

HyperMiner: Topic Taxonomy Mining with Hyperbolic Embedding

Yishi Xu¹ Dongsheng Wang¹ Bo Chen¹ Ruiying Lu¹ Zhibin Duan¹ Mingyuan Zhou²

¹ Xidian University

² The University of Texas at Austin

Source code: <https://github.com/NoviceStone/HyperMiner>

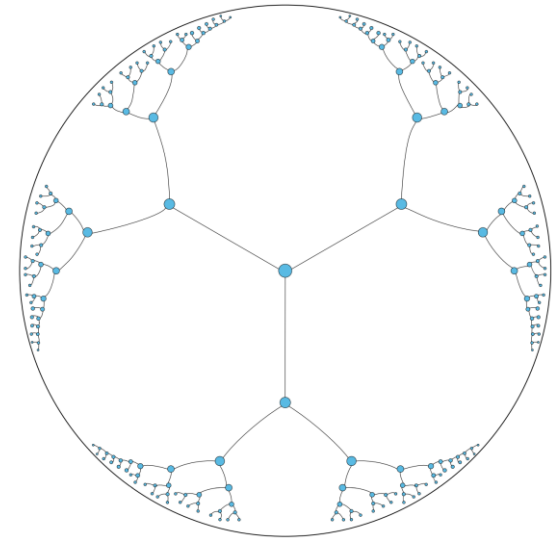
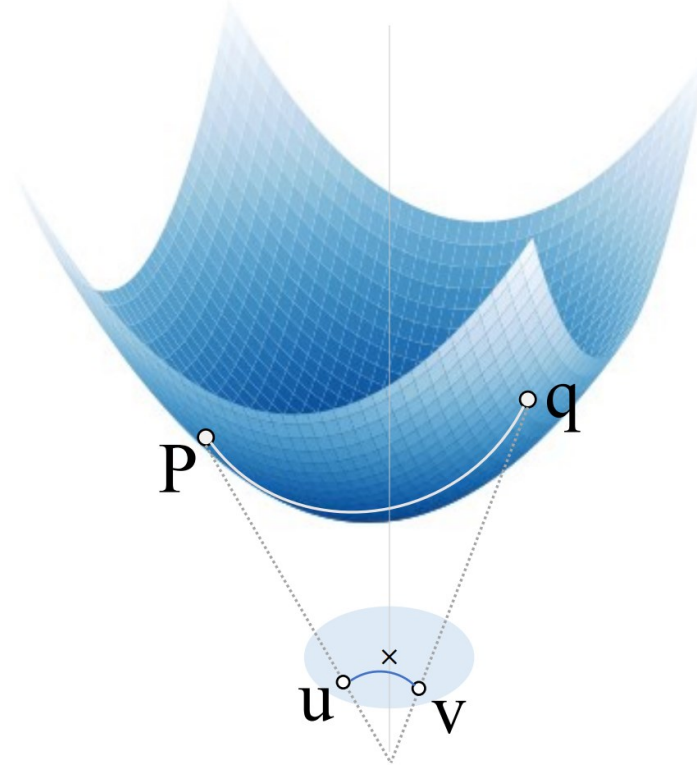
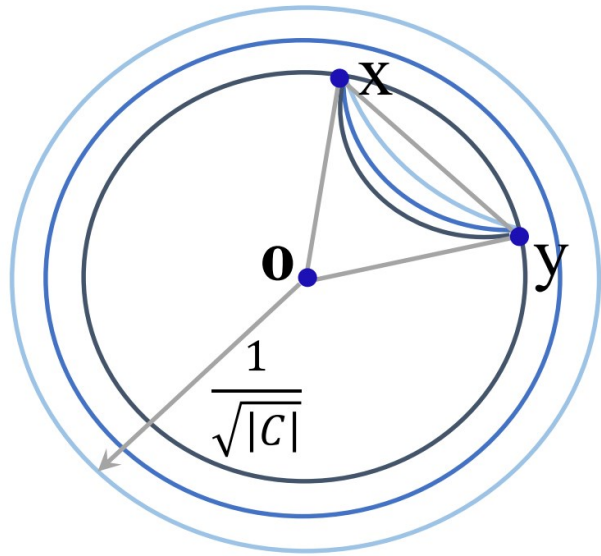
Motivation

Existing embedded topic models generally hold the Euclidean embedding space assumption, leading to a fundamental limitation in capturing hierarchical relationships given that

- The lexical hierarchy naturally exists for the words of vocabulary
- A semantic hierarchy is also expected between topics and words

Motivation

Hyperbolic geometry has shown superior performance in modeling hierarchical data, with the tree-likeness property of its distance metric.




Approach

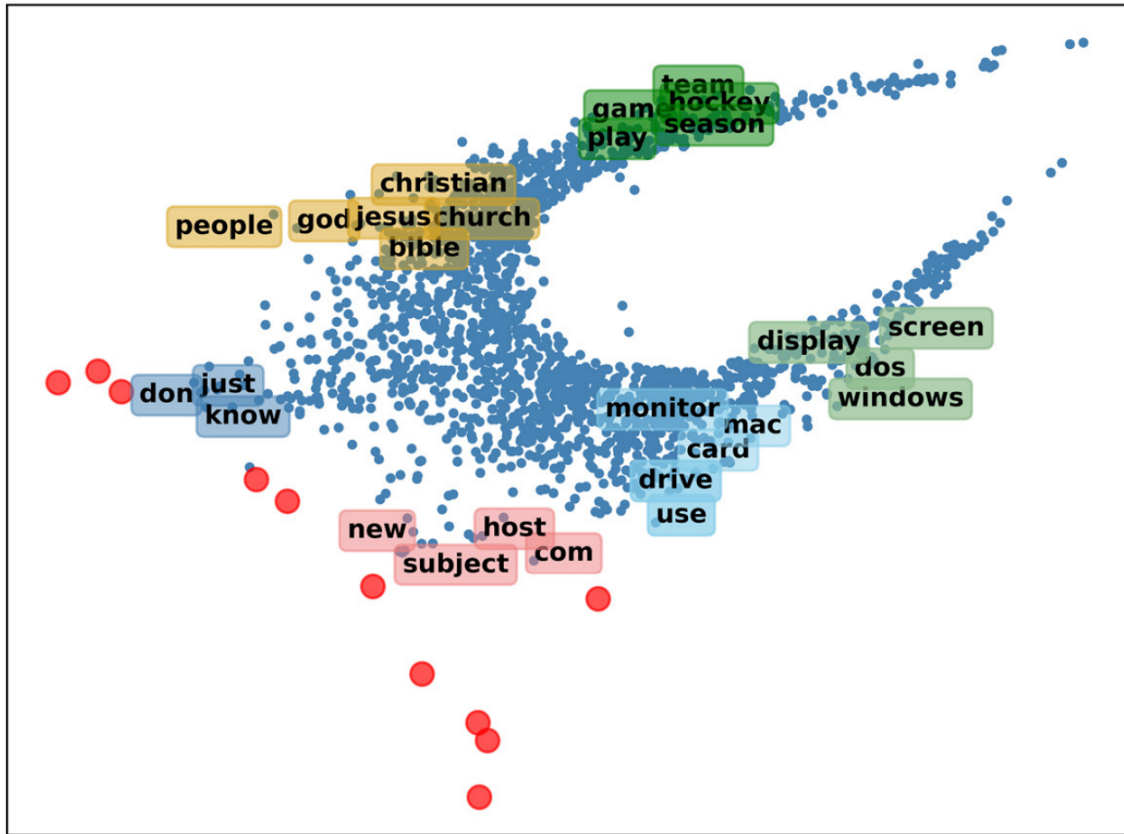
Embed both words and topics into a shared hyperbolic space instead of Euclidean space, so that the hyperbolic distance metric can be used to measure the semantic similarity between topic and words. Note that

- $\beta_k \in \mathbb{R}^V$: distribution of words for the k-th topic
- $\alpha_k \in \mathbb{R}^D$: vector representation of the k-th topic
- $\rho \in \mathbb{R}^{D \times V}$: word embeddings of the vocabulary

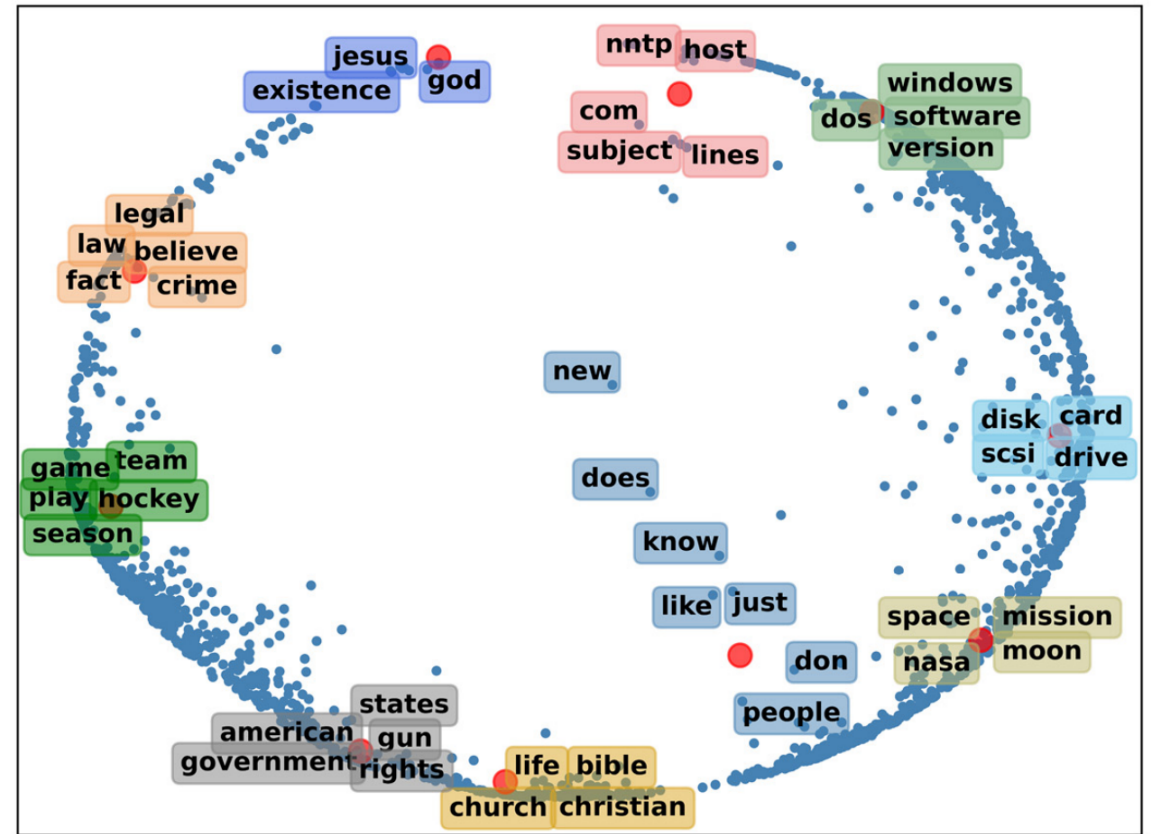
$\beta_k = \text{Softmax}(\text{dist}(\rho, \alpha_k))$ Inner product for Euclidean embeddings, replace it with distance metric for hyperbolic embeddings.



Implicit semantic hierarchy mining

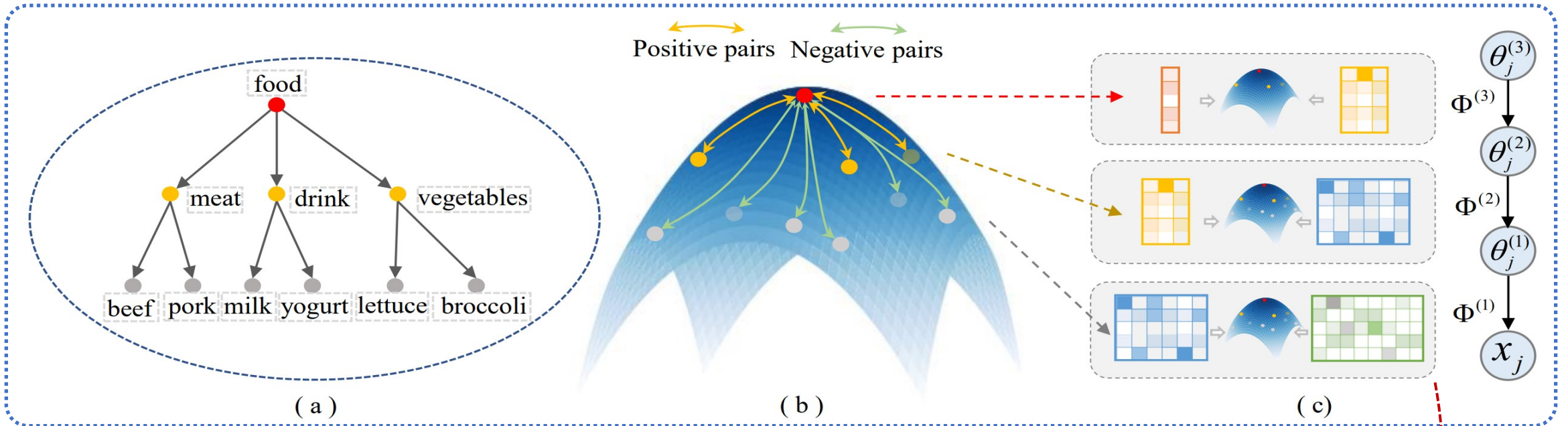


2D Euclidean embedding space
learned by ETM



2D hyperbolic embedding space
learned by HyperETM

Knowledge-guided topic modeling



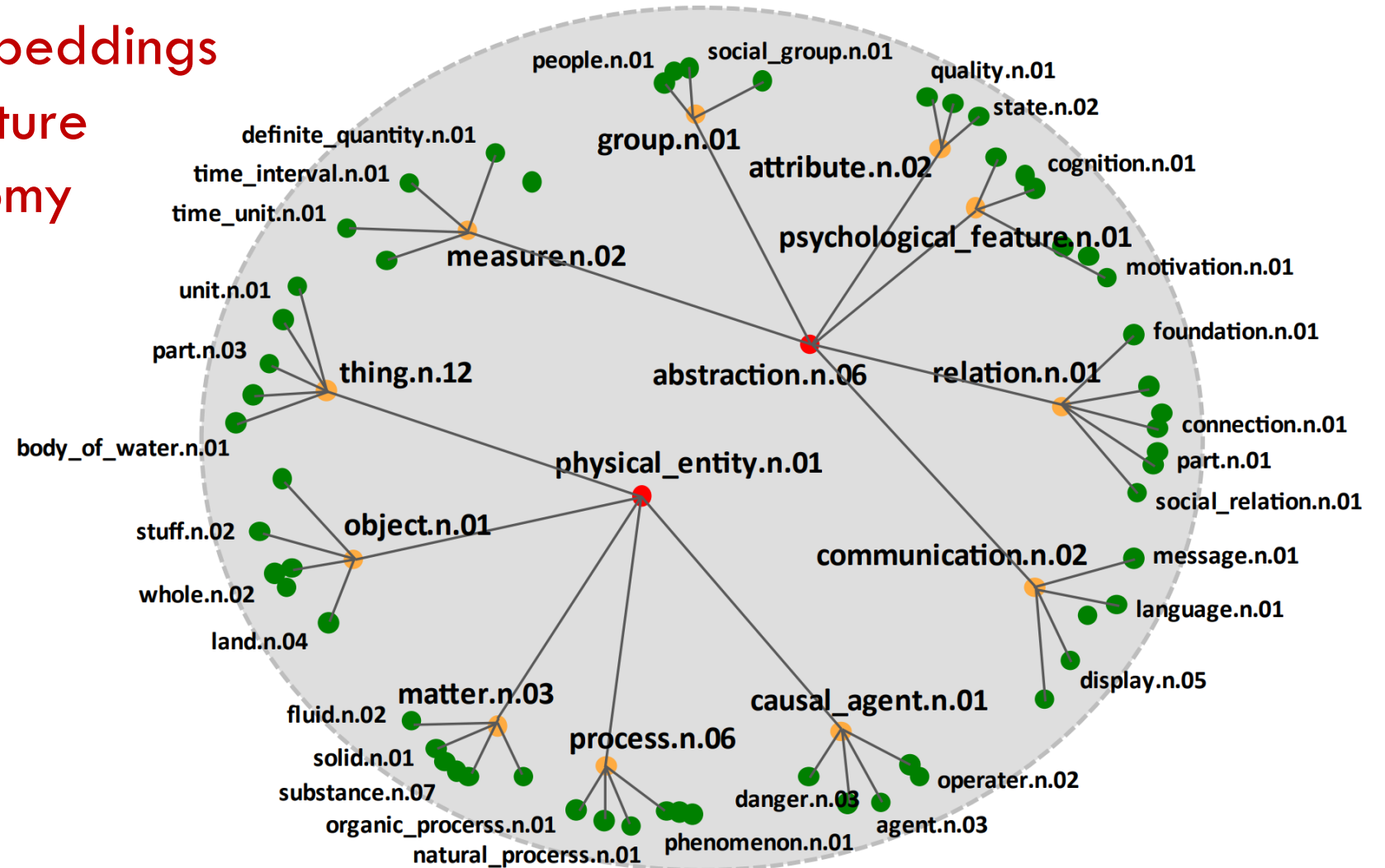
Injecting tree-structure knowledge via contrastive loss

$$\mathcal{L}_{\text{Contra}} = \mathbb{E}_{\alpha_i \in \mathcal{T}} \left[-\log \frac{\exp(\text{dist}(\alpha_i, \alpha_i^+)/\tau)}{\exp(\text{dist}(\alpha_i, \alpha_i^+)/\tau) + \sum_{\alpha_i^- \in \mathcal{Q}(\alpha_i)} \exp(\text{dist}(\alpha_i, \alpha_i^-)/\tau)} \right]$$

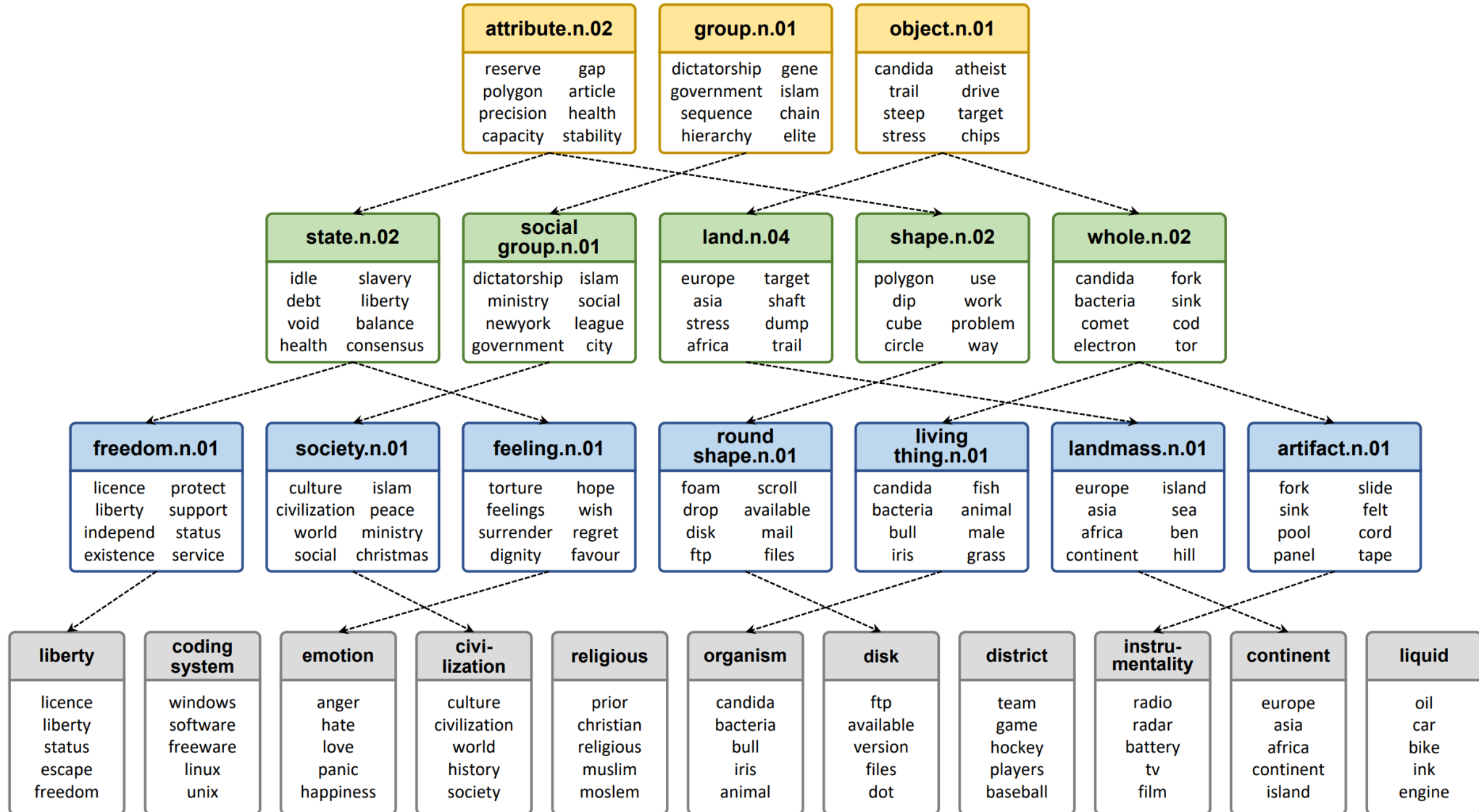
$$\Phi^{(l)} = \text{Softmax}(\text{dist}(\alpha^{(l-1)}, \alpha^{(l)})) \quad \alpha^{(l)} : \text{topic embeddings at layer } l$$

Topic hierarchy visualization

Distribution of topic embeddings well preserves the structure of prior concept taxonomy



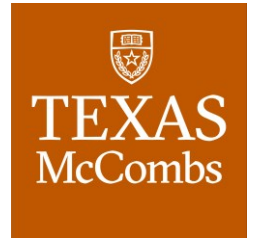
Mined topic taxonomy



Clustering performance

- Our approach yields better document representations

Method	20NG		TMN	
	km-Purity	km-NMI	km-Purity	km-NMI
LDA [7]	38.43 \pm 0.52	35.98 \pm 0.39	48.17 \pm 0.69	30.96 \pm 0.78
ProdLDA [16]	39.21 \pm 0.63	36.52 \pm 0.51	55.28 \pm 0.67	35.57 \pm 0.72
ETM [27]	42.68 \pm 0.71	37.72 \pm 0.64	59.35 \pm 0.74	38.75 \pm 0.86
WHAI [17]	40.89 \pm 0.35	38.90 \pm 0.27	58.06 \pm 0.45	37.34 \pm 0.48
SawETM [28]	43.36 \pm 0.48	41.59 \pm 0.62	61.13 \pm 0.56	40.78 \pm 0.63
TopicNet [31]	42.94 \pm 0.41	40.76 \pm 0.53	60.52 \pm 0.50	40.09 \pm 0.54
HyperETM	43.63 \pm 0.51	39.06 \pm 0.64	61.22 \pm 0.62	40.52 \pm 0.71
HyperMiner	<u>44.37</u> \pm 0.38	<u>42.83</u> \pm 0.45	<u>62.96</u> \pm 0.48	<u>41.93</u> \pm 0.52
HyperMiner-KG	45.16 \pm 0.35	43.65 \pm 0.39	63.84 \pm 0.43	42.81 \pm 0.47



Thank you.

Please feel free to contact us by e-mail

xuyishi@stu.xidian.edu.cn bchen@mail.xidian.edu.cn

mingyuan.zhou@mcombs.utexas.edu

Paper can be downloaded from

