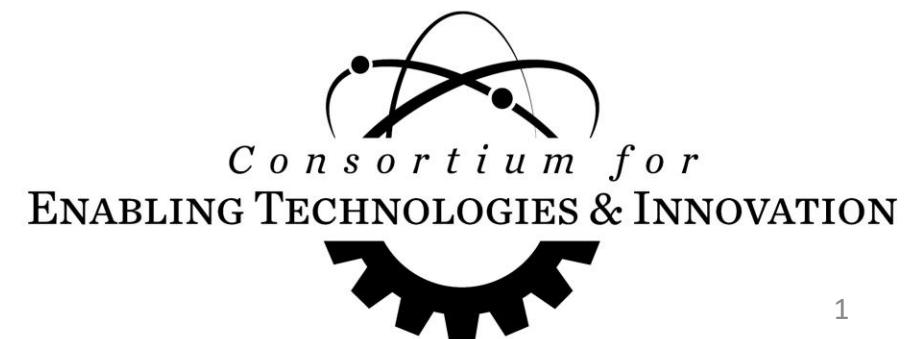


AToMS: Author-Topic Manifold Summarization for Interpretable Author-Collaboration Forecasting

Conrad D. Hougen[§], Karl T. Pazdernik^{*}, Alfred O. Hero[§]

University of Michigan[§], Pacific Northwest National Laboratory^{*}

Date 02/08/23

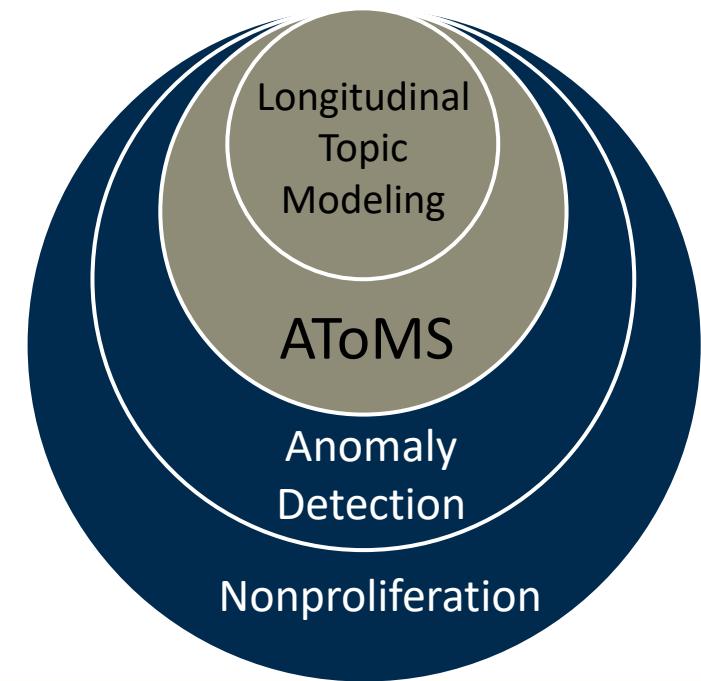




Thrust 1: Computer and Engineering Sciences for Nonproliferation

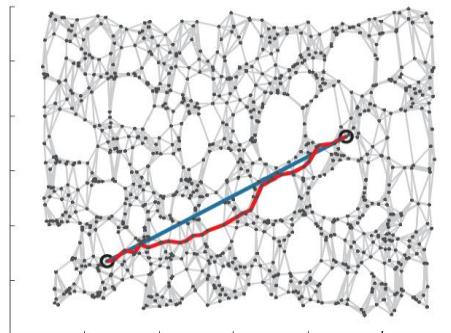
Thrust 2

Thrust 3

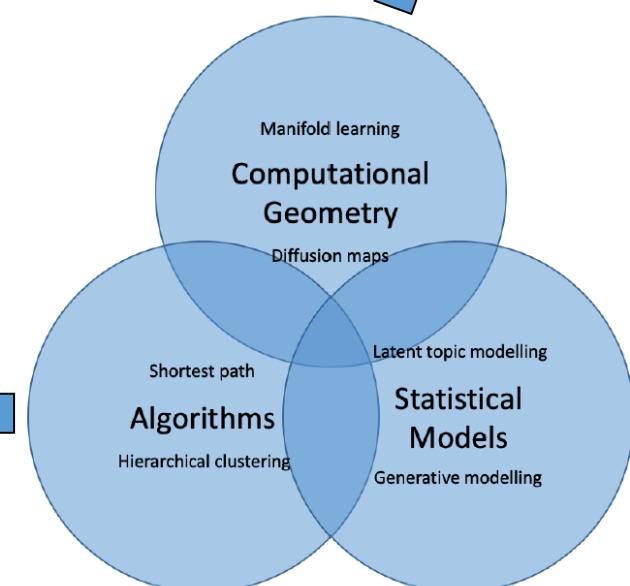
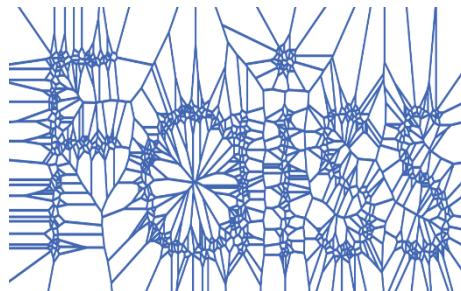


>> Foundations

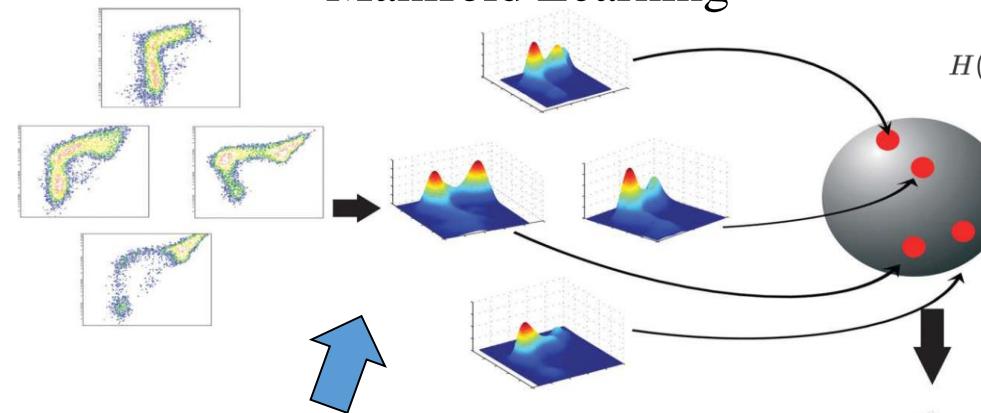
Dijkstra's Algorithm



Similarity Search

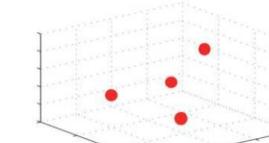


Manifold Learning

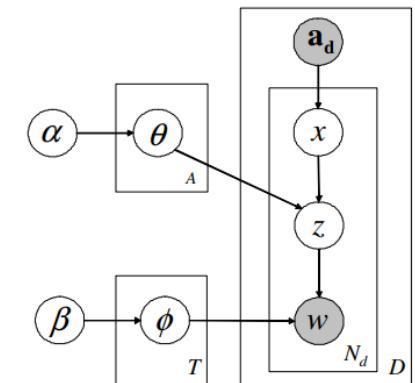
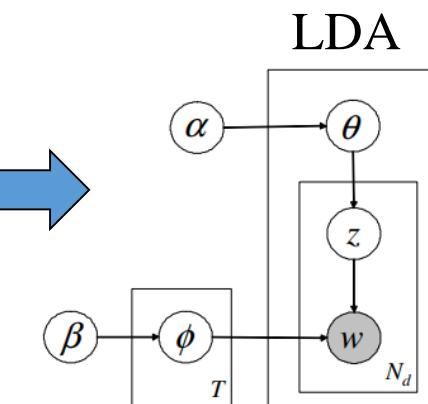


$$H(P, Q) = \frac{1}{\sqrt{2}} \sqrt{\sum_{n=1}^N (\sqrt{p_n} - \sqrt{q_n})^2}$$

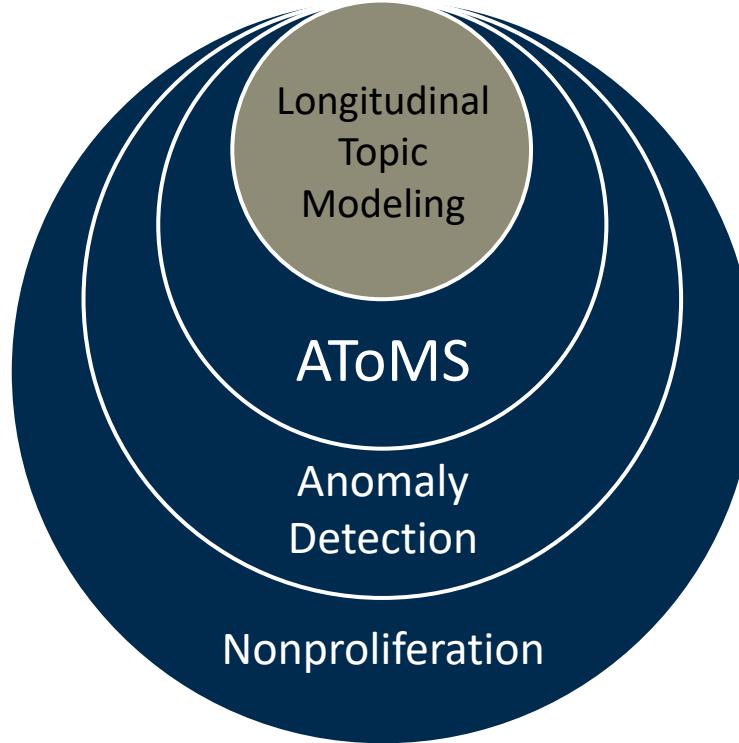
Low-Dimensional Embedding



Author-Topic Model



» GDLTM: A Geometry-Driven Longitudinal Topic Model



[1] Wang, Y., Hougen, C., Oselio, B., Dempsey, W., & Hero, A. "A Geometry-Driven Longitudinal Topic Model." Harvard Data Science Review (2021).

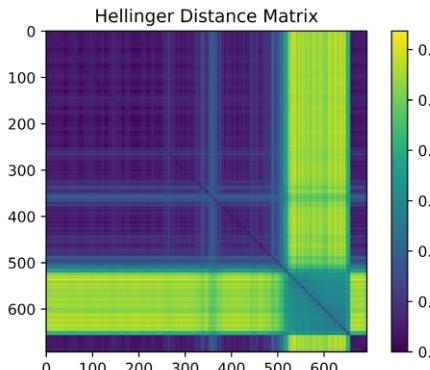
» GDLTM: A Geometry-Driven Longitudinal Topic Model

Temporal Smoothing

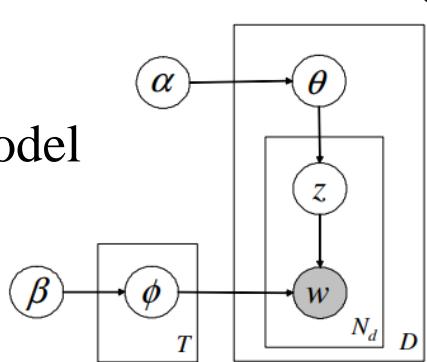
	d_1	d_2	d_3	d_4	d_5
r_1	Doc1	Doc6	Doc11	Doc16	Doc21
r_2	Doc2	Doc7	Doc12	Doc17	Doc22
r_3	Doc3	Doc8	Doc13	Doc18	Doc23
r_4	Doc4	Doc9	Doc14	Doc19	Doc24
r_5	Doc5	Doc10	Doc15	Doc20	Doc25

	s_1	s_2	s_3	s_4	s_5
s_1	Doc1	Doc6	Doc11	Doc17	Doc23
s_2	Doc2	Doc7	Doc12	Doc19	Doc24
s_3	Doc3	Doc8	Doc13	Doc20	Doc25
s_4	Doc4	Doc9	Doc15		
s_5	Doc5	Doc10			

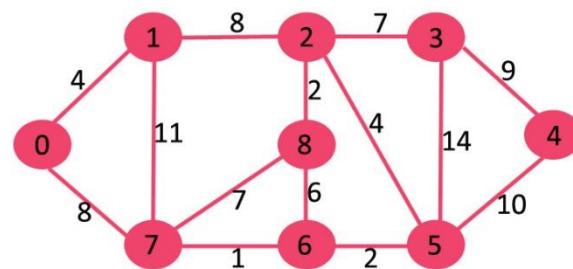
Distance Matrix



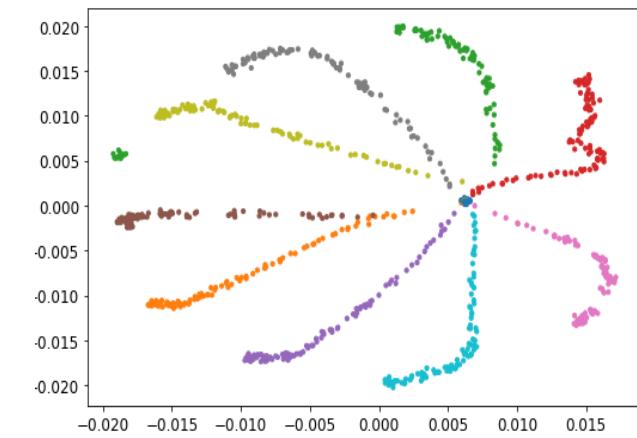
Topic Model



Shortest Path Algorithm



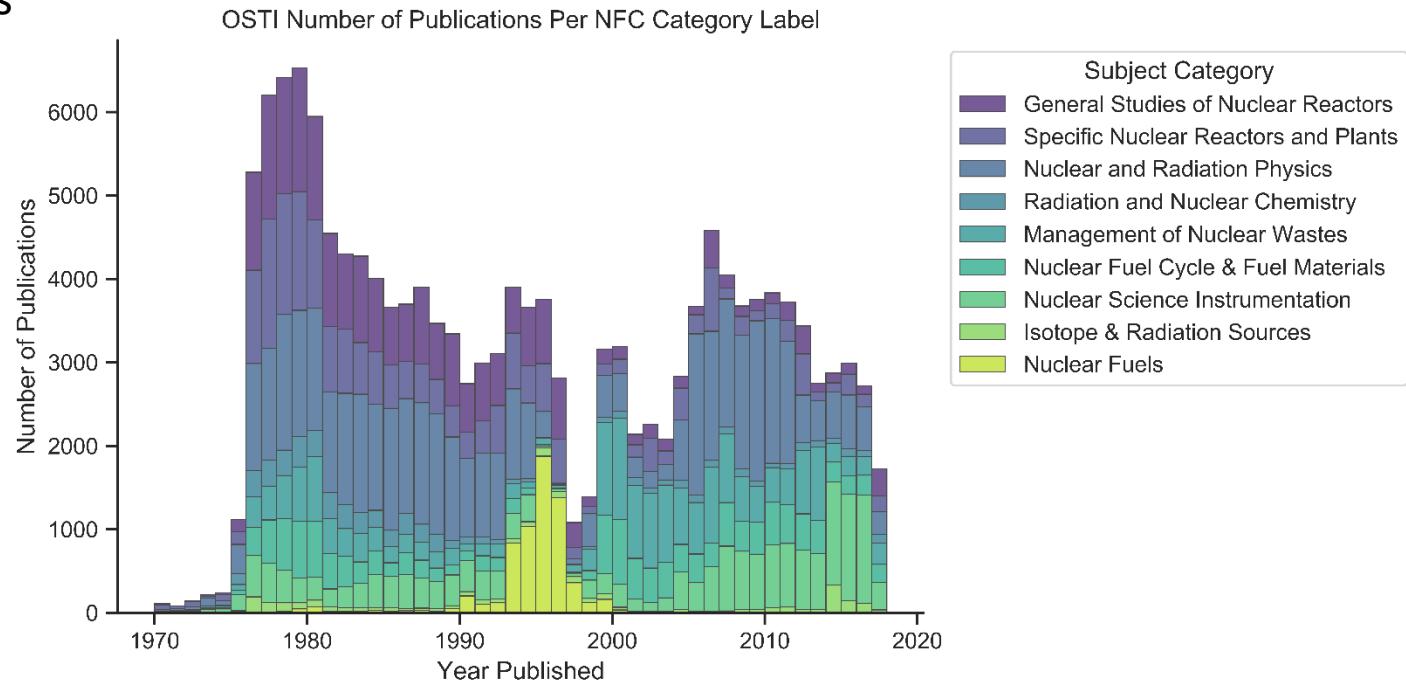
Low-Dimensional Embedding



[1] Wang, Y., Hougen, C., Oselio, B., Dempsey, W., & Hero, A. "A Geometry-Driven Longitudinal Topic Model." Harvard Data Science Review (2021).

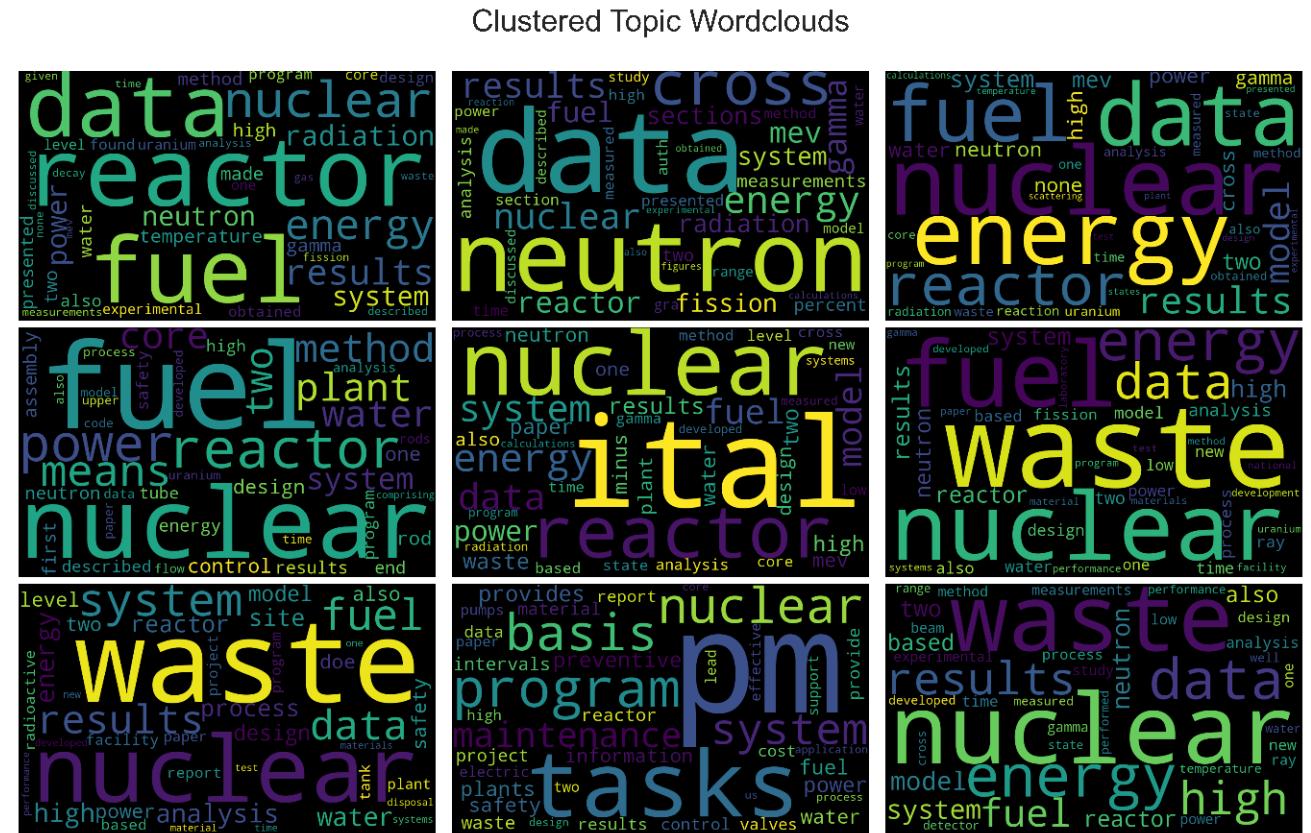
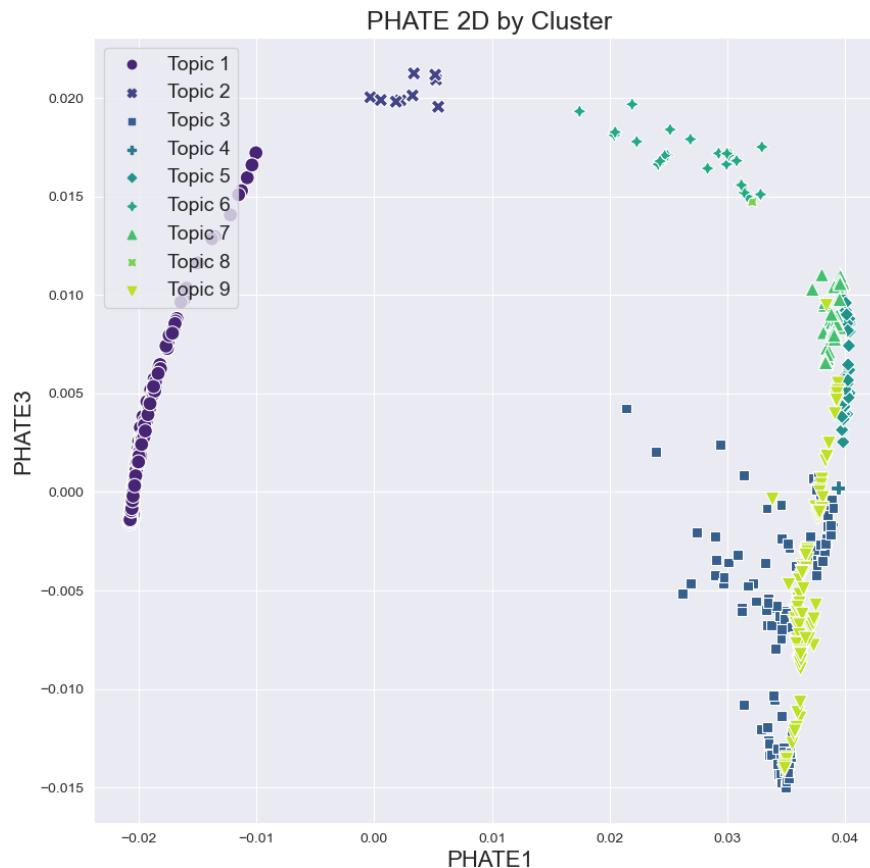
» GDLTM: A Geometry-Driven Longitudinal Topic Model

- GDLTM on OSTI dataset of manuscript abstracts
 - Filtered on **Nuclear Fuel Cycle** labeled papers
 - ~150,000 manuscripts
 - ~360,000 authors



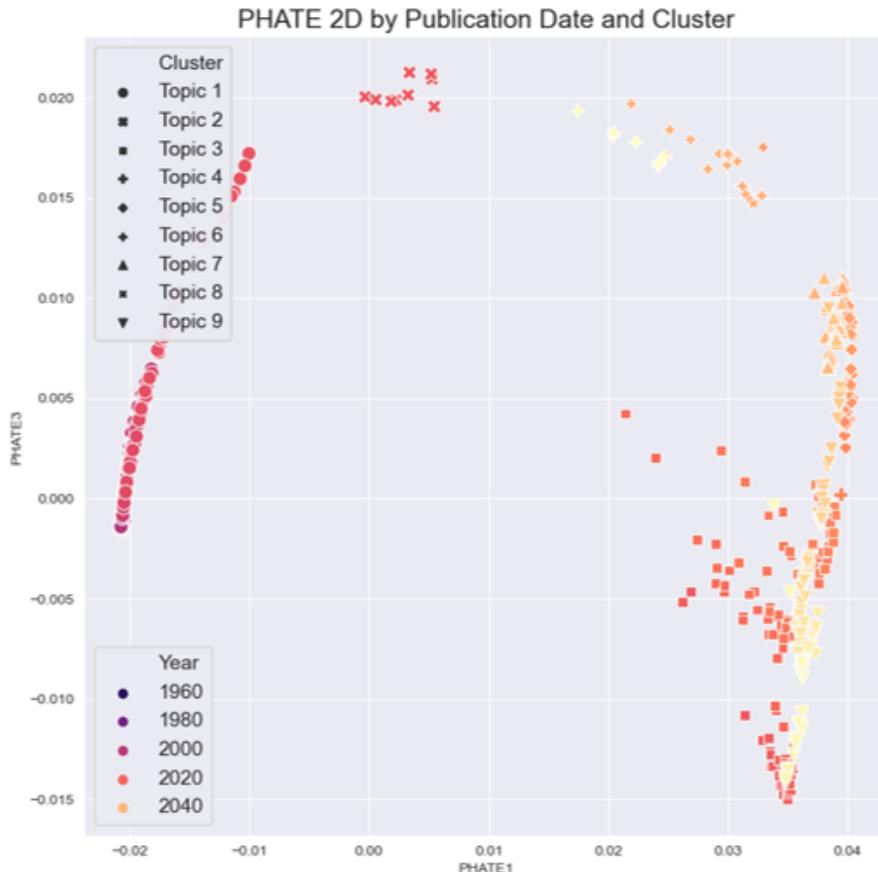
[1] Wang, Y., Hogen, C., Oselio, B., Dempsey, W., & Hero, A. "A Geometry-Driven Longitudinal Topic Model." Harvard Data Science Review (2021).

>> GDLTM: A Geometry-Driven Longitudinal Topic Model

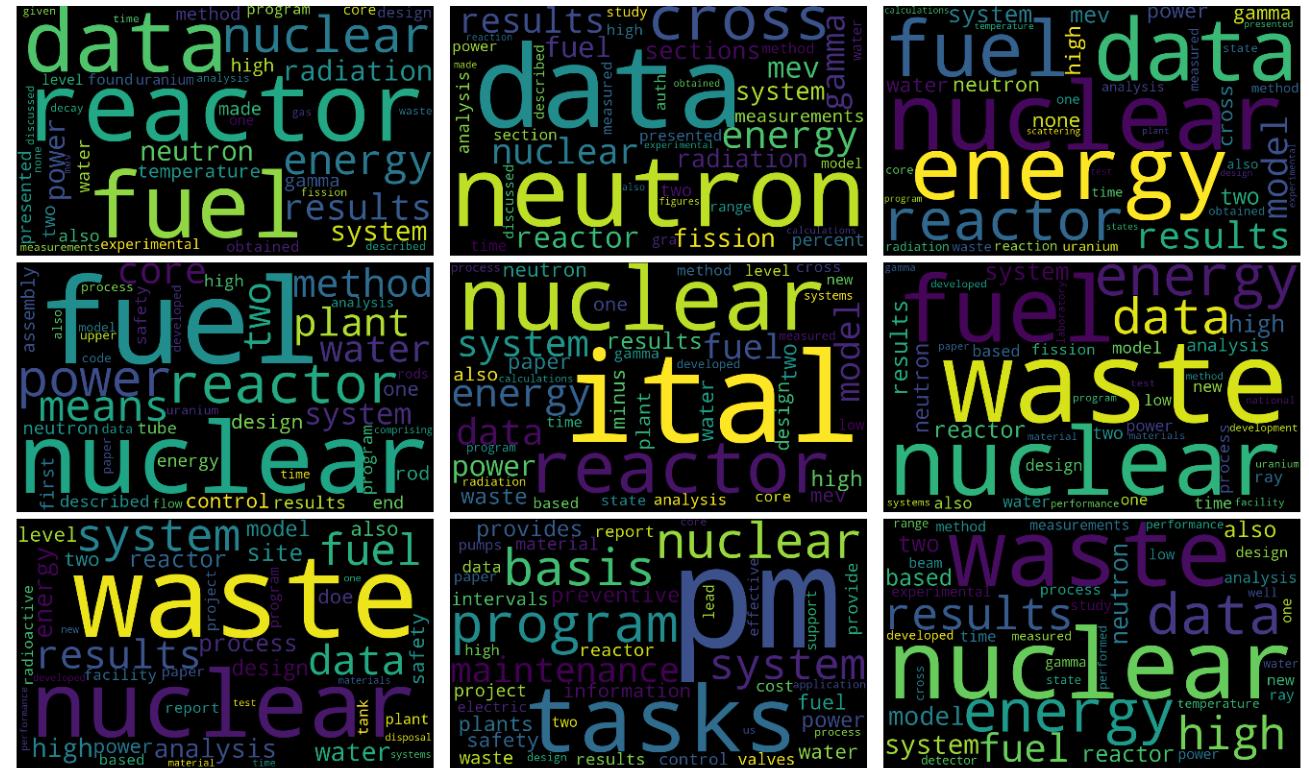


[1] Wang, Y., Hougen, C., Oselio, B., Dempsey, W., & Hero, A. "A Geometry-Driven Longitudinal Topic Model." Harvard Data Science Review (2021).

>> GDLTM: A Geometry-Driven Longitudinal Topic Model



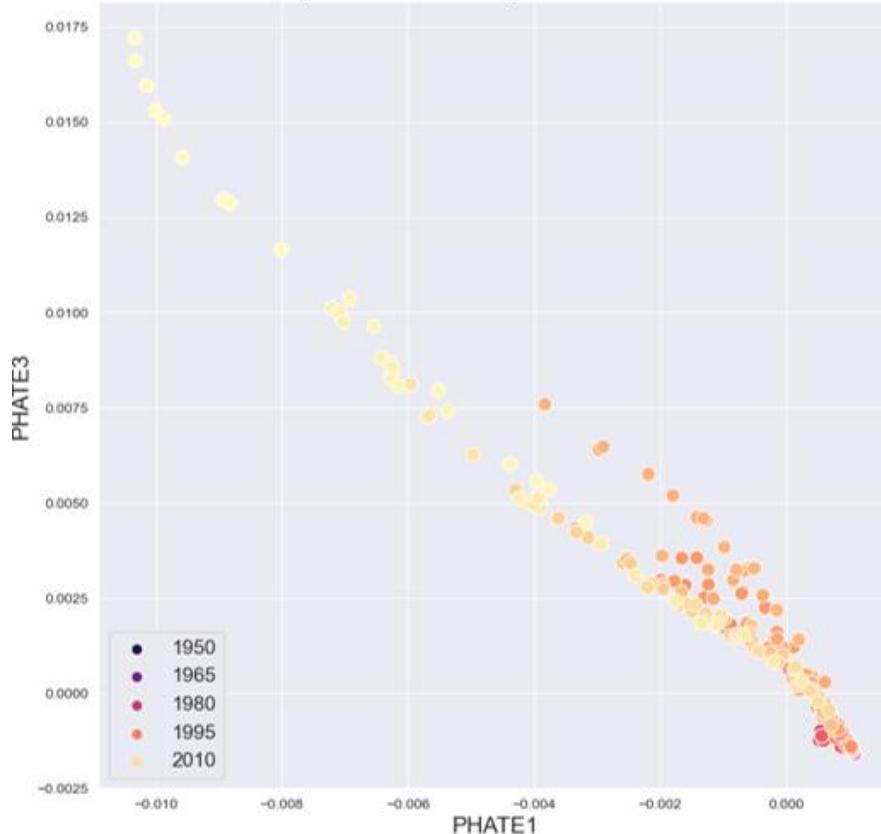
Clustered Topic Wordclouds



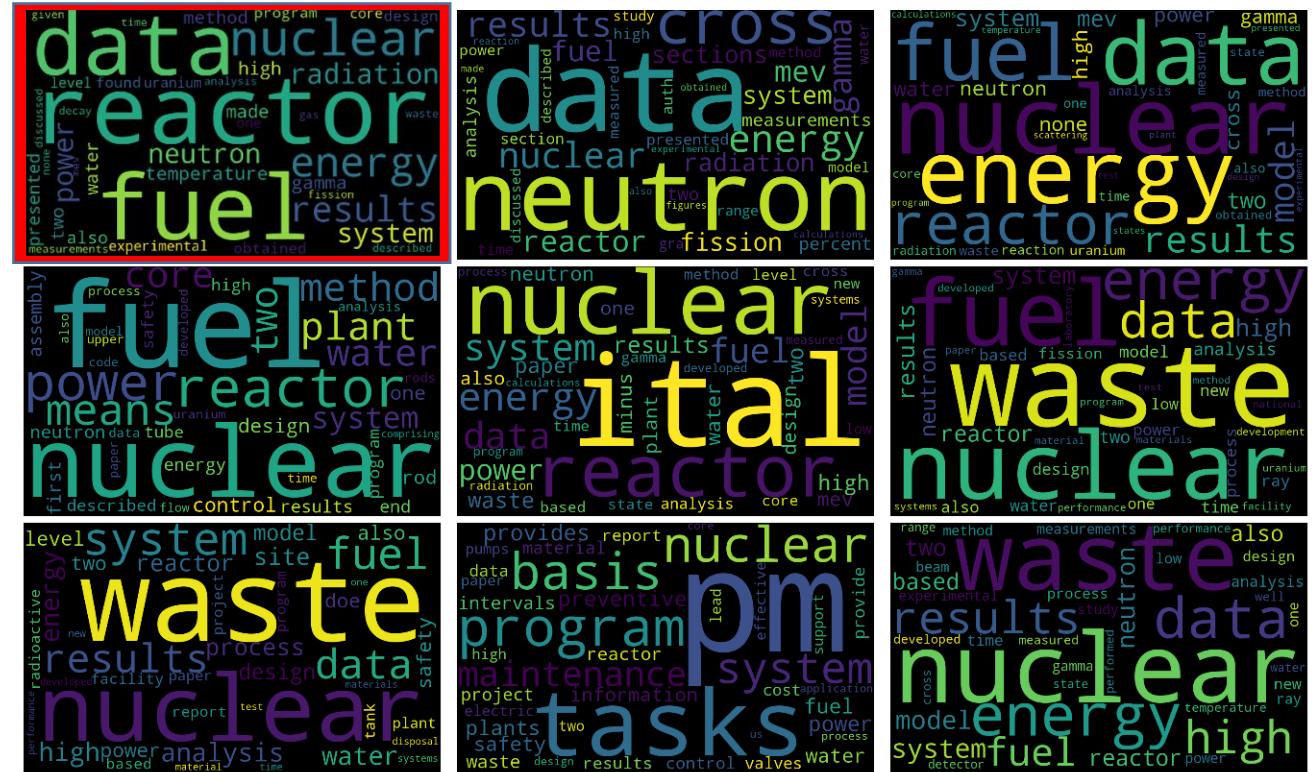
[1] Wang, Y., Hougen, C., Oselio, B., Dempsey, W., & Hero, A. "A Geometry-Driven Longitudinal Topic Model." Harvard Data Science Review (2021).

>> GDLTM: A Geometry-Driven Longitudinal Topic Model

Topic 1 PHATE 2D by Publication Date

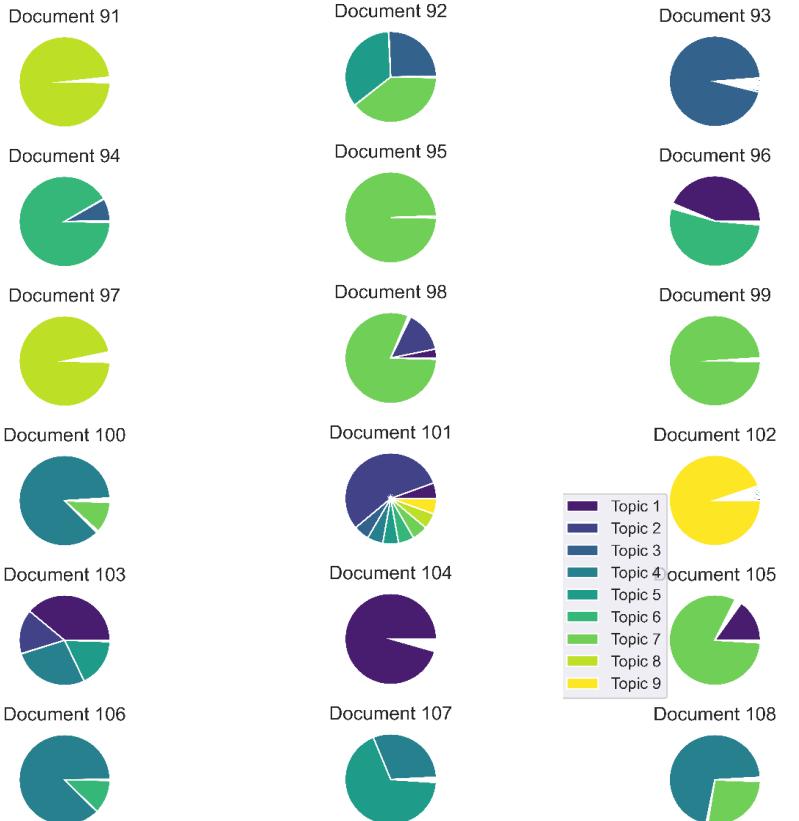


Clustered Topic Wordclouds

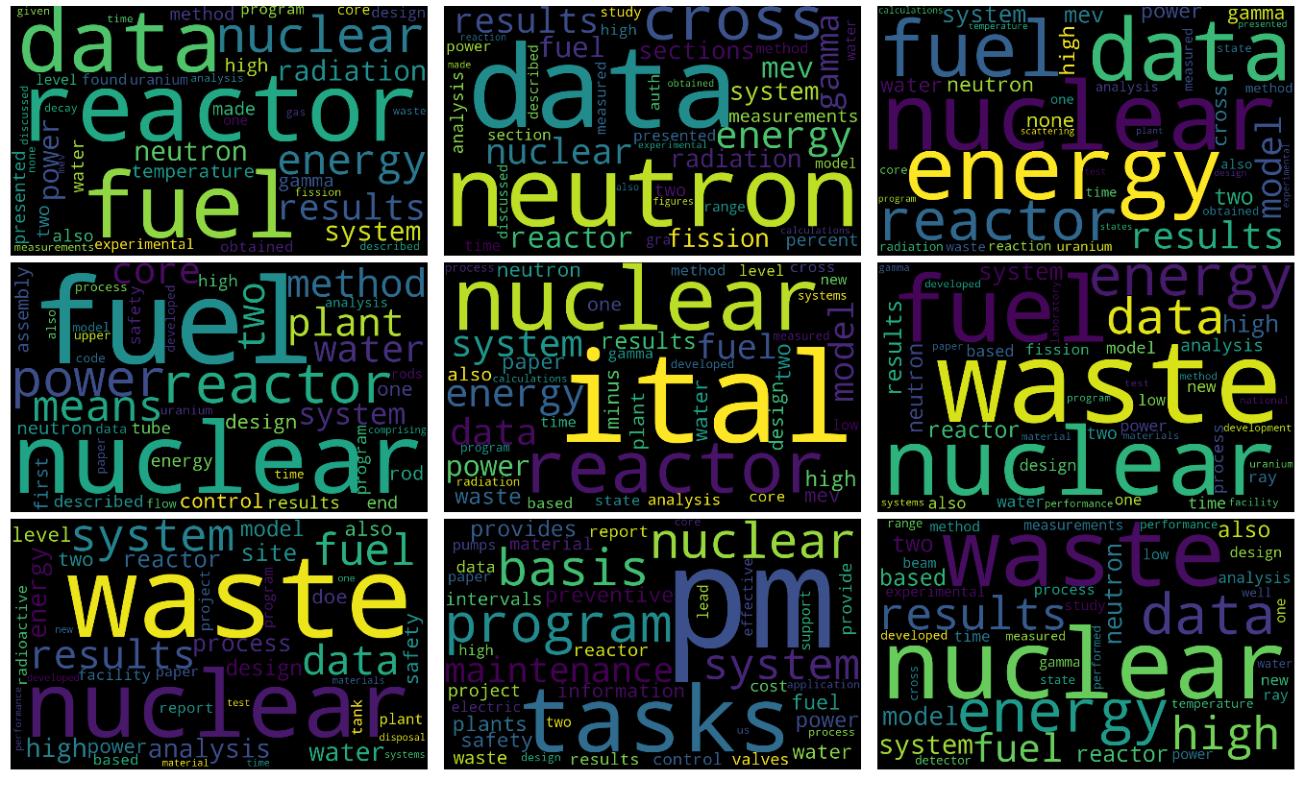


[1] Wang, Y., Hougen, C., Oselio, B., Dempsey, W., & Hero, A. "A Geometry-Driven Longitudinal Topic Model." Harvard Data Science Review (2021).

>> GDLTM: A Geometry-Driven Longitudinal Topic Model

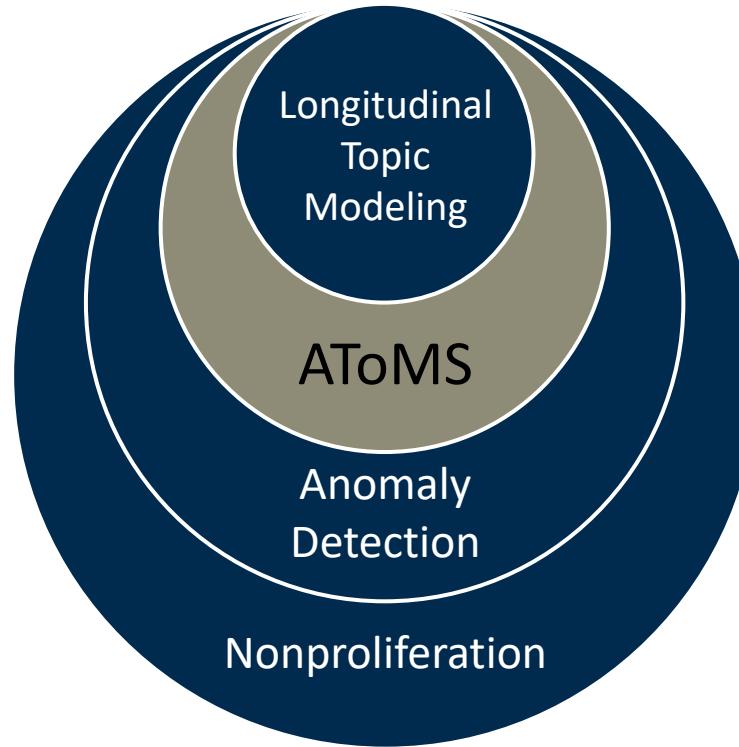


Clustered Topic Wordclouds

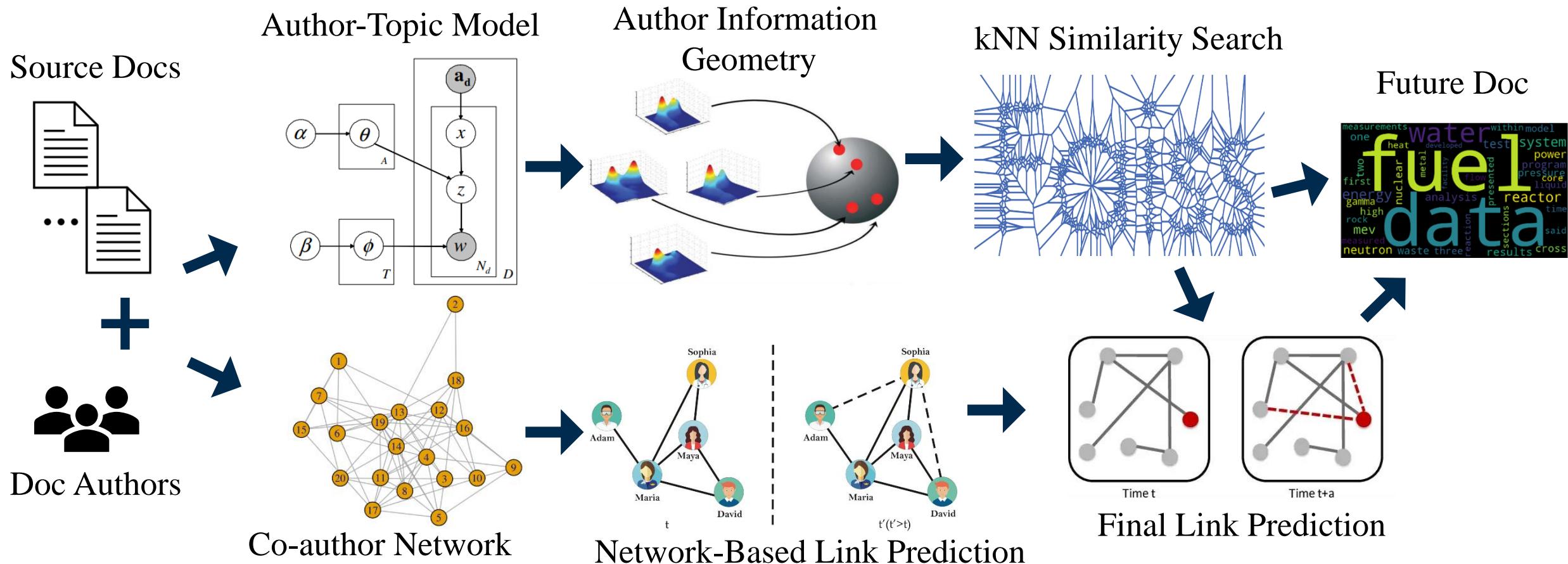


[1] Wang, Y., Hougen, C., Oselio, B., Dempsey, W., & Hero, A. "A Geometry-Driven Longitudinal Topic Model." Harvard Data Science Review (2021).

» AToMS: Author Topic Manifold Summarization



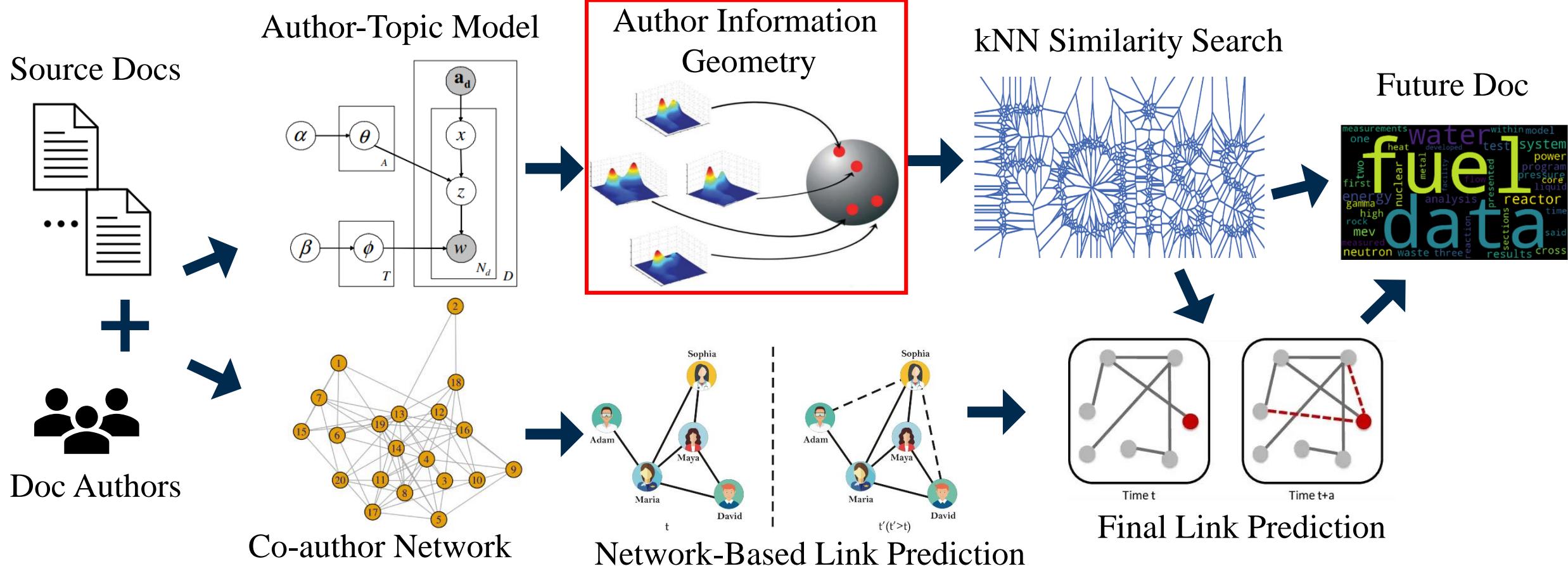
» AToMS: Author Topic Manifold Summarization



[2] Ahmad, I., Akhtar, M.U., Noor, S. *et al.* Missing Link Prediction using Common Neighbor and Centrality based Parameterized Algorithm. *Sci Rep* 10, 364 (2020).

[3] Michal Rosen-Zvi, Thomas Griffiths, Mark Steyvers, and Padhraic Smyth. 2004. The author-topic model for authors and documents. In Proceedings of the 20th conference on Uncertainty in artificial intelligence (UAI '04). AUAI Press, Arlington, Virginia, USA, 487–494.

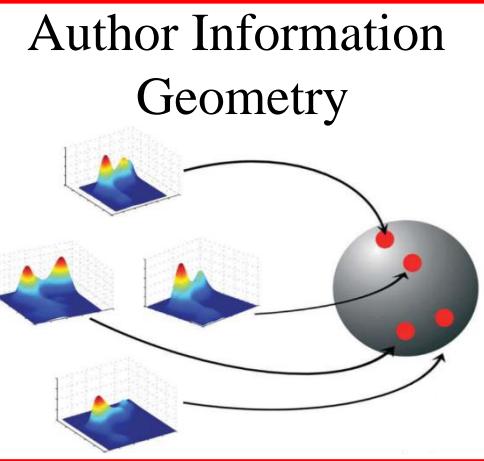
» AToMS: Author Topic Manifold Summarization



[2] Ahmad, I., Akhtar, M.U., Noor, S. *et al.* Missing Link Prediction using Common Neighbor and Centrality based Parameterized Algorithm. *Sci Rep* 10, 364 (2020).

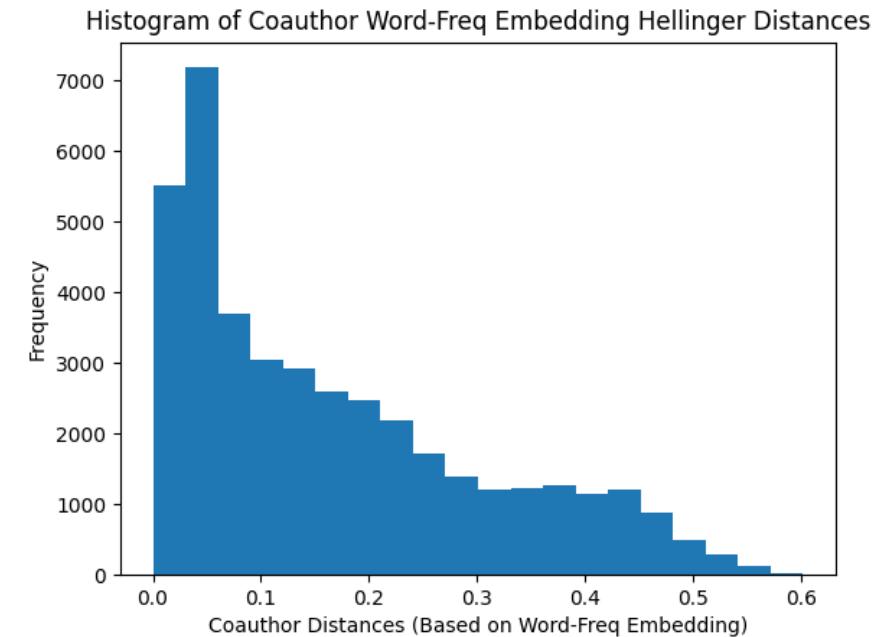
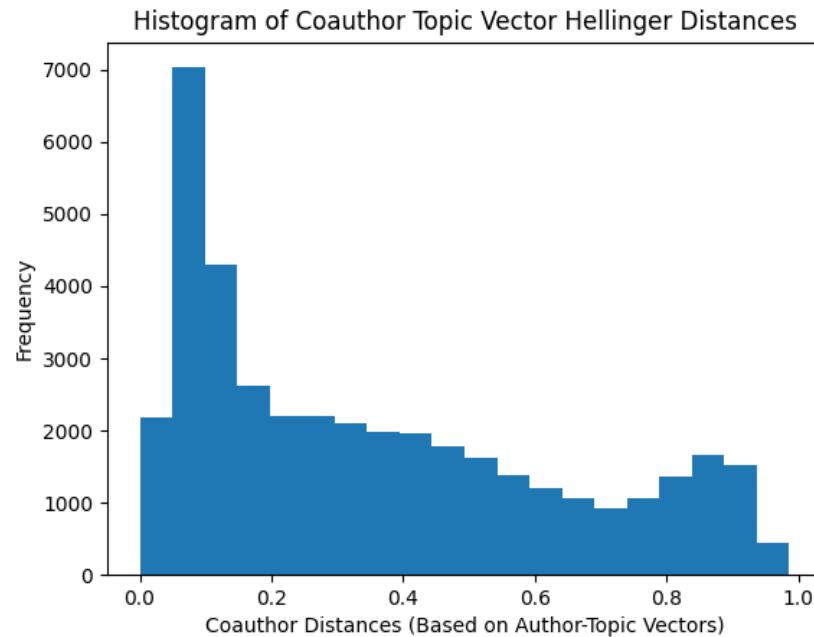
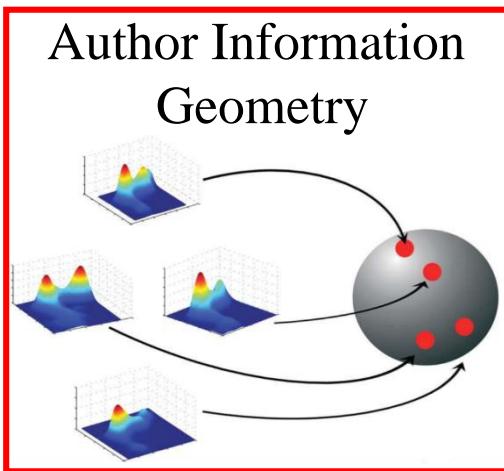
[3] Michal Rosen-Zvi, Thomas Griffiths, Mark Steyvers, and Padhraic Smyth. 2004. The author-topic model for authors and documents. In Proceedings of the 20th conference on Uncertainty in artificial intelligence (UAI '04). AUAI Press, Arlington, Virginia, USA, 487–494.

» AToMS: Author Topic Manifold Summarization



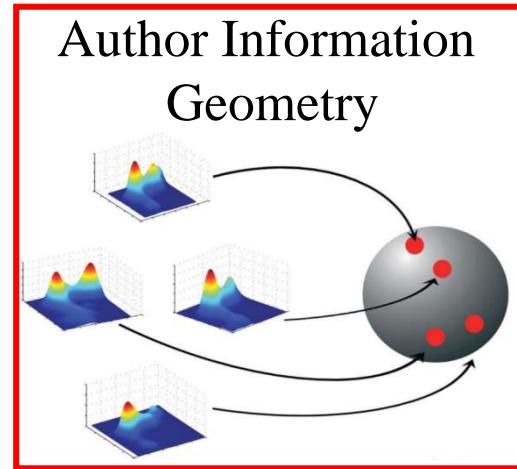
- Important Question: How should we represent authors?
 - Representation depends on distance metric
 - Hellinger distance: $H(P, Q) = \frac{1}{\sqrt{2}} \sqrt{\sum_{n=1}^N (\sqrt{p_n} - \sqrt{q_n})^2}$
 - Topic vectors are not always well-separated
- Alternative symmetric distance metrics may have disadvantages
 - E.g. Wasserstein, cosine, etc.
 - Depends on the author embedding

» AToMS Author Vectors

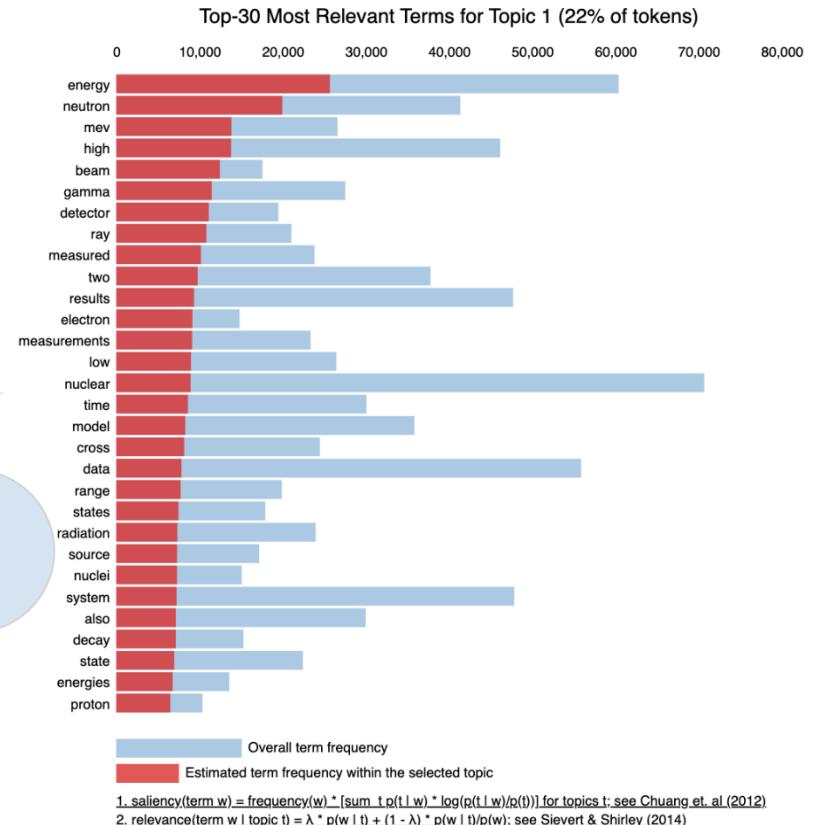


- Topic vectors are not always well-separated

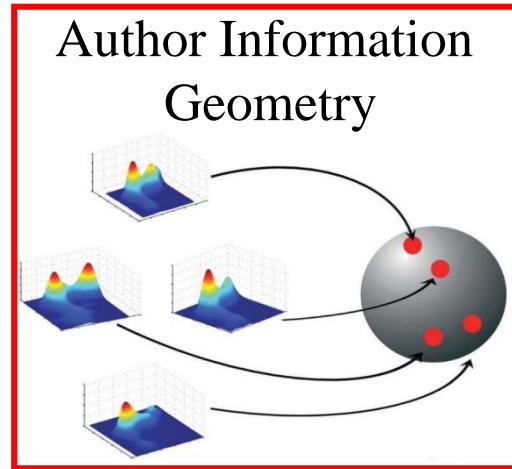
» AToMS: Author Topic Manifold Summarization



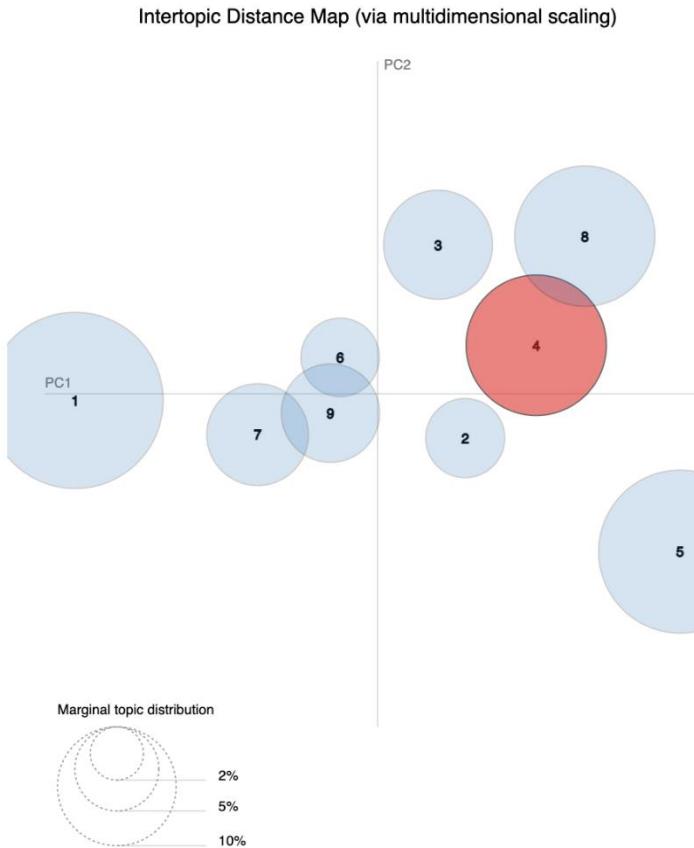
- Topic vectors are not always well-separated



» AToMS: Author Topic Manifold Summarization



- Topic vectors are not always well-separated

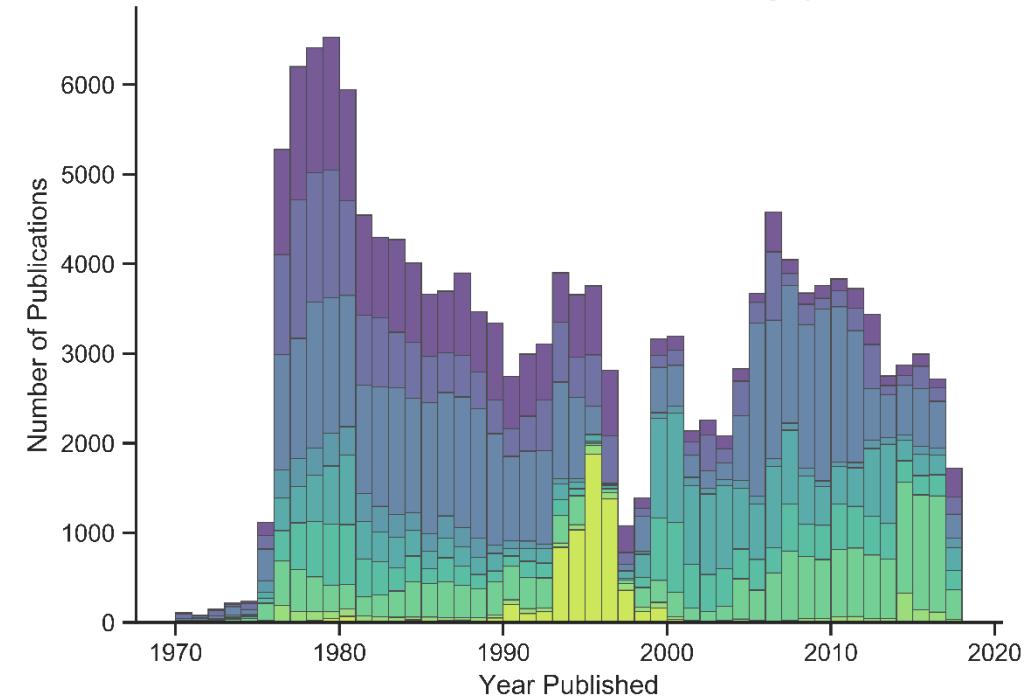


>> Experimental Setup

- OSTI (DOE Office of Scientific and Technical Information)
 - Database of articles from DOE-funded R&D
 - Title, author, OSTI subject label, abstract text, etc.
- Nuclear Fuel Cycle Documents
 - Filtered from OSTI using subject labels
 - 9 OSTI subject labels are considered NFC related
- Preprocessing
 - Filter out “stop words” – standard practice in NLP
 - Author entity disambiguation
 - Remove documents with missing abstract or authors
- Experimental Setup
 - Reduce corpus to specific subsets for train and test
 - Consider only authors that exist in the training set
 - Predict links between pairs of authors

Subject Category	
General Studies of Nuclear Reactors	
Specific Nuclear Reactors and Plants	
Nuclear and Radiation Physics	
Radiation and Nuclear Chemistry	
Management of Nuclear Wastes	
Nuclear Fuel Cycle & Fuel Materials	
Nuclear Science Instrumentation	
Isotope & Radiation Sources	
Nuclear Fuels	

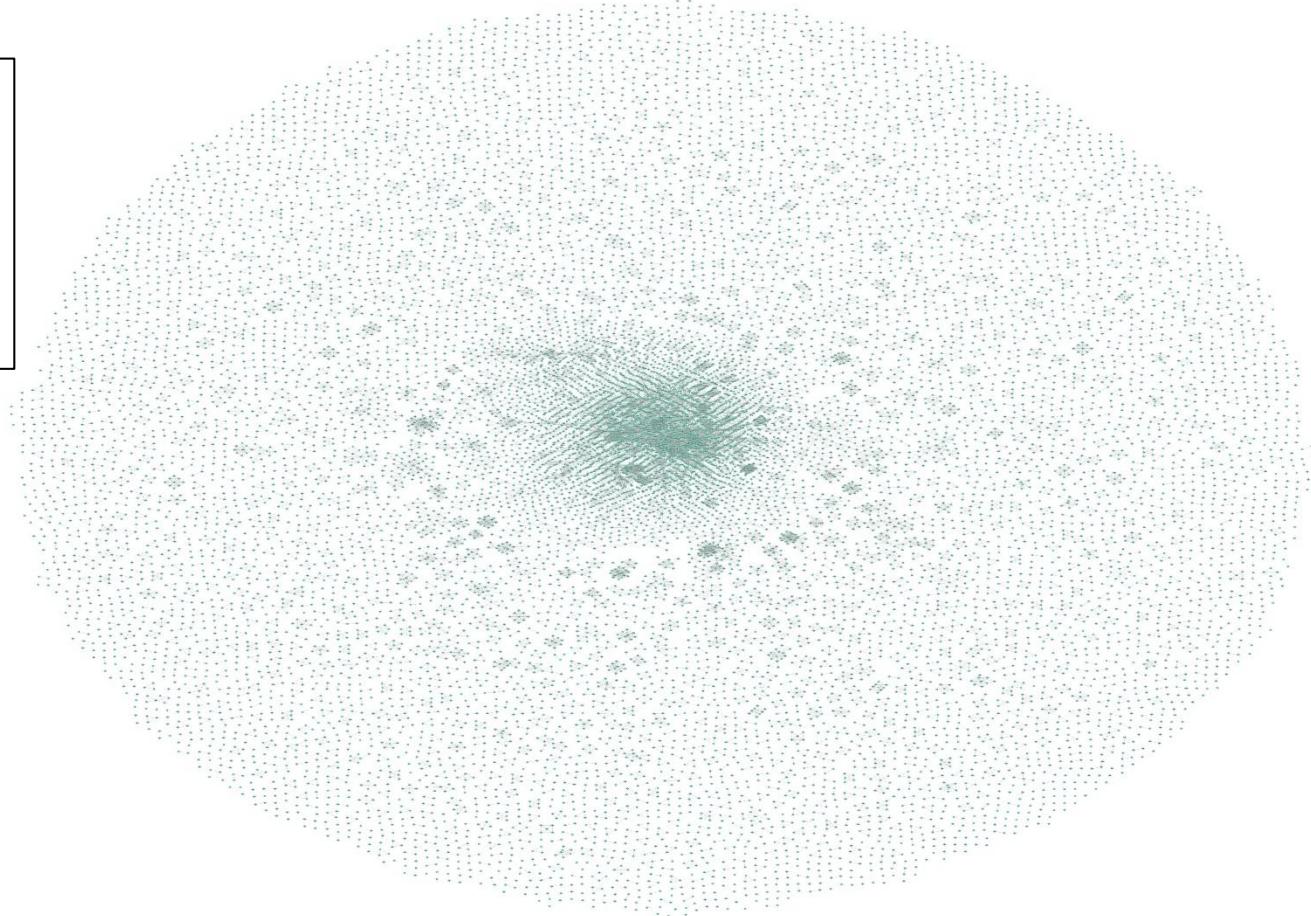
OSTI Number of Publications Per NFC Category Label

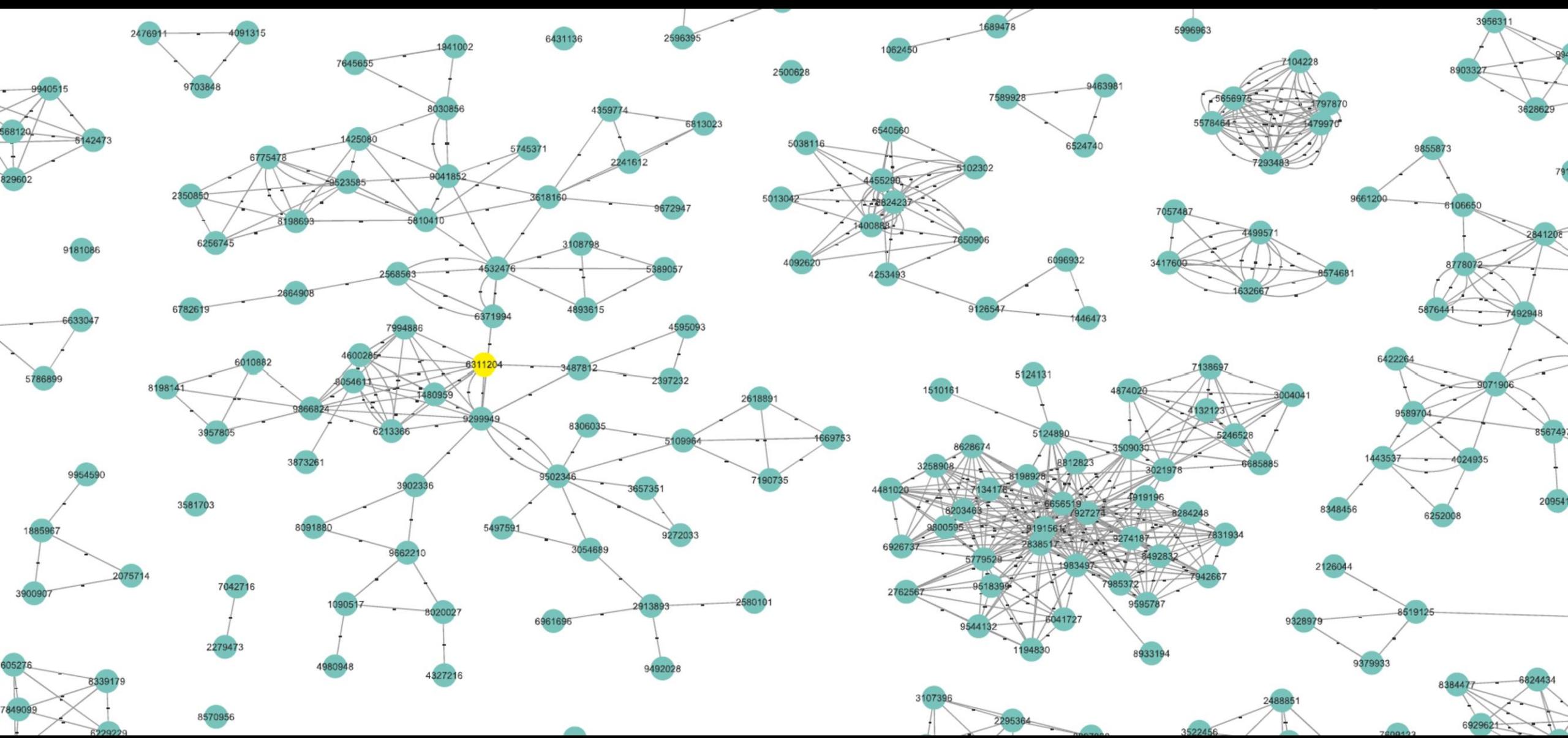


» OSTI NFC: ['81-'82][⁹⁰][15,471][36,282][52,779]

- Training Data
 - OSTI NFC documents from 1981 through 1982
 - 15,471 authors
 - 36,282 edges in coauthorship network
 - 8,846 manuscript abstracts

Property Name	Value
Avg. Number of Neighbors	8.602
Network Diameter	29
Clustering Coefficient	0.767
Network Density	0.002
Connected Components	4251



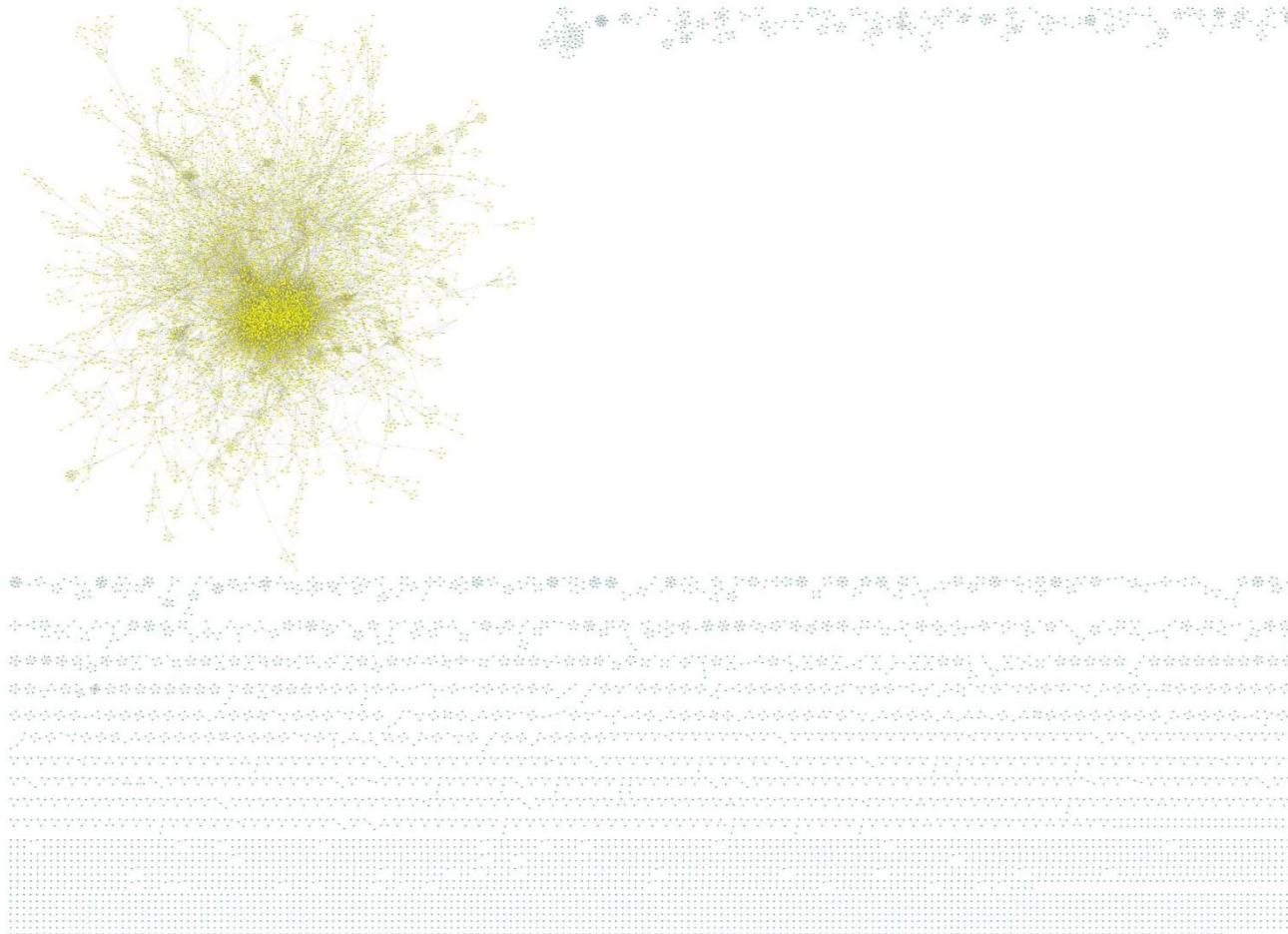


PNNL-SA-181770

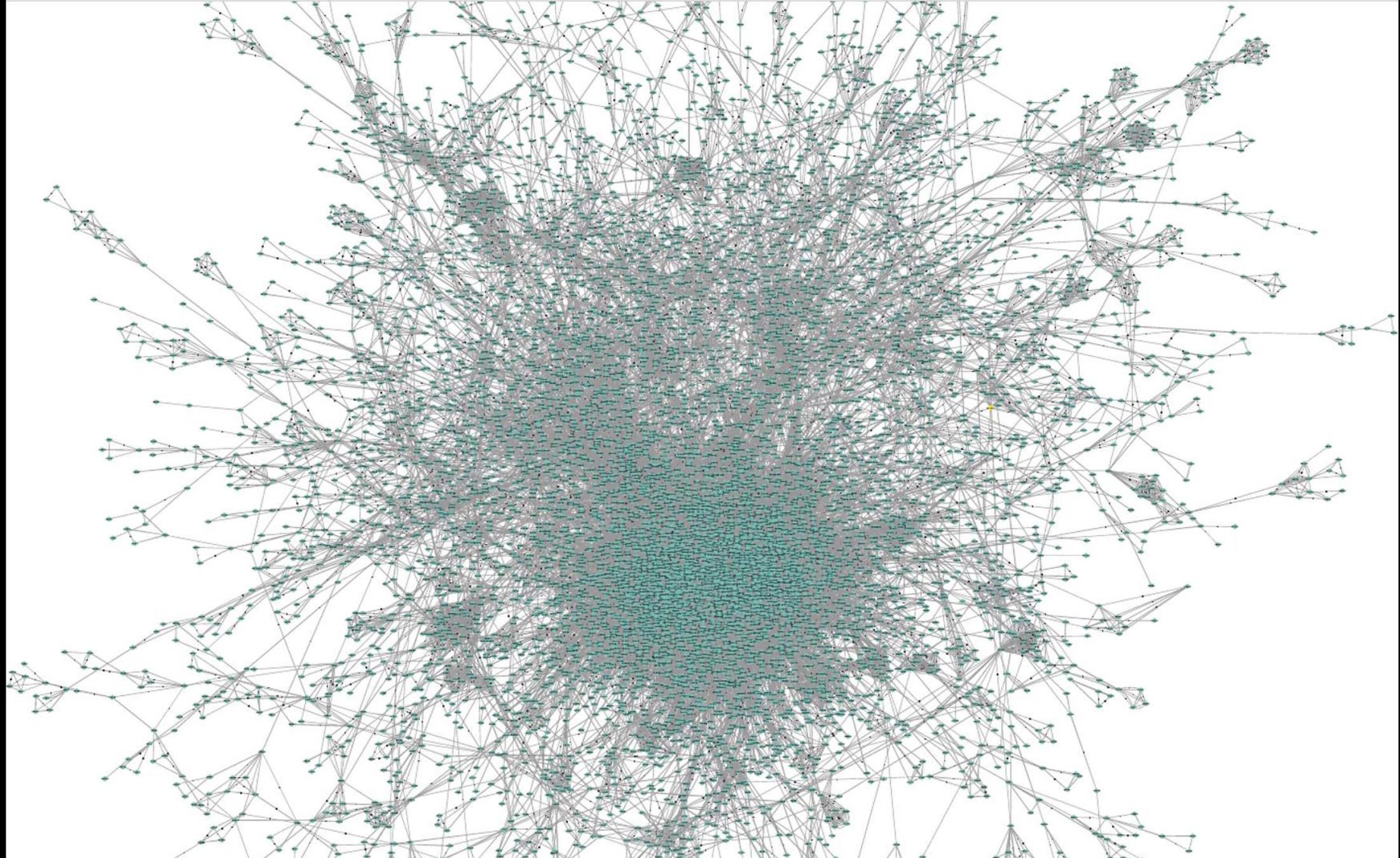
» OSTI NFC: ['81-'82]['90][15,471][36,282][52,779]

- Test Data
 - OSTI NFC documents from 1983 through 1990
 - Restricted to same 15,471 authors
 - $52,779 - 36,282 = 16,497$ new edges
 - 29,100 new manuscript abstracts

Property Name	Value
Avg. Number of Neighbors	8.864
Network Diameter	25
Clustering Coefficient	0.615
Network Density	0.001
Connected Components	3119



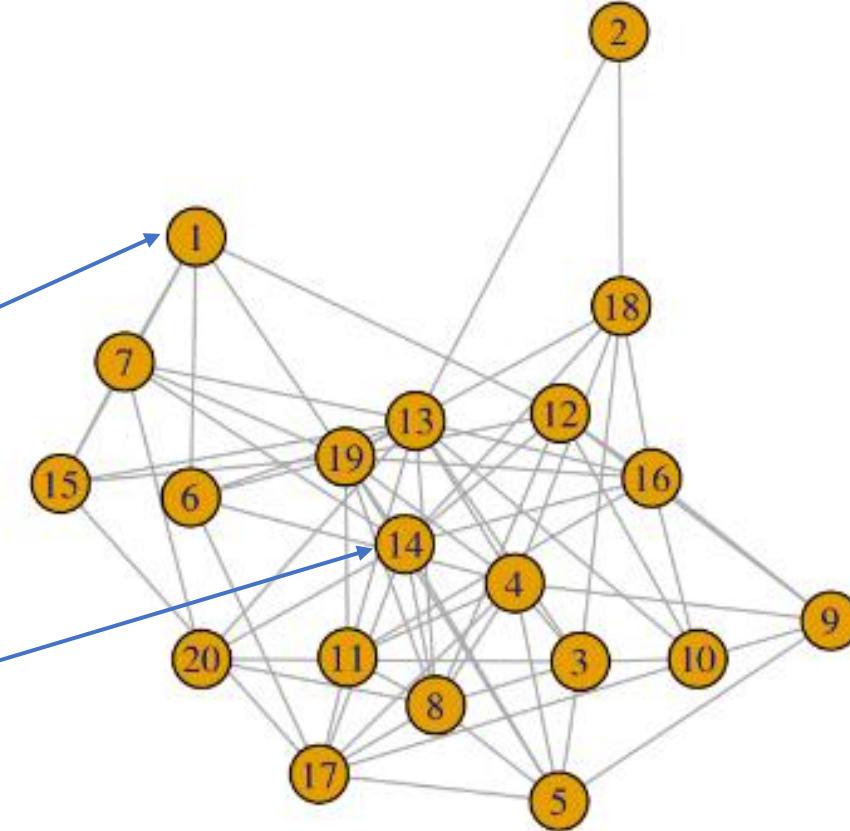
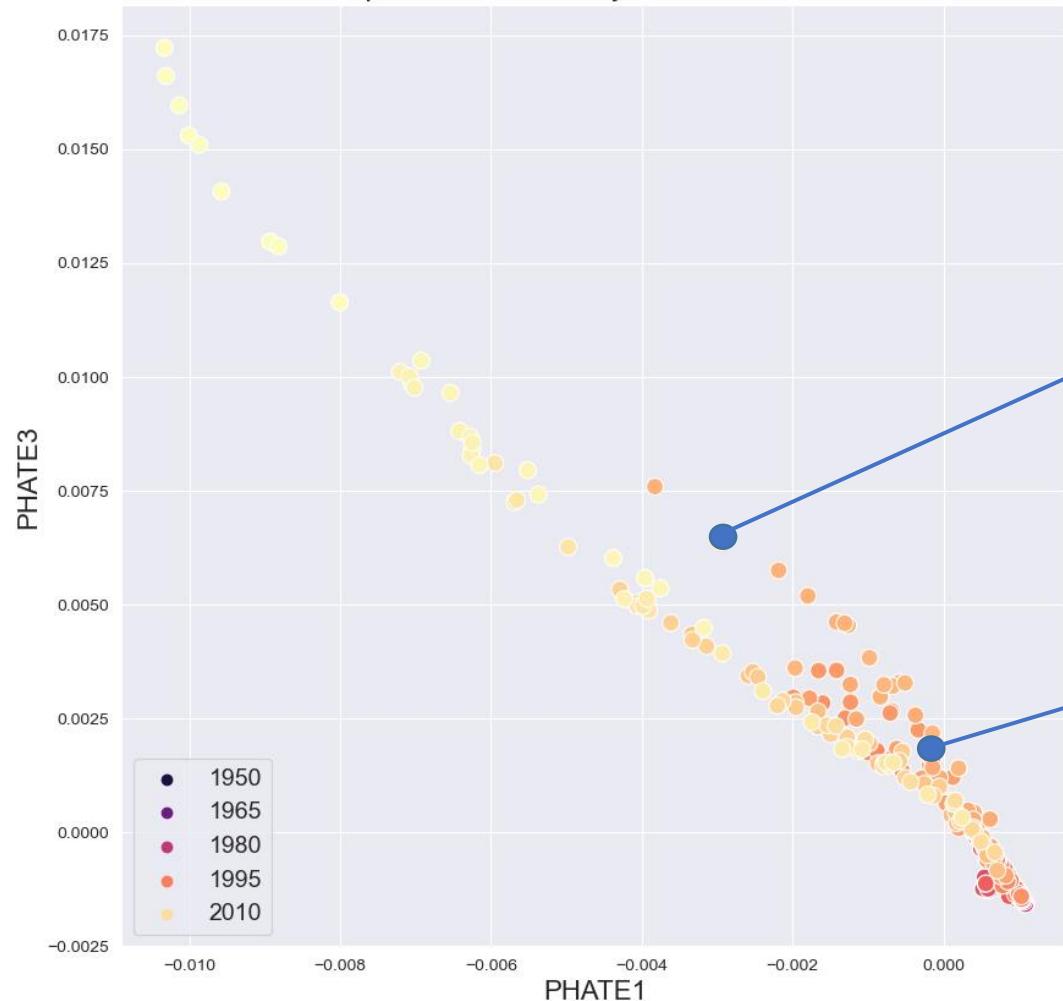
0



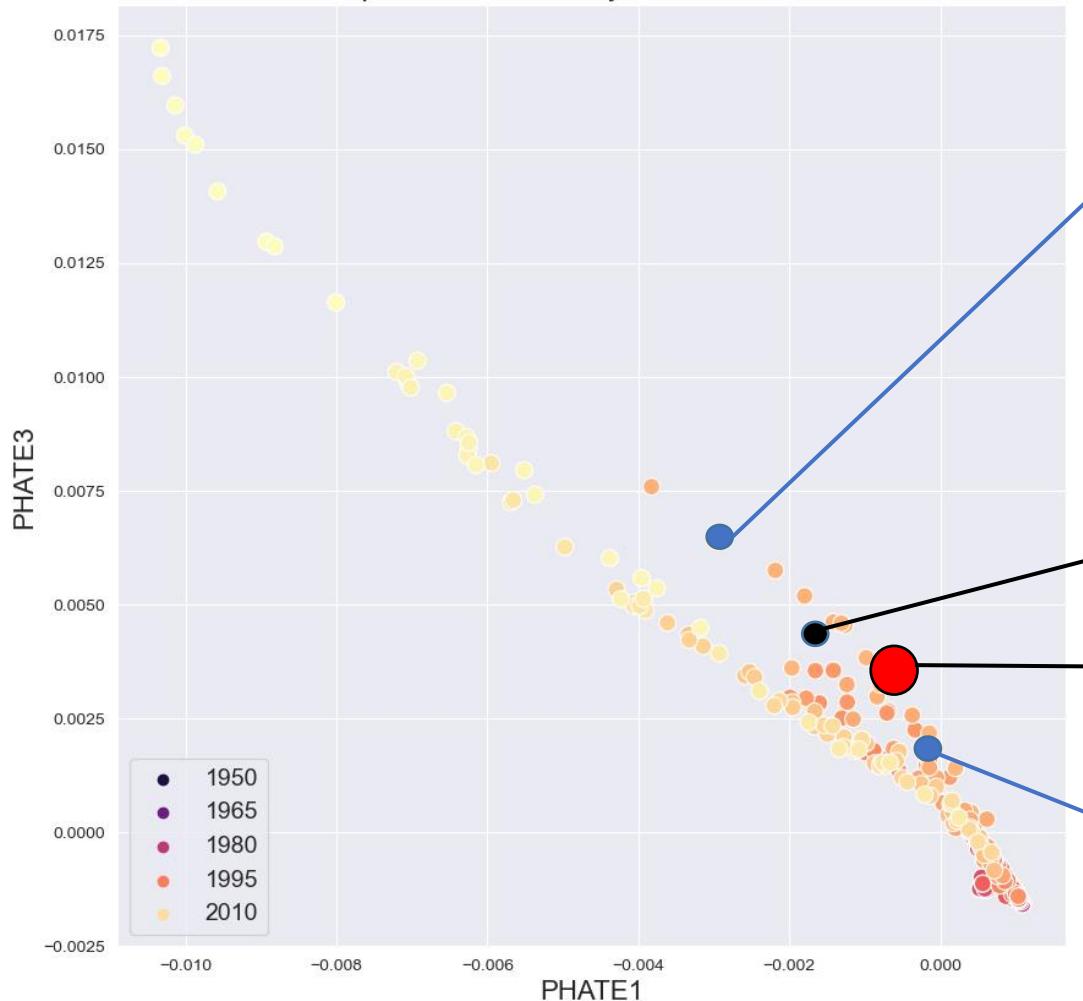
» OSTI NFC: ['81-'82]['90][15,471][36,282][52,779]

Baseline Method	TP	FP	FN	Precision	Recall	F1-Score
Graph Distance 2 (GD2)	426	17,715	16,071	0.0235	0.0258	0.0246
Common Neighbors (CN)	939	17,202	15,558	0.0518	0.0569	0.0542
Adamic-Adar (AA)	905	17,236	15,592	0.0499	0.0549	0.0523
Jaccard Coefficient (JC)	619	17,522	15,878	0.0341	0.0375	0.0357
Resource Allocation (RA)	774	17,367	15,723	0.0427	0.0469	0.0447
AToMS Method	TP	FP	FN	Precision	Recall	F1-Score
Common Neighbors/Hellinger	944	17,197	15,553	0.0520	0.0572	0.0545
Adamic-Adar/Hellinger	907	17,234	15,590	0.0500	0.0550	0.0524
Jaccard Coefficient/Hellinger	527	17,614	15,970	0.0291	0.0319	0.0304
Resource Allocation/Hellinger	781	17,360	15,716	0.0431	0.0473	0.0450

Topic 1 PHATE 2D by Publication Date



Topic 1 PHATE 2D by Publication Date



nuclear effects high model may
waste fission time measured reactor mev
surface analysis obtained two measurements
ground found reactor also energy system
three test report neutron temperature
cross sections state samples radiation results solution

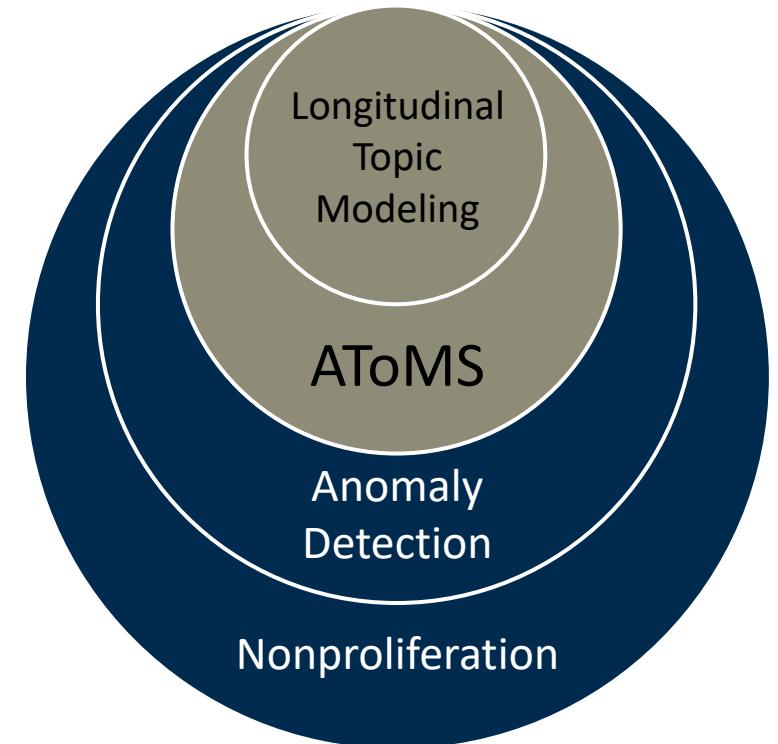
tests high rate made also test spent discussed
data one presented decay storage system water energy new
well model materials gamma radioactive
analysis results found reactor program
nu clear uranium two method

neutron system analysis temperature process
uranium presented two potential field reactor mev
fuel state high water test made code
design program well found states also
model energy cross sections experimental waste flow gas power
thermal sections analysis reaction report

measurements one heat developed test system power
within model first energy nuclear reactor pressure core
water nuclear facility analysis sections reactor time
fuel metal flow presented said sections cross

>> Conclusions and Future Work

- Link Prediction
 - Challenging problem in bibliographic networks
 - Networks are extremely sparse
 - AToMS outperforms “network-only” algorithms
 - Common neighbors is surprisingly effective
- Interpretability
 - AToMS provides predictive document modeling
 - Topic modeling reveals author interests
- Longitudinal Author-Topic Modeling
 - AToMS can be improved by using temporal information
 - Combine GDLTM longitudinal topics with AToMS model
- Nonproliferation Objective
 - Computational tools assist understanding bibliographic data
 - Can signal knowledge-based threats



ACKNOWLEDGEMENTS

This material is based upon work supported by the Department of Energy / National Nuclear Security Administration under Award Number(s) DE-NA0003921.



OSTI NFC 1.1 (897 Authors, 1275 edges, 560 edge budget)

Method Name	TP	FP	TN	FN	Precision	Recall	F1-Score
AToMSv2_pa_he	0	560	-	155	0.0000	0.0000	0.0000
AToMSv2_aa_he	16	544	-	139	0.0286	0.1032	0.0448
AToMSv2_jc_he	20	540	-	135	0.0357	0.1290	0.0559
AToMSv2_ra_he	17	543	-	138	0.0304	0.1097	0.0476
AToMSv2_cn_he	15	545	-	140	0.0268	0.0968	0.0420
AToMSv2_jc_co	15	545	-	140	0.0268	0.0968	0.0420
PA	1	559	-	154	0.0018	0.0065	0.0028
AA	15	545	-	140	0.0268	0.0968	0.0420
JC	16	544	-	139	0.0286	0.1032	0.0448
RA	16	544	-	139	0.0286	0.1032	0.0448
CN	15	545	-	140	0.0268	0.0968	0.0420
GD2	4	152	-	151	0.0256	0.0258	0.0257