

Embedding as a Tool for Algorithm Design

Le Song

Center for Machine Learning

College of Computing

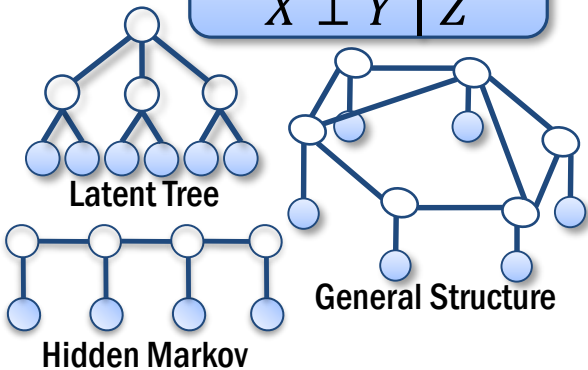
Georgia Institute of Technology

More Structure

Less

Graphical Models

$$X \perp Y | Z$$



Squeeze more info. out of big data

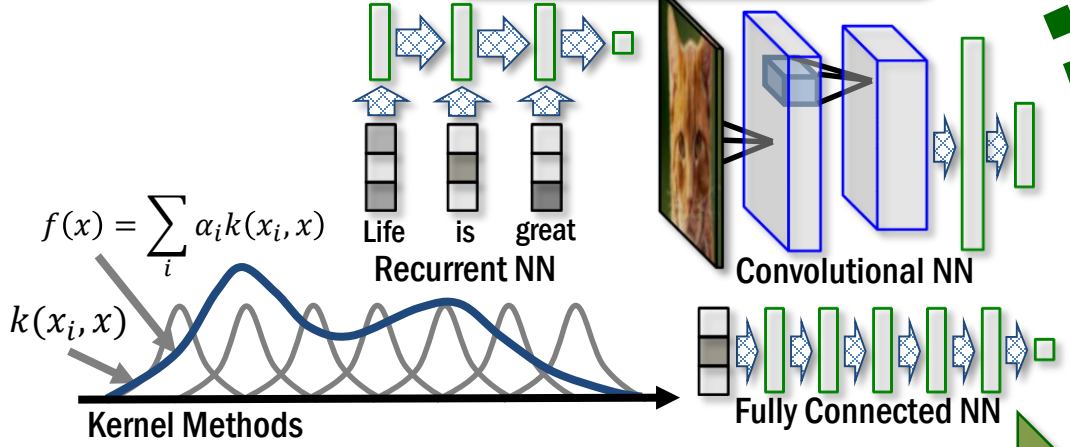
Something here?

Pareto Frontier

Machine Learning Algorithms = Model + Algorithm

Function Approximation

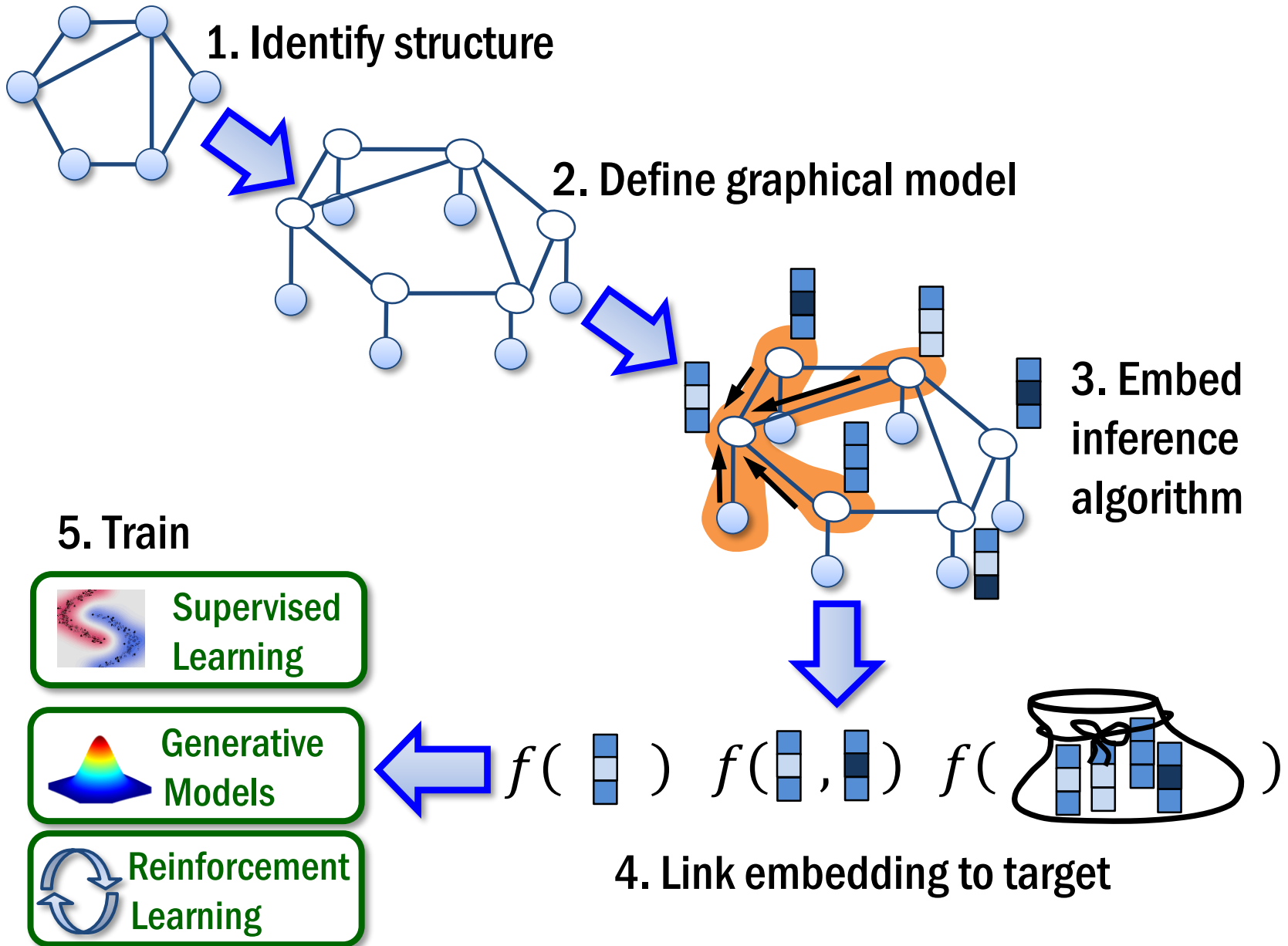
$$y = f(x)$$



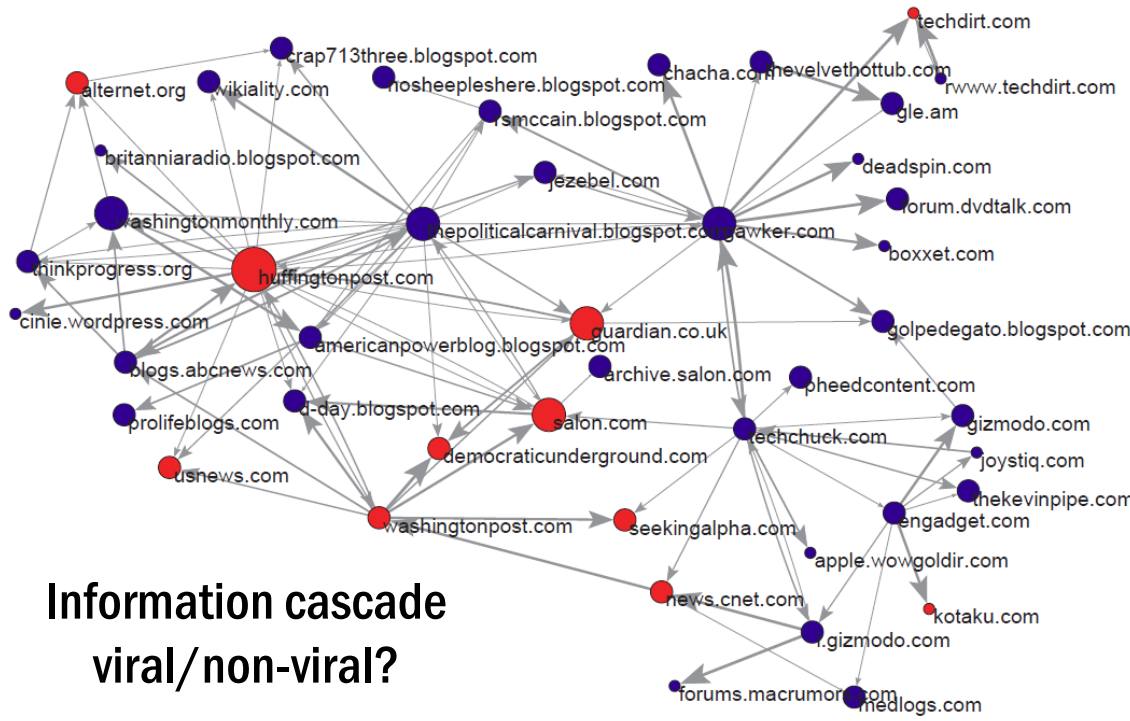
Less

More Scalable

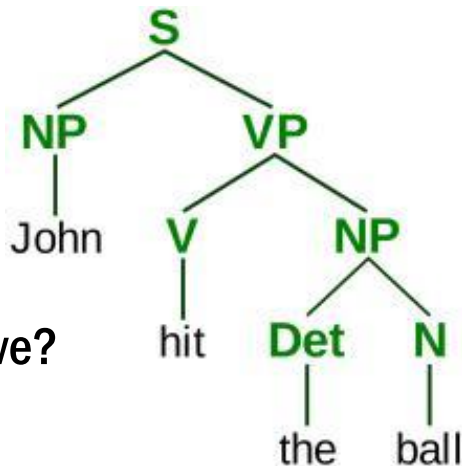
Embedding algorithms



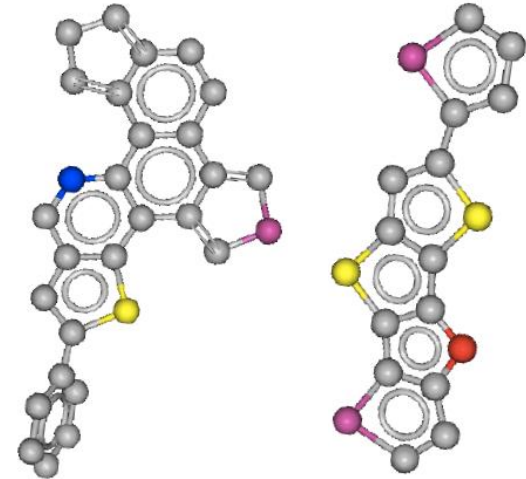
Motivation 1: Prediction for structured data



Natural language
positive/negative?



Drug/materials
effective/ineffective?



```

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
  
```

code graphs
benign/
malicious?

```

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]
  
```

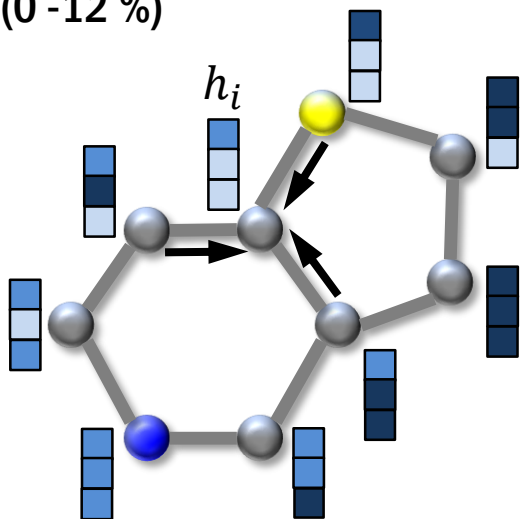
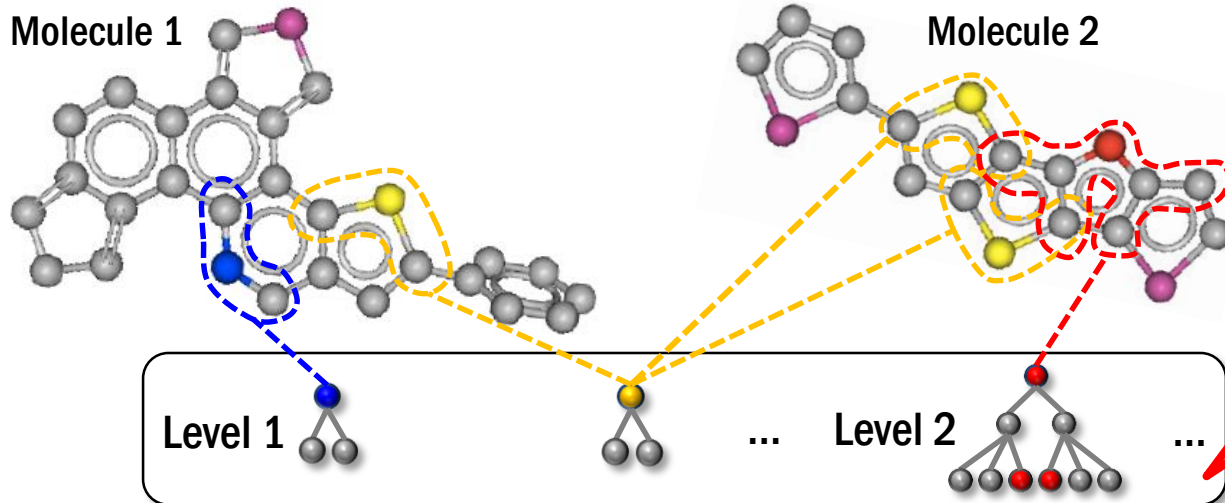
Big dataset, explosive feature space

2.3 M organic materials

“Bag of structures” representation

Predict

Efficiency (PCE)
(0 - 12 %)



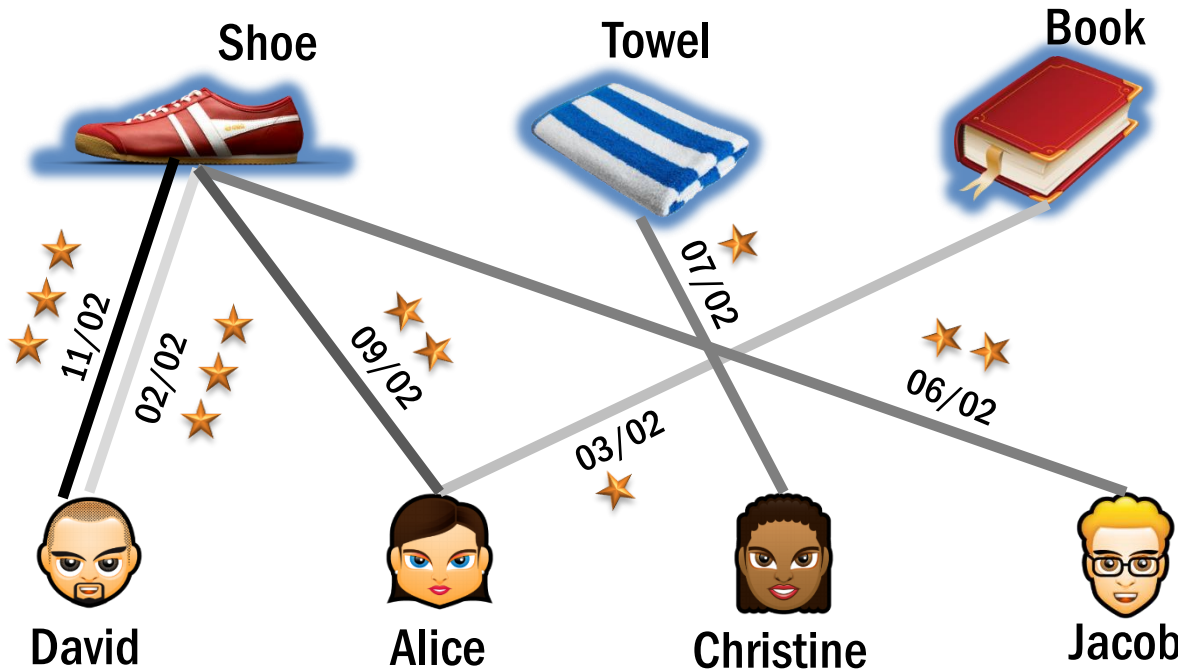
Weisfeiler-Lehman algorithm

1. $h_i \leftarrow \text{Hash}(\text{node type}), \forall i$
2. Iterate T times:
 $h_i \leftarrow \text{Hash}(h_i + \sum_{j \in \mathcal{N}(i)} h_j), \forall i$
3. Aggregate $\sum_{\forall i} h_i$

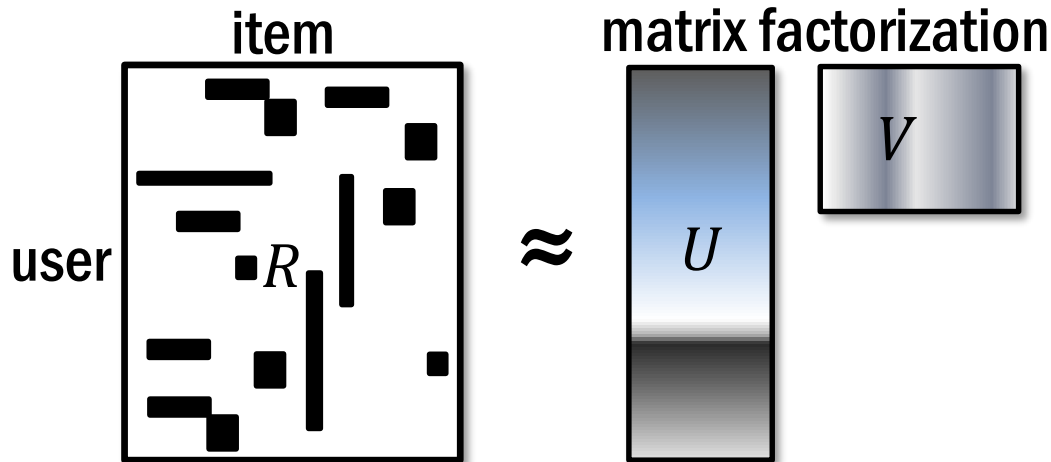
Hash manually designed, need 100 million param.

Embedding reduces model size by 1000 times !

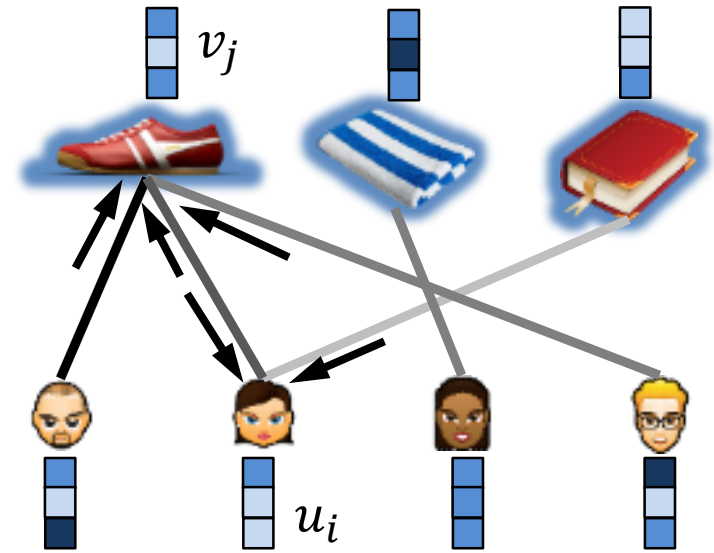
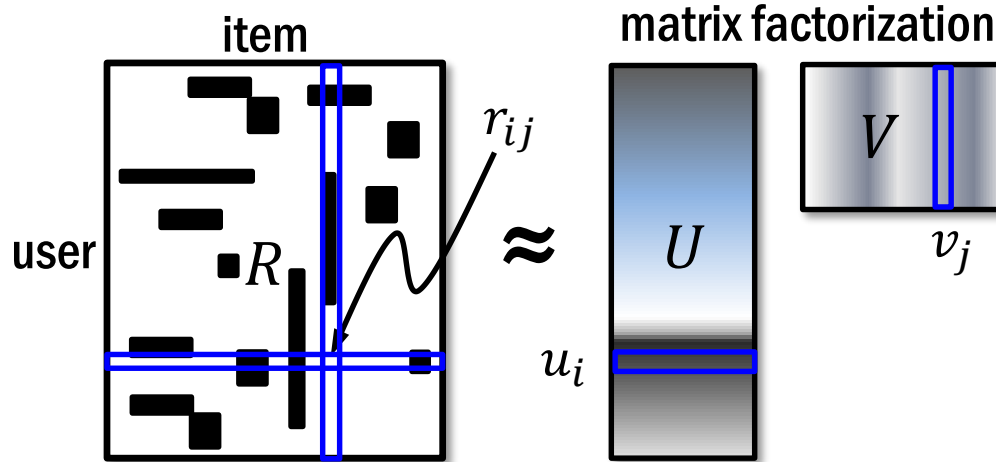
Motivation 2: Dynamic processes over networks



who will do
what and when?



Complex behavior not well captured



Alternating least square

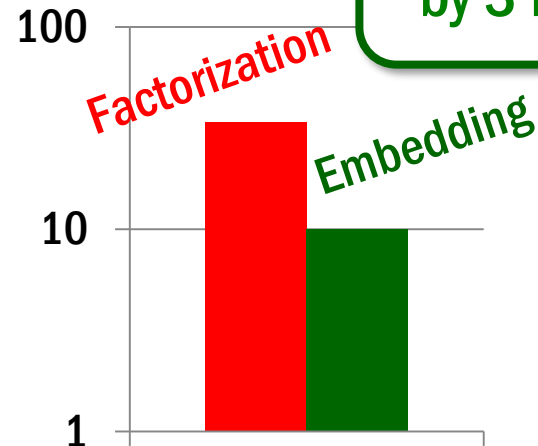
1. Initialize $u_i, v_j, \forall i, j$
2. Iterate T times

$$u_i \leftarrow \operatorname{argmin}_u \sum_{j \in \mathcal{N}(i)} (r_{ij} - u \cdot v_j)^2, \forall i$$

$$v_j \leftarrow \operatorname{argmin}_v \sum_{i \in \mathcal{N}(j)} (r_{ij} - u_i \cdot v)^2, \forall j$$

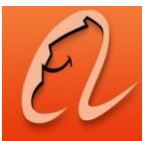
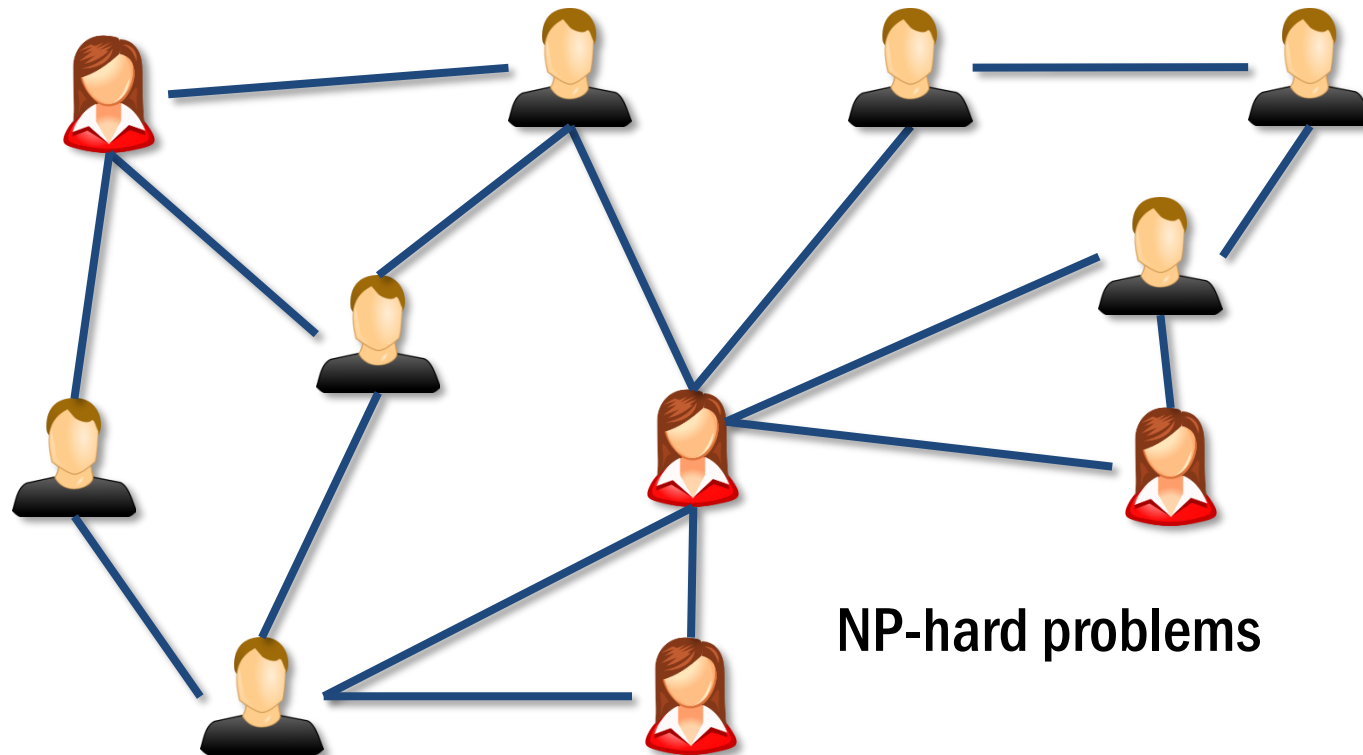
temporal / sequential
information not modeled

Return Time
MAE (hour)



Motivation 3: Combinatorial optimizations over graphs

Application	Optimization problem
Advertisers: influence maximization Analysts: community discovery Platforms: resource scheduling	Minimum vertex/set cover Maximum cut Traveling salesman



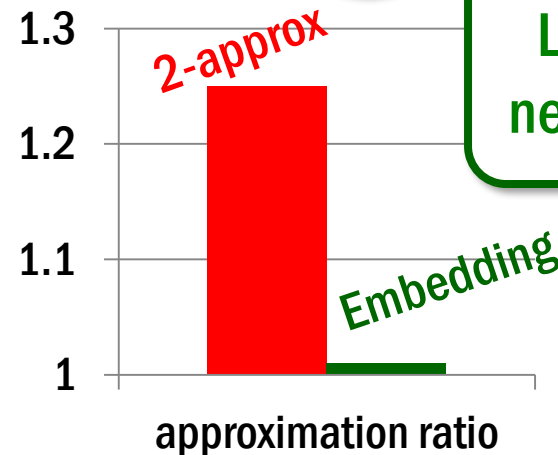
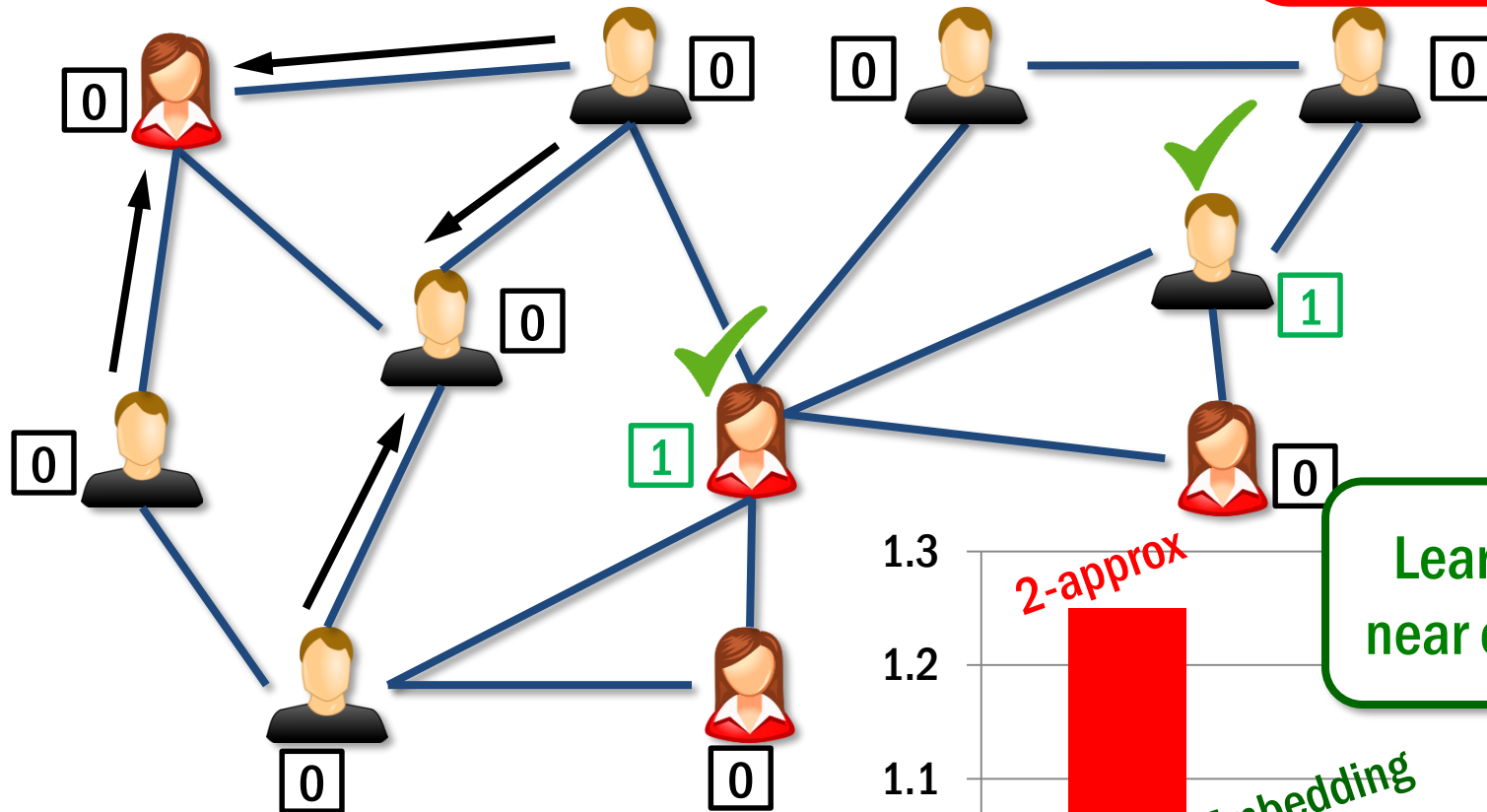
Simple heuristics do not exploit data

2 - approximation for minimum vertex cover

Repeat till all edges covered:

- Select uncovered edge with **largest total degree**

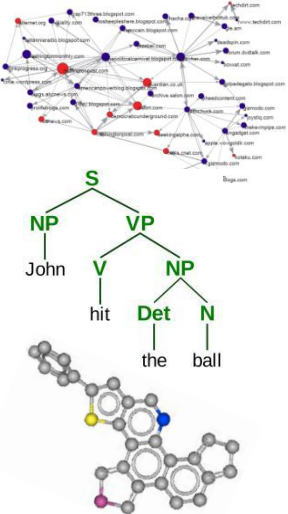
Manually designed rule.
Can we learn from data?



Learn to be near optimal!

Key observation & fundamental question

Algorithm = Structured composition of **manually** designed operation

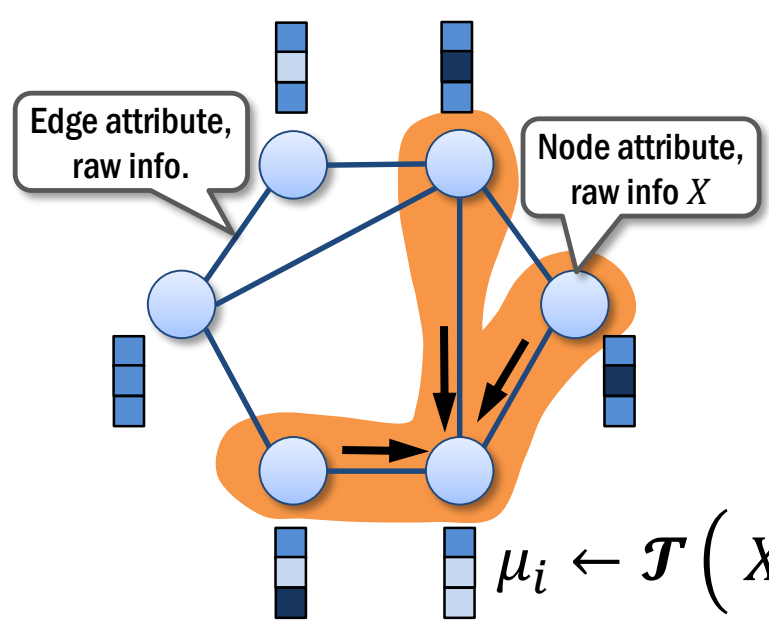
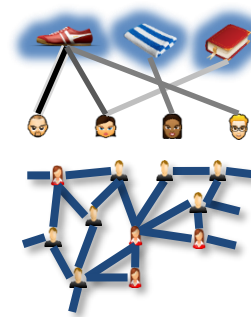


```

mov  [reg=Chovar_05], edi
mov  [reg=Chovar_10], edi
mov  [reg=Chovar_2C], edi
mov  edi, [edi]
mov  [reg=Chovar_4], 0
mov  [reg=Chovar_1], edi
call  ecx

; uc_MOC1818
cmp  byte_1, uc_MOC1818
jz   short uc_MOC1818

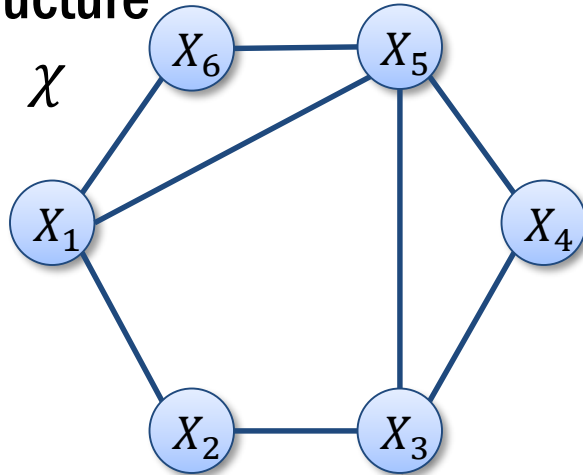
mov  esi, edi
mov  [reg=Chovar_10], esi
mov  [reg=Chovar_24], esi
call  CRYPTD_md5cx
;
mov  [reg=Chovar_1], esi
mov  [reg=Chovar_20], esi
call  esi, [reg=Chovar_20]
    
```



Design in a unified framework?
Learn these algorithms?

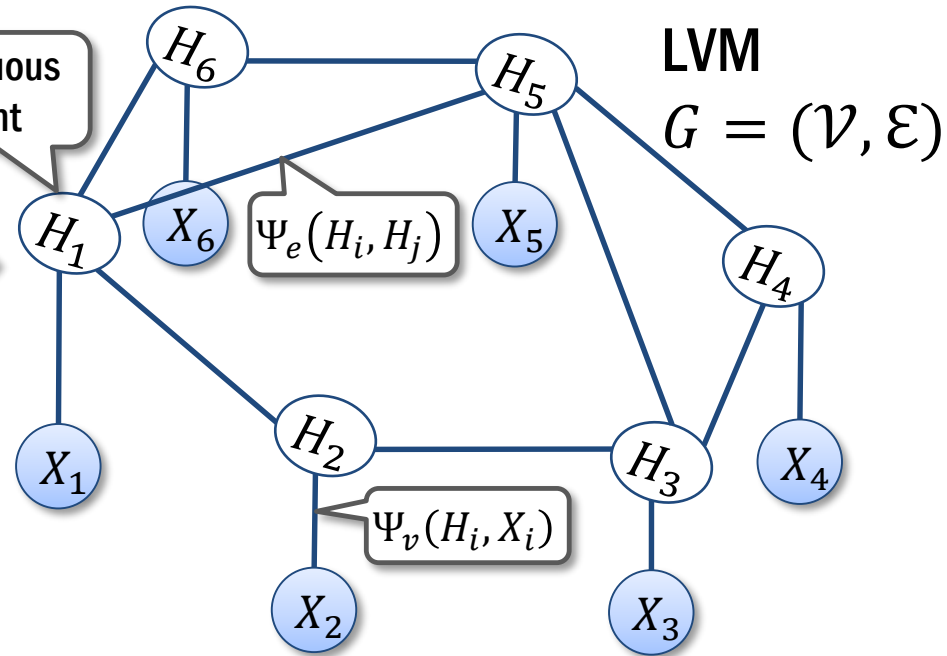
Represent structure as latent variable model (LVM)

Structure



Represent

Continuous Latent



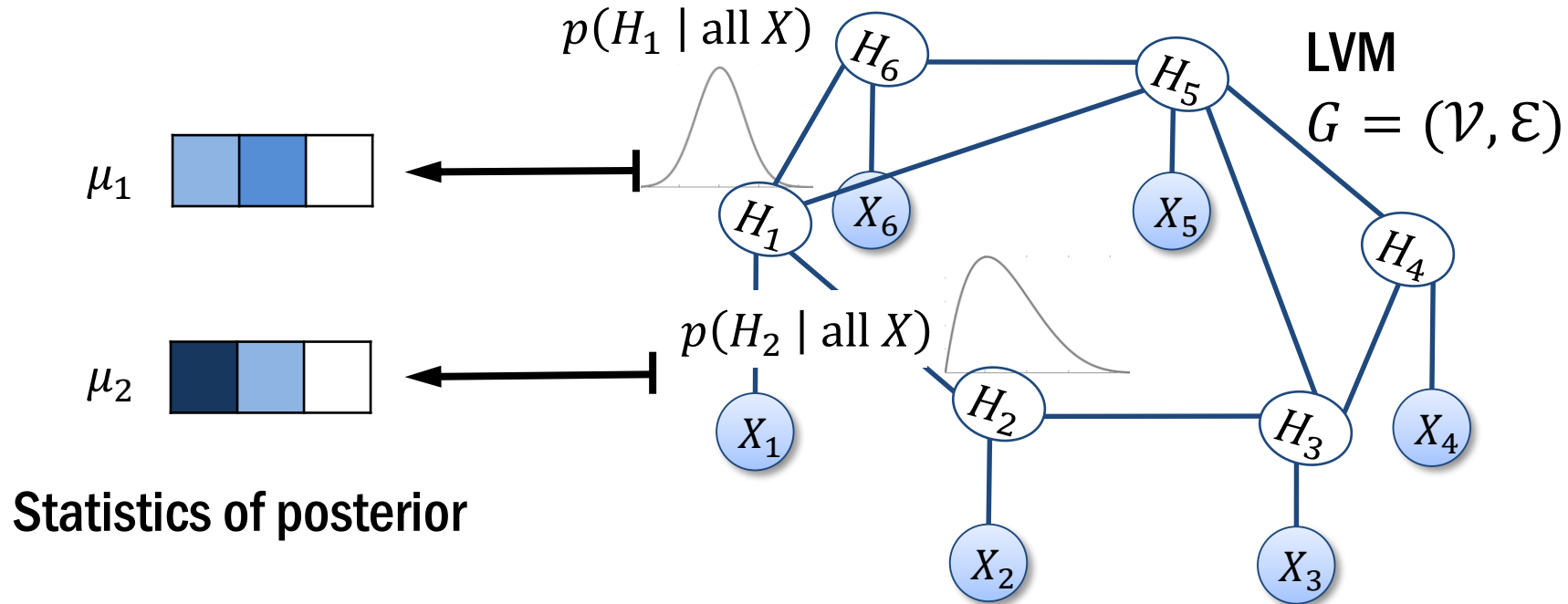
Joint likelihood of hidden variables

$$p(\text{all } H \mid \text{all } X) \propto \prod_{i \in \mathcal{V}} \underbrace{\Psi_v(H_i, X_i \mid \theta_v)}_{\text{Nonnegative node potential}} \prod_{(i,j) \in \mathcal{E}} \underbrace{\Psi_e(H_i, H_j \mid \theta_e)}_{\text{Nonnegative edge potential}}$$

Nonnegative
node potential

Nonnegative
edge potential

Posterior distribution as features



Statistics of posterior

$$p(H_i | \text{all } X) = \sum_{\text{all } H_j \text{ except } H_i} p(\text{all } H | \text{all } X)$$

Capture both nodal and topological info.

Aggregate information from distant nodes

Mean field algorithm aggregates information

Approximate posterior

$$p(H_i | \text{all } X) \approx q_i(H_i)$$

via fixed point iteration:

1. Initialize $q_i(H_i), \forall i$

2. Iterate T times

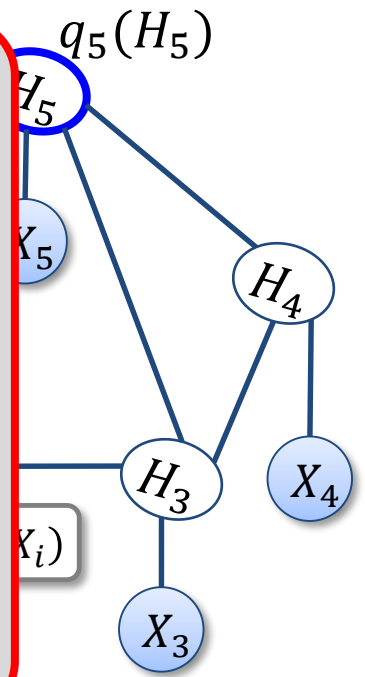
$$q_i(H_i) \leftarrow \Psi_v(H_i, X_i)$$

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

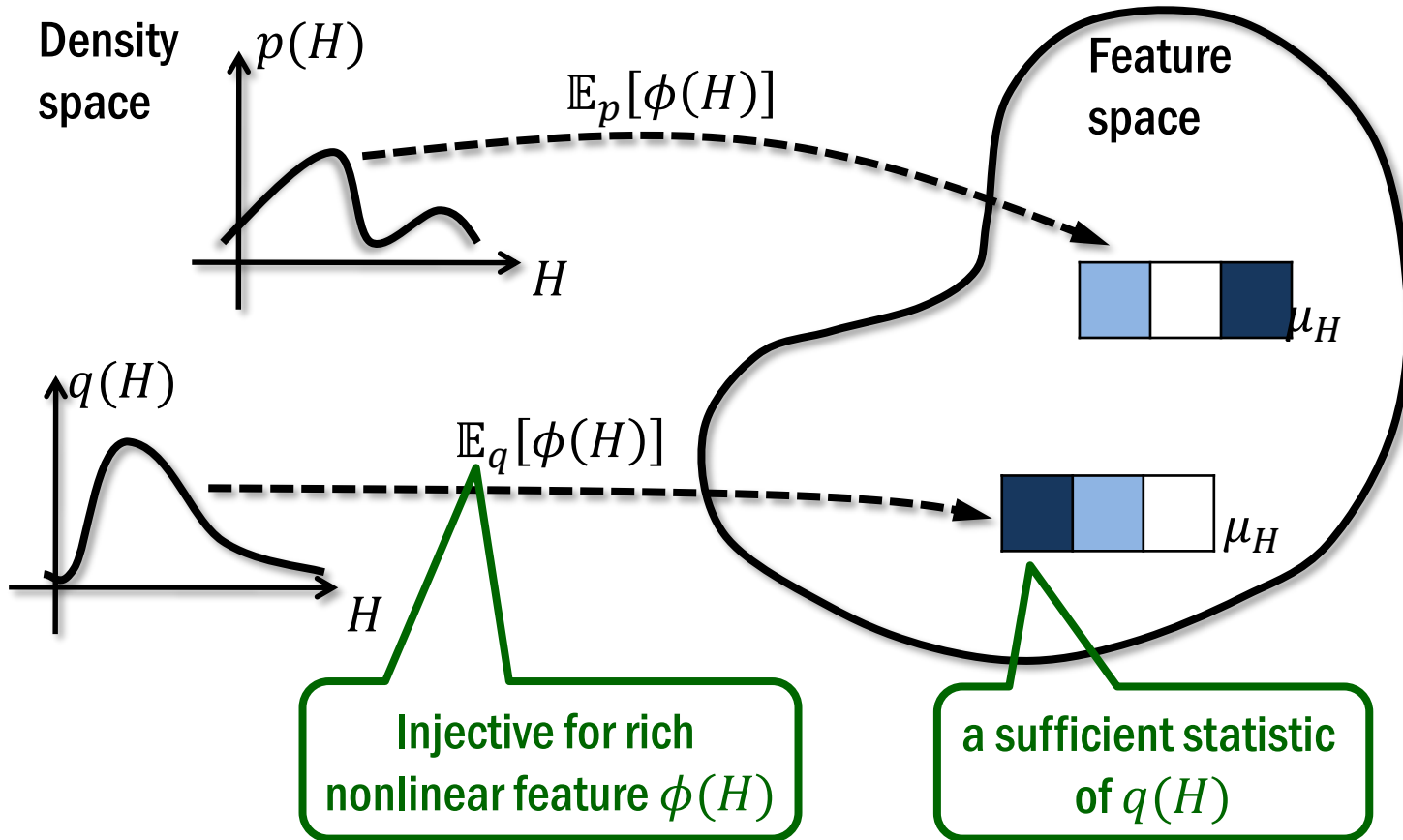
$$\mathcal{T} \left(X_i, \{q_j(H_j)\}_{j \in \mathcal{N}(i)} \right)$$

Need to perform integration

Need to learn Ψ_v and Ψ_e



What's embedding?



Example:

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



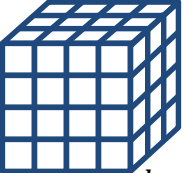
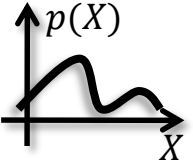

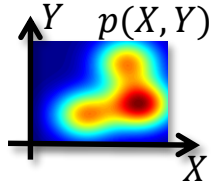

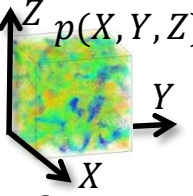
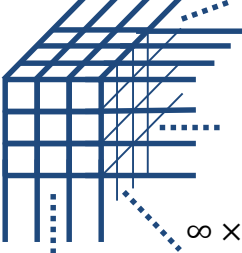
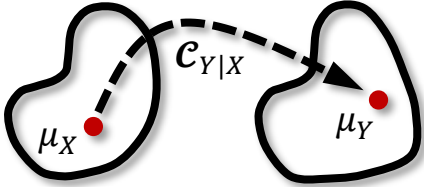
$$\mu_H =$$

$$\left[\begin{array}{c} \text{Mean,} \\ \hline \text{Variance,} \\ \hline \text{higher} \\ \text{order} \\ \text{moment} \\ \vdots \end{array} \right]$$

Equivalent Operation

$$\mathcal{J}(q(H)) = \tilde{\mathcal{J}}(\mu_H)$$

Learning via embedding

	Distributions			Probabilistic Operations
Discrete	$\mathbb{P}(X)$  $d_X \times 1$	$\mathbb{P}(X, Y)$  $d_X \times d_Y$	$\mathbb{P}(X, Y, Z)$  $d_X \times d_Y \times d_Z$	Sum Rule: $\mathbb{P}(Y) = \sum_X \mathbb{P}(Y X)\mathbb{P}(X)$ Product Rule: $\mathbb{P}(Y, X) = \mathbb{P}(Y X)\mathbb{P}(X)$ Bayes Rule: $\mathbb{P}(X y) = \frac{\mathbb{P}(y X)\mathbb{P}(X)}{\mathbb{P}(y)}$
Embedding	 $\mu_X := \mathbb{E}_X[\phi(X)]$  $\infty \times 1$	 $\mathcal{C}_{XY} := \mathbb{E}_{XY}[\phi(X) \otimes \phi(Y)]$  $\infty \times \infty$	 $\mathcal{C}_{XYZ} := \mathbb{E}_{XYZ}[\phi(X) \otimes \phi(Y) \otimes \phi(Z)]$  $\infty \times \infty \times \infty$	 Sum Rule: $\mu_Y = \mathcal{C}_{Y X}\mu_X$ Product Rule: $\mathcal{C}_{YX} = \mathcal{C}_{Y X}\mathcal{C}_{XX}$ Bayes Rule: $\mu_{X y} = \mathcal{C}_{X Y}\phi(y)$

Divergence & Independence measure

- Feature selection
- Clustering
- Reduction
- Transfer

Embedding graphical models

- Spectral HMM
- Kernel belief propagation
- Latent tree & junction tree

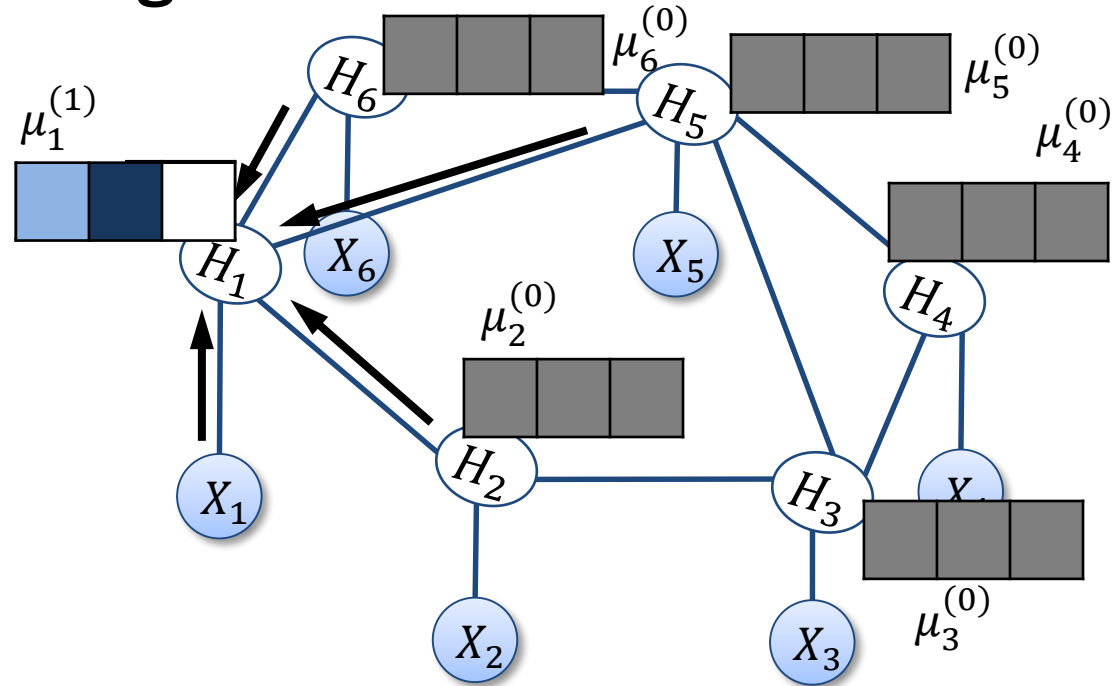
Embedding mean field

Approximate embedding of

$$p(H_i | \text{all } X) \mapsto \mu_i$$

via fixed point update

1. Initialize $\mu_i, \forall i$
2. Iterate T times



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

Embedding mean field

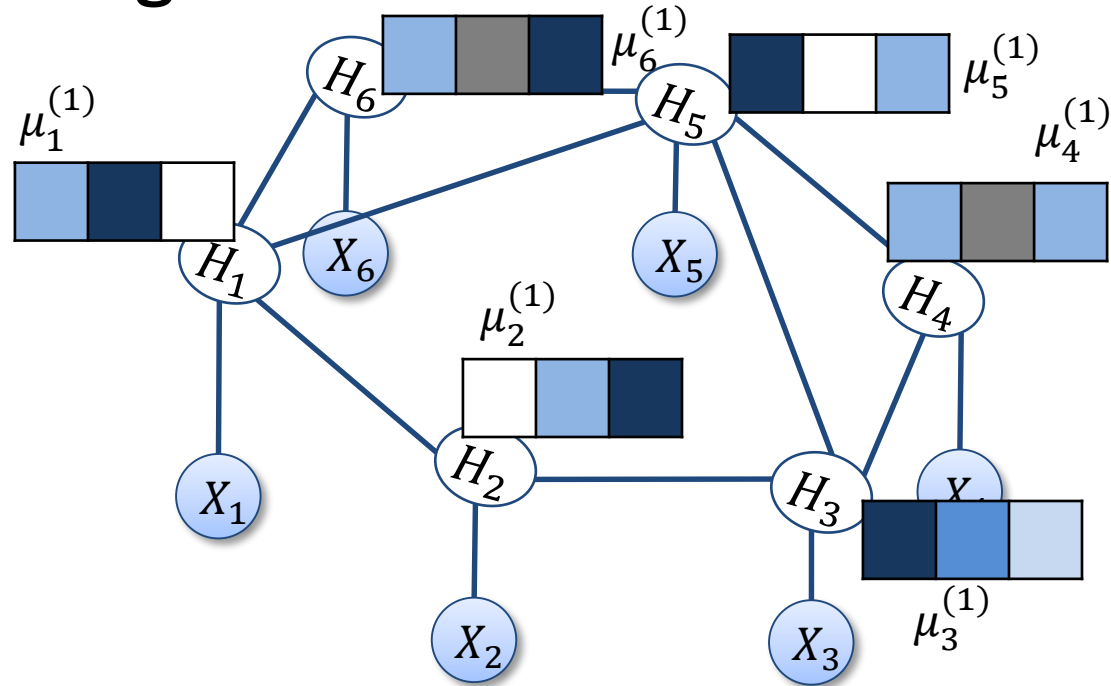
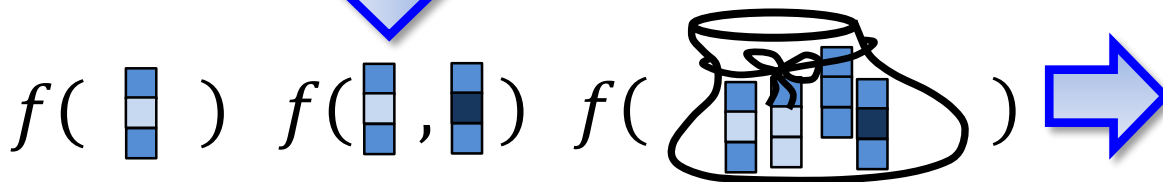
Approximate embedding of

$$p(H_i | \text{all } X) \mapsto \mu_i$$

via fixed point update

1. Initialize $\mu_i, \forall i$
2. Iterate T times

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



Supervised Learning

Generative Models

Reinforcement Learning

Directly parameterize nonlinear mapping

$$\mu_i \leftarrow \tilde{\mathcal{F}} \left(X_i, \{\mu_j\}_{j \in \mathcal{N}(i)} \right)$$

Use any universal function approximator, eg. kernel function

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

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

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

$d \times n$ matrix $d \times d$ matrix

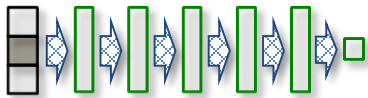
Will be learned

Embedded algorithm is flexible yet structured

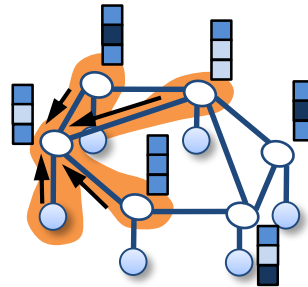
Embedded algorithm = Structured composition of nonlinear functions

Model Space

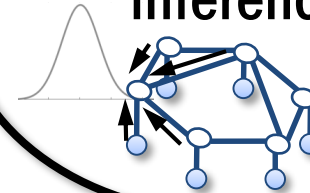
Generic function approximator



Embedded inference algorithm



Graphical model inference



Benefit of the new view: belief propagation

Approximate posterior

$$p(H_i | \{x_j\}) = \Psi_v(H_i, X_i) \Psi_e(H_i, H_j)$$

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

via fixed point iteration

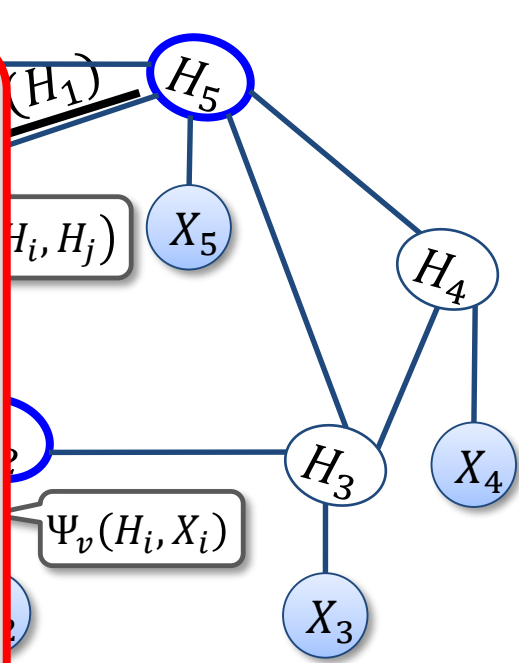
1. Initialize m_{ij}
2. Iterate T times

$$m_{ij}(H_j) \leftarrow \int_{\mathcal{H}} \underbrace{\Psi_v(H_i, X_i) \Psi_e(H_i, H_j)}_{\mathcal{T}(X_i, \{m_{\ell i}(H_i)\}_{\ell \in \mathcal{N}(i) \setminus j})} \cdot \prod_{\ell \in \mathcal{N}(i) \setminus j} m_{\ell i}(H_i) dH_i, \forall i, j$$

$$\mathcal{T}(X_i, \{m_{\ell i}(H_i)\}_{\ell \in \mathcal{N}(i) \setminus j})$$

Need to perform integration

Need to learn Ψ_v and Ψ_e



Embed belief propagation

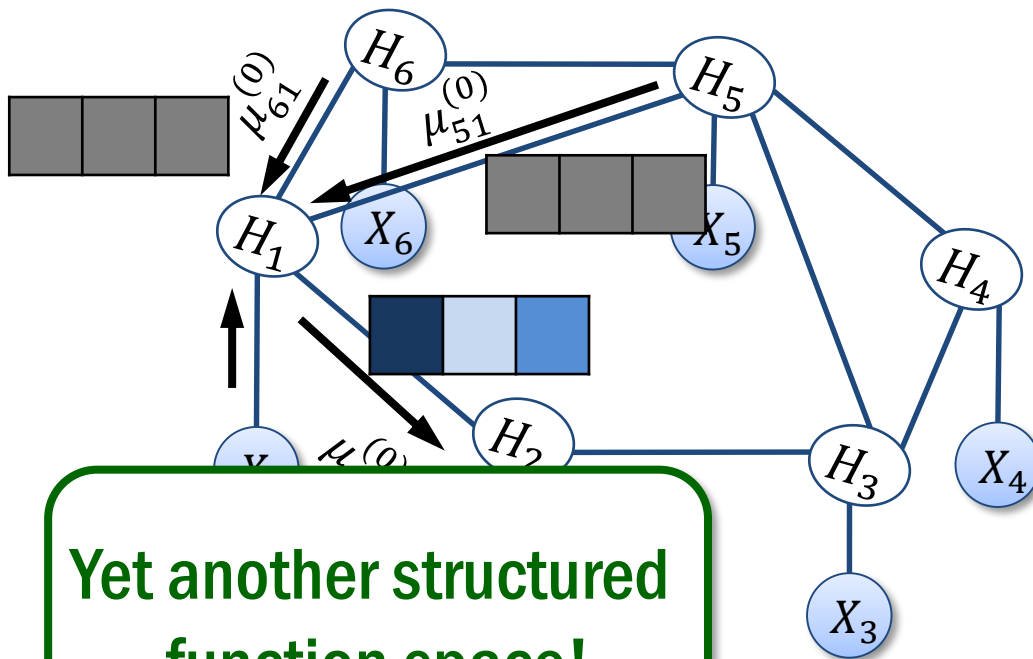
Approximate embedding of

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

via fixed point update

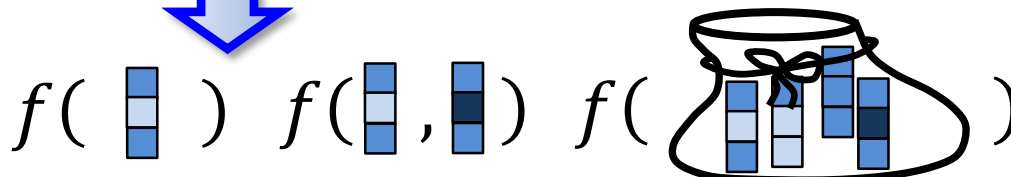
1. Initialize $\mu_{ij}, \forall (i, j)$
2. Iterate T times

Yet another structured function space!

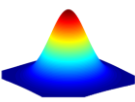


$$\mu_{ij} \leftarrow \tilde{\mathcal{F}}(X_i, \{\mu_{\ell i}\}_{\ell \in \mathcal{N}(i) \setminus j}), \forall (i, j)$$

3. Aggregate $\mu_i = \tilde{\mathcal{F}}(\{\mu_{\ell i}\}_{\ell \in \mathcal{N}(i)}), \forall i$



 Supervised Learning

 Generative Models

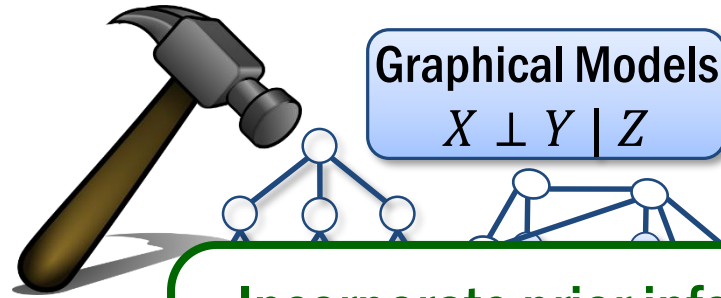
 Reinforcement Learning

New tools for algorithm design

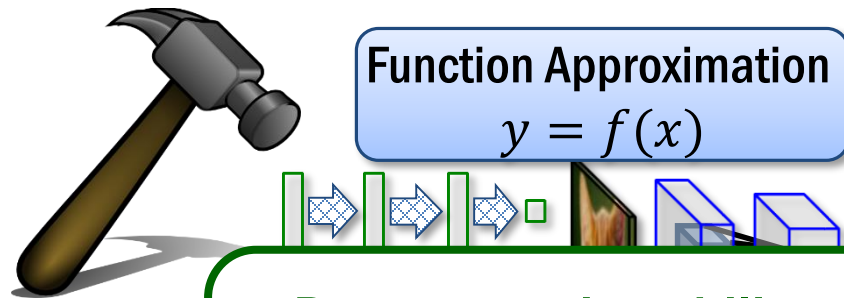
Manual algorithm design

(city=Atlanta) AND (age=40)
(browser=IE) XOR (system=Linux)
(bought=car) OR (usage<3 years)

Explosive combinations!

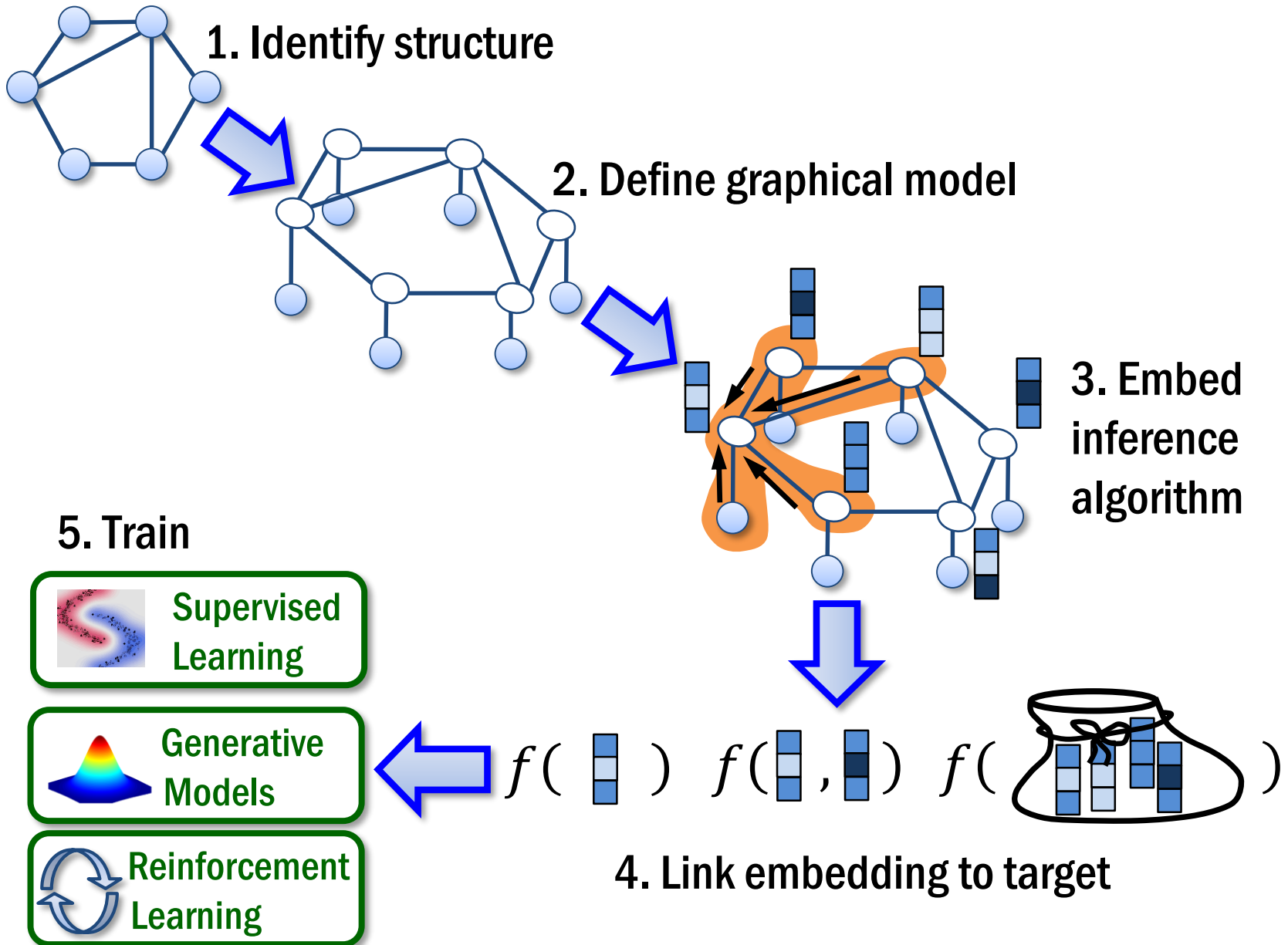


Incorporate prior info.
Reason about structure
Inference algorithm

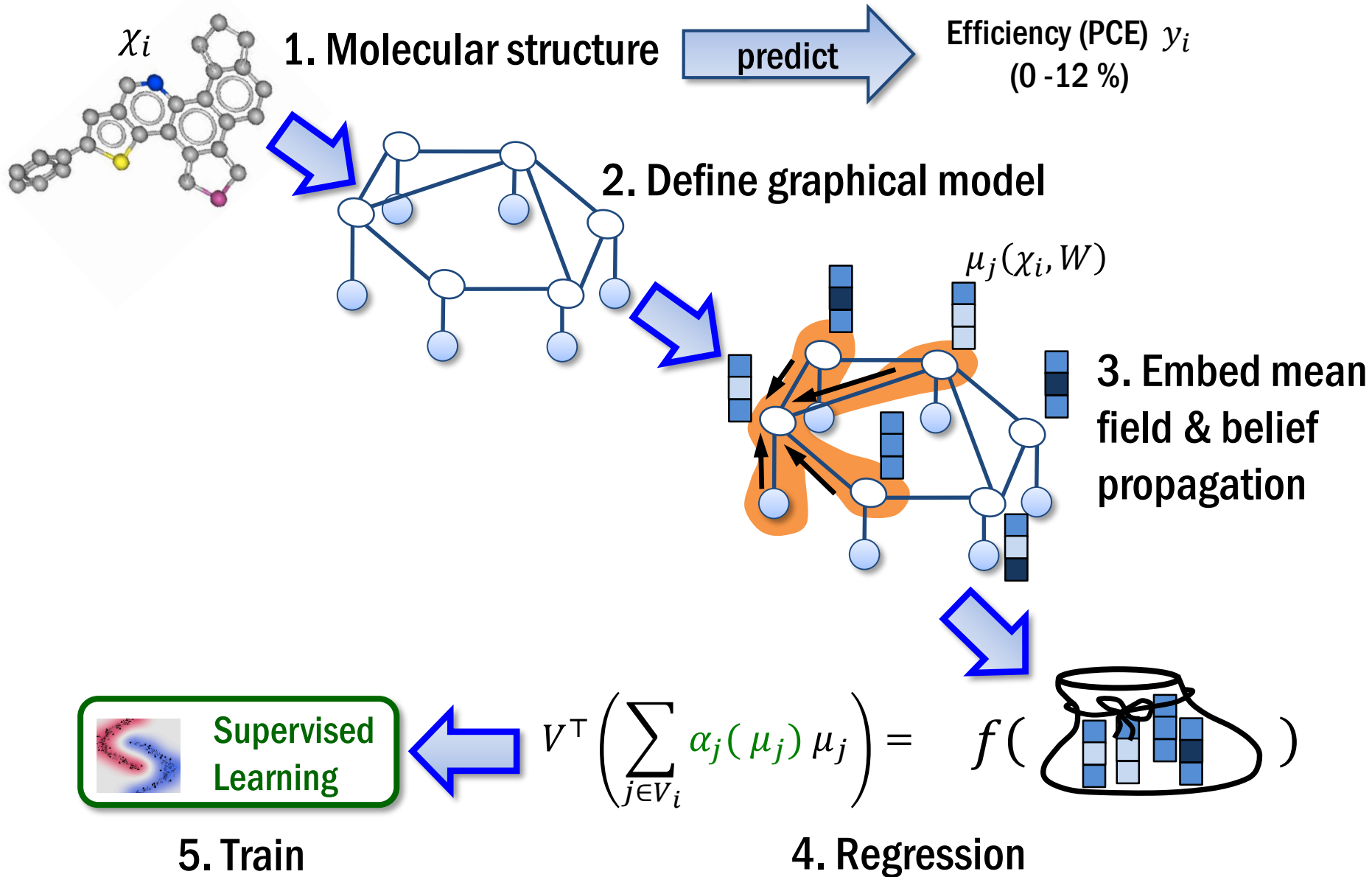


Representation ability
Statistical complexity
Generalization ability

Embedding algorithms

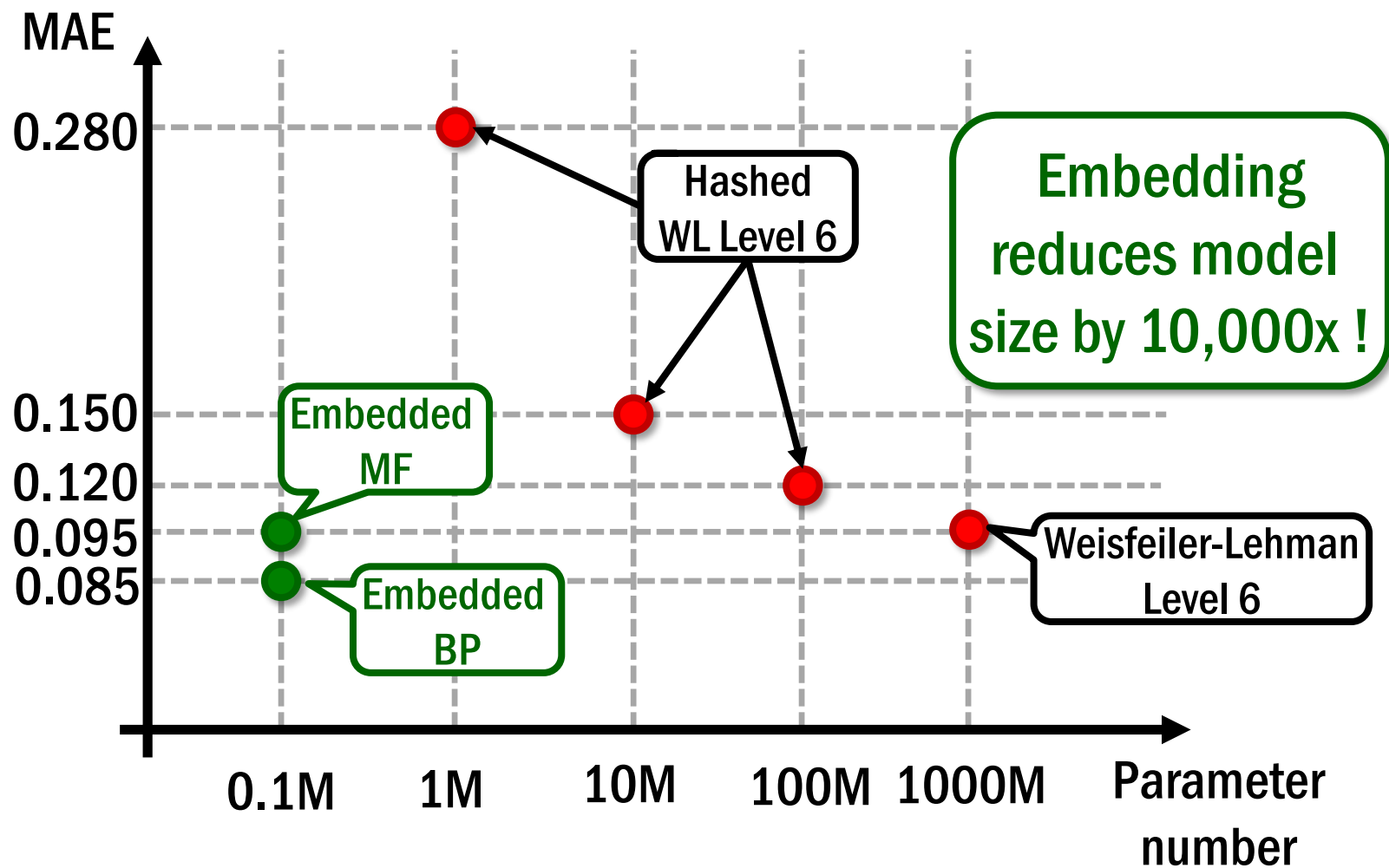


Scenario 1: Prediction for structured data

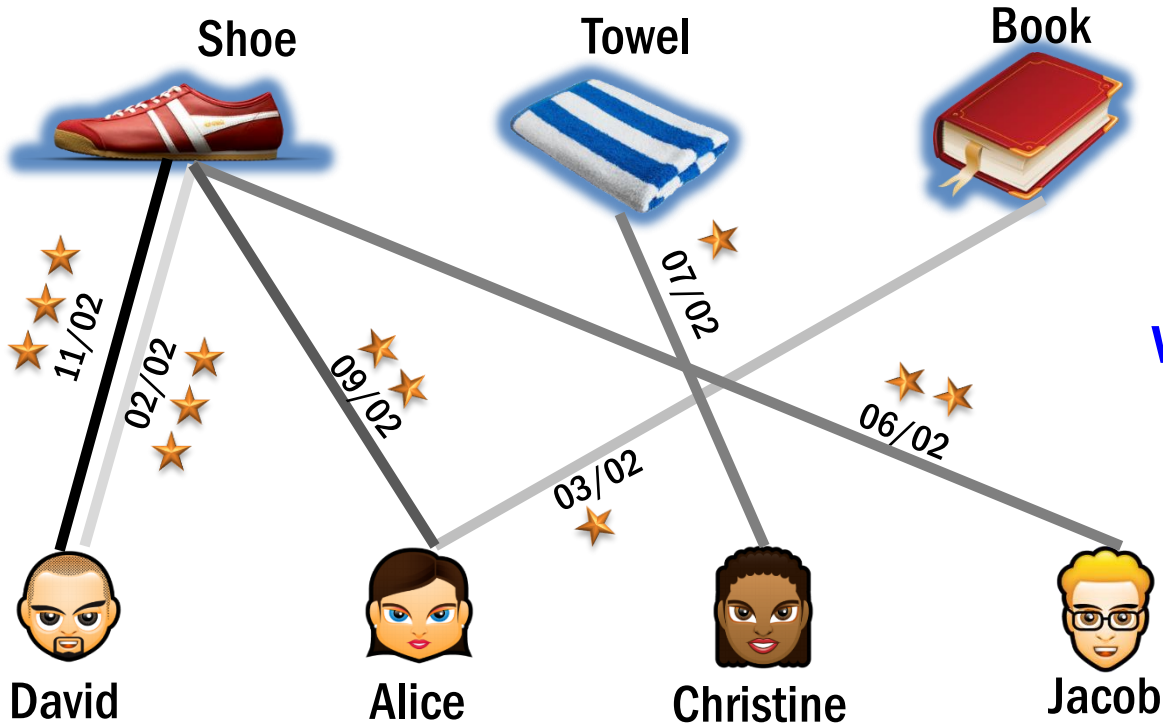


More compact model and lower error

Harvard clean energy dataset, 2.3 million organic molecules,
predict power conversion efficiency (0 -12 %)

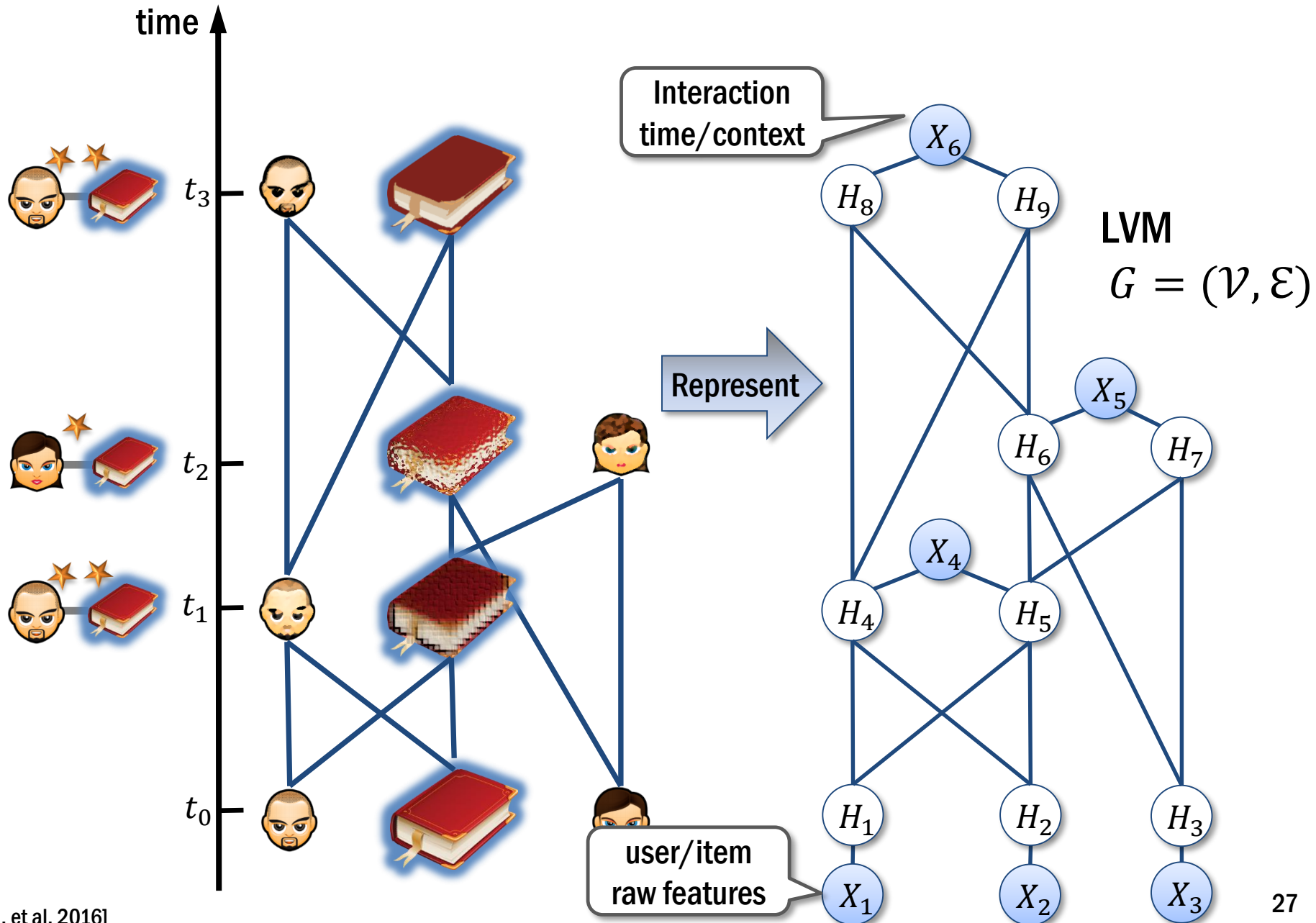


Motivation 2: Dynamic processes over networks

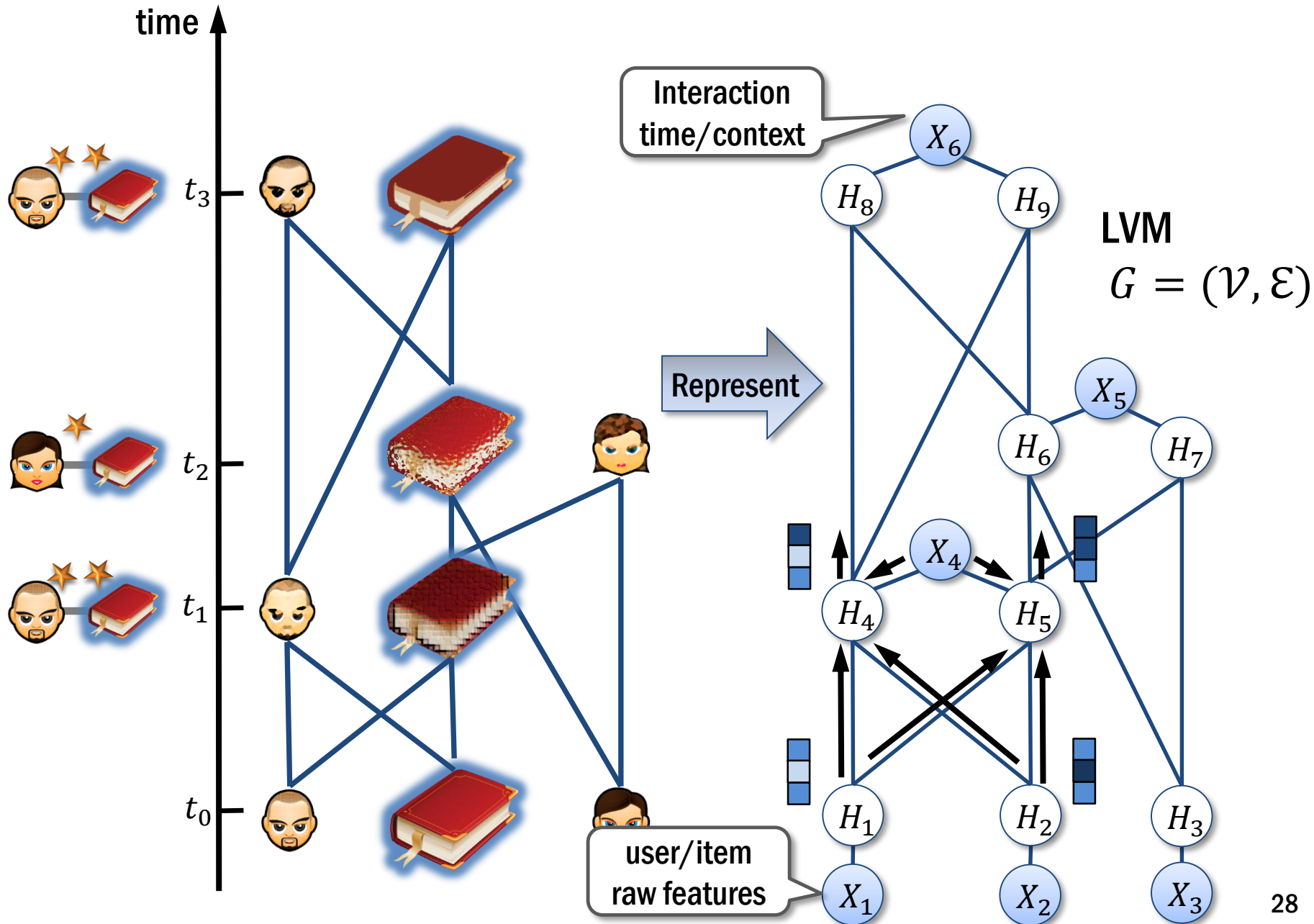


who will do
what and when?

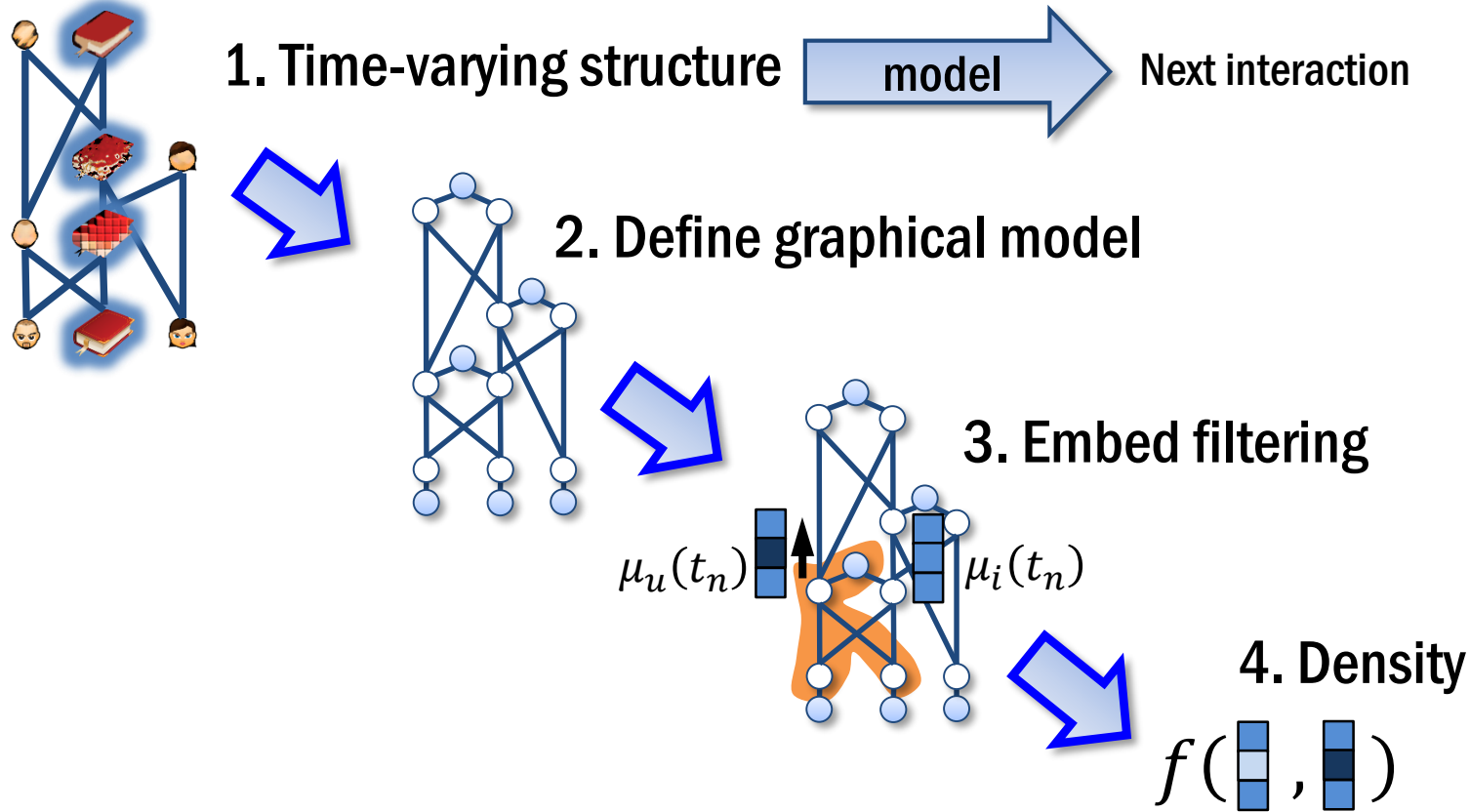
Unroll: time-varying dependency structure



Embed filtering/forward message passing



Embedding algorithm for building generative model

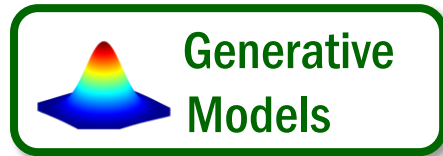


Compatibility between user u and item i

$$\alpha_{ui} = \exp(\mu_u^\top(t_n)\mu_i(t_n))$$

Likelihood of next event time $p_{ui}(t)$

$$\alpha_{ui}(t - t_n) \exp\left(-\frac{\alpha_{ui}(t - t_n)^2}{2}\right)$$



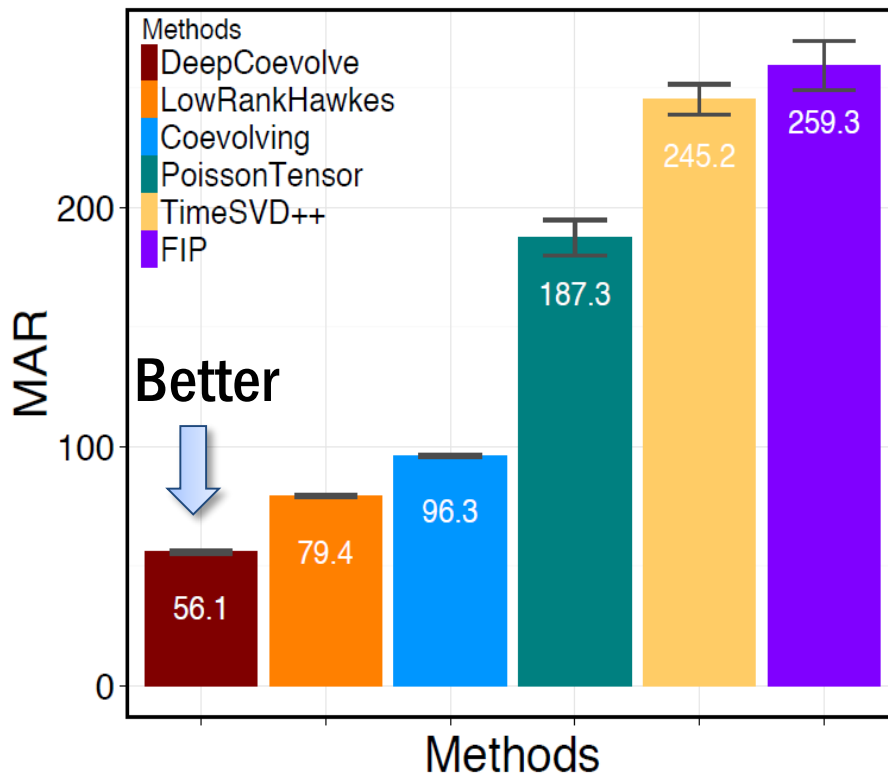
5. Train with
MLE or GAN

IPTV dataset

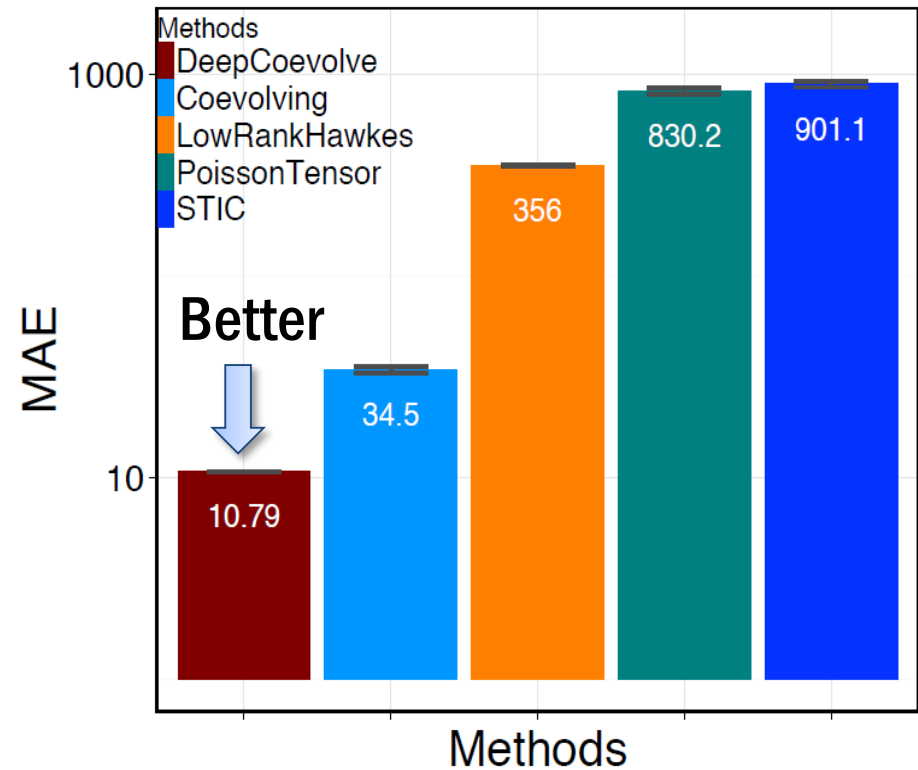
7,100 users, 436 programs, ~2M views

MAR: mean absolute rank difference

MAE: mean absolute error (hours)

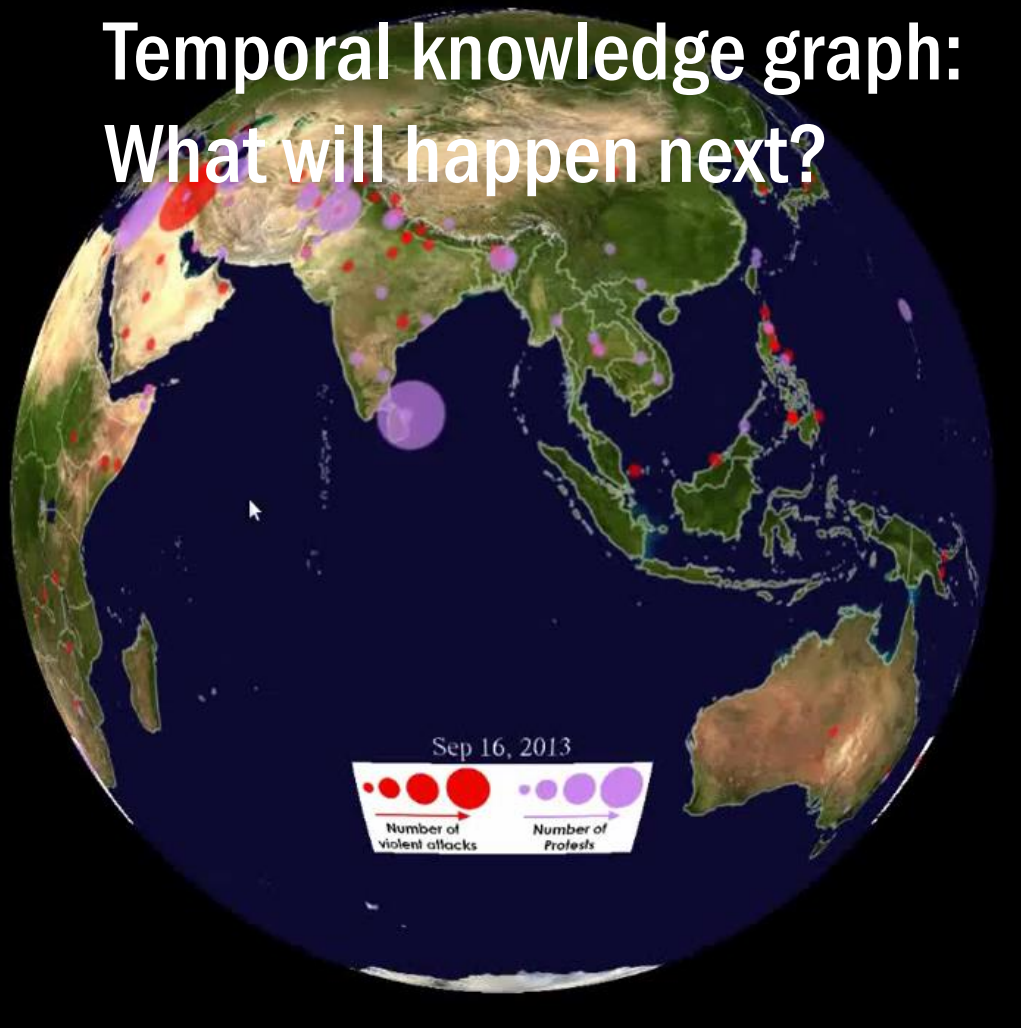


Next item prediction



Return time prediction

Temporal knowledge graph: What will happen next?

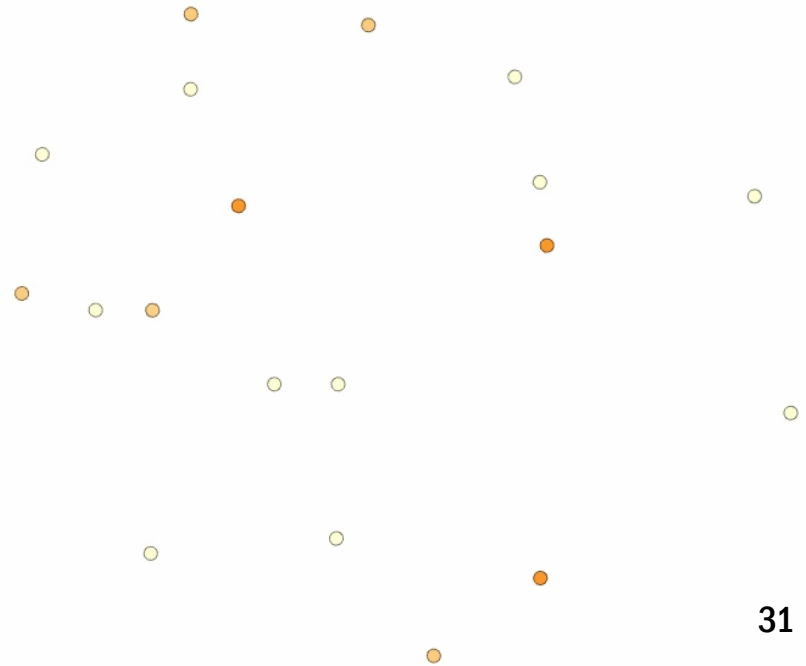


GDELT database

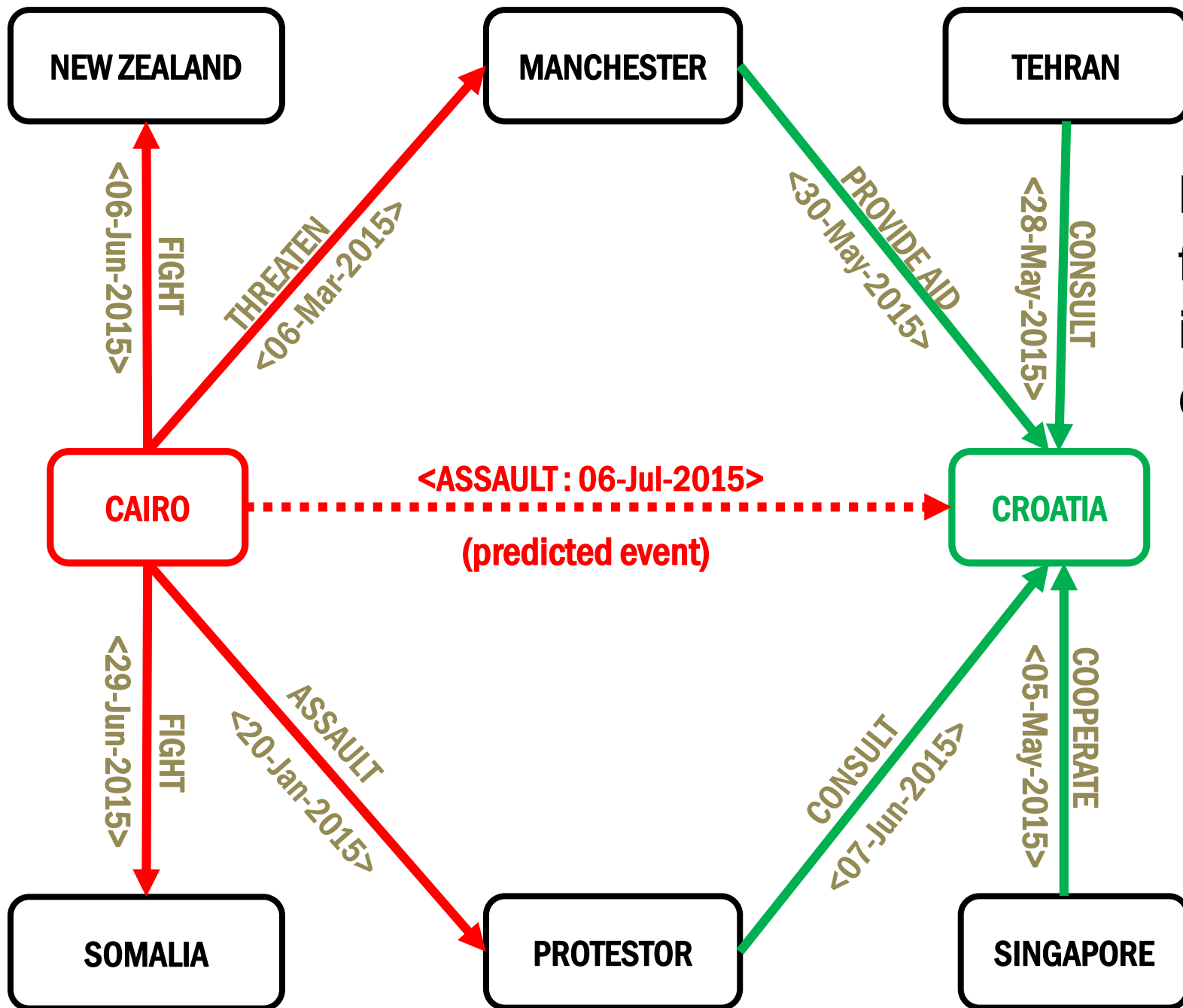
Events in news media

subject – relation – object
and time

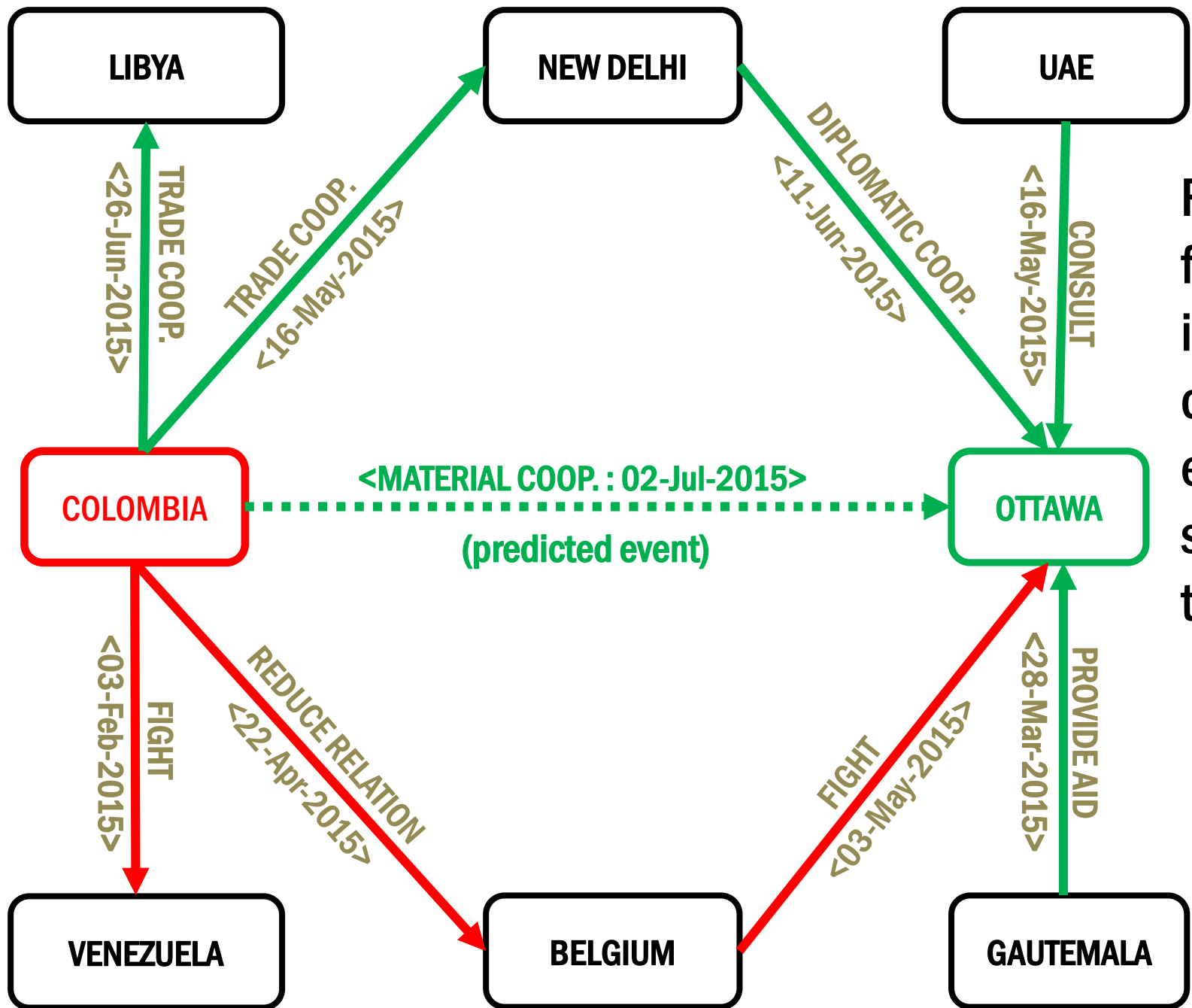
Total archives span >215 years,
trillion of events



Time-varying dependency structure



Enemy's friend is an enemy

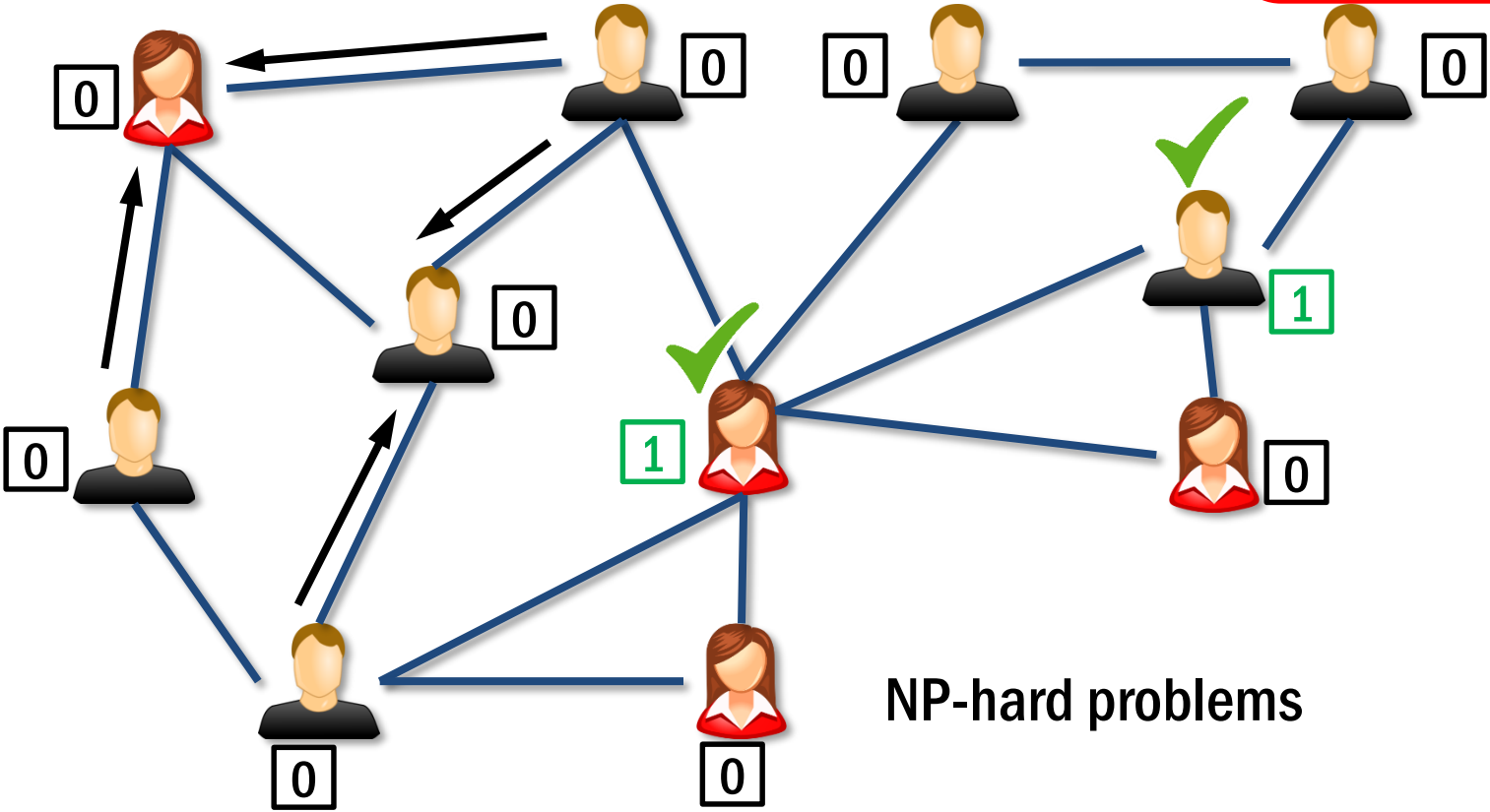


Friends' friend is a friend, common enemy strengthen the tie

Scenario 3: Combinatorial optimization over graph

2 - approximation for minimum vertex cover
Repeat till all edges covered:
• Select uncovered edge with **largest total degree**

Manually designed rule.
Can we learn from data?



NP-hard problems



Greedy algorithm as Markov decision process

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. $x_i + x_j \geq 1, \forall (i, j) \in \mathcal{E}$

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

Reward: $r^t = -1$

State S : current selected nodes

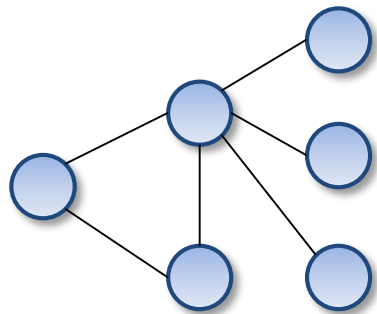
Action value function: $\hat{Q}(S, i)$

Greedy policy:

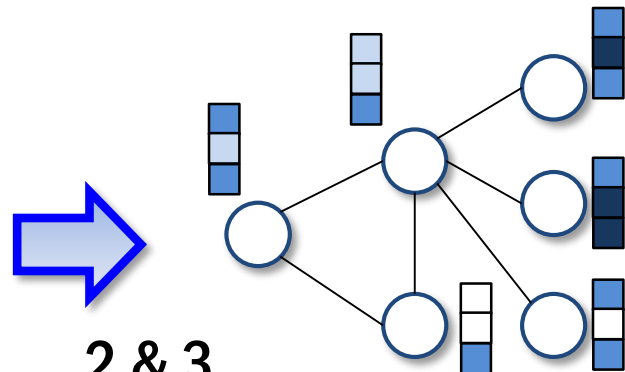
$$i^* = \operatorname{argmax}_i \hat{Q}(S, i)$$

Update state S

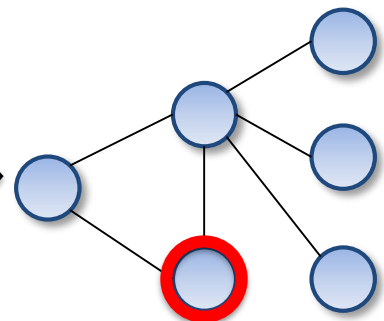
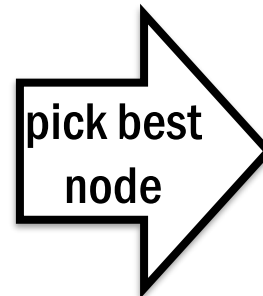
Embedding for state-action value function



1. Problem graph



2 & 3.
Model &
Embed MF



Greedy action
 $i^* = \operatorname{argmax}_i \hat{Q}(S, i)$



State-action value function

4. Q-function

$$\hat{Q}(S, i)$$

$$= \theta_1 \sigma(\theta_2 \sum_{j \in V} \mu_j + \theta_3 \mu_i)$$

aggregated
embedding ↑

↑ individual
embedding

5. Train

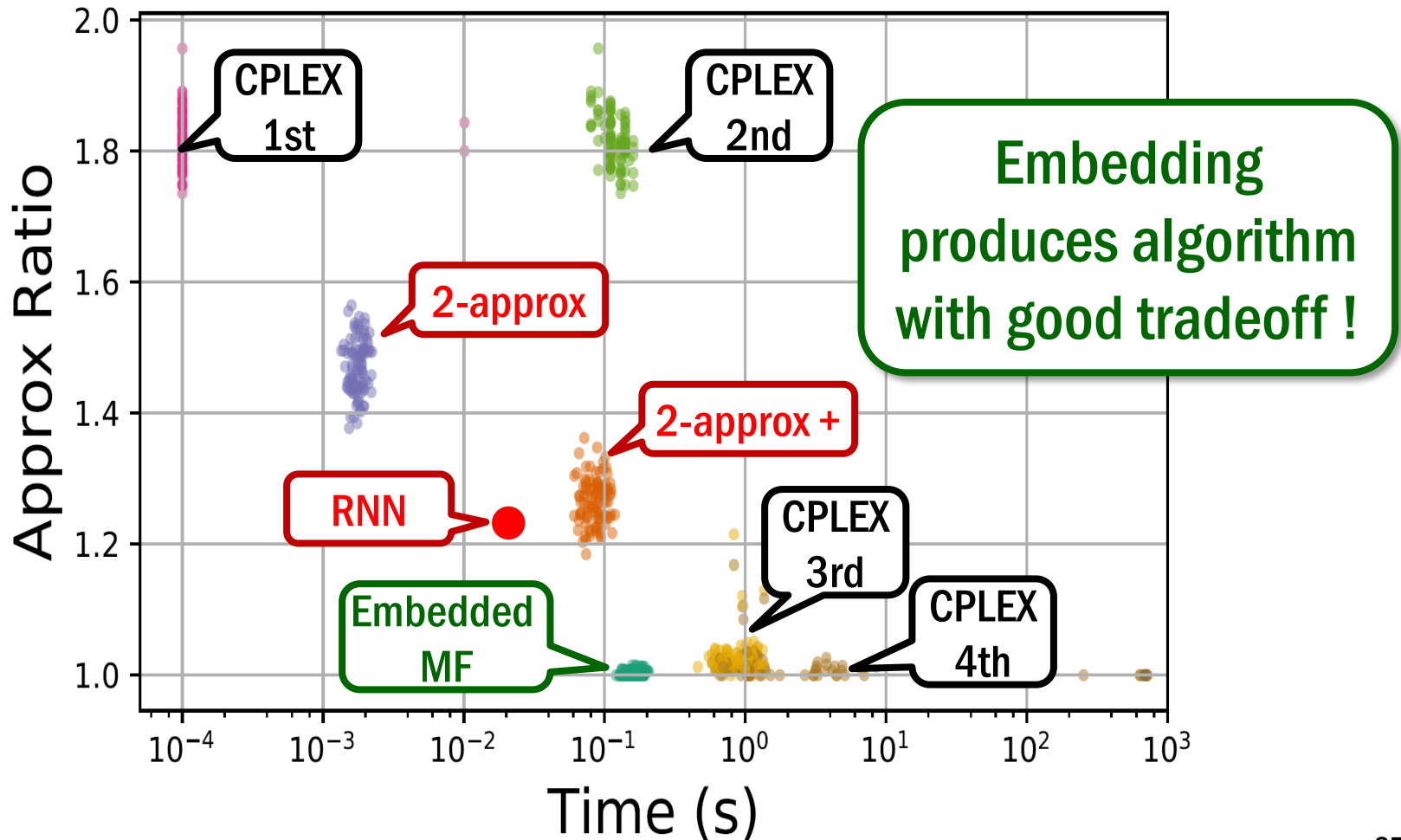


Runtime quality trade-off

Generate 200 Barabasi-Albert networks with 300 nodes

Let CPLEX produces 1st, 2nd, 3rd, 4th feasible solutions

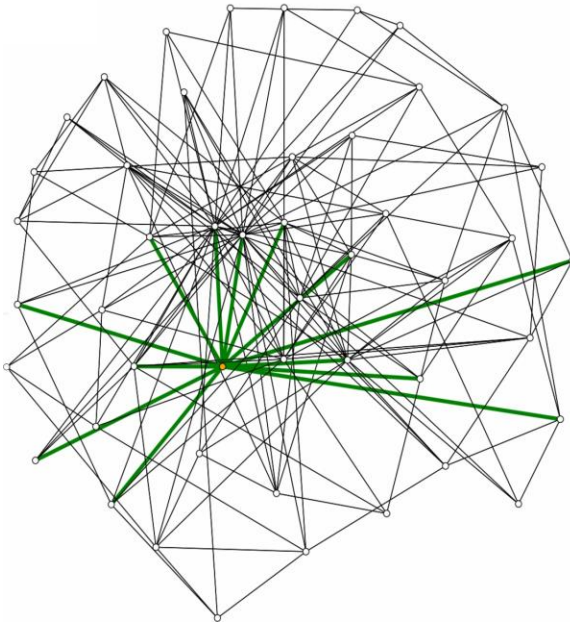
MVC Barabasi-Albert



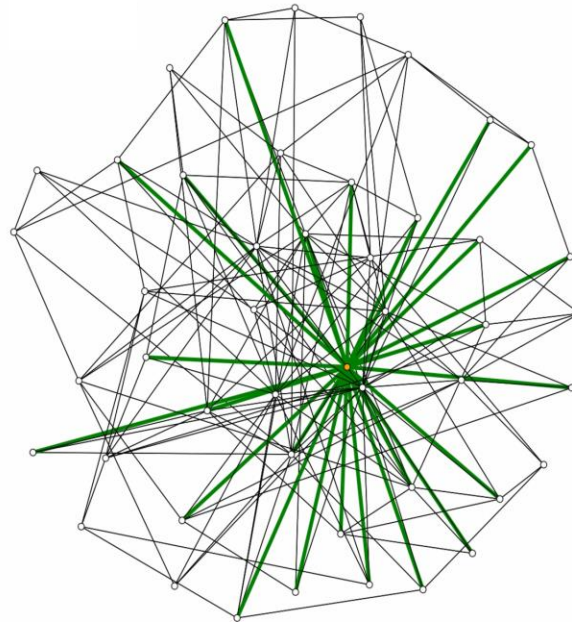
What algorithm is learned?

Learned algorithm balances between

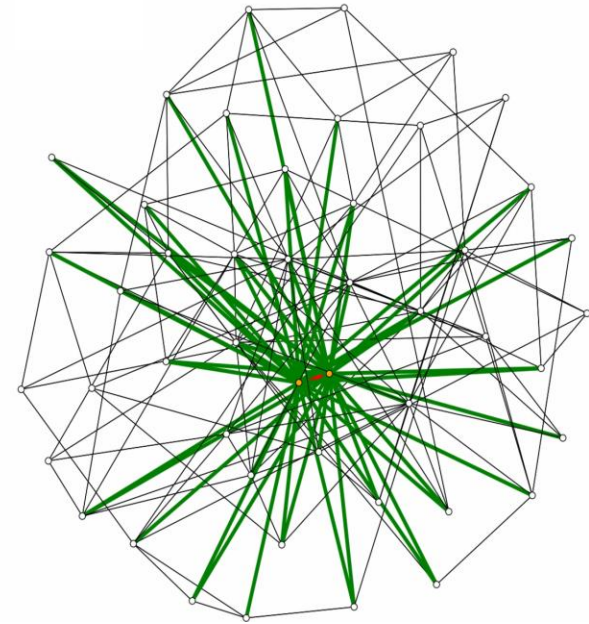
- degree of the picked node and
- fragmentation of the graph



Embedding



Node greedy



Edge greedy

Program with perception and uncertain components

```
result = Operation(a, b)
```

```
result.clear(), carry = 0
```

```
For i in range(len(a)):
```

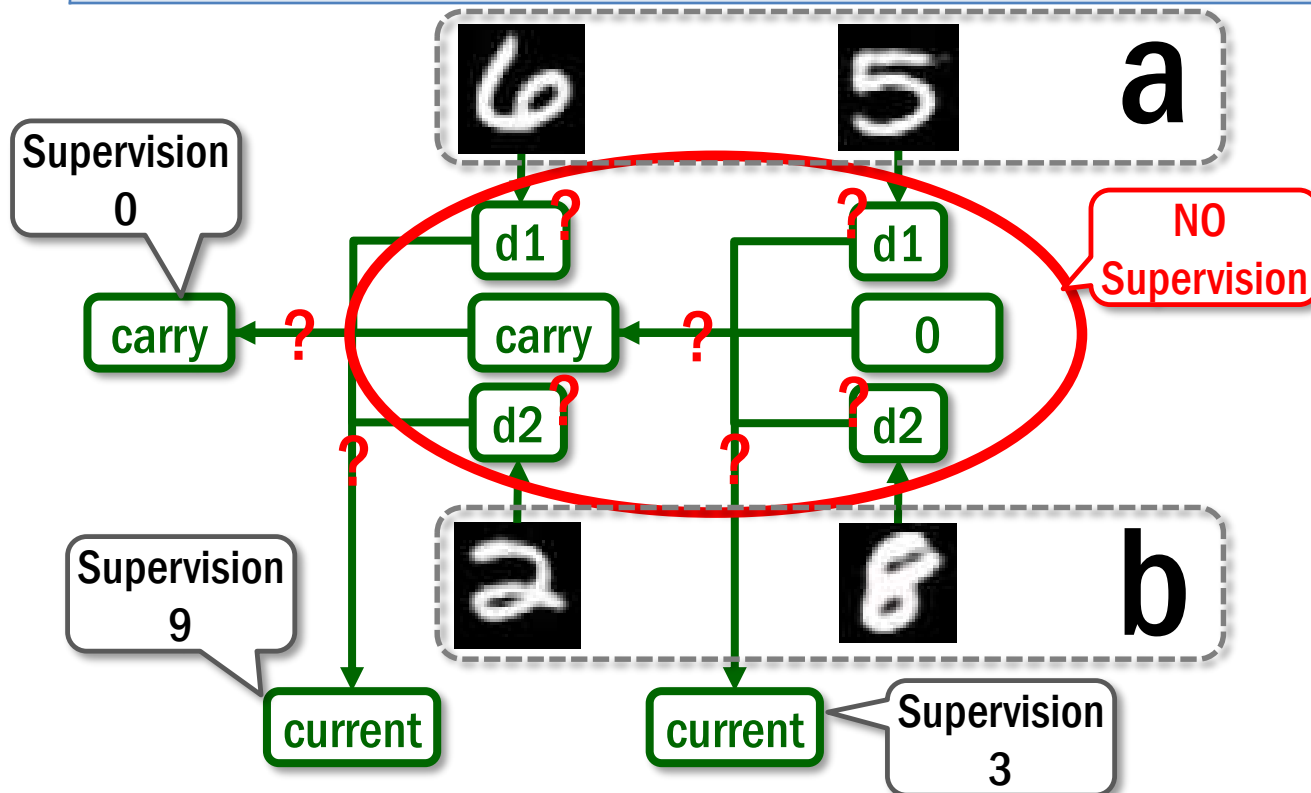
```
    d1 = Recognize(a[i]), d2 = Recognize(b[i])
```

```
    current = Func1(d1, d2, carry), carry = Func2(d1, d2, carry)
```

```
    result.append(current)
```

```
result.append(carry)
```

$$\begin{array}{r}
 \begin{array}{|c|c|} \hline 6 & 5 \\ \hline \end{array} \\
 ? \begin{array}{|c|c|} \hline 2 & 8 \\ \hline \end{array} \\
 \hline
 = \begin{array}{|c|c|} \hline 9 & 3 \\ \hline \end{array}
 \end{array}$$



Algorithm
=
Function
structure

Embedding as a tool for algorithm design

