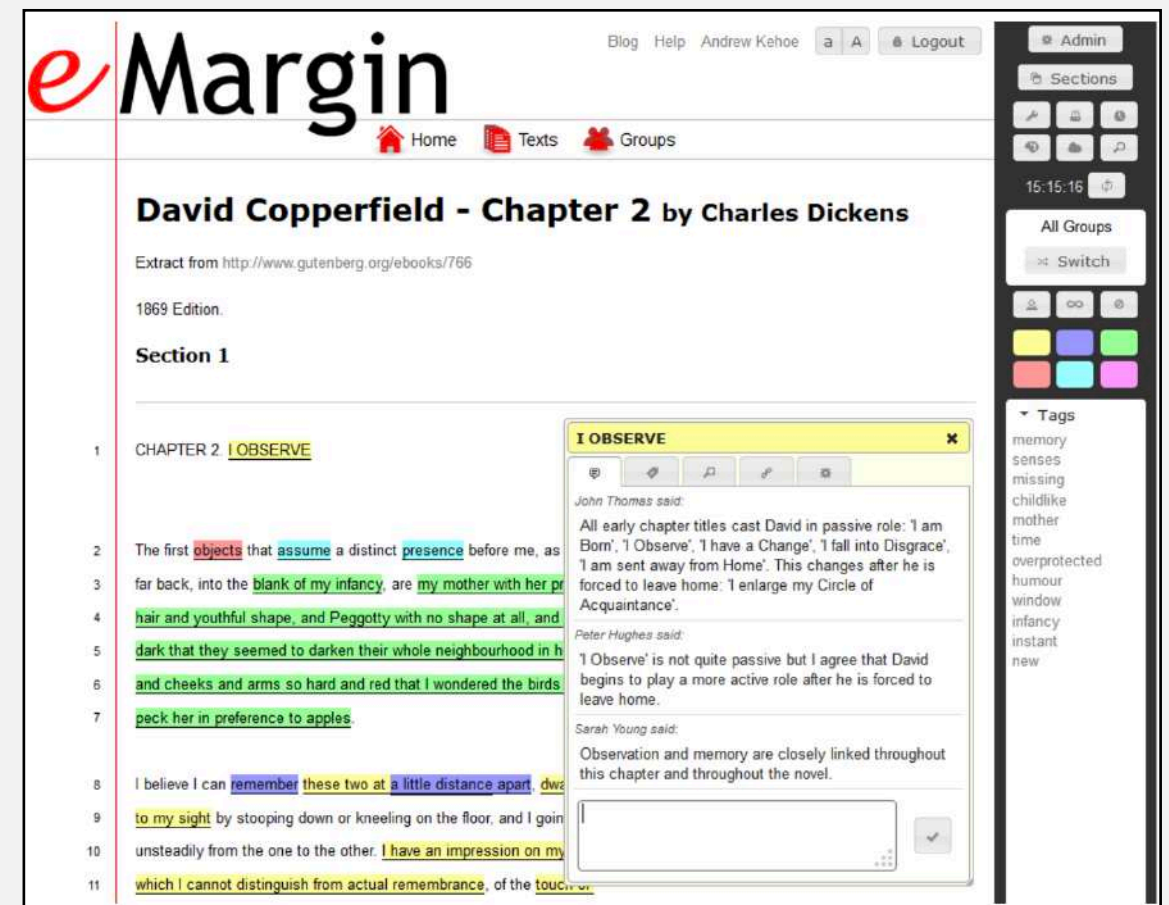# Visually Exploring Documents: Topic Modeling

# Working with text

- "Text visualization" can mean many things

  - fonts, sizes, colors, kerning, typography in general

- We should consider text visualization from the more general definition of visualization: *amplify cognition*

  - Help people improve comprehension of a piece of text

  - Understand themes in text corpora

  - Locate/search relevant textual documents

# Tasks in Text Visualization

- The way we approach text visualization is highly dependent on tasks.

- Close reading: deeply comprehending text, going beyond the words on the page **[Jänicke et al. 2015]**

# Close Reading Designs



**Color**            [Alexander et al. 2014]

**Size**    [Walsh et al. 2014]

**Node-link Diagrams**

[Coles et al. 2014]

# Distant Reading

- Understand structure, relationships, themes, connections within a document.



**[Posavec 2007]**

# Document Exploration

- A "document": sequence of words.

  - A book, a wikipedia article, a paper abstract, a tweet

- Can view as generic high-dimensional data, apply techniques we've discussed so far for visual exploration

- But documents are unique:

  - Really high-dimensional (e.g. dimensionality = size of English vocabulary)

  - Usually sparse

  - Yet, each dimension - word - is interpretable

# Document Exploration Tasks

- Suppose you were provided hundreds of articles relevant to your respective research backgrounds.

- What tasks are relevant to you for gaining an understanding - and advancing your research - on the document corpus?

# Exploring Documents: Dimensionality Reduction

- We represent each document as a point in a high-dimensional space.

  - Each dimension is a word

  - The value of a dimension is the word count: the number of times the word appeared in the document

    **[Anderson et al. 2014]**



**Limitations with this approach?**

# Matrix Factorization

- Dimensionality reduction is … too reductive. We seek better **representations** of documents.

- Let's consider **matrix factorization**.

$$W \in \mathbb{R}^{n \times d} \longrightarrow \|W - \tilde{U}\tilde{V}^T\| \qquad \tilde{U} \in \mathbb{R}^{n \times k} \quad \tilde{V} \in \mathbb{R}^{d \times k}$$

$\mathbf{w}_i \approx \tilde{\mathbf{u}}_i \tilde{V}^T$ a reconstruction of word counts for document $i$

$\tilde{\mathbf{u}}_i$ assigns importance to each column of $\tilde{V}$

$\tilde{V}$ each column vector: a set of weights, one for each word $\longrightarrow$ **shared across documents**

$k$ determines approximation quality

$k \ll \min(n, d)$ we share limited information across documents for reconstruction

**Ideally:** columns represent meaningful, latent factors shared across documents

# Optimistic, but unrealistic, example

- We set *k* to 2

$$\tilde{V} = [\tilde{\mathbf{v}}_1 \, \tilde{\mathbf{v}}_2] \qquad \tilde{\mathbf{v}}_1 \in \mathbb{R}^d \qquad \tilde{\mathbf{v}}_2 \in \mathbb{R}^d$$

- Then we can represent each document with 2 numbers

$$\tilde{\mathbf{w}}_i = [w_{i1}, w_{i2}]$$

$$\tilde{\mathbf{w}}_i = [1,0] \quad \text{take on words from } \tilde{\mathbf{v}}_1$$

$$\tilde{\mathbf{w}}_i = [0,1] \quad \text{take on words from } \tilde{\mathbf{v}}_2$$

- Does this look familiar? What are we doing with documents in this manner?

- Clustering! (kind of) 

$\tilde{\mathbf{v}}_1 \in \mathbb{R}^d$ **set of words that describe one cluster**

$\tilde{\mathbf{v}}_2 \in \mathbb{R}^d$ **set of words that describe other cluster**

# Matrix Factorization: SVD

- Clustering interpretation only meaningful if the reconstruction is good!

$$\| W - \tilde{U}\tilde{V}^{\mathrm{T}} \|_F$$

- We can find the global minimum of this energy via the singular value decomposition (SVD)

$$W = U\Lambda V^{T} \, , \, U \in \mathbb{R}^{n \times n} \, , \, \Lambda \in \mathbb{R}^{n \times d} \, , \, V \in \mathbb{R}^{d \times d} \qquad U^{T}U = I \, , \, V^{T}V = I$$

- Truncate the SVD to the top k singular values

$$\tilde{U} = U_{1:k}\sqrt{\Lambda_{1:k}} \qquad \tilde{V} = V_{1:k}\sqrt{\Lambda_{1:k}}$$

- Can write approximation as the following expansion:

$$W \approx \sum_{i=1}^{k} \lambda_i \mathbf{u}_i \mathbf{v}_i^{T} \quad \textbf{best rank-k approximation}$$

# Nonnegative Matrix Factorization

- Limitation: the document and latent factors can be negative - not easily interpretable! Bad for visualization.

- So, then, let's enforce nonnegativity!

$$\|W - UV^T\|_F^2 \, , \quad s.t. \; U, V \geq 0$$

- We can then say "latent factors $V$ contribute $u$ amount to the reconstruction"

- Challenge: fixing $U$, energy is convex in $V$. Fixing $V$, energy is convex in $U$. But not convex in both! (due to nonnegativity constraint - cannot apply Eckart-Young theorem)

# Algorithm: Multiplicative Updates

$$\|W - UV^T\|_F^2, \quad s.t. \ U, V \geq 0$$

- Alternate between the following updates:

$$V_{ab} \leftarrow V_{ab} \frac{(U^T W)_{ab}}{(U^T U V^T)_{ab}} \qquad U_{ab} \leftarrow U_{ab} \frac{(WV)_{ab}}{(UV^T V)_{ab}}$$

- Update can be seen as a particular (per-element) step size chosen for gradient descent:

$$E(V) = \|W - UV^T\|_F^2 = tr((W - UV^T)^T(W - UV^T))$$

$$= tr(W^T W - 2VU^T W + VU^T UV^T)$$

- Take derivative of trace:

$$\frac{\mathbf{d}E}{\mathbf{d}V^T} = 2(U^T UV^T - U^T W)$$

# Algorithm: Multiplicative Updates continued...

$$V_{ab} \leftarrow V_{ab} - \eta_{ab}(U^T U V^T - U^T W)_{ab} \quad , \quad \eta_{ab} = \frac{V_{ab}}{(U^T U V^T)_{ab}}$$

$$V_{ab} \leftarrow V_{ab} - \frac{V_{ab}}{(U^T U V^T)_{ab}}(U^T U V^T - U^T W)_{ab} = V_{ab} - V_{ab} + V_{ab}\frac{(U^T W)_{ab}}{(U^T U V^T)_{ab}}$$

- Starting from an initial guess for *U* and *V* that are both nonnegative, this scheme ensures:

  - Both will remain nonnegative

  - Energy decreases at each iteration

  - Will arrive at *some* fixed-point solution (local minimum)

**[Lee & Seung 2001]**

# Still, some limitations

- Weights are unbounded:

  - Some latent factors could dominate others.

  - Thus, document weights become hard to interpret.

- At this point, might start introducing regularization terms.

- However, consider the following probabilistic interpretation:

$$w_{ij} \approx \mathbf{u}_i^T \mathbf{v}_j = \sum_{l=1}^{k} u_{il} v_{jl} \longrightarrow \sum_{l} p(z_l \,|\, \theta) p(w \,|\, z_l, \beta)$$

**probability of latent factor, given a document**

**probability of word, given latent factor**

# Latent Dirichlet Allocation

- Next, let's consider *all* words in a document:

$$p(\theta, \mathbf{w} \mid \alpha, \beta) = \underline{p(\theta \mid \alpha)} \prod_n \sum_{z_n} p(z_n \mid \theta) p(w_n \mid z_n, \beta)$$

Probability of latent factors - or **topics**
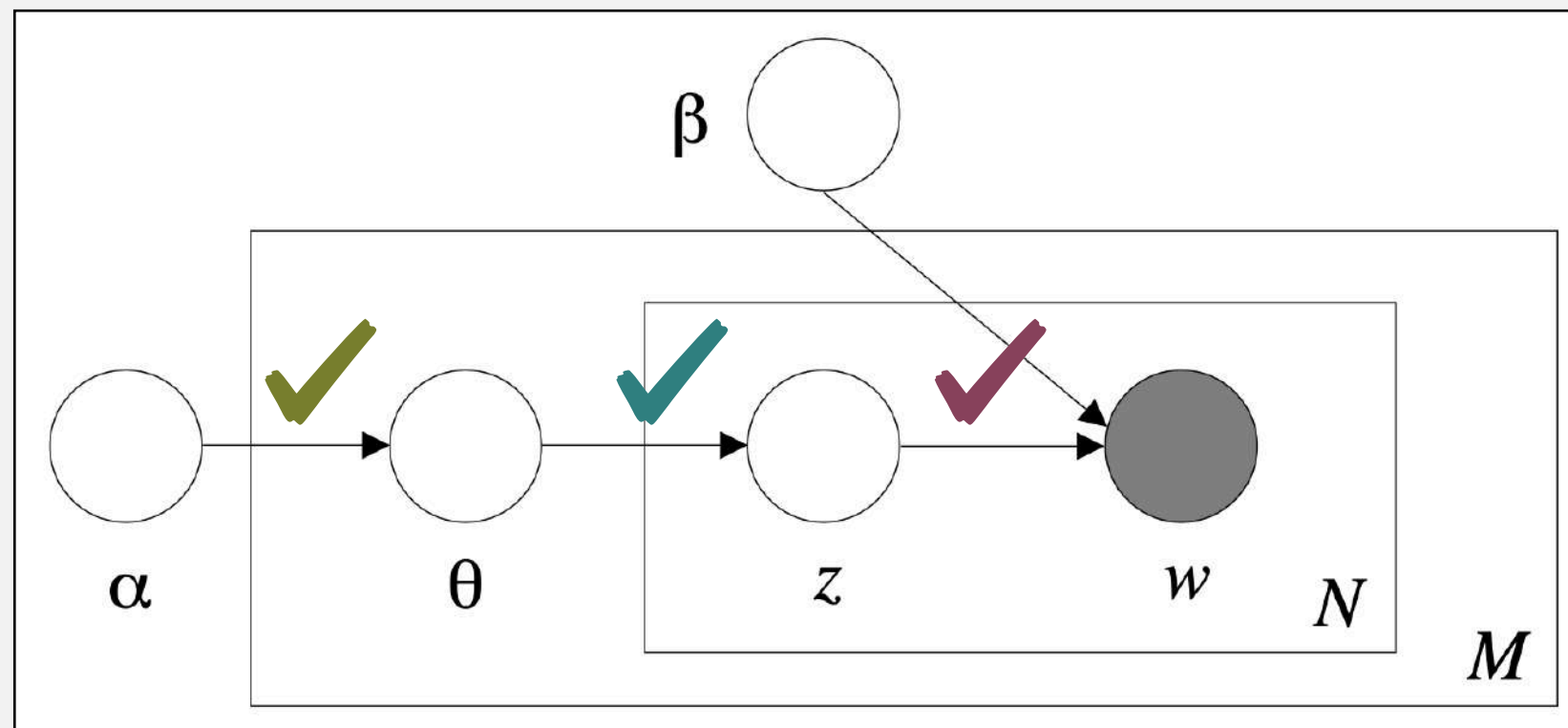
- We marginalize out $\theta$

$$p(\mathbf{w} \mid \alpha, \beta) = \int p(\theta \mid \alpha) \left( \prod_n \sum_{z_n} p(z_n \mid \theta) p(w_n \mid z_n, \beta) \right) d\theta$$

- Last, consider *all* documents:

$$p(\mathbf{w} \mid \alpha, \beta) = \prod_{d=1}^{M} \int p(\theta_d \mid \alpha) \left( \prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn} \mid \theta_d) p(w_{dn} \mid z_{dn}, \beta) \right) d\theta_d$$

# Probabilistic Model

$$p(\mathbf{w} \mid \alpha, \beta) = \prod_{d=1}^{M} \int p(\theta_d \mid \alpha) \left( \prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn} \mid \theta_d) p(w_{dn} \mid z_{dn}, \beta) \right) d\theta_d$$
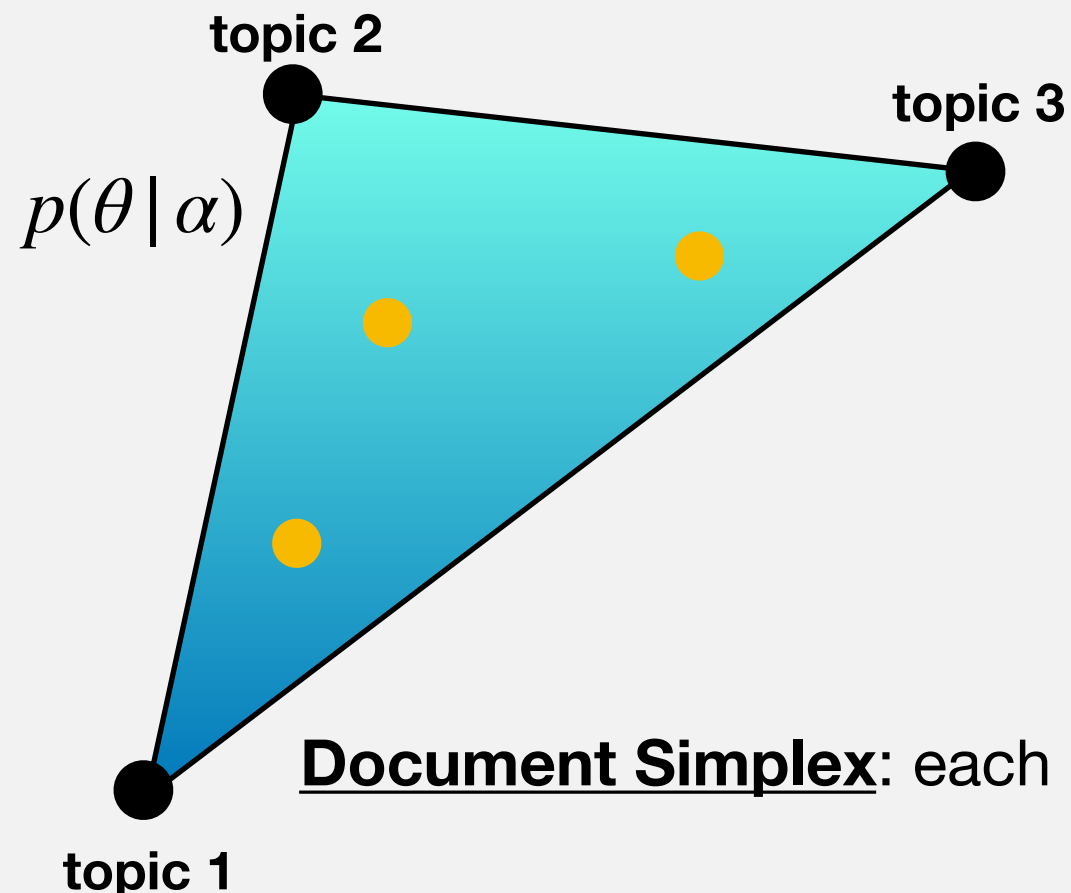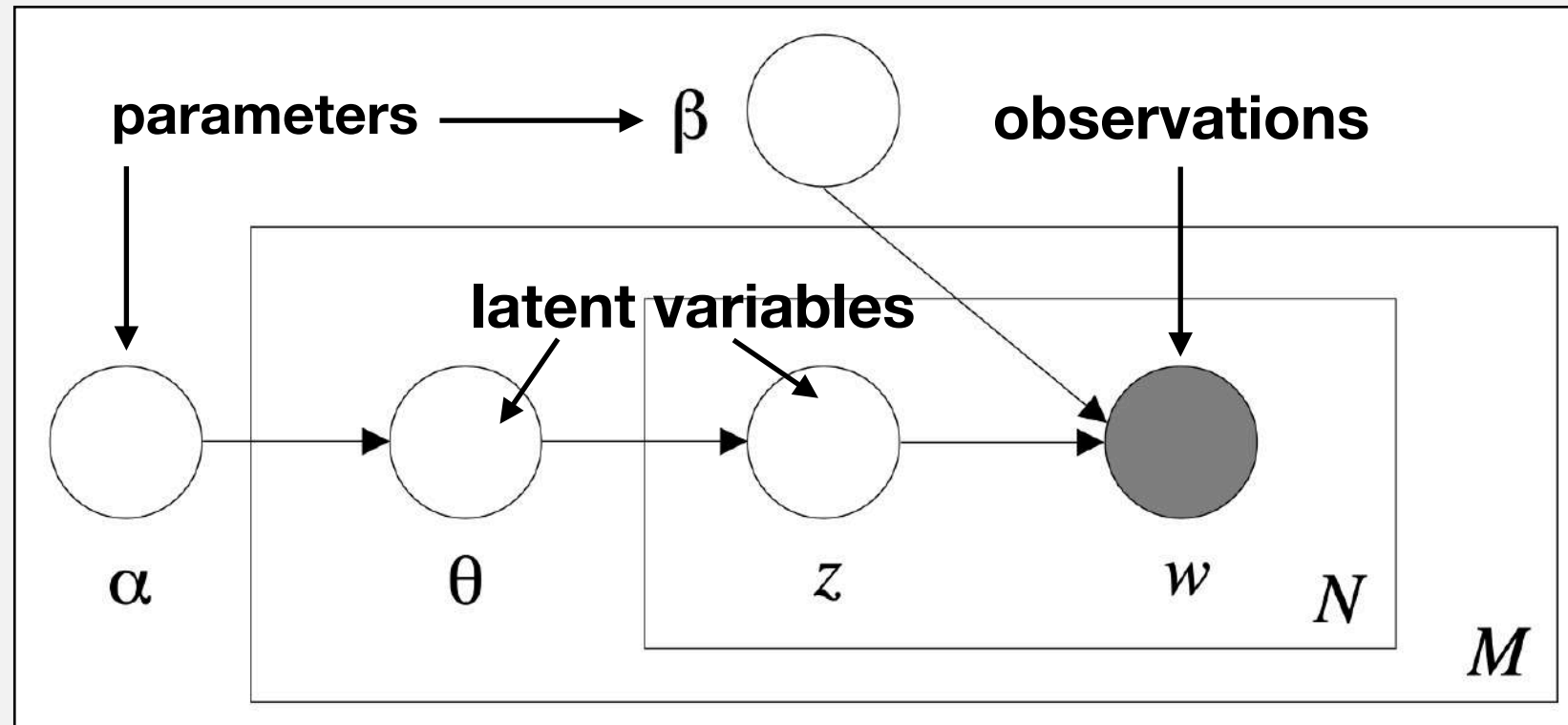


For a given document, draw a mixture of topics

To generate a word, first draw a random topic, given the mixture
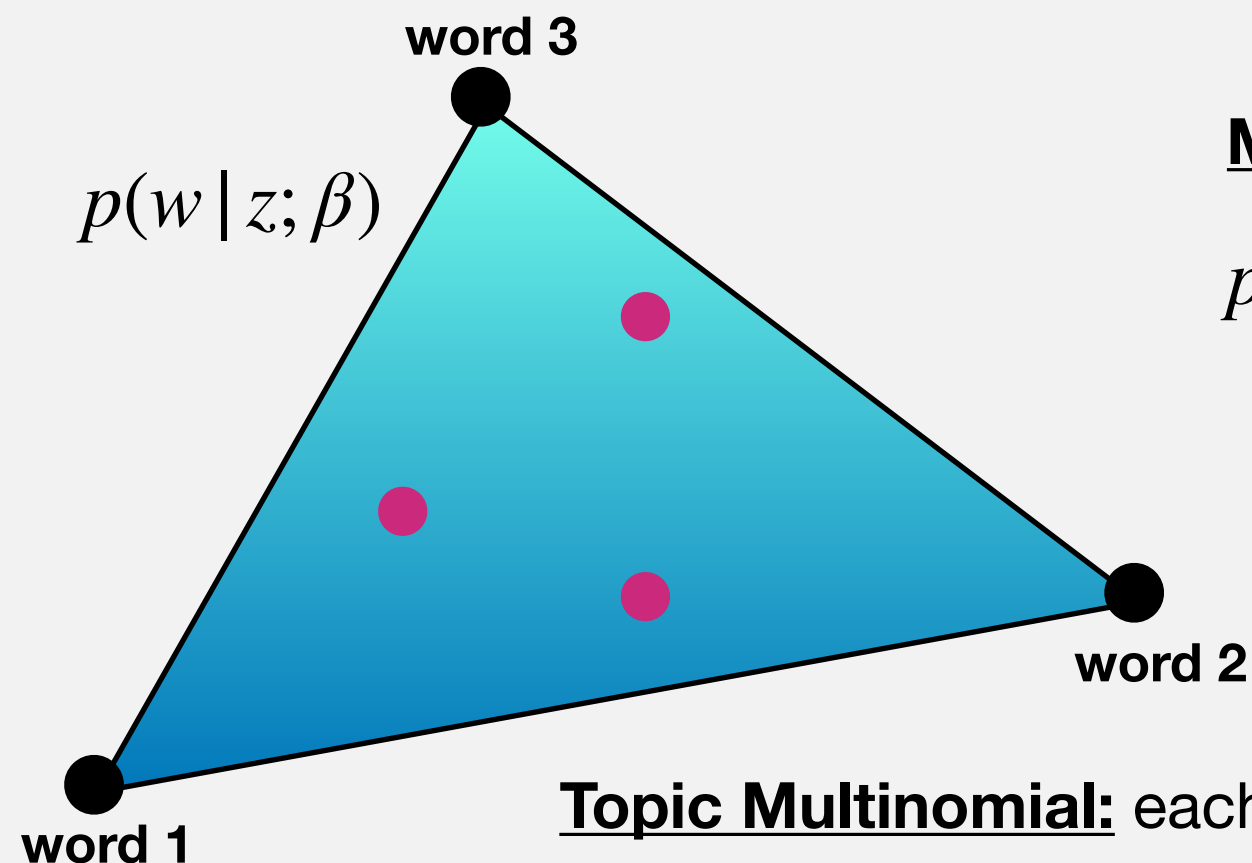
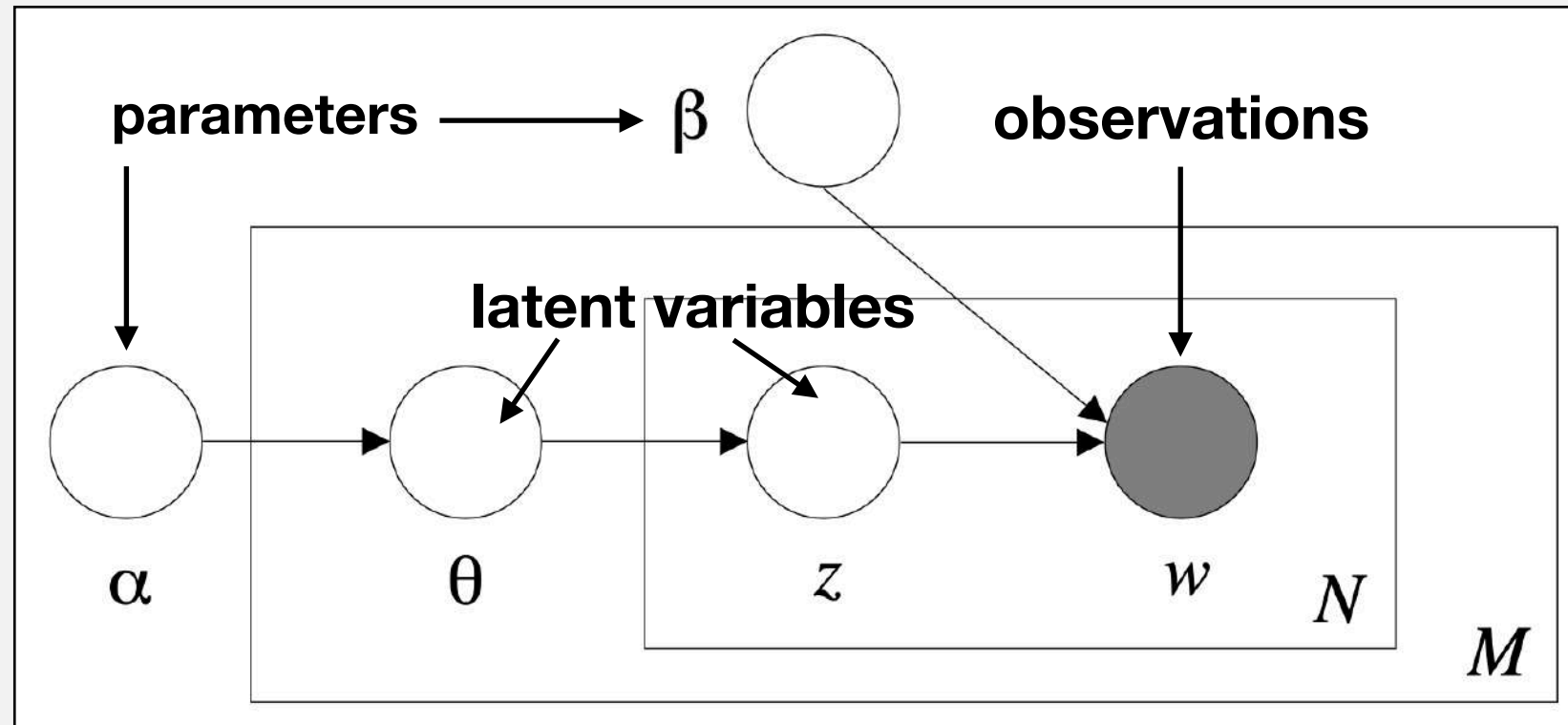Last, sample a word from the topic

# Probabilistic Model



**Dirichlet distribution**

$$p(\theta \mid \alpha) = Z(\alpha) \prod_i \theta_i^{\alpha_i - 1}$$

**Document Simplex**: each point is a mixture over topics

# Probabilistic Model



$$p(w \mid z; \beta)$$

**Multinomial distribution**

$$p(w \mid z; \beta) = Z(\beta_{z_i}) \prod_i w_i^{\beta_{z_i}}$$

**word 3**

**word 2**

**word 1**

**Topic Multinomial:** each point is a distribution over words

# Generative Model

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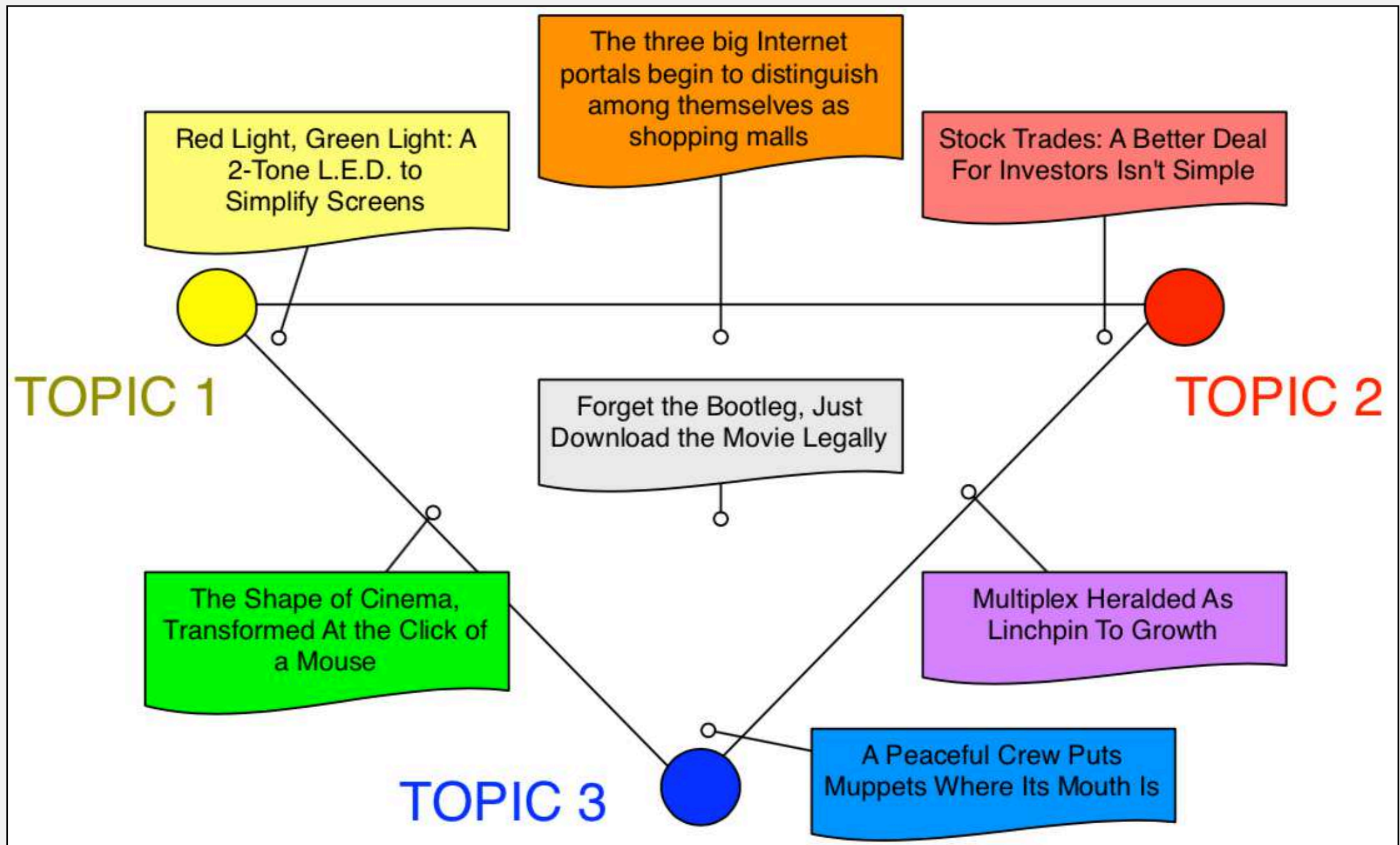computer, technology, system, service, site, phone, internet, machine

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

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

**(slide from David Mimno)**

# Generative Model

# Generative Model

computer, technology, system, service, site, phone, internet, machine

sell, sale, store, product, business, advertising, market, consumer

play, film, movie, theater, production, star, director, stage

Hollywood studios are preparing to let people download and buy electronic copies of movies over the Internet, much as record labels now sell songs for 99 cents through Apple Computer's iTunes music store and other online services ...

**(slide from David Mimno)**

# Generative Model



computer, technology, system, service, site, phone, internet, machine

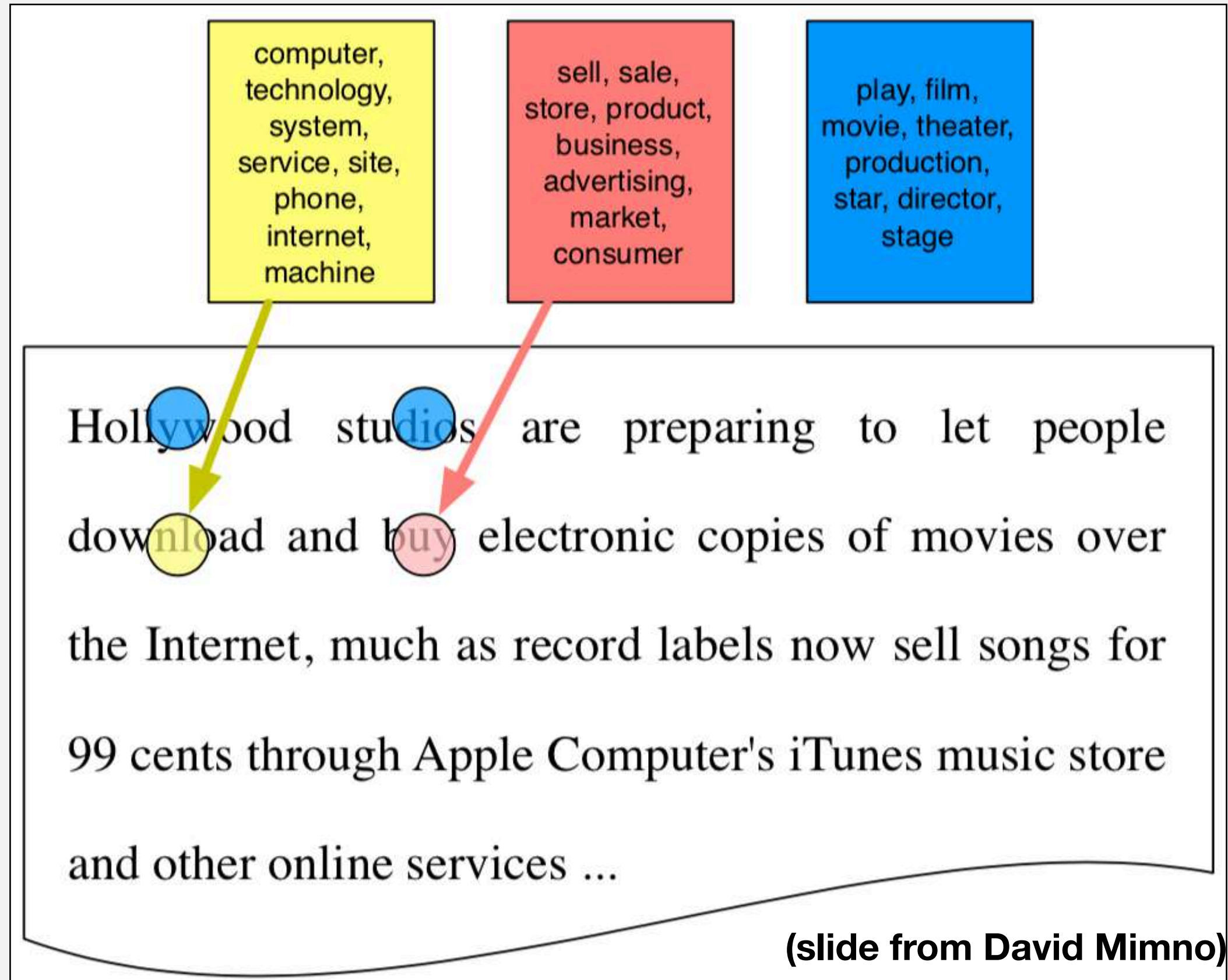sell, sale, store, product, business, advertising, market, consumer

play, film, movie, theater, production, star, director, stage

Hollywood studios are preparing to let people download and buy electronic copies of movies over the Internet, much as record labels now sell songs for 99 cents through Apple Computer's iTunes music store and other online services ...

**(slide from David Mimno)**

# Generative Model

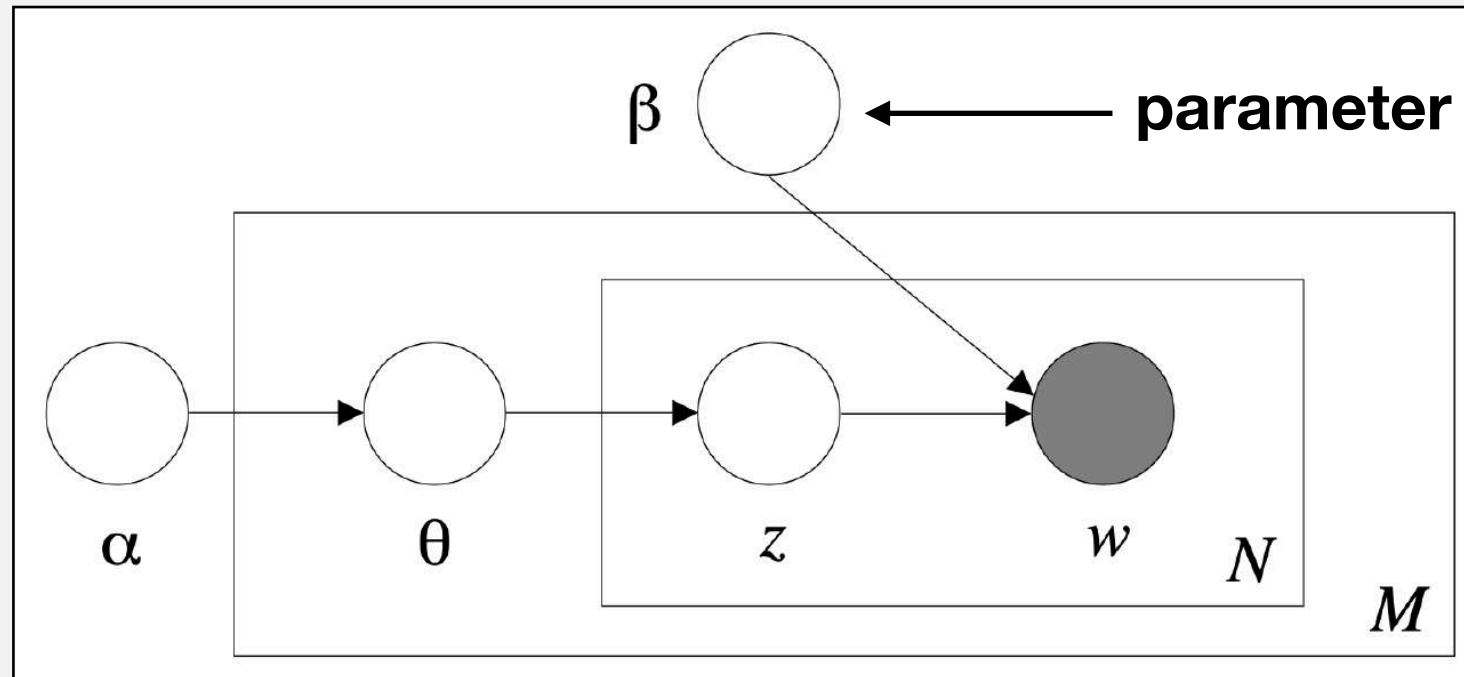computer, technology, system, service, site, phone, internet, machine

sell, sale, store, product, business, advertising, market, consumer

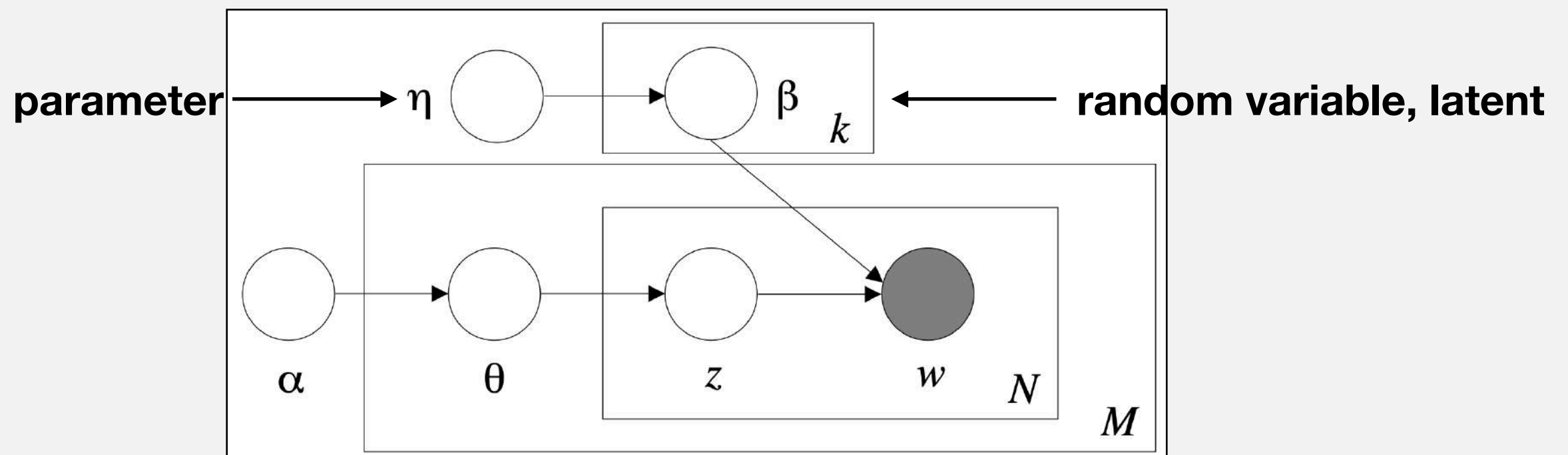play, film, movie, theater, production, star, director, stage

Hollywood studios are preparing to let people download and buy electronic copies of movies over the Internet, much as record labels now sell songs for 99 cents through Apple Computer's iTunes music store and other online services ...

**(slide from David Mimno)**

# Generative Model



**(slide from David Mimno)**

# Problem with Sparsity



- Documents might be similar, but have no words in common!

# Inference in LDA

- Need to compute the posterior distribution of hidden variables - intervals, point estimates, etc..

- However, this is intractable: need to evaluate a really complicated integral

- **Variational Inference**: we approximate the distribution we would like, with one that permits tractable optimization

$$q(\beta, \mathbf{z}, \theta \,|\, \lambda, \phi, \gamma) = \prod_{i=1}^{k} \mathbf{Dir}(\beta_i \,|\, \lambda_i) \prod_{d=1}^{M} \mathbf{Dir}(\theta_d \,|\, \gamma_d) \prod_{n=1}^{N} \mathbf{Mult}(z_n \,|\, \phi_n)$$

**Probability distribution for a topic (over words)**

**Probability distribution for a document (over topics)**

# Relating *q* and *p*

- We maximize the Evidence Lower BOund (ELBO):

$$\log p(\mathbf{w} \,|\, \alpha, \eta) \geq \mathscr{L}(\mathbf{w}, \lambda, \phi, \gamma) = \mathbb{E}_q[\log p(\mathbf{w}, \beta, \mathbf{z}, \theta \,|\, \alpha, \eta)] - \mathbb{E}_q[\log q(\beta, \mathbf{z}, \theta \,|\, \lambda, \phi, \gamma)]$$

- Log-likelihood of document

- Latent variables are shared between distributions

- Differences: Dirichlet and Multinomial parameters

- A consequence of Jensen's inequality: a concave (e.g. *log*) function of an expectation is lower-bound by the expectation of the concave function

# Main Optimization

$$\mathscr{L} = \sum_d \mathbb{E}_q[\log p(\mathbf{w}_d | \beta_d, \mathbf{z}_d, \theta)] + \mathbb{E}_q[\log p(\mathbf{z}_d | \theta_d)] + \mathbb{E}_q[\log p(\theta_d | \alpha)] + \mathbb{E}_q[\log p(\beta | \eta)]$$

$$- \mathbb{E}_q[\log q(\mathbf{z}_d | \phi_d)] - \mathbb{E}_q[\log q(\theta_d | \gamma_d)] - \mathbb{E}_q[\log q(\beta | \lambda)]$$

- Can solve for variational parameters via coordinate ascent:

$$\phi_{dwk} \propto \exp\{\mathbb{E}_q[\log \theta_{dk}] + \mathbb{E}_q[\log \beta_{kw}]\}$$ **topic-relevance**

**if *k* is related to *d and* w, then the topic, document, and word are all related**

$$\gamma_{dk} = \alpha + \sum_w n_{dw}\phi_{dwk}$$

**a topic is relevant to a document, if document-conditioned words are topic-relevant**

$$\lambda_{kw} = \eta + \sum_d n_{dw}\phi_{dwk}$$

**a topic is relevant to a word, if documents that contain the word are topic-relevant**

# What does inference give us?

- We have now estimated our variational parameters.

  - One describes a probability distribution for each topic

  - The other describes a probability distribution for each document

**Topics**

$\mathbf{Dir}(\beta_i \,|\, \lambda_i)$

**Documents**

$\mathbf{Dir}(\theta_d \,|\, \gamma_d)$

- We can draw random samples from each distribution … or, if we just want a single realization, we take expectations:

$$\mathbb{E}[\mathbf{Dir}(\beta_i \,|\, \lambda_i)] = \frac{\lambda_i}{\sum_{j=1}^{n} \lambda_{ij}}$$

$$\mathbb{E}[\mathbf{Dir}(\theta_d \,|\, \gamma_d)] = \frac{\gamma_d}{\sum_{j=1}^{k} \gamma_{dj}}$$

# Tasks in Visualizing Topic Models

- Comparing documents

- Comparing topics

- Understanding a topic

- Understanding a document, *in terms of* topics

- Other data:

  - Time? Document Categories?

**Tasks determine how we prioritize visual encodings and interactions!**

# Visualizing a topic?

- Easy? choose its highest-probability words.

- **Visualizing multiple topics**: some decisions need to be made…



Topic 02 graphics virtual simulation interaction pis visualization surface visual human-machine physical haptic touch imagine realistic interactive force tracking

Topic 07 recognition speech sign musical music signed signing sound speaker computer-based automatic auditory emotional processing channel synthesis communication

Topic 22 language text tagging linguistic natural categorization machine relation processing message meaning nlp corpora extraction sentence translation word training

Topic 23 mining discovery dataset massive machine network scientific detection statistical pattern novel knowledge complex field developing time source

Topic 28 network security privacy response service communication distributed emergency policy collaborative justice wireless criminal released internet private fire sharing

Topic 01 parallel database query relational management processing http performance estimate optimizer spatio-temporal implementation answer operation hardware

Topic 04 model reduction performance dimensionality existing space statistical measure optimization selection approach based novel popular method machine

Topic 13 reasoning planning complex decision theory causal intelligence uncertainty computational domain real-world probabilistic knowledge graphical

Topic 19 undergraduate graduate course education program educational computing engineering university curriculum interdisciplinary project school women underrepresented

Topic 00 creative creativity computational media scratch children interactive reading designer content artist study technology animation collaboration

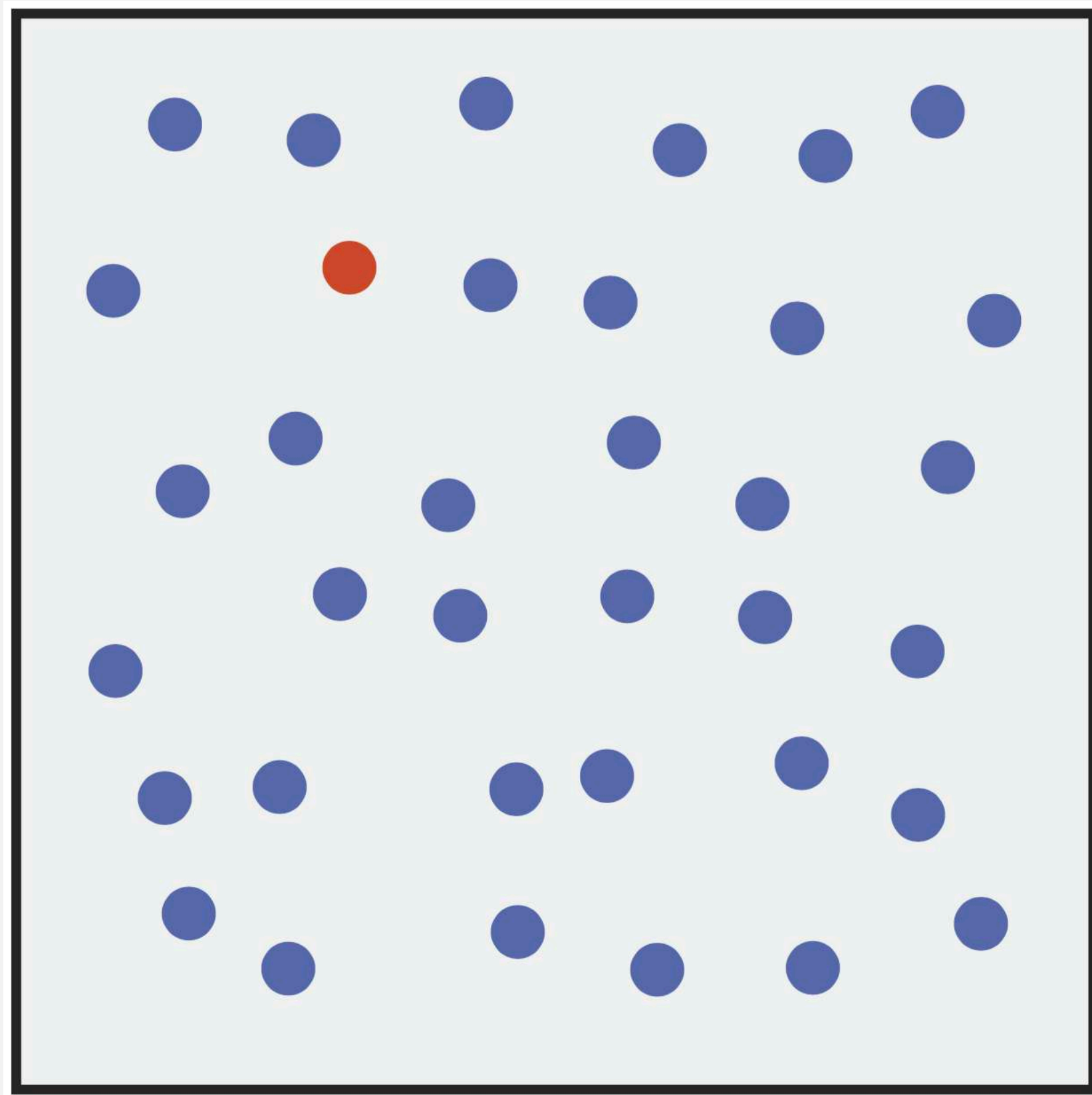Topic 05 user create help potential available goal process people generate ability current solution set building example enable cost knowledge difficult complex

Topic 11 intelligence cognitive agent people human behavior intelligent strategy interaction individual learn environment ability understanding machine

Topic 12 visual computational neuroscience brain cortex stimuli memory response activity understanding mechanism neuronal natural neural movement

Topic 09 biology computational biological network sequencing high-throughput sequence interaction protein gene proposed bioinformatics evolutionary

**[Dou et al. 2011]**

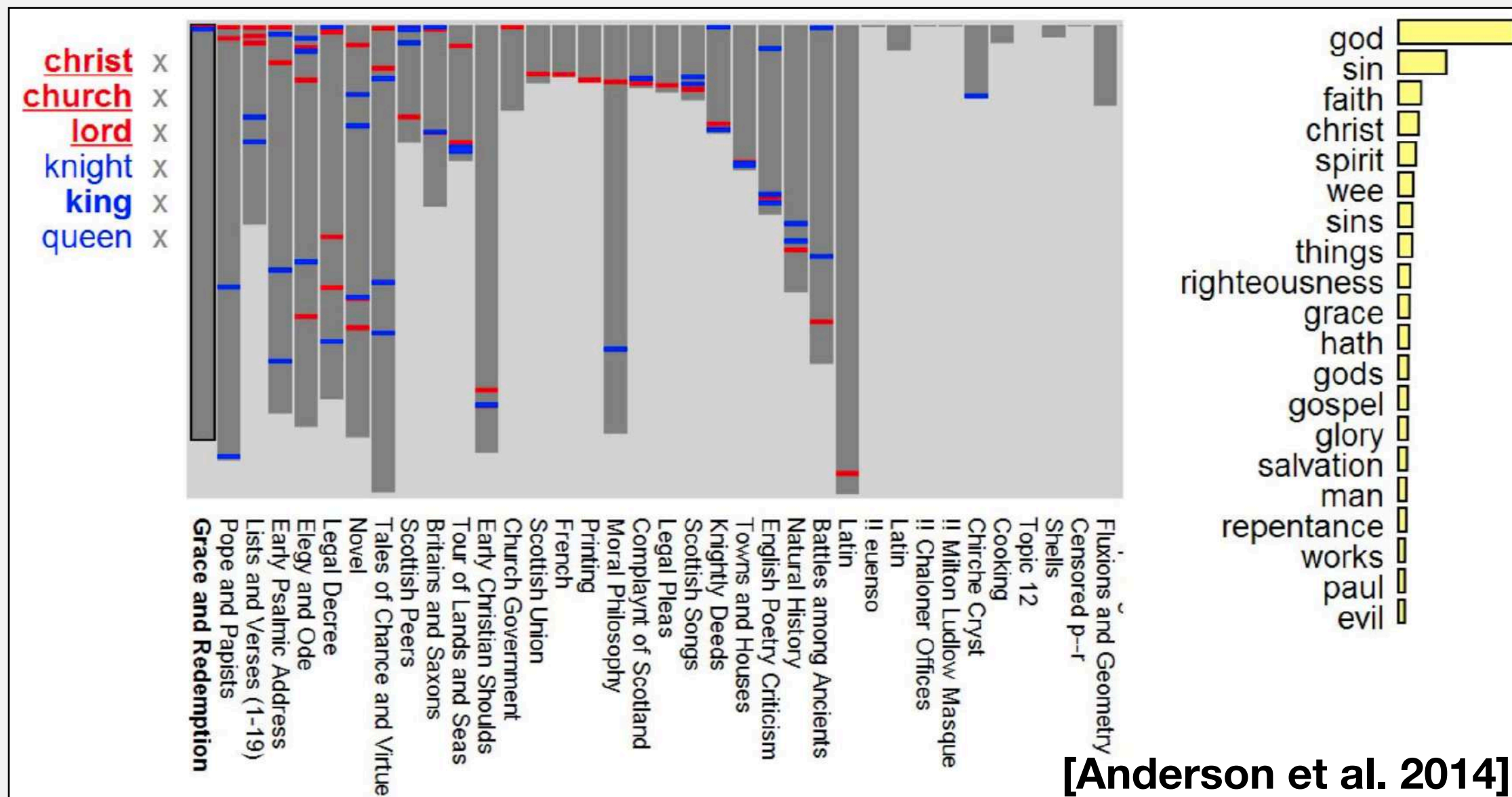# Preattentive Processing

# Implications for Text Visualization

- Showing a list/scatterplot of text, without consideration of visual channels, can result in <u>serial processing</u> - *slow*!

- If possible, use other visual channels to style text.

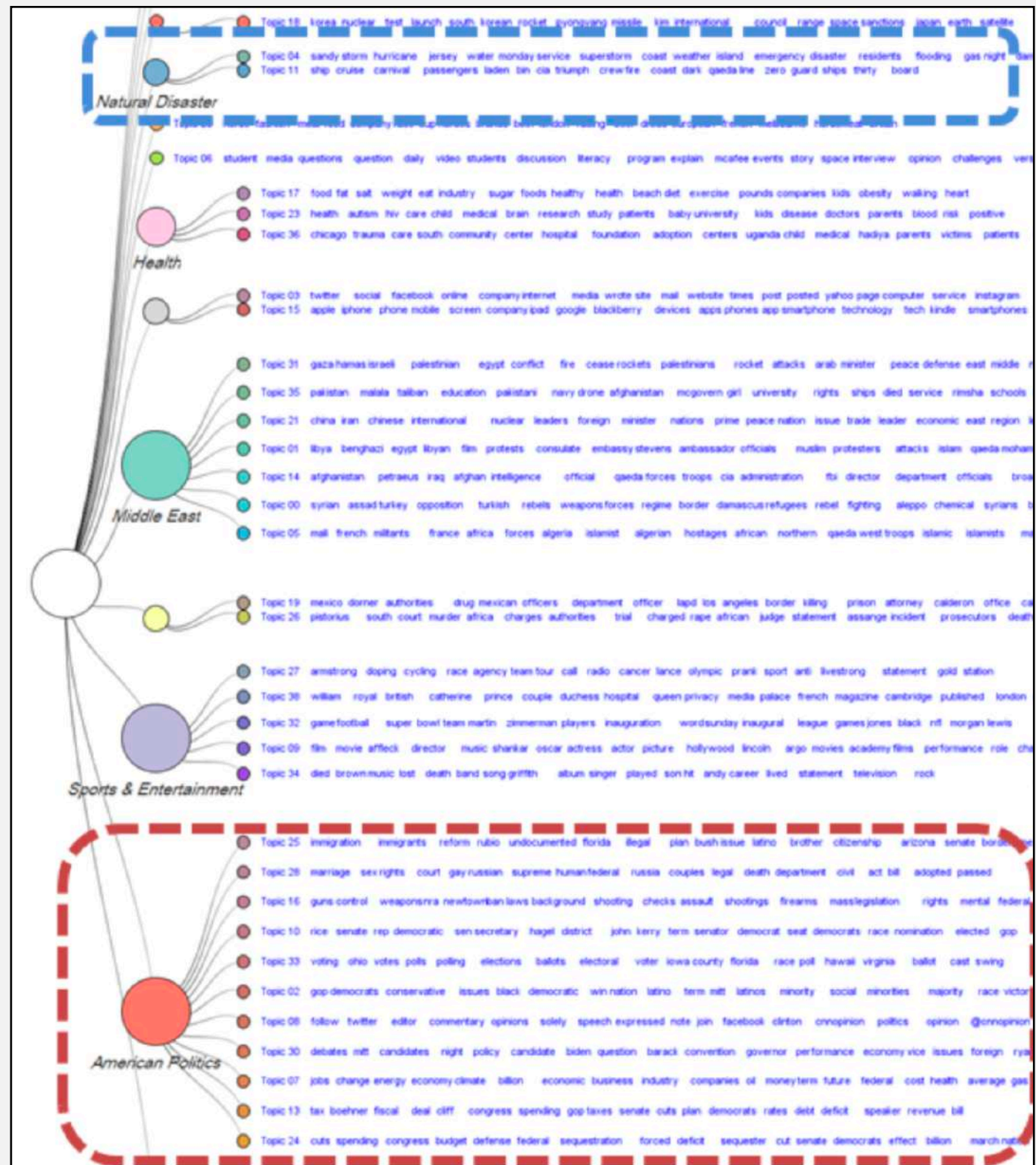- Topic modeling gives us additional information which we may use to visually style text data - **so we should use it**.

# Topic Vis, Pt. 2

- An alternative:



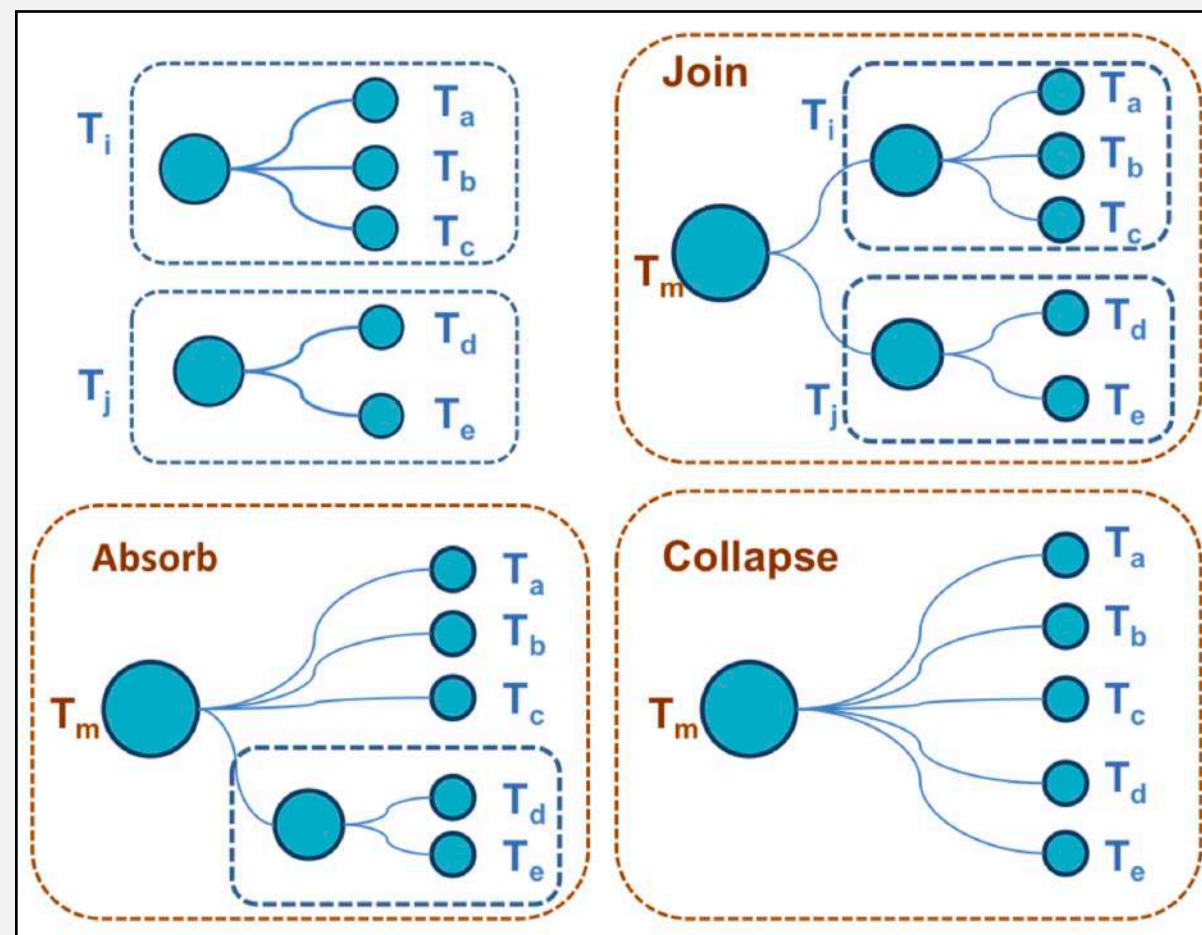[Anderson et al. 2014]

# Topic Vis, Pt. 3

- Lots of topics?

- Hierarchies!

[Dou et al. 2013]

# Building the Hierarchy

- Start from a set of topics, treat each as leaf nodes in a tree, repeat:

  - Consider the following types of operations for a pair of subtrees:

  

# Building the Hierarchy

- Find the operation that gives a pair of subtrees the "lowest cost". Cost? Distance between topics?

$$d_H(t_i, t_j) = \sum_{v=1}^{N} (\sqrt{t_{i,v}} - \sqrt{t_{j,v}})^2$$

- Non-leaf nodes? Average their distributions…

- Sidebar: any potential issues with this distance?

- Algorithm continued: merge the two nodes with lowest cost, repeat until we reach the root!

# Hierarchy: Good?

- Might be imperfect, so allow the user to adjust by exposing these operations:

# Visualizing a document?

- Easy? Show topic assignments.

- Visualizing multiple topics? Parallel coordinates!



**[Dou et al. 2011]**

# Document Vis, Pt. 2

- Poly-lines can quickly become a source of clutter!

- An alternative:

[Anderson et al. 2014]

# Document Vis, Pt. 3

- Order matters!

- An alternative:



**[Dou et al. 2011]**

# Text Vis?

- Show the document directly:



[Dou et al. 2011]

# Text Vis, Pt. 2

- Use topic model to aid in showing text!



**[Anderson et al. 2014]**

# Handling Time

- How do topics vary over time?



**[Dou et al. 2011]**

# What is this time series?

- Need to derive a time-dependent measure of **topic relevance**.

- Method:

  - Fix a time interval unit

  - Take a temporal sliding window over all time steps, where the window length is this interval.

  - For a given temporal window, gather all documents:

$$\tau_D(i) = \sum_{d \in D} d_i$$

# Hierarchical Temporal Evolution

- Top-level nodes:



[Dou et al. 2013]

# Hierarchical Temporal Evolution

- Interactive! Click on a node, expand to its children:

# Topic Modeling for Exploring Conversations

- Conversation data

  - Set of speakers

  - Each speaker delivers an <u>utterance</u>: a set of statements, at a particular time.

  - Treat each utterance as a document: topic modeling!

# ConToVi: Multi-Party Conversation Exploration

- Objective: show the progression of a conversation, in the context of topic membership

- Necessitates different views:

**[El-Assady et al. 2016]**

# Topic Glyph

- Explicitly visually encode topic membership:



- Arc angles: map to the global view (previous slide)

- Brightness: document-topic weight

- Other visual encodings?

# View 2: Glyphs



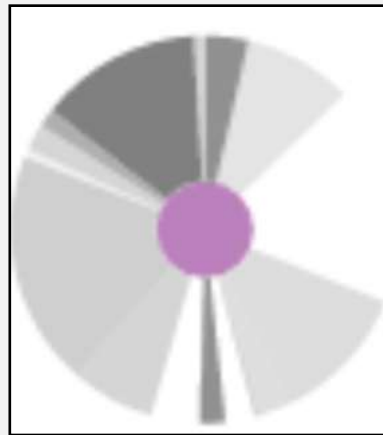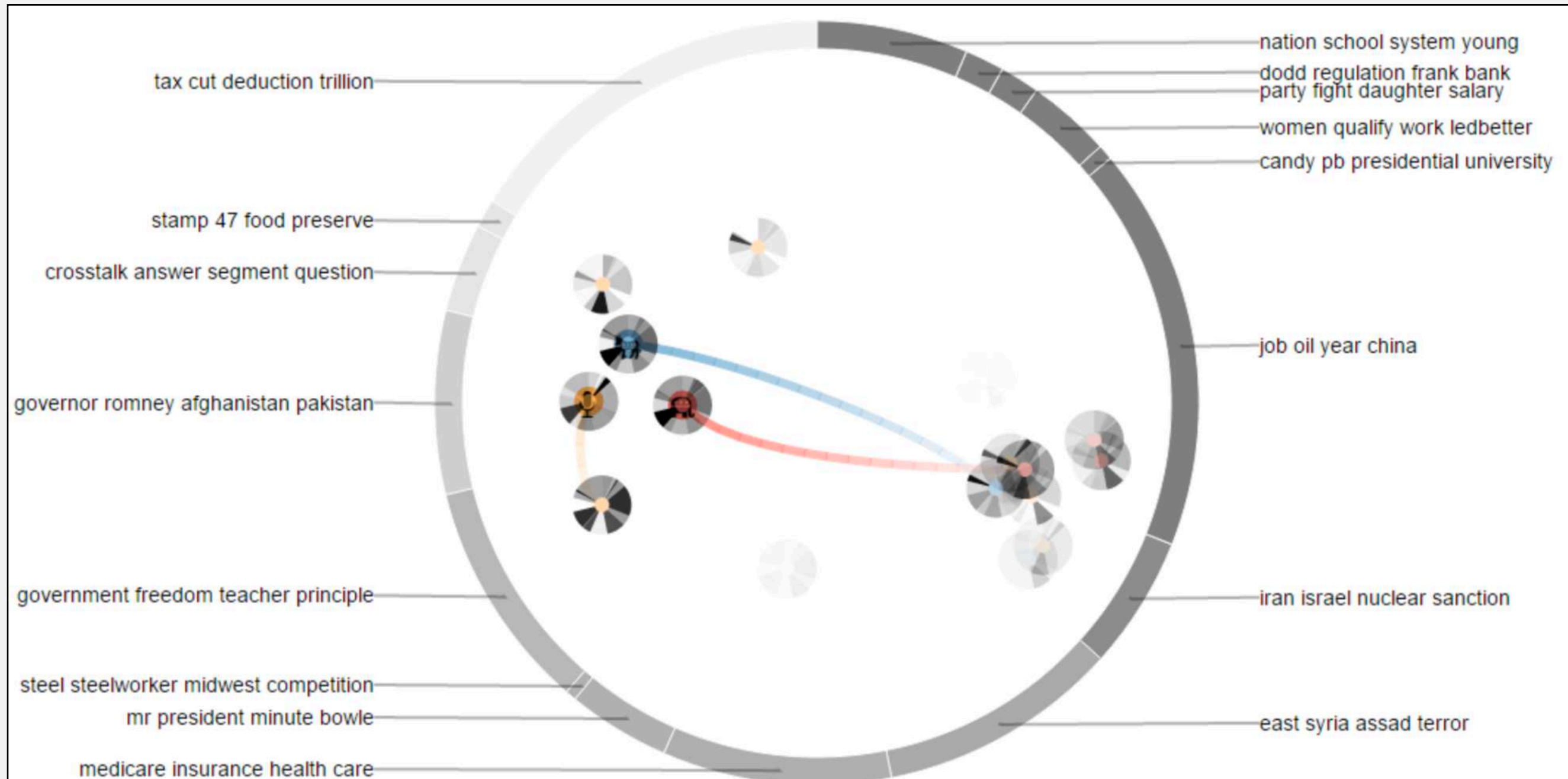tax cut deduction trillion

stamp 47 food preserve

crosstalk answer segment question

governor romney afghanistan pakistan

government freedom teacher principle

steel steelworker midwest competition

mr president minute bowle

medicare insurance health care

nation school system young

dodd regulation frank bank

party fight daughter salary

women qualify work ledbetter

candy pb presidential university

job oil year china

iran israel nuclear sanction

east syria assad terror

**Animation with Context!**

# View 3: "Sedimentation"



party fight daughter occasional
east syria assad terror
women work advocacy ledbetter
teacher state poor medicaid
government prosperous responsibility 47
mr president minute bowle
governor romney afghanistan pakistan
dodd regulation frank bank
tax cut deduction 5

medicare insurance health care
candy pb presidential university
job oil year china
crosstalk answer segment question
iran pressure israel nuclear
nation school system young

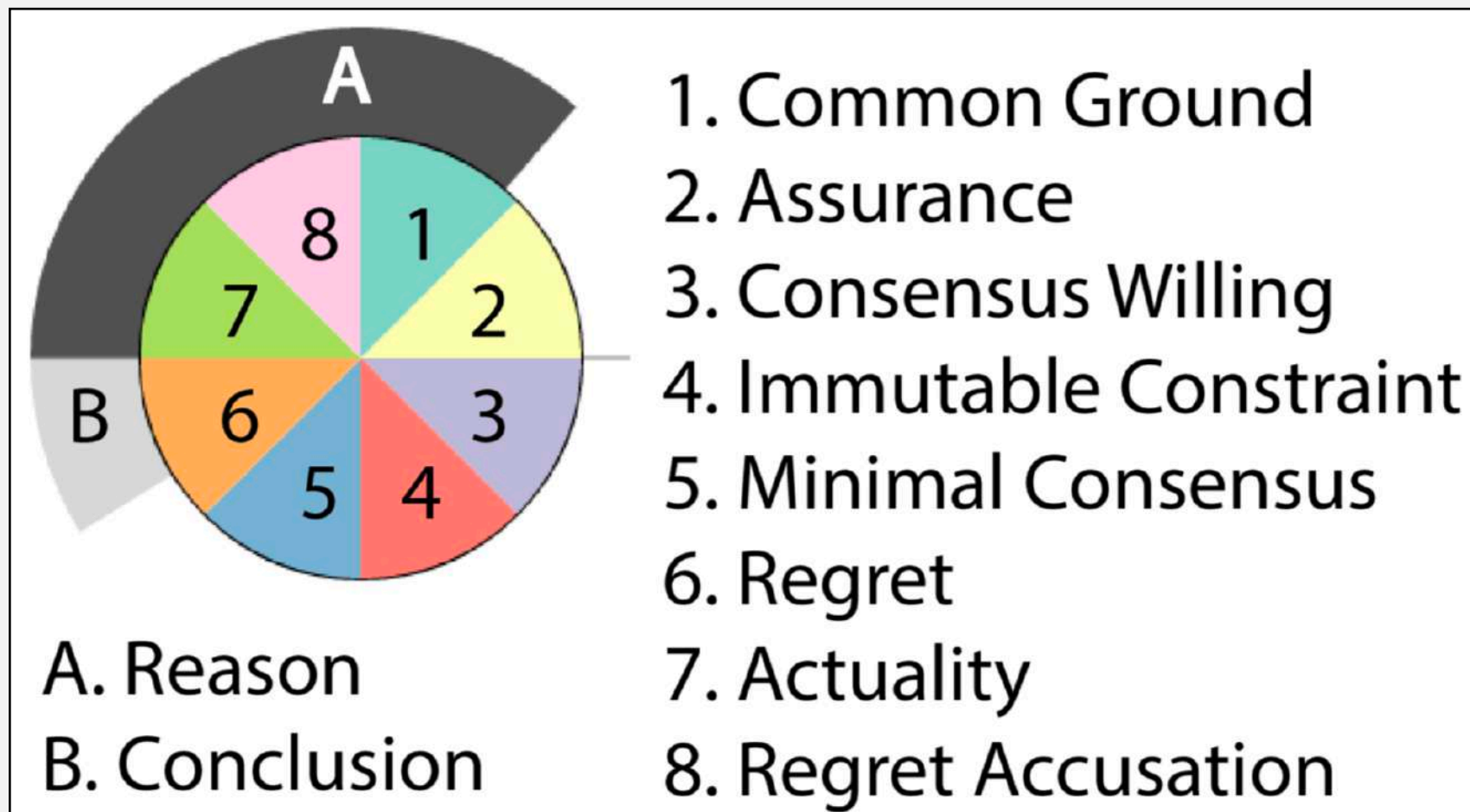**Animation with Context!**

# View 4: Speaker Paths



**Potential Issues? Improvements?**

# Argumentation Patterns

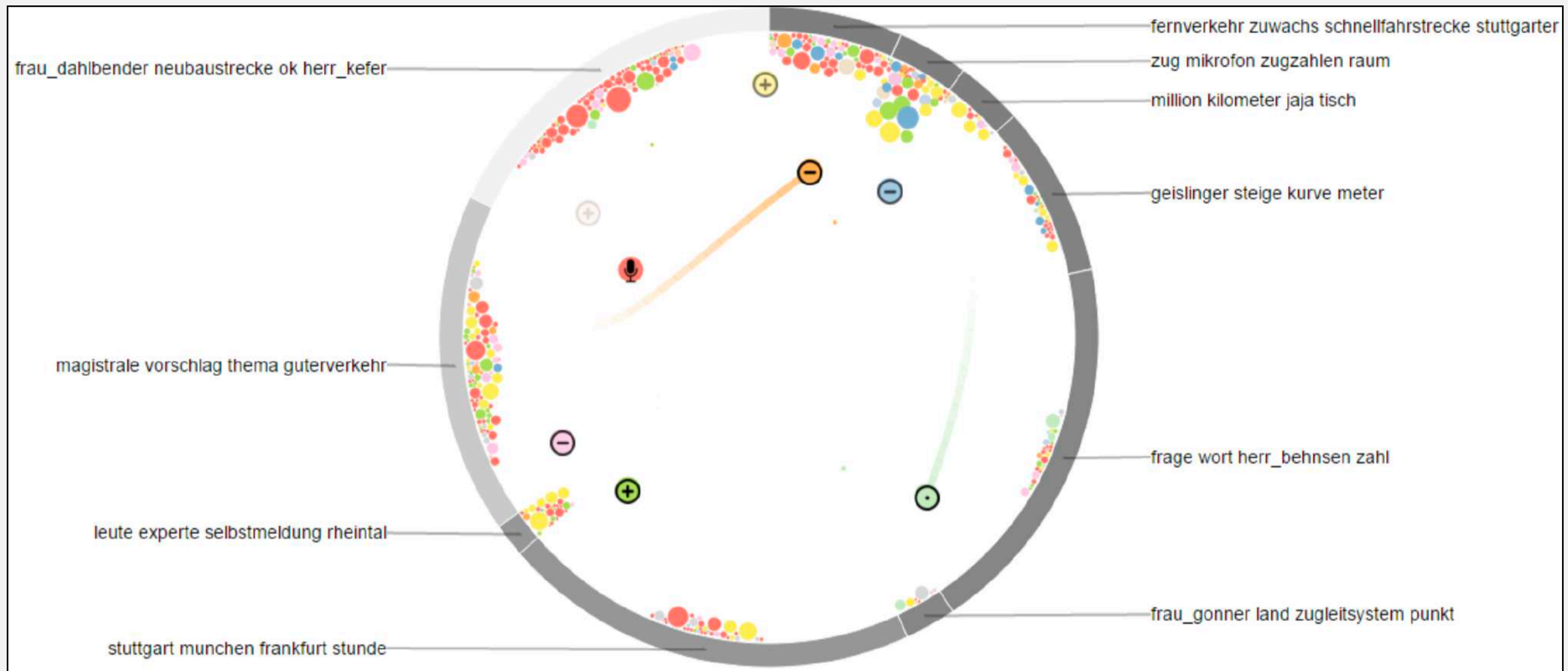- Speaker patterns are extracted, used as additional semantics to understand utterances:



A. Reason
B. Conclusion

1. Common Ground
2. Assurance
3. Consensus Willing
4. Immutable Constraint
5. Minimal Consensus
6. Regret
7. Actuality
8. Regret Accusation
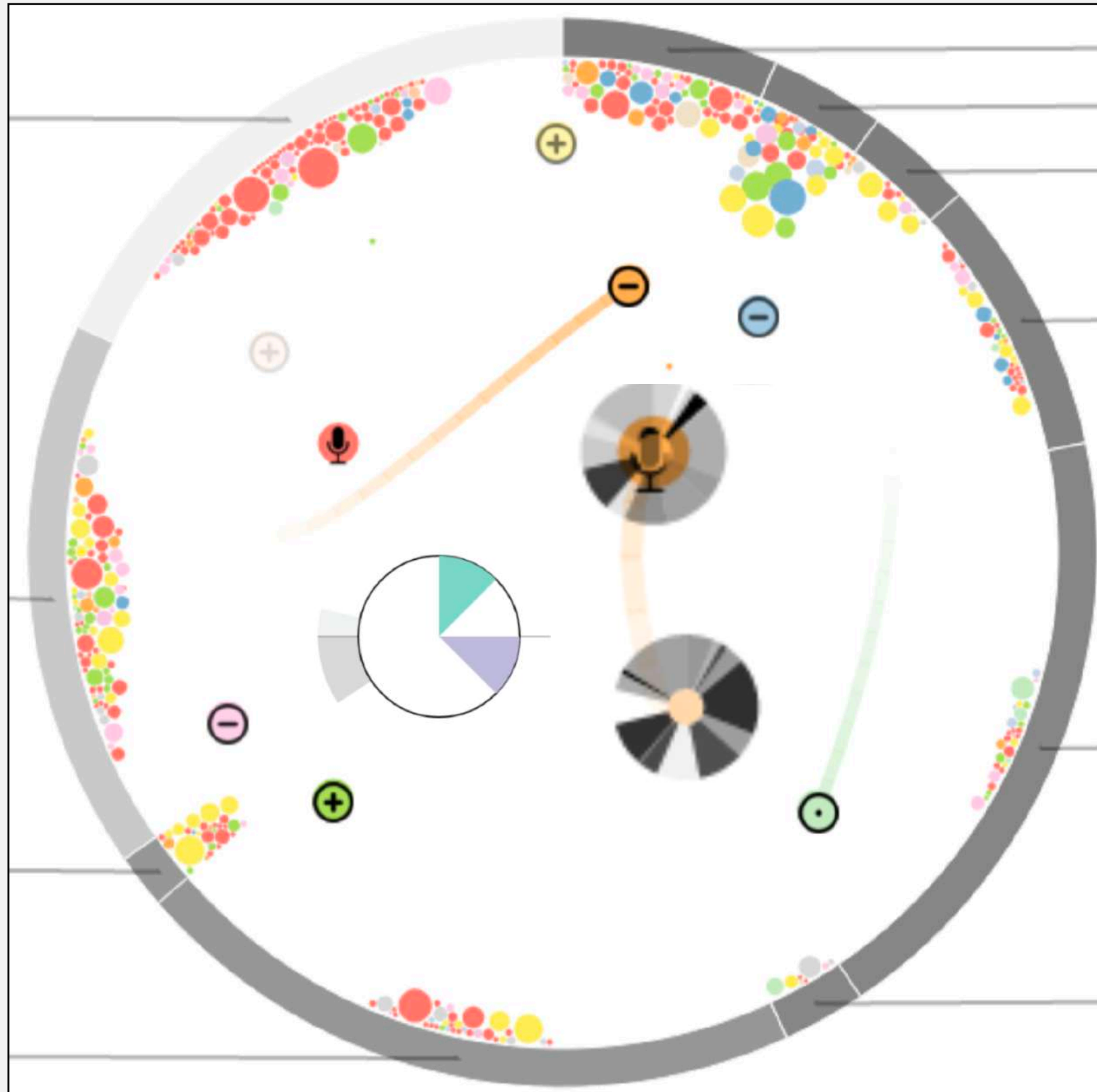
# Comparing Stances

# Detail-on-Demand

# Putting it together

# Sidebar: Encoding Overload?

# Sidebar: Encoding Overload?

- Equally important to designing visual encodings is the ability for us (humans) to *decode* the visualization.

- **Decode**: going from visual channel back to data

- Be careful about *prioritizing* visual encodings:

  - Most important aspects of data should get mapped to channels, or combination of channels, that a human can easily decode (e.g. *channel effectiveness*)