# Probabilistic Embedding Models

## Robert Bamler, UC Irvine

Internal Blockchain Meeting
7 October 2019

UCI

# Who Are Your **Peers**?

# Word Embeddings



word meanings
(semantic representation)

UCI

# Word Embeddings

"king"    "man"    "woman"    "queen"

▶ Capture **semantic relations**:  −  +  ≈ 

▶ Used for **transfer learning** in natural language processing, e.g., for **sentiment analysis**:

vs.

Robert Bamler

4

UCI

# Distributional Hypothesis

**Assumption:** Words that appear in similar contexts are similar in meaning.

[...] we can not **dedicate** — we can not **consecrate** — we can not **hallow** , this ground [...]

(A. Lincoln, 1863)

## Semantic relations between words:

When Seth had lived 105 years , he became the father of Enosh . [...]

When Enosh had lived 90 years , he became the father of Kenan .

(Gen 5, 6 and Gen 5, 9)

UCI

# "Neural" Word Embeddings: word2vec

**Idea:** Map each word in the vocabulary to a vector in $\mathbb{R}^k$.

Pull words that appear in similar contexts together.

To prevent overall collapse, push random words apart.

classifier cl. cl. fier cl. ? cl. fier cl. cl. cl. cl. cl.

negative negat... negative positive examples exampl... examples examples exampl... examples

... ... ...
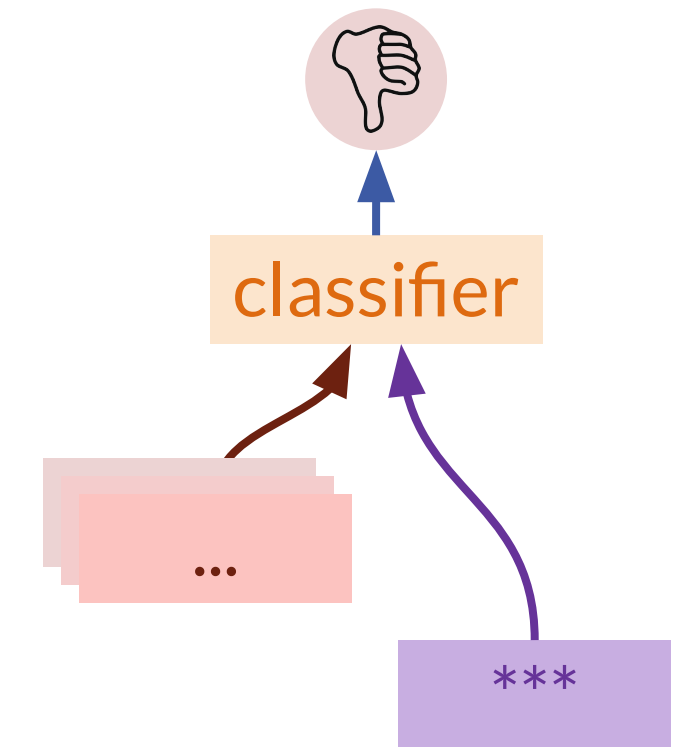
... ... *** *** *** ... ...

UCI

# "**Neural**" Word Embeddings: word2vec

[Mikolov et al., ICLR 2013 & NIPS 2013]

Intuitively:

Pull words that appear in similar **contexts** together.

To prevent overall collapse, push random words apart.

Minimize loss function: $\ell = -\sum_{(i,j)\in\text{pos.}} \log \sigma(u_i^\top v_j) - \sum_{(i,j)\in\text{neg.}} \log \sigma(-u_i^\top v_j)$

**word embedding**
(vector in $\mathbb{R}^k$)

**context embedding**
(vector in $\mathbb{R}^k$)

UCI

# Our Extension: Dynamic Word Embeddings

[Bamler & Mandt, ICML 2017]

"Computer" **in 1961**



© 20th century FOX

"Computer" **today**

UCI

# Detecting **Subtle Changes** Over Time

[Bamler & Mandt, ICML 2017]

**Naive idea:**
Fit individual embedding vector for each word and each year.

**Problem:**
Only few data **per word & year**.
→ *small signal/noise ratio*

Our Solution:

1999

**probabilistic** version of word embeddings

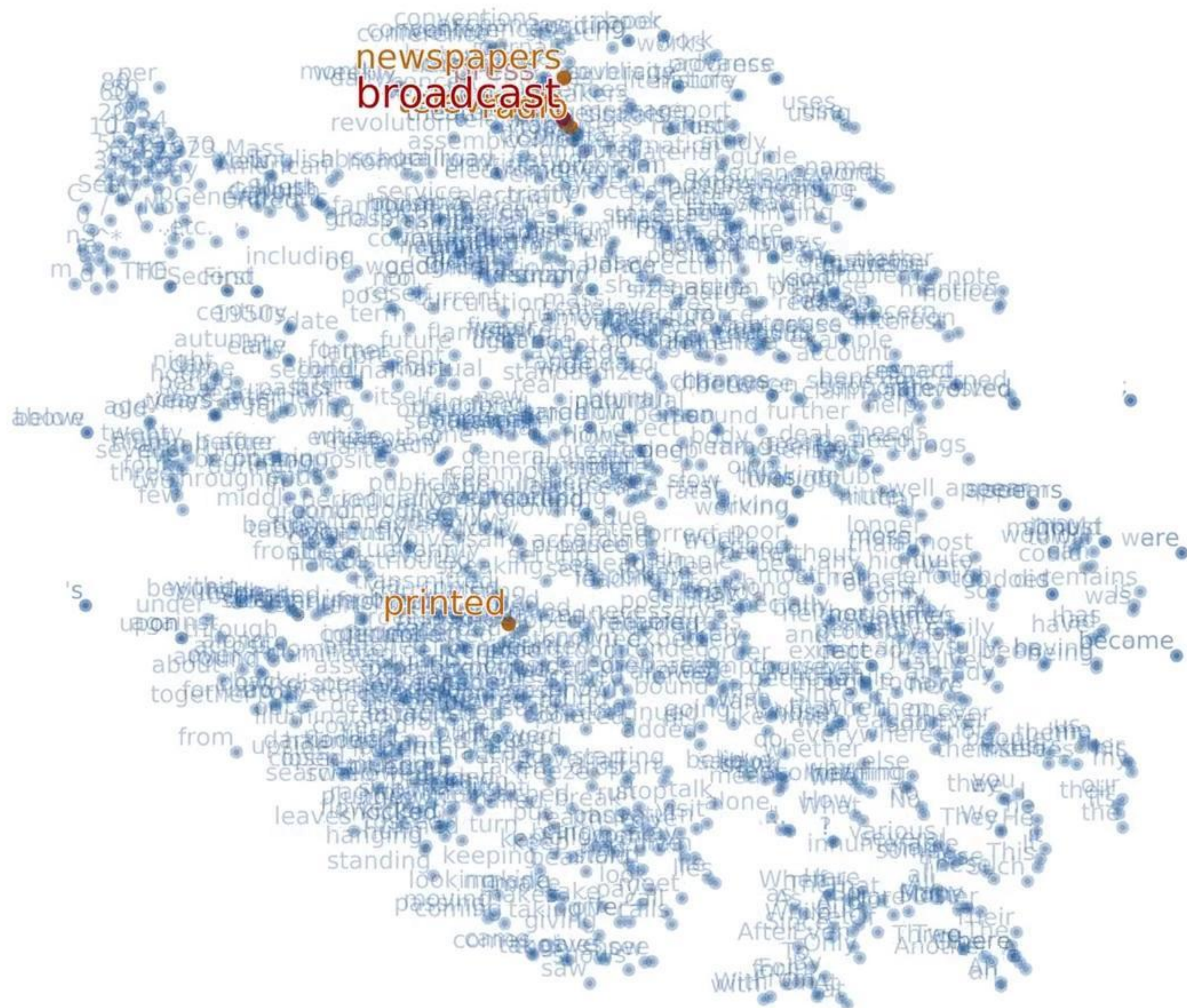**+** **probabilistic** model of the dynamics

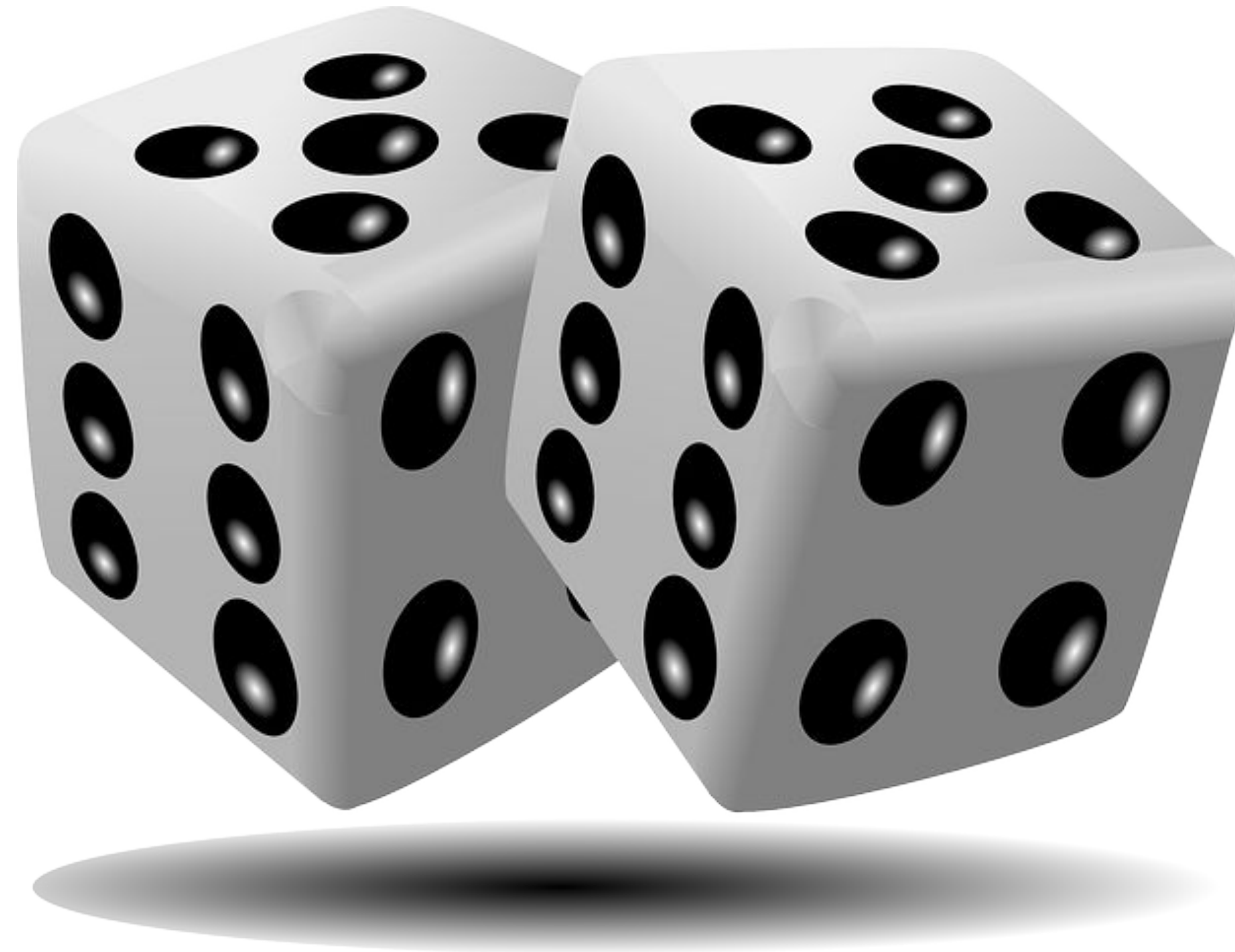**+** scalable approximate Bayesian inference

UCI

198**b**

broadcast
radio

television
newspapers
TV
printed
press

# Bear With Me: Probabilistic Models & Inference

UCI

# Probabilistic Models & Bayesian Inference

Notation: *x* = observations (data)
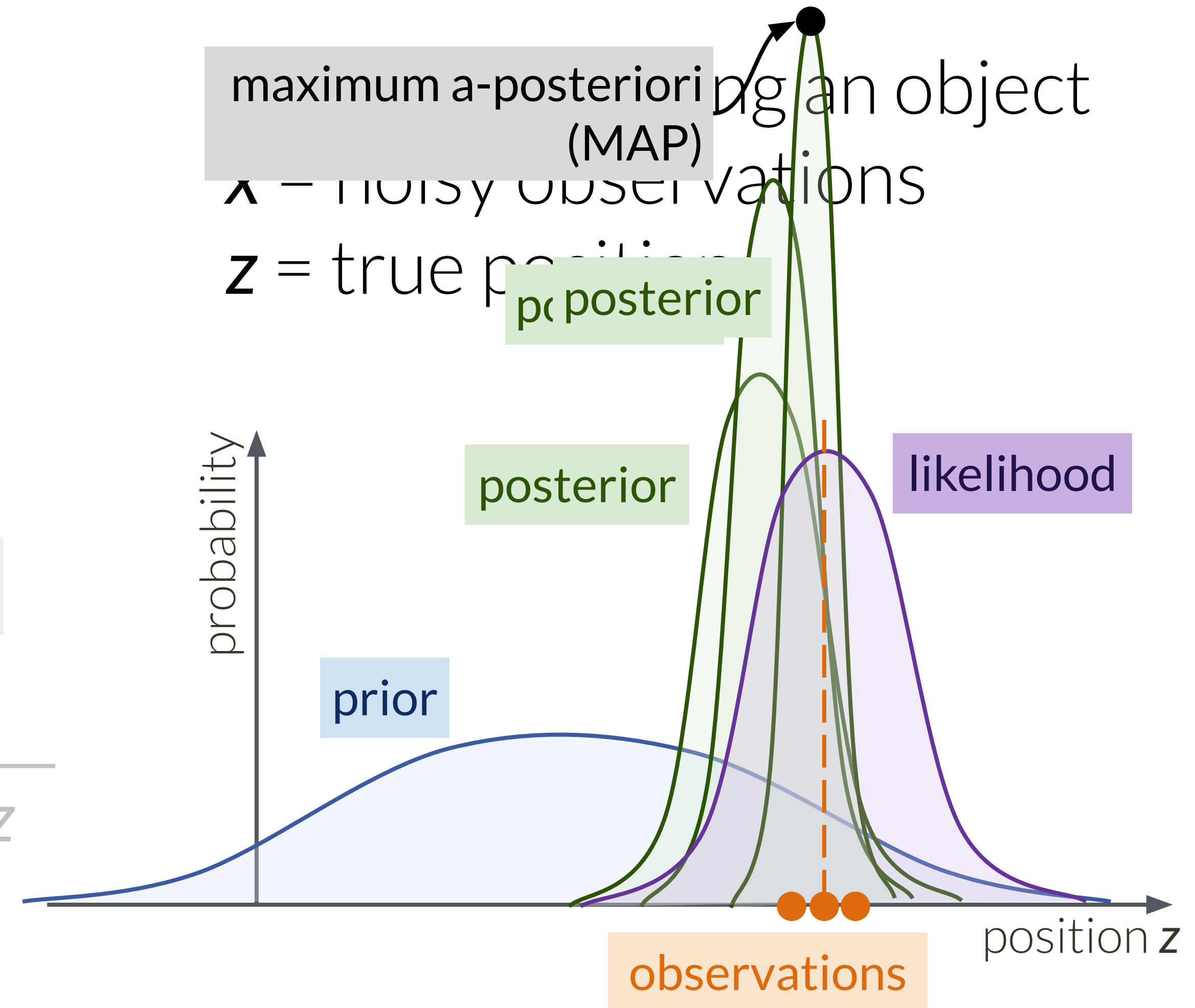　　　　　 *z* = latent (i.e., unknown)
　　　　　 variables that caused *x*

Probabilistic model: $p(x, z) = p(z)\, p(x|z)$

prior　　likelihood

Inference: find **posterior** $p(z|x) = \dfrac{p(x, z)}{\int p(x, z)\, dz}$

probability that latent variables *z* explain the observed data *x*.

maximum a-posteriori ng an object
(MAP)

**X** = noisy observations

*z* = true position

posterior

posterior

likelihood

probability

prior

position *z*

observations

UCI

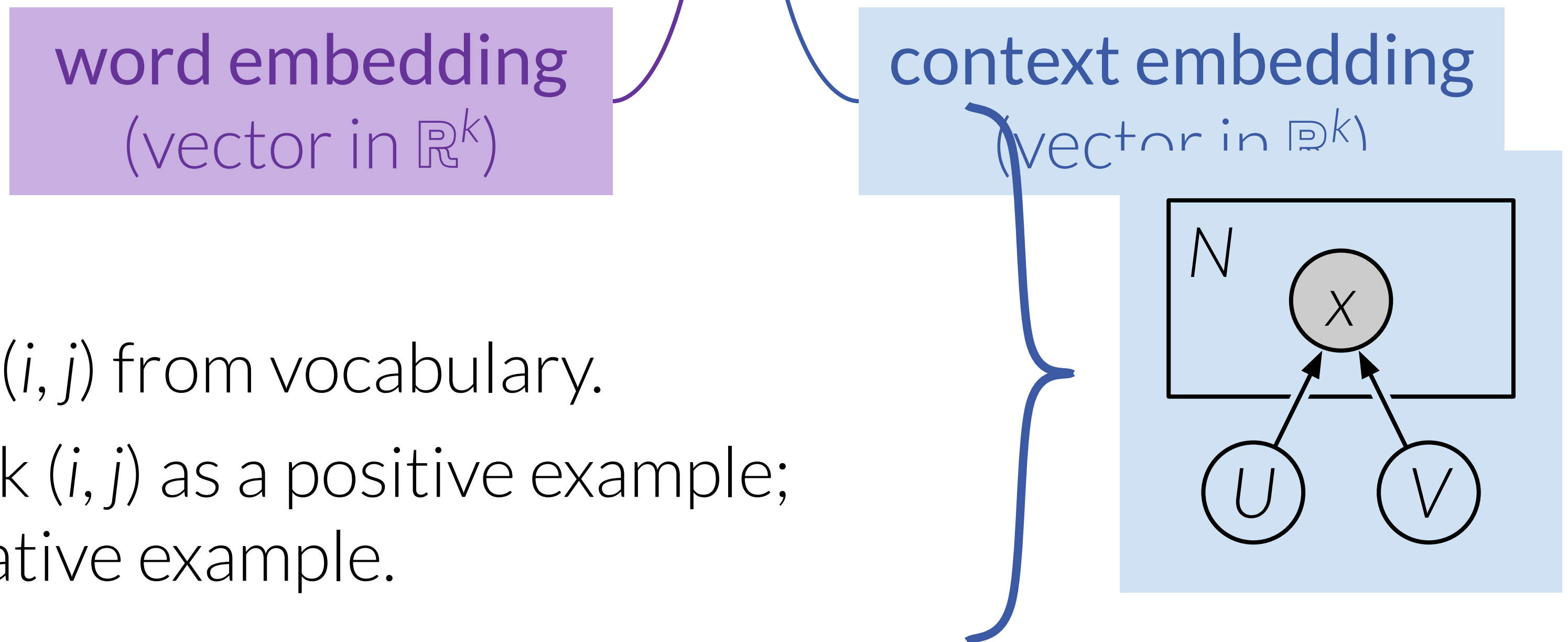# Example 1: Probabilistic Variant of word2vec  [Barkan AAAI 2017]

**Reminder:** word2vec minimizes loss  $\ell = -\sum_{(i,j)\in\text{pos.}} \log \sigma(u_i^\top v_j) - \sum_{(i,j)\in\text{neg.}} \log \sigma(-u_i^\top v_j)$

**Observation:** $\ell = -\log p(x|z)$

**word embedding**
(vector in $\mathbb{R}^k$)

**context embedding**
(vector in $\mathbb{R}^k$)

**Generative process:**

*repeat:*

- Draw a random pair of words $(i, j)$ from vocabulary.

- **With probability $\sigma(u_i^\top v_j)$:** mark $(i, j)$ as a positive example; **otherwise:** mark $(i, j)$ as a negative example.
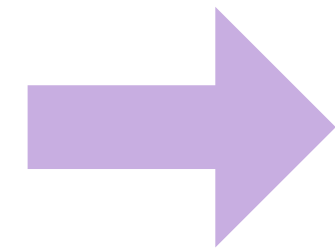


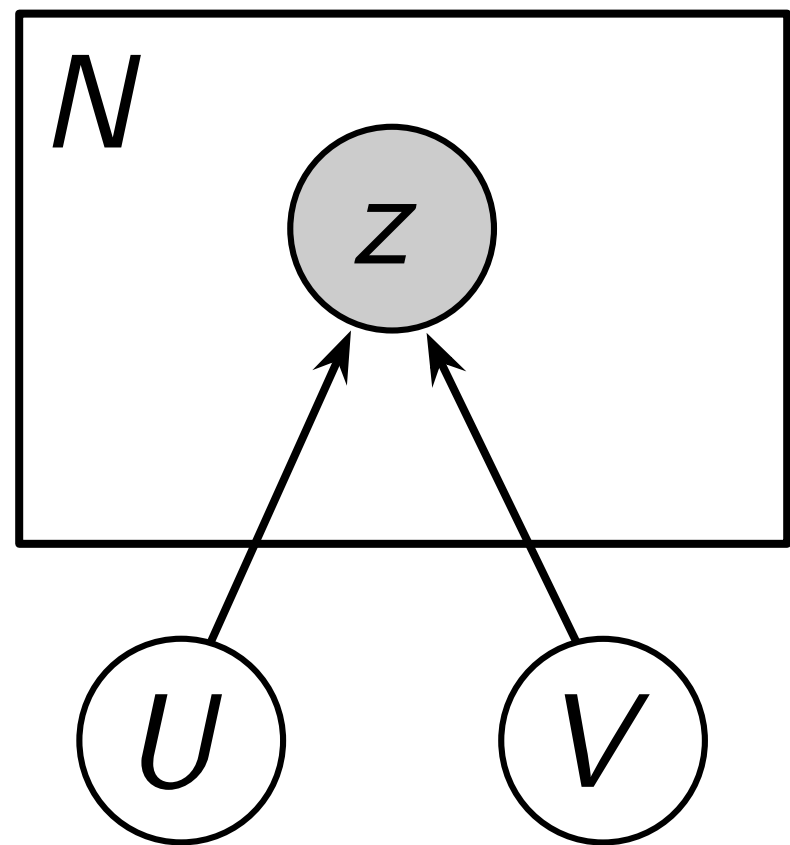UCI

# Enough Equations, **Back to Pretty Pictures**
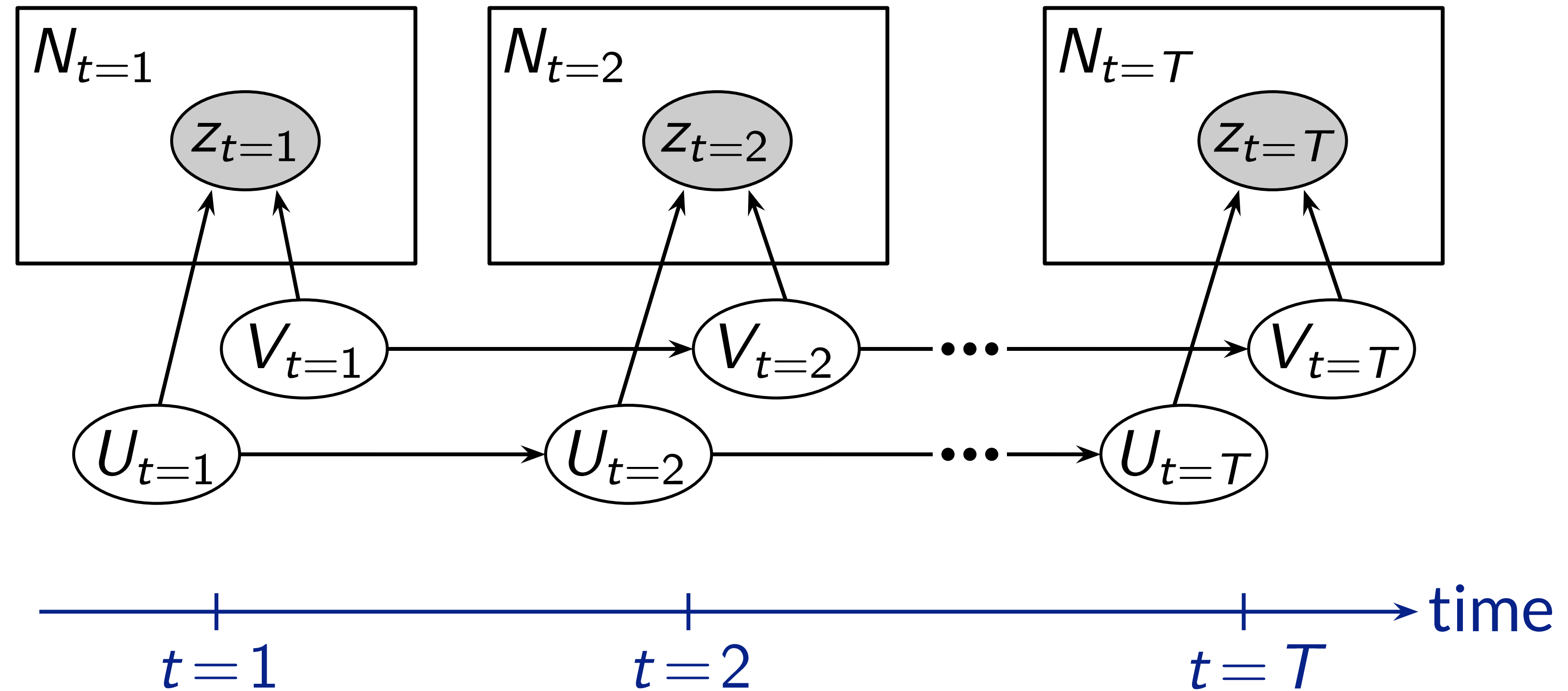
Nymphenburg Park, **Munich** #nofilter

UCI

# Dynamic Word Embeddings Model

[Bamler & Mandt, ICML 2017]

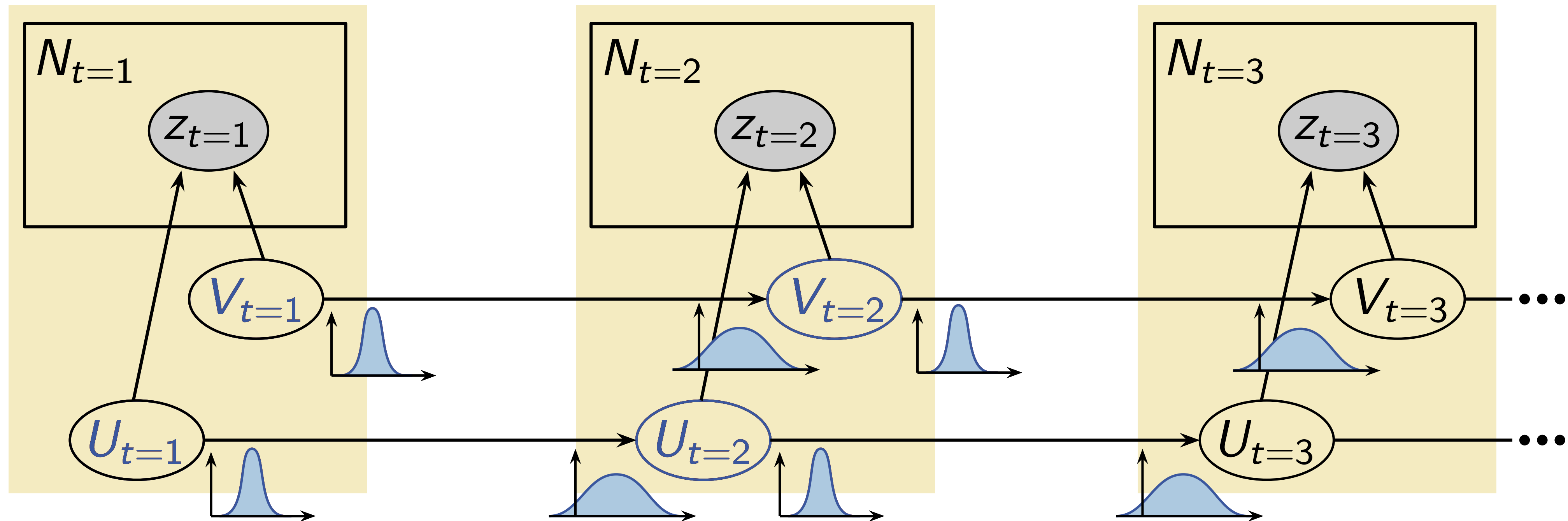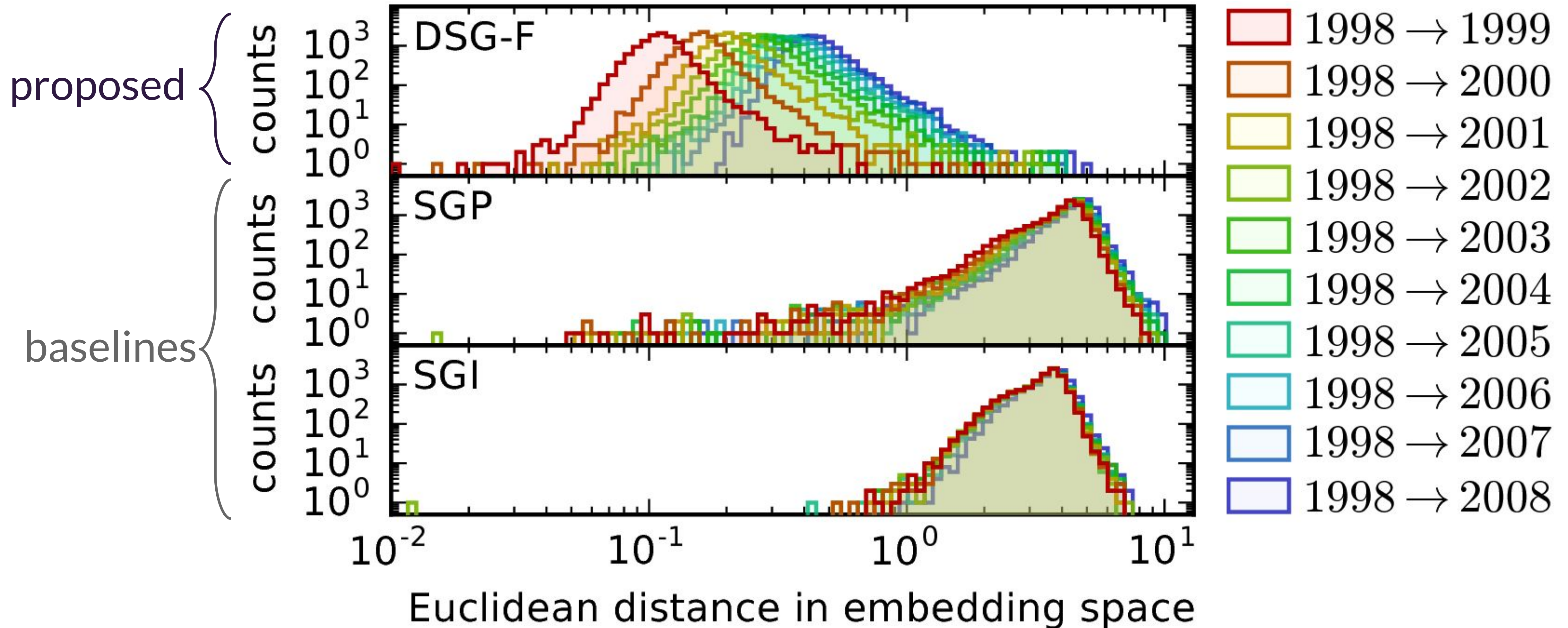static model $\Rightarrow$ Dynamic Model

UCI

# On-line Learning With Variational Inference

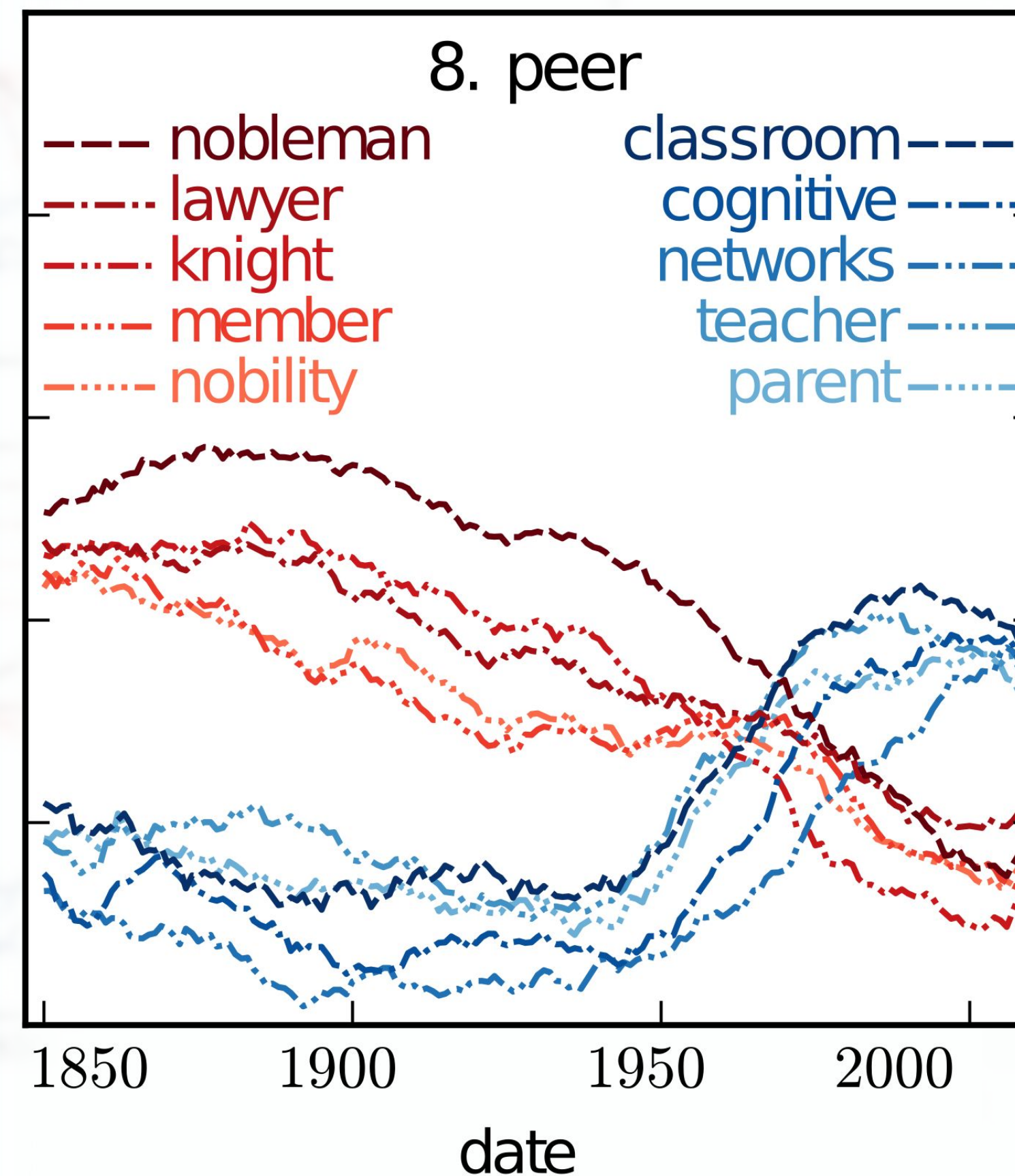[see "Dynamic Word Embeddings", Bamler & Mandt, ICML 2017]

UCI

# Results III: Smooth Trajectories in Embedding Space

(Training data: Google Books corpus [Michel et al., 2011])

UCI

# Results: Top 10 Most Mobile Words (1850→2008)

(Training data: Google Books corpus [Michel et al., 2011])

UCI

# Time Travel

word embeddings
for year **1800**

word embeddings
for year **1801**

word embeddings
for year **2008**

**"car"**

Kalman filter

books written
in **1800**

books written
in **1801**

books written
in **2008**

Robert Bamler

UCI

# Results: Word Aging with Goldstone GD

[Bamler & Mandt, ICML 2018]

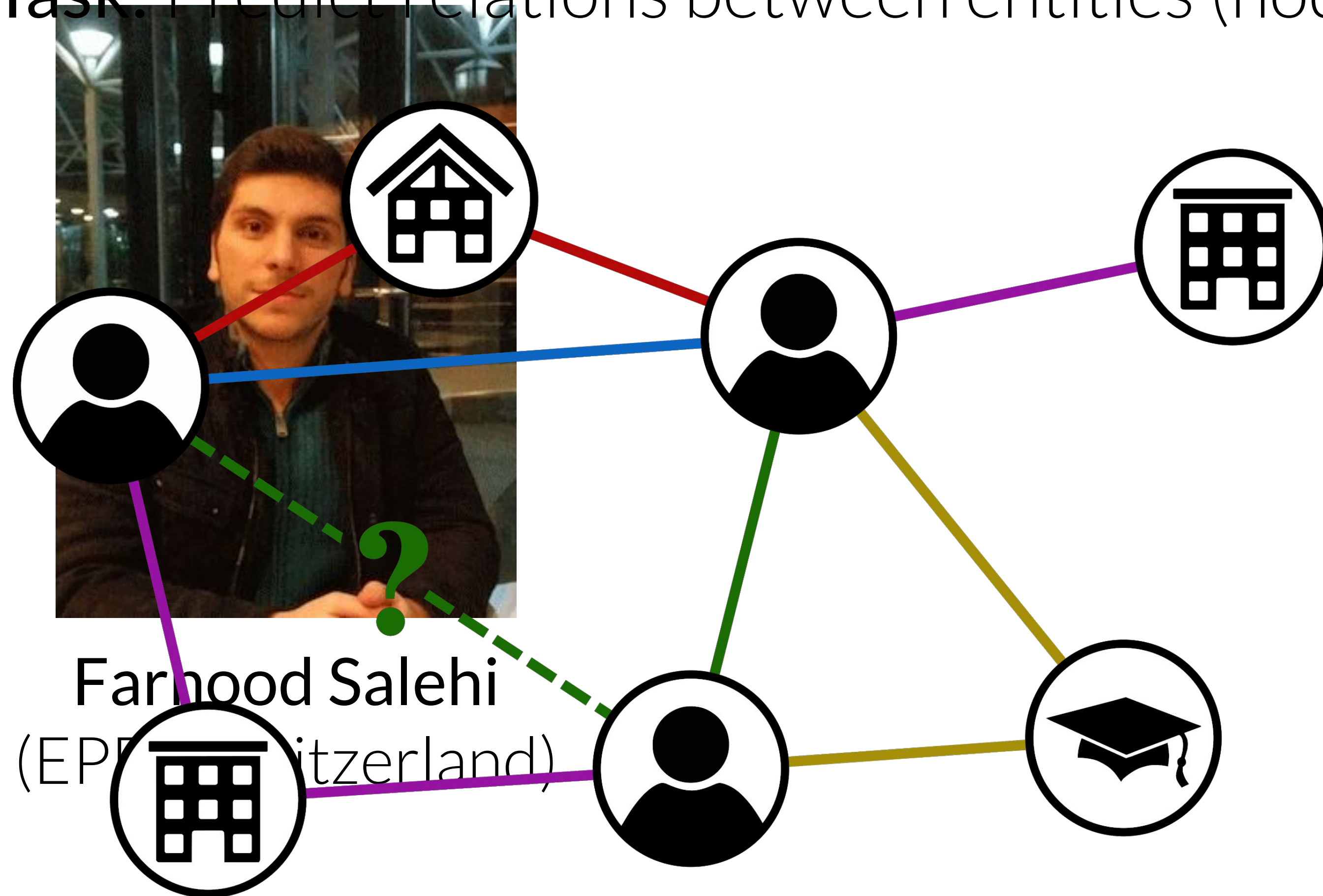| 2008 | 1800 |
|---|---|
| *car* | boat, saddle, canoe, wagon, box |
| *DNA* | potassium, chemical, sodium, molecules, displacement |
| *tuberculosis* | chronic, paralysis, irritation, disease, vomiting |

UCI

# Example 2: Probabilistic Knowledge Graphs

[Bamler, Salehi & Mandt, UAI 2019]

**Task:** Predict relations between entities (nodes) in a knowledge graph.



Farhood Salehi
(EPFL, Switzerland)

UCI

# Example 2: Probabilistic Knowledge Graphs

[Bamler, Salehi & Mandt, UAI 2019]

**Task:** Predict relations between entities (nodes) in a knowledge graph.

**State of the Art:**
Learn embeddings for entities and relation types.

**Problem:**
Highly sensitive to hyperparameters
[Kadlec et al., 2017]



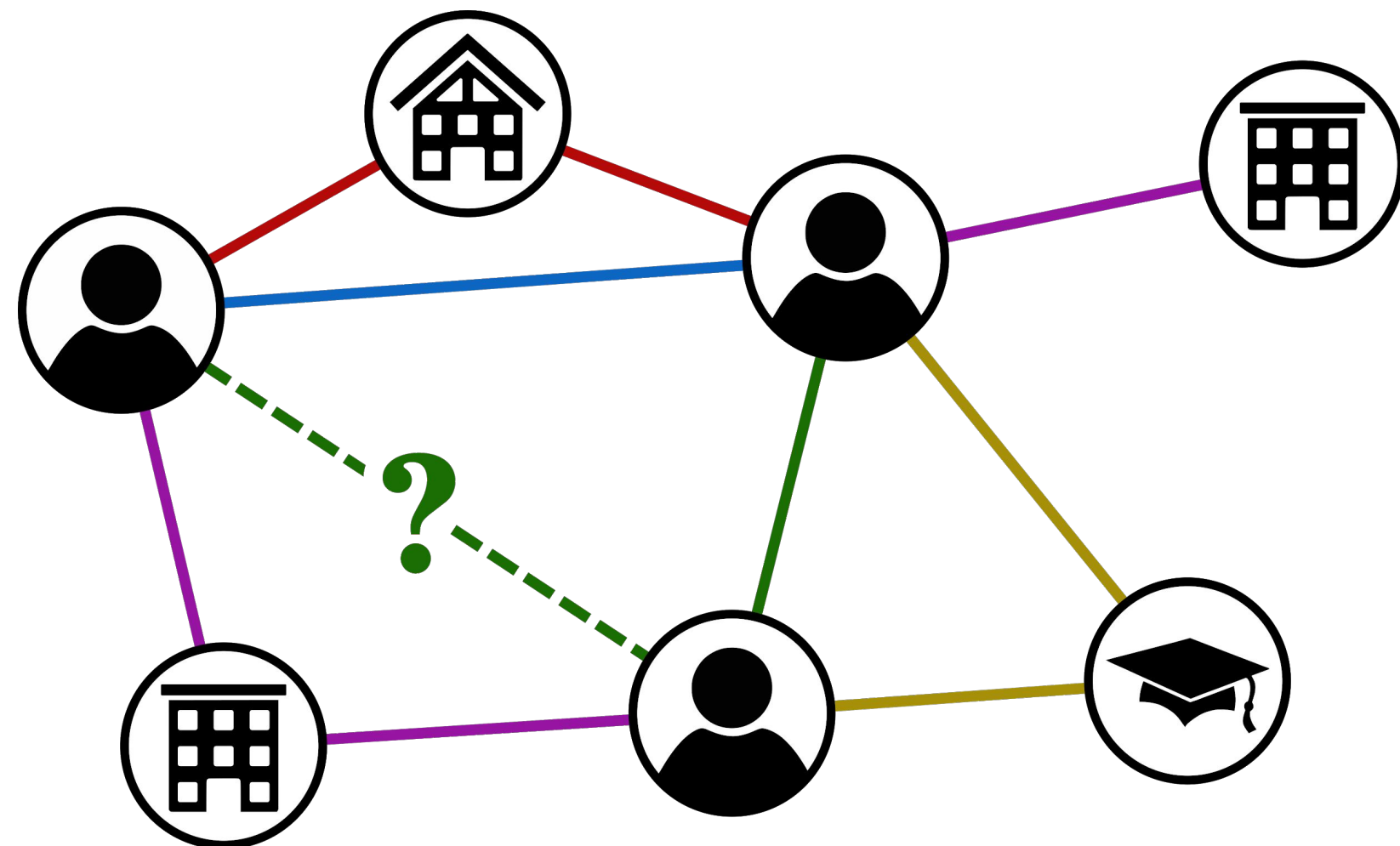**Reason:** Many entities are supported only by few data points.

UCI

# Example 2: Probabilistic Knowledge Graphs

[Bamler, Salehi & Mandt, UAI 2019]

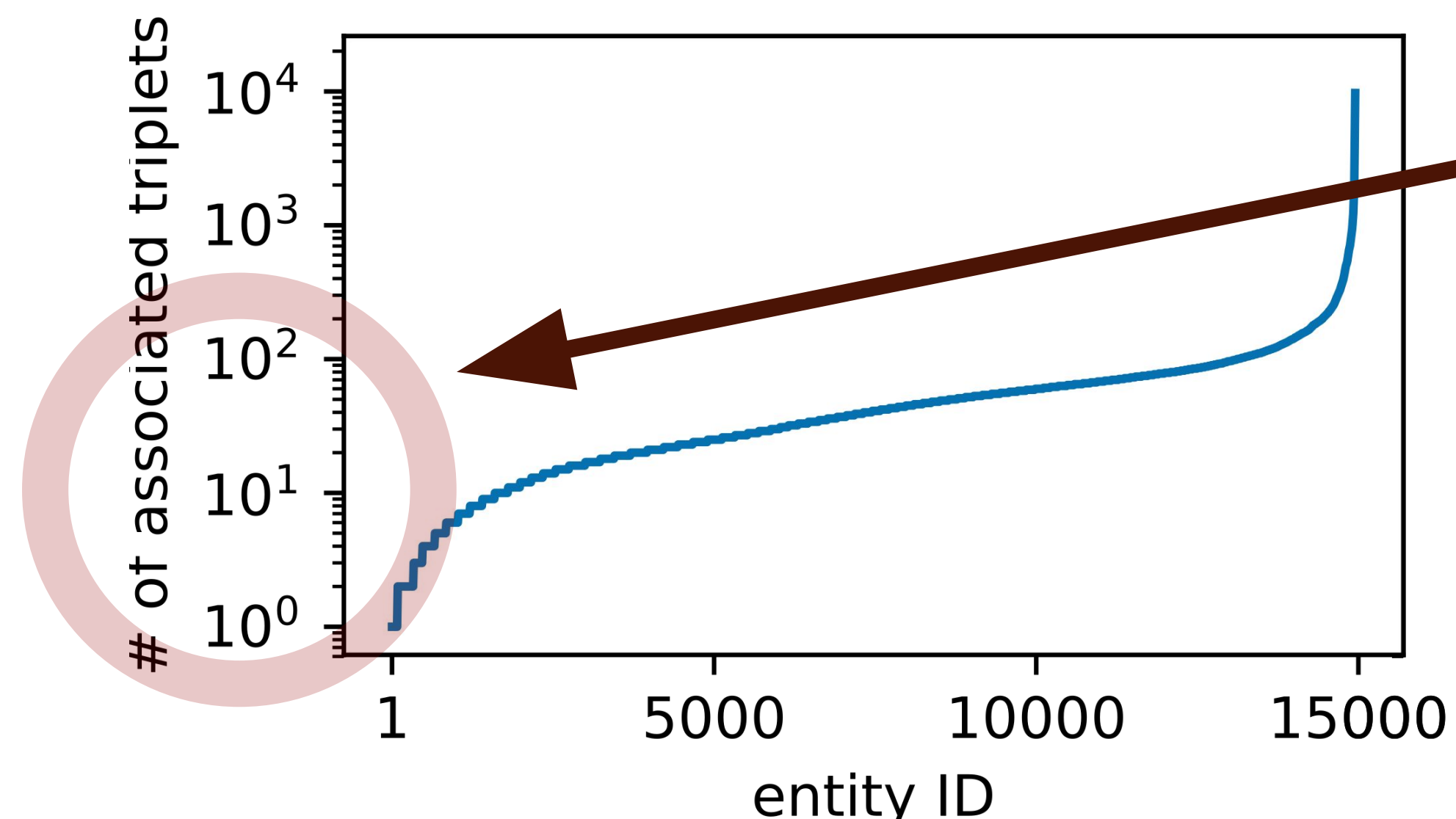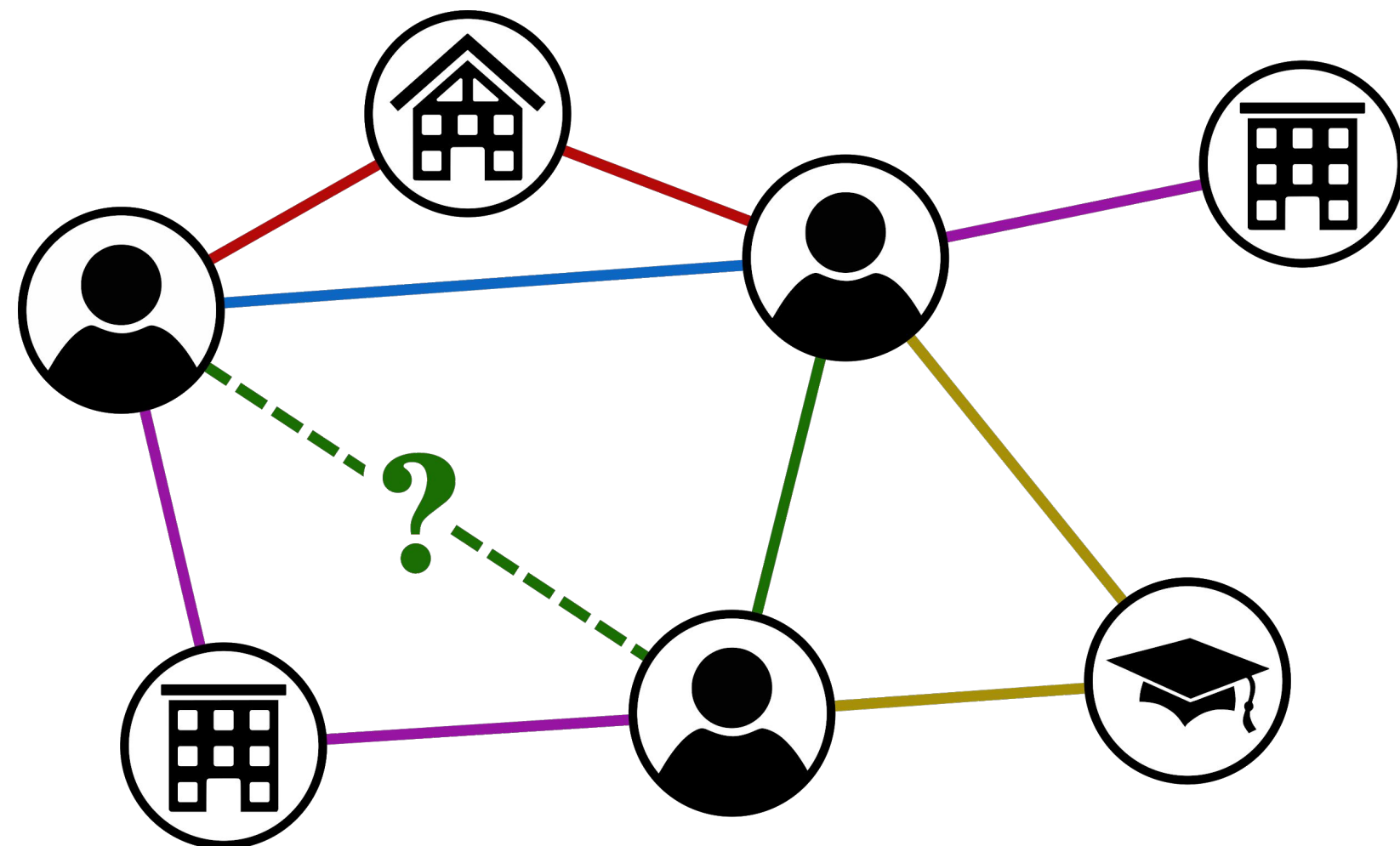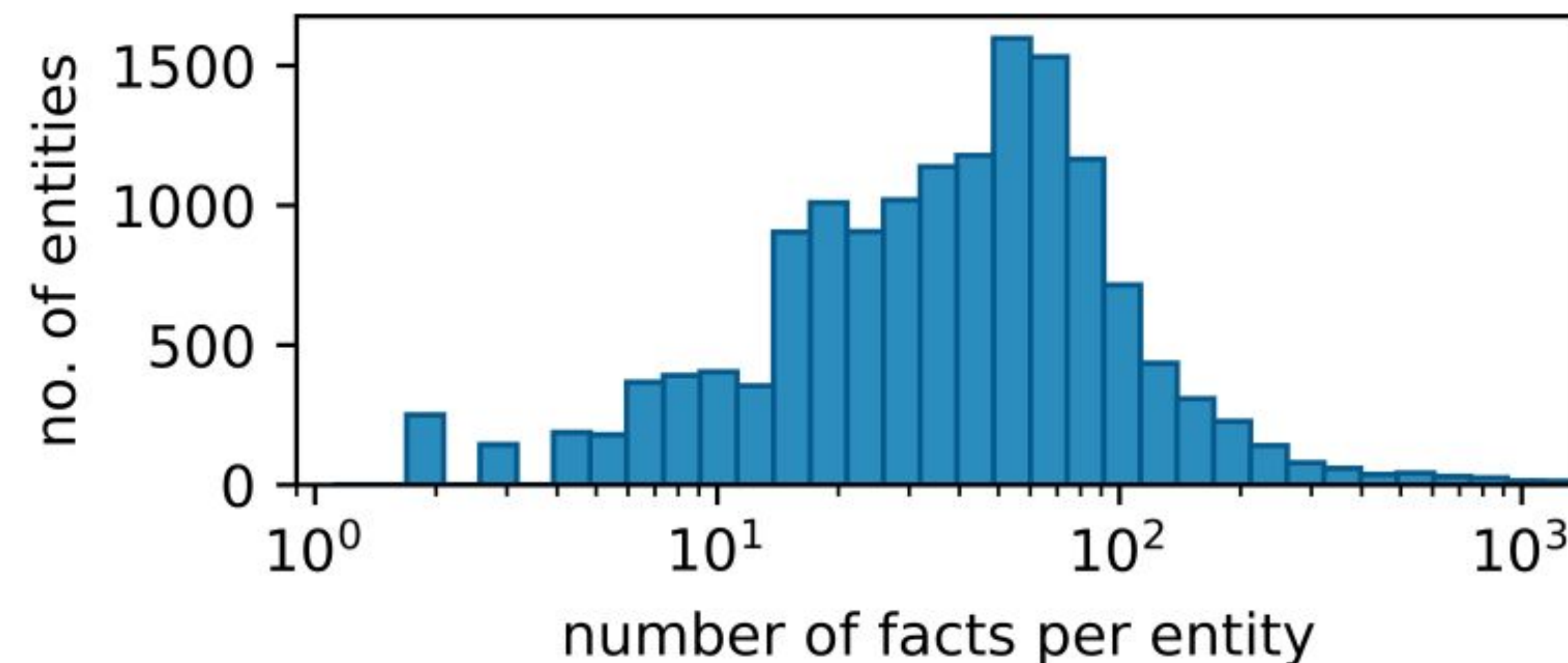**Task:** Predict relations between entities (nodes) in a knowledge graph.

## State of the Art:

Learn embeddings for entities and relation types.

## Problem:

Highly sensitive to hyperparameters
[Kadlec et al., 2017]
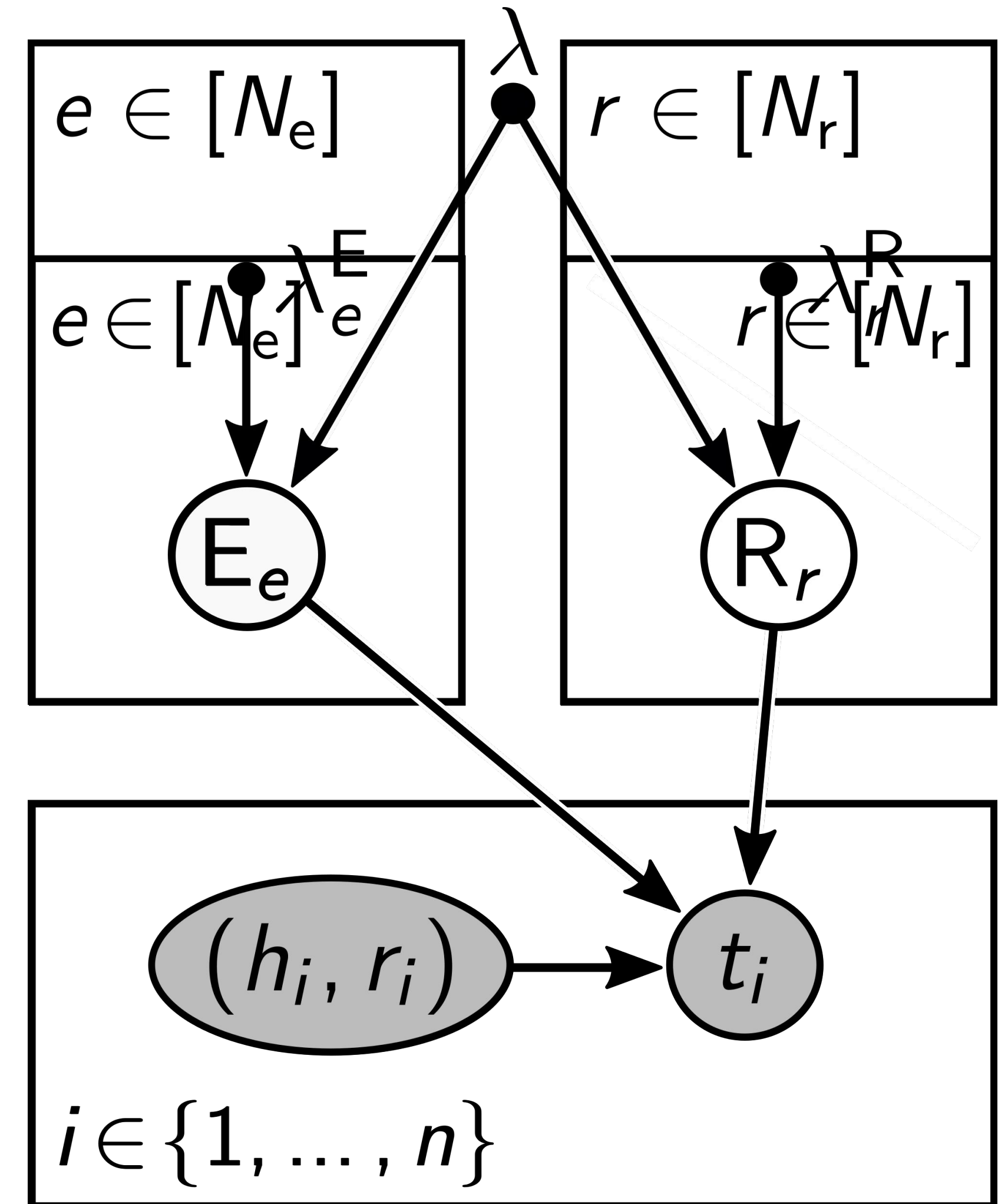
UCI

# Example 2: Probabilistic Knowledge Graphs

[Bamler, Salehi & Mandt, UAI 2019]

## Our Solution:

→ Reinterpret existing models as **probabilistic generative models** of relational facts (**h**ead, **r**elation, **t**ail).

→ Introduce macroscopic number of *local* hyperparameters **(> 10,000)**.

→ Tune hyperparameters efficiently with **variational expectation maximization**.
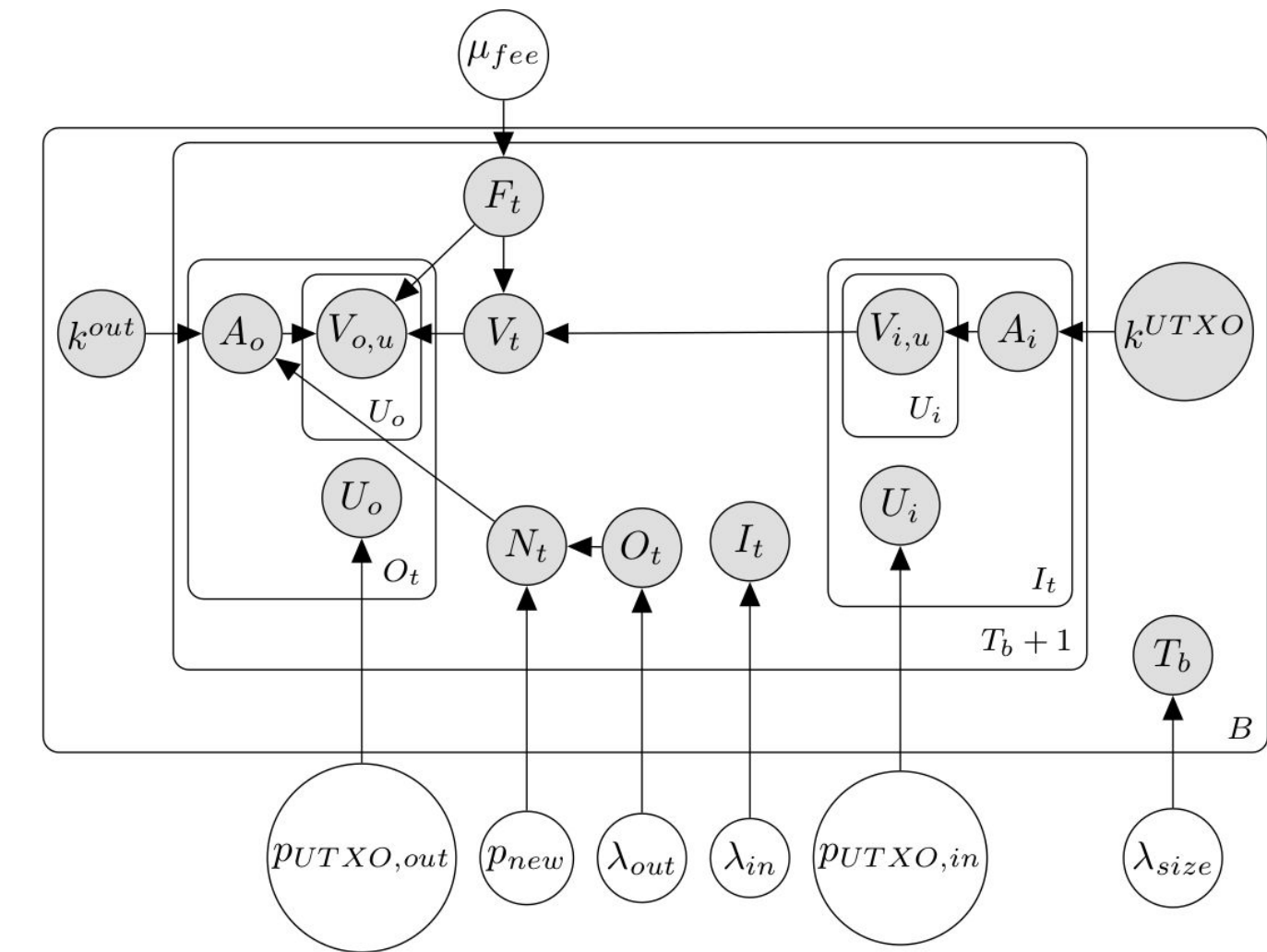
UCI

# Results: Probabilistic Knowledge Graph Embeddings

[Bamler, Salehi & Mandt, UAI 2019]

Link prediction outperforms previous state of the art.

| | data set → | WN18RR | | WN18 | | FB15K-237 | | FB15K | |
|---|---|---|---|---|---|---|---|---|---|
| ↓ model | ↓ variant | MRR | Hits@10 | MRR | Hits@10 | MRR | Hits@10 | MRR | Hits@10 |
| DistMult | Yang et al. [2015] (orig.) | – | – | 0.83 | 0.942 | – | – | 0.35 | 0.577 |
| DistMult | Kadlec et al. [2017] | – | – | 0.790 | 0.950 | – | – | 0.837 | 0.904 |
| DistMult | Dettmers et al. [2018] | 0.43 | 0.49 | 0.822 | 0.936 | 0.241 | 0.419 | 0.654 | 0.824 |
| DistMult | Ours (after variational EM) | **0.455** | **0.544** | **0.911** | **0.961** | **0.357** | **0.548** | **0.841** | **0.914** |
| ComplEx | Trouillon et al. [2016] (orig.) | – | – | 0.941 | 0.947 | – | – | 0.692 | 0.840 |
| ComplEx | Lacroix et al. [2018]* | 0.478 | 0.569 | 0.952 | 0.963 | 0.364 | 0.555 | **0.857** | 0.909 |
| ComplEx | Ours (after variational EM) | **0.486** | **0.579** | **0.953** | **0.964** | **0.365** | **0.560** | 0.854 | **0.915** |

UCI

# Discussion: Embedding the Blockchain

▶ **Analogy to knowledge graphs:**
transaction ≈ relational fact; token ≈ relation

▶ **Temporal component:**
≈ Dynamic Word Embeddings, but more ephemeral.



[Jourdan et al., 2019]

▶ **Analysis (ideas):**
▷ semantic analysis of users & tokens using embeddings
▷ predict transactions?

▶ **Action Items:**
▷ *Get a data set.*
▷ Discuss existing literature.
▷ Come up with more concrete analysis questions (maybe as we look at the data).

UCI