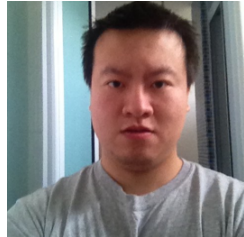


External Evaluation of Topic Models: A Graph Mining Approach

Hau Chan*



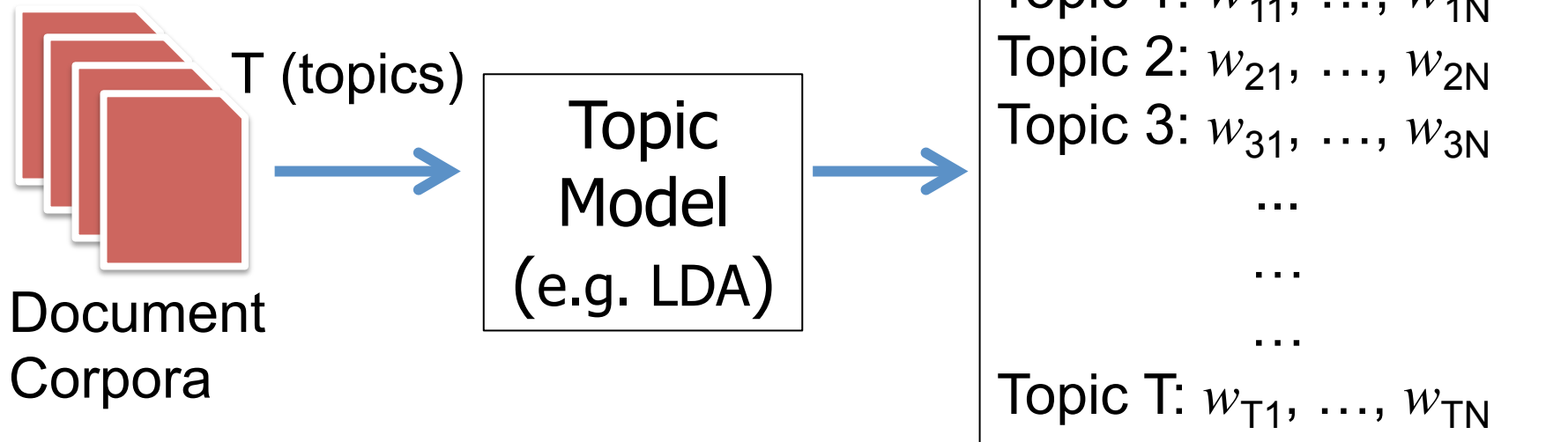
Leman Akoglu





Stony Brook University

ICDM 2013 December 7-10, 2013

Topic Models



EXAMPLE TOPICS T1 (HIGH-QUALITY) AND T2 (LOW-QUALITY) OF A TOPIC MODEL.

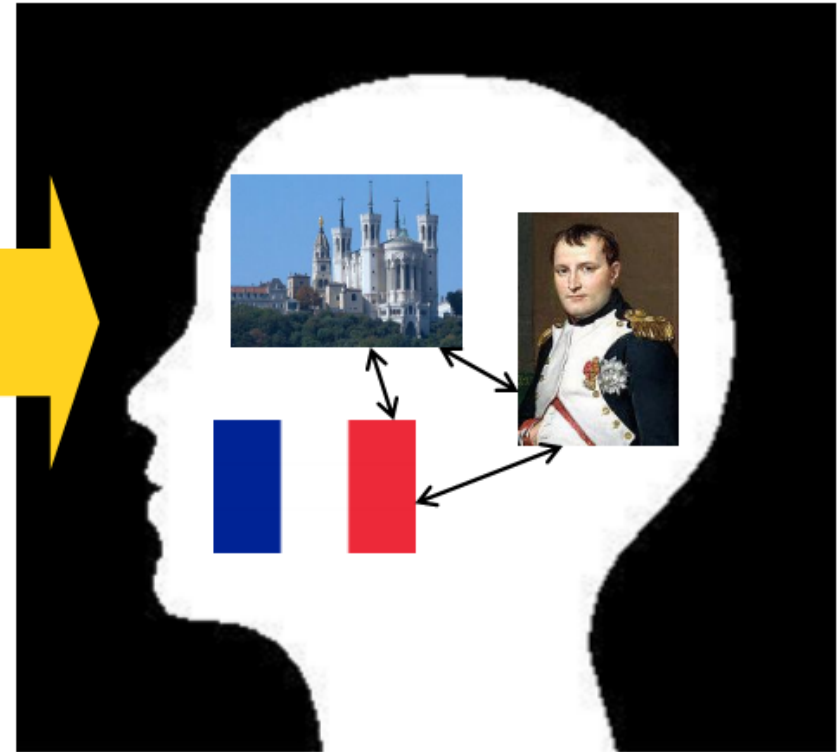
	T1: steam, engine, valve, piston, cylinder, pressure, boiler, air, pump, pipe
	T2: cut, system, capital, pointed, opening, building, character, round, france, paris

Our goal: distinguish **good** topics from **poor** topics!

Motivation

- Negative correlation (!) between
 - **Human** evaluation & **Statistical** evaluation
[Chang+ NIPS'09]
- Applications where **human-perceived quality** essential:
 - doc-doc similarity (via topic distribution)
 - word-sense disambiguation
 - multi-doc summarization

Main idea: exploit Wikipedia



Understand how humans
navigate Wikipedia

Get an idea of how
people connect concepts

[West-Leskovec, 2012]

Main idea: exploit Wikipedia

WIKIPEDIA
The Free Encyclopedia

Steam

From Wikipedia, the free encyclopedia

For other uses, see [Steam \(disambiguation\)](#).

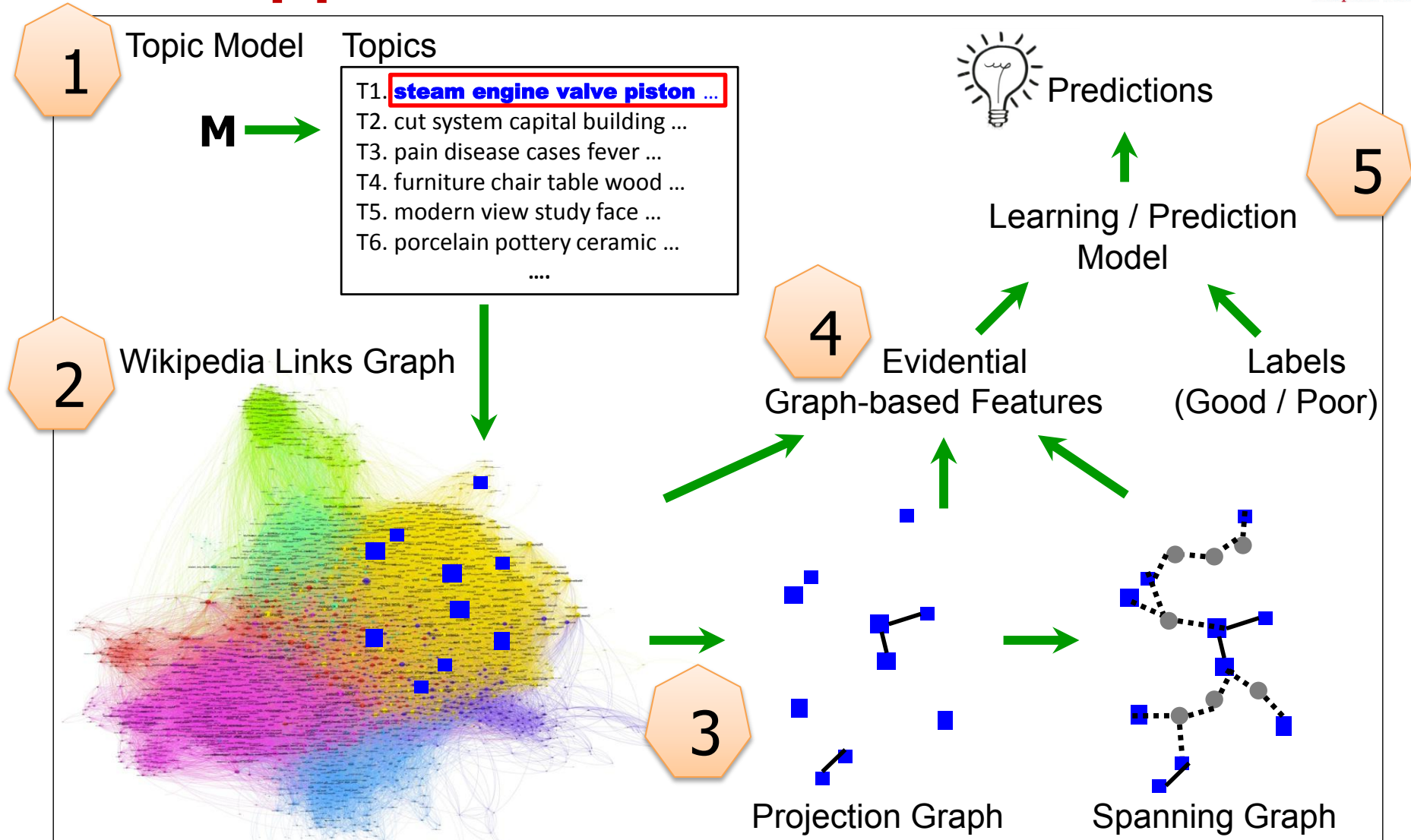
Steam is the technical term for [water vapor](#), the [gaseous phase](#) of [water](#) which is formed when water [boils](#). Technically speaking, in terms of the chemistry and physics, steam is invisible and cannot be seen; however, in common language it is often used to refer to the visible [mist](#) of water droplets formed as this water vapor [condenses](#) in the presence of (cooler) [air](#). At lower pressures, such as in the [upper atmosphere](#) or at the top of high mountains water boils at a lower temperature than the nominal 100 °C (212 °F) at [standard temperature and pressure](#). If heated further it becomes [superheated steam](#).

Wikipedia is a graph!

node = entity & (directed) edge = pagelink

**Intuition: Good topic words are conceptually
“coherent” → “close-by” in WikiLinks**

Our Approach



Projection Graph

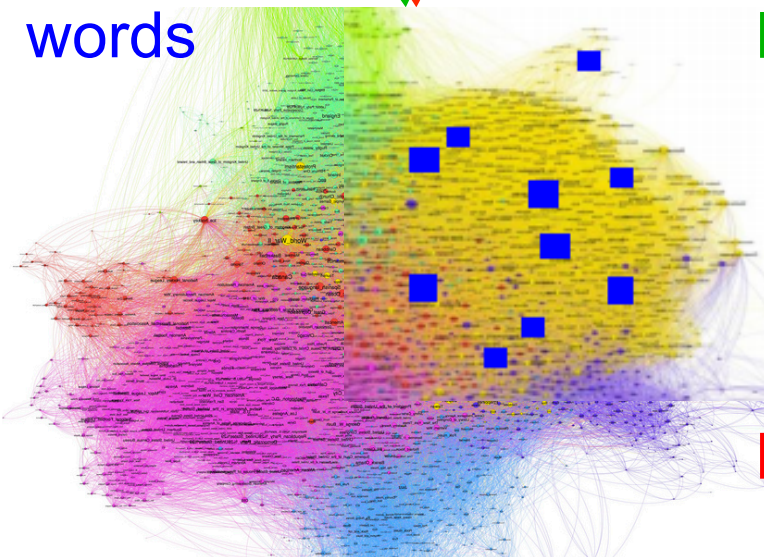
EXAMPLE TOPICS T1 (HIGH-QUALITY) AND T2 (LOW-QUALITY) OF A TOPIC MODEL.

T1: steam, engine, valve, piston, cylinder, pressure, boiler, air, pump, pipe

T2: cut, system, capital, pointed, opening, building, character, round, france, paris

match T2
words

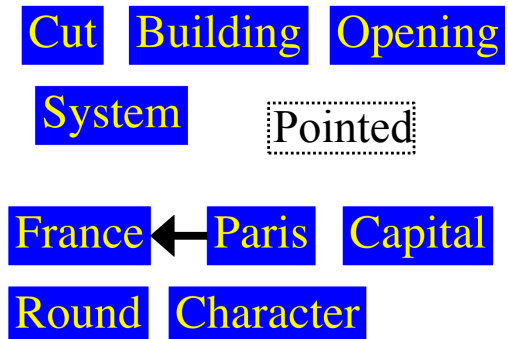
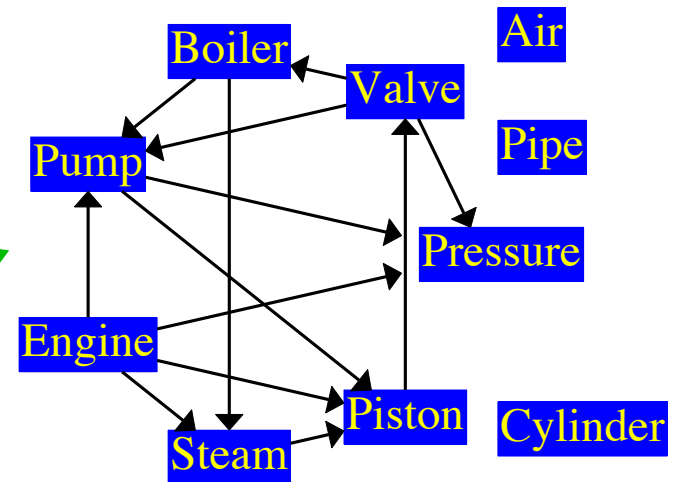
match T1
words



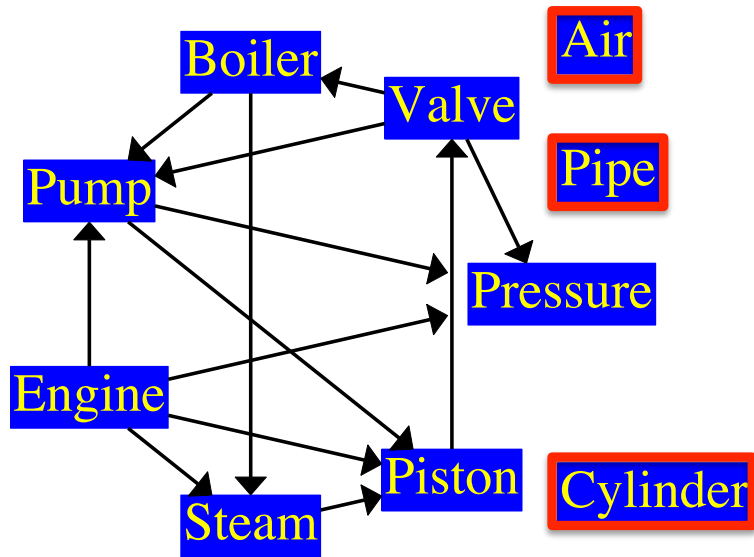
WikiLinks Graph

PROJ T1

PROJ T2



Spanning Graph

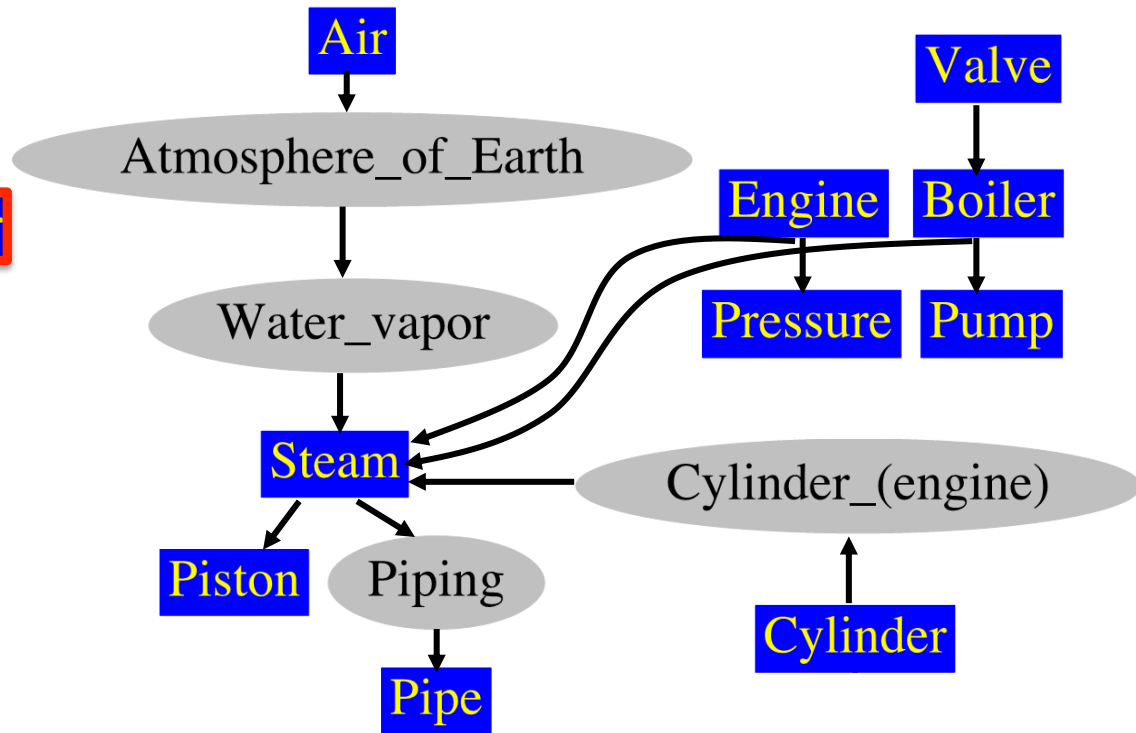


PROJ T1

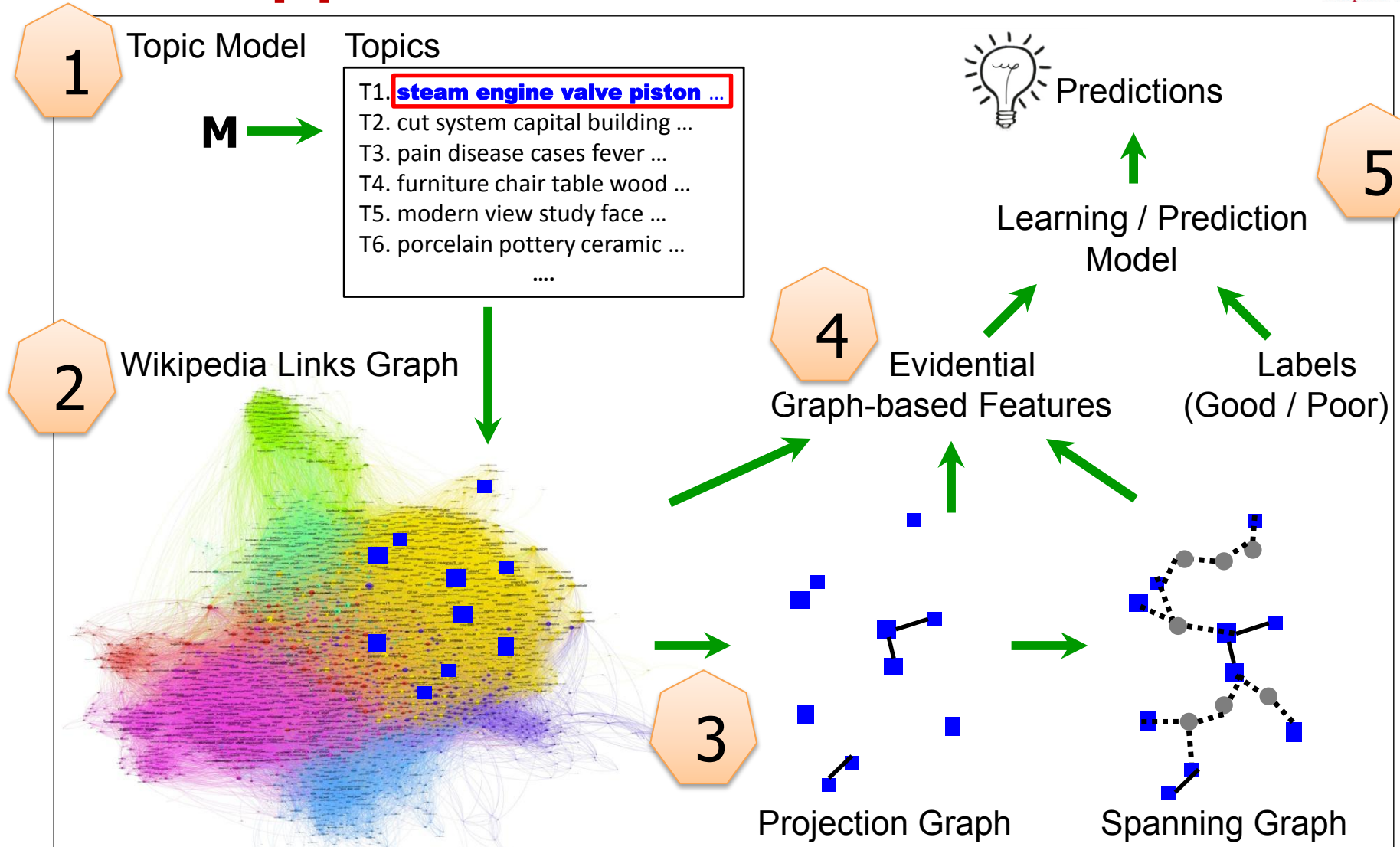
WikiLinks induced on topic-words

SPAN T1

Connection subgraph of topic-words

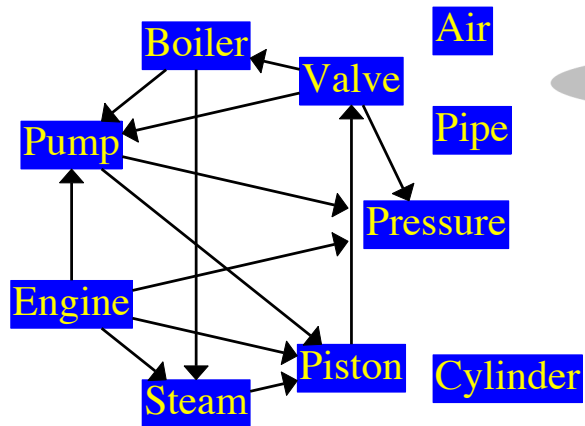


Our Approach



All Features: 3 groups

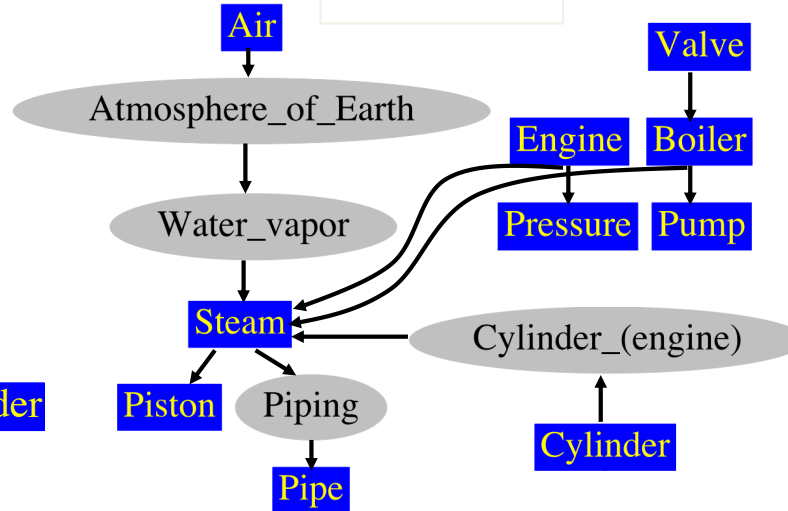
PROJ



4
features

e.g
Largest CC Size

SPAN



9
features

e.g
#Connector Nodes

SP

	Air	...	Pipe
Air	1		5
	Shortest Path		
Pipe	NP		1

9
features

e.g
Average SP

Labels: Good vs. Poor Topics

Dataset	# Documents	# Topics	Labels
Press	2,246	100	No
Brain Injury	10,000	200	No
Books	12,000	120	Yes
News	55,000	117	Yes

Relative Quality Prediction (Labels – Generated)

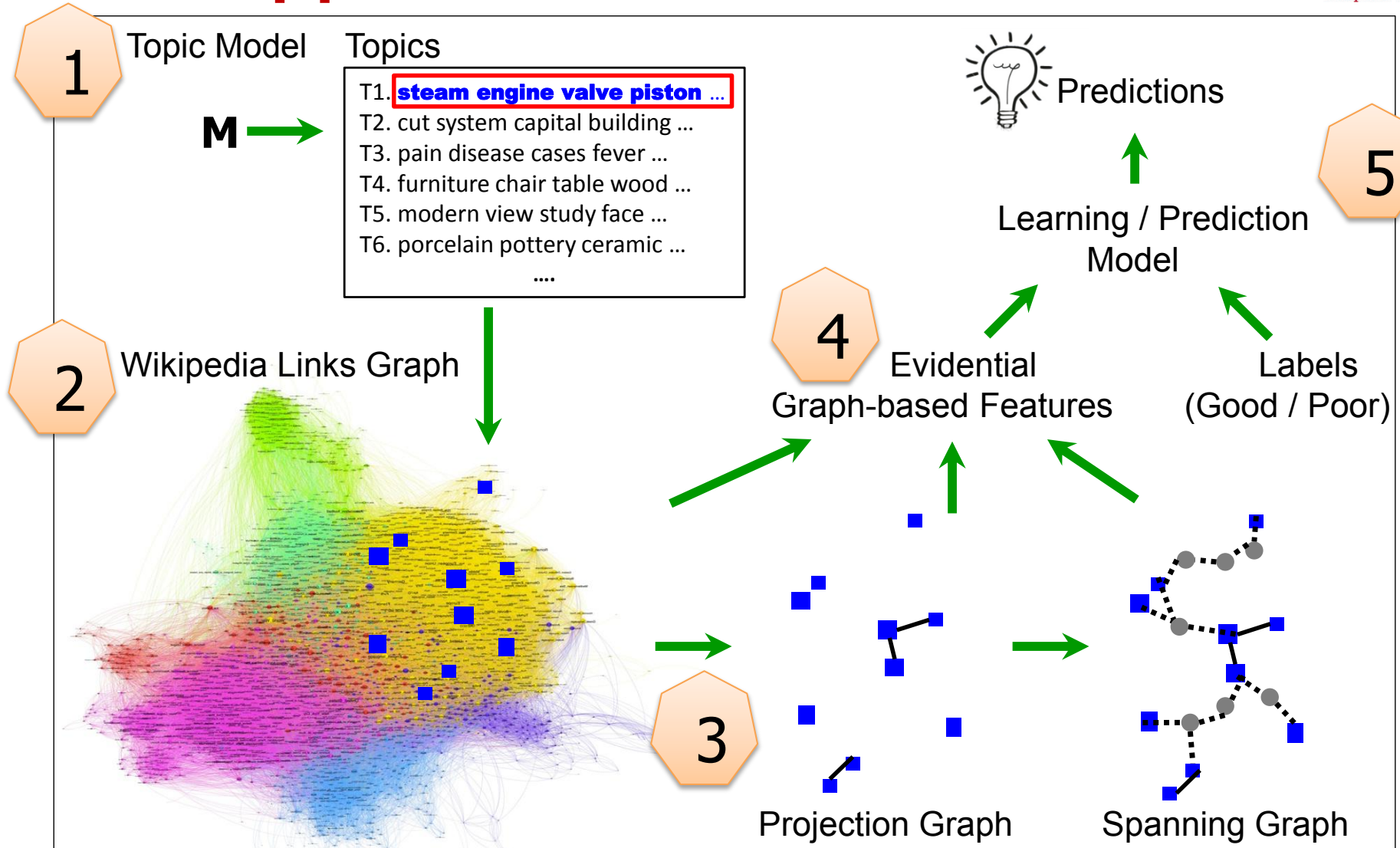
T topics; W words; **Top 10** vs **Top X-Y**

$w_1, w_2, \dots, w_{10}, \dots, w_{11}, \dots, w_{20}, \dots, w_{31}, \dots, w_{40}, \dots, w_{91}, \dots, w_{100}$

Absolute Quality Prediction (Labels – Human Annotators)

*We thank David Newman and his group for sharing these data

Our Approach





Relative Quality Prediction

 w_1, w_2, \dots, w_{10} w_{11}, \dots, w_{20} w_{31}, \dots, w_{40} w_{91}, \dots, w_{100}

Feature set	top-10 vs.	top-[11-20]		top-[31-40]		top-[91-100]	
		PRESS	BRAIN	PRESS	BRAIN	PRESS	BRAIN
BASELINE-MAJORITY		0.500	0.500	0.500	0.500	0.500	0.500
PROJ		0.505	0.622	0.715	0.705	0.765	0.725
D-SPAN		0.650	0.687	0.760	0.740	0.805	0.762
D-SP		0.605	0.665	0.710	0.760	0.750	0.790
PROJ+D-SPAN		0.650	0.687	0.745	0.722	0.790	0.777
PROJ+D-SP		0.650	0.672	0.710	0.752	0.815	0.800
PROJ+D-SPAN+D-SP		0.660	0.687	0.735	0.752	0.810	0.807

>15%

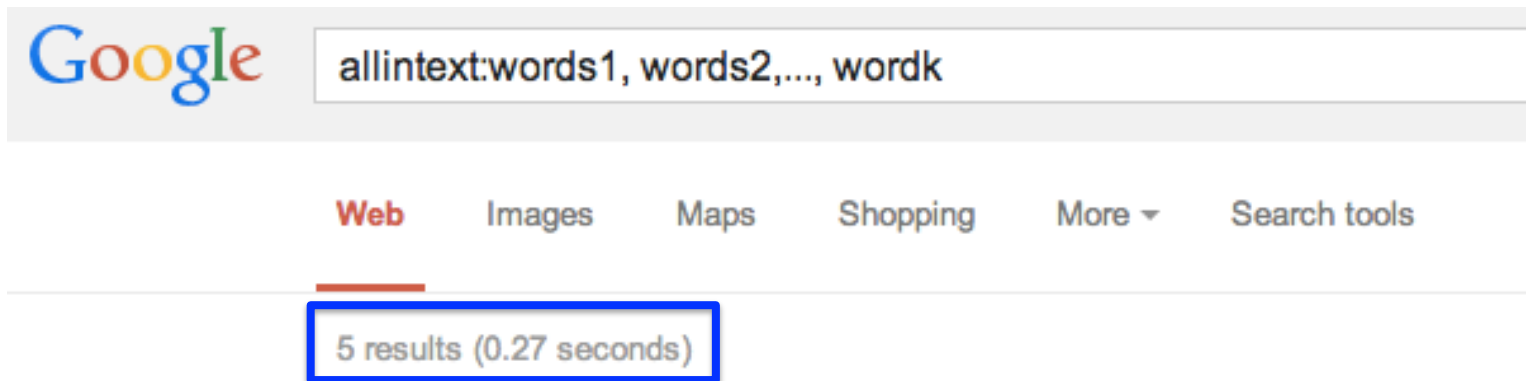
>23%

>30%

improvement over random baseline

Consistent higher accuracy for easier relative tasks

Absolute Quality Prediction – Baselines



Baseline (1)

Google Features:

- (1) allintext
- (2) intitle
- (3) inanchor
- (4) inurl

Baseline (2)

Personalize PageRank:

- (1) average pairwise score
- (2) median pairwise score
- (3) average pairwise rank
- (4) median pairwise rank

Absolute Quality Prediction

Feature set	BOOKS	NEWS	BOOKS +NEWS
BASELINE-MAJORITY	0.610	0.521	0.549
BASELINE-GOOGLE	0.642	0.624	0.629
BASELINE-PPR	0.842	0.735	0.785
PROJ	0.875	0.812	0.848
D-SPAN	0.892	0.769	0.844
D-SP	0.883	0.786	0.852
PROJ+D-SPAN	0.883	0.795	0.844
PROJ+D-SP	0.892	0.795	0.848
PROJ+D-SPAN+D-SP	0.900	0.821	0.831

6% – 9%
improvement over ALL baselines

Cross-Domain Prediction

Train \ Test	BOOKS	NEWS
	BOOKS	NEWS
BOOKS	0.900	0.769
NEWS	0.867	0.821

Our **graph-centric features** are **domain-independent**
(only based on “graph closeness”)

Learned Coefficients

Selected Feature	Coef: BOOKS	Coef: NEWS
<i>g_MNumMiss</i>	0.0626	0.0918
<i>g_SRatioC</i>	0.2940	0.5909
<i>g_MMaxDeg</i>	-0.2921	-0.4541
<i>g_MSizeMaxComp</i>	-0.8667	— — — —
<i>g_SAvgMSTWeight</i>	— — — —	0.2598
<i>NumSP2</i>	-0.9685	— — — —

Good topics:

- Fewer missing (matched) words on WikiLinks
- Fewer connector nodes (in spanning graph)
- Higher maximum degree (in projection graph)

Thank you!

- SBU Office of the Vice President for Research
- NSF Graduate Research Fellowship
- ICDM 2013 Travel Grants

hauchan@cs.stonybrook.edu

<http://www.cs.sunysb.edu/~hauchan>



ICDM 2013

IEEE International Conference on Data Mining
Dallas, Texas / December 7-10, 2013

4. Prediction

- For **relative quality prediction**, we consider topics with 20, 40, and 100 words (ordered by descending probability of describing the topics) where the topics are generated by LDA
 - Positive Class: top 10 words for each topics
 - Negative Class: Bottom 11-20, Bottom 31-40, and Bottom 91-100 words for each topics
- For **absolute quality prediction**
 - Positive Class: Good topics
 - Negative Class: Bad topics

5. Learning Model

- Logistics Regression Classifier
 - L1-Norm Regularization
- We report Leave-One-Out Cross-Validation (LOOCV)