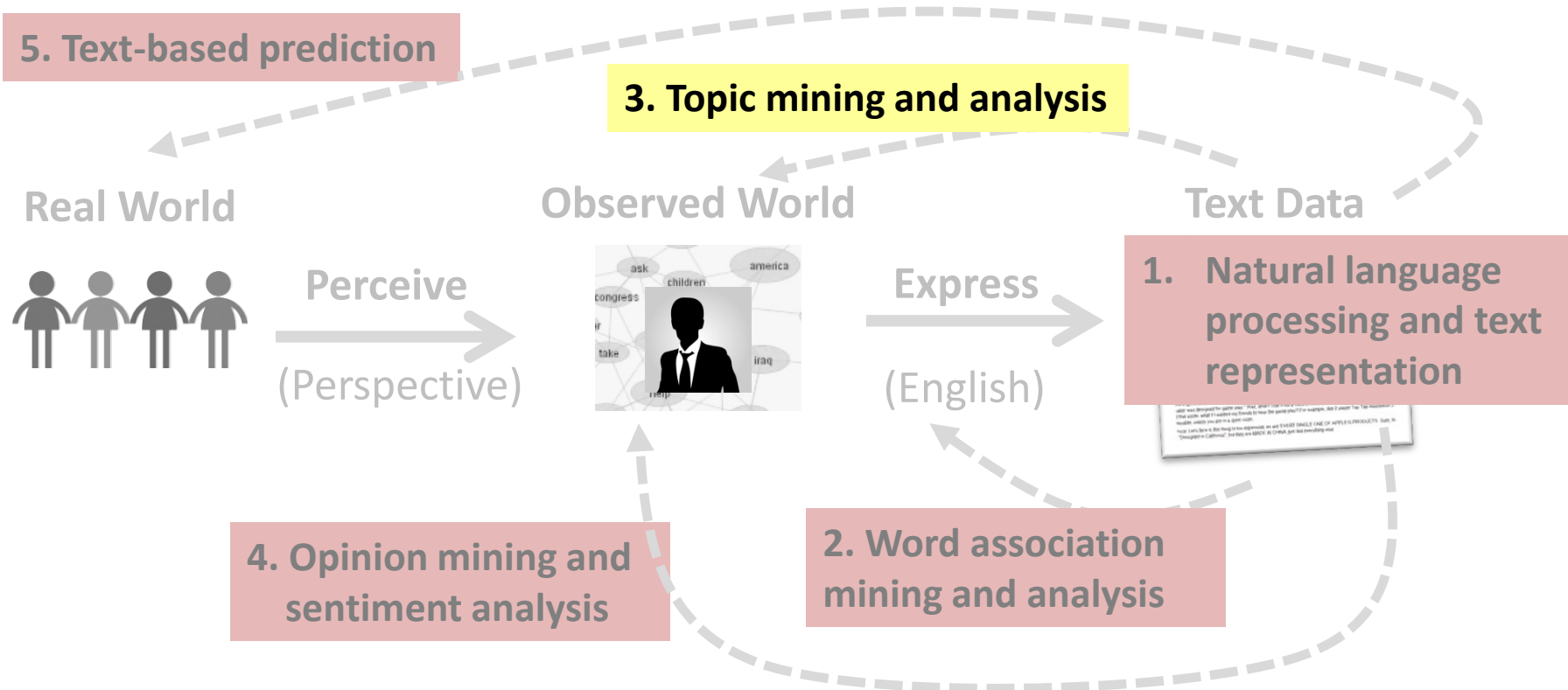




Probabilistic Topic Models: Mixture of Unigram Language Models

ChengXiang “Cheng” Zhai
Department of Computer Science
University of Illinois at Urbana-Champaign

Probabilistic Topic Models: Mixture of Unigram LMs



Factoring out Background Words

d

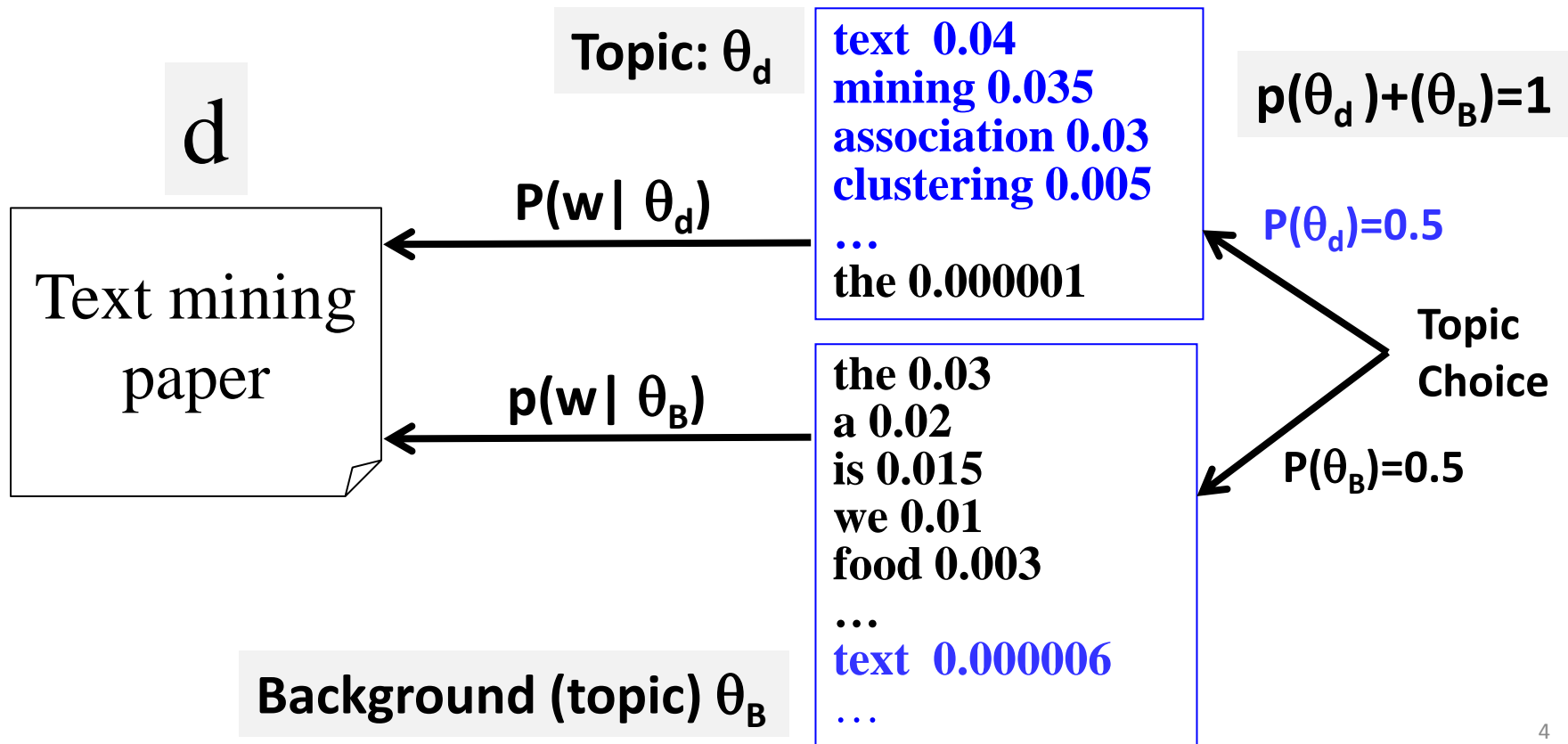
Text mining
paper

$p(w | \theta)$

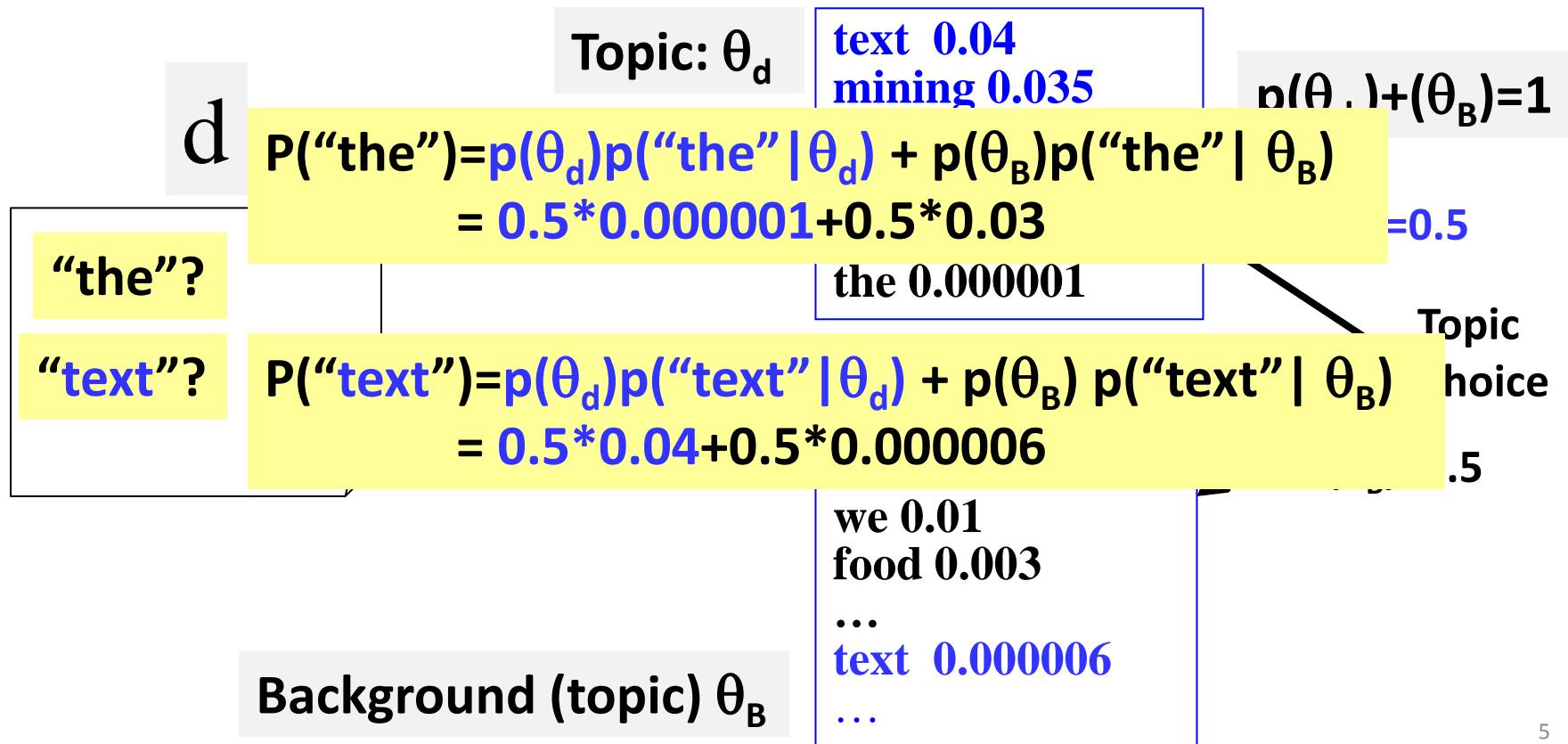
the 0.031
a 0.018
...
text 0.04
mining 0.035
association 0.03
clustering 0.005
computer 0.0009
...
food 0.000001
...

How can we get rid of
these common words?

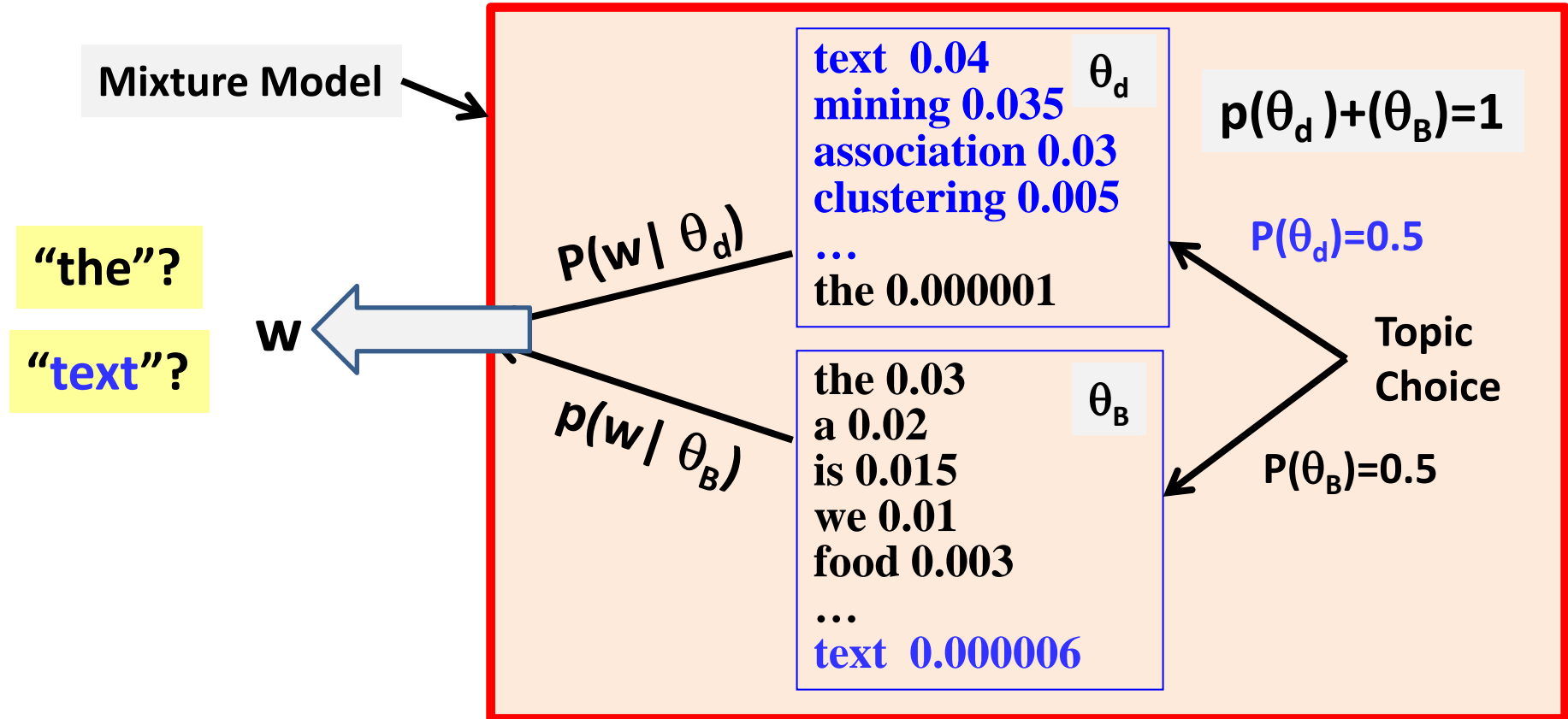
Generate d Using Two Word Distributions



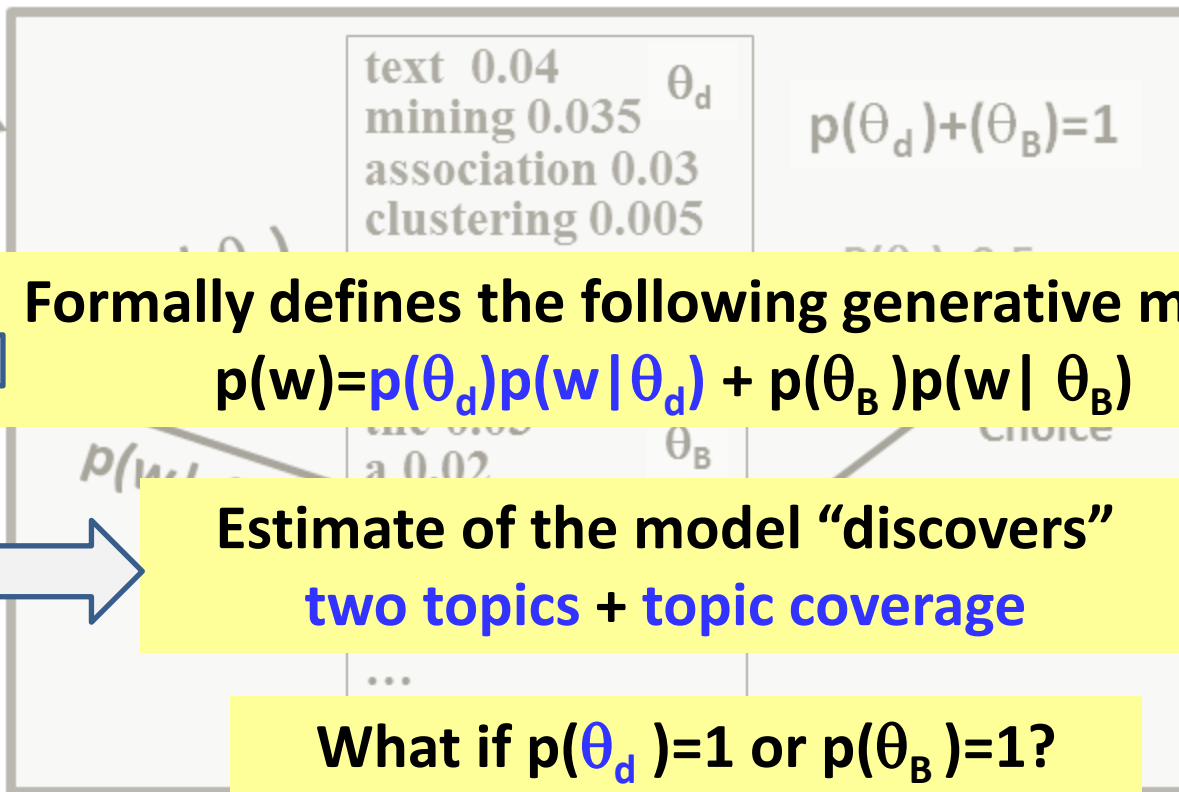
What's the probability of observing a word w ?



The Idea of a Mixture Model



As a Generative Model...



Mixture of Two Unigram Language Models

- **Data:** Document d
- **Mixture Model: parameters** $\Lambda = (\{p(w | \theta_d)\}, \{p(w | \theta_B)\}, p(\theta_B), p(\theta_d))$
 - Two unigram LMs: θ_d (the topic of d); θ_B (background topic)
 - Mixing weight (topic choice): $p(\theta_d) + p(\theta_B) = 1$

- **Likelihood function:**

$$\begin{aligned} p(d | \Lambda) &= \prod_{i=1}^{|d|} p(x_i | \Lambda) = \prod_{i=1}^{|d|} [p(\theta_d)p(x_i | \theta_d) + p(\theta_B)p(x_i | \theta_B)] \\ &= \prod_{i=1}^M [p(\theta_d)p(w_i | \theta_d) + p(\theta_B)p(w_i | \theta_B)]^{c(w,d)} \end{aligned}$$

- **ML Estimate:** $\Lambda^* = \arg \max_{\Lambda} p(d | \Lambda)$

$$\text{Subject to} \quad \sum_{i=1}^M p(w_i | \theta_d) = \sum_{i=1}^M p(w_i | \theta_B) = 1 \quad p(\theta_d) + p(\theta_B) = 1$$