

# Who, What, When, Where, and Why? A Computational Approach to Understanding Historical Events Using State Department Cables

Allison J.B. Chaney  
Princeton University

Joint work with  
David Blei, Hanna Wallach, and  
the History Lab at Columbia

We can do nothing but scrutinize historical events themselves if we want to discover what they are.

– Dean W.R. Matthews, *What is an Historical Event?*

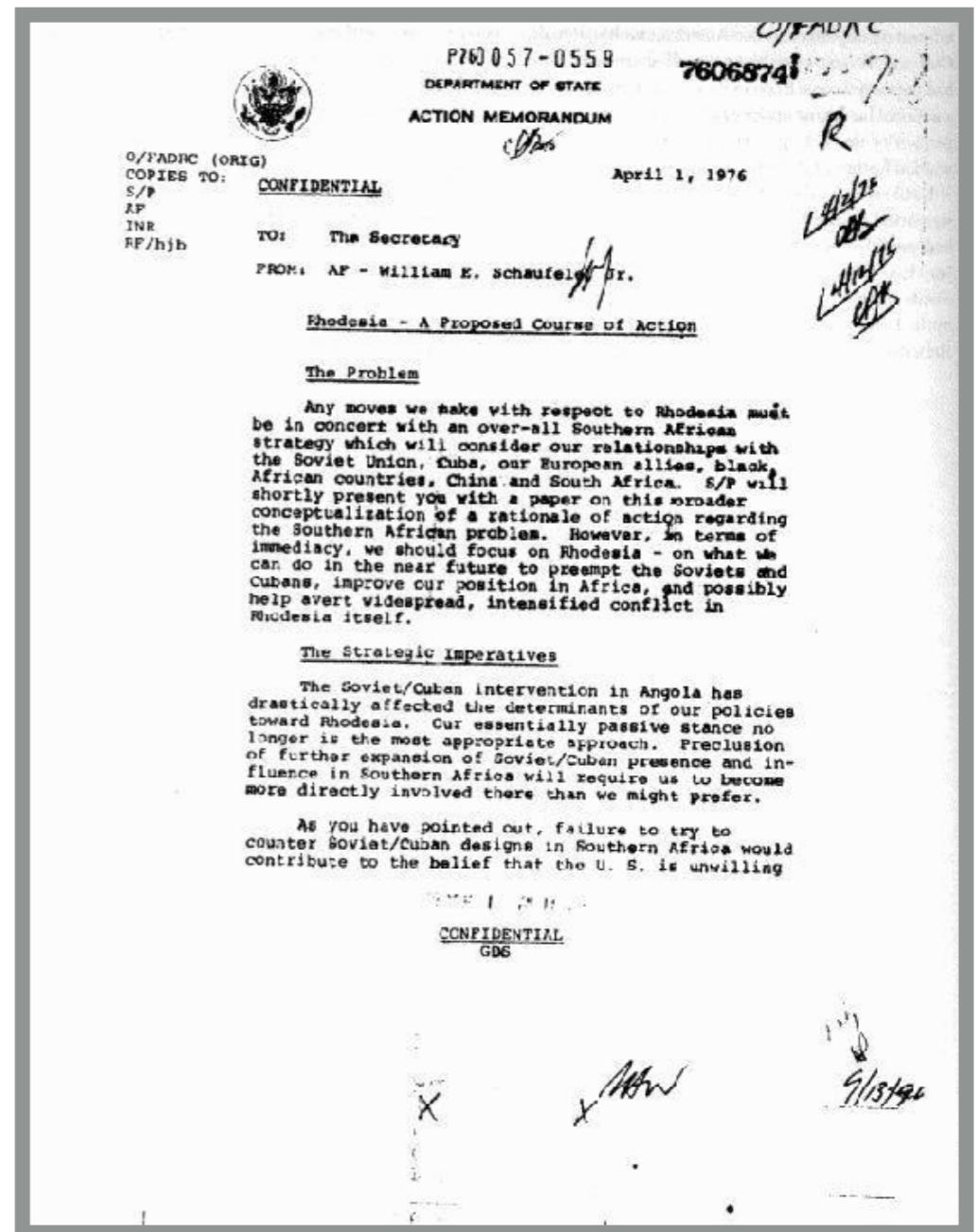


Matthew Connelly's  
History Lab at Columbia



# Matthew Connelly's History Lab at Columbia

## U.S. State Department Cables

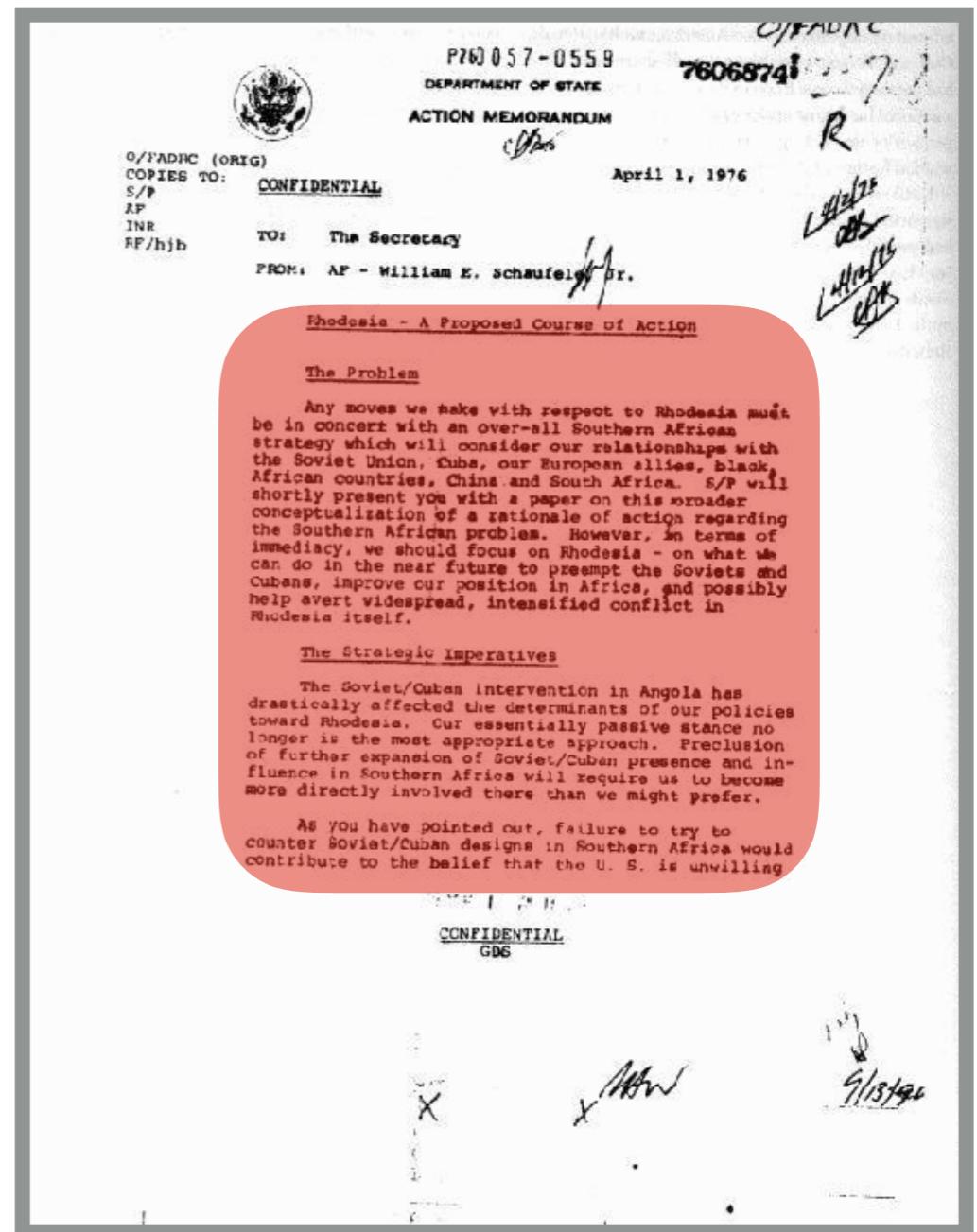




# Matthew Connelly's History Lab at Columbia

## U.S. State Department Cables

### Message content (text)



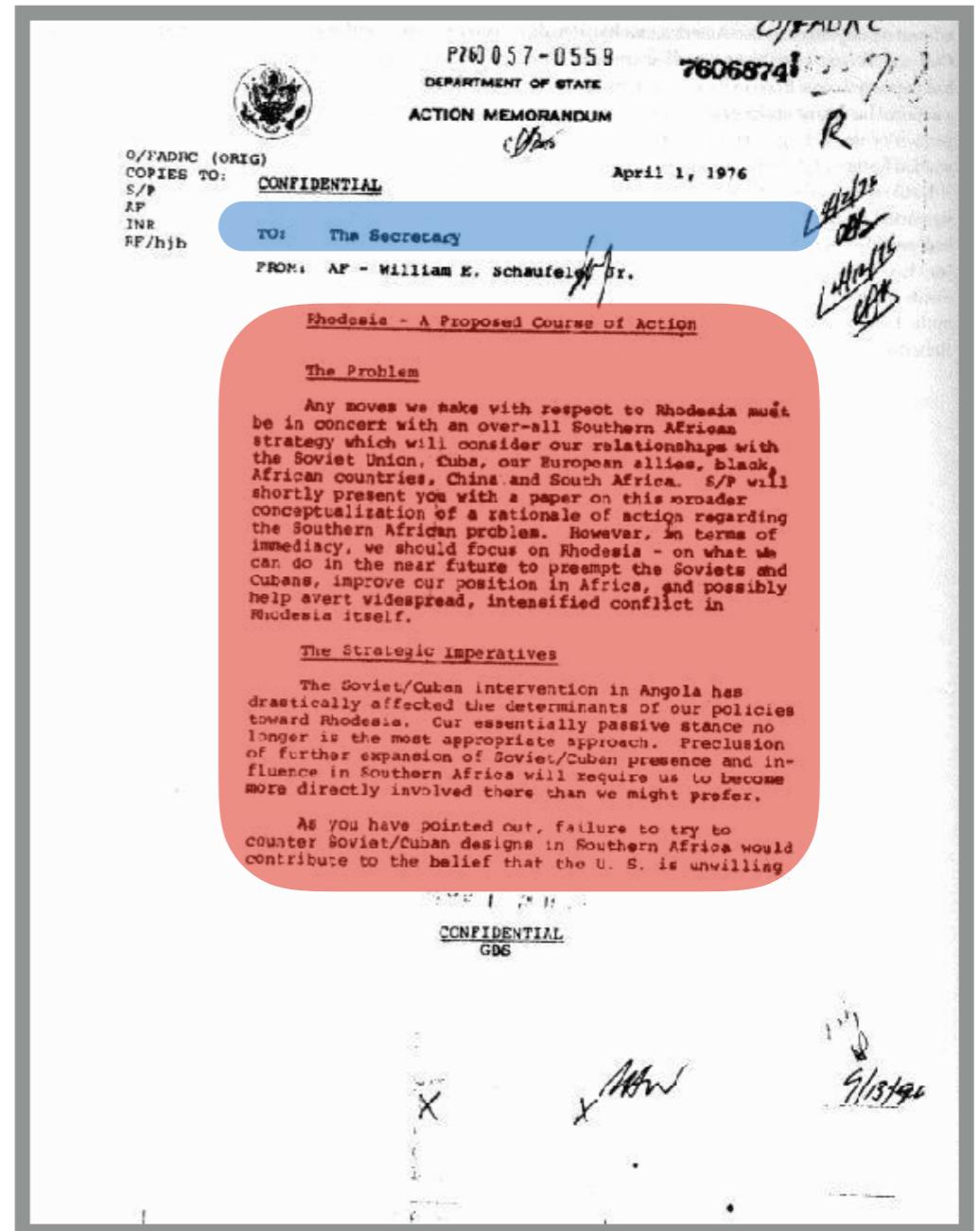


# Matthew Connelly's History Lab at Columbia

## U.S. State Department Cables

Message content (text)

Sending entity (embassy,  
department, or individual)





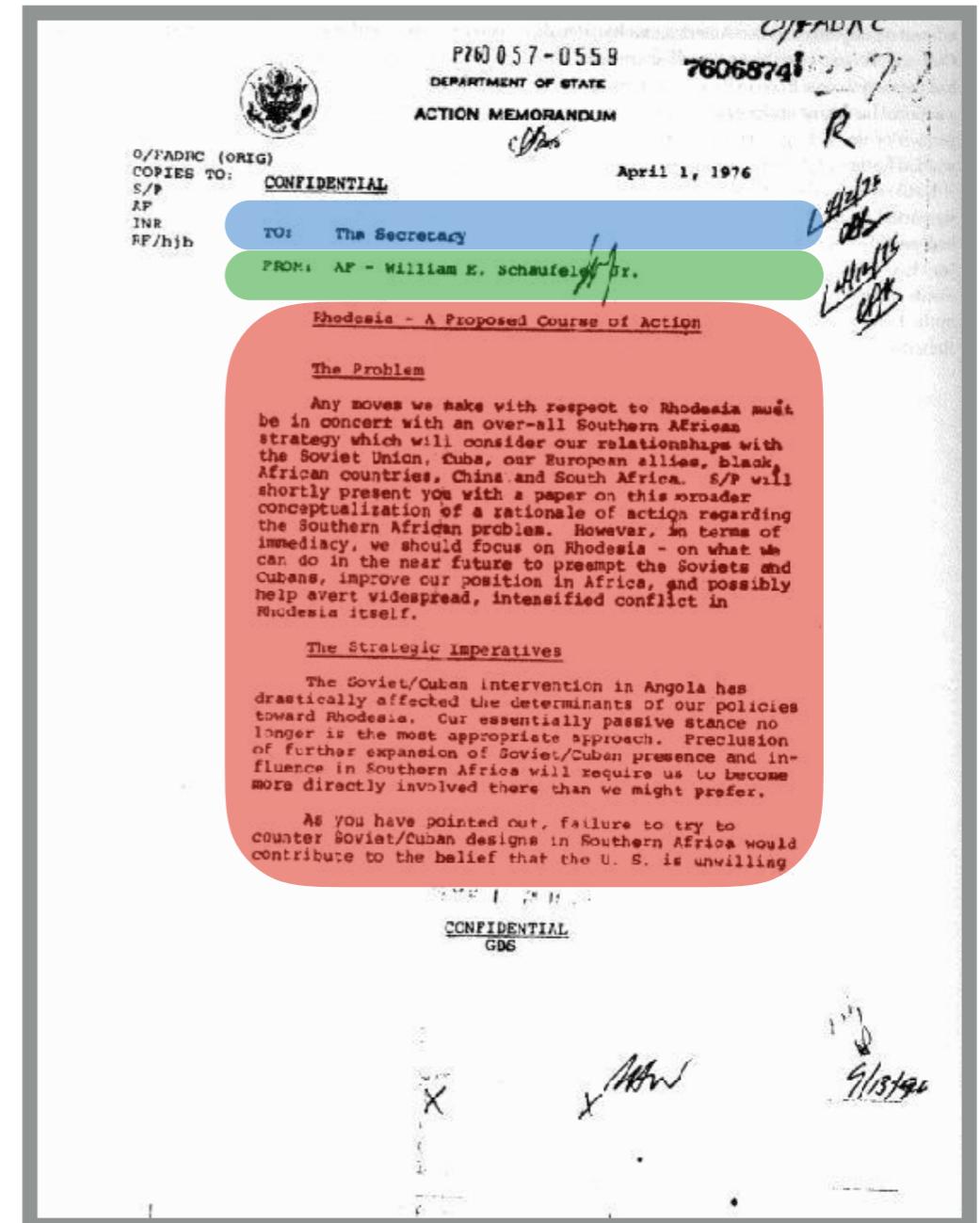
# Matthew Connelly's History Lab at Columbia

## U.S. State Department Cables

Message content (text)

Sending entity (embassy,  
department, or individual)

One or more receiving entities





# Matthew Connelly's History Lab at Columbia

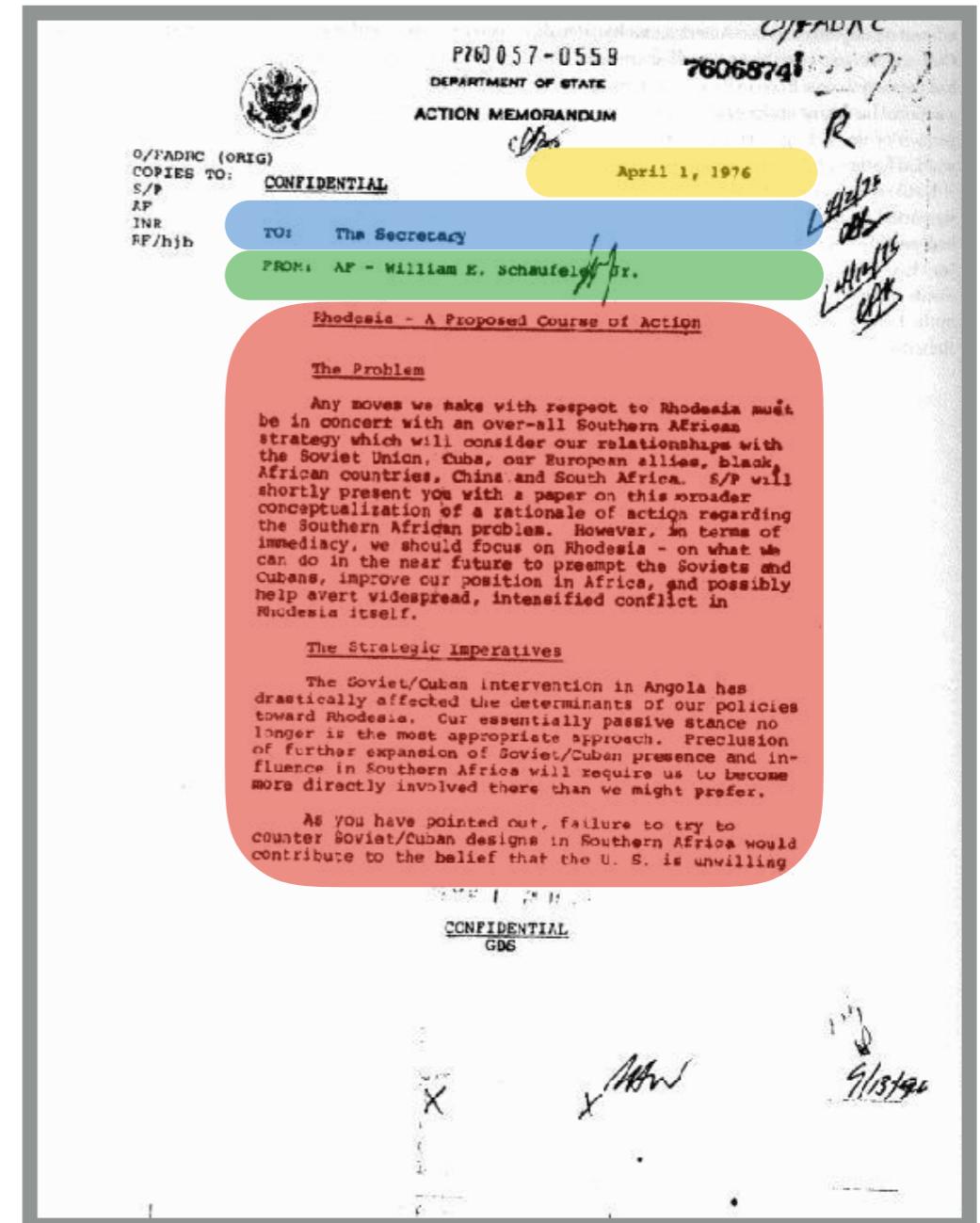
## U.S. State Department Cables

Message content (text)

Sending entity (embassy,  
department, or individual)

One or more receiving entities

Send date





# Matthew Connelly's History Lab at Columbia

## U.S. State Department Cables

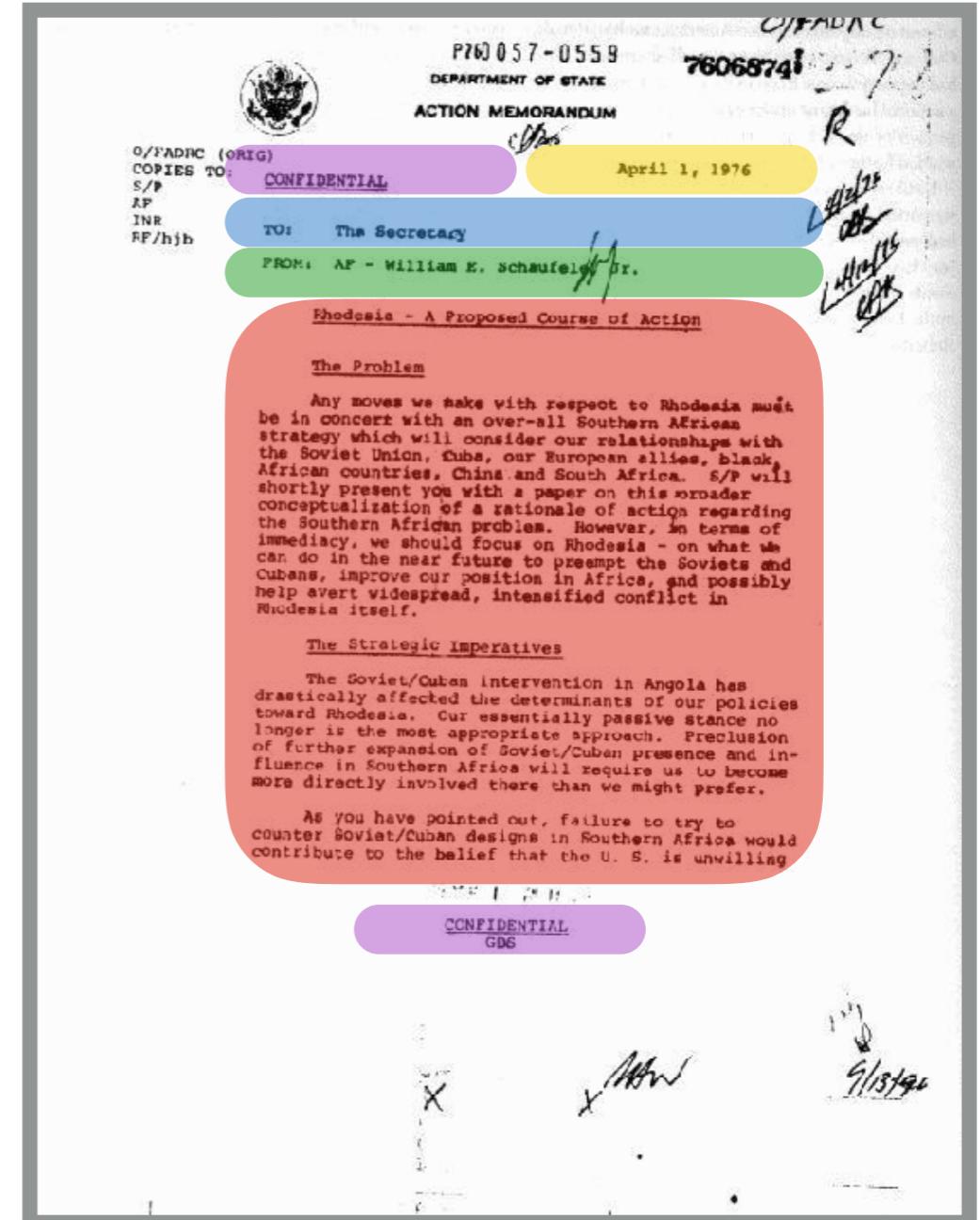
Message content (text)

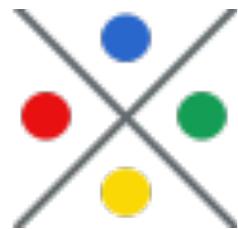
Sending entity (embassy,  
department, or individual)

One or more receiving entities

Send date

Classification level





Matthew Connelly's  
History Lab at Columbia

## U.S. State Department Cables

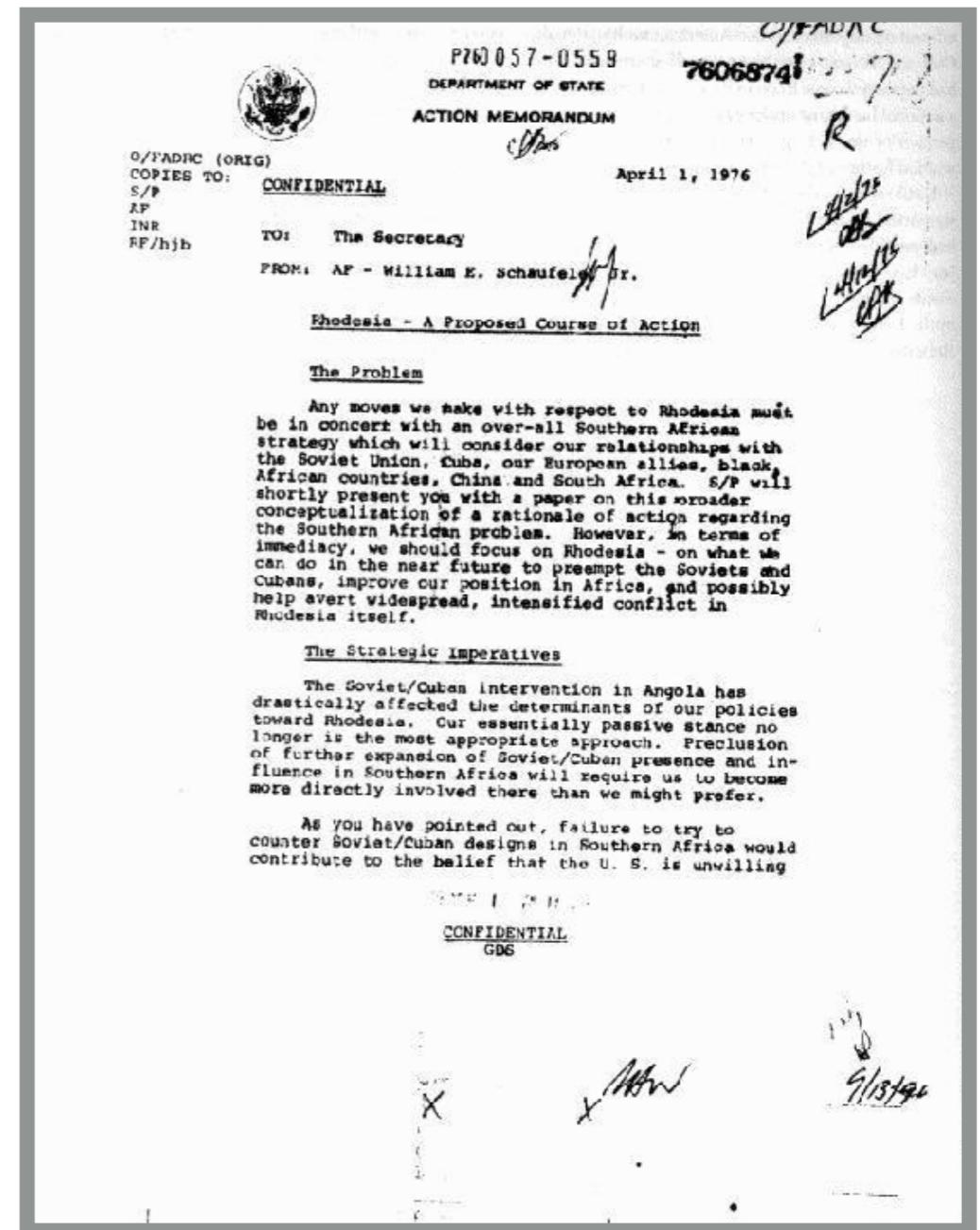
2,674,486 messages

sent between 1973 and 1978

34,204 unique sending entities

23.4% sent from State Dept.

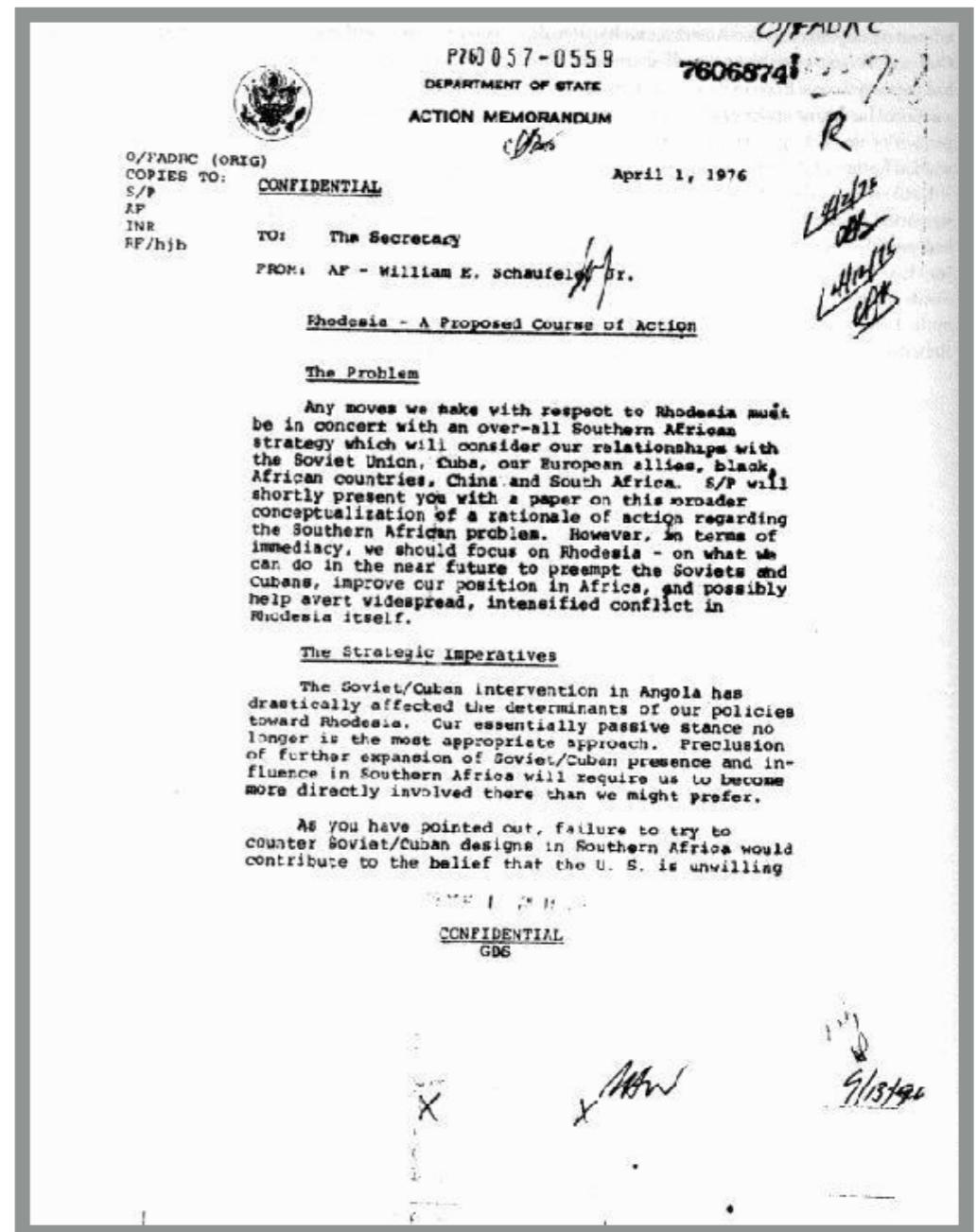
41.6% sent to State Dept.





# Matthew Connelly's History Lab at Columbia

## U.S. State Department Cables

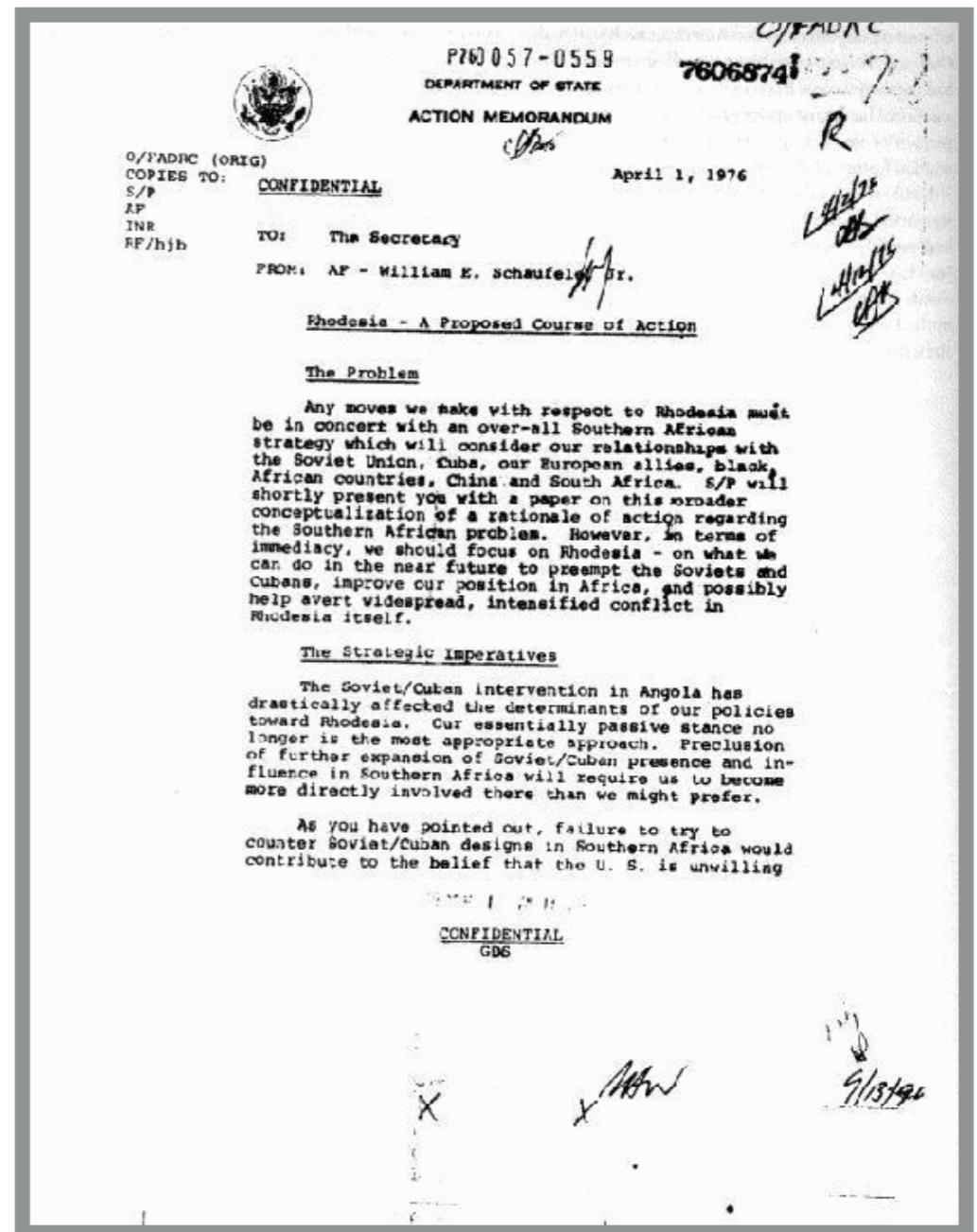




# Matthew Connelly's History Lab at Columbia

## U.S. State Department Cables

What interesting events can be found in these messages?



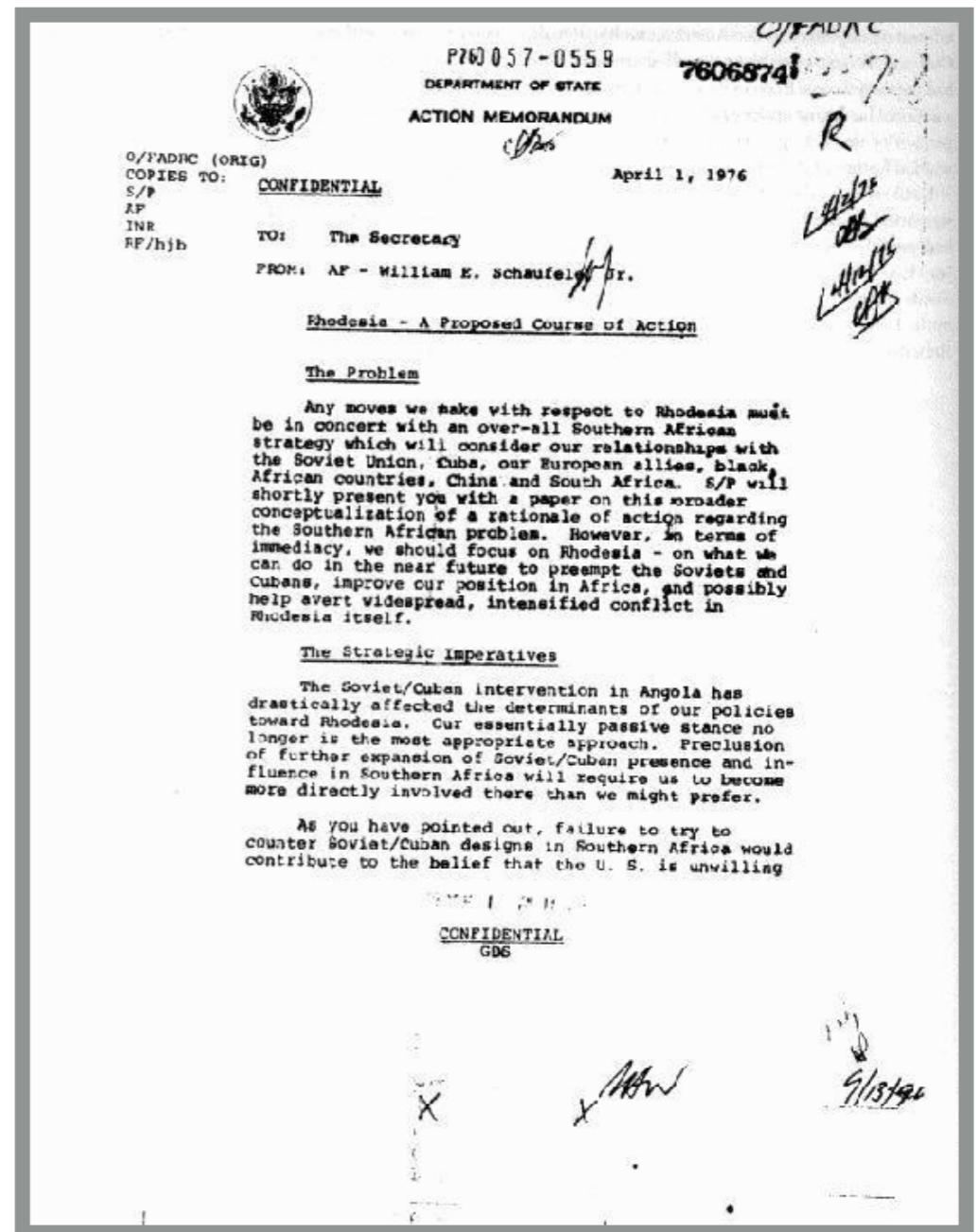


# Matthew Connelly's History Lab at Columbia

## U.S. State Department Cables

What interesting events can be found in these messages?

How can we describe events?





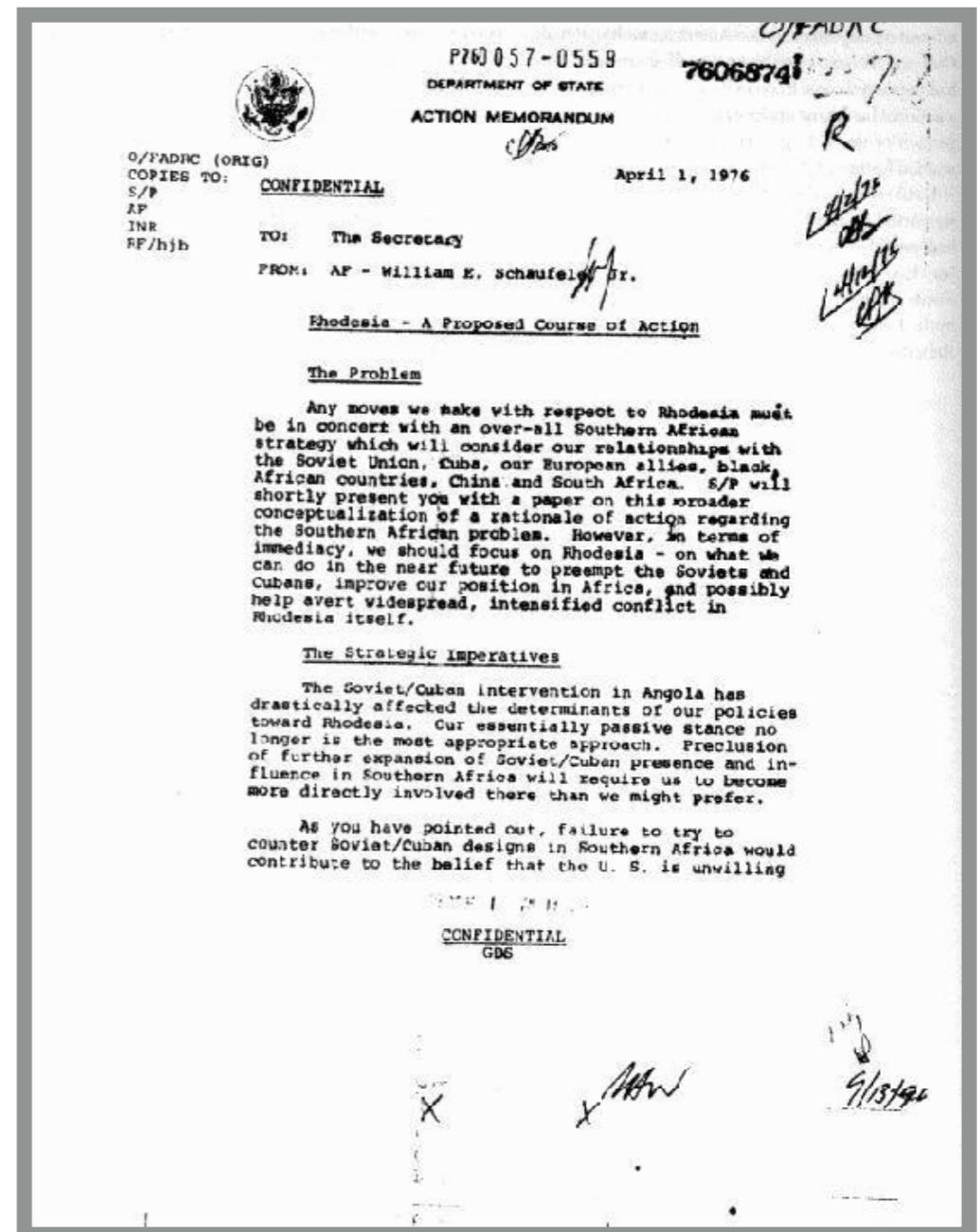
# Matthew Connelly's History Lab at Columbia

## U.S. State Department Cables

What interesting events can be found in these messages?

How can we describe events?

What relationships do different entities have?





# Matthew Connelly's History Lab at Columbia

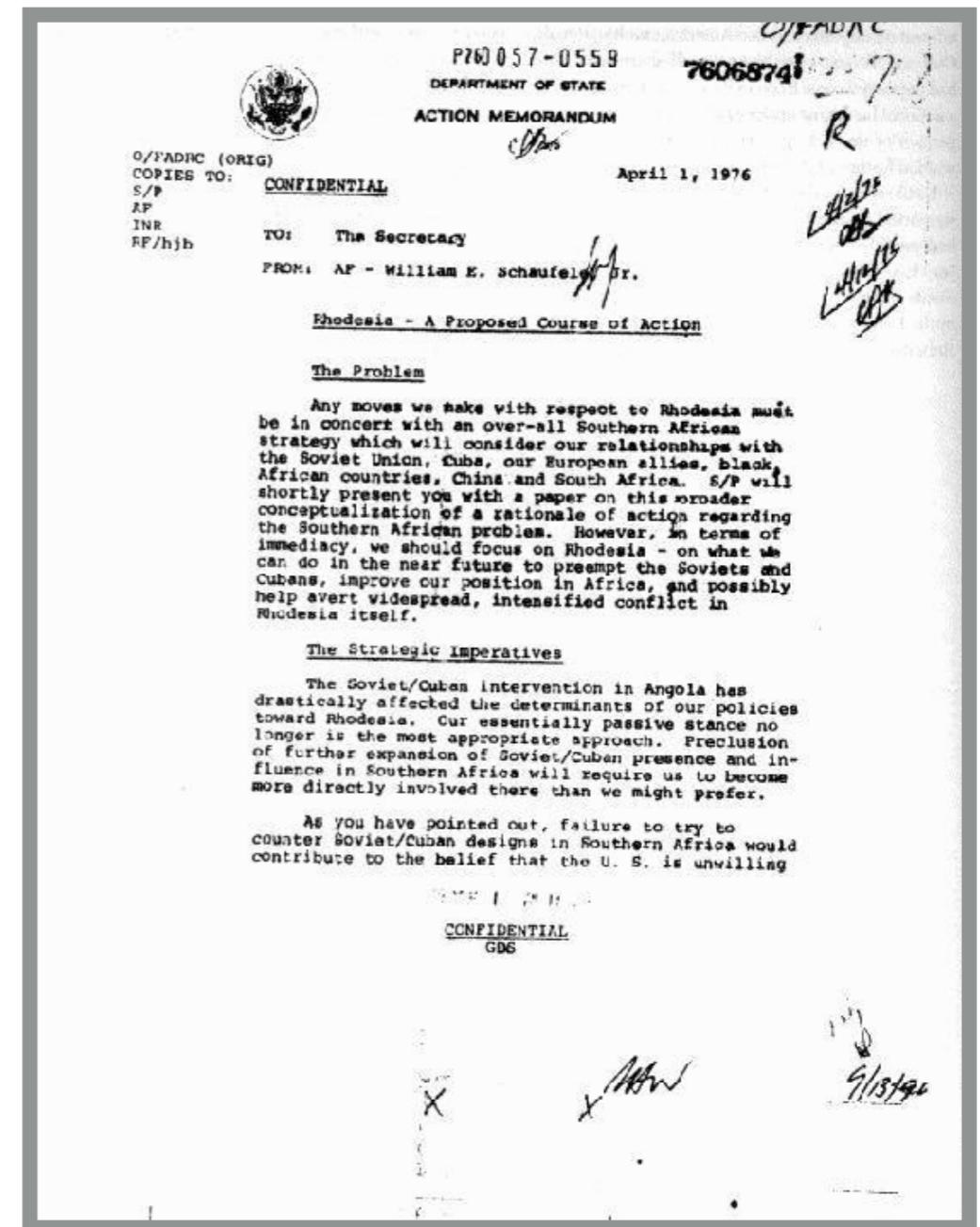
## U.S. State Department Cables

What interesting events can be found in these messages?

How can we describe events?

What relationships do different entities have?

What are the typical concerns of different entities?



Events are **unobserved**.

Events are **unobserved**.

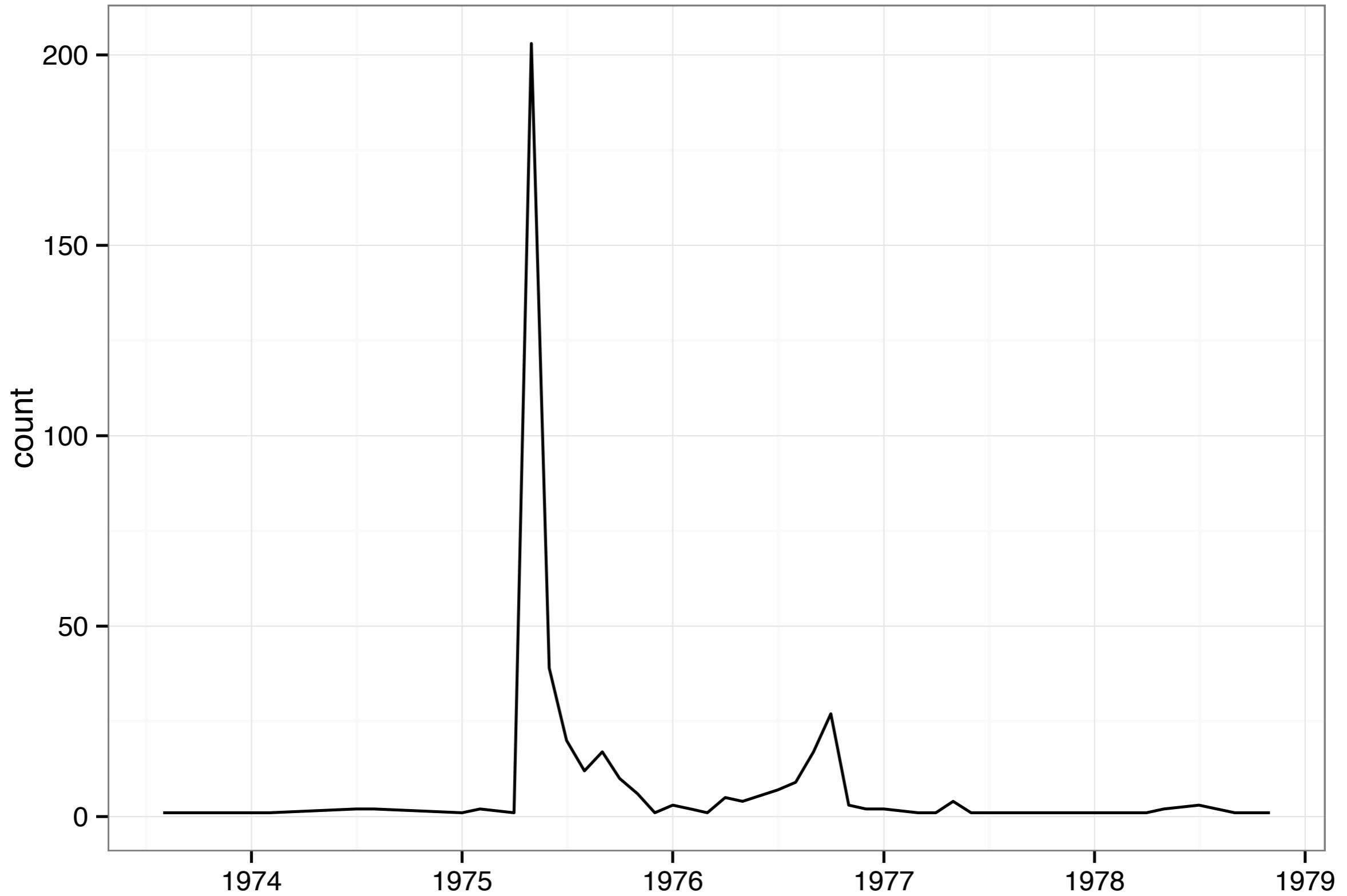
What are **observed** ways to characterize events?

Events are **unobserved**.

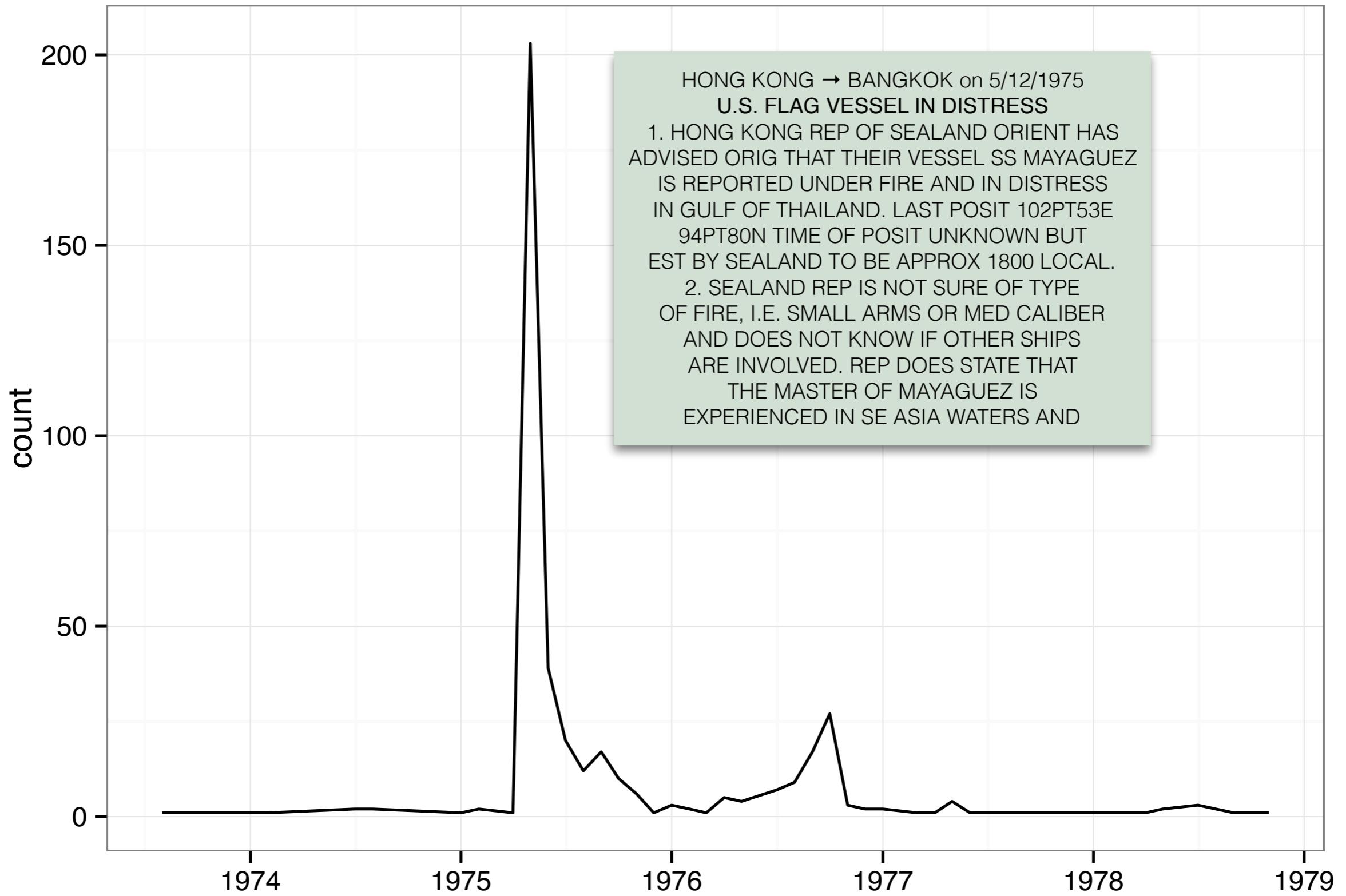
What are **observed** ways to characterize events?

Temporary shifts away  
from business as usual  
in message content

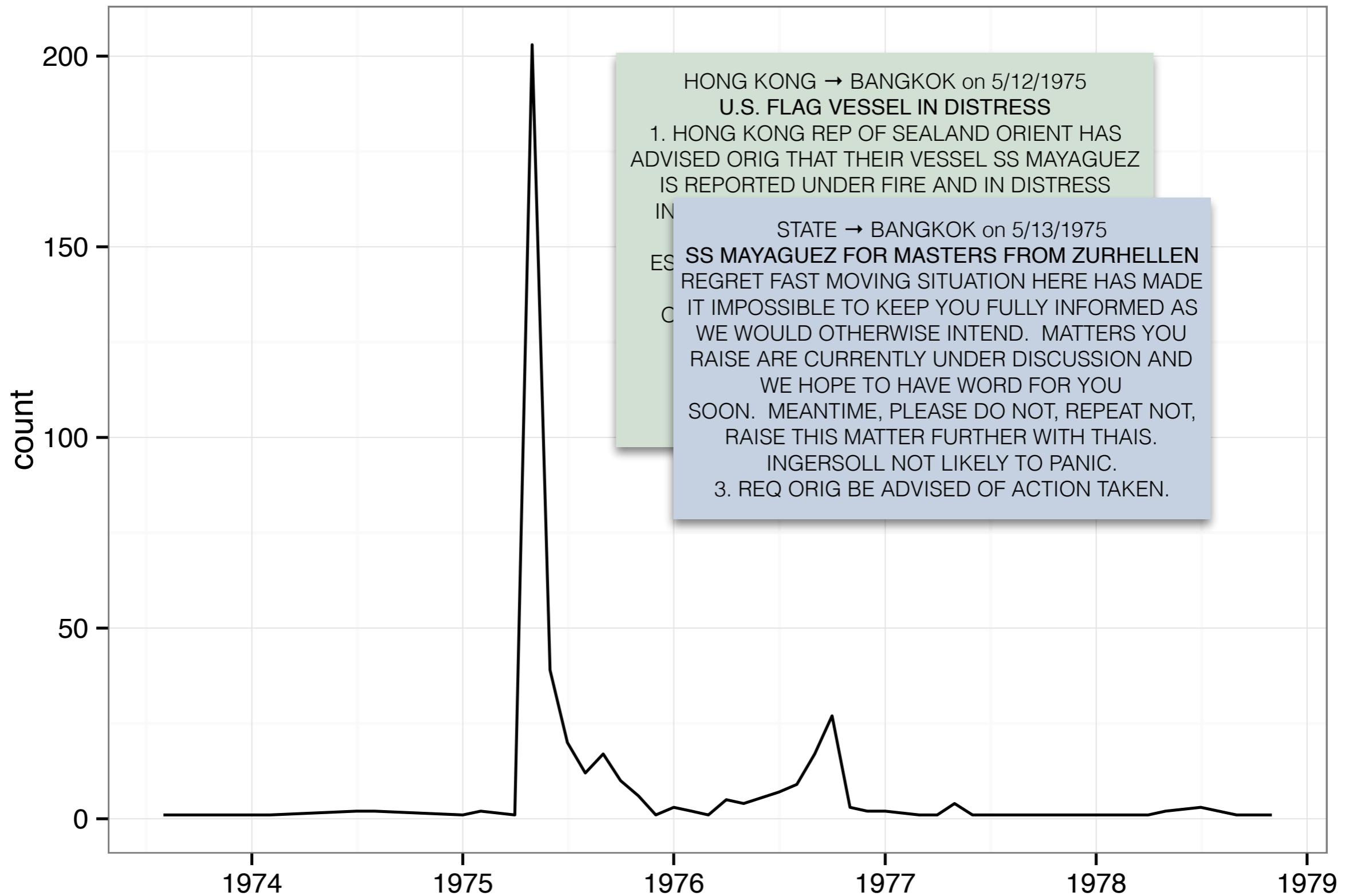
# Mayaguez Incident



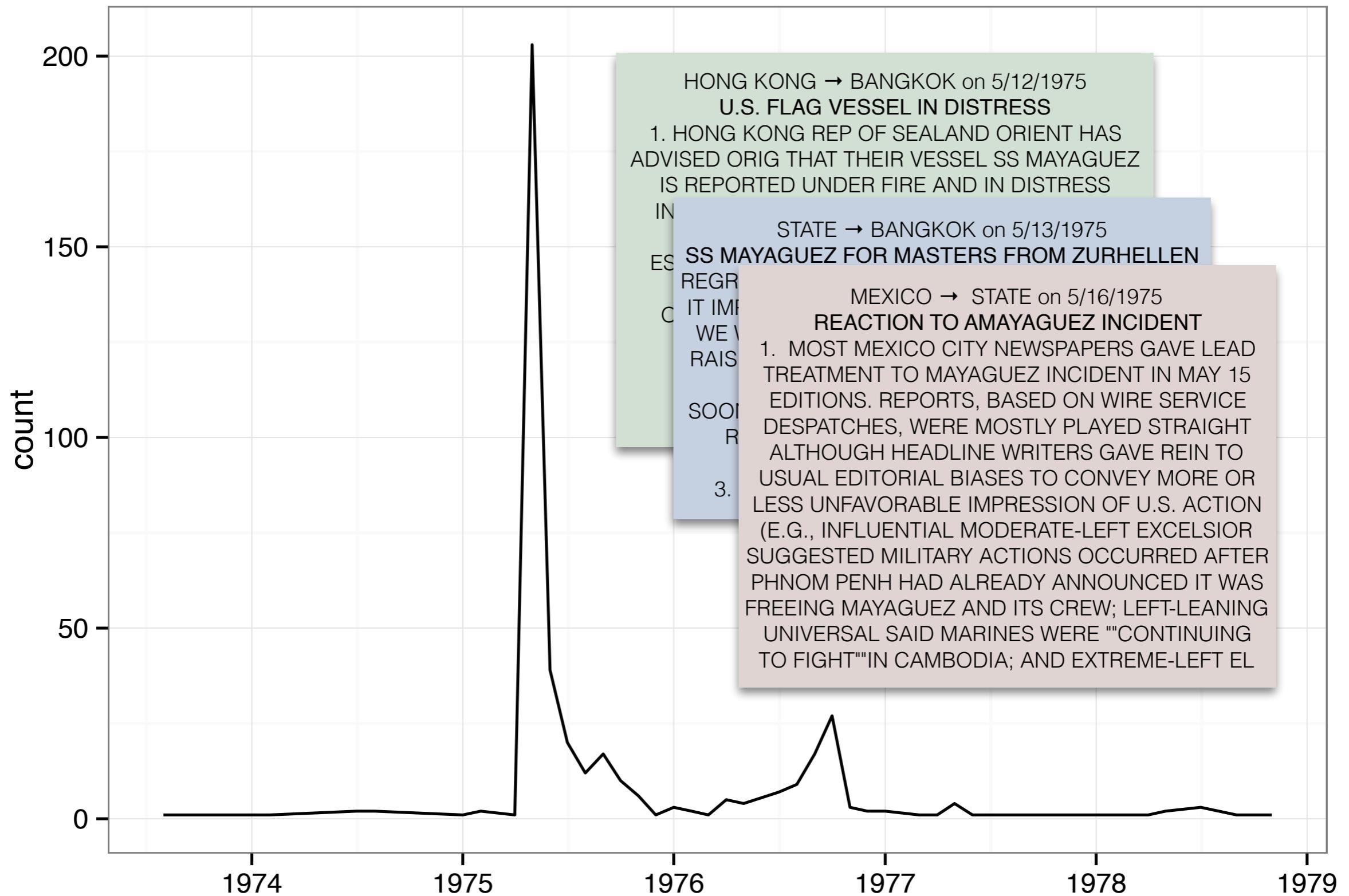
# Mayaguez Incident



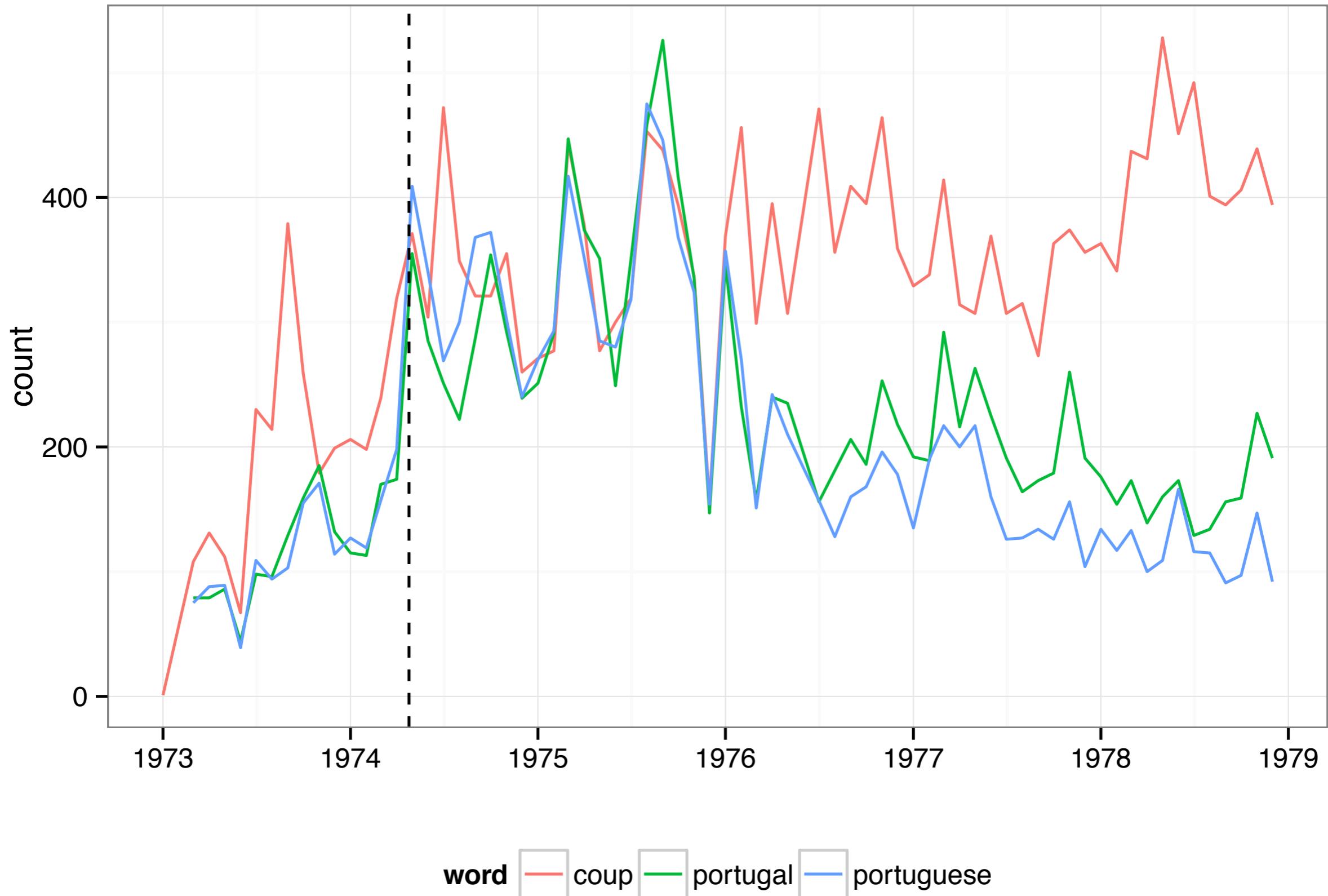
# Mayaguez Incident

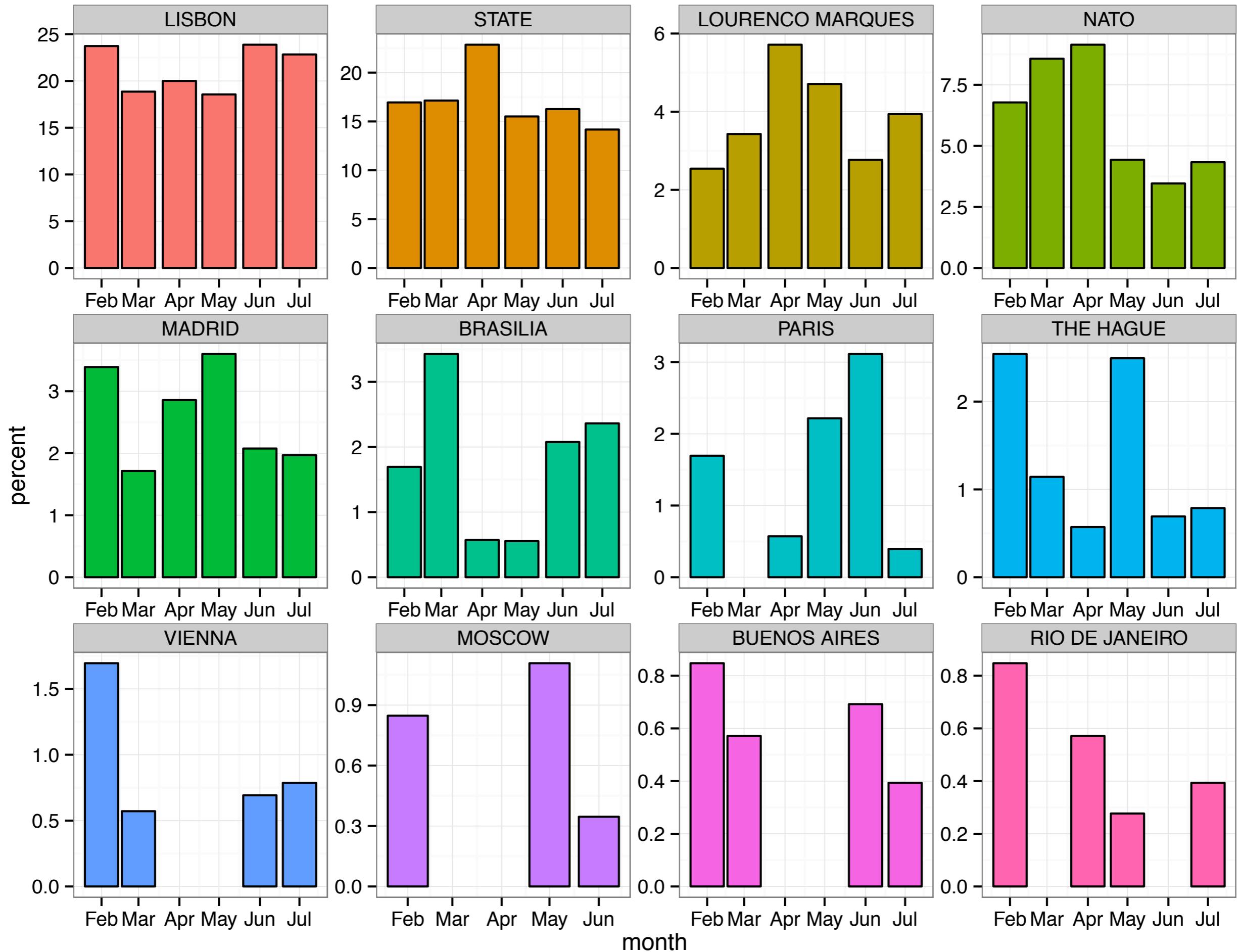


# Mayaguez Incident



# Carnation Revolution





What are the key actors in  
constructing our model?

# cables

HONG KONG → BANGKOK on 5/12/1975  
U.S. FLAG VESSEL IN DISTRESS

1. HONG KONG REP OF SEALAND ORIENT HAS ADVISED ORIG THAT THEIR VESSEL SS MAYAGUEZ IS REPORTED UNDER FIRE AND IN DISTRESS

IN

STATE → BANGKOK on 5/13/1975

ES SS MAYAGUEZ FOR MASTERS FROM ZURHELLEN

REGR

IT IMP

WE

RAIS

SOON

R

3.

MEXICO → STATE on 5/16/1975

REACTION TO AMAYAGUEZ INCIDENT

1. MOST MEXICO CITY NEWSPAPERS GAVE LEAD TREATMENT TO MAYAGUEZ INCIDENT IN MAY 15 EDITIONS. REPORTS, BASED ON WIRE SERVICE DESPATCHES, WERE MOSTLY PLAYED STRAIGHT ALTHOUGH HEADLINE WRITERS GAVE REIN TO USUAL EDITORIAL BIASES TO CONVEY MORE OR LESS UNFAVORABLE IMPRESSION OF U.S. ACTION (E.G., INFLUENTIAL MODERATE-LEFT EXCELSIOR SUGGESTED MILITARY ACTIONS OCCURRED AFTER PHNOM PENH HAD ALREADY ANNOUNCED IT WAS FREEING MAYAGUEZ AND ITS CREW; LEFT-LEANING UNIVERSAL SAID MARINES WERE ""CONTINUING TO FIGHT"" IN CAMBODIA; AND EXTREME-LEFT EL

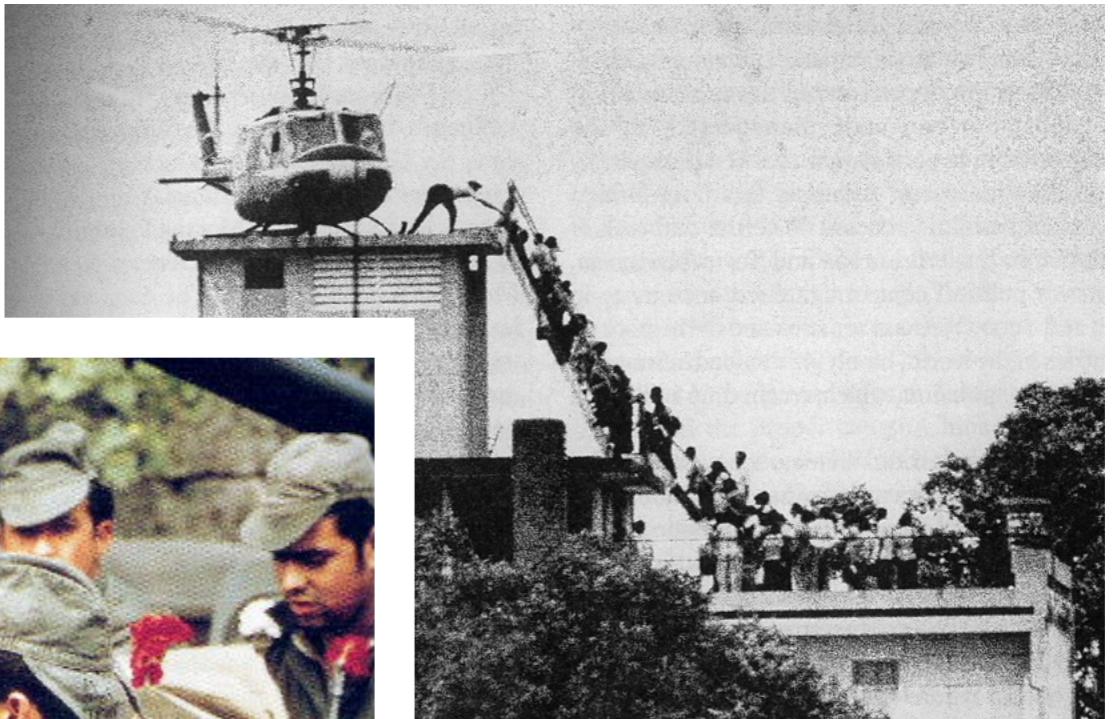
# entities



# cables

HONG KONG →  
BANGKOK on 5/12/1975  
U.S. FLAG VESSEL IN  
STATF → BANGKOK on  
1 MEXICO → STATE on  
5/16/1975  
REACTION TO  
AMAYAGUEZ INCIDENT  
1. MOST MEXICO CITY  
NEWSPAPERS GAVE  
LEAD TREATMENT TO

# events



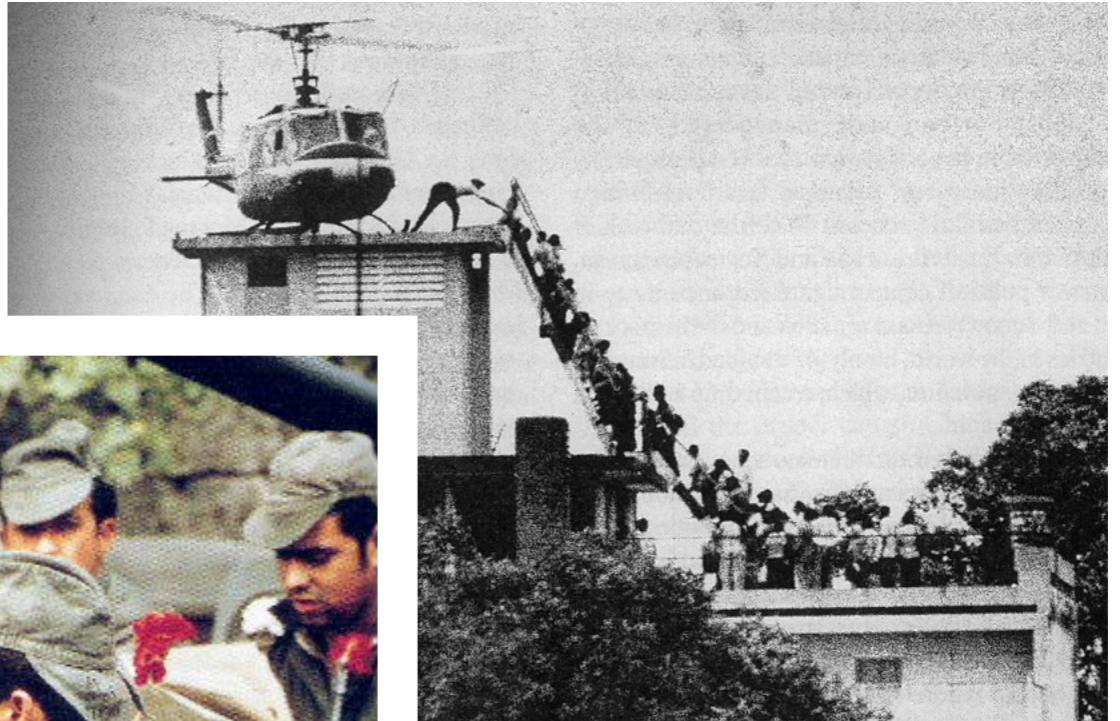
## cables

HONG KONG →  
BANGKOK on 5/12/1975  
U.S. FLAG VESSEL IN  
STATF → BANGKOK on  
1  
S  
MEXICO → STATE on  
5/16/1975  
REACTION TO  
AMAYAGUEZ INCIDENT  
1. MOST MEXICO CITY  
NEWSPAPERS GAVE  
LEAD TREATMENT TO

## entities



# events



## cables

HONG KONG →  
BANGKOK on 5/12/1975  
U.S. FLAG VESSEL IN  
STATF → BANGKOK on  
1 MEXICO → STATE on  
5/16/1975  
REACTION TO  
AMAYAGUEZ INCIDENT  
1. MOST MEXICO CITY  
NEWSPAPERS GAVE  
LEAD TREATMENT TO

## entities



# representing cables with topics

**11/28/197: EUROPEAN MINISTERS OF EDUCATION CONFERENCE BUCHAREST**  
SECOND AND THIRD PLENARY SESSIONS OF CONFERENCE CONTINUED NOVEMBER  
26 WITH SPEECHES BY 18 ADDITIONAL DELEGATIONS AND EIGHT OBSERVER  
DELEGATIONS ON DEVELOPMENTS IN NATIONAL EDUCATION SINCE 1967. OF THREE  
COUNTRIES REPRESENTED AS OBSERVERS (U.S., CANADA, ISRAEL), ONLY CANADIAN  
DELEGATION ADDRESSED PLENARY WITH DESCRIPTION OF MUTUAL CANADIAN-  
EUROPEAN EDUCATIONAL PROBLEMS AND CONCERNS. CANADIAN DELEGATE  
CONCLUDED WITH EXPRESSION OF CANADIAN DESIRE FOR FULL MEMBERSHIP IN  
EUROPEAN REGION. FOLLOWING DEPARTMENT INSTRUCTIONS, U.S. DEL DID NOT  
ADDRESS PLENARY ALTHOUGH DELEGATION PLANS PARTICIPATION AS  
APPROPRIATE IN WORKING SESSIONS BEGINNING NOVEMBER 27.



*Latent Dirichlet allocation. Blei, Ng, and Jordan, 2003.*

# representing cables with topics

11/28/197: EUROPEAN MINISTERS OF EDUCATION CONFERENCE BUCHAREST  
SECOND AND THIRD PLENARY SESSIONS OF CONFERENCE CONTINUED NOVEMBER  
26 WITH SPEECHES BY 18 ADDITIONAL DELEGATIONS AND EIGHT OBSERVER  
DELEGATIONS ON DEVELOPMENTS IN NATIONAL EDUCATION SINCE 1967. OF THREE  
COUNTRIES REPRESENTED AS OBSERVERS (U.S., CANADA, ISRAEL), ONLY CANADIAN  
DELEGATION ADDRESSED PLENARY WITH DESCRIPTION OF MUTUAL CANADIAN-  
EUROPEAN EDUCATIONAL PROBLEMS AND CONCERNS. CANADIAN DELEGATE  
CONCLUDED WITH EXPRESSION OF CANADIAN DESIRE FOR FULL MEMBERSHIP IN  
EUROPEAN REGION. FOLLOWING DEPARTMENT INSTRUCTIONS, U.S. DEL DID NOT  
ADDRESS PLENARY ALTHOUGH DELEGATION PLANS PARTICIPATION AS  
APPROPRIATE IN WORKING SESSIONS BEGINNING NOVEMBER 27.



*Latent Dirichlet allocation.*  
Blei, Ng, and Jordan, 2003.

**Advantage:** Good for discovering general, interpretable themes useful for representing *entities' typical concerns*

**Disadvantage:** Does not capture word-level shifts in language or subject  
not useful for representing *event descriptions*

# representing cables with words

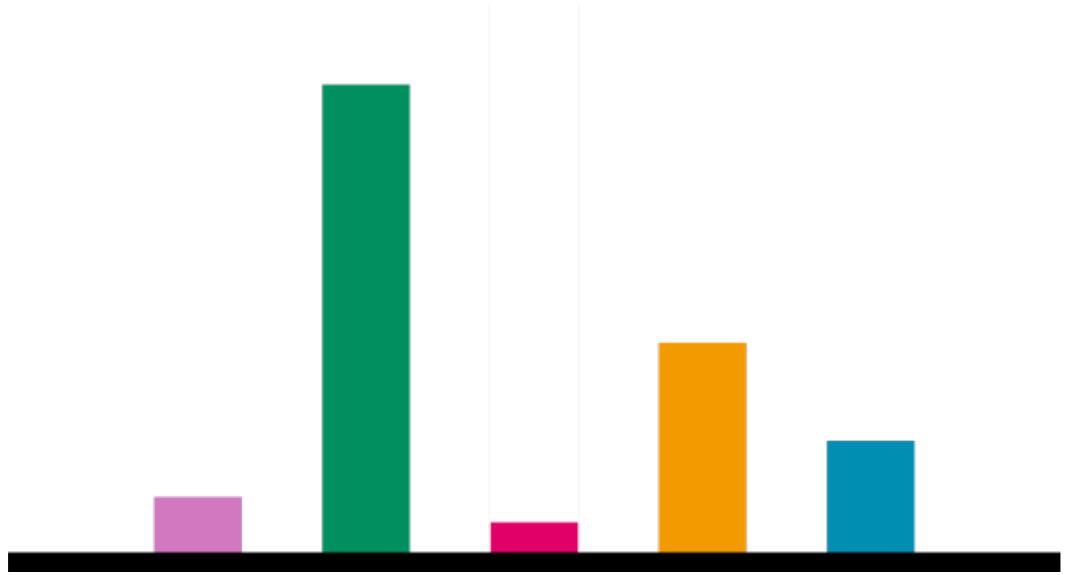
**11/28/197: EUROPEAN MINISTERS OF EDUCATION CONFERENCE BUCHAREST**  
SECOND AND THIRD PLENARY SESSIONS OF CONFERENCE CONTINUED NOVEMBER 26 WITH SPEECHES BY 18 ADDITIONAL DELEGATIONS AND EIGHT OBSERVER DELEGATIONS ON DEVELOPMENTS IN NATIONAL EDUCATION SINCE 1967. OF THREE COUNTRIES REPRESENTED AS OBSERVERS (U.S., CANADA, ISRAEL), ONLY CANADIAN DELEGATION ADDRESSED PLENARY WITH DESCRIPTION OF MUTUAL CANADIAN-EUROPEAN EDUCATIONAL PROBLEMS AND CONCERNS. CANADIAN DELEGATE CONCLUDED WITH EXPRESSION OF CANADIAN DESIRE FOR FULL MEMBERSHIP IN EUROPEAN REGION. FOLLOWING DEPARTMENT INSTRUCTIONS, U.S. DEL DID NOT ADDRESS PLENARY ALTHOUGH DELEGATION PLANS PARTICIPATION AS APPROPRIATE IN WORKING SESSIONS BEGINNING NOVEMBER 27.

**Advantage:** Good for capturing shifts in specific subjects and terminology  
useful for representing *event descriptions*

**Disadvantage:** harder to interpret, may cause scalability issues

not useful for representing *entities' typical concerns*

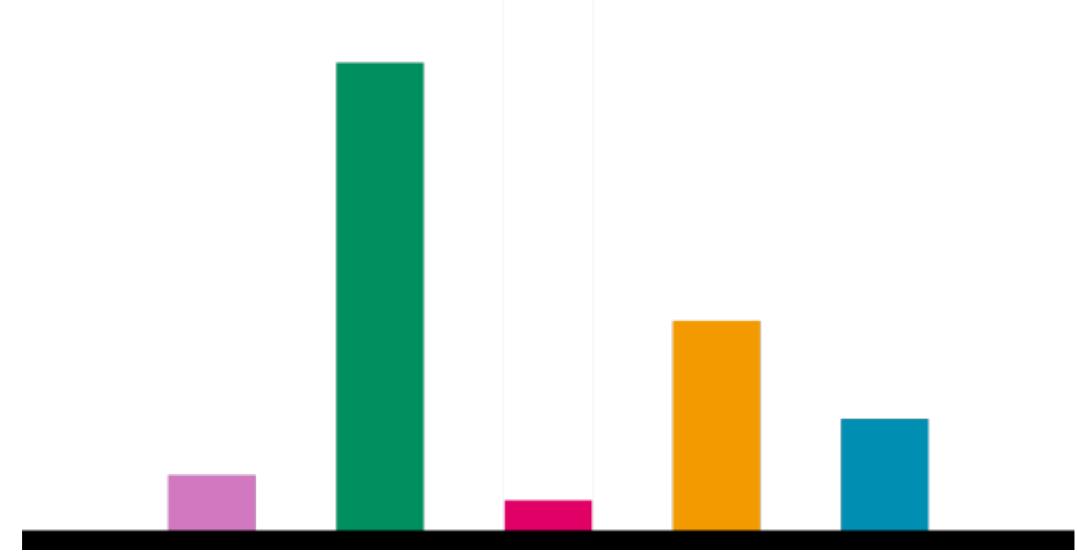
# modeling entities



typical concerns  
of the Bangkok Embassy



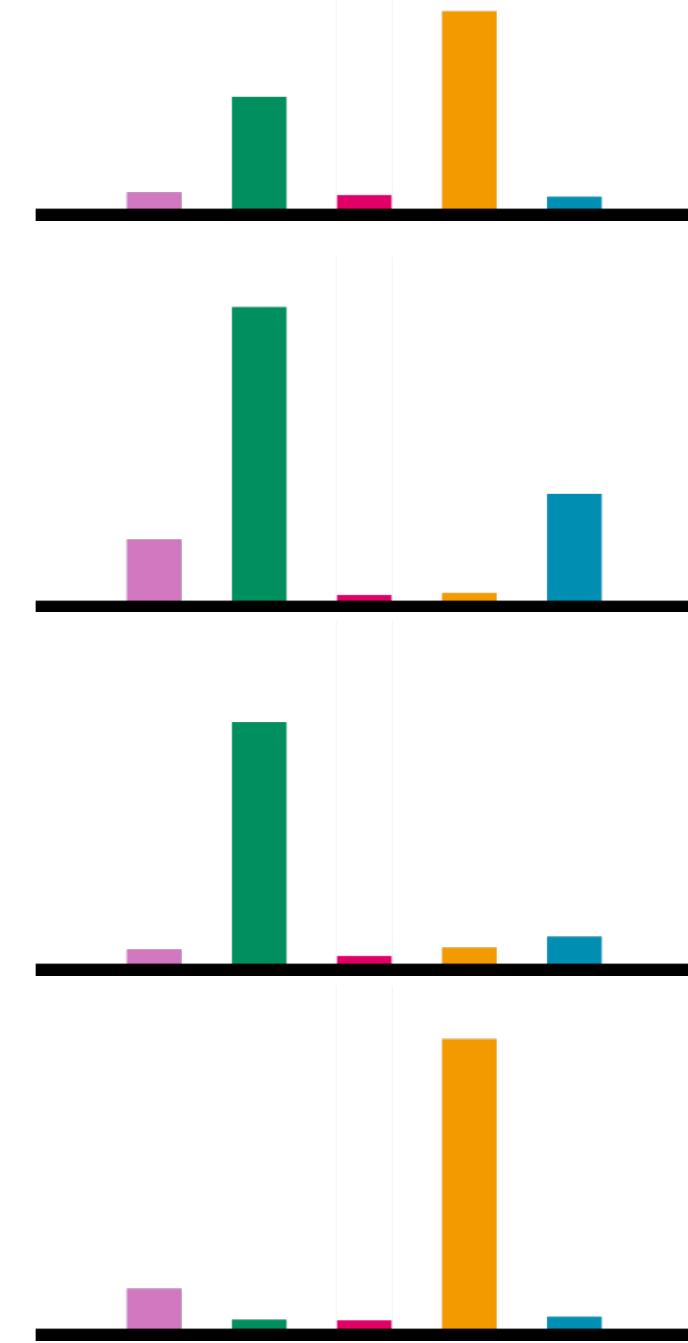
# modeling entities



typical concerns  
of the Bangkok Embassy



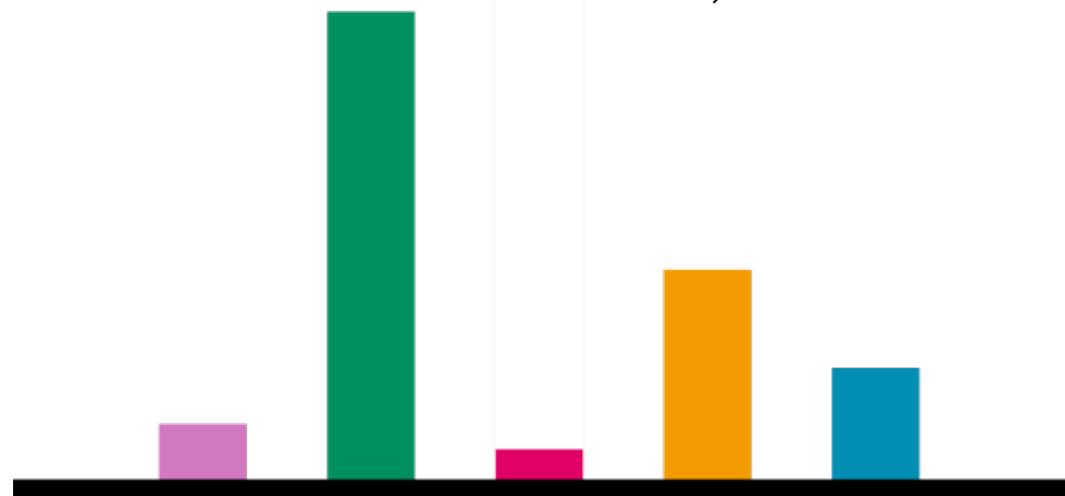
cables sent



# modeling entities

for each entity  $i$  and topic  $k$ , draw typical concerns:

$$\phi_{i,k} \sim \text{Gamma}(\alpha_\phi, \beta_\phi)$$



typical concerns  
of entity  $i$

# entity-only model of cables

for each entity  $i$  and topic  $k$ , draw typical concerns:

$$\phi_{i,k} \sim \text{Gamma}(\alpha_\phi, \beta_\phi)$$

# entity-only model of cables

for each entity  $i$  and topic  $k$ , draw typical concerns:

$$\phi_{i,k} \sim \text{Gamma}(\alpha_\phi, \beta_\phi)$$

for each topic  $k$  and vocabulary term  $v$ , draw topics:

$$\theta_{k,v} \sim \text{Gamma}(\alpha_\theta, \beta_\theta)$$

# entity-only model of cables

for each entity  $i$  and topic  $k$ , draw typical concerns:

$$\phi_{i,k} \sim \text{Gamma}(\alpha_\phi, \beta_\phi)$$

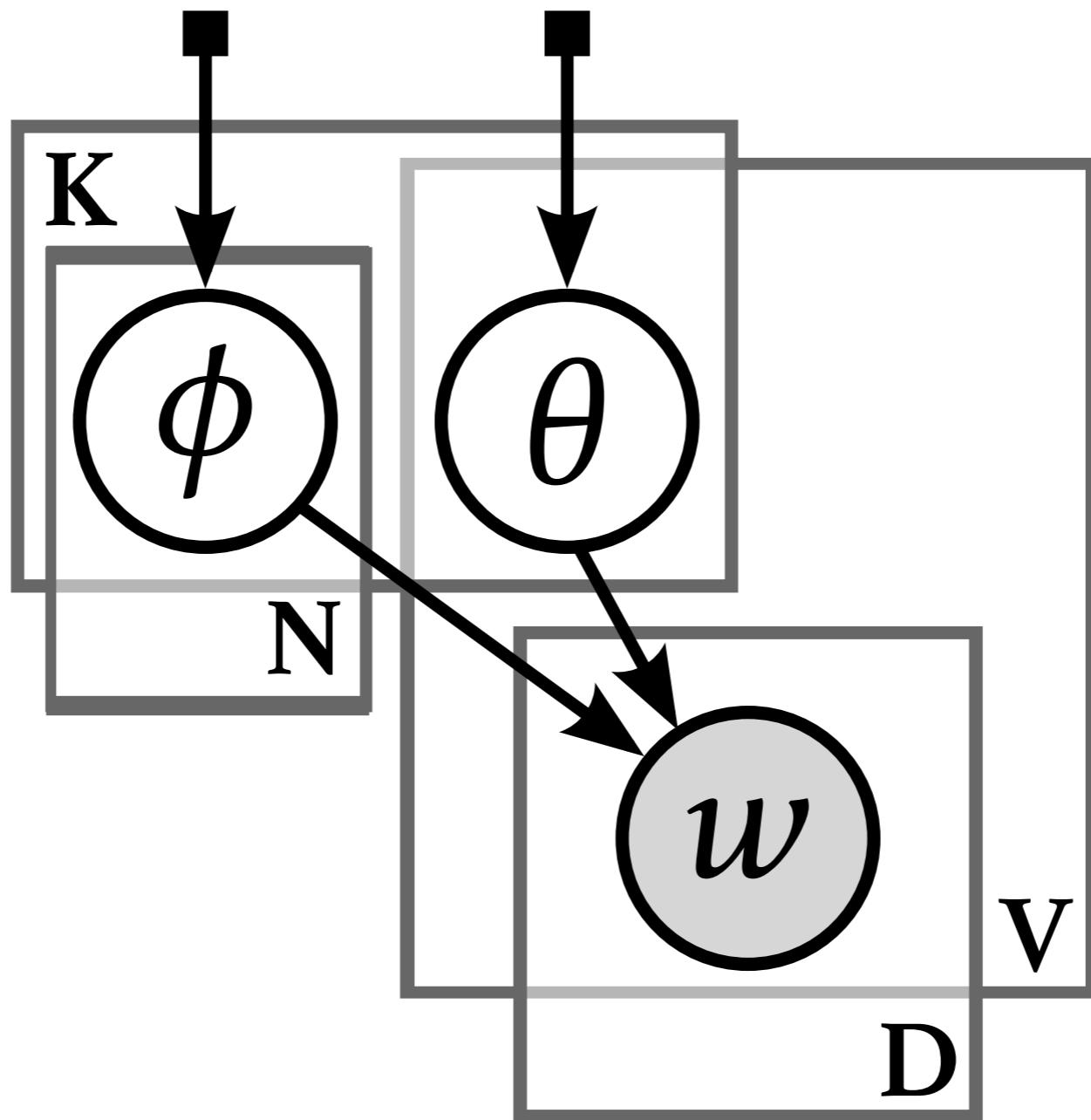
for each topic  $k$  and vocabulary term  $v$ , draw topics:

$$\theta_{k,v} \sim \text{Gamma}(\alpha_\theta, \beta_\theta)$$

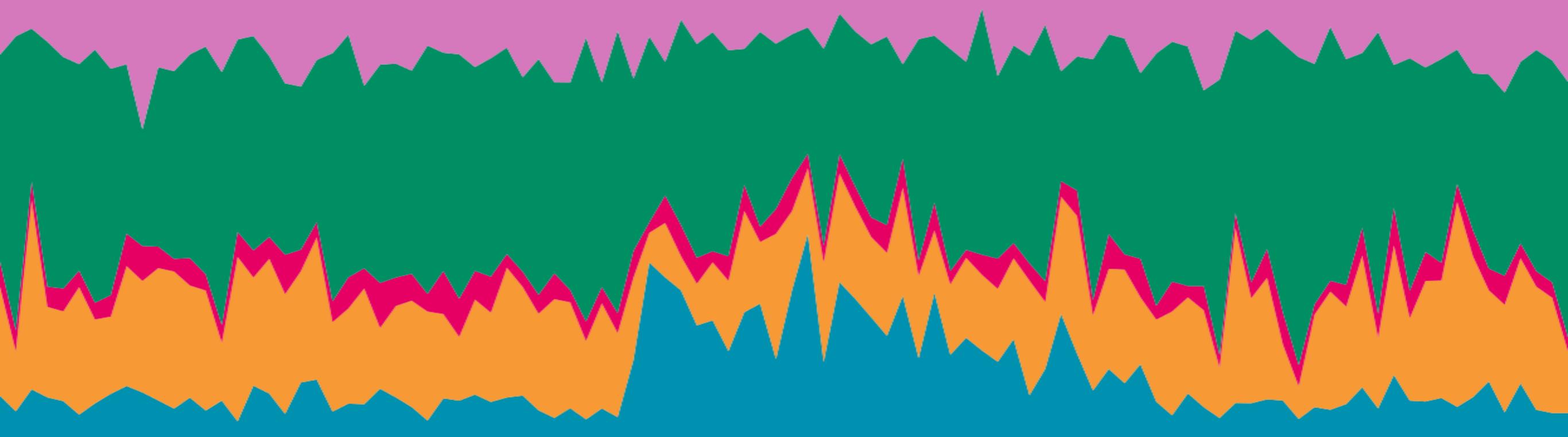
for each cable  $n$  (sent by entity  $i$ ) and vocabulary term  $v$ :

$$w_{n,v} \sim \text{Poisson} \left( \sum_k \phi_{i,k} \theta_{k,v} \right)$$

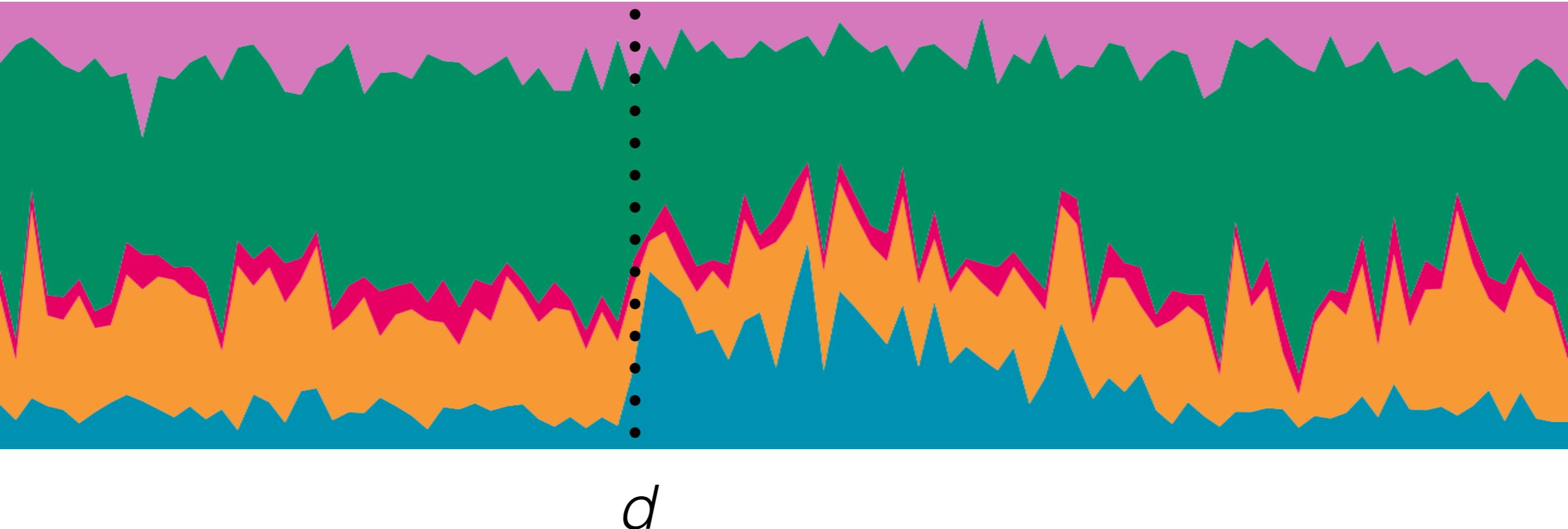
# entity-only model of cables



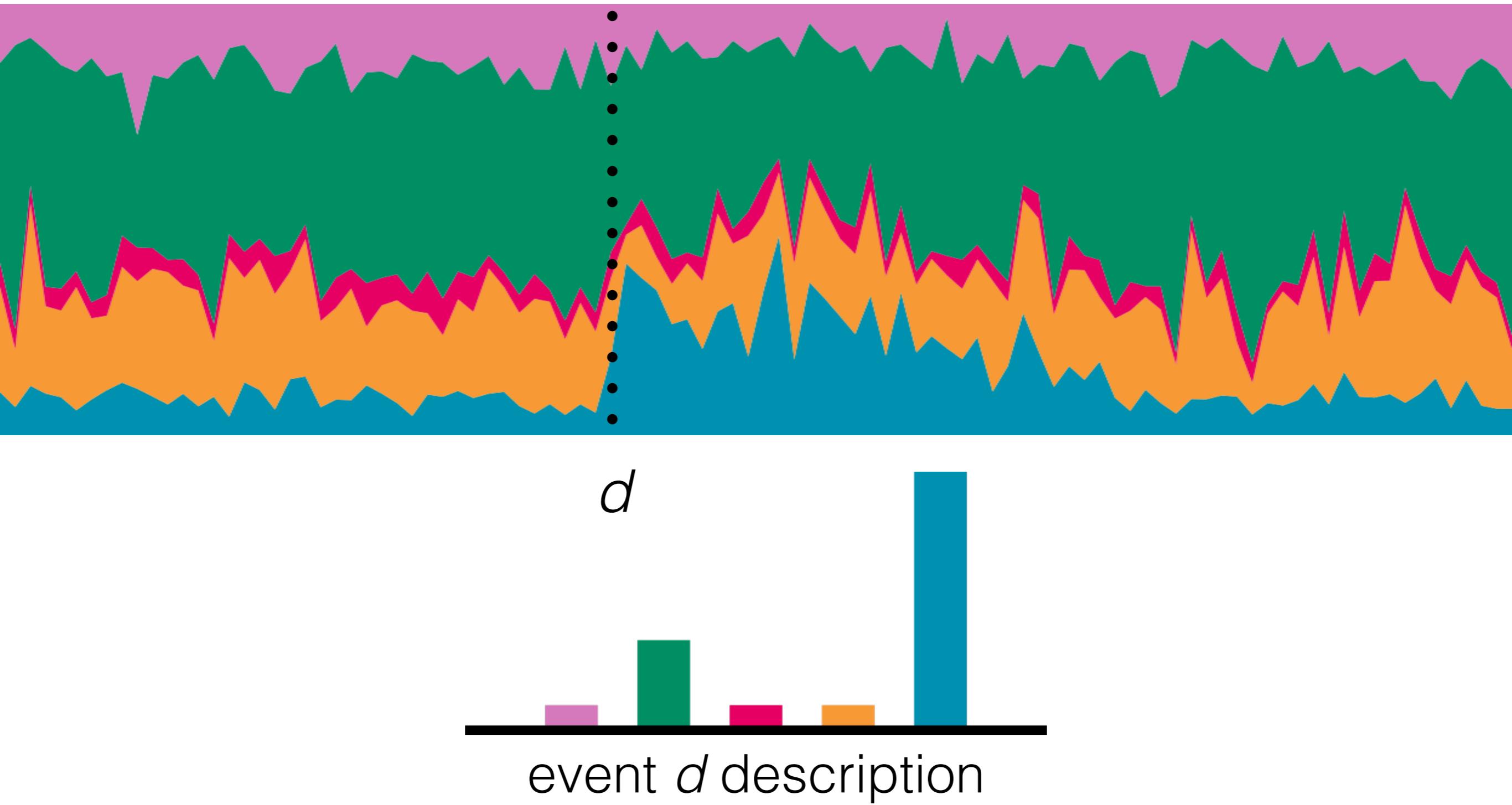
# modeling events



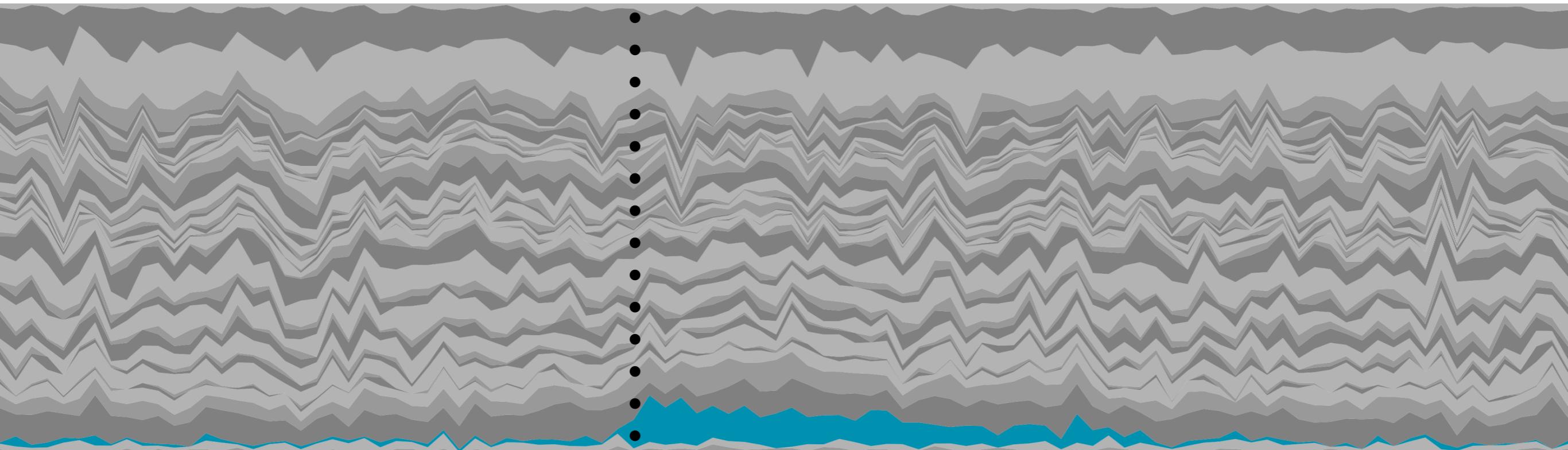
# modeling events



# modeling events



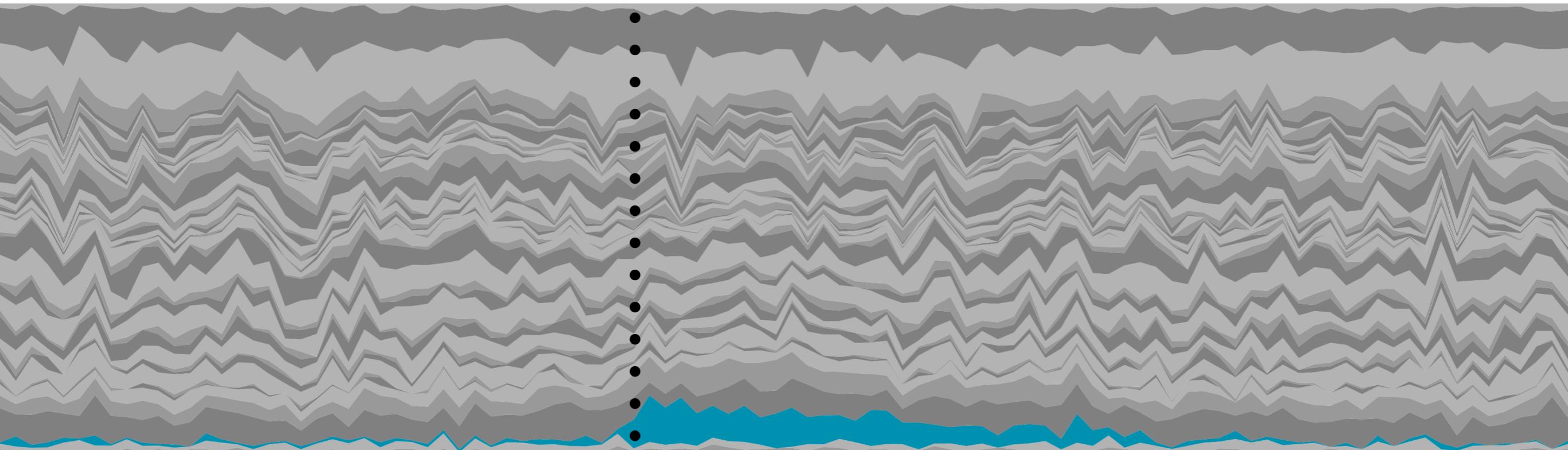
# modeling events



$d$



# modeling events



*d*



event *d* description

# event-only model of cables

for each time  $t$ , draw event strength:

$$\epsilon_t \sim \text{Gamma}(\alpha_\epsilon, \beta_\epsilon)$$

# event-only model of cables

for each time  $t$ , draw event strength:

$$\epsilon_t \sim \text{Gamma}(\alpha_\epsilon, \beta_\epsilon)$$

for each time  $t$  and vocab term  $v$ , draw event description:

$$\pi_{t,v} \sim \text{Gamma}(\alpha_\pi, \beta_\pi)$$

# event-only model of cables

for each time  $t$ , draw event strength:

$$\epsilon_t \sim \text{Gamma}(\alpha_\epsilon, \beta_\epsilon)$$

for each time  $t$  and vocab term  $v$ , draw event description:

$$\pi_{t,v} \sim \text{Gamma}(\alpha_\pi, \beta_\pi)$$

for each cable  $n$  (sent at time  $d$ ) and vocabulary term  $v$ :

$$w_{n,v} \sim \text{Poisson} \left( \sum_t f(t, d) \epsilon_t, \pi_{t,v} \right)$$

# event-only model of cables

for each time  $t$ , draw event strength:

$$\epsilon_t \sim \text{Gamma}(\alpha_\epsilon, \beta_\epsilon)$$

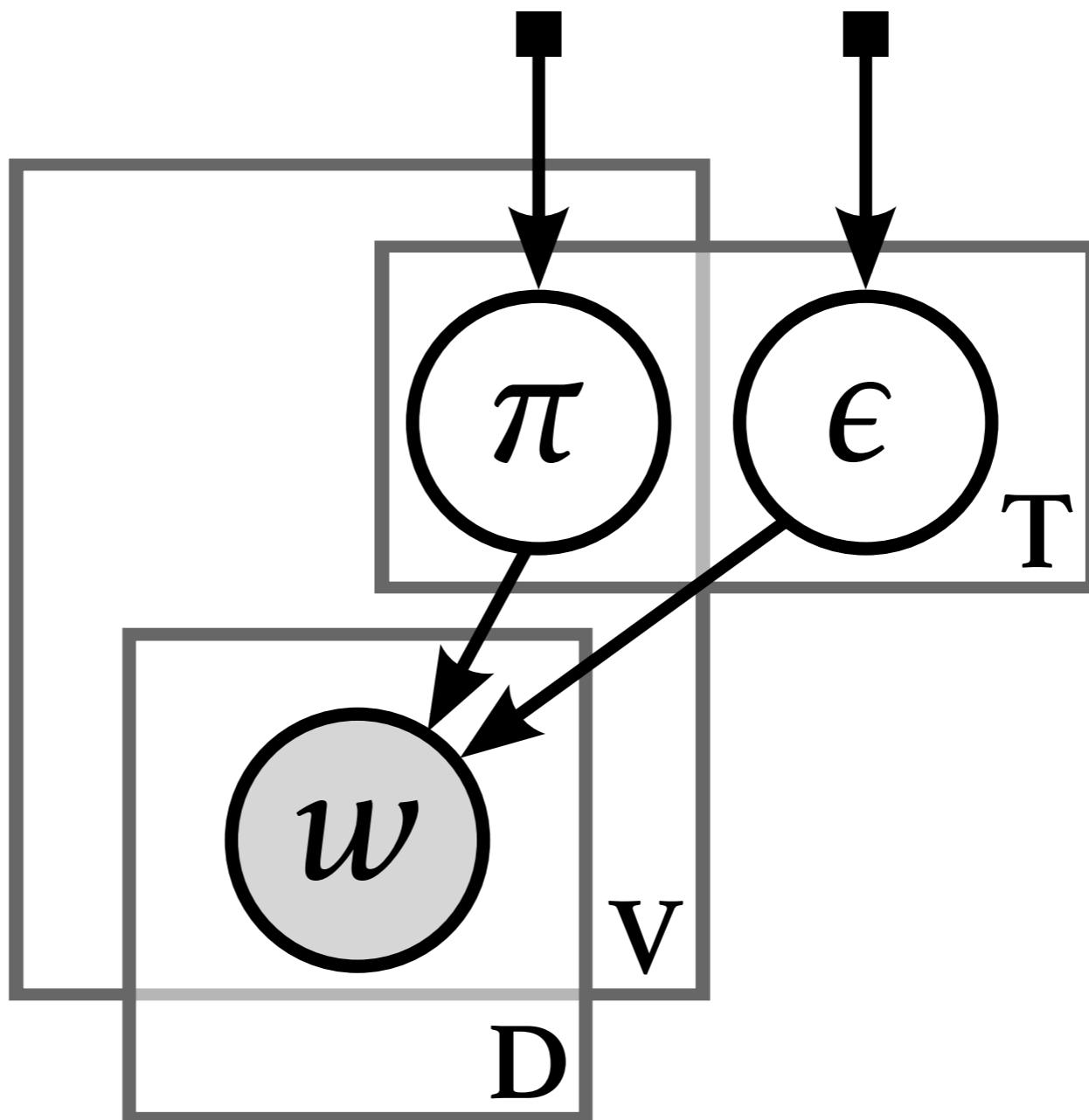
for each time  $t$  and vocab term  $v$ , draw event description:

$$\pi_{t,v} \sim \text{Gamma}(\alpha_\pi, \beta_\pi)$$

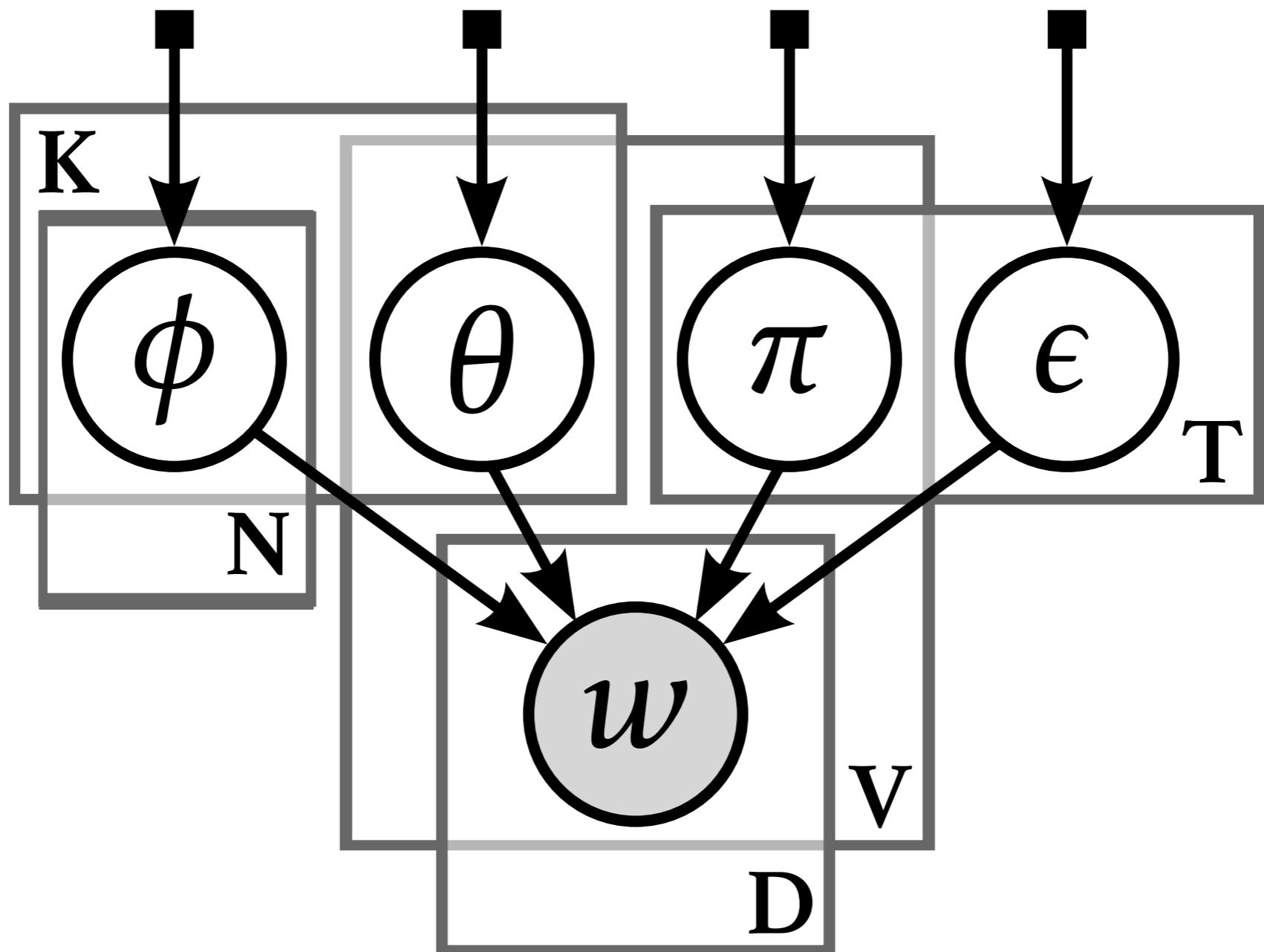
for each cable  $n$  (sent at time  $d$ ) and vocabulary term  $v$ :

$$w_{n,v} \sim \text{Poisson} \left( \sum_t f(t, d) \epsilon_t, \pi_{t,v} \right)$$

# event-only model of cables



# full model of cables



# full model of cables

for each cable  $n$  (sent by entity  $i$  at time  $d$ ) and vocab term  $v$ :

$$w_{n,v} \sim \text{Poisson} \left( \sum_k \phi_{i,k} \theta_{k,v} + \sum_t f(t, d) \epsilon_t, \pi_{t,v} \right)$$

# full model of cables

for each cable  $n$  (sent by entity  $i$  at time  $d$ ) and vocab term  $v$ :

sum over topics

$$w_{n,v} \sim \text{Poisson} \left( \sum_k \phi_{i,k} \theta_{k,v} + \sum_t f(t, d) \epsilon_t, \pi_{t,v} \right)$$


# full model of cables

for each cable  $n$  (sent by entity  $i$  at time  $d$ ) and vocab term  $v$ :

sum over topics

$$w_{n,v} \sim \text{Poisson} \left( \sum_k \phi_{i,k} \theta_{k,v} + \sum_t f(t,d) \epsilon_t, \pi_{t,v} \right)$$

The diagram illustrates the components of the Poisson distribution. On the left, there is a horizontal axis labeled "typical concerns of entity  $i$ ". Above this axis, there is a bar chart with four bars of different heights and colors: purple, green, orange, and teal. An arrow points from this bar chart towards the summation term  $\sum_k \phi_{i,k} \theta_{k,v}$  in the Poisson formula. Another arrow points from the entire formula towards the right side of the equation.

# full model of cables

for each cable  $n$  (sent by entity  $i$  at time  $d$ ) and vocab term  $v$ :

sum over topics

$$w_{n,v} \sim \text{Poisson} \left( \sum_k \phi_{i,k} \theta_{k,v} + \sum_t f(t,d) \epsilon_t, \pi_{t,v} \right)$$

typical concerns of entity  $i$

topics for word  $v$

# full model of cables

for each cable  $n$  (sent by entity  $i$  at time  $d$ ) and vocab term  $v$ :

$$w_{n,v} \sim \text{Poisson} \left( \sum_k \phi_{i,k} \theta_{k,v} + \sum_t f(t, d) \epsilon_t, \pi_{t,v} \right)$$

# full model of cables

for each cable  $n$  (sent by entity  $i$  at time  $d$ ) and vocab term  $v$ :

1973 1978  
sum over all time steps

$$w_{n,v} \sim \text{Poisson} \left( \sum_k \phi_{i,k} \theta_{k,v} + \sum_t f(t, d) \epsilon_t, \pi_{t,v} \right)$$

# full model of cables

for each cable  $n$  (sent by entity  $i$  at time  $d$ ) and vocab term  $v$ :

$$w_{n,v} \sim \text{Poisson} \left( \sum_k \phi_{i,k} \theta_{k,v} + \sum_t f(t, d) \epsilon_t, \pi_{t,v} \right)$$

1973      1978  
sum over all time steps

decay of relevancy

# full model of cables

for each cable  $n$  (sent by entity  $i$  at time  $d$ ) and vocab term  $v$ :

$$w_{n,v} \sim \text{Poisson} \left( \sum_k \phi_{i,k} \theta_{k,v} + \sum_t f(t, d) \epsilon_t, \pi_{t,v} \right)$$

event  $t$   
strength

decay of relevancy

The diagram illustrates the components of the Poisson distribution formula. At the top, a horizontal dashed line with arrows at both ends spans from '1973' to '1978'. Below this line, the text 'sum over all time steps' is centered. To the right of the formula, an arrow points from the term  $\sum_t f(t, d) \epsilon_t$  to the text 'event  $t$  strength'. Below the formula, a green triangle is shown on a horizontal baseline, representing the 'decay of relevancy'.

# full model of cables

for each cable  $n$  (sent by entity  $i$  at time  $d$ ) and vocab term  $v$ :

$$w_{n,v} \sim \text{Poisson} \left( \sum_k \phi_{i,k} \theta_{k,v} + \sum_t f(t, d) \epsilon_t, \pi_{t,v} \right)$$

event  $t$   
strength

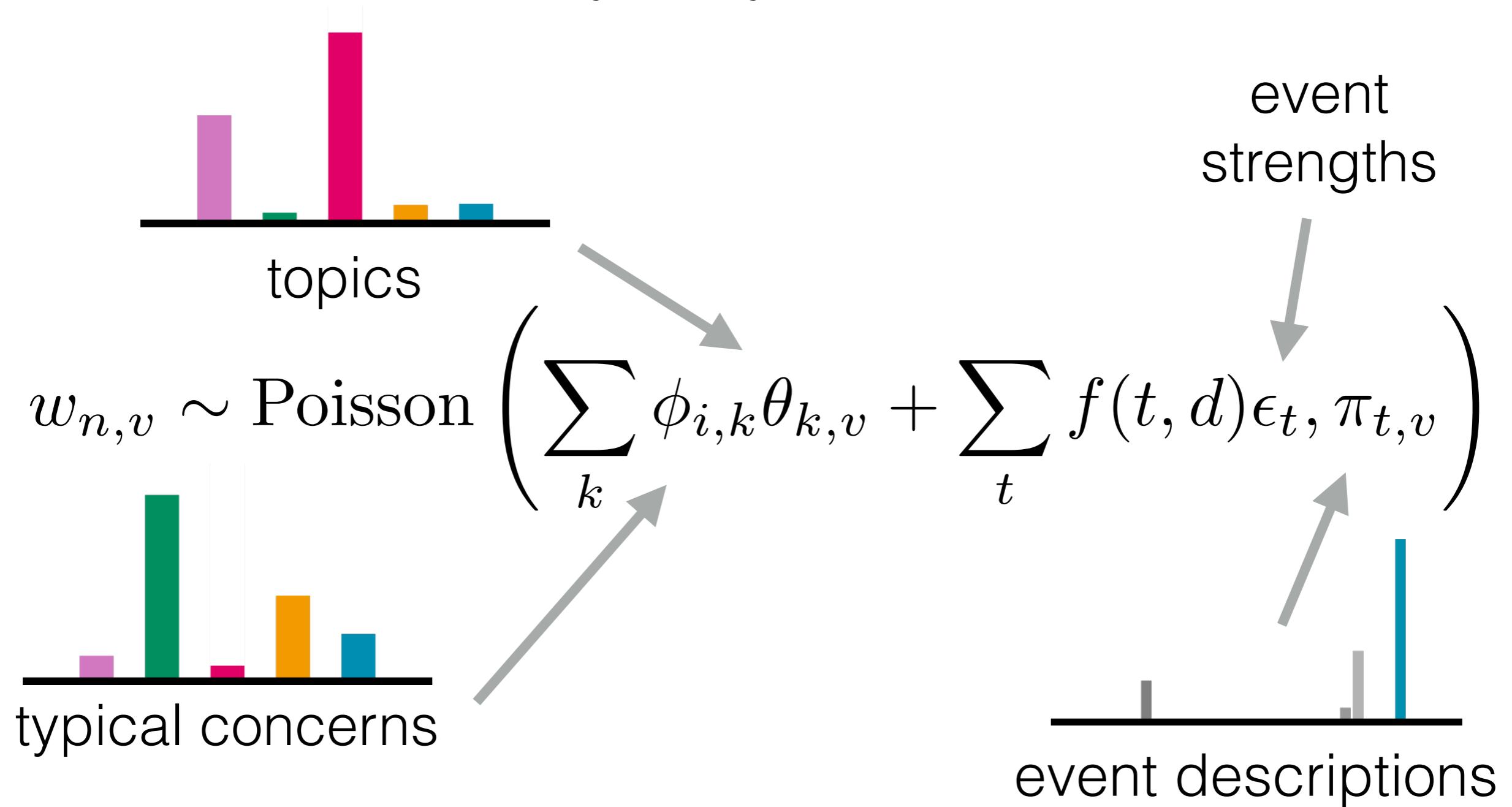
decay of relevancy

event descriptions  
in vocabulary space

1973      1978  
sum over all time steps

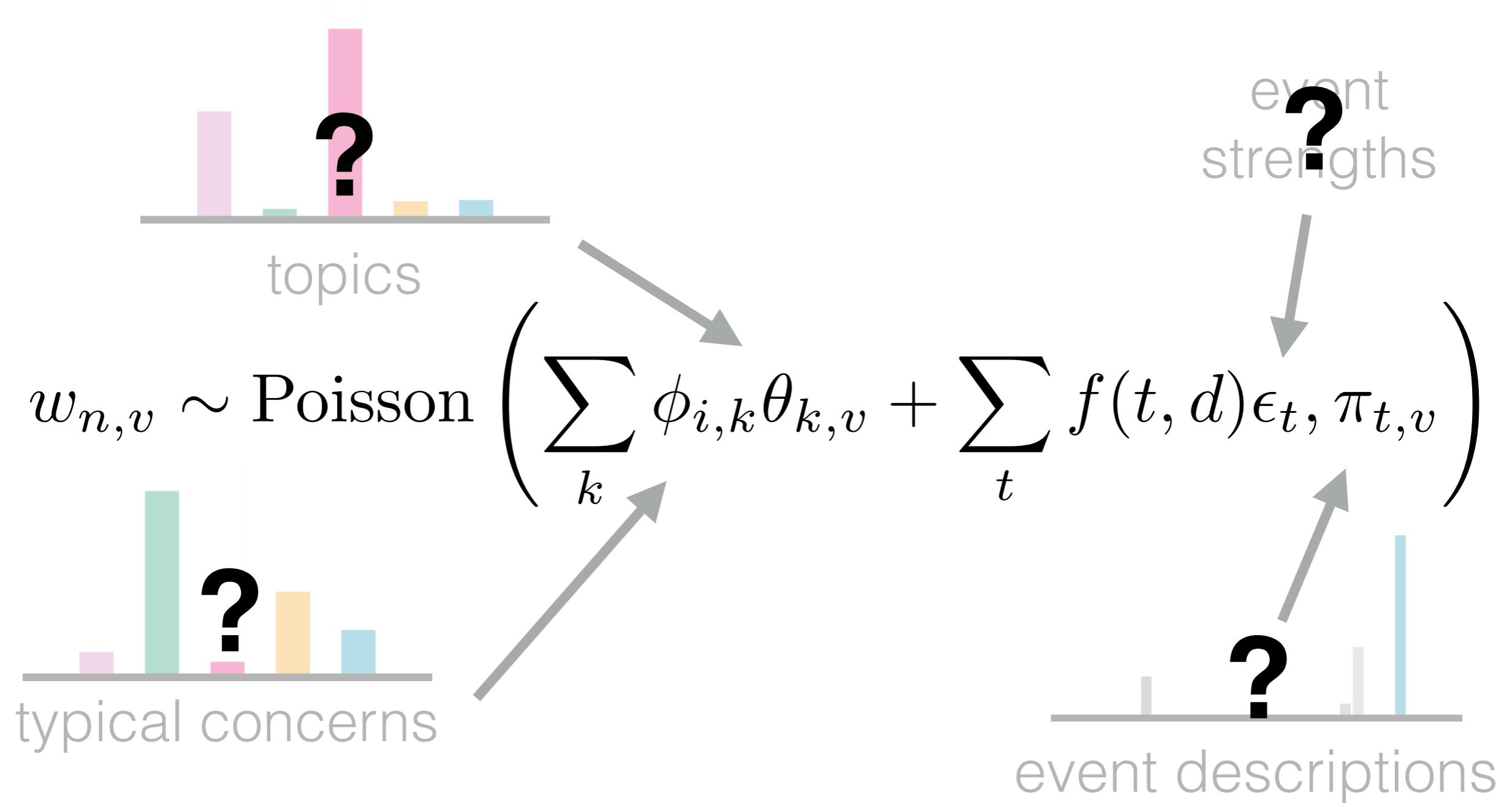
# full model of cables

for each cable  $n$  (sent by entity  $i$  at time  $d$ ) and vocab term  $v$ :

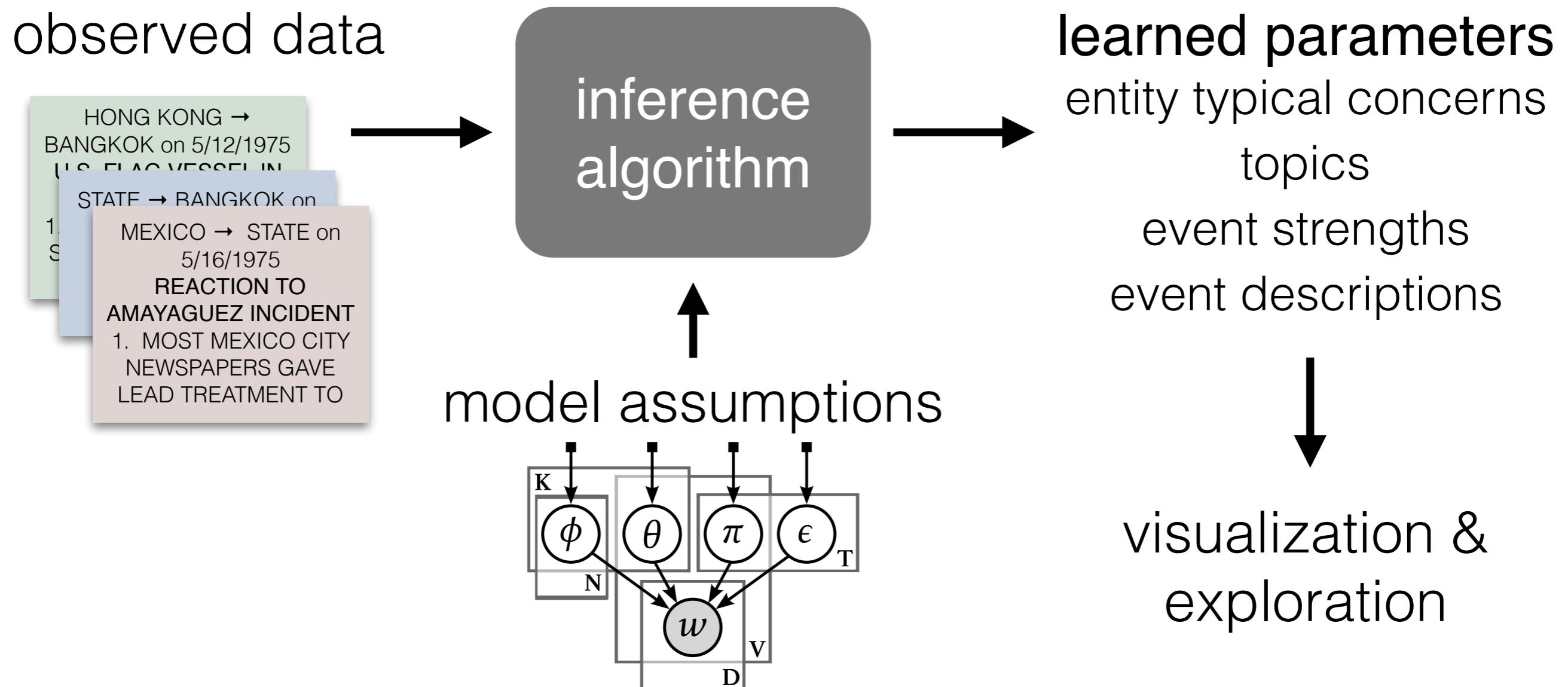


# full model of cables

for each cable  $n$  (sent by entity  $i$  at time  $d$ ) and vocab term  $v$ :



# How do we find the values of the hidden parameters that best fit the data?



# Posterior Distribution

$$p(\phi, \theta, \epsilon, \pi \mid w, \alpha, \beta) = \frac{p(\phi, \theta, \epsilon, \pi, w \mid \alpha, \beta)}{\int_{\phi} \int_{\theta} \int_{\epsilon} \int_{\pi} p(\phi, \theta, \epsilon, \pi, w \mid \alpha, \beta)}$$

latent model parameters

easy to compute

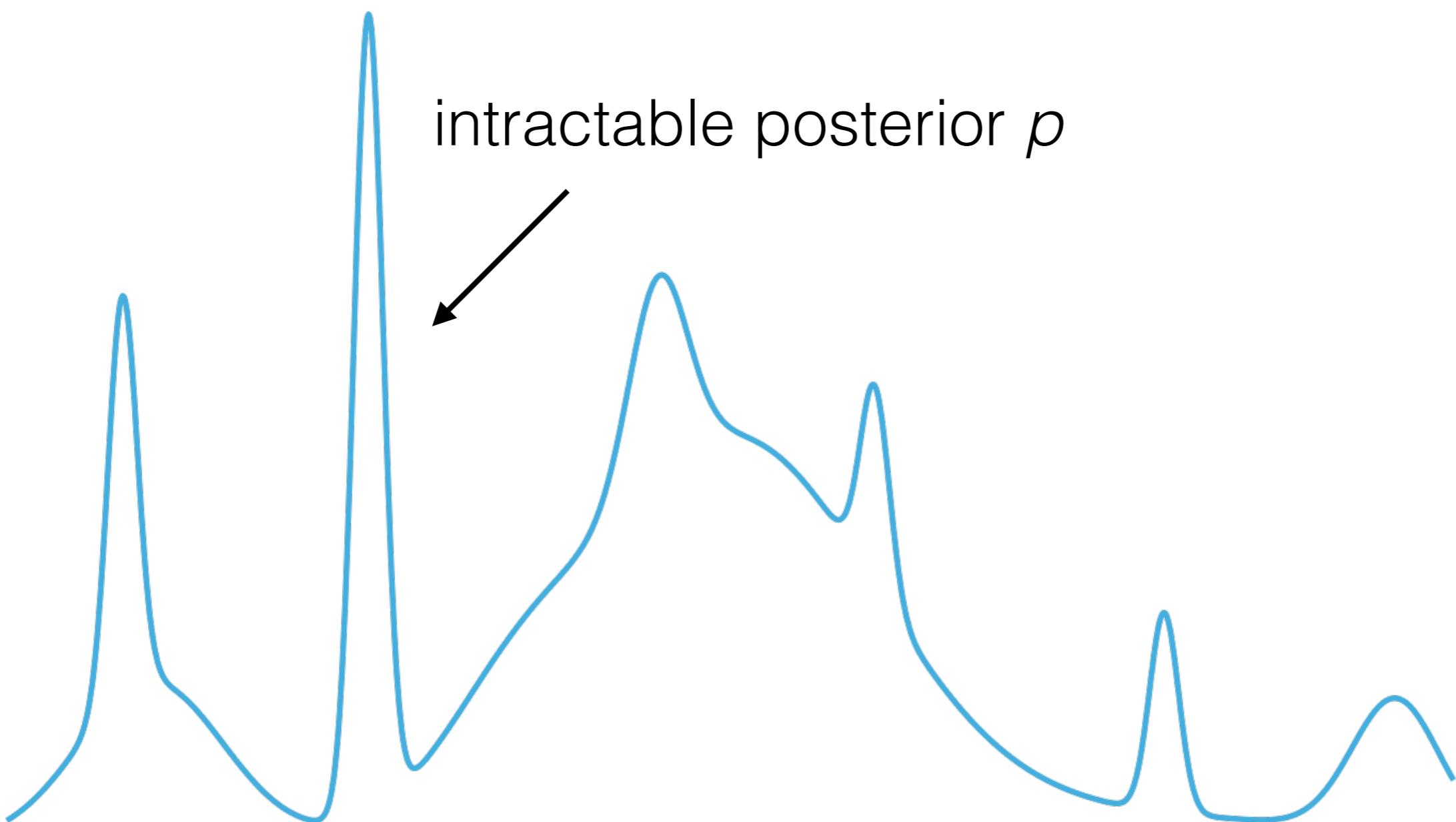
observed data

model hyperparameters

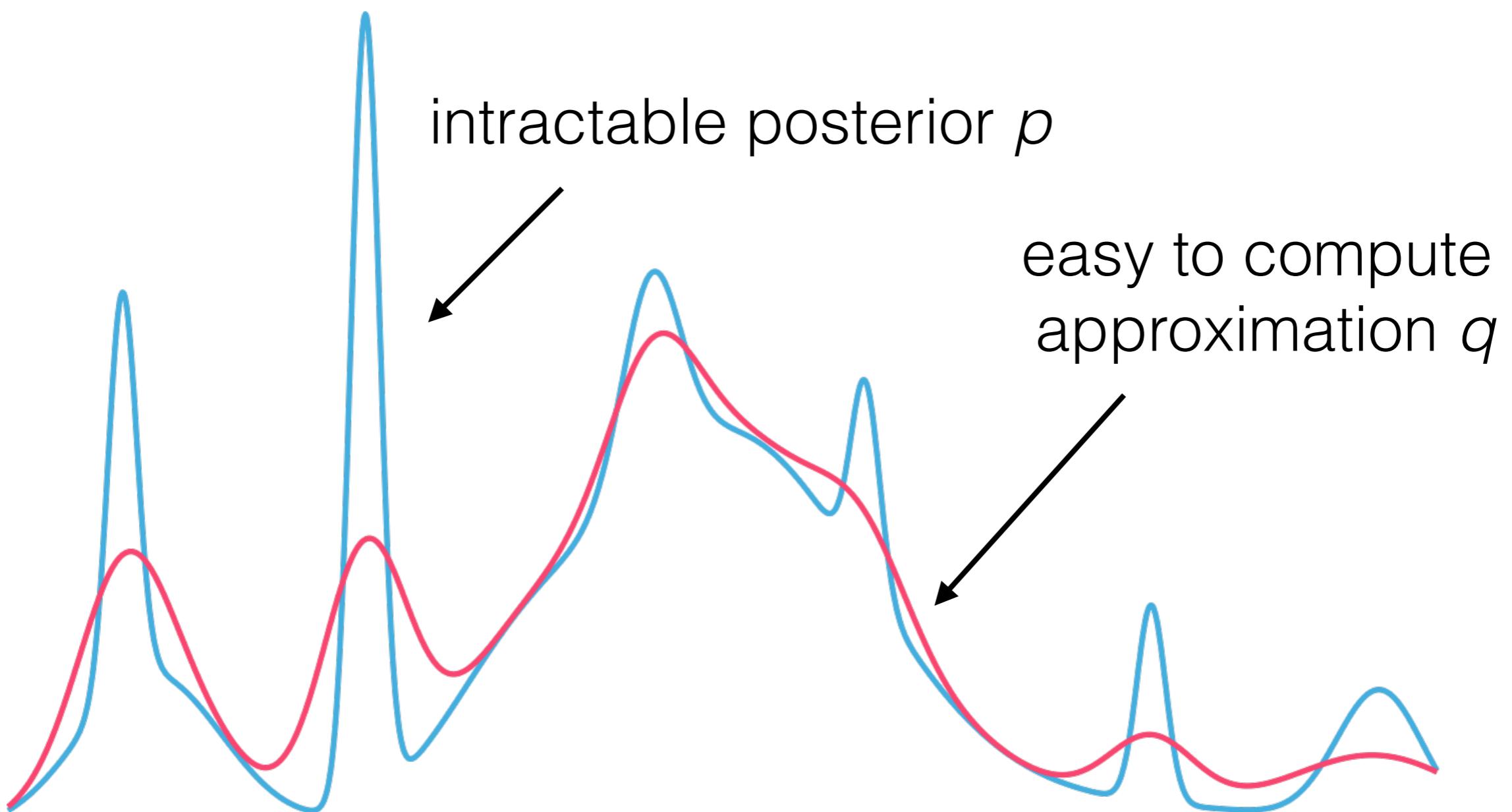
intractable

The diagram illustrates the posterior distribution formula. It features a central fraction where the numerator is  $p(\phi, \theta, \epsilon, \pi, w \mid \alpha, \beta)$  and the denominator is  $\int_{\phi} \int_{\theta} \int_{\epsilon} \int_{\pi} p(\phi, \theta, \epsilon, \pi, w \mid \alpha, \beta)$ . Four arrows point to specific parts of the formula: one from the left labeled "latent model parameters" points to the set of parameters  $(\phi, \theta, \epsilon, \pi)$  in the numerator; another from the right labeled "easy to compute" points to the same set in the denominator; a third arrow from the bottom left labeled "observed data" points to the variable  $w$ ; and a fourth arrow from the bottom right labeled "intractable" points to the denominator integral.

# Variational Inference



# Variational Inference



# Variational Inference

Standard conjugate variational inference

- model-specific mathematical derivations of closed-form updates
- uses coordinate ascent: iteratively updates each parameter while holding the others fixed

*Black box variational inference.* Ranganath, Gerrish, and Blei, 2014.

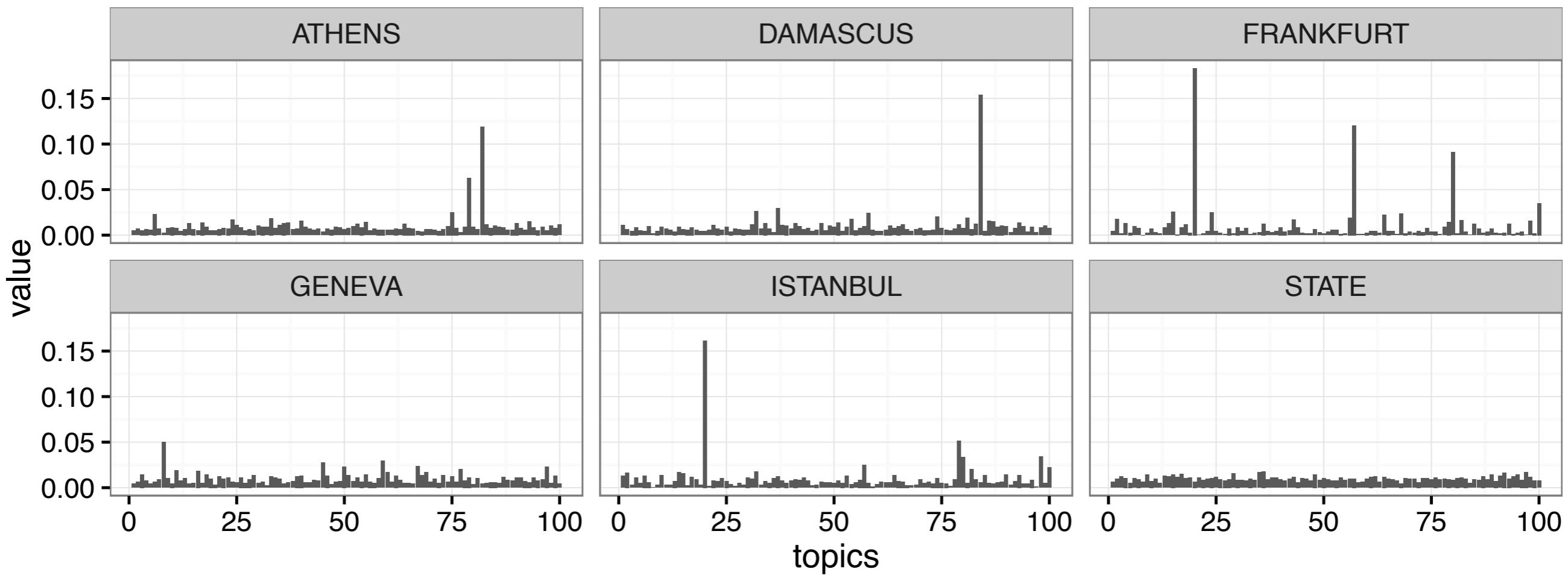
- reduces model-specific mathematical derivations
- uses stochastic optimization
  - noisy gradient is computed from Monte Carlo samples from the variational distribution
- downsides: looser fit, slower to converge

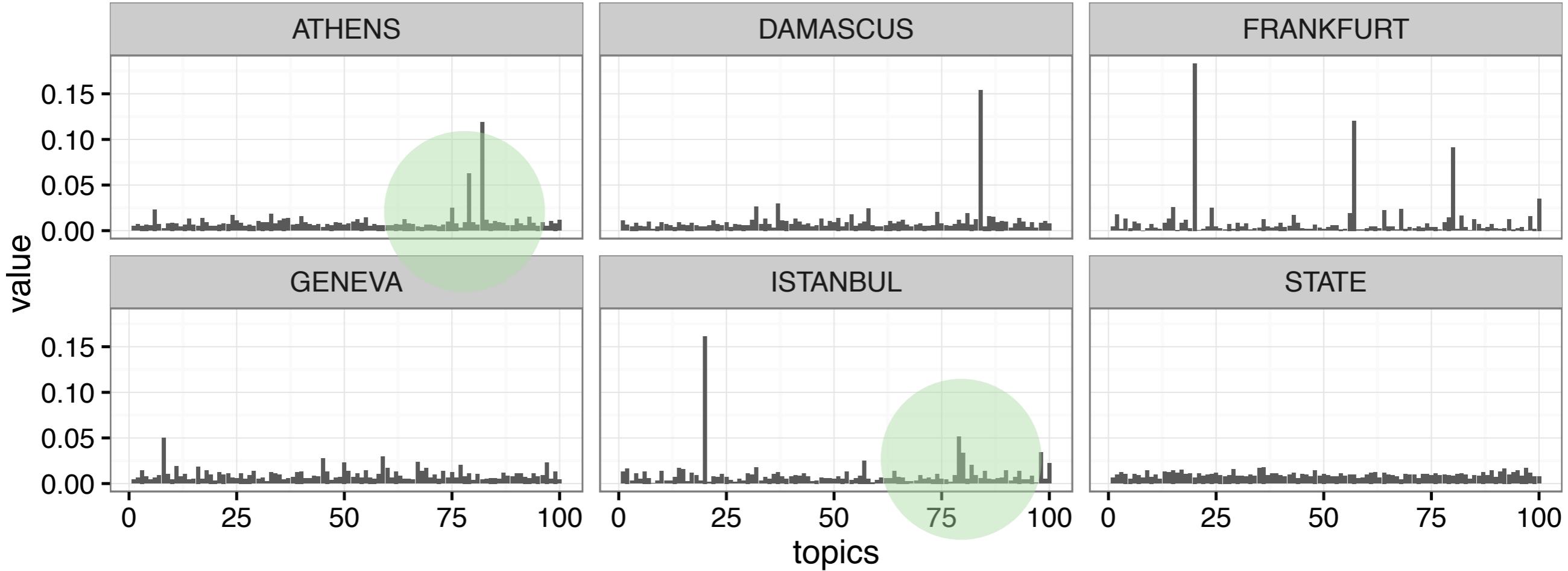
# Validation

- compare discovered events to manually collected examples of known historical events (and corresponding cables)
  - How many of the known events are recovered?
  - How does the average distributions of the known cables compares to the discovered event distribution?
- present the discovered events (date, topic distribution, and entities involved) to an expert historian

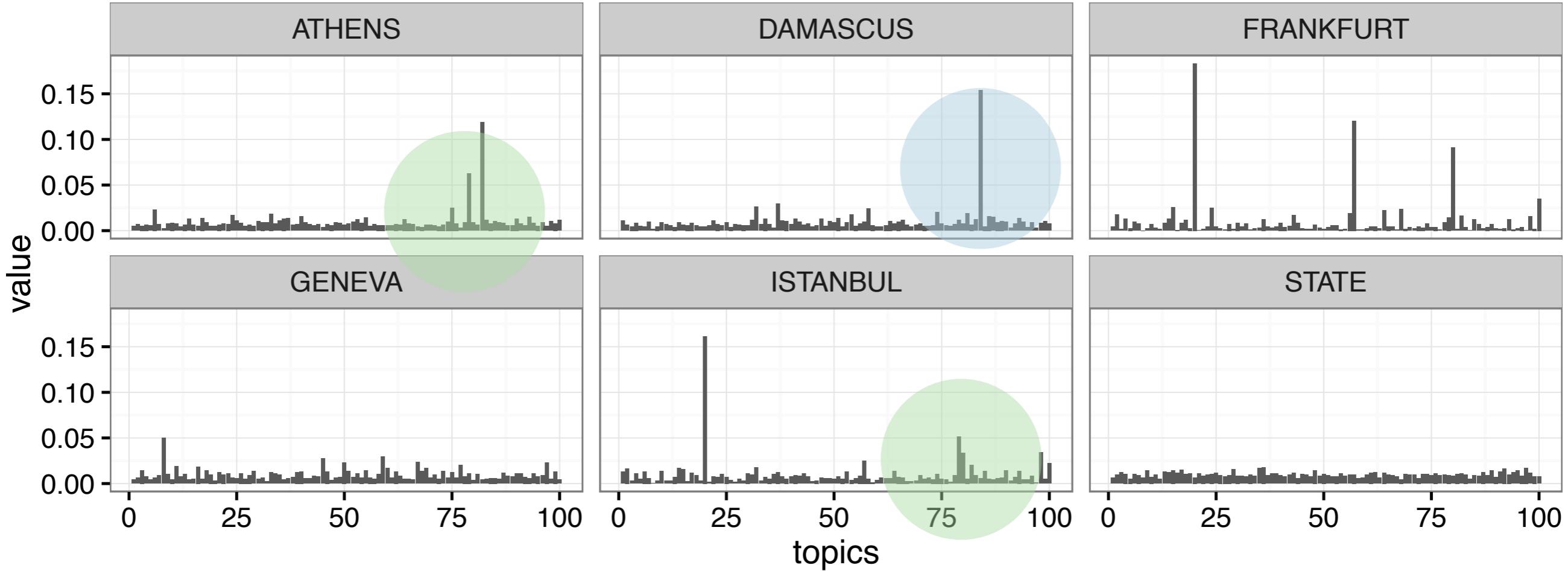
# Exploratory Results

## 1976



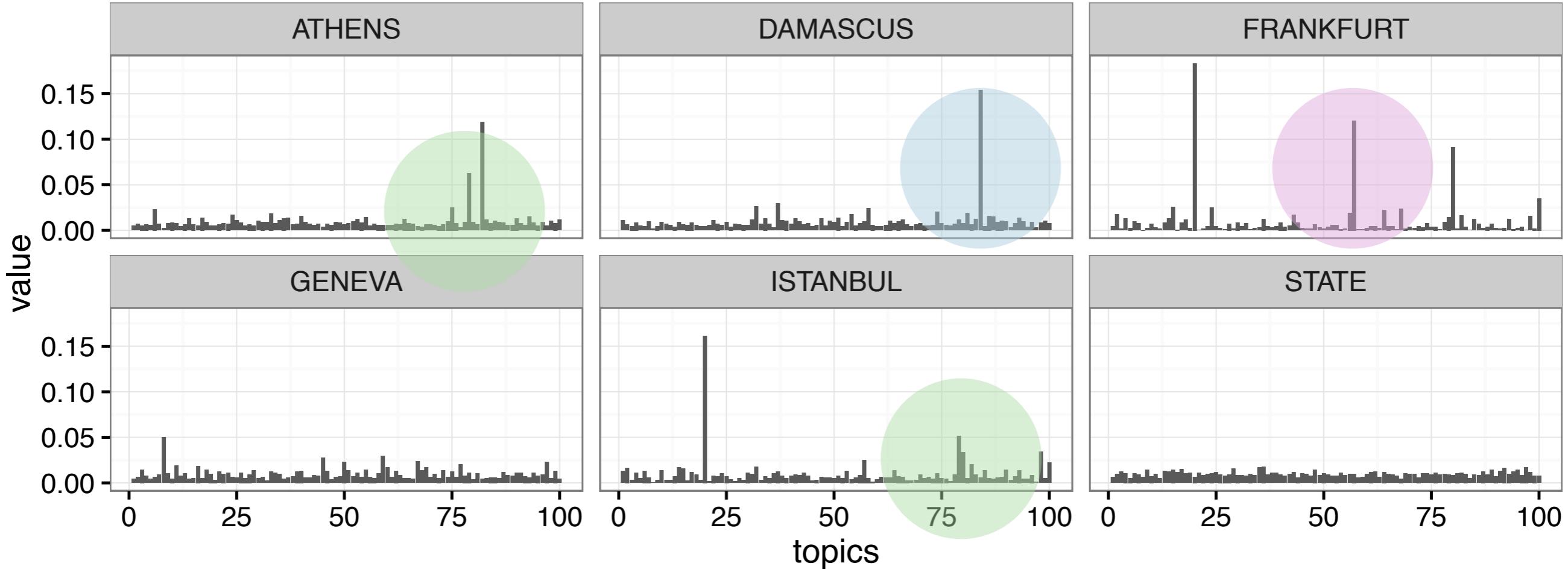


TURKISH  
TURKEY  
GREEK  
ANKARA  
GREECE



TURKISH  
TURKEY  
GREEK  
ANKARA  
GREECE

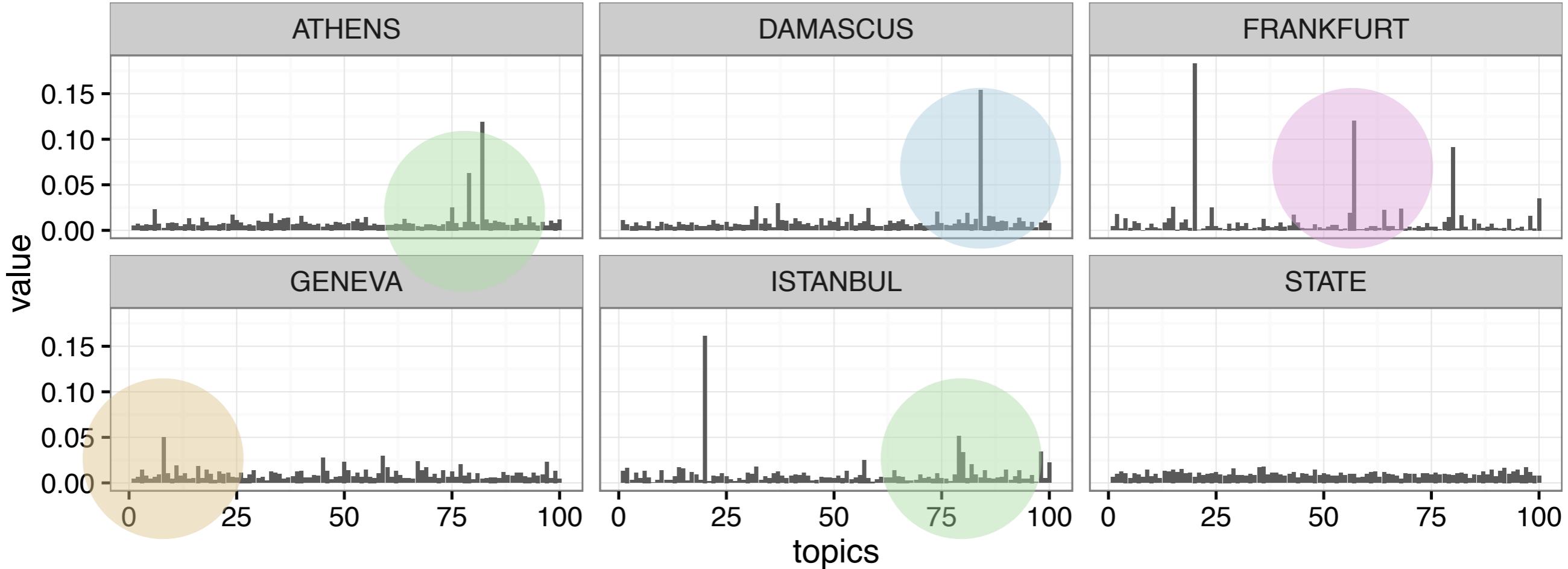
BEIRUT  
DAMASCUS  
SYRIAN  
LEBANON  
ARAB



TURKISH  
TURKEY  
GREEK  
ANKARA  
GREECE

BEIRUT  
DAMASCUS  
SYRIAN  
LEBANON  
ARAB

BONN  
GERMAN  
FRG  
GERMANY  
BERLIN



TURKISH  
TURKEY  
GREEK  
ANKARA  
GREECE

BEIRUT  
DAMASCUS  
SYRIAN  
LEBANON  
ARAB

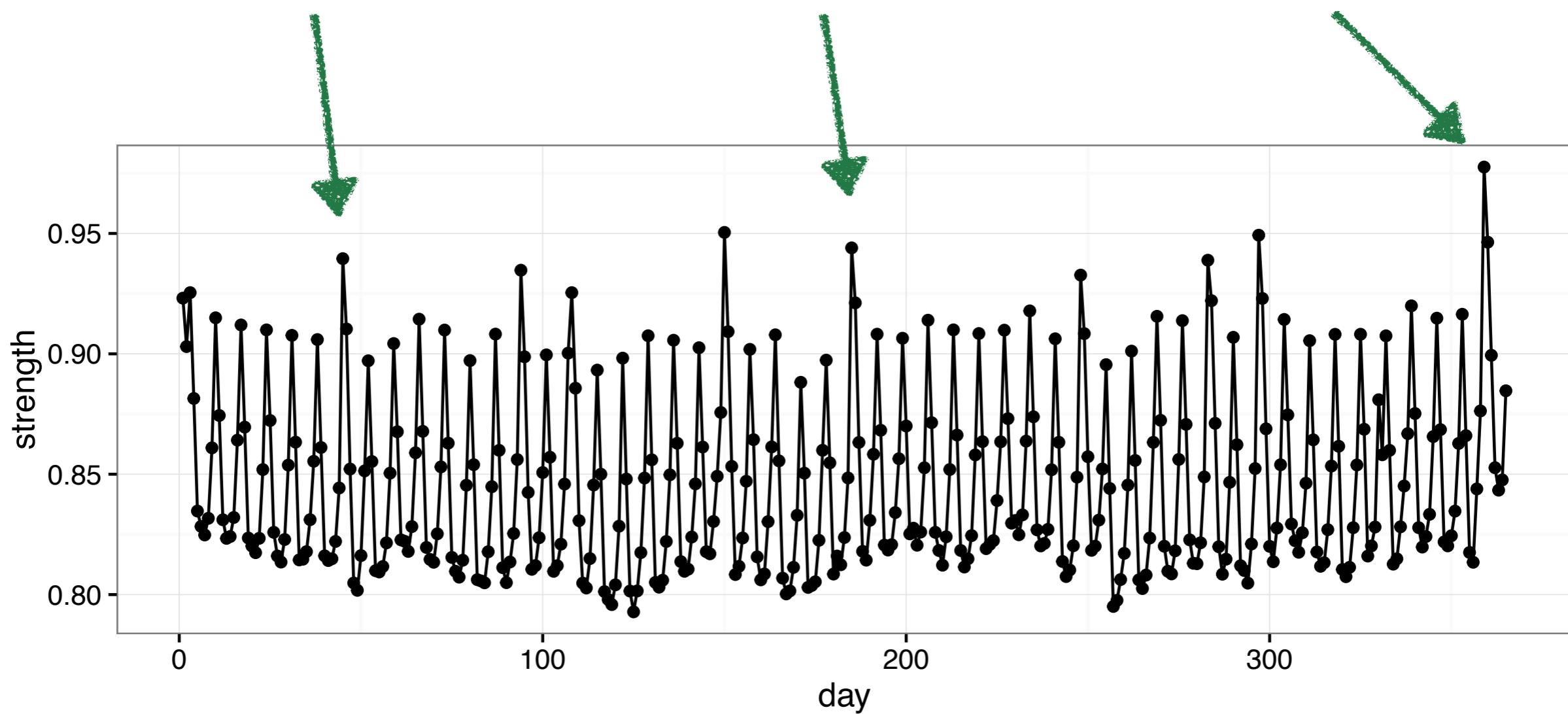
BONN  
GERMAN  
FRG  
GERMANY  
BERLIN

USUN  
GENEVA  
SECRETARIAT  
DRAFT  
SESSION

FOLLOW  
DHABI  
REQUEST  
RTC  
FRANKFURT

ACTIVITY  
HOPE  
COMPLICATE  
PROMOTION  
APPEAR

AWARE  
FLY  
DO  
GSDR  
DOUBT



# Current Work

# Current Work

- Goals:

# Current Work

- Goals:
  - show that our model outperforms heuristic baselines on simulated data

# Current Work

- Goals:
  - show that our model outperforms heuristic baselines on simulated data
  - explore model sensitivity to hyperparameters

# Current Work

- Goals:
  - show that our model outperforms heuristic baselines on simulated data
  - explore model sensitivity to hyperparameters
  - explore more real-world data with the model

# Current Work

- Goals:
  - show that our model outperforms heuristic baselines on simulated data
  - explore model sensitivity to hyperparameters
  - explore more real-world data with the model
- Model extensions

# Current Work

- Goals:
  - show that our model outperforms heuristic baselines on simulated data
  - explore model sensitivity to hyperparameters
  - explore more real-world data with the model
- Model extensions
  - hierarchical model: document specific parameters

# Current Work

- Goals:
  - show that our model outperforms heuristic baselines on simulated data
  - explore model sensitivity to hyperparameters
  - explore more real-world data with the model
- Model extensions
  - hierarchical model: document specific parameters
  - include interactions between entities

# Current Work

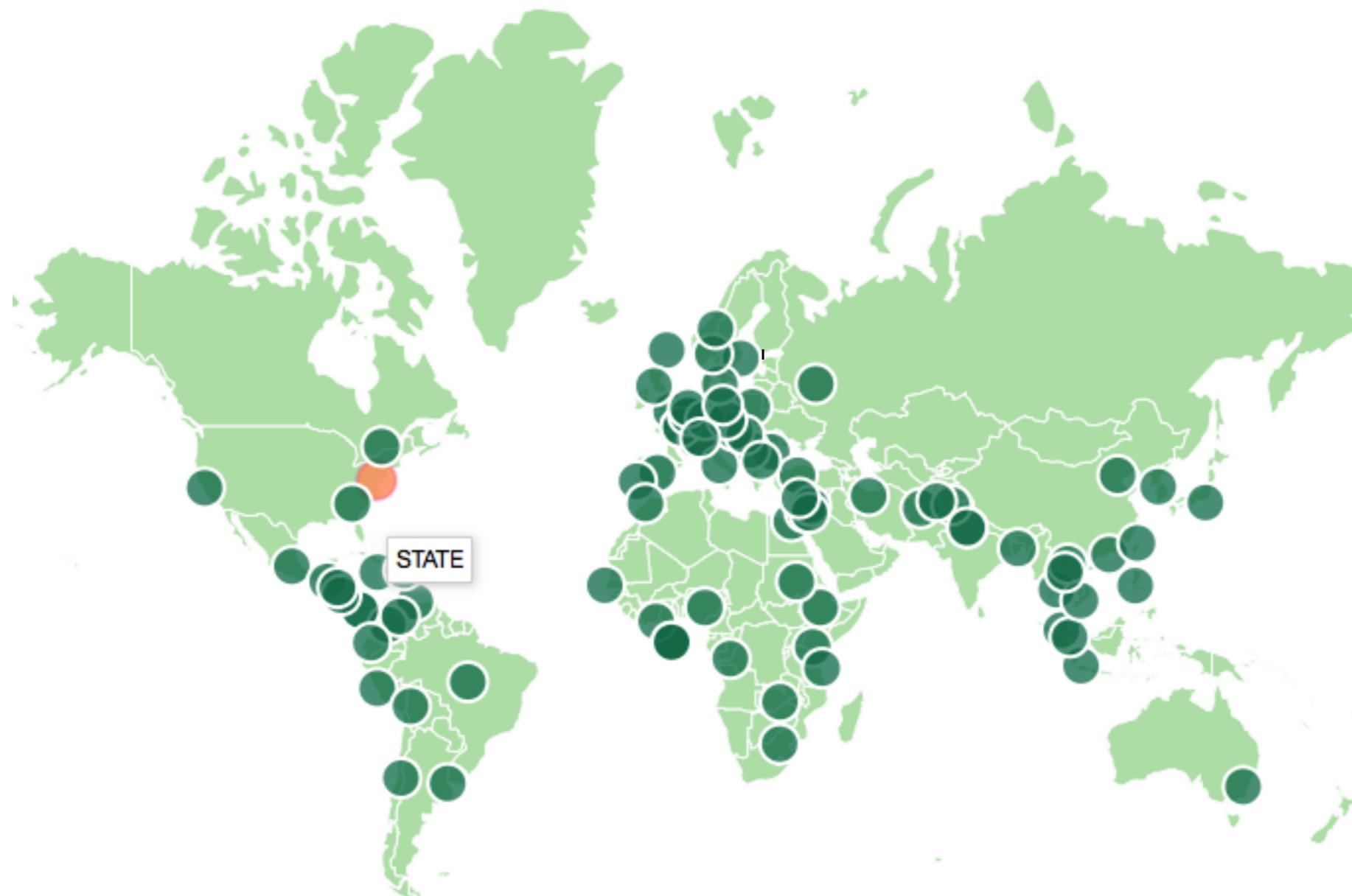
- Goals:
  - show that our model outperforms heuristic baselines on simulated data
  - explore model sensitivity to hyperparameters
  - explore more real-world data with the model
- Model extensions
  - hierarchical model: document specific parameters
  - include interactions between entities
  - learn event duration

# Current Work

- Goals:
  - show that our model outperforms heuristic baselines on simulated data
  - explore model sensitivity to hyperparameters
  - explore more real-world data with the model
- Model extensions
  - hierarchical model: document specific parameters
  - include interactions between entities
  - learn event duration
  - explore different event decay shapes

# Visualization

## Entities Overview



**Thank you!**  
Questions and suggestions welcome.

# Complete Conditionals

$$\phi_{i,k} \mid \theta, \epsilon, \pi, \mathbf{W} \sim \text{Gamma} \left( \alpha_\phi + \sum_{v,n \in N_i} z_{n,v,k}^{entity}, \beta_\phi + |N_i| \sum_v \theta_{k,v} \right)$$

$$\theta_{k,v} \mid \phi, \epsilon, \pi, \mathbf{W} \sim \text{Gamma} \left( \alpha_\theta + \sum_{v,n} z_{n,v,k}^{entity}, \beta_\theta + \sum_i |N_i| \phi_{i,k} \right)$$

$$\epsilon_t \mid \phi, \theta, \pi, \mathbf{W} \sim \text{Gamma} \left( \alpha_\epsilon + \sum_{v,n \in N_t} z_{n,v,t}^{event}, \beta_\epsilon + \sum_v \sum_{d=0}^{\delta} |N_{t+d}| f(t+d, t) \sum_k \pi_{t,k} \right)$$

$$\pi_{t,v} \mid \phi, \theta, \epsilon, \mathbf{W} \sim \text{Gamma} \left( \alpha_\pi + \sum_{v,n \in N_t} z_{n,v,t}^{event}, \beta_\pi + \sum_v \sum_{d=0}^{\delta} |N_{t+d}| f(t+d, t) \epsilon_d \right)$$