

Multi-Modal Bayesian Embeddings for Learning Social Knowledge Graphs

Zhilin Yang¹², Jie Tang¹, William W. Cohen²

¹Tsinghua University

²Carnegie Mellon University

AMiner: academic social network



data mining

14



Jiawei Han (韩家炜) ✓

h-index: 134 | **#Paper:** 819 | **#Citation:** 81713

Professor

Department of Computer Science, University of Illinois at Urbana-Champaign

Similar

Data Mining | Data Cube | Data Analysis | Information Retrieval | Indexation

33792 views

4



Philip S. Yu ✓

h-index: 134 | **#Paper:** 911 | **#Citation:** 79798

Professor and Wexler Chair in Information Technology

Department of Computer Science, University of Illinois Chicago

Similar

Data Mining | Indexation | Social Network | Internet | Data Analysis

1682 views

5



Christos Faloutsos ✓

h-index: 103 | **#Paper:** 550 | **#Citation:** 55079

Professor

Dept. of Computer Science Carnegie Mellon University

Similar

Data Mining | Social Network | Large Graph | Power Law | Graph Theory

1223 views

Research interests



Text-Based Approach

818
Entity Linking with a Knowledge Base: Issues, Techniques, and Solutions
Wei Shen, Jianyong Wang, **Jiawei Han**
Knowledge and Data Engineering, IEEE Transactions (2015)
Bibtex <http://dx.doi.org/10.1109/TKDE.2014.2327028>

817
Power-Based Diagnosis of Node Silence in Remote High-End Sensing Systems
Yong Yang, Lu Su, Mohammad Maifi Hasan Khan, Michael LeMay, Tarek F. Abdelzaher, **Jiawei Han**
TOSN (2015)
Bibtex <http://doi.acm.org/10.1145/2661639>

816
A Framework of Mining Trajectories from Untrustworthy Data in Cyber-P
Lu An Tang, Ying Xu, Quanqun Gu, **Jiawei Han**, Guofei Jiang, Aïme Loupa, Thomas F. La



Jiawei Han (韩家炜)

H-Index: 133 | #Paper: 818 | #Citation: 111164

Department of Computer Science, University of Illinois at Urbana-Champaign

Professor

Similar

Data Mining | Information Extraction | Data Analysis | Machine Learning

List of publications



Infer

Research interests

Text-Based Approach

818
Entity Linking with a Knowledge Base: Issues, Techniques, and Solutions
Wei Shen, Jianyong Wang, **Jiawei Han**
Knowledge and Data Engineering, IEEE Transactions (2015)
Bibtex <http://dx.doi.org/10.1109/TKDE.2014.2327028>

817
Power-Based Diagnosis of Node Silence in Remote High-End Sensing Systems
Yong Yang, Lu Su, Mohammad Maifi Hasan Khan, Michael LeMay, Tarek F. Abdelzaher, **Jiawei Han**
TOSN (2015)
Bibtex <http://doi.acm.org/10.1145/2661639>

816
A Framework of Mining Trajectories from Untrustworthy Data in Cyber-P...
Lu An Tang, Ying Yu, Quanqun Gu, **Jiawei Han**, Guofei Jiang, Ailee Leung, Thomas F. Le...



Jiawei Han (韩家炜) ✓

H-Index: 133 | **#Paper:** 818 | **#Citation:** 111164

Department of Computer Science, University of Illinois at Urbana-Champaign

Professor

Similar

Data Mining | Information Extraction | Data Analysis | Machine Learning

Term Frequency => “challenging problem”
TF-IDF => “line drawing”

818
Entity Linking with a Knowledge Base: Issues, Techniques, and Solutions
 Wei Shen, Jianyong Wang, **Jiawei Han**
 Knowledge and Data Engineering, IEEE Transactions (2015)
 Bibtext <http://dx.doi.org/10.1109/TKDE.2014.2327028>

817
Power-Based Diagnosis of Node Silence in Remote High-End Sensing Systems
 Yong Yang, Lu Su, Mohammad Maifi Hasan Khan, Michael LeMay, Tarek F. Abdelzaher, **Jiawei Han**
 TOSN (2015)
 Bibtext <http://doi.acm.org/10.1145/2661639>

816
A Framework of Mining Trajectories from Untrustworthy Data in Cyber-P...
 Lu An Tang, Yiqiao Yu, Quanquan Gu, **Jiawei Han**, Guofei Jiang, Alice Leung, Thomas F. Le...

Knowledge-Driven Approach

List of publications



Infer

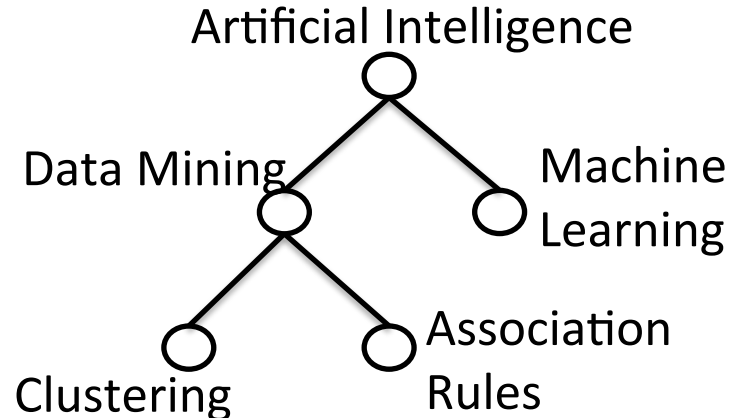


Jiawei Han (韩家炜) ✓
H-Index: 133 | **#Paper:** 818 | **#Citation:** 111164
 Department of Computer Science, University of Illinois at Urbana-Champaign
 Professor

Similar

Data Mining | Information Extraction | Data Analysis | Machine Learning

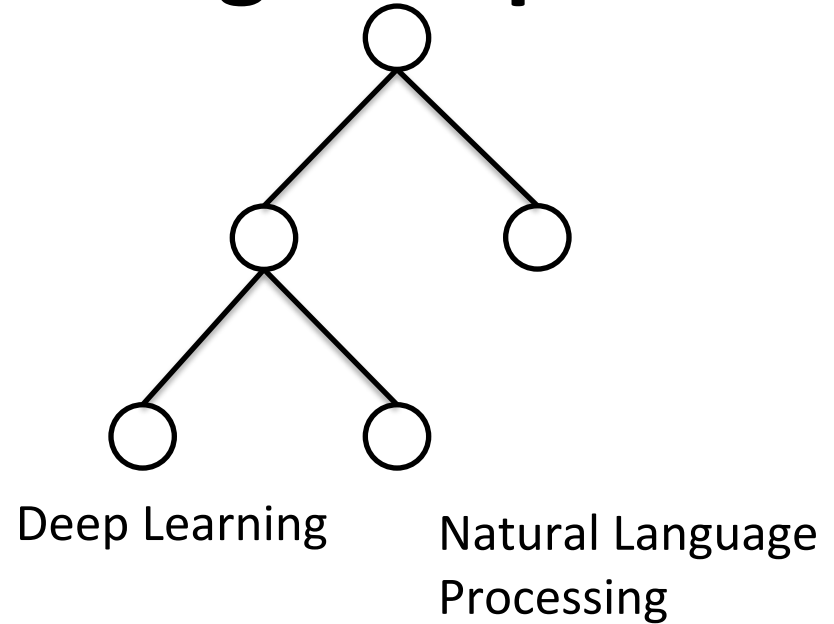
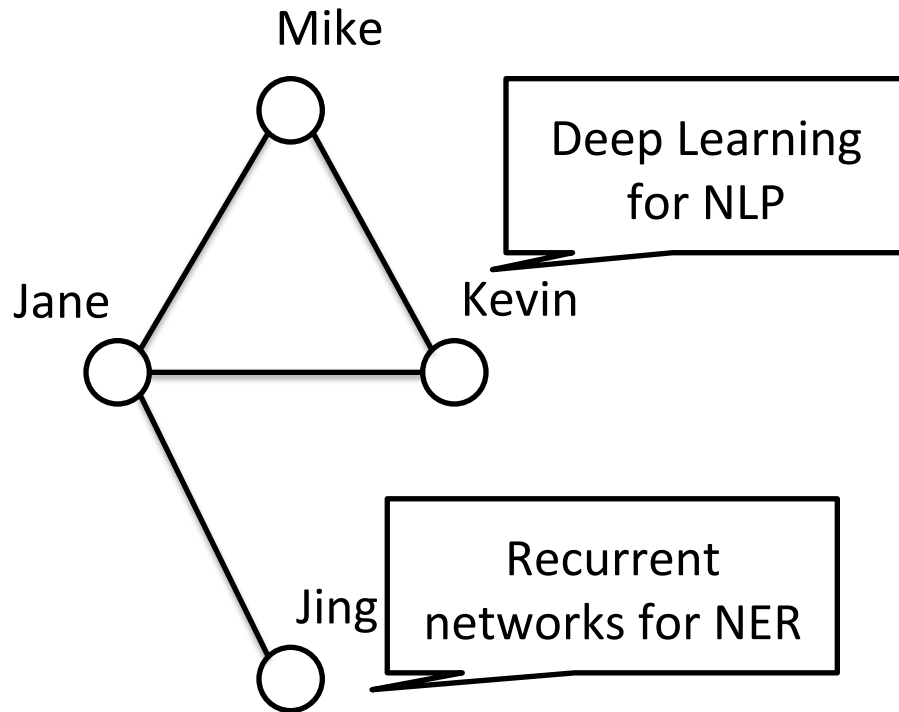
Research interests



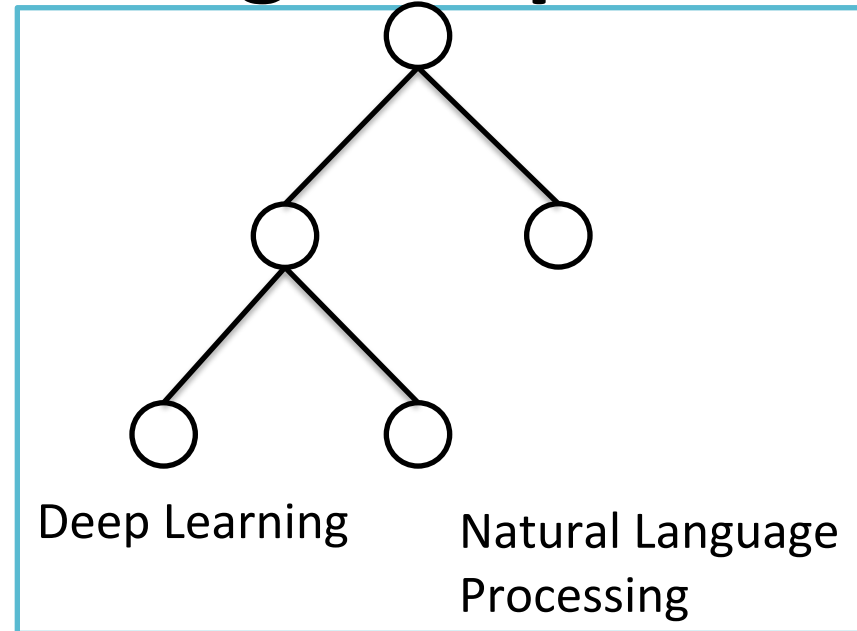
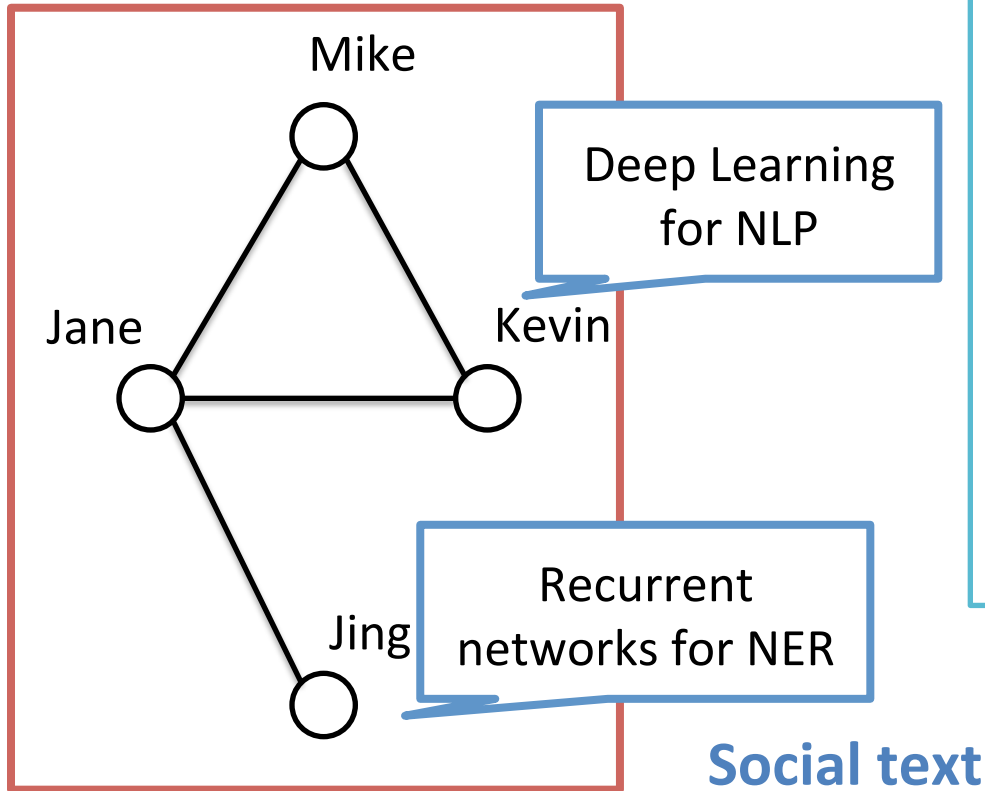
Knowledge bases

Problem:

Learning Social Knowledge Graphs



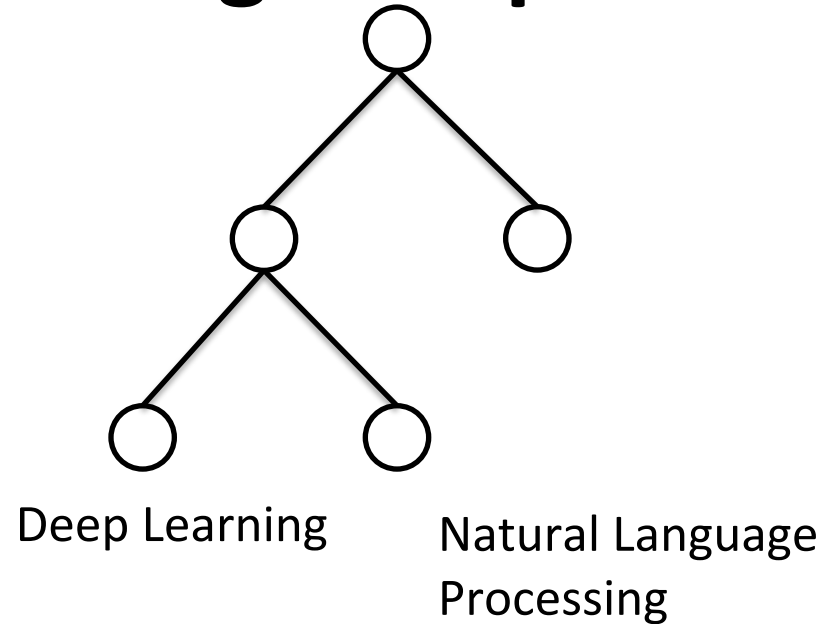
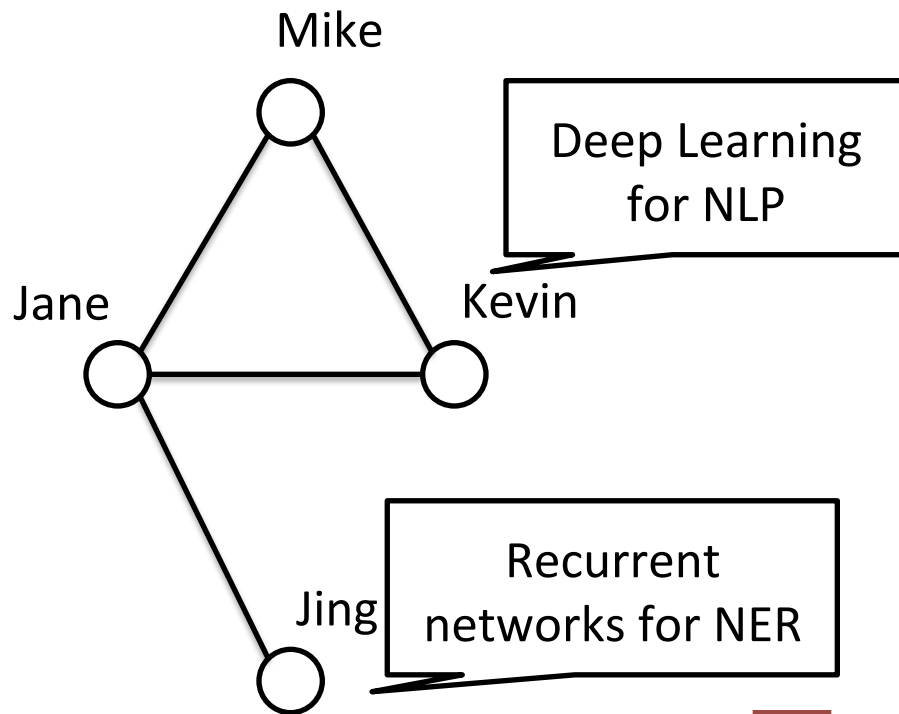
Problem: Learning Social Knowledge Graphs



Knowledge base

Social network structure

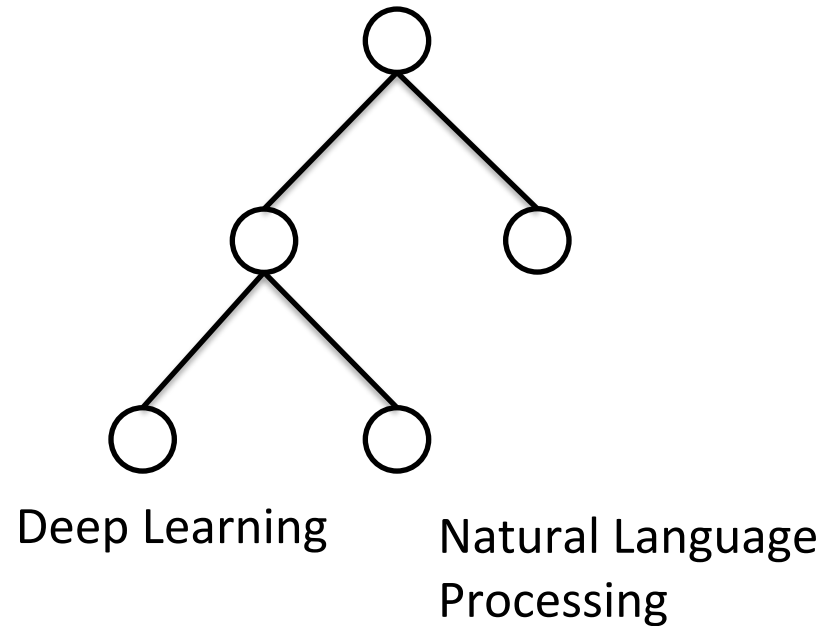
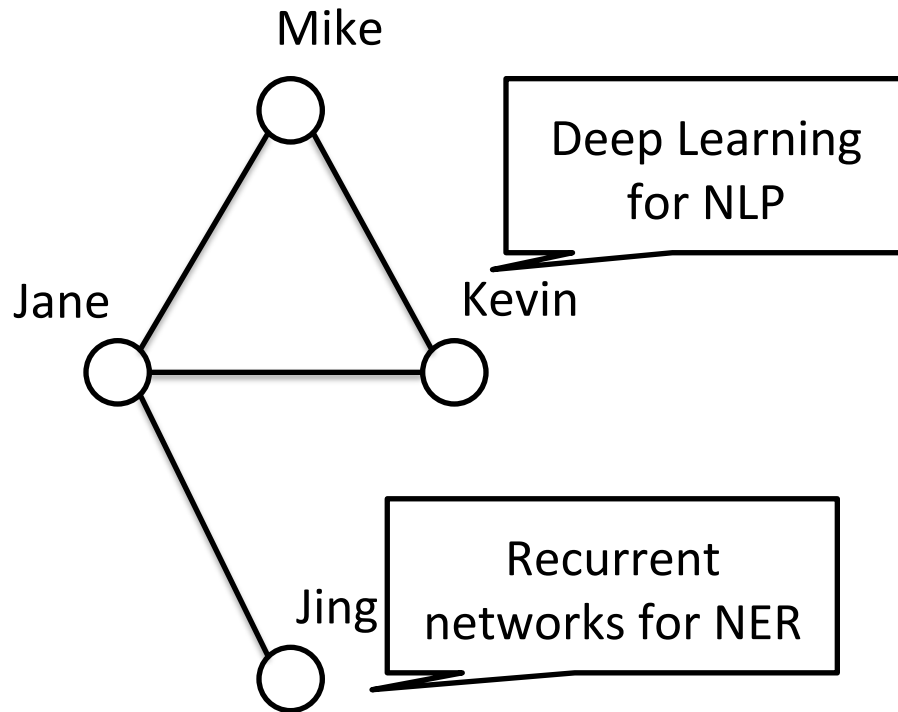
Problem: Learning Social Knowledge Graphs



Infer a ranked list of concepts

Kevin: Deep Learning, Natural Language Processing
Jing: Recurrent Networks, Named Entity Recognition

Challenges

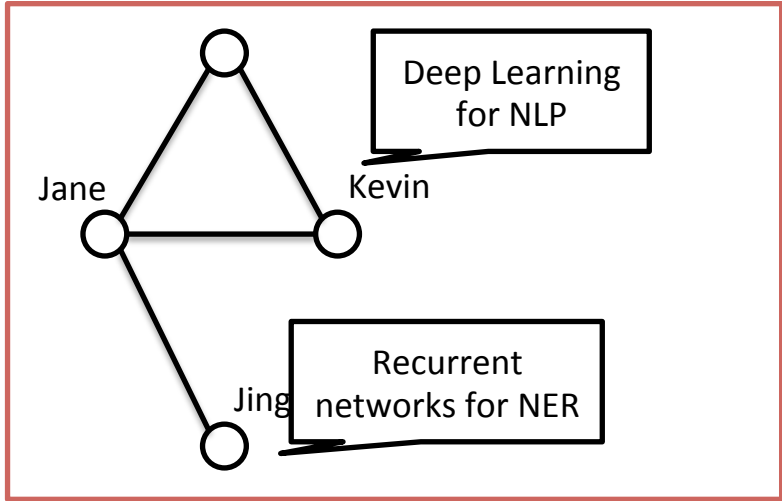


Two modalities – users and concepts

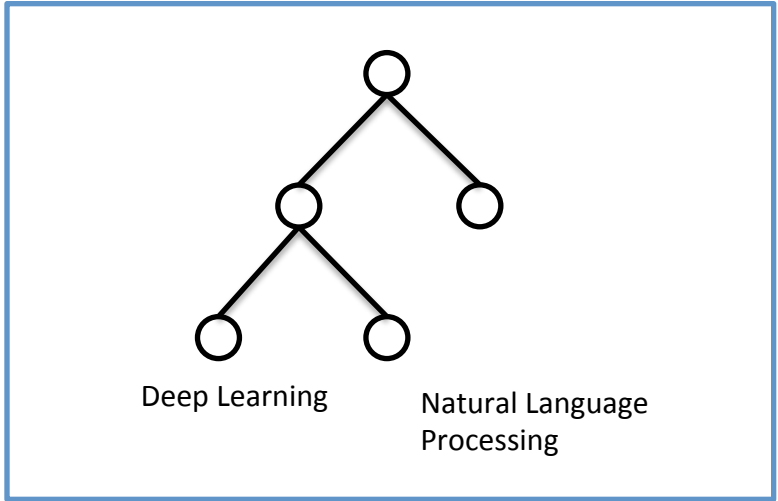
How to leverage information from both modalities?

How to connect these two modalities?

Approach



Learn user embeddings

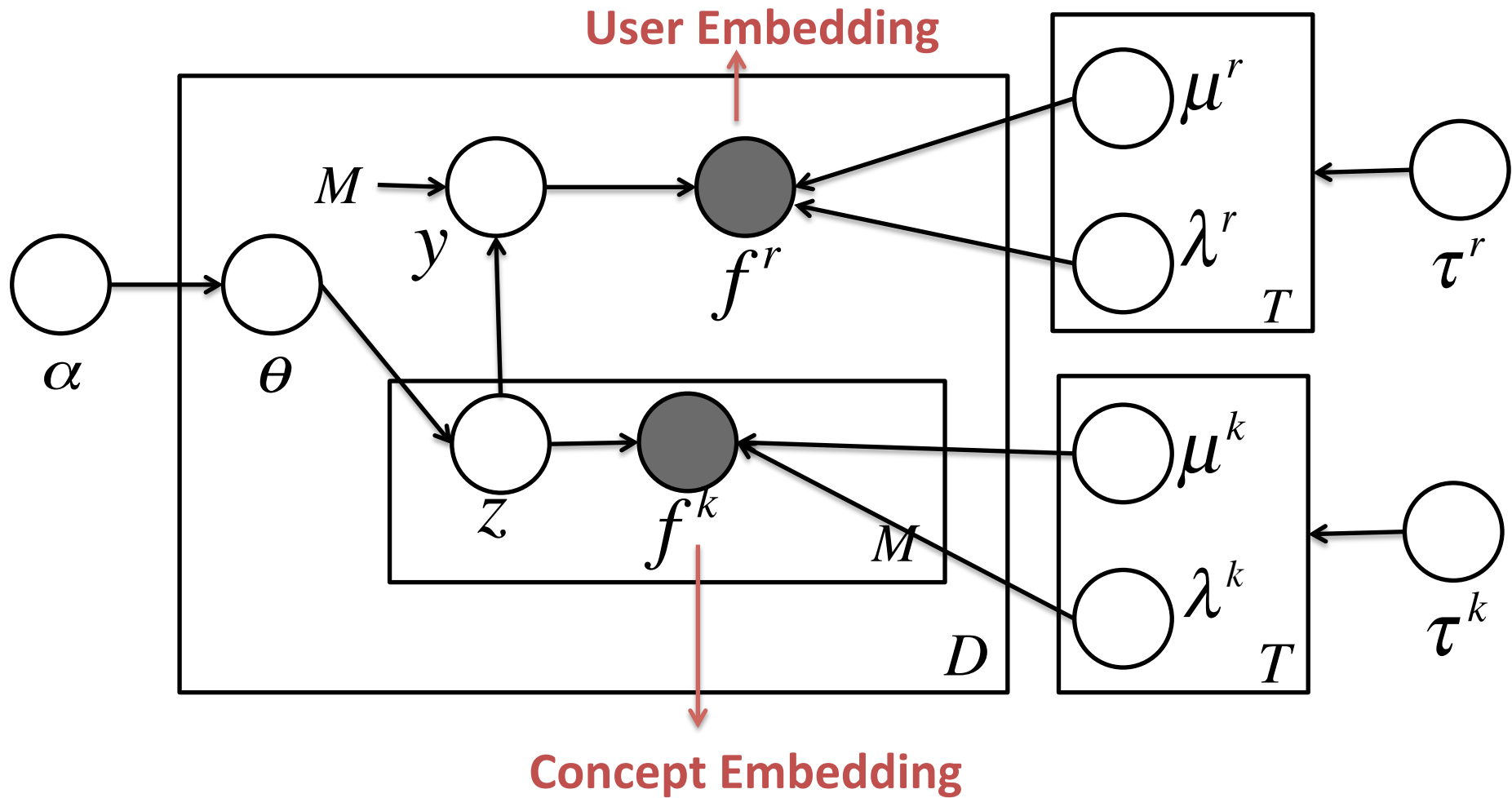


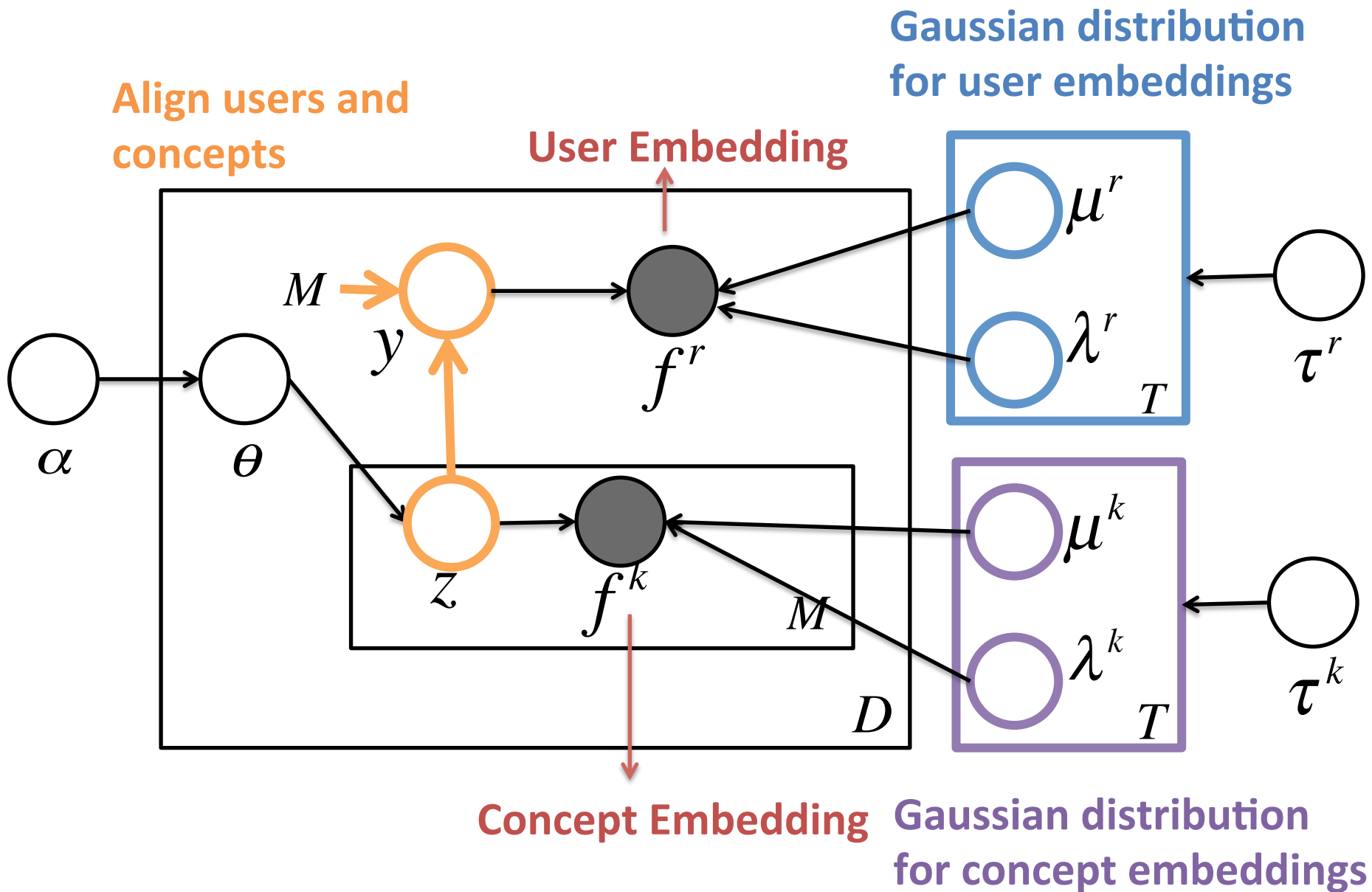
Learn concept embeddings

Model

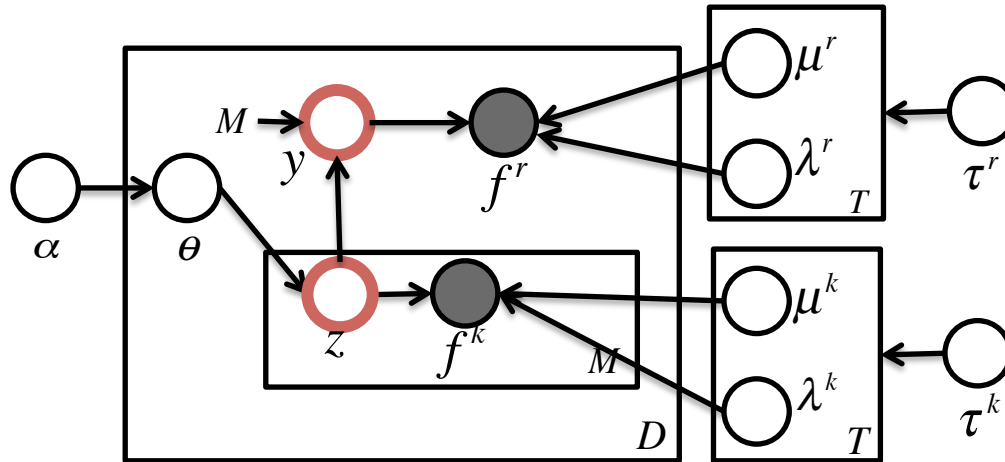
Social KG

Model





Inference and Learning

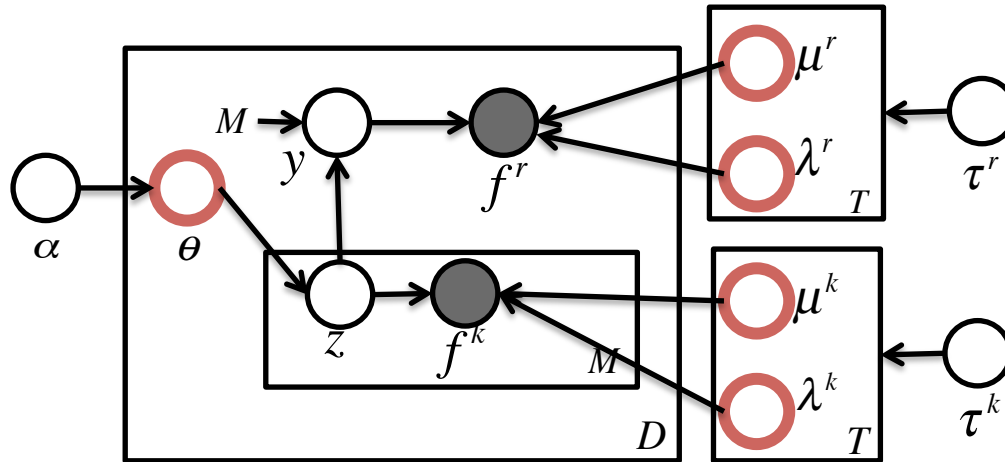


Collapsed Gibbs sampling

Iterate between:

1. **Sample latent variables**

Inference and Learning

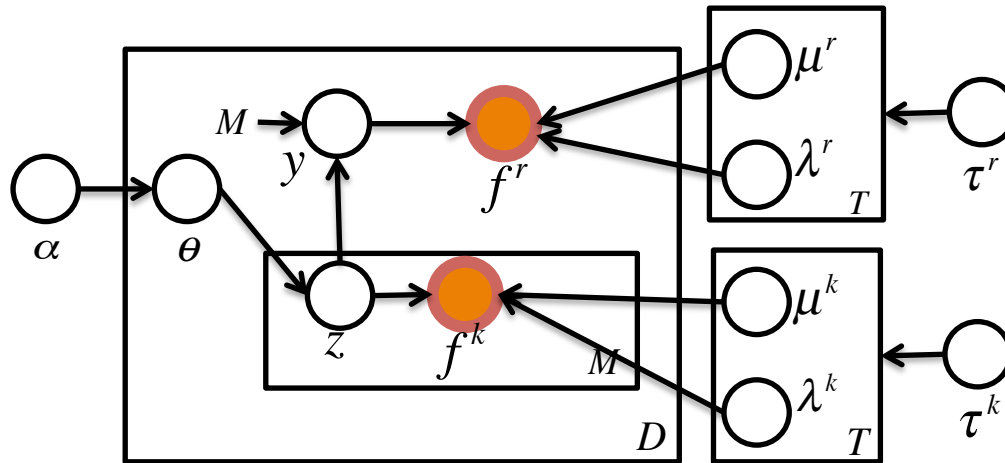


Collapsed Gibbs sampling

Iterate between:

1. Sample latent variables
2. **Update parameters**

Inference and Learning



Collapsed Gibbs sampling

Iterate between:

1. Sample latent variables
2. Update parameters
3. **Update embeddings**

AMiner Research Interest Dataset

- 644,985 researchers
- Terms in these researchers' publications
 - Filtered with Wikipedia
- Evaluation
 - Homepage matching
 - 1,874 researchers
 - Using homepages as ground truth
 - LinkedIn matching
 - 113 researchers
 - Using LinkedIn skills as ground truth

Code and data available:

<https://github.com/kimiyoung/genvector>

Homepage Matching

Using homepages as ground truth.

Method	Precision@5
GenVector	78.1003%
GenVector-E	77.8548%
Sys-Base	73.8189%
Author-Topic	74.4397%
NTN	65.8911%
CountKG	54.4823%

GenVector

Our model

GenVector-E

Our model w/o embedding update

Sys-Base

AMiner baseline: key term extraction

CountKG

Rank by frequency

Author-topic

Classic topic models

NTN

Neural tensor networks

LinkedIn Matching

Using LinkedIn skills as ground truth.

Method	Precision@5
GenVector	50.4424%
GenVector-E	49.9145%
Author-Topic	47.6106%
NTN	42.0512%
CountKG	46.8376%

GenVector Our model
GenVector-E Our model w/o embedding update

CountKG Rank by frequency
Author-topic Classic topic models
NTN Neural tensor networks

Error Rate of Irrelevant Cases

Manually label terms that are clearly NOT research interests, e.g., challenging problem.

Method	Error rate
GenVector	1.2%
Sys-Base	18.8%
Author-Topic	1.6%
NTN	7.2%

GenVector
Sys-Base

Our model
AMiner baseline: key term extraction

Author-topic
NTN

Classic topic models
Neural tensor networks

Qualitative Study: Top Concepts within Topics

GenVector

Query expansion
Concept mining
Language modeling
Information extraction
Knowledge extraction
Entity linking
Language models
Named entity recognition
Document clustering
Latent semantic indexing

Author-Topic

Speech recognition
Natural language
***Integrated circuits**
Document retrieval
Language models
Language model
***Microphone array**
Computational linguistics
***Semidefinite programming**
Active learning

Qualitative Study: Top Concepts within Topics

GenVector

Image processing
Face recognition
Feature extraction
Computer vision
Image segmentation
Image analysis
Feature detection
Digital image processing
Machine learning algorithms
Machine vision

Author-Topic

Face recognition
***Food intake**
Face detection
Image recognition
***Atmospheric chemistry**
Feature extraction
Statistical learning
Discriminant analysis
Object tracking
***Human factors**

Qualitative Study: Research Interests

GenVector

Feature extraction
Image segmentation
Image matching
Image classification
Face recognition

Sys-Base

Face recognition
Face image
***Novel approach**
***Line drawing**
Discriminant analysis

Qualitative Study: Research Interests

GenVector

Unsupervised learning
Feature learning
Bayesian networks
Reinforcement learning
Dimensionality reduction

Sys-Base

***Challenging problem**
Reinforcement learning
***Autonomous helicopter**
***Autonomous helicopter flight**
Near-optimal planning

Online Test

Does Jiawei Han have these skills or expertise? ✕

Information Extraction ✕ Efficient Algorithm ✕ Data Mining ✕ Large Databases ✕ Data Cube ✕

Data Integration ✕ Type another area of expertise

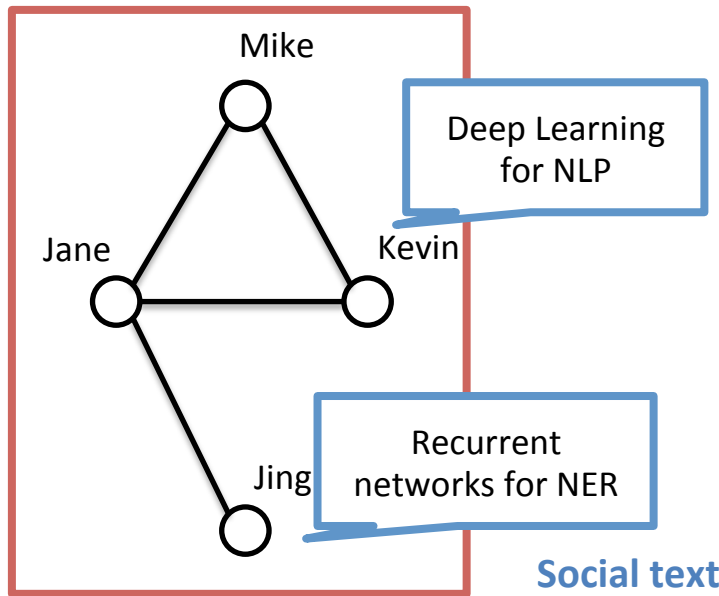
Skip Vote

A/B test with live users

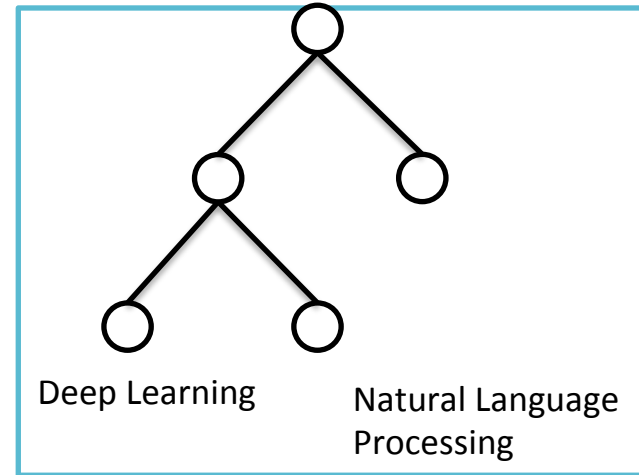
- Mixing the results with Sys-Base

Method	Error rate
GenVector	3.33%
Sys-Base	10.00%

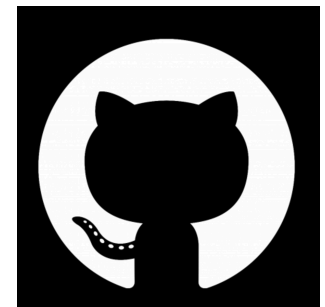
Other Social Networks?



Social network structure



Knowledge base



Conclusion

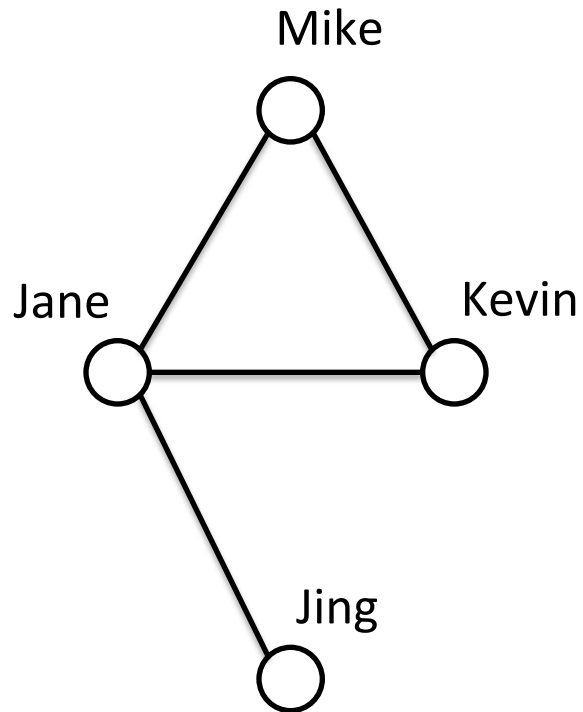
- Study a novel problem
 - Learning social knowledge graphs
- Propose a model
 - Multi-modal Bayesian embedding
 - Integrate embeddings into graphical models
- AMiner research interest dataset
 - 644,985 researchers
 - Homepage and LinkedIn matching as ground truth
- Online deployment on AMiner

Thanks!

Code and data:

<https://github.com/kimiyoung/genvector>

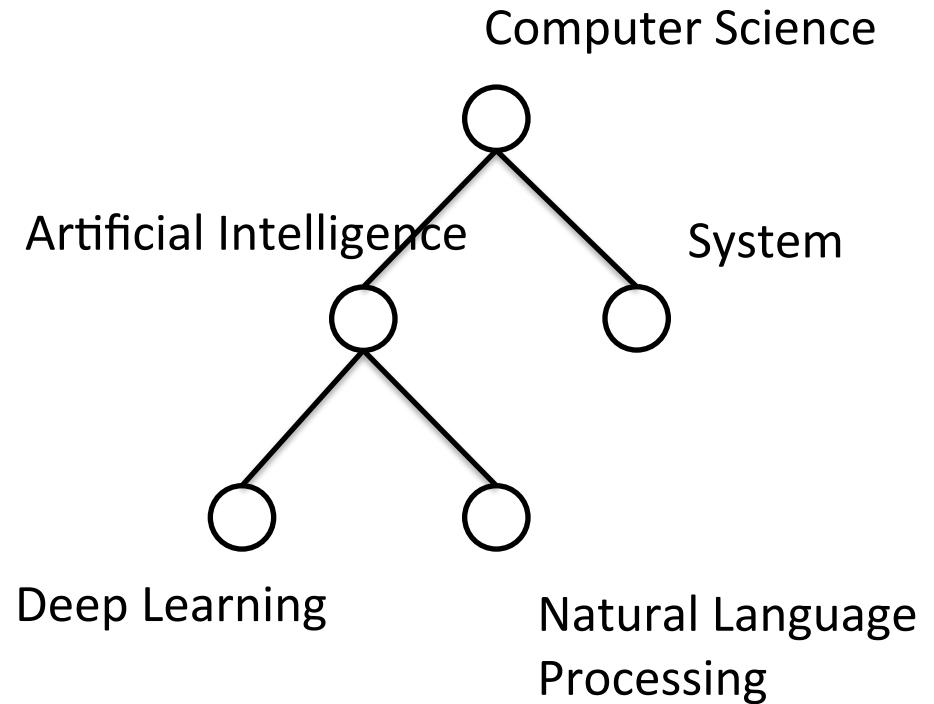
Social Networks



AMiner, Facebook, Twitter...

Huge amounts of information

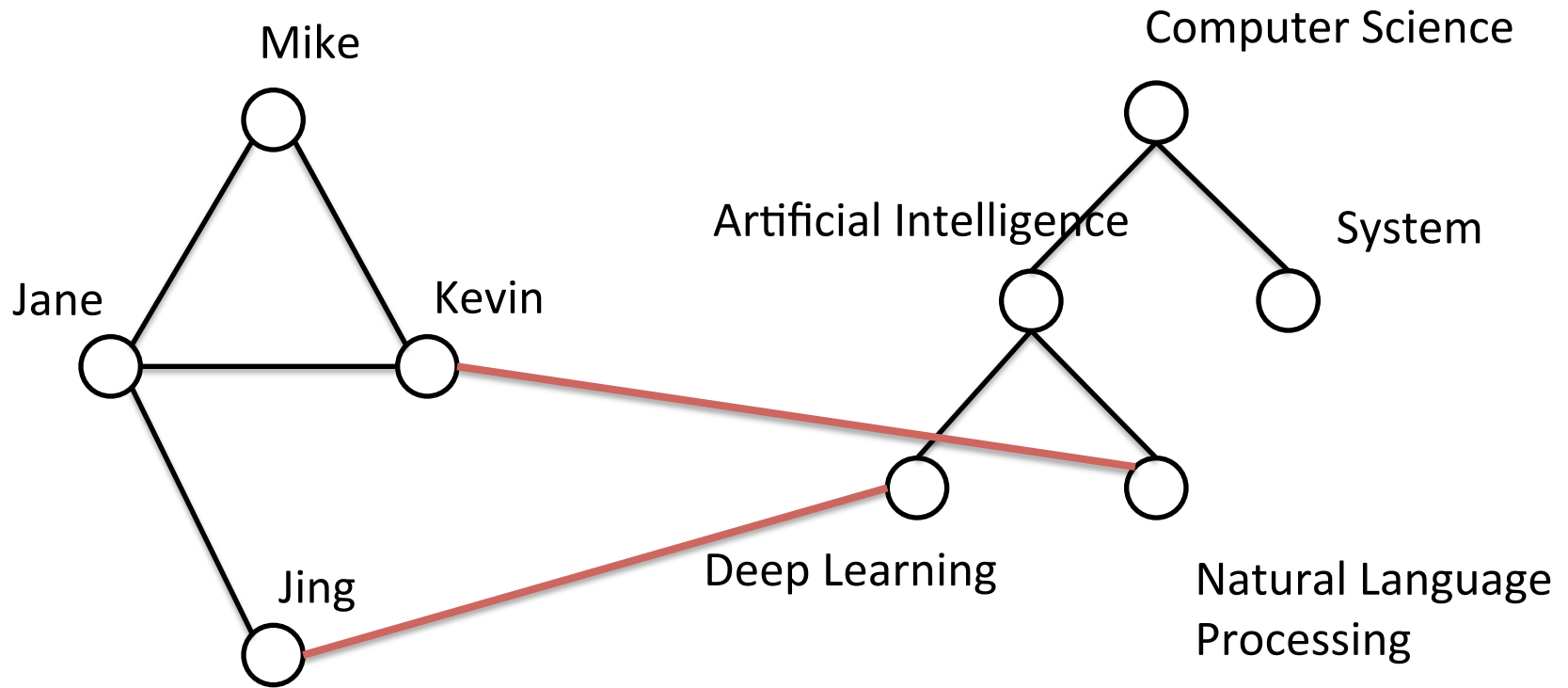
Knowledge Bases



Wikipedia, Freebase, Yago, NELL...

Huge amounts of knowledge

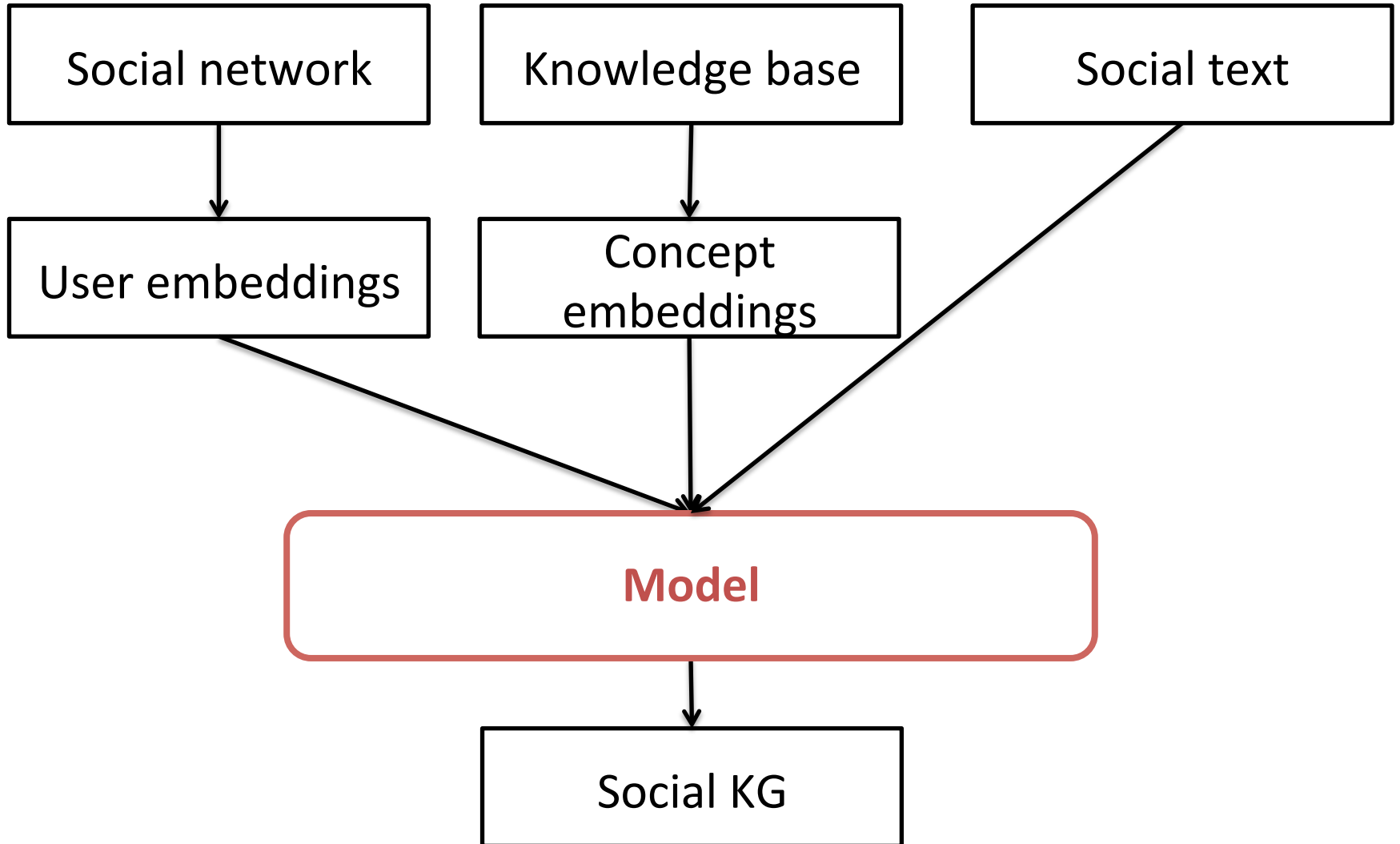
Bridge the Gap



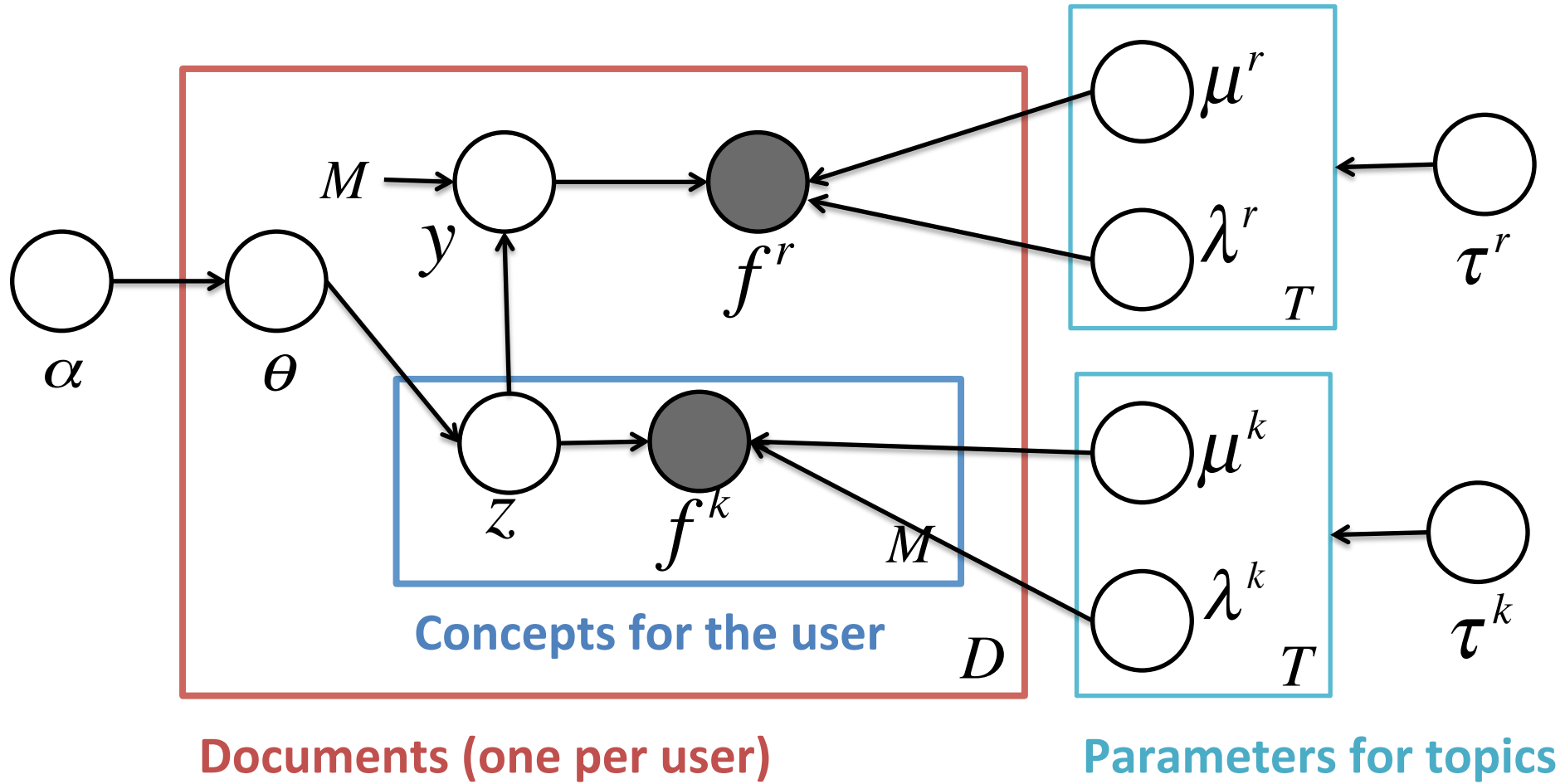
Better user understanding
e.g. mine research interests on AMiner

Approach

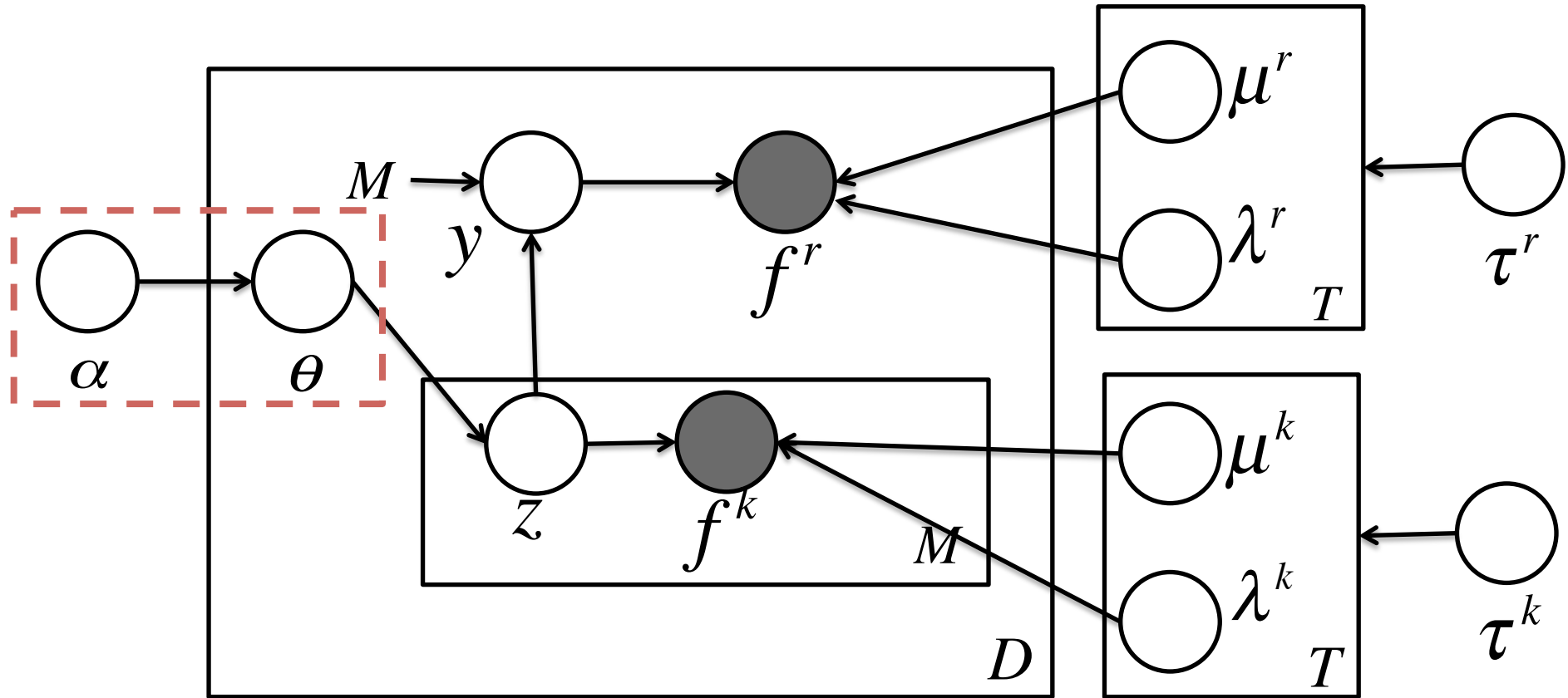
Copy picture



Model

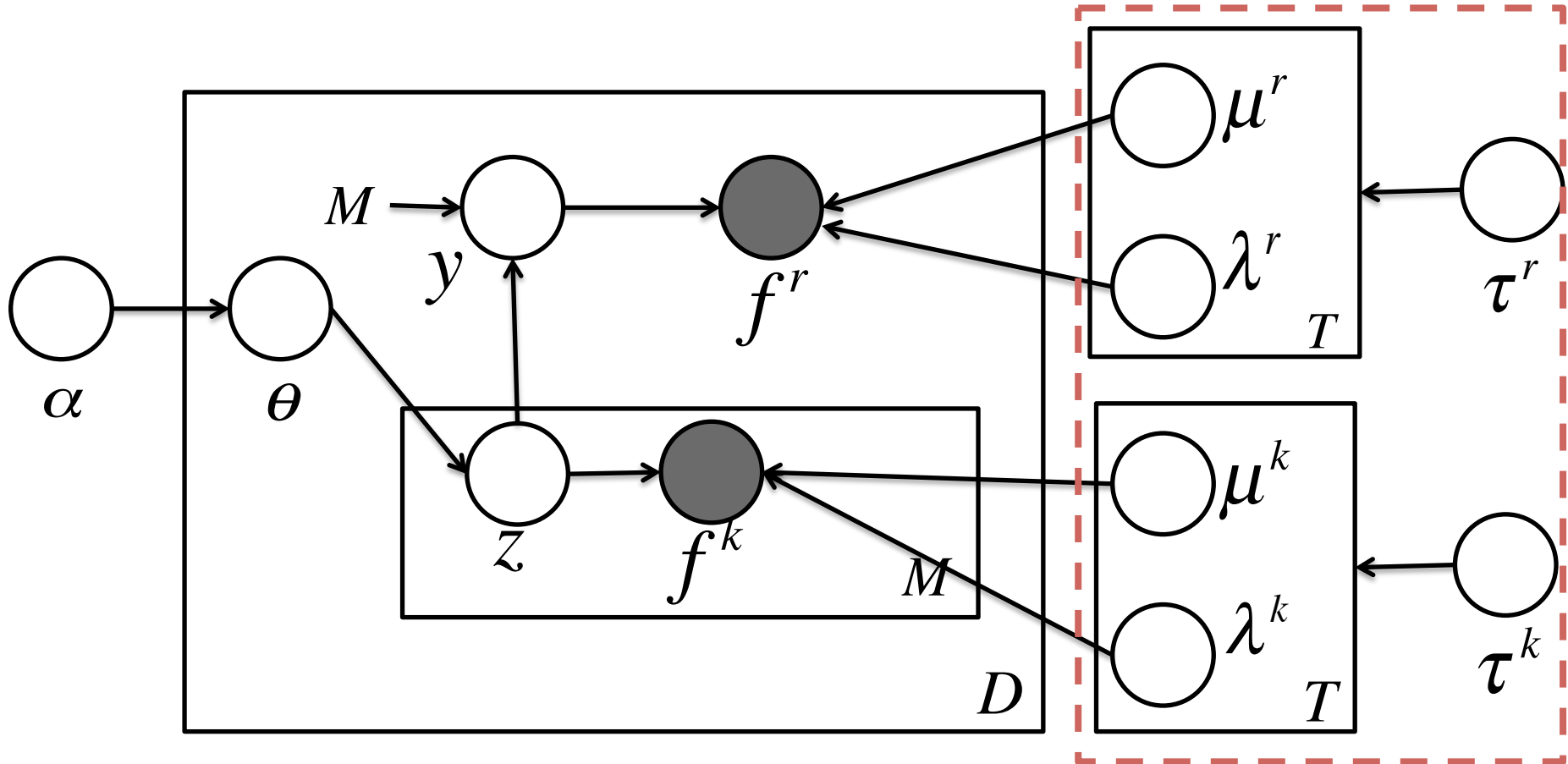


Model



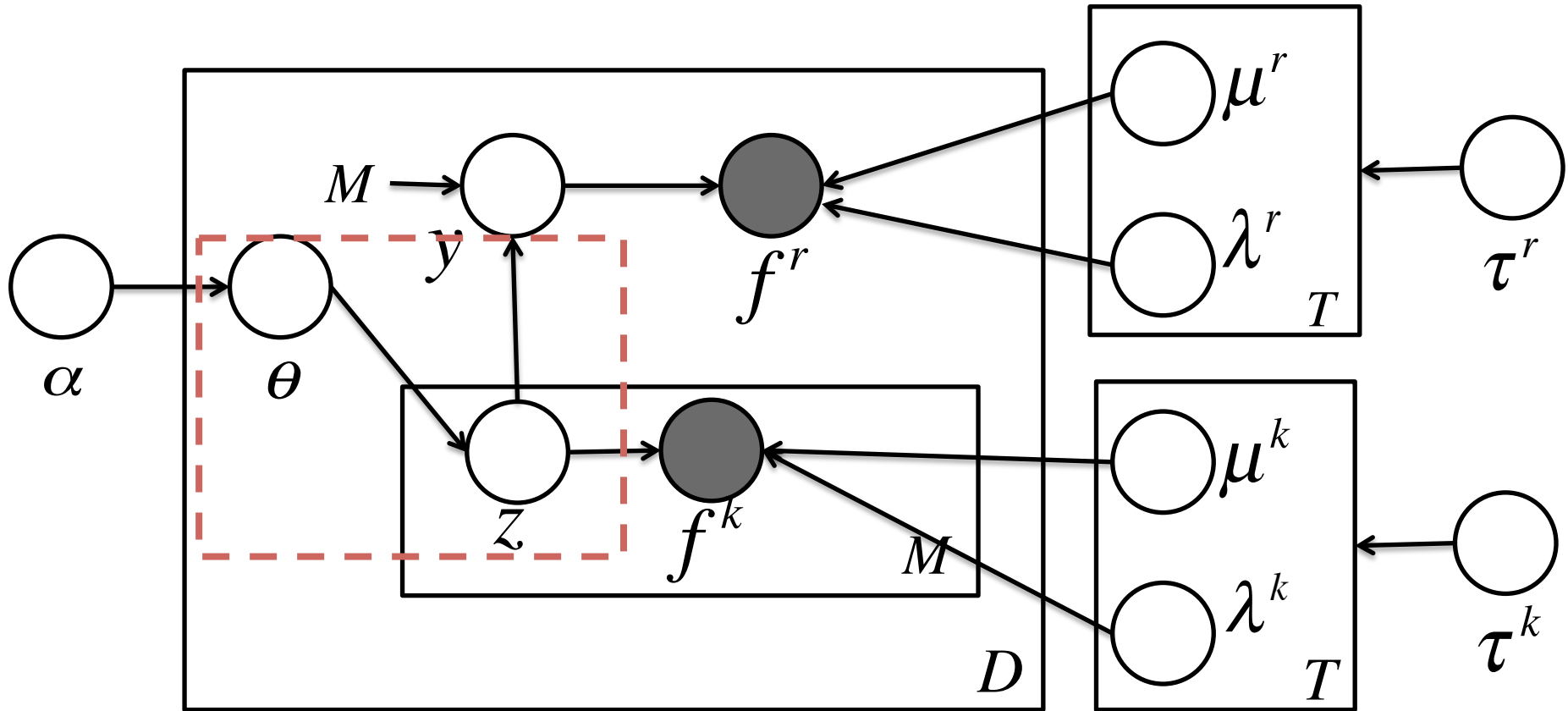
Generate a topic distribution for each document (from a Dirichlet)

Model



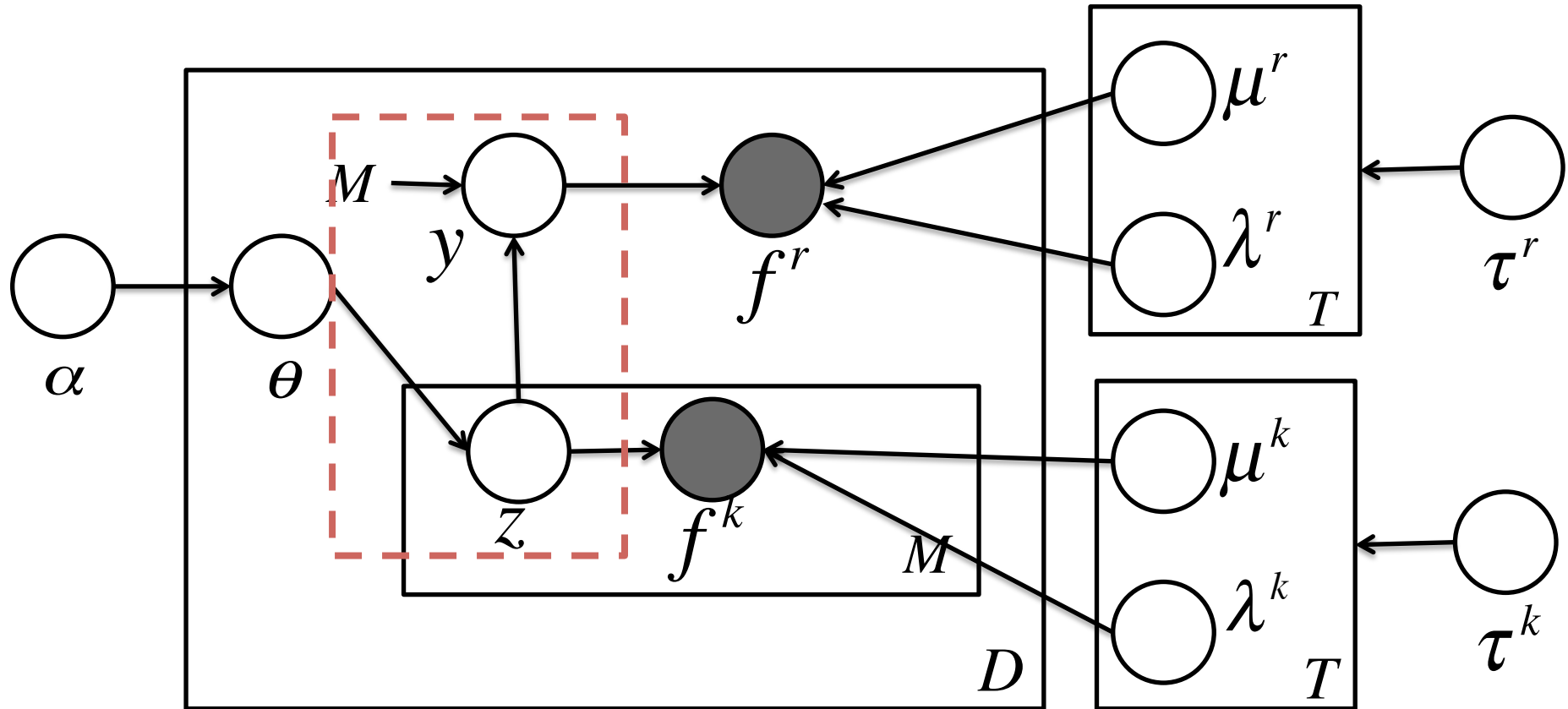
Generate Gaussian distribution for each embedding space (from a Normal Gamma)

Model



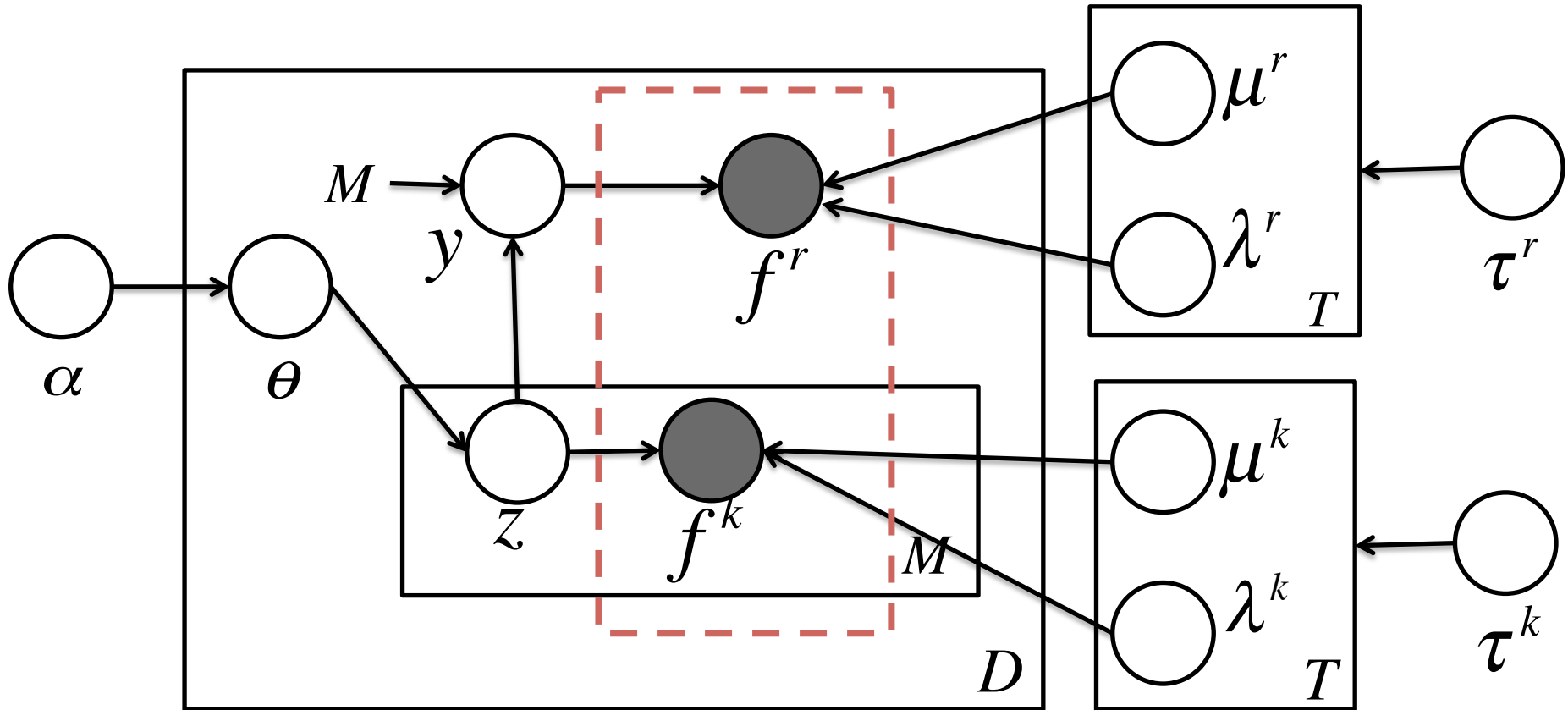
Generate the topic for each concept (from a Multinomial)

Model



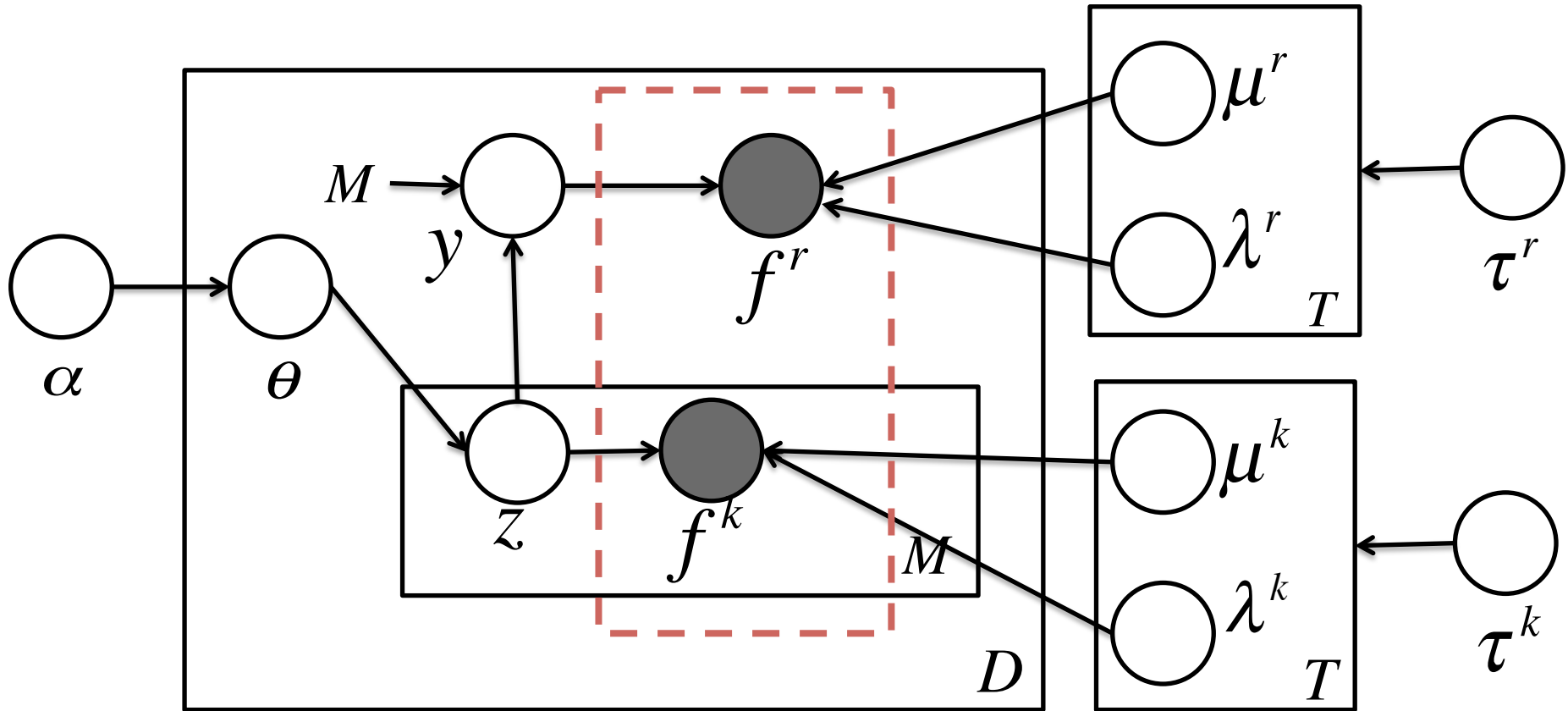
Generate the topic for each user (from a Uniform)

Model



Generate embeddings for users and concepts (from a Gaussian)

Model



General

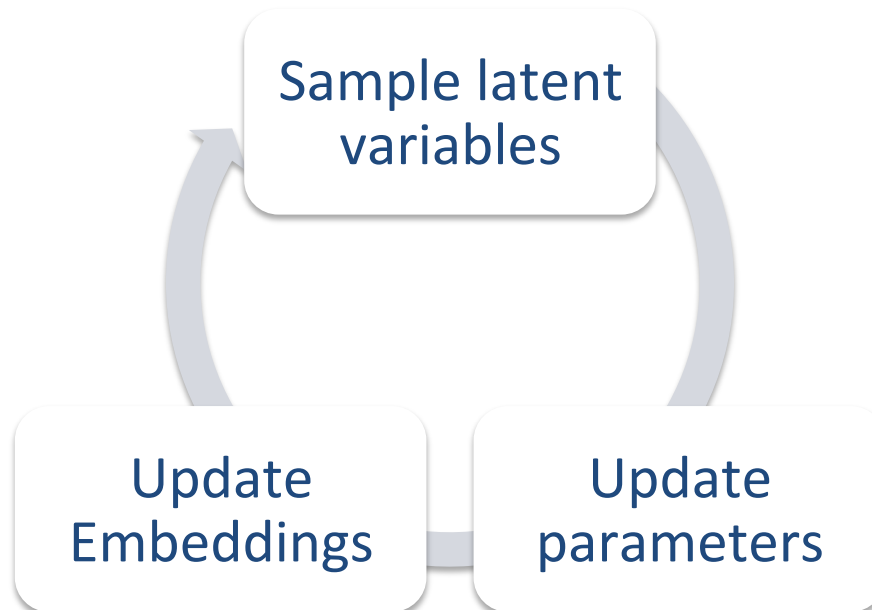
Inference and Learning

Collapsed Gibbs sampling for inference

Add picture

Update the embedding during learning

Different from LDAs with discrete observed variables



Methods for Comparison

Method	Description
GenVector	Our model
GenVector-E	Our model w/o embedding update
Sys-Base	AMiner baseline: key term extraction
CountKG	Rank by frequency
Author-topic	Classic topic models
NTN	Neural tensor networks