# Dynamic Topic Modeling for Monitoring Market Competition from Online Text and Image Data

## Hao Zhang[1], Gunhee Kim[2], Eric P. Xing[1]
### [1]: Carnegie Mellon University, [2]: Seoul National University

KDD2015

## Problem Statement

The increasing pervasiveness of Internet has lead to a wealth of consumer-created data over a multitude of online platforms

twitter · yelp · Google Reviews · 大众点评 dianping.com · BLOG:Sphere · Blogger

### What can we learn?
- 🙂 General public's opinion towards different companies' products and service
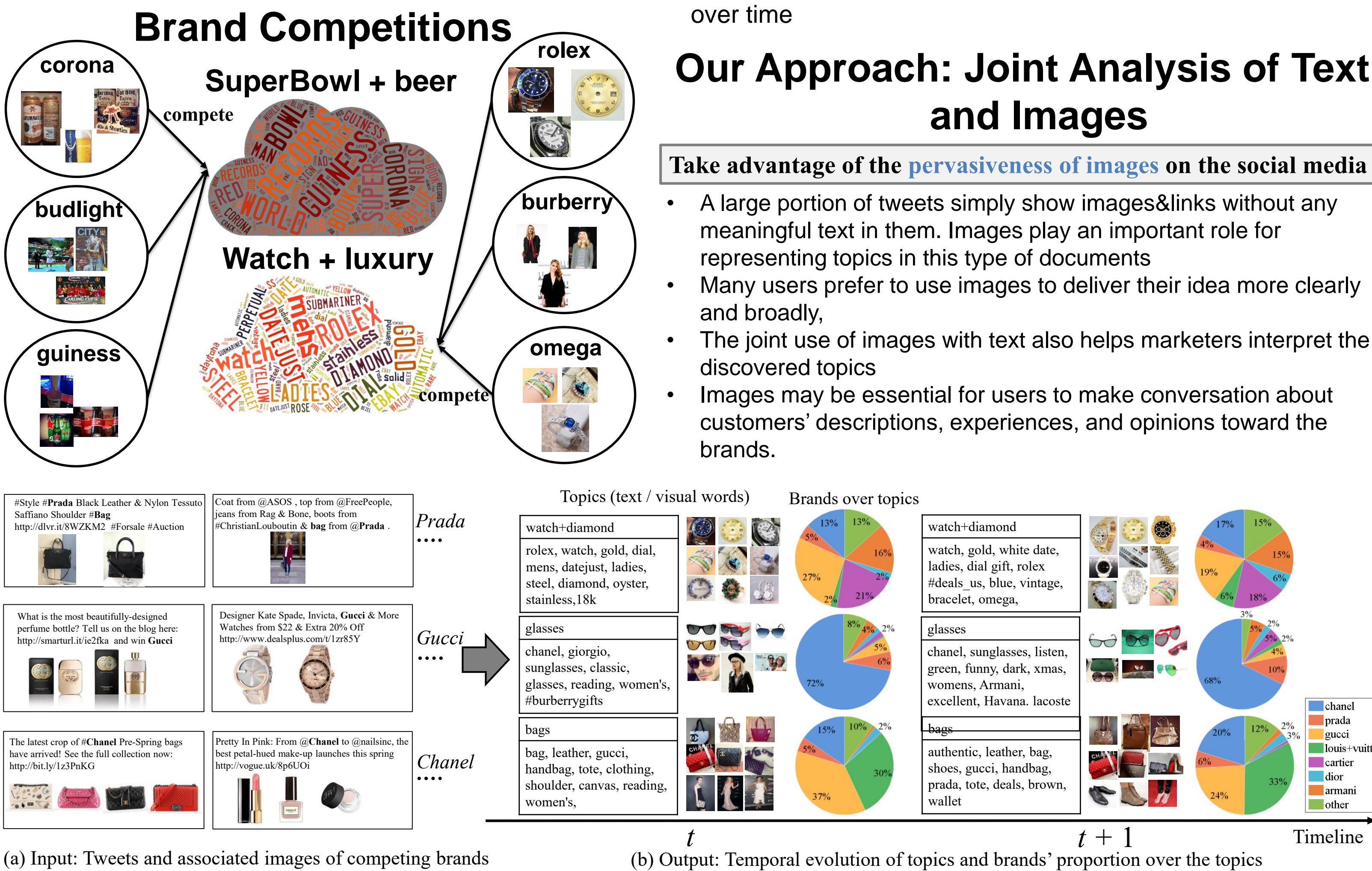- 🙂 Performance evaluations in different market conditions (time, location etc.)

### Brand Competitions

corona · rolex
SuperBowl + beer — compete
budlight · burberry
Watch + luxury
guiness · omega — compete

### What does marketers want to see?
- **Detection:** Listen in consumers' opinions towards their products and their competitors
- **Summarization:** Summarize/visualize how a shared market is occupied by different brands
- **Dynamics:** Monitoring the changes of market competition over time

### Our Approach: Joint Analysis of Text and Images

**Take advantage of the pervasiveness of images on the social media**
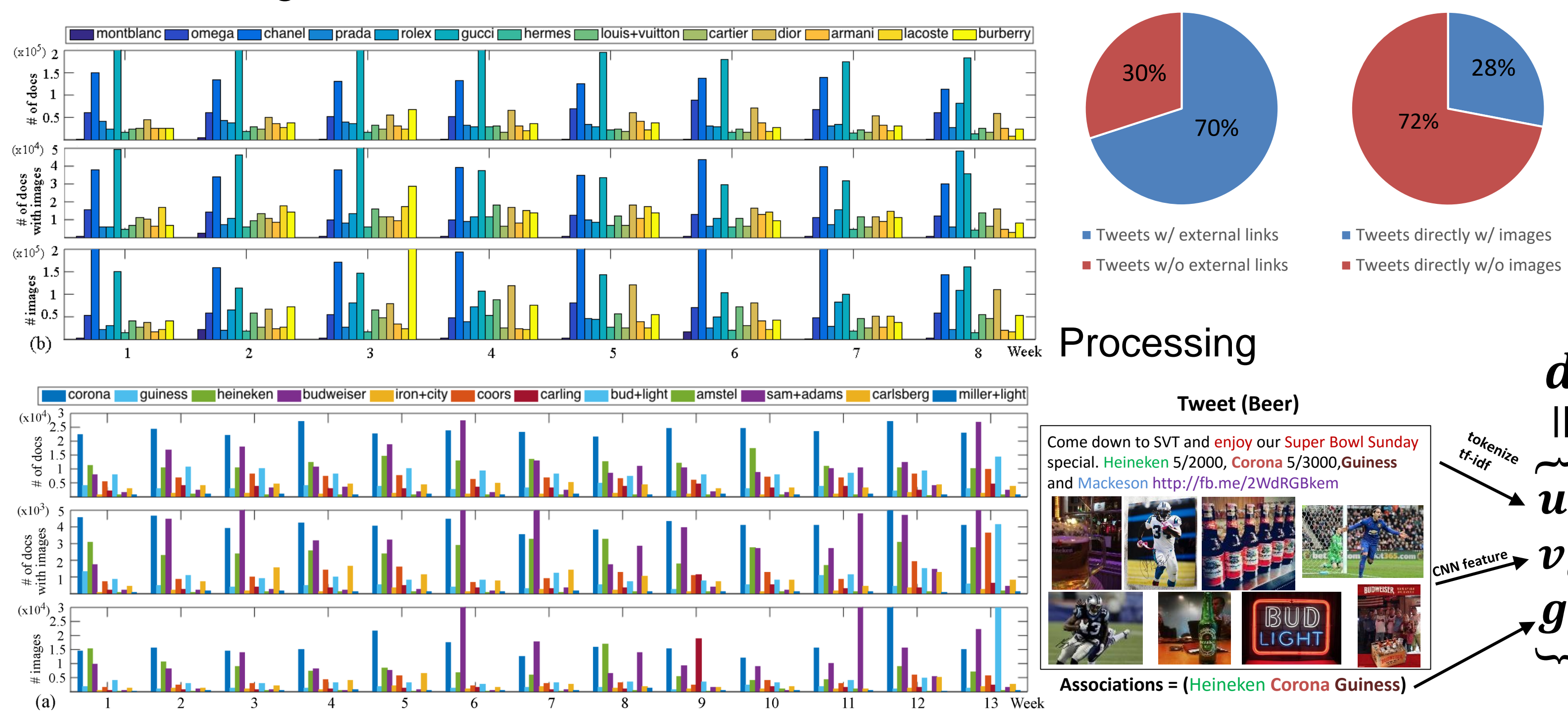- A large portion of tweets simply show images&links without any meaningful text in them. Images play an important role for representing topics in this type of documents
- Many users prefer to use images to deliver their idea more clearly and broadly,
- The joint use of images with text also helps marketers interpret the discovered topics
- Images may be essential for users to make conversation about customers' descriptions, experiences, and opinions toward the brands.



(a) Input: Tweets and associated images of competing brands
(b) Output: Temporal evolution of topics and brands' proportion over the topics

## Collecting Data

Crawling raw tweets and associated Images using the **REST** API
- 2 groups of brands: **Luxury** (13 brands) **Beer** (12 brands)
- **6.6M** tweets and **7.5M** images from **twitter** and external links
- Time range: **10/20/2014 to 02/01/2015**



### Processing

Tweet (Beer): Come down to SVT and enjoy our Super Bowl Sunday special. Heineken 5/2000, Corona 5/3000,Guiness and Mackeson http://fb.me/2WdRGBkem

$$d = \{u_d, v_d, g_d\}$$

CNN feature

Associations = {Heineken Corona Guiness}

## Competitive Dynamic Multi-view STC (cdSTC)

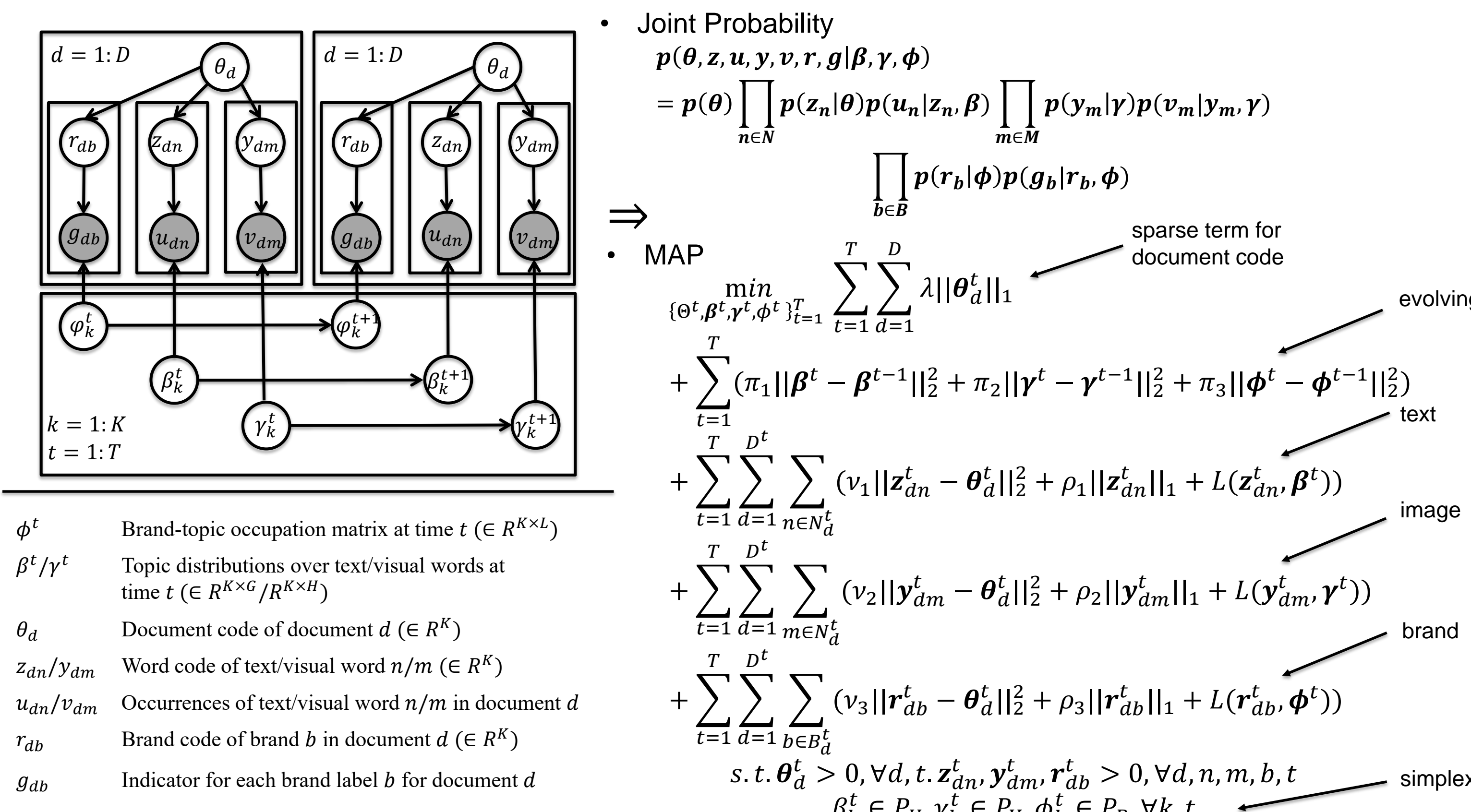The model aims to address 3 major challenges

1. **Multi-view** — • Modeling of multi-view representations of text and images

2. **Competition** — • Modeling of latent topics that are competitively shared by multiple brands

3. **Dynamic** — • Tracking temporal evolution of the topics and competitions



- Joint Probability
$$p(\theta, z, u, v, r, g | \beta, \gamma, \phi) = p(\theta) \prod_{n \in N} p(z_n|\theta) p(u_n|z_n,\beta) \prod_{m \in M} p(y_m|\gamma) p(v_m|y_m,\gamma) \prod_{b \in B} p(r_b|\phi) p(g_b|r_b,\phi)$$

- MAP
$$\min_{\{\theta^t,\beta^t,\gamma^t,\phi^t\}_{t=1}^T} \sum_{t=1}^T \sum_{d=1}^D \lambda \|\theta_d^t\|_1$$ — sparse term for document code

$$+ \sum_{t=1}^T (\pi_1 \|\beta^t - \beta^{t-1}\|_2^2 + \pi_2 \|\gamma^t - \gamma^{t-1}\|_2^2 + \pi_3 \|\phi^t - \phi^{t-1}\|_2^2)$$ — evolving chain

$$+ \sum_{t=1}^T \sum_{d=1}^{D^t} \sum_{n \in N_d^t} (\nu_1 \|z_{dn}^t - \theta_d^t\|_2^2 + \rho_1 \|z_{dn}^t\|_1 + L(z_{dn}^t, \beta^t))$$ — text

$$+ \sum_{t=1}^T \sum_{d=1}^{D^t} \sum_{m \in M_d^t} (\nu_2 \|y_{dm}^t - \theta_d^t\|_2^2 + \rho_2 \|y_{dm}^t\|_1 + L(y_{dm}^t, \gamma^t))$$ — image

$$+ \sum_{t=1}^T \sum_{d=1}^{D^t} \sum_{b \in B_d^t} (\nu_3 \|r_{db}^t - \theta_d^t\|_2^2 + \rho_3 \|r_{db}^t\|_1 + L(r_{db}^t, \phi^t))$$ — brand

$$s.t.\ \theta_d^t > 0, \forall d, t. z_{dn}^t, y_{dm}^t, r_{db}^t > 0, \forall d, n, m, b, t$$
$$\beta_k^t \in P_U, \gamma_k^t \in P_V, \phi_k^t \in P_B, \forall k, t$$ — simplex

| Symbol | Description |
|---|---|
| $\phi^t$ | Brand-topic occupation matrix at time $t$ ($\in R^{K \times L}$) |
| $\beta^t/\gamma^t$ | Topic distributions over text/visual words at time $t$ ($\in R^{K \times G}/R^{K \times H}$) |
| $\theta_d$ | Document code of document $d$ ($\in R^K$) |
| $z_{dn}/y_{dm}$ | Word code of text/visual word $n/m$ ($\in R^K$) |
| $u_{dn}/v_{dm}$ | Occurrences of text/visual word $n/m$ in document $d$ |
| $r_{db}$ | Brand code of brand $b$ in document $d$ ($\in R^K$) |
| $g_{db}$ | Indicator for each brand label $b$ for document $d$ |

## Evaluation: Topic Quality

amazon mechanical turk

Argument 1: Lower perplexity ≠ higher quality [J. Chang 2009]
Argument 2: Perplexity is not a fair metric for models with different distributions

– Define the **Coherence Measure (CM)** and the **Validity Measure (VM)**:

$$CM = \frac{\#\ of\ relevant\ words}{\#\ of\ words\ in\ valid\ topics} \qquad VM = \frac{\#\ of\ valid\ topics}{\#\ of\ topics}$$

- **Average VM/CM on text topics**

|  | VM (Beer / Luxury) | CM (Beer / Luxury) |
|---|---|---|
| dLDA | 0.53 / 0.68 | 0.55 / 0.52 |
| STC + dyn | 0.44 / 0.66 | 0.57 / 0.57 |
| cdSTC + multi | 0.51 / 0.70 | **0.63 / 0.59** |
| cdSTC + text | **0.605 / 0.71** | 0.61 / **0.59** |

- **Average VM/CM on visual topic**

|  | VM (Beer / Luxury) | CM (Beer / Luxury) |
|---|---|---|
| Kmeans | 0.39 / 0.56 | 0.59 / 0.64 |
| LDA + multi | **0.57** / 0.63 | 0.51 / 0.69 |
| cdSTC + multi | **0.57 / 0.65** | **0.66 / 0.71** |

## Evaluation: Prediction

- **Task I**: Given a novel tweet, can we predict its most associated brand?

novel tweets → Model → infer → Gucci ⟺
$$\max_{\{\theta^t, \mathcal{M}^t, \eta^t\}} \sum_{t=1}^T f(\theta^t, \mathcal{M}^t, D^t) + CR(\theta^t, \eta^t) + \frac{1}{2}\|\eta^t\|_2^2$$
$$s.t.\ \theta_d^t > 0, \forall d, t. z_{dn}^t, y_{dm}^t > 0, \forall d, n, m, t$$
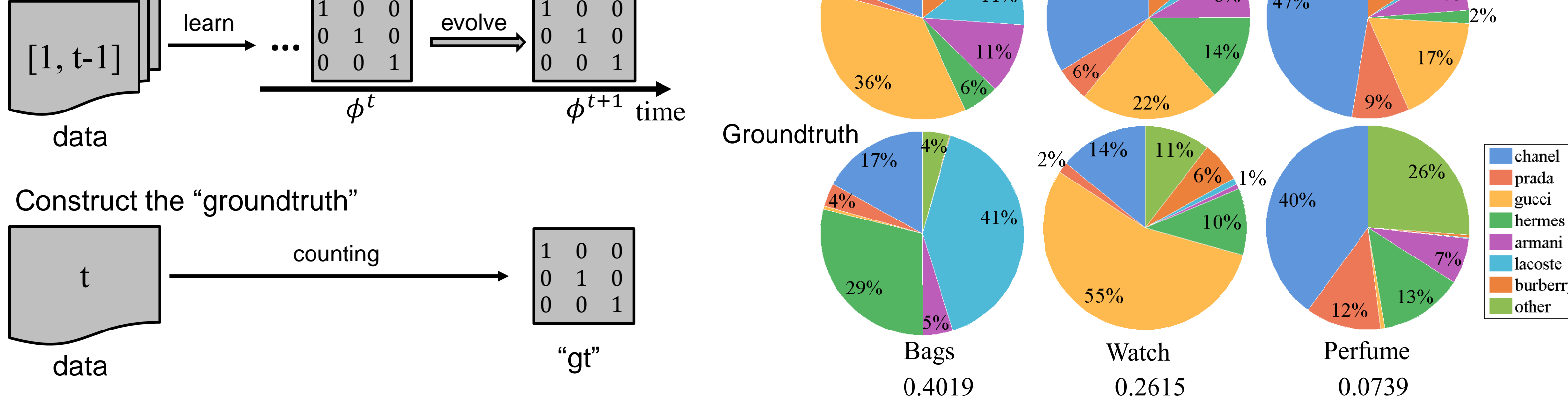$$\beta_k^t \in P_U, \gamma_k^t \in P_V, \forall k, t$$

**Task I-I:** Randomly split data in every time slice into 90% for training and 10% for testing

**Task I-II:** Use the data in $[1, t-1]$ for training, $[t-1, t]$ for testing



- **Task II**: Given an unseen past document, can we predict its timestamp?

past tweets → locate → Sent at this time point

$$\max_t p(d|\mathcal{M}^t), where$$
$$p(d|\mathcal{M}^t) = \prod_{n \in N_d} p(u_n|\beta^t) \prod_{m \in M_d} p(v_m|\gamma^t) \prod_{b \in B_d} p(g_b|\phi^t)$$



- **Task III**: Can we predict future competition trends using past data?

Evolve the competition matrix

$[1, t-1]$ data → learn → ... → evolve → $\phi^{t+1}$ time

Construct the "groundtruth"

$t$ data → counting → "gt"

Prediction / Groundtruth

Bags 0.4019 · Watch 0.2615 · Perfume 0.0739

## Brand Competition Monitoring

### Objective
- How brands occupy the market in every time slice?
- How each textual/visual topic evolves over time?
- How each brand's occupation changes over time?
- How's the competition trends between multi-brands like over time?

easy → difficult