

Improvements to The Bayesian Topic N-gram Models

Hiroshi Noji^{1,3}

Daichi Mochihashi^{2,3}

Yusuke Miyao^{1,3}

1. National Institute of Informatics, Tokyo
2. The Institute of Statistical Mathematics, Tokyo
3. Graduate University for Advanced Studies

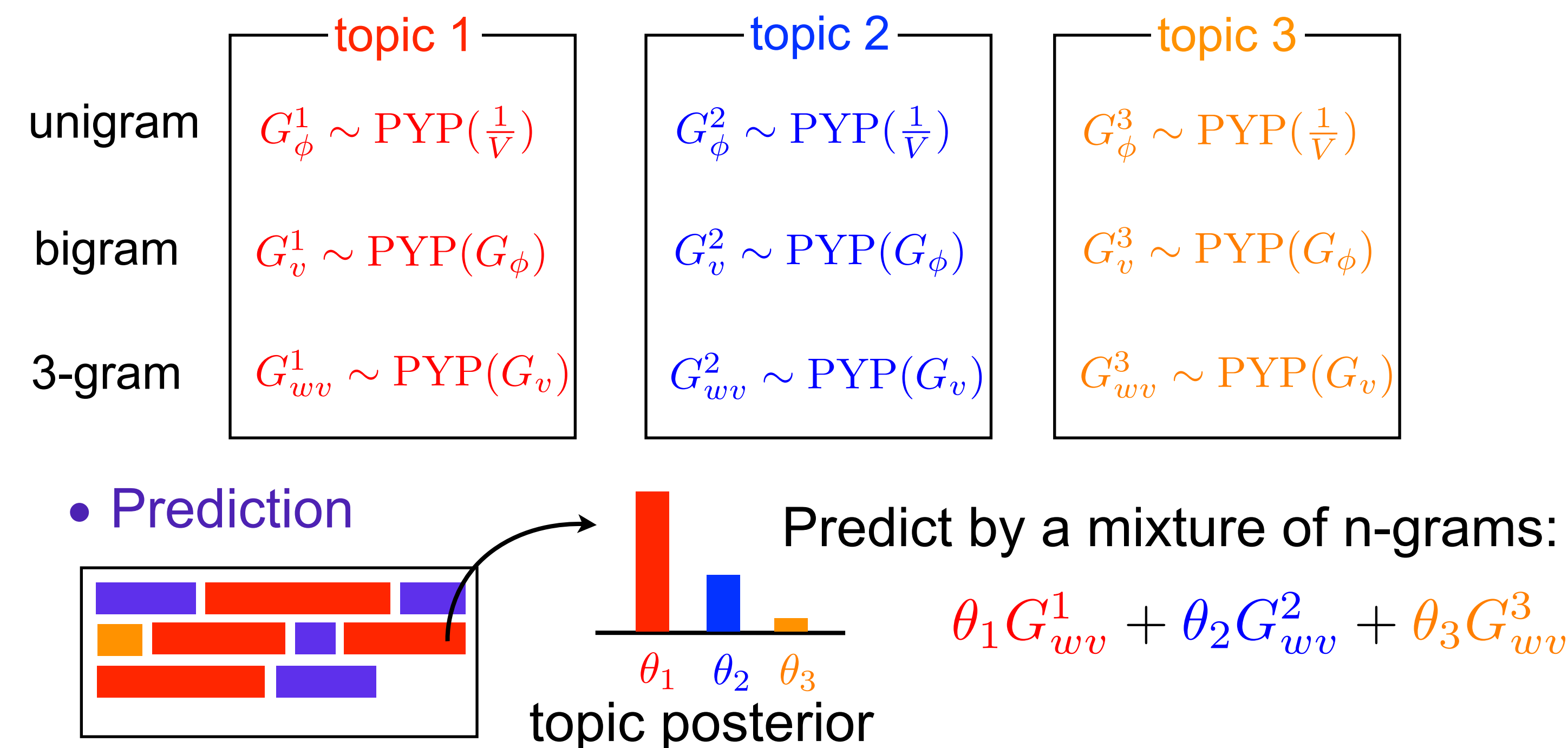
Overview

- The predictions of n-gram language models are very *local*



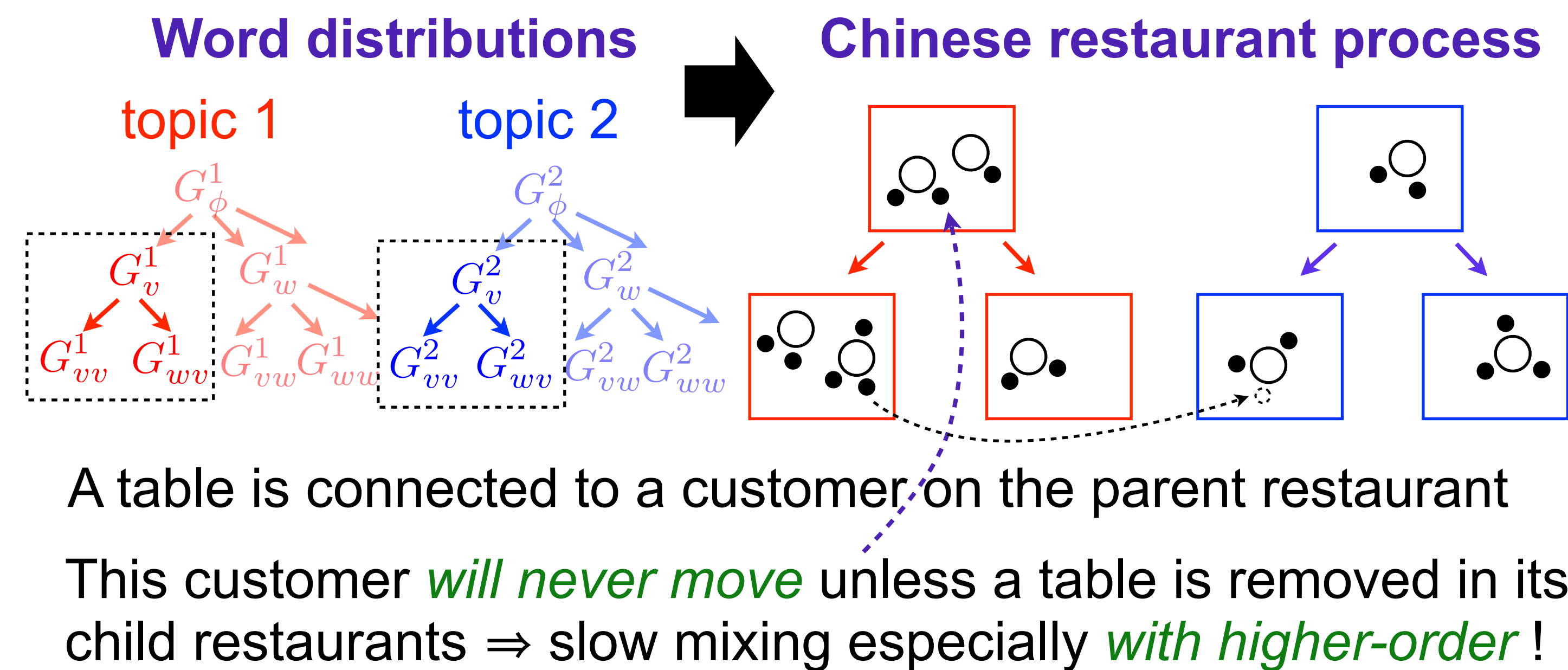
- Problems of the previous n-gram + topic models:
 - Prediction get much sparse *with higher order n-grams*
 - Inference (local Gibbs) is not very efficient
 - Resolve these with *hierarchical prior* + *blocked sampling*

Basic model (Wallach' 06)

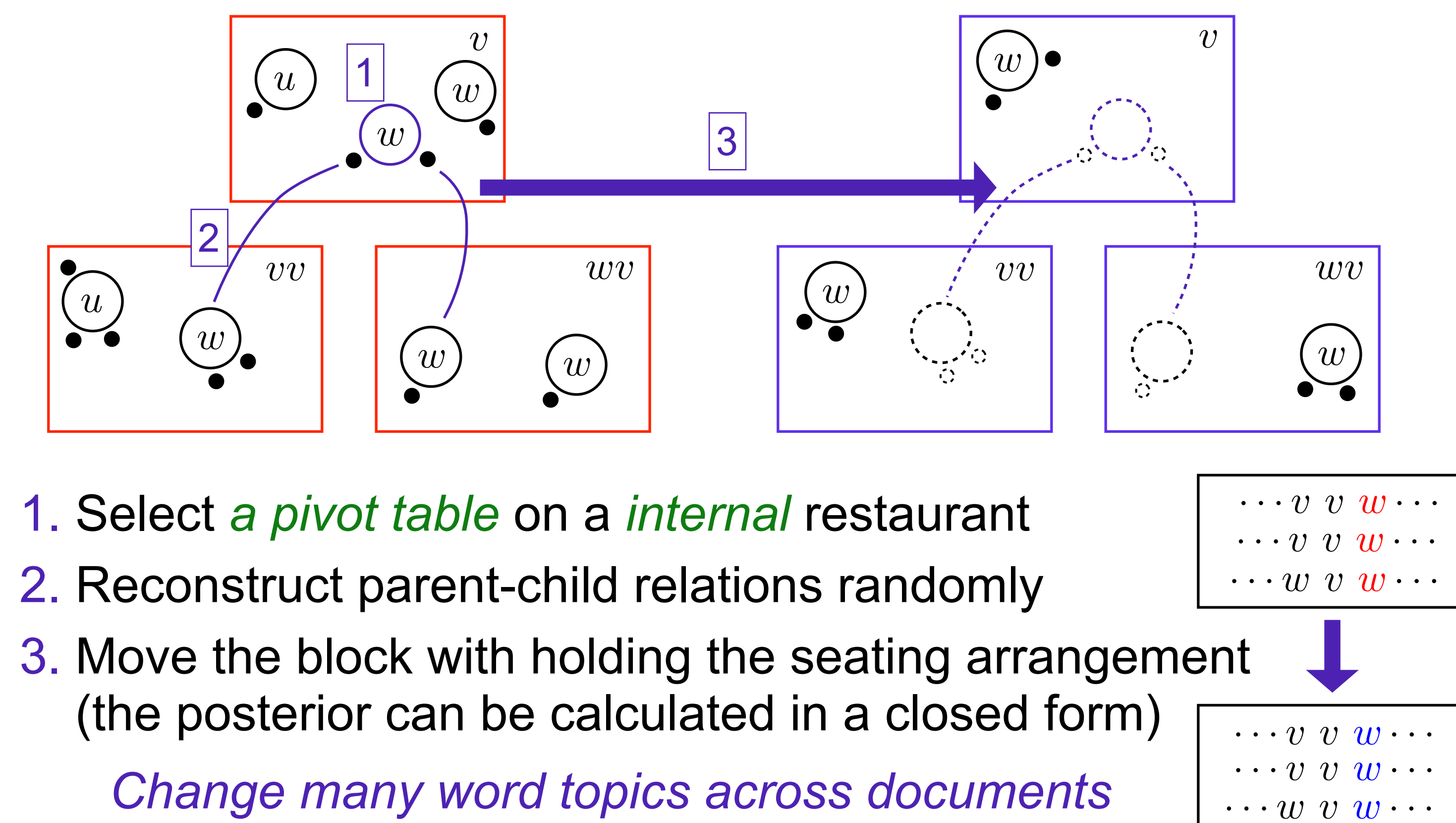


- Prediction

Problem of the naive Gibbs sampler

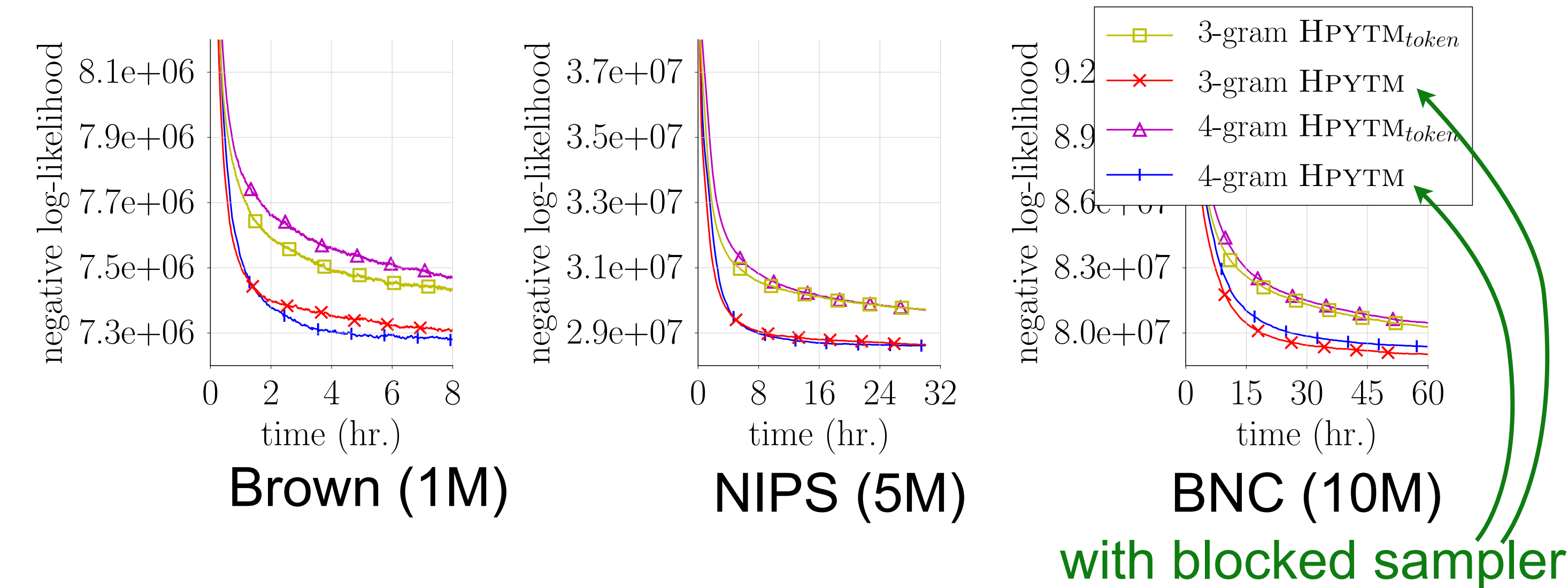


Proposed blocked sampler



Effects of the blocked sampler

- Use Wallach's model for simplicity
- It can be applied to our extended models with a little effort



Perplexity results

- BNC (10M) with 4-gram models
- # topics=100 (not so sensitive)
- Our models are *normalization-free*:
 - much faster prediction than rescaling methods
- Conclusion to the model design:
 - Flat structure* between the global and topics is better than *hierarchical structure* => *but why?*

HPYLM (no topic)	169.2
Wallach	140.4
Wallach + block	133.1
Unigram rescaling	130.3
Hierarchical	129.0
Switching	125.5

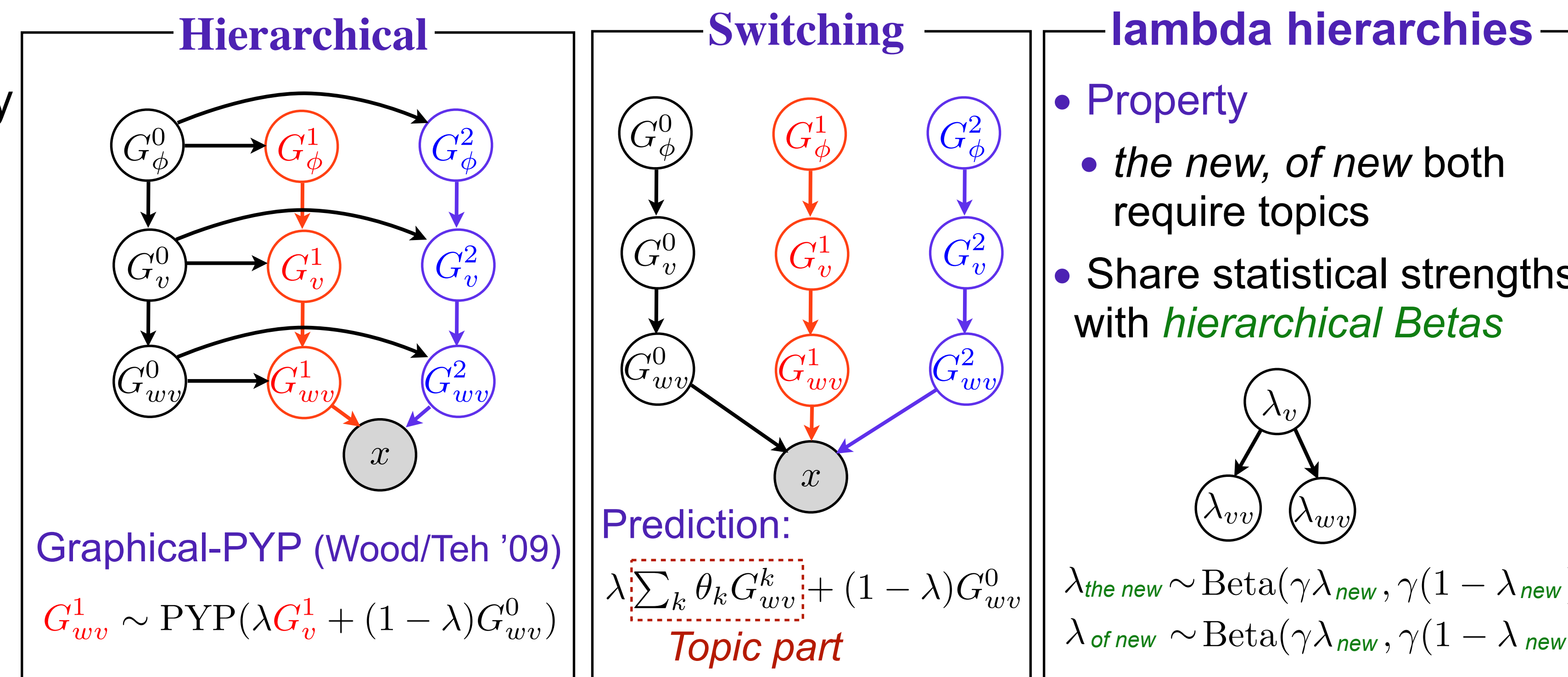
Posterior inspections

- Switching** assigns only some part of words topics
there has been much recent work on measuring image statistics and on learning probability distributions on images. we observe that the mapping from images to statistics is many-to-one and show it can be quantified by a phase space factor.
- lighter words are assigned the global model (topic 0)
- might led to more *accurate topic prediction*
- Hierarchical** also learns differences of contexts, but all words are assigned topics

λ_h	h
0.0–0.1	in spite, were unable, a sort, on behalf, . regardless
0.5–0.6	assumed it, rand mines, plans was, other excersises
0.9–1.0	that the, the existing, the new, their own, and spatial

Hierarchical priors for ease of sparseness

- Motivation
 - We don't want to assign all n-grams topics equally
- local context candidates**
 - ... in order => to, that } *not require* topics
 - ... would like => you, to }
 - ... state of => **Washington, the** } *require* topics
 - ... the new => **york, algorithm** }
- Comparing two models encoding this difference
 - Hierarchical**: the global model is used as a prior
 - Switching**: two models are switched



References

Hanna M. Wallach. 2006. Topic modeling: beyond bag- of-words. In *Proc of ICML '06*, pages 977–984.
Frank Wood and Yee Whye Teh. 2009. A hierarchical nonparametric Bayesian approach to statistical language model domain adaptation. In *Proc of AISTATS*, volume 12.