# Reading Tea Leaves: How Humans Interpret Topic Models

# Topic Models

▷ **Used to identify the main themes in a collection of documents.**

▷ **Documents are a collection of topics. Topics are a distribution over words.**

**TOPIC 1**
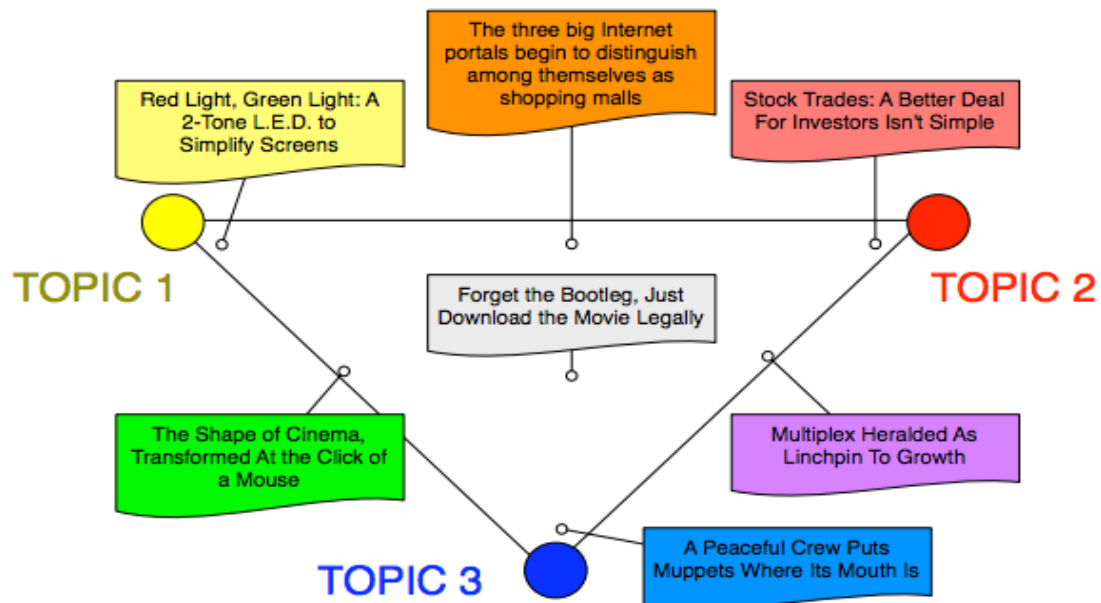computer, technology, system, service, site, phone, internet, machine

**TOPIC 2**
sell, sale, store, product, business, advertising, market, consumer

**TOPIC 3**
play, film, movie, theater, production, star, director, stage

(a) Topics

Red Light, Green Light: A 2-Tone L.E.D. to Simplify Screens

The three big Internet portals begin to distinguish among themselves as shopping malls

Stock Trades: A Better Deal For Investors Isn't Simple

**TOPIC 1**

Forget the Bootleg, Just Download the Movie Legally

**TOPIC 2**

The Shape of Cinema, Transformed At the Click of a Mouse

Multiplex Heralded As Linchpin To Growth
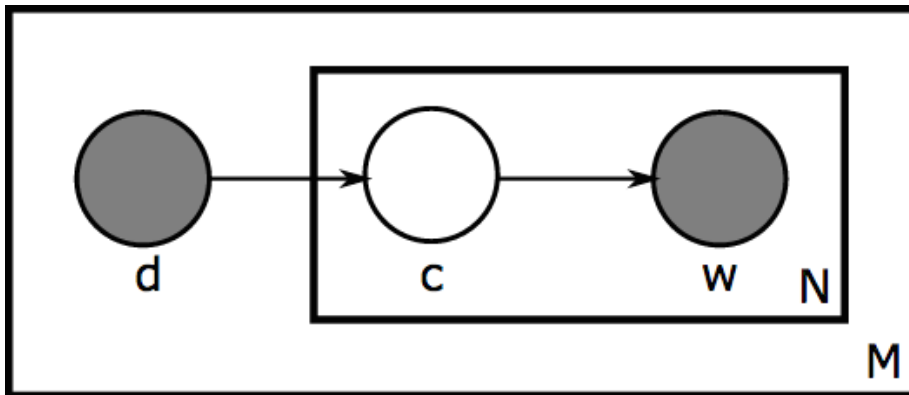
**TOPIC 3**

A Peaceful Crew Puts Muppets Where Its Mouth Is
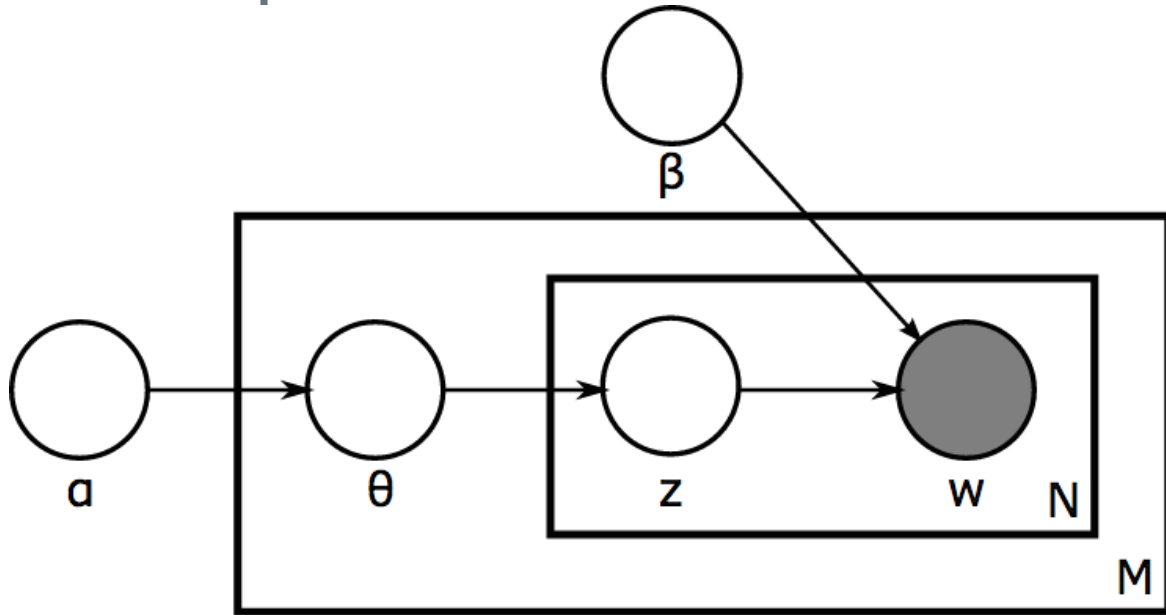
(b) Document Assignments to Topics

# Probabilistic Latent Semantic Indexing (pLSI)

Probability of each co-occurrence is modelled as a mixture of conditionally independent multinomial distributions
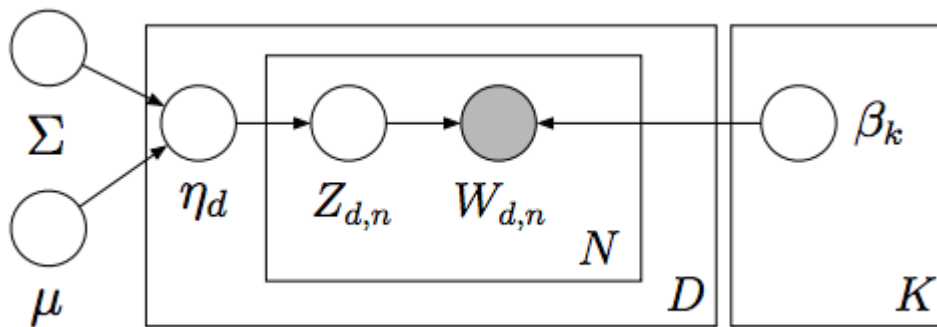
# Latent Dirichlet Analysis (LDA)

**Topic distribution assumed to have a Dirichlet prior.**

# Correlated Topic Model (CTM)

Allows for richer covariance structure between topic proportions. Uses a logistic normal prior over topic mixture proportions.

**(a) Topics**

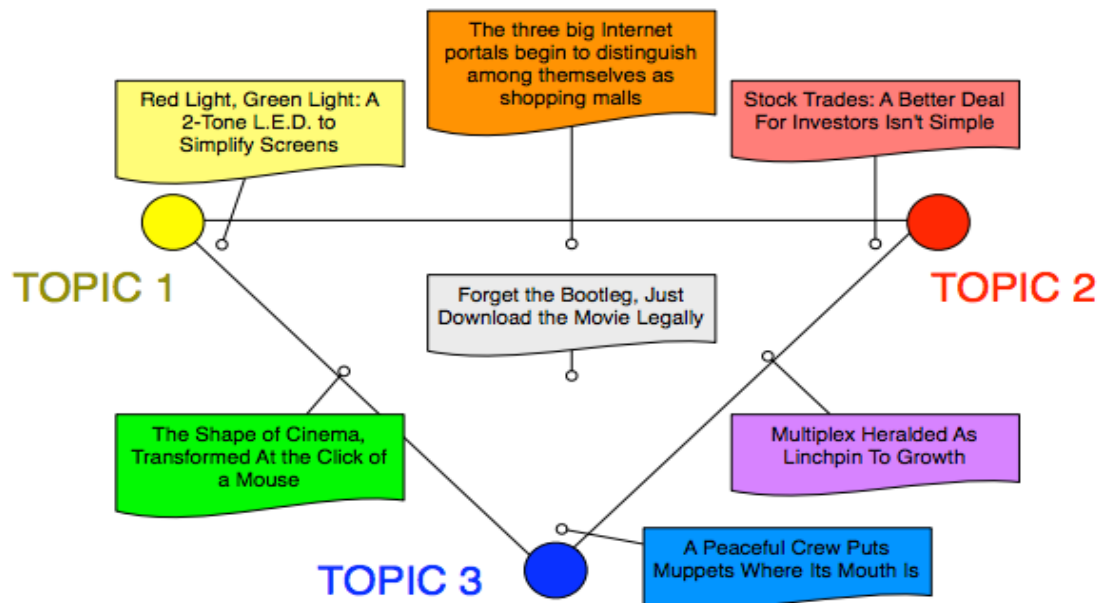TOPIC 1: computer, technology, system, service, site, phone, internet, machine

TOPIC 2: sell, sale, store, product, business, advertising, market, consumer

TOPIC 3: play, film, movie, theater, production, star, director, stage

**(b) Document Assignments to Topics**

Red Light, Green Light: A 2-Tone L.E.D. to Simplify Screens

The three big Internet portals begin to distinguish among themselves as shopping malls

Stock Trades: A Better Deal For Investors Isn't Simple

TOPIC 1

TOPIC 2

Forget the Bootleg, Just Download the Movie Legally

The Shape of Cinema, Transformed At the Click of a Mouse

Multiplex Heralded As Linchpin To Growth

TOPIC 3

A Peaceful Crew Puts Muppets Where Its Mouth Is

# Goals & Motivation

▷ Previously, no measure of interpretability of this latent space.

▷ Present a method for measuring interpretability of topic models using human evaluation tasks.

# Word Intrusion

**1 / 10**

floppy   alphabet   computer   processor   memory   disk

**2 / 10**

molecule   education   study   university   school   student

**3 / 10**

linguistics   actor   film   comedy   director   movie

**4 / 10**

islands   island   bird   coast   portuguese   mainland
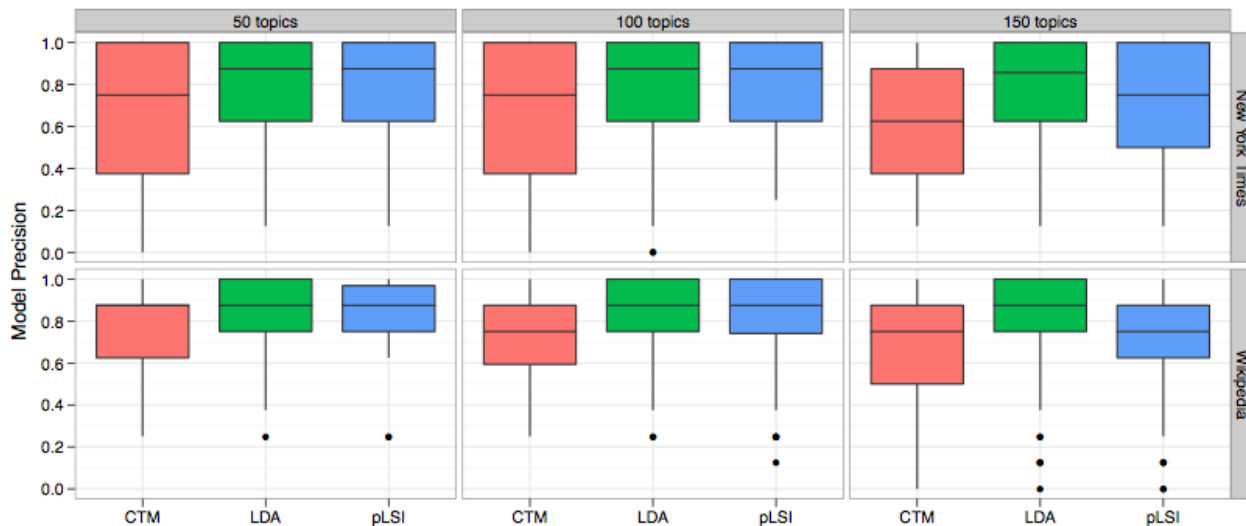
# Topic Intrusion

## DOUGLAS_HOFSTADTER

Douglas Richard Hofstadter (born February 15, 1945 in New York, New York) is an American academic whose research focuses on consciousness, thinking and creativity. He is best known for ", first published in
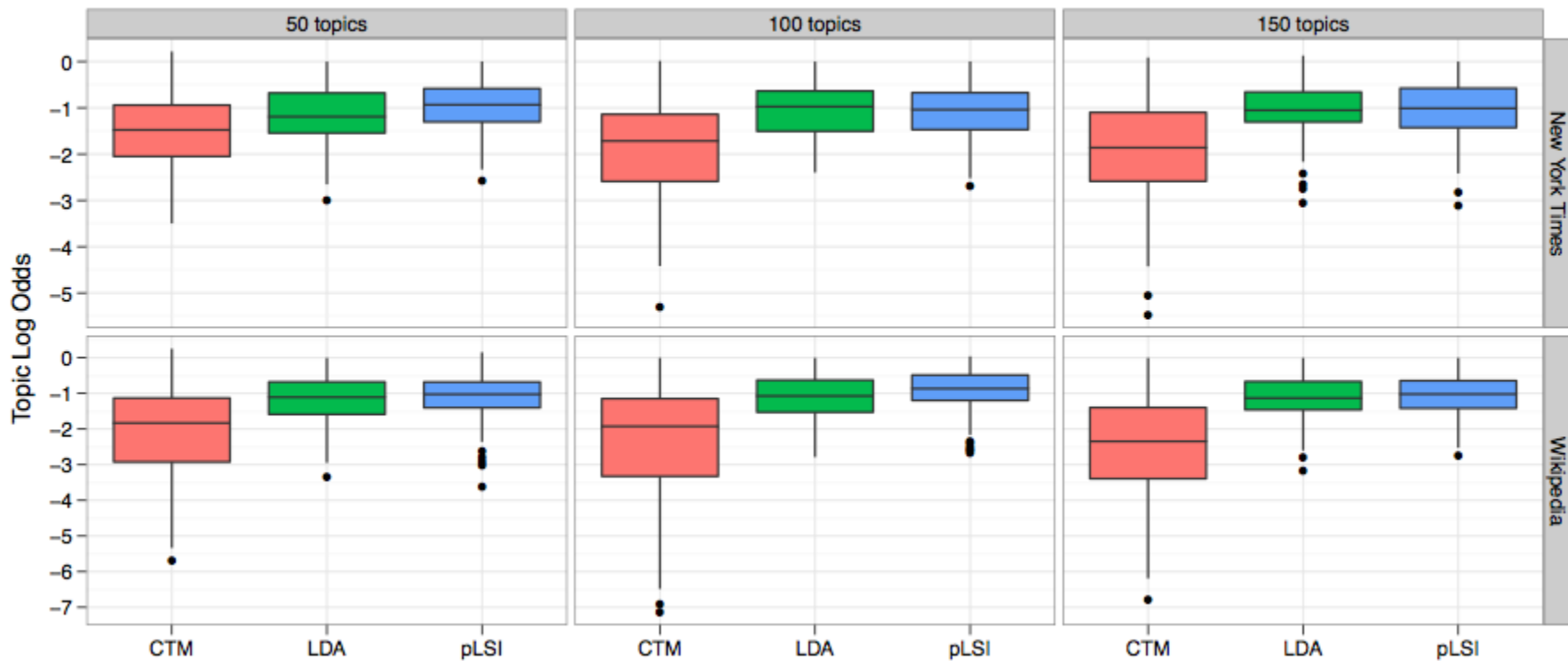
Show entire excerpt

| student | school | study | education | research | university | science | learn |
|---------|--------|-------|-----------|----------|------------|---------|-------|
| human | life | scientific | science | scientist | experiment | work | idea |
| play | role | good | actor | star | career | show | performance |
| write | work | book | publish | life | friend | influence | father |

# Results - Word Intrusion

| CORPUS | TOPICS | LDA | CTM | pLSI |
|--------|--------|-----|-----|------|
| NEW YORK TIMES | 50 | **-7.3214 / 784.38** | -7.3335 / 788.58 | -7.3384 / 796.43 |
| | 100 | -7.2761 / 778.24 | **-7.2647 / 762.16** | -7.2834 / 785.05 |
| | 150 | -7.2477 / 777.32 | -7.2467 / **755.55** | **-7.2382** / 770.36 |
| WIKIPEDIA | 50 | **-7.5257** / 961.86 | -7.5332 / **936.58** | -7.5378 / 975.88 |
| | 100 | -7.4629 / 935.53 | **-7.4385 / 880.30** | -7.4748 / 951.78 |
| | 150 | -7.4266 / 929.76 | **-7.3872 / 852.46** | -7.4355 / 945.29 |

# Results - Topic Intrusion

# Conclusion

▷ Traditional metrics of evaluation do not capture whether topics are coherent.

▷ When developing topic models we should now focus on evaluations which depend on real-world task performance.

# References

Blei, D., & Lafferty, J. (2006). Correlated topic models. *Advances in neural information processing systems*, *18*, 147. (CTM Figure)

Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent dirichlet allocation. *the Journal of machine Learning research*, *3*, 993-1022. (LDA Figure)

Chang, J., Gerrish, S., Wang, C., Boyd-graber, J. L., & Blei, D. M. (2009). Reading tea leaves: How humans interpret topic models. In *Advances in neural information processing systems* (pp. 288-296).

"Plsi 1" by Bkkbrad, EduardoValle - http://en.wikipedia.org/wiki/File:Plsi.svg. Licensed under CC BY-SA 3.0 via Wikimedia Commons

# Further Reading

Blei, D. M. (2012). Probabilistic topic models. *Communications of the ACM,55*(4), 77-84.