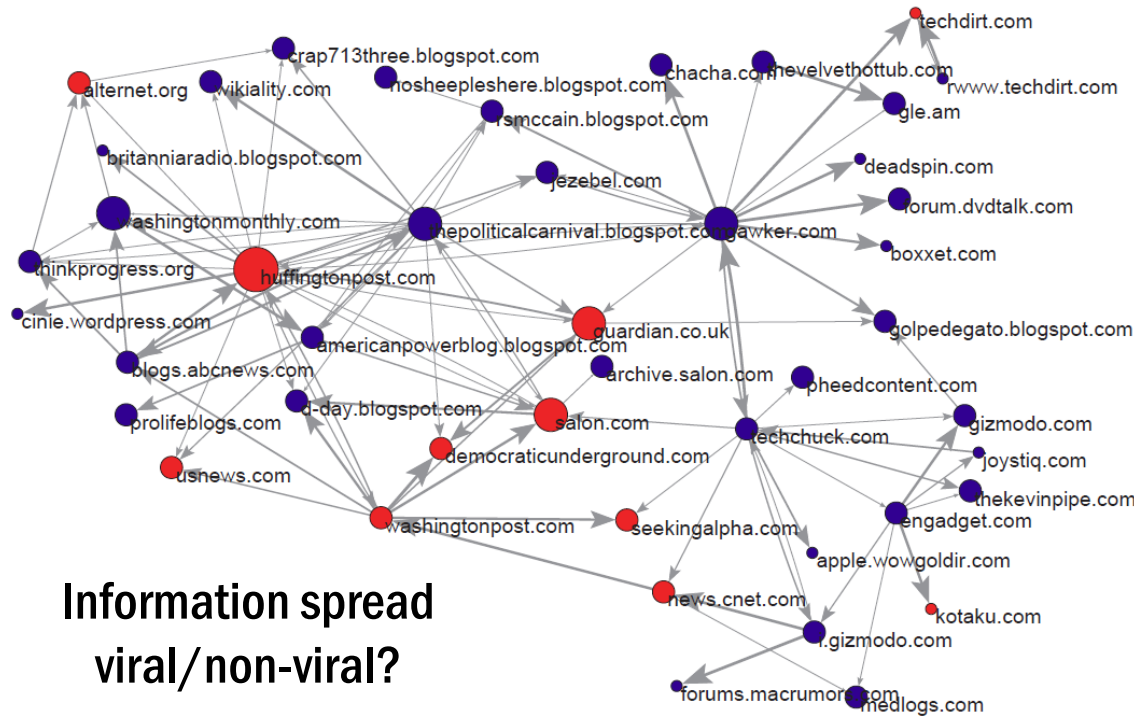# What is machine learning (ML)

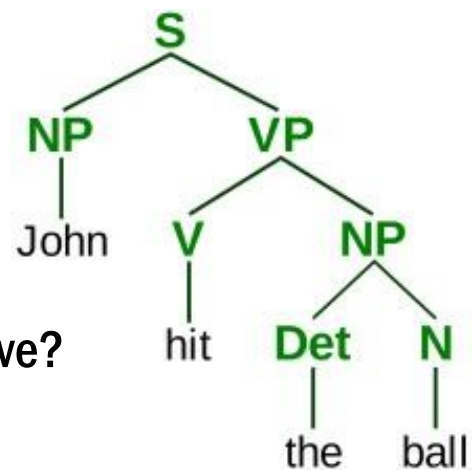Design algorithms and systems that can improve their <u>performance</u> with <u>data</u>



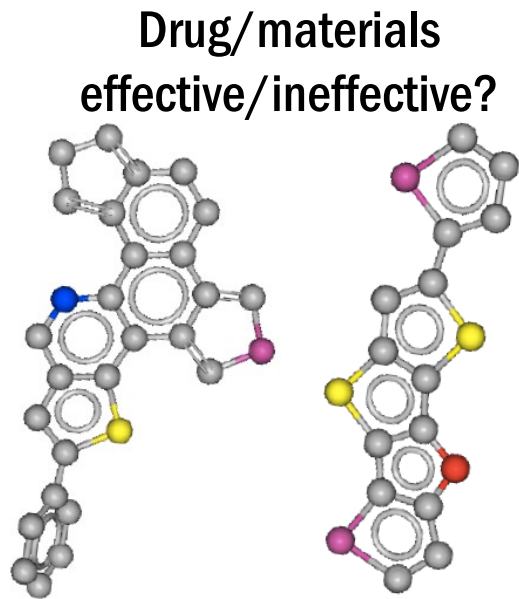The best design pattern for big data?

Embedding structures

# Ex 1: Prediction for structured data

**Drug/materials effective/ineffective?**



**Information spread viral/non-viral?**



**Natural language positive/negative?**



**code graphs benign/ malicious?**



```
mov    [esp+4Ch+var_40], edi
mov    [esp+4Ch+n], 18h
mov    [esp+4Ch+var_3C], edx
mov    edx, [esi]
mov    [esp+4Ch+dest], 0
mov    [esp+4Ch+src], edx
call   eax
```

```
loc_80C1B2B:
cmp    bp, 1
jz     short loc_80C1B88
```

```
xor    eax, eax
cmp    bp, 2
jz     short loc_80C1B48
```

```
loc_80C1B48:
cmp    ebx, 12h
movzx  edx, byte ptr [edi+3]
movzx  ecx, byte ptr [edi+4]
jnz    short loc_80C1B39
```

```
lea    eax, [ebx+13h]
...
mov    [esp+4Ch+src], offset aD1_both_c ;
mov    [esp+4Ch+dest], eax
mov    [esp+4Ch+var_24], eax
call   CRYPTO_malloc
...
mov    [esp+4Ch+dest], ecx ; dest
mov    [esp+4Ch+src], edi ; src
mov    [esp+4Ch+var_20], ecx
call   _memcpy
mov    ecx, [esp+4Ch+var_20]
```
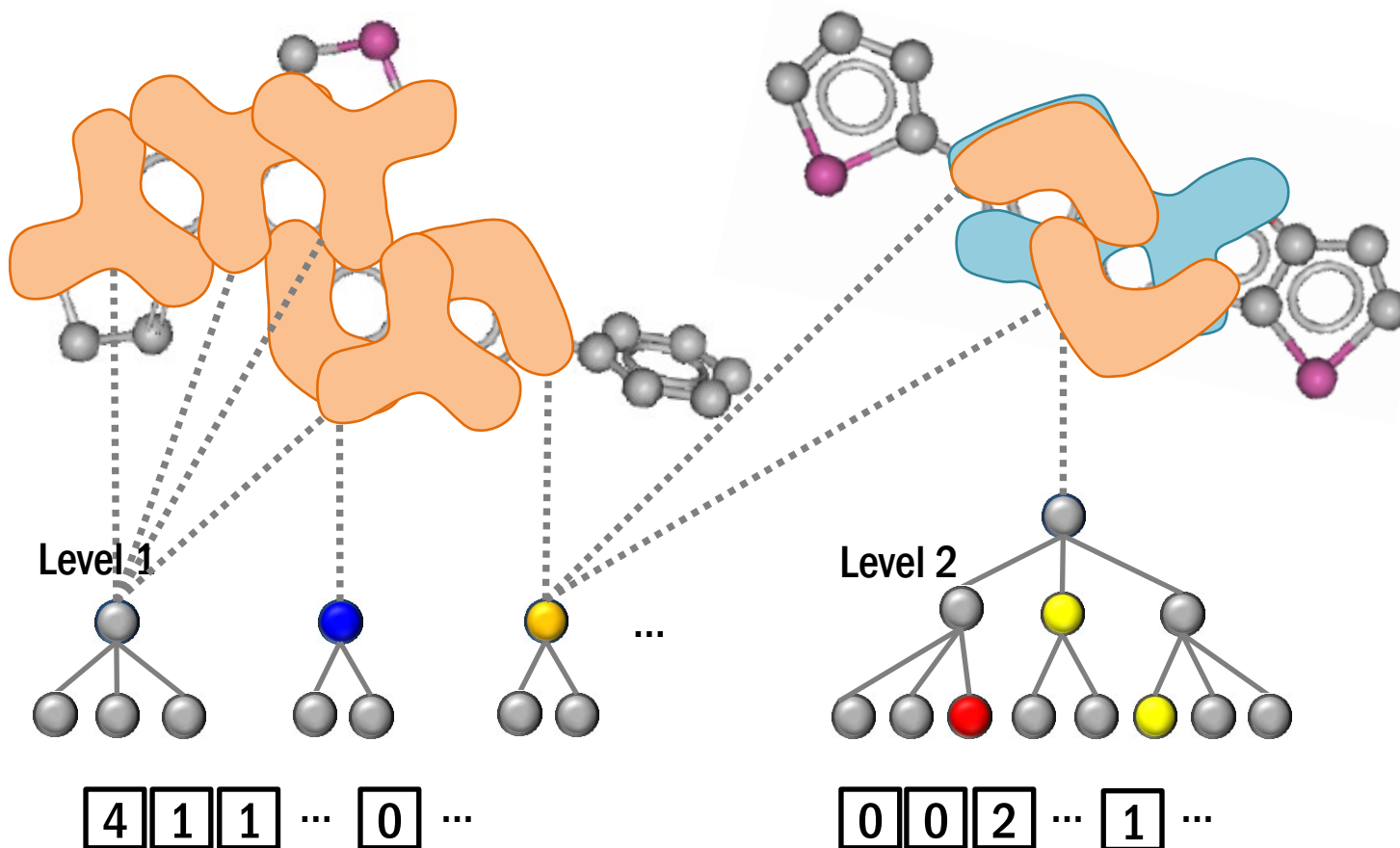
3

# Big dataset, explosive feature space



2.3 million organic materials

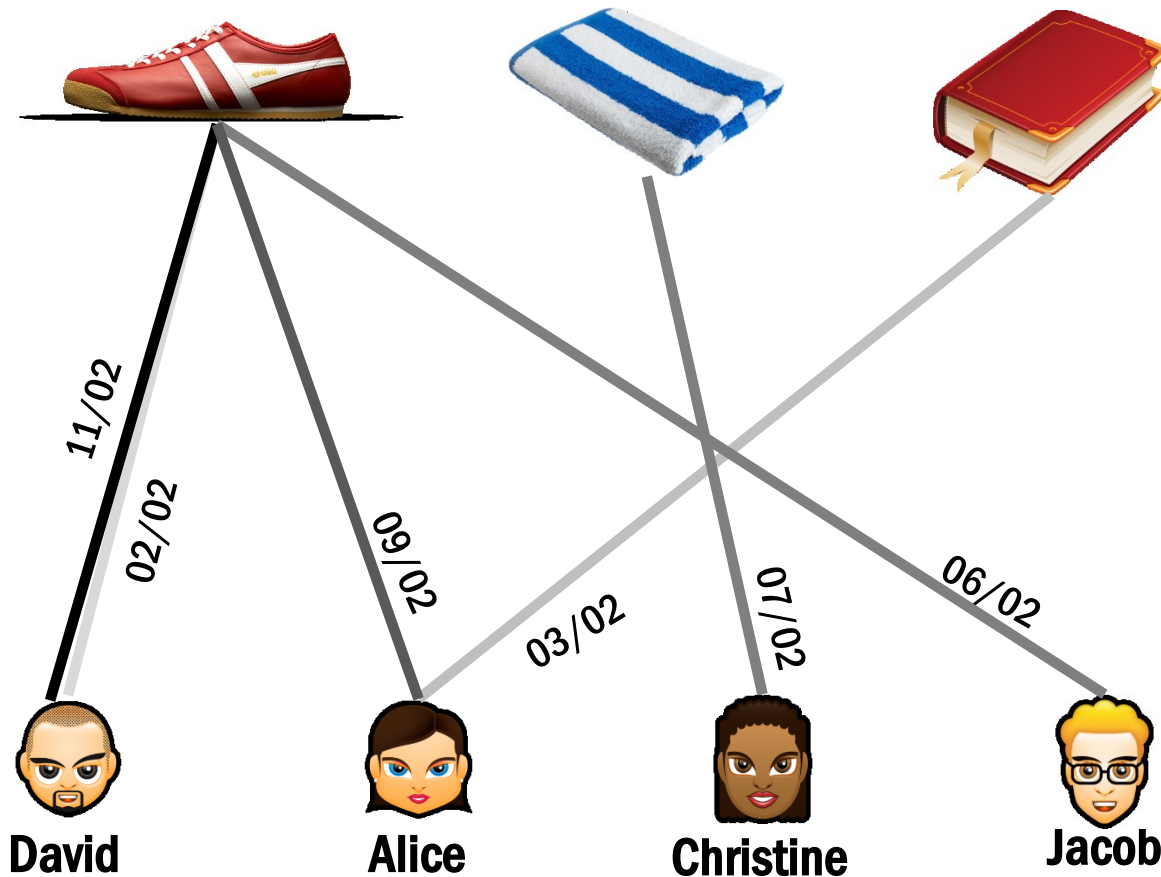Structure elements

Feature vector

Predict

Efficiency (PCE) (0 -12 %)

Level 1

Level 2

4 1 1 ... 0 ...

0 0 2 ... 1 ...

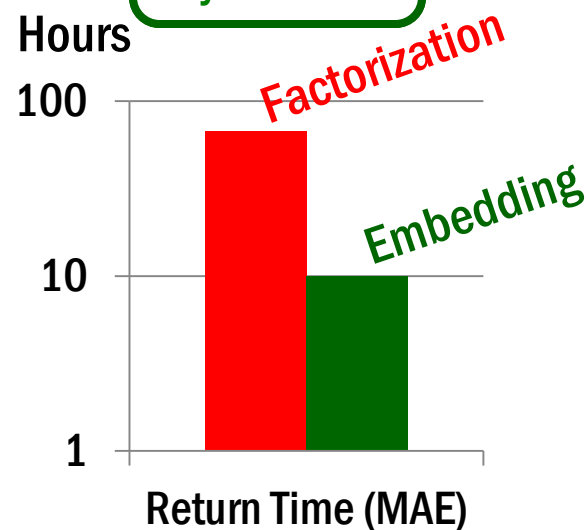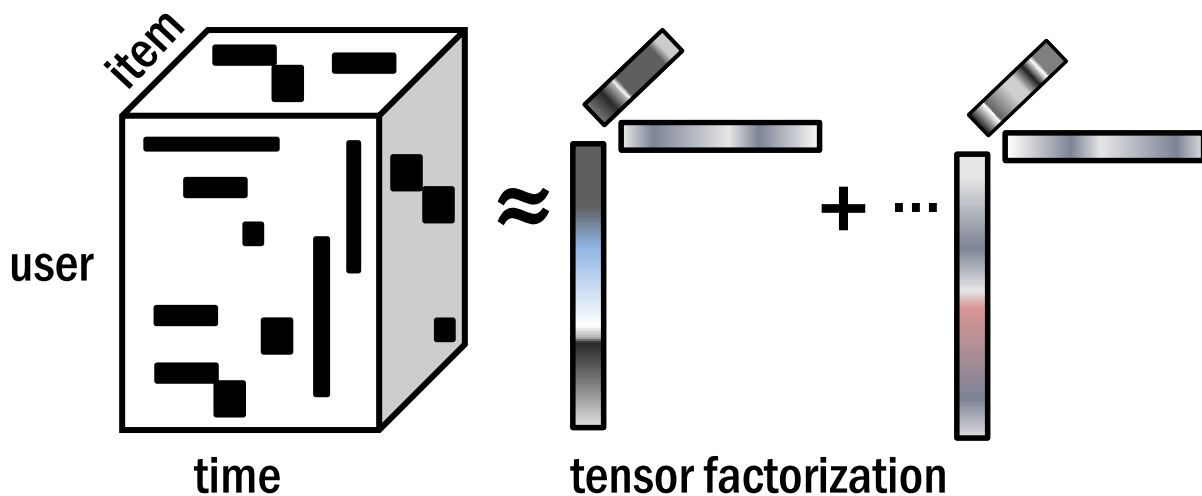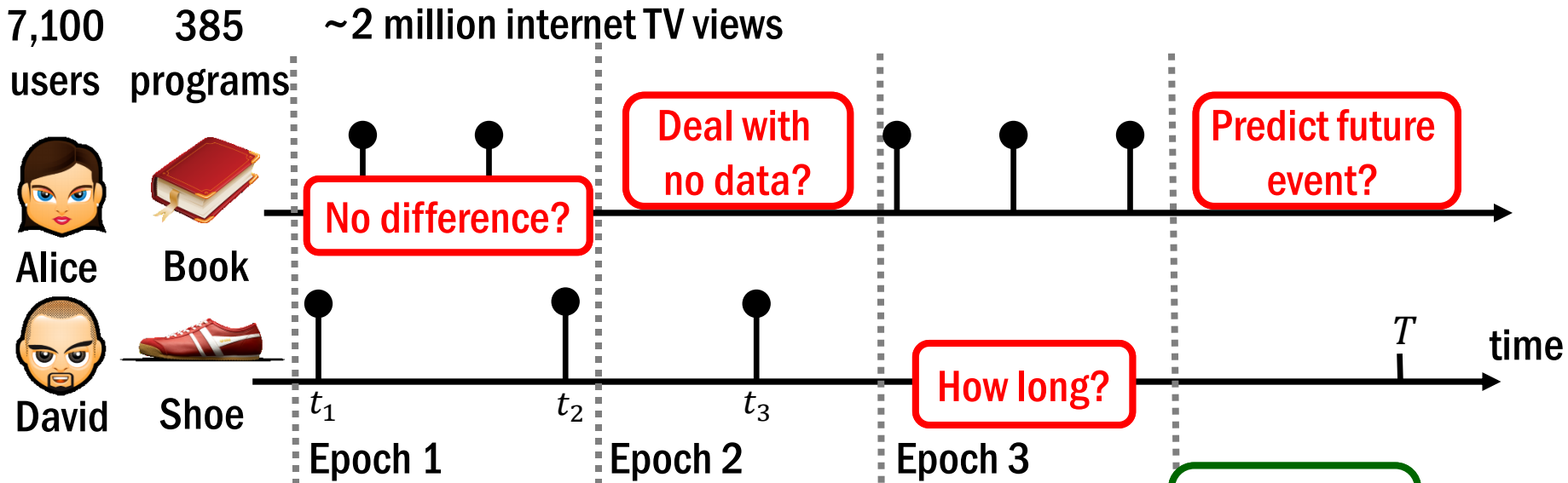| method | dimension | MAE |
|---|---|---|
| Level 6 | 1.3 billion | 0.096 |
| Embedding | 0.1 million | 0.085 |

Reduce model size by 10,000 times!

4

# Ex 2: Social information network modeling

who and when
will do what?

11/02

02/02

09/02

03/02

07/02

06/02

David

Alice

Christine

Jacob

amazon spotify PANDORA last.fm reddit yelp WeChat

# Complex behavior not well modeled

7,100 users    385 programs    ~2 million internet TV views

Alice    Book

David    Shoe

No difference?

Deal with no data?

Predict future event?

How long?

$t_1$    $t_2$    $t_3$    $T$    time

Epoch 1    Epoch 2    Epoch 3

Reduce error by 5 folds!

item    user    time

$\approx$    $+ \cdots$

tensor factorization

Hours

100

10

1

Factorization

Embedding

Return Time (MAE)
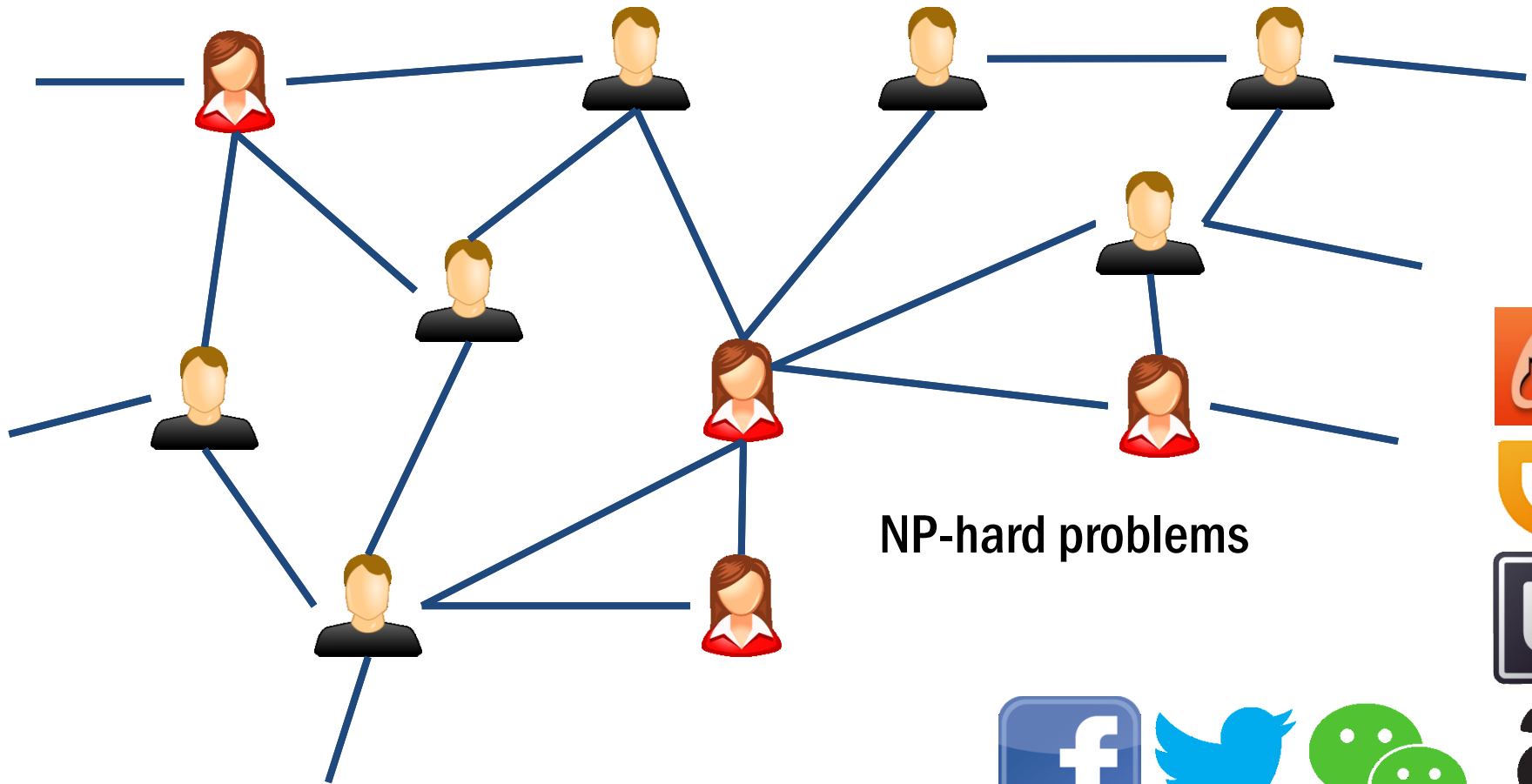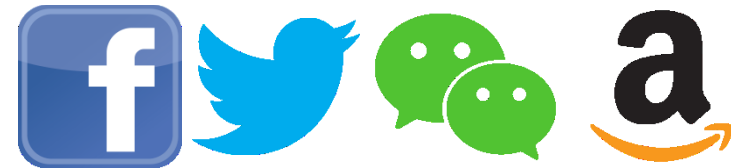
6

# Ex 3: Combinatorial optimizations over graphs

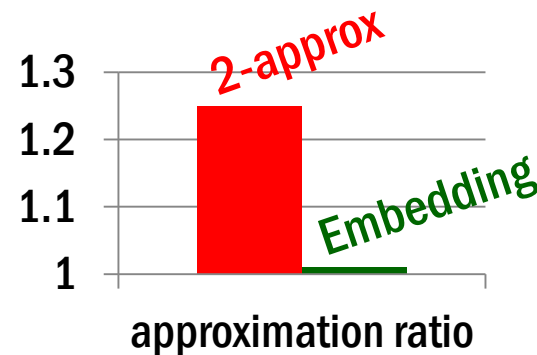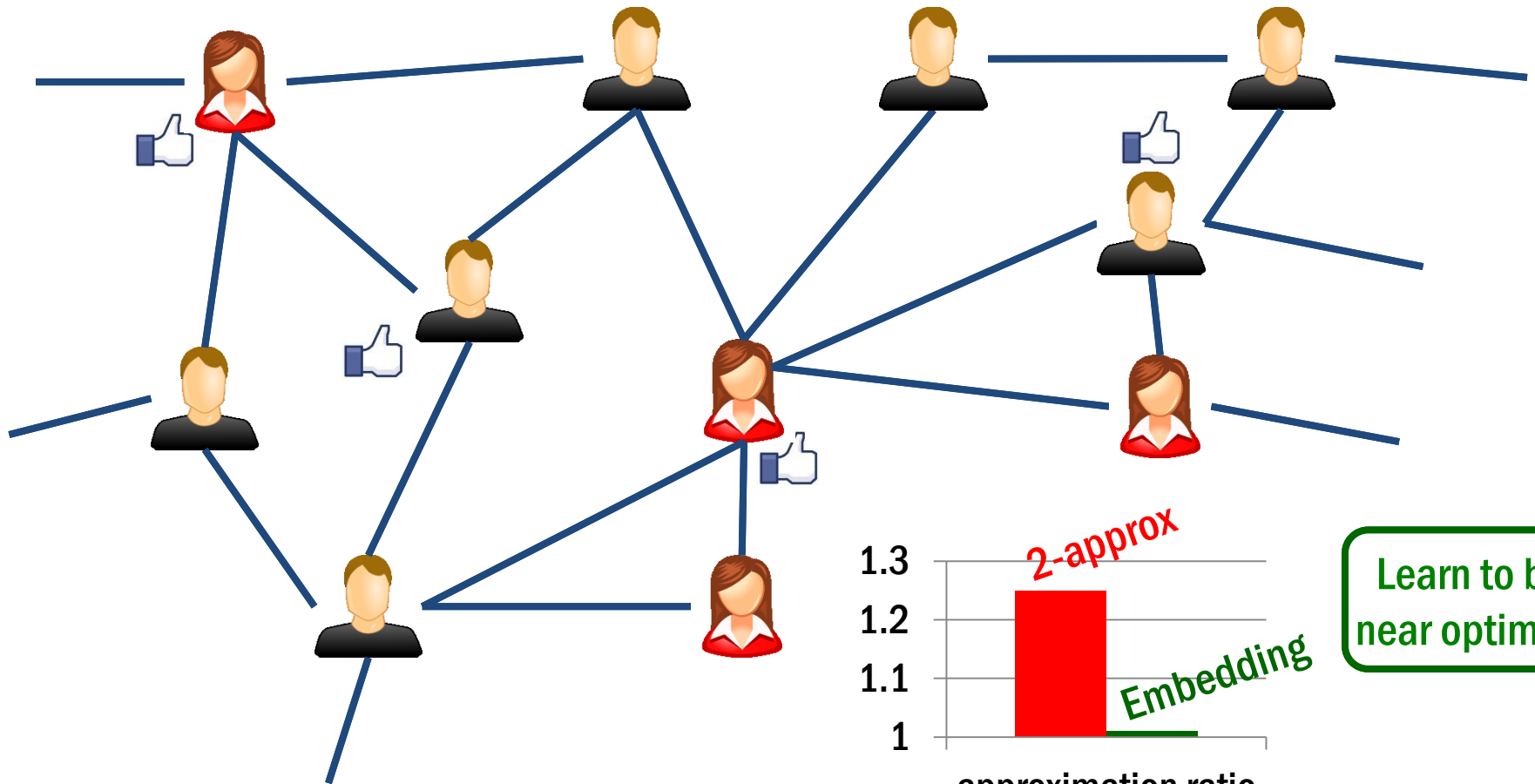| Application | Optimization Problem |
|---|---|
| Influence maximization | Minimum vertex/set cover |
| Community discovery | Maximum cut |
| Resource scheduling | Traveling salesman |



NP-hard problems

# Simple heuristics do not exploit data

**2 - approximation for minimum vertex cover**

Repeat till all edges covered:
1.  Select uncovered edge with largest total degree
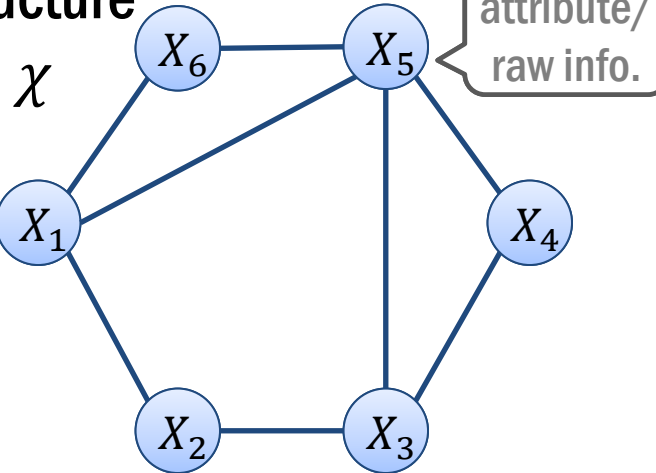
Decision not data-driven. Can we learn from data?



2-approx

Embedding

1.3
1.2
1.1
1

approximation ratio

Learn to be near optimal!

# Fundamental problems

Structure

$\chi$

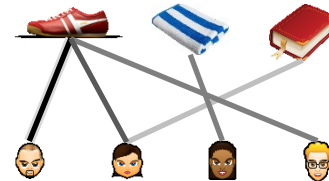attribute/ raw info.

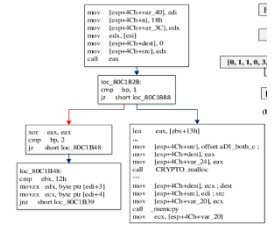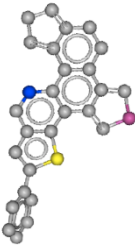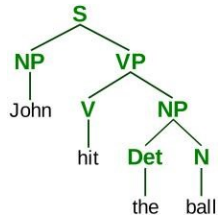Nodes: $X_6$, $X_5$, $X_1$, $X_4$, $X_2$, $X_3$
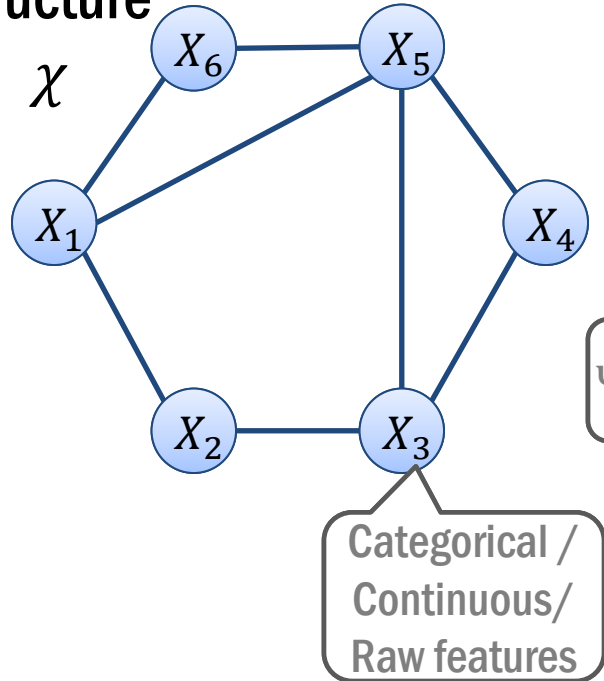
How to describe node?

How to describe entire structure?

How to incorporate various info.?

How to do it efficiently?

# Represent structure as latent variable model (LVM)

**Structure**

$\chi$



**Represent**

**LVM**

$G = (\mathcal{V}, \mathcal{E})$

Continuous Latent

$\Psi_e(H_i, H_j)$

$\Psi_v(H_i, X_i)$

Categorical / Continuous/ Raw features

## Joint likelihood

$$p(\{H_i\}, \{X_i\}) \propto \prod_{i \in \mathcal{V}} \Psi_v(H_i, X_i | \theta_v) \prod_{(i,j) \in \mathcal{E}} \Psi_e(H_i, H_j | \theta_e)$$

Nonnegative node potential

Nonnegative edge potential

[Dai, Dai & Song 2016]

# Posterior distribution as features

**Features of nodes**

$$\mu_1(\chi, W)$$ 

$$+$$

$$\mu_2(\chi, W)$$ 

$$+$$

$$\vdots$$

$$= \mu^a(\chi, W)$$ 

**Features of the entire structure**

$$p(H_1|\{x_j\})$$

$$p(H_2|\{x_j\})$$

LVM

$$G = (\mathcal{V}, \mathcal{E})$$

posterior



$$p(H_i|\{x_j\}) = \frac{\sum_{all\ H_j\ except\ H_i} p(\{H_j\}, \{x_j\})}{p(\{x_j\})}$$

Capture both nodal and topological info.
Aggregate information from distant nodes

[Dai, Dai & Song 2016]

# Mean field algorithm aggregates information

**Approximate posterior**

$$p(H_i|\{x_j\}) \approx q_i(H_i)$$

**via fixed point update**
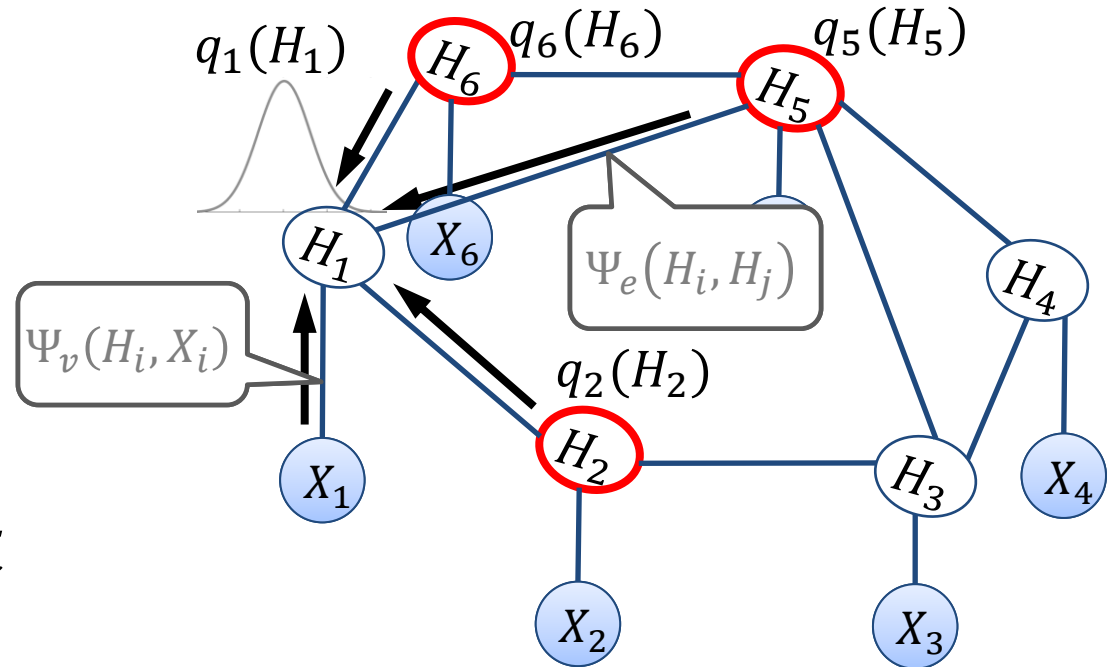
1. Initialize $q_i(H_i), \forall\, i$

2. Iterate many times

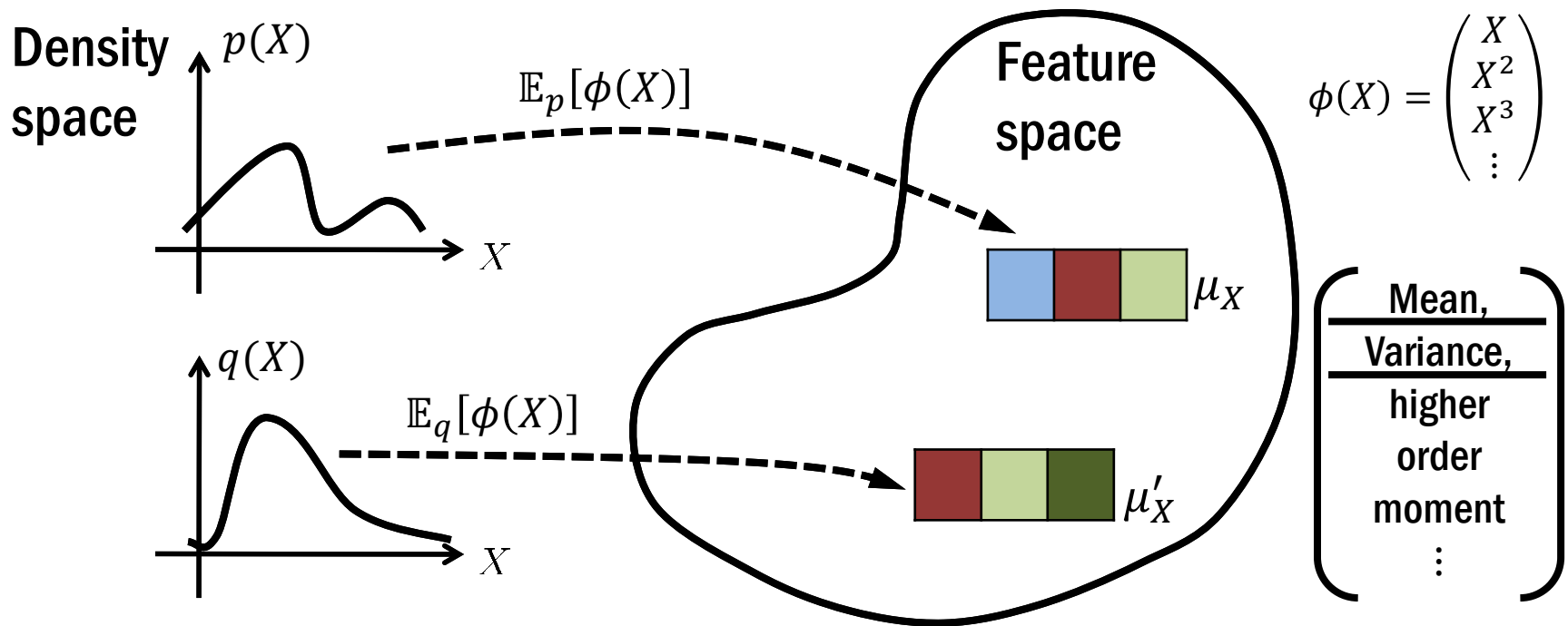$$q_i(H_i) \leftarrow \Psi_v(H_i, X_i) \cdot$$

$$\prod_{j \in \mathcal{N}(i)} \exp\left( \int_{\mathcal{H}} q_j(H_j) \log(\Psi_e(H_i, H_j))\, dH_j \right), \forall\, i$$

$$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad}$$

$$\mathcal{T} \circ \left( X_i, \{q_j(H_j)\}_{j \in \mathcal{N}(i)} \right)$$

$q_1(H_1)$   $q_6(H_6)$   $q_5(H_5)$
$H_6$   $H_5$
$\Psi_e(H_i, H_j)$
$H_1$   $X_6$   $H_4$
$\Psi_v(H_i, X_i)$
$q_2(H_2)$
$X_1$   $H_2$   $H_3$   $X_4$
$X_2$   $X_3$

# Embedding of distribution



**Density space**

$p(X)$

$X$

$\mathbb{E}_p[\phi(X)]$

$q(X)$

$X$

$\mathbb{E}_q[\phi(X)]$

**Feature space**

$\mu_X$

$\mu'_X$

$$\phi(X) = \begin{pmatrix} X \\ X^2 \\ X^3 \\ \vdots \end{pmatrix}$$

$$\begin{pmatrix} \text{Mean,} \\ \text{Variance,} \\ \text{higher} \\ \text{order} \\ \text{moment} \\ \vdots \end{pmatrix}$$

**Injective for rich nonlinear feature $\phi(x)$**

$\mu_X$ **is a sufficient statistic of** $p(X)$

Operator View
$$\mathcal{T} \circ p(x) = \widetilde{\mathcal{T}} \circ \mu_X$$
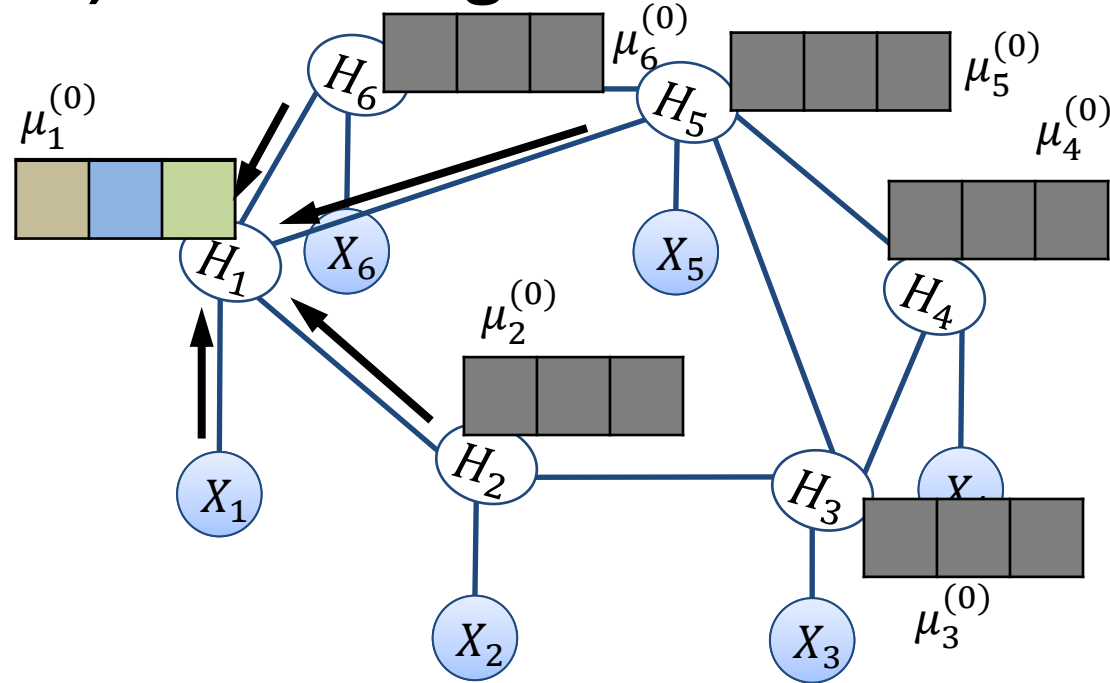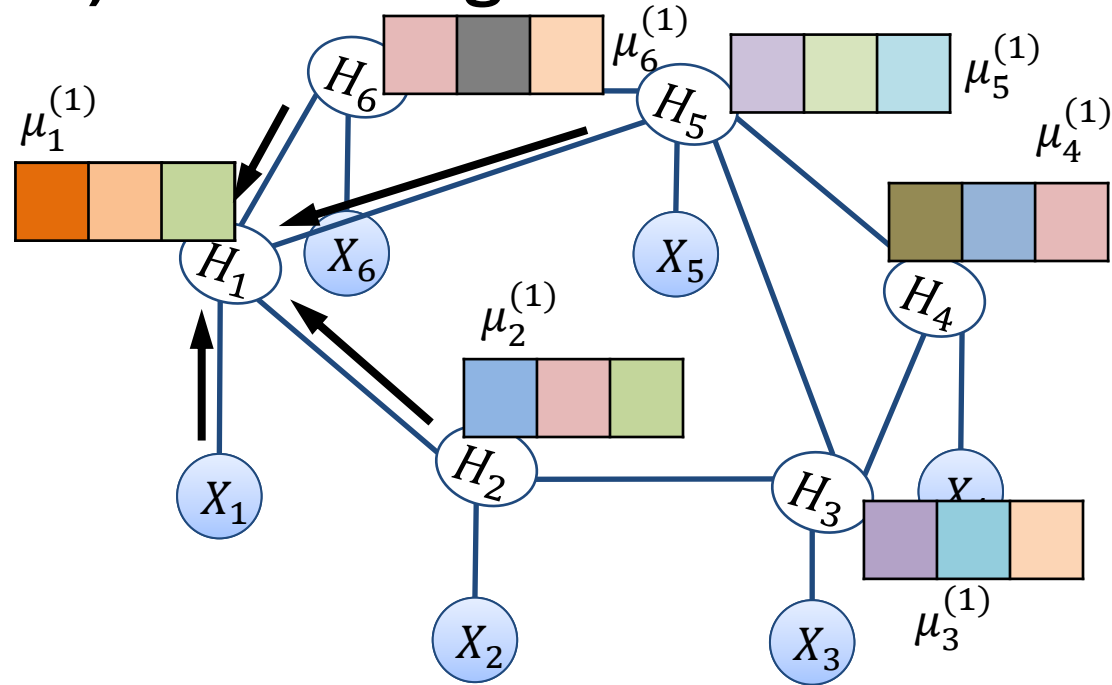
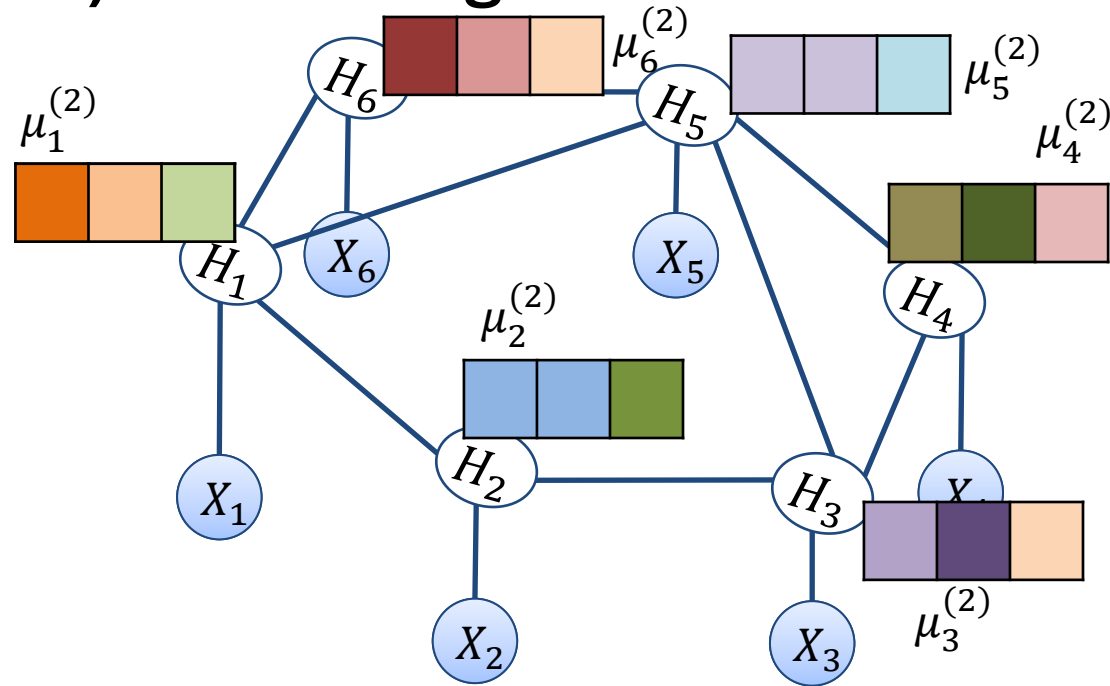[Smola, Gretton, Song & Scholkopf. 2007]

# Structure2vec (S2V): embedding mean field

Approximate embedding of

$$p(H_i | \{x_j\}) \mapsto \mu_i$$

via fixed point update

1. Initialize $\mu_i, \forall\, i$

2. Iterate many times



$$\mu_i \leftarrow \widetilde{\mathcal{T}} \circ \left( X_i, \{\mu_j\}_{j \in \mathcal{N}(i)} \right), \forall\, i$$

# Structure2vec (S2V): embedding mean field

Approximate embedding of

$$p\big(H_i \big| \{x_j\}\big) \mapsto \mu_i$$

via fixed point update

1. Initialize $\mu_i, \forall\, i$

2. Iterate many times

$$\mu_i \leftarrow \widetilde{\mathcal{T}} \circ \Big( X_i, \{\mu_j\}_{j \in \mathcal{N}(i)} \Big), \forall\, i$$

# Structure2vec (S2V): embedding mean field

Approximate embedding of

$$p\big(H_i\big|\{x_j\}\big) \mapsto \mu_i$$

via fixed point update

1. Initialize $\mu_i, \forall\, i$

2. Iterate many times

$$\mu_i \leftarrow \widetilde{\mathcal{T}} \circ \Big( X_i, \{\mu_j\}_{j\in\mathcal{N}(i)} \Big), \forall\, i$$

How to parametrize $\widetilde{\mathcal{T}}$?
Depends on unknown $\Psi_v(H_i, X_i)$ and $\Psi_e(H_i, H_j)$

# Directly parameterize nonlinear mapping

$$\mu_i \leftarrow \widetilde{\mathcal{T}} \circ \left( X_i, \{\mu_j\}_{j \in \mathcal{N}(i)} \right)$$

**Any universal nonlinear function will do**

**Eg. assume $\mu_i \in \mathcal{R}^d, X_i \in \mathcal{R}^n$, neural network parameterization**

$$\mu_i \leftarrow \sigma \left( W_1 X_i + W_2 \sum_{j \in \mathcal{N}(i)} \mu_j \right)$$

max$\{0, \cdot\}$
tanh$(\cdot)$
sigmoid$(\cdot)$

$d \times n$
matrix

$d \times d$
matrix

Learn with supervision, unsupervised learning, or reinforcement learning

# Embedding belief propagation

Approximate $p\big(H_i\big|\{x_j\},\theta\big)$ as

$$q_i(H_i) = \Psi_v(H_i, x_i|\theta) \cdot$$

$$\prod_{j \in \mathcal{N}(i)} m_{ji}(H_i)$$

$G = (\mathcal{V}, \mathcal{E})$

$\Psi_v(H_i, X_i)$
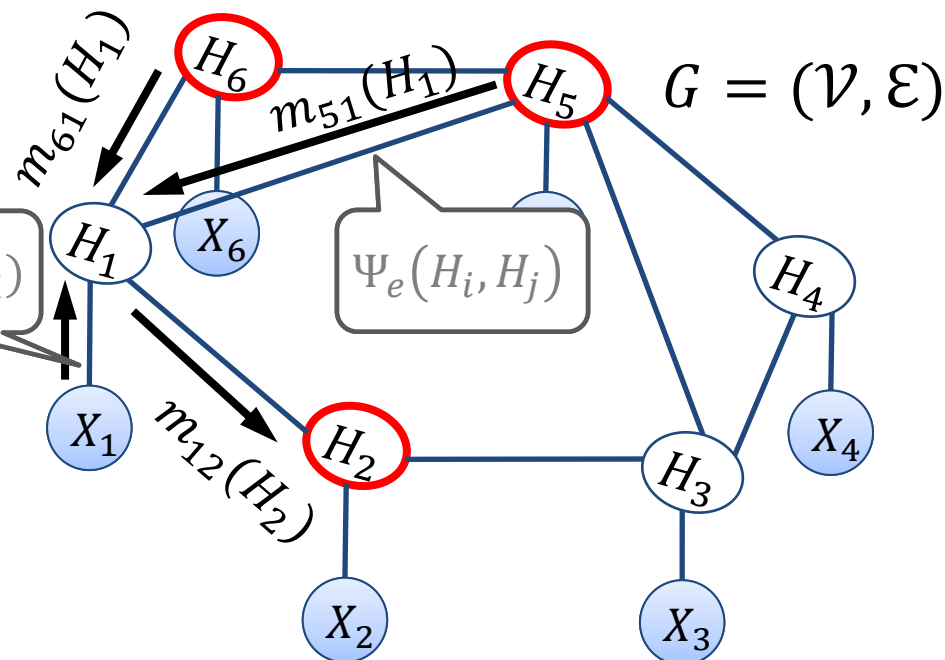
$\Psi_e(H_i, H_j)$

$\boxed{\mathcal{T}' \circ \big( X_i, \{m_{\ell i}(H_i)\}_{\ell \in \mathcal{N}(i)} \big)}$
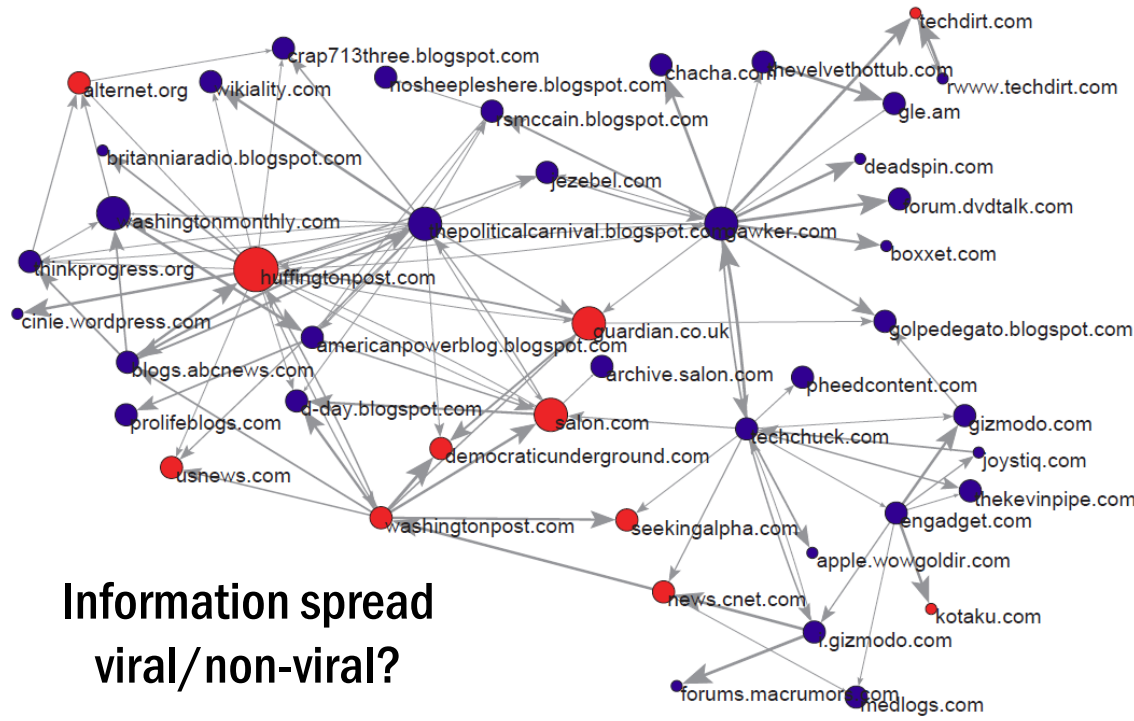
1. Initialize $m_{ij}(H_j), \forall i, j$

2. Iterate many times

$$m_{ij}(H_j) \leftarrow \int_{\mathcal{H}} \Psi_v(H_i, X_i|\theta)\Psi_e(H_i, H_j|\theta) \cdot \prod_{\ell \in \mathcal{N}(i)\backslash j} m_{\ell i}(H_i)\, dH_i, \forall i, j$$

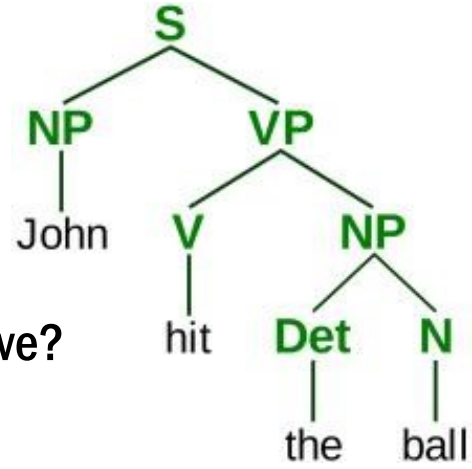$$\boxed{\mathcal{T} \circ \big( X_i, \{m_{\ell i}(H_i)\}_{\ell \in \mathcal{N}(i)\backslash j} \big)}$$



[Song et al. 11a,b]
[Song et al. 10a,b]

# Ex 1: Prediction for structured data

**Drug/materials effective/ineffective?**



**Information spread viral/non-viral?**



**Natural language positive/negative?**



**code graphs benign/malicious?**

# Algorithm learning

Given $m$ data points $\{\chi_1, \chi_2, \dots, \chi_m\}$

And their labels $\{y_1, y_2, \dots, y_m\}$

Estimate parameters $W$ and $V$ via

$$\min_{V,W} \; L(V, W) := \sum_{i=1}^{m} (y_i - V^\top \mu^a(W, \chi_i))^2$$

| Computation | Operation | Similar to |
|---|---|---|
| Objective $L(V, W)$ | A sequence of nonlinear mappings over graph | Graphical model inference |
| Gradient $\dfrac{\partial L}{\partial W}$ | Chain rule of derivatives in reverse order | Back propagation in deep learning |

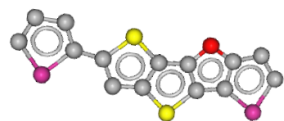# 10,000x smaller model but accurate prediction

Harvard clean energy project: predict material efficiency (0-12)
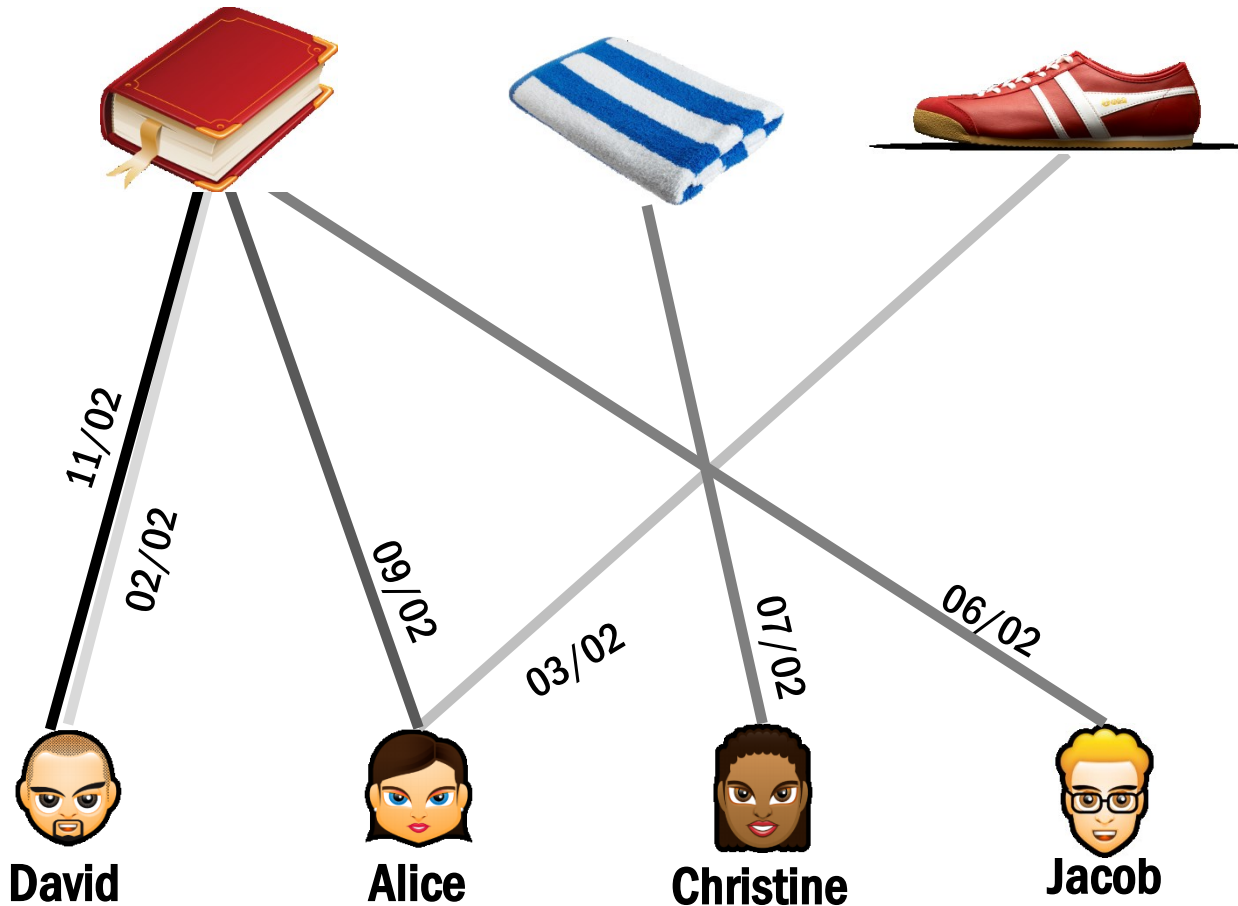2.3 million organic molecules
90% for training, 10% data for testing

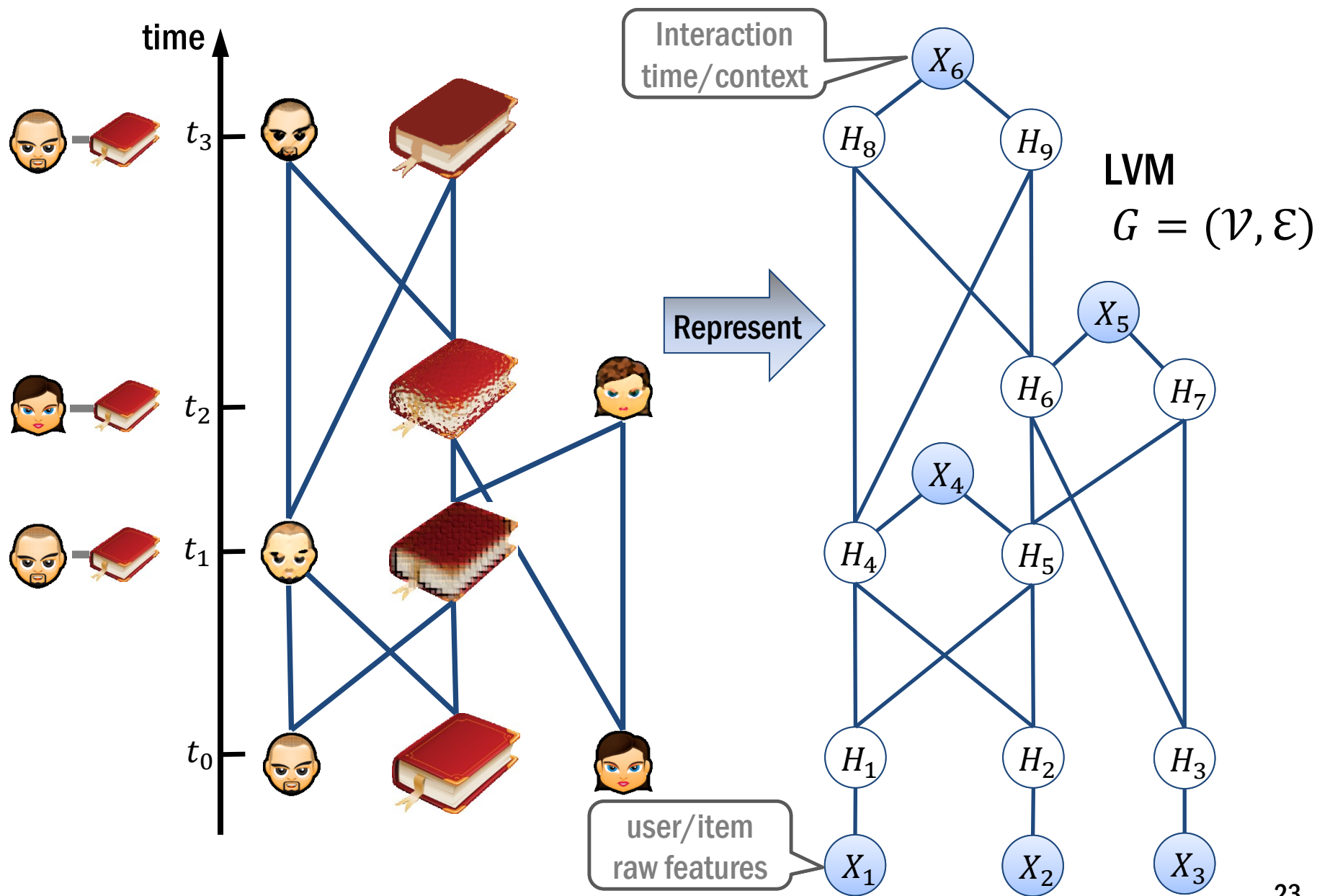| | Test MAE | Test RMSE | # parameters |
|---|---|---|---|
| Mean predictor | 1.986 | 2.406 | 1 |
| WL level-3 | 0.143 | 0.204 | 1.6 m |
| WL level-6 | 0.096 | 0.137 | 1.3 b |
| S2V-MF | 0.091 | 0.125 | 0.1 m |
| S2V-BP | 0.085 | 0.117 | 0.1 m |

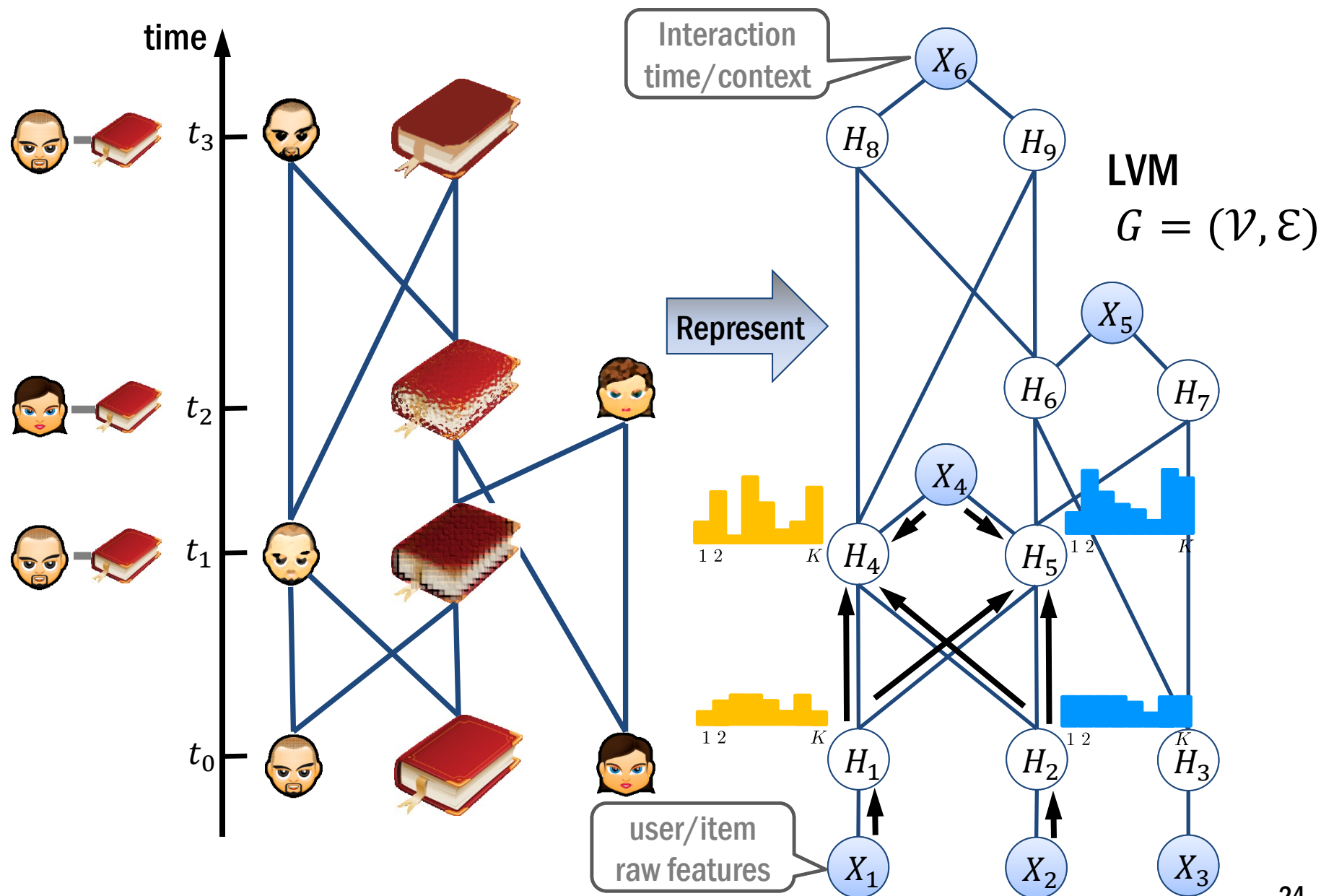~4% relative error

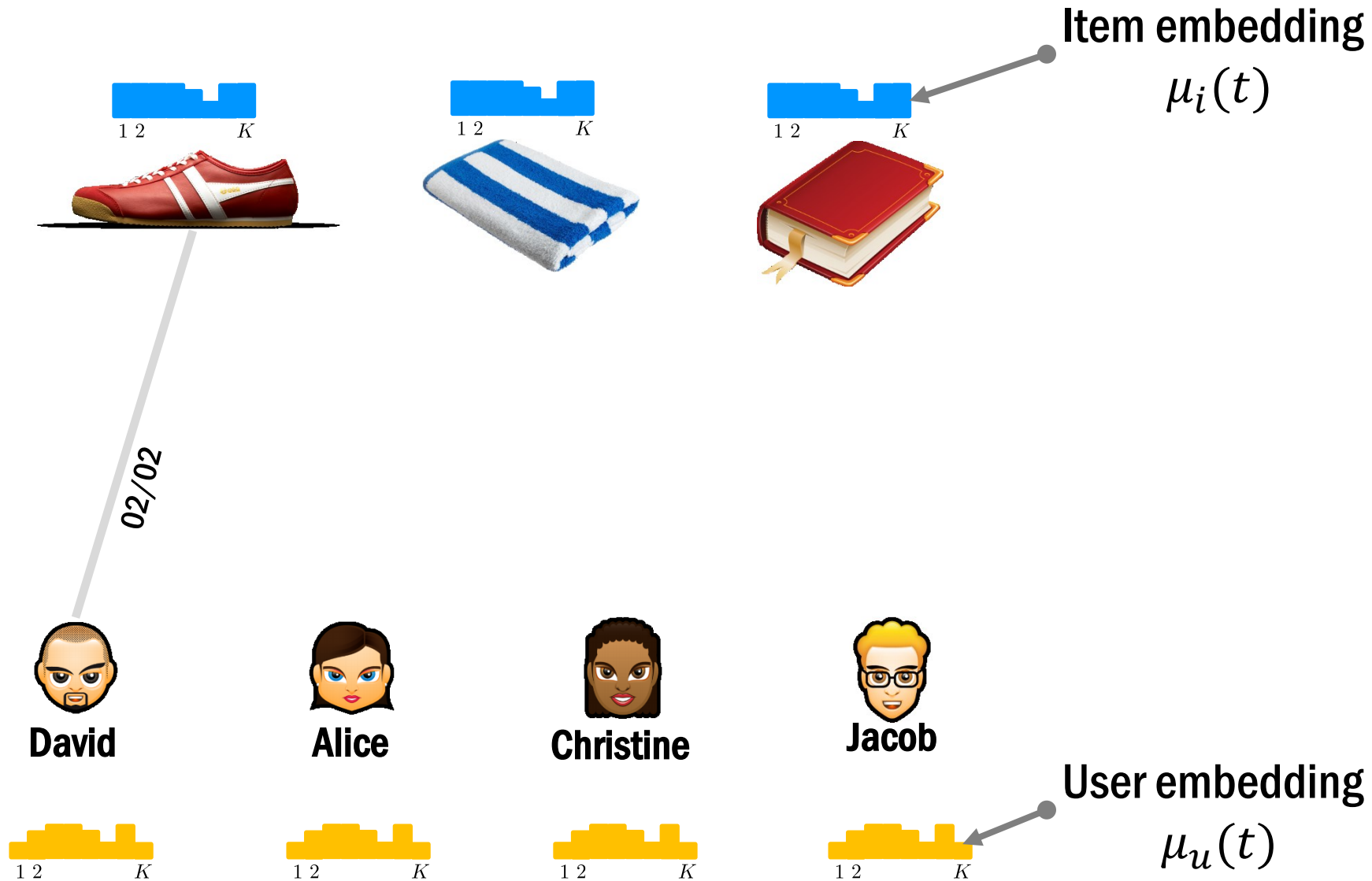# Ex 2: Social information network modeling



who and when
will do what?

11/02

02/02

09/02

03/02

07/02

06/02

David

Alice

Christine

Jacob

# Unroll: time-varying dependency structure



time

$t_3$

$t_2$

$t_1$

$t_0$

Interaction time/context

$X_6$

$H_8$  $H_9$

**LVM**

$G = (\mathcal{V}, \mathcal{E})$

Represent

$X_5$

$H_6$  $H_7$

$X_4$

$H_4$  $H_5$

$H_1$  $H_2$  $H_3$

user/item raw features

$X_1$  $X_2$  $X_3$

23

# Embed filtering/forward belief propagation



time

$t_3$

$t_2$

$t_1$

$t_0$

Represent

Interaction time/context

$X_6$

$H_8$    $H_9$

LVM

$G = (\mathcal{V}, \mathcal{E})$

$X_5$

$H_6$    $H_7$

$X_4$

$H_4$    $H_5$

1 2    K      1 2    K

$H_1$    $H_2$    $H_3$

1 2    K      1 2    K

user/item raw features

$X_1$    $X_2$    $X_3$

# Co-evolutionary embedding



Item embedding $\mu_i(t)$

02/02

David   Alice   Christine   Jacob

User embedding $\mu_u(t)$

# Co-evolutionary embedding



Item embedding $\mu_i(t)$

02/02

03/02

David  Alice  Christine  Jacob

User embedding $\mu_u(t)$

# Co-evolutionary embedding



Item embedding
$\mu_i(t)$

User embedding
$\mu_u(t)$

David  Alice  Christine  Jacob

02/02  03/02  06/02

# Co-evolutionary embedding



Item embedding $\mu_i(t)$

02/02

03/02

07/02

06/02

David

Alice

Christine

Jacob

User embedding $\mu_u(t)$

# Co-evolutionary embedding



Item embedding $\mu_i(t)$

User embedding $\mu_u(t)$

02/02

09/02

03/02

07/02

06/02

David

Alice

Christine

Jacob

# Co-evolutionary embedding



Item embedding
$\mu_i(t)$

Update
U→I

Update
I→U

02/02

09/02

03/02

07/02

06/02

David

Alice

Christine

Jacob

User embedding
$\mu_u(t)$

# From embedding to next interaction time

Link embedding with interaction data using generative model



Intensity of interaction determined by compatibility and time-lapse

$$\lambda_{ui}(t|t_n) = \boxed{\exp\left(\mu_u^\top(t_n)\mu_i(t_n)\right)} \cdot \boxed{(t - t_n)}$$

**Density function**

$$p_{ui}(t|t_n) = \lambda_{ui}(t|t_n)\, S_{ui}(t|t_n)$$

**Survival function**

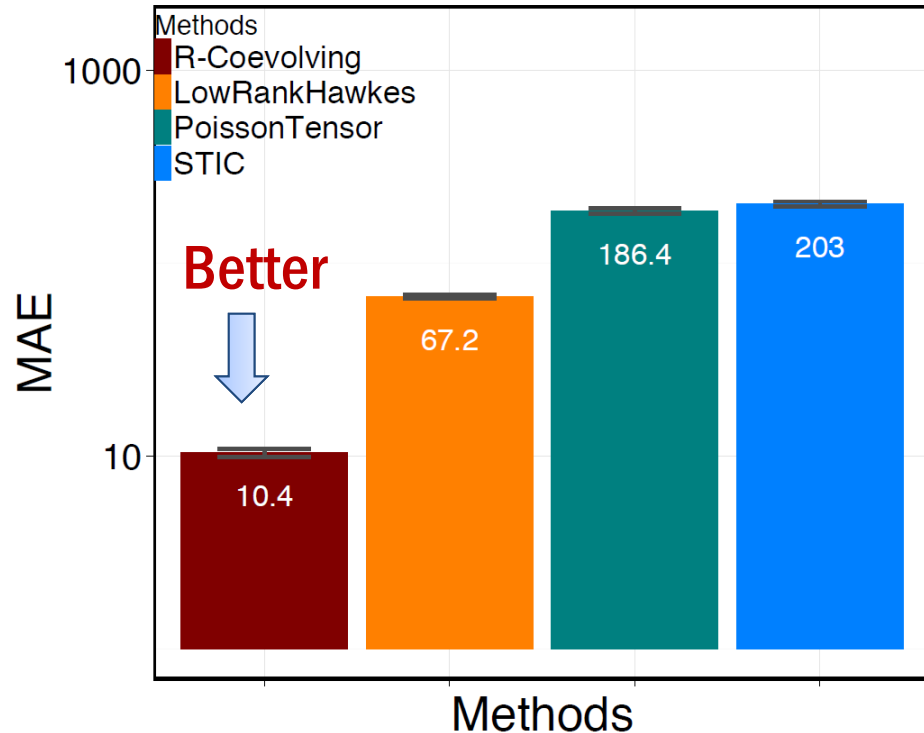$$S_{ui}(t|t_n) = \exp\left(-\int_{t_n}^{t} \lambda_{ui}(\tau)d\tau\right)$$

# Embedding leads to better prediction

Reddit dataset: prediction of discussion forum participation
1,000 users, 1403 groups, ~10K interactions
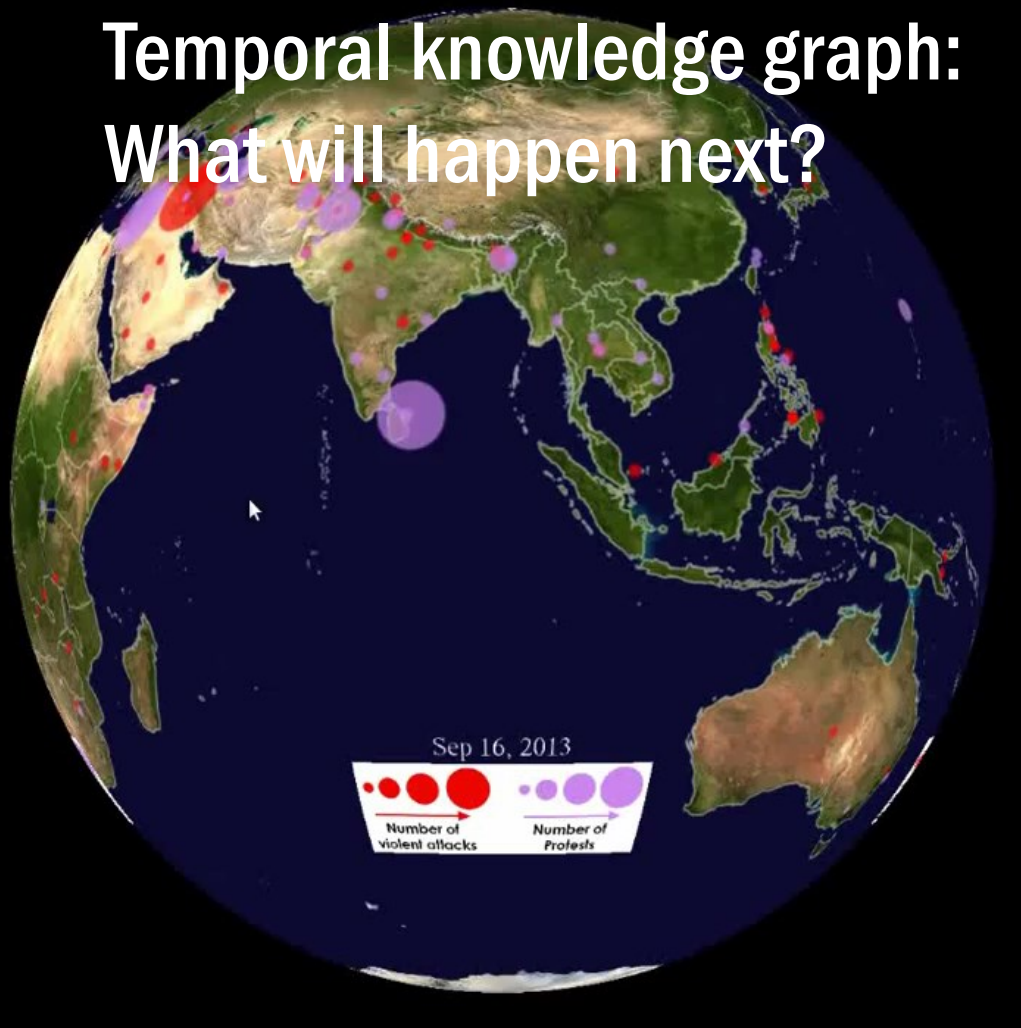


**Next item prediction**
**MAR: mean absolute rank difference**

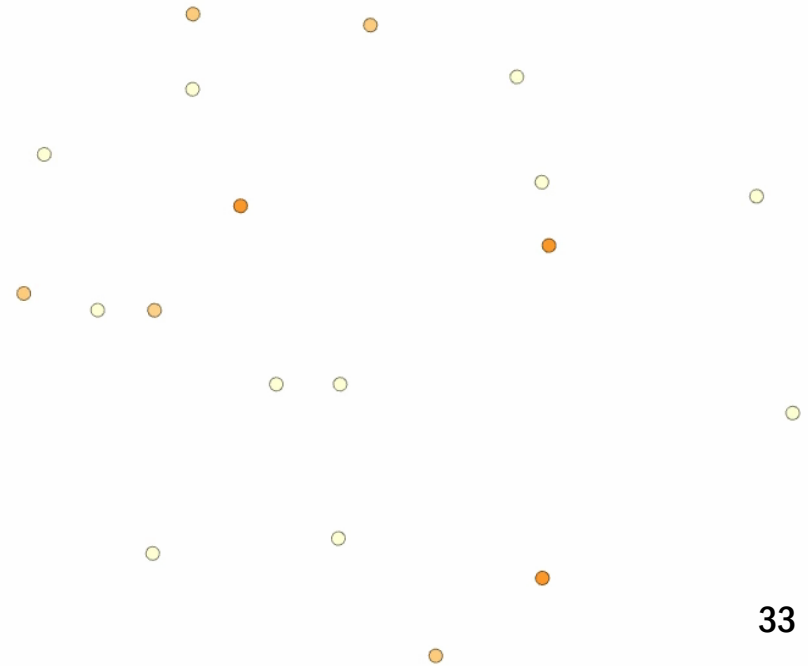**Return time prediction**
**MAE: mean absolute error (hours)**

# Temporal knowledge graph: What will happen next?



Sep 16, 2013

Number of violent attacks — Number of Protests

GDELT database:

Events in news media

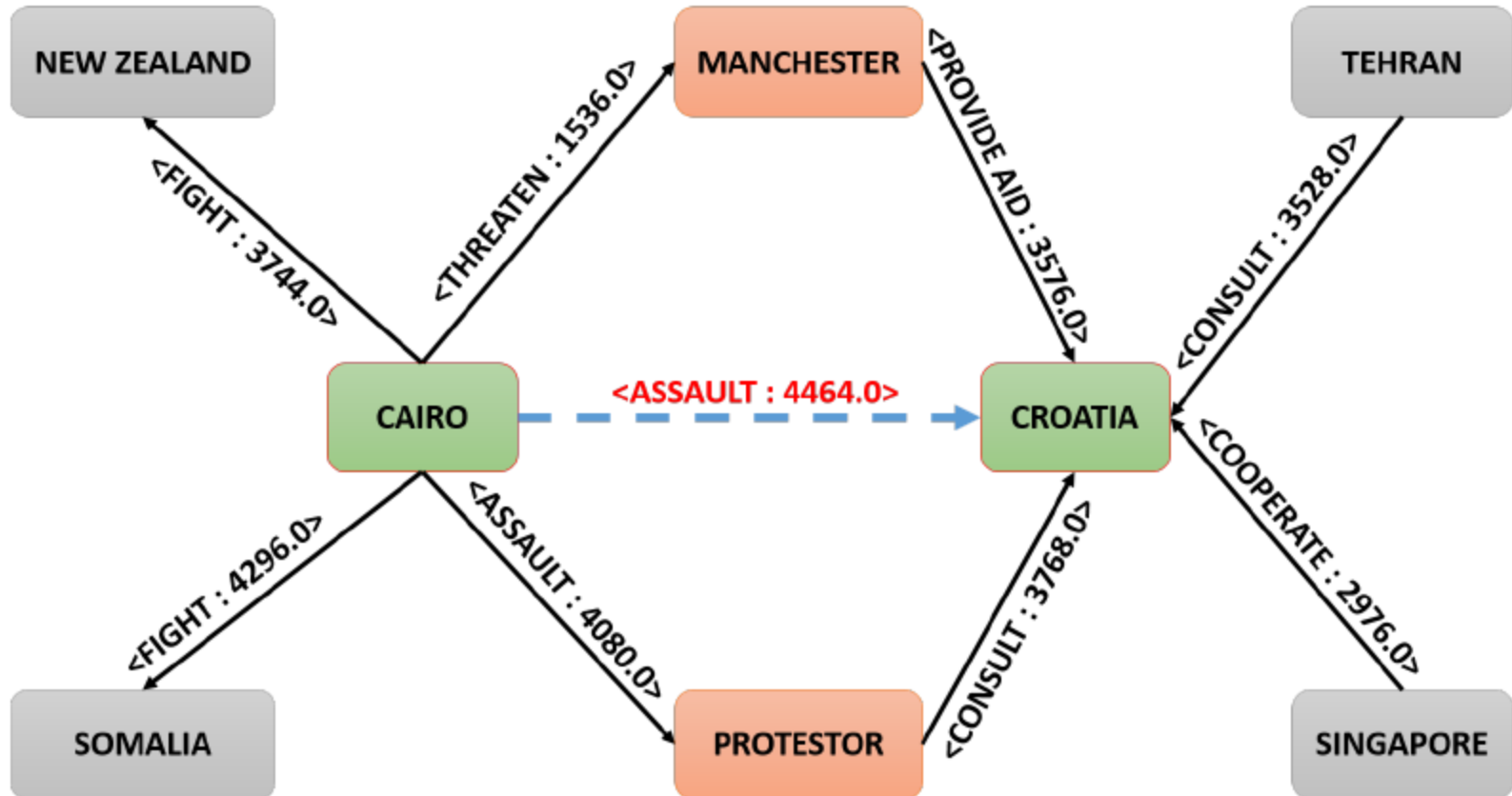Total archives span >215 years, trillion of events

Event (knowledge item):
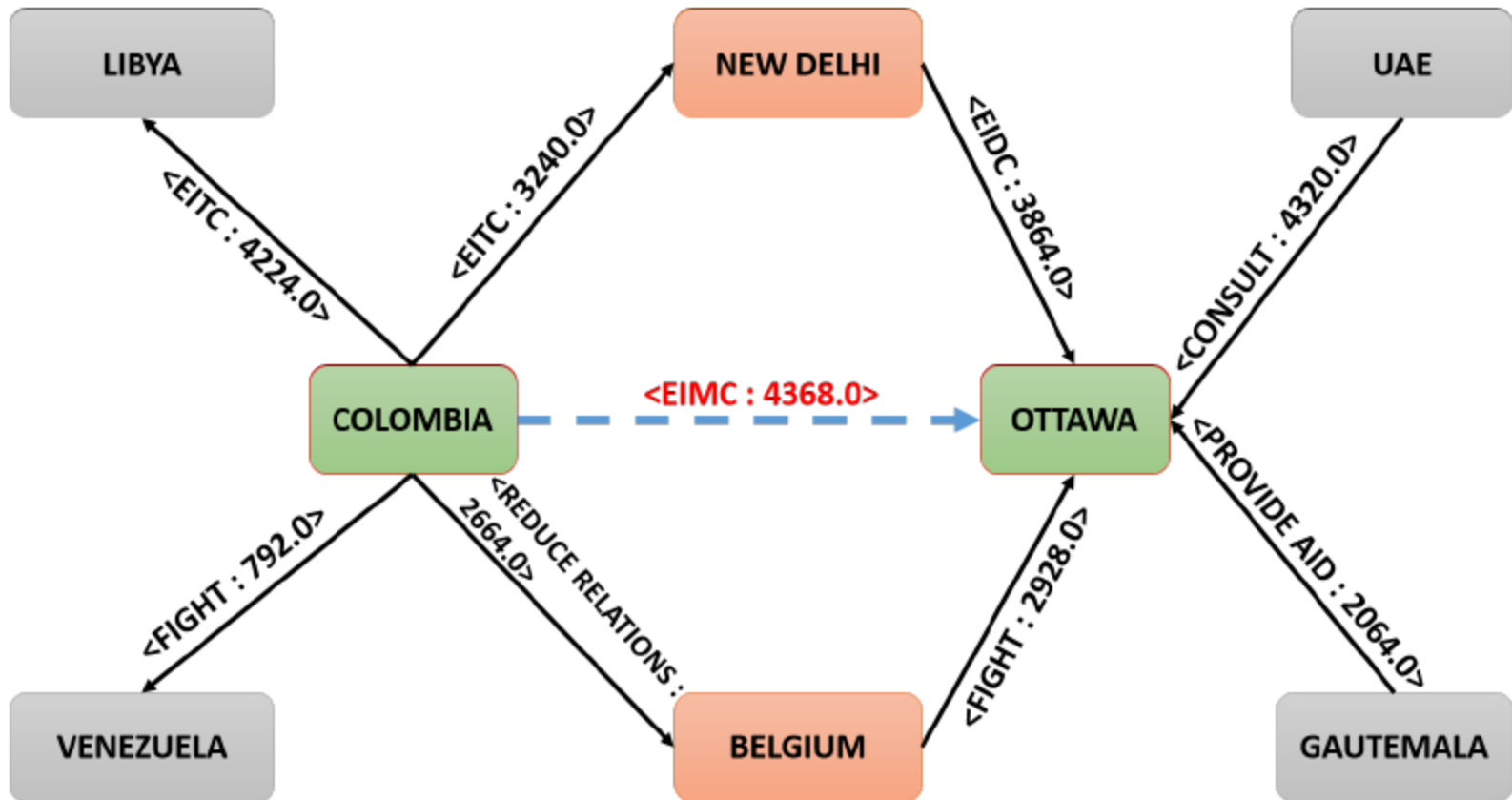- Subject --- relation --- object
- Time

# Reasoning over time I
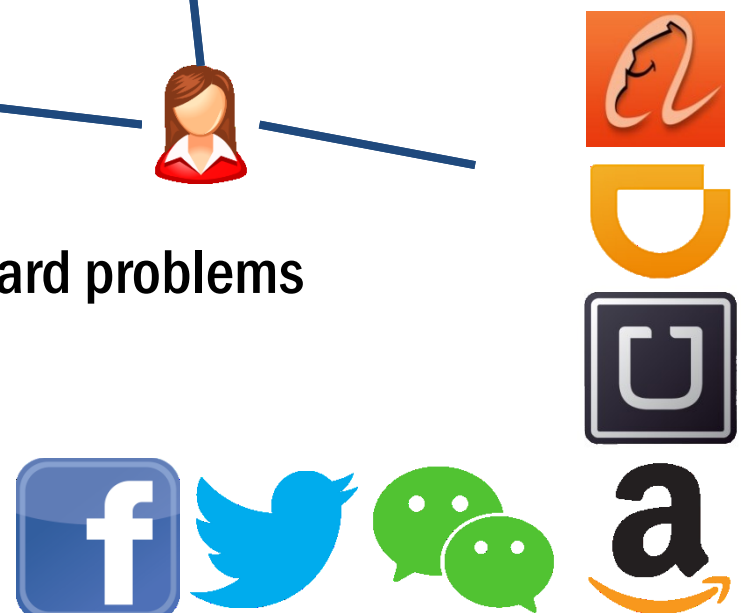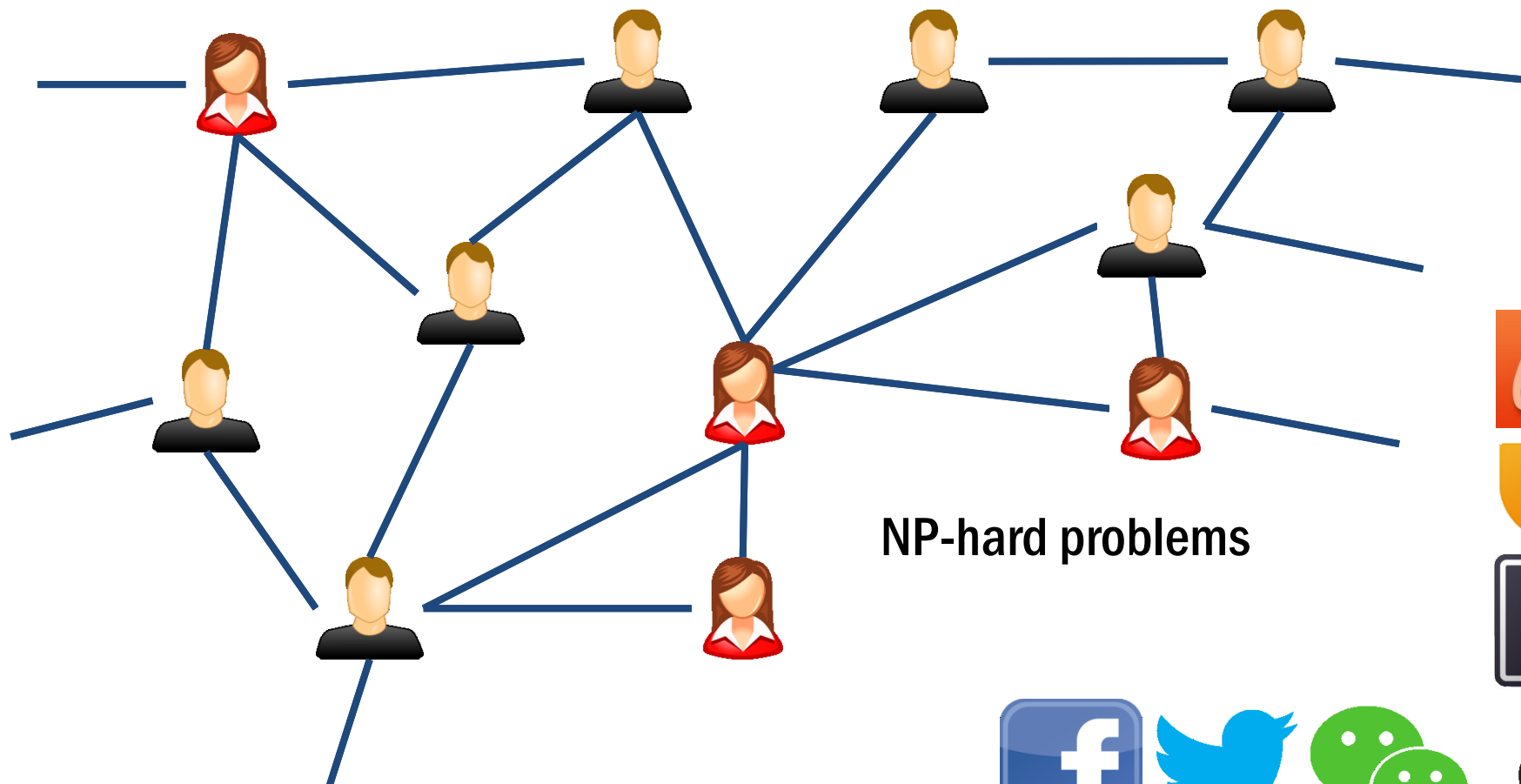
Enemy's friend is enemy

# Reasoning over time II

Friends' friend is a friend, common enemy improves bond

EITC / EIDC / EIMC: some form of cooperation

# App 3: Combinatorial optimizations over graphs

| Application | Optimization Problem |
|---|---|
| Influence maximization | Minimum vertex/set cover |
| Community discovery | Maximum cut |
| Resource scheduling | Traveling salesman |

NP-hard problems

# Combinatorial optimization as MDP

Minimum vertex cover: smallest number of nodes to cover all edges

$$\min_{x_i \in \{0,1\}} \sum_{i \in \mathcal{V}} x_i$$

$$s.t. \, x_i + x_j > 0, \forall (i,j) \in \mathcal{E}$$

Repeat:

1. Compute total degree of each uncovered edge

2. Select both ends of uncovered edge with largest total degree

Until all edges are covered

multistage decision making problem
$$r^t = \sum_{i \in \mathcal{V}} x_i^t - \sum_{i \in \mathcal{V}} x_i^{t+1} = -1$$
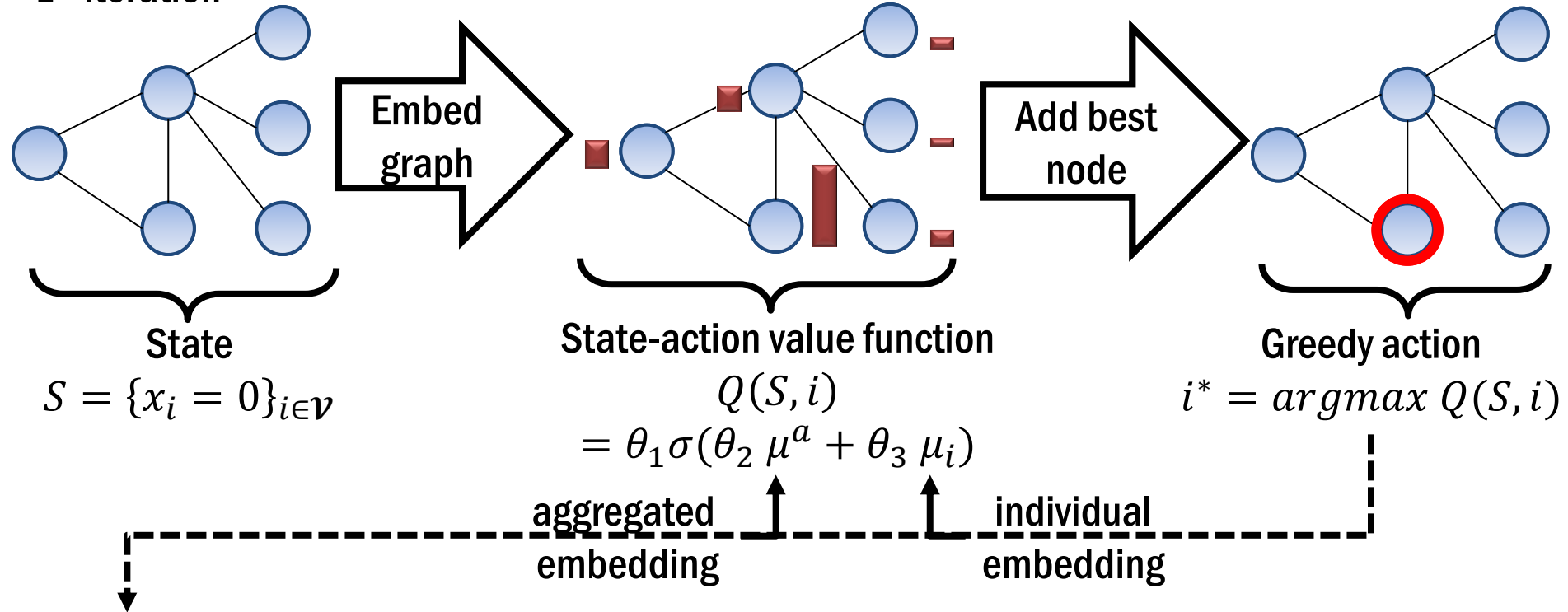
State $S$: current set of nodes selected

Action value function: $Q(S,i)$

Greedy policy:
$$i^* = argmax_i \, Q(S,i)$$

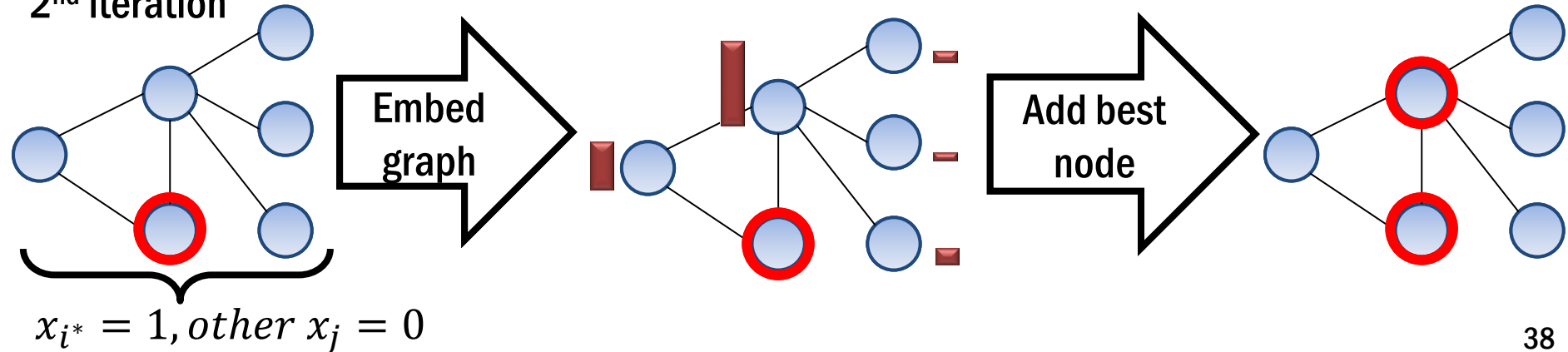Update state $S$

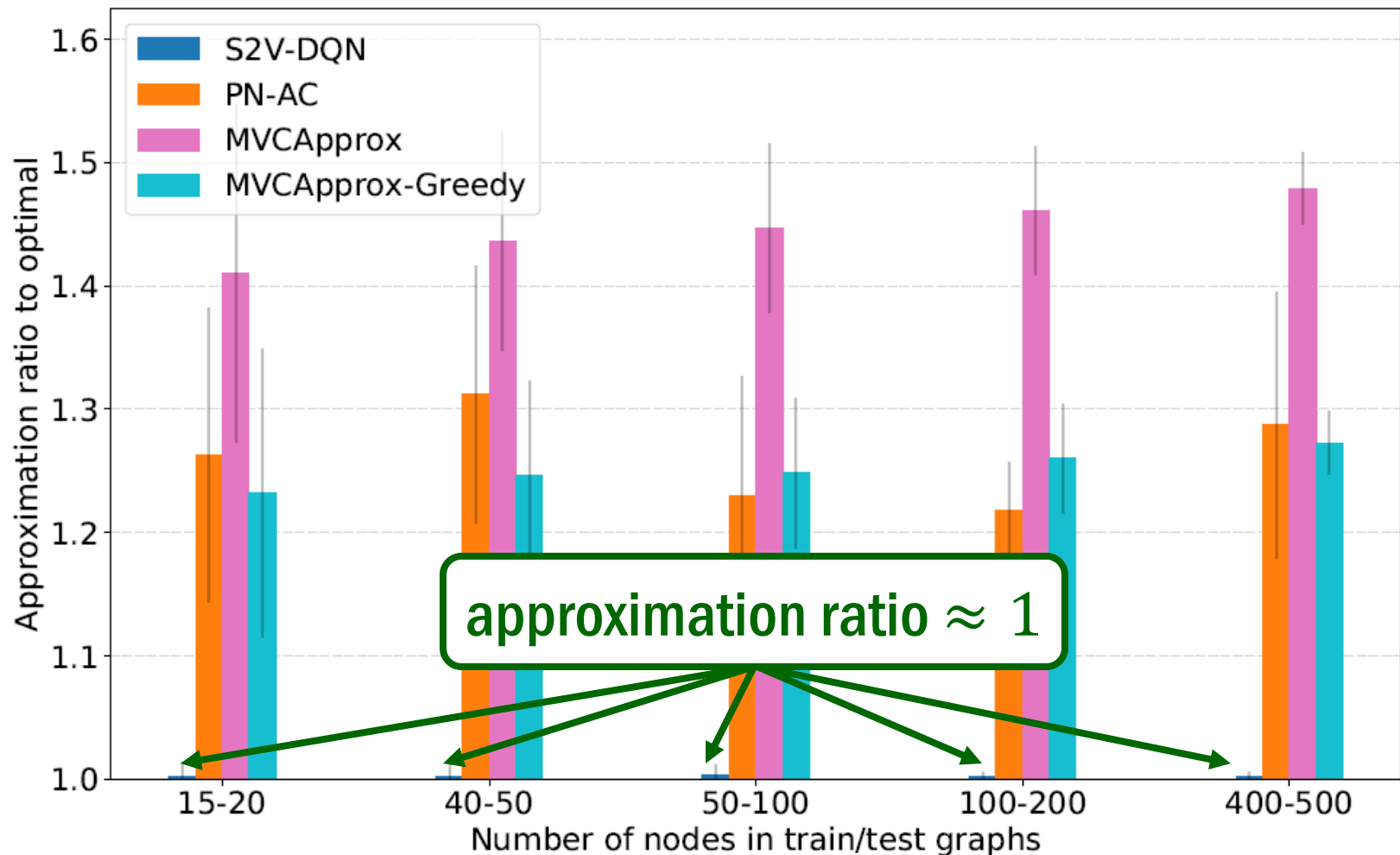# Graph embedding for state-action value function



**1st iteration**

Embed graph

**State**
$$S = \{x_i = 0\}_{i \in \mathcal{V}}$$

**State-action value function**
$$Q(S, i)$$
$$= \theta_1 \sigma(\theta_2\, \mu^a + \theta_3\, \mu_i)$$

Add best node

**Greedy action**
$$i^* = argmax\, Q(S, i)$$

aggregated embedding

individual embedding

**2nd iteration**

Embed graph

Add best node

$$x_{i^*} = 1, other\, x_j = 0$$

# Embedding leads to better heuristic algorithm

Minimum vertex cover: smallest number of nodes to cover all edges
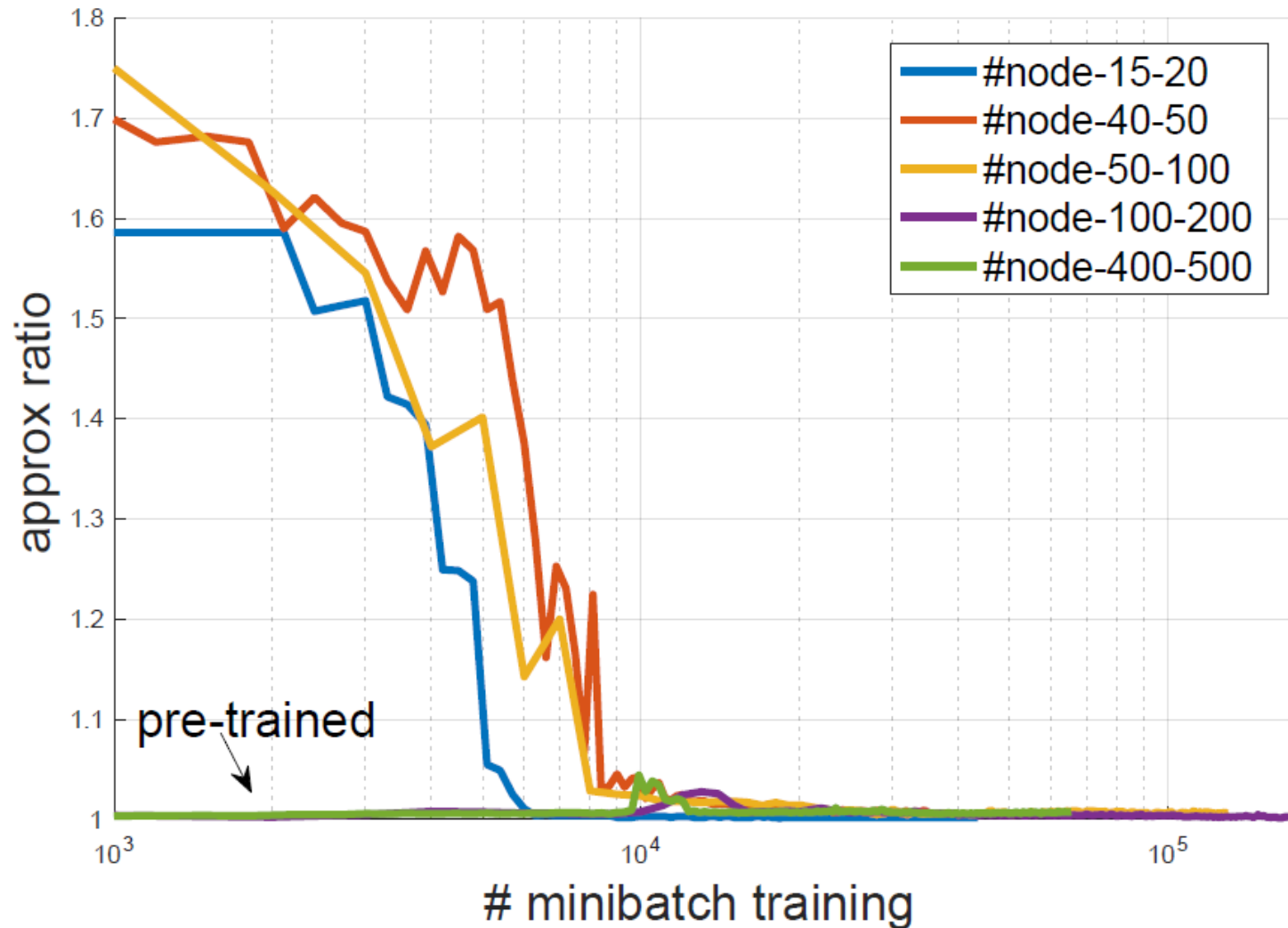A distribution of scale free networks
Optimal approximated by running CPLEX for 1 hour
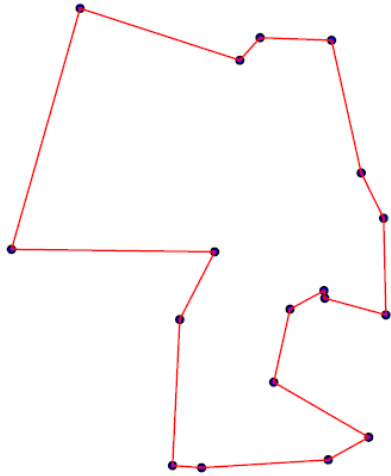


approximation ratio ≈ 1

# Training converge quite fast

Pre-training: initialize embedding parameters with ones trained with smaller networks
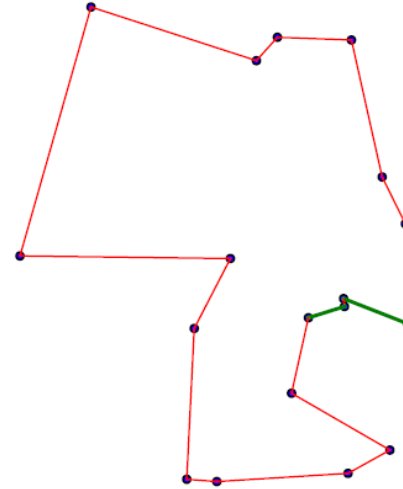
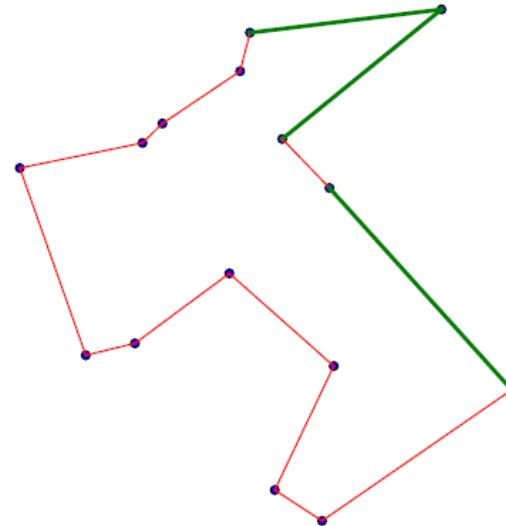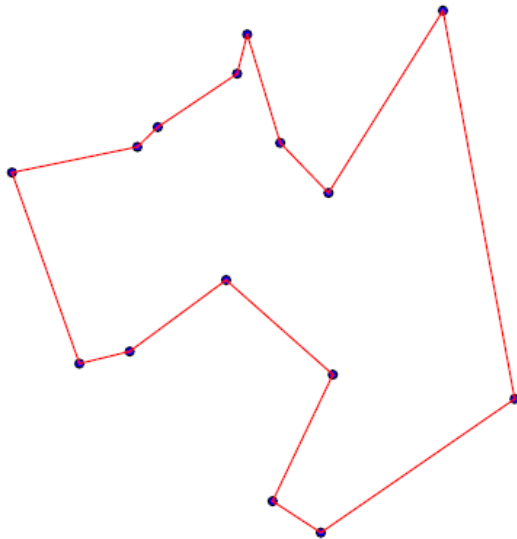# Also good for traveling salesman problem
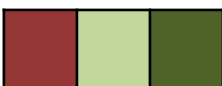
Optimal

Embedding

0.07% longer

0.5% longer

# Embedding as a tool for algorithm design

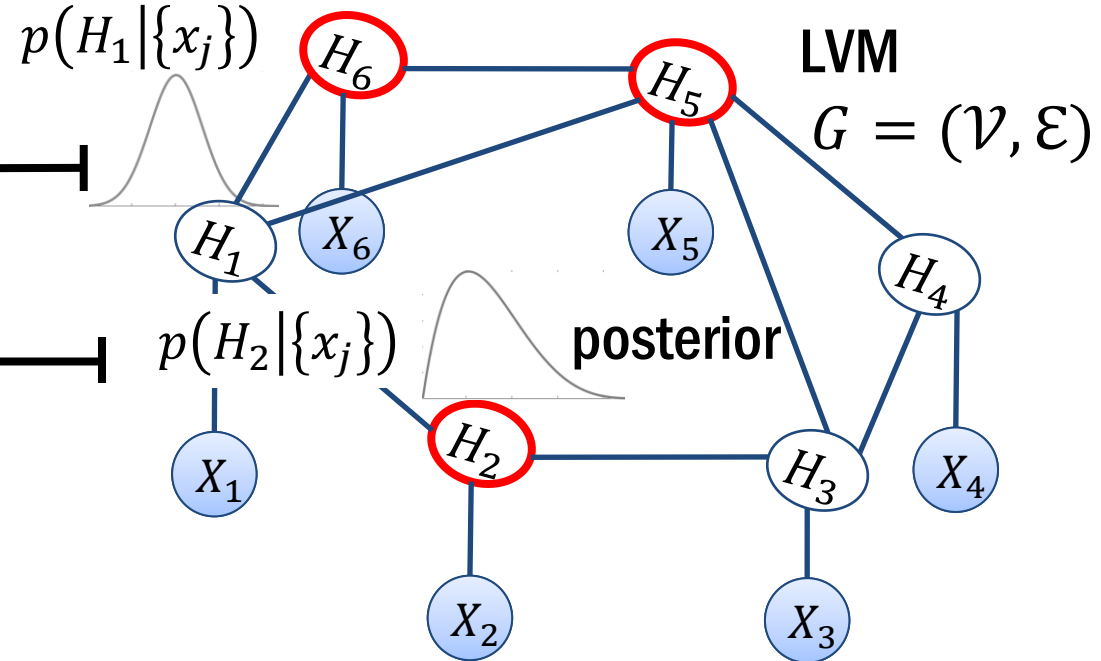**Embedding of node**

$$\mu_1(\chi, W)$$

$$+$$

$$\mu_2(\chi, W)$$

$$+$$

$$\vdots$$

$$= \mu^a(\chi, W)$$

**Embedding of entire structure**

$p(H_1|\{x_j\})$

$H_6$    $H_5$    LVM

$G = (\mathcal{V}, \mathcal{E})$

$H_1$   $X_6$    $X_5$    $H_4$

$p(H_2|\{x_j\})$    **posterior**

$X_1$

$H_2$    $H_3$    $X_4$

$X_2$    $X_3$

- **Embedding structures**

- Learn better? Nonconvex & RL?

- New system & programming language?

[Dai, Dai & Song 2016]