

## Motivation

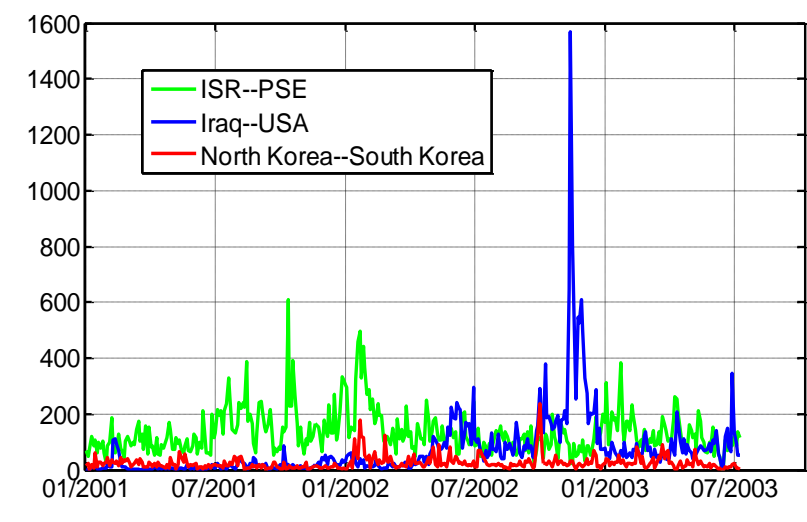
### Real world applications:

International Relation Study; Recommender System...

### Time varying counts

➤ Capture complex dependencies across time steps

➤ Infer interpretable latent structure to analyze and predict



## Related Work

### Poisson Gamma Dynamical Systems (PGDS) [2]

$$y_v^{(t)} \sim \text{Pois}(\delta^{(t)} \sum_{k=1}^K \phi_{vk} \theta_k^{(t)}) \text{ and } \theta_k^{(t)} \sim \text{Gam}(\tau_0 \sum_{k_2=1}^K \pi_{kk_2} \theta_{k_2}^{(t-1)}, \tau_0)$$

✓ Gamma-Poisson construction, and supports expressive latent transition structures

✓ Shallow model may still have shortcomings in capturing long-range temporal dependencies

### Deep Temporal Sigmoid Belief Networks [3]

$$p_{\theta}(V, H) = p(h_1)p(v_1|h_1) \cdot \prod_{t=2}^T p(h_t|h_{t-1}, v_{t-1}) \cdot p(v_t|h_t, v_{t-1})$$

$$p(h_{jt} = 1|h_{t-1}, v_{t-1}) = \sigma(w_{1j}^T h_{t-1} + w_{3j}^T v_{t-1} + b_j), p(v_t|h_t, v_{t-1}) = \prod_{m=1}^M y_{mt}^{v_{mt}}$$

✓ A sequential stack of sigmoid belief networks (SBNs)

✓ How the layers in DTSBN are related to each other lacks an intuitive interpretation

✓ Modelling count data with Replicated Softmax Model

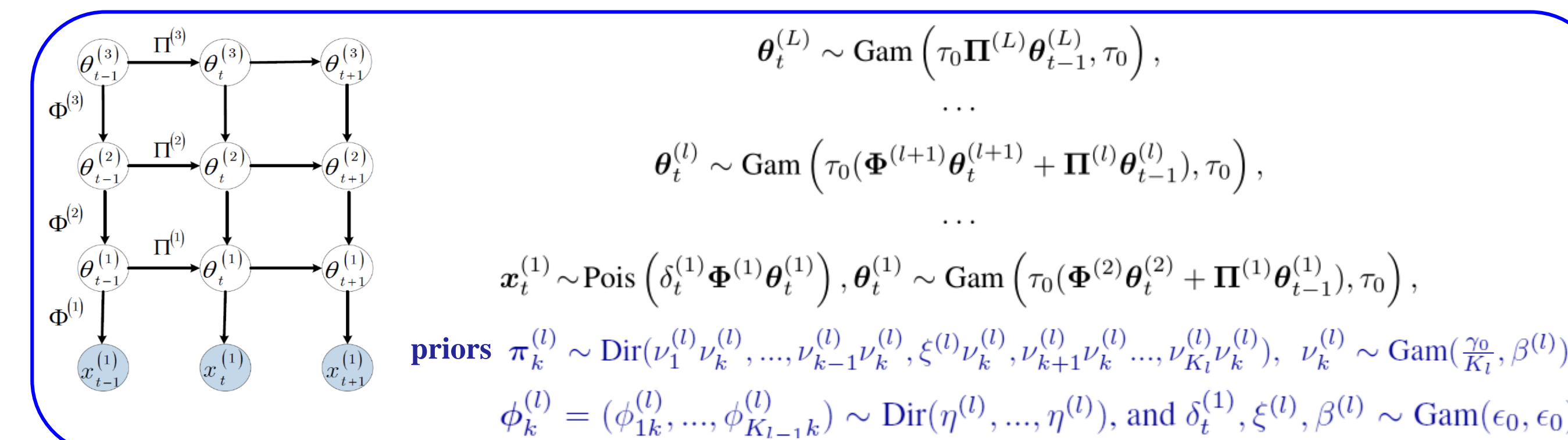
## Contributions

➤ Improve previously proposed models by mining deep hierarchical latent structure from the data, and capturing long-range temporal dependencies.

➤ Derive closed-form Gibbs sampling update equations.

➤ Develop stochastic gradient MCMC inference that is scalable to big dataset.

### Deep Poisson-Gamma Dynamical Systems



➤ Hierarchical structure

$$\mathbb{E}[x_t^{(i)} | \theta_t^{(i)}, \{\Phi^{(p)}\}_{p=1}^L] = [\prod_{p=1}^L \Phi^{(p)}] \theta_t^{(i)}$$

➤ Long-range temporal dependence

$$\mathbb{E}[x_t^{(1)} | \theta_{t-1}^{(1)}, \theta_{t-2}^{(1)}, \theta_{t-3}^{(1)} / \delta_t^{(1)}] = \Phi^{(1)} \Pi^{(1)} \theta_{t-1}^{(1)} + \Phi^{(1)} \Phi^{(2)} [\Pi^{(2)}]^2 \theta_{t-2}^{(1)} + \Phi^{(1)} \Phi^{(2)} (\Pi^{(2)} \Phi^{(3)} + \Phi^{(3)} \Pi^{(3)}) [\Pi^{(3)}]^2 \theta_{t-3}^{(1)}$$

## Inference

### Backward and upward propagation of latent counts

Techniques:

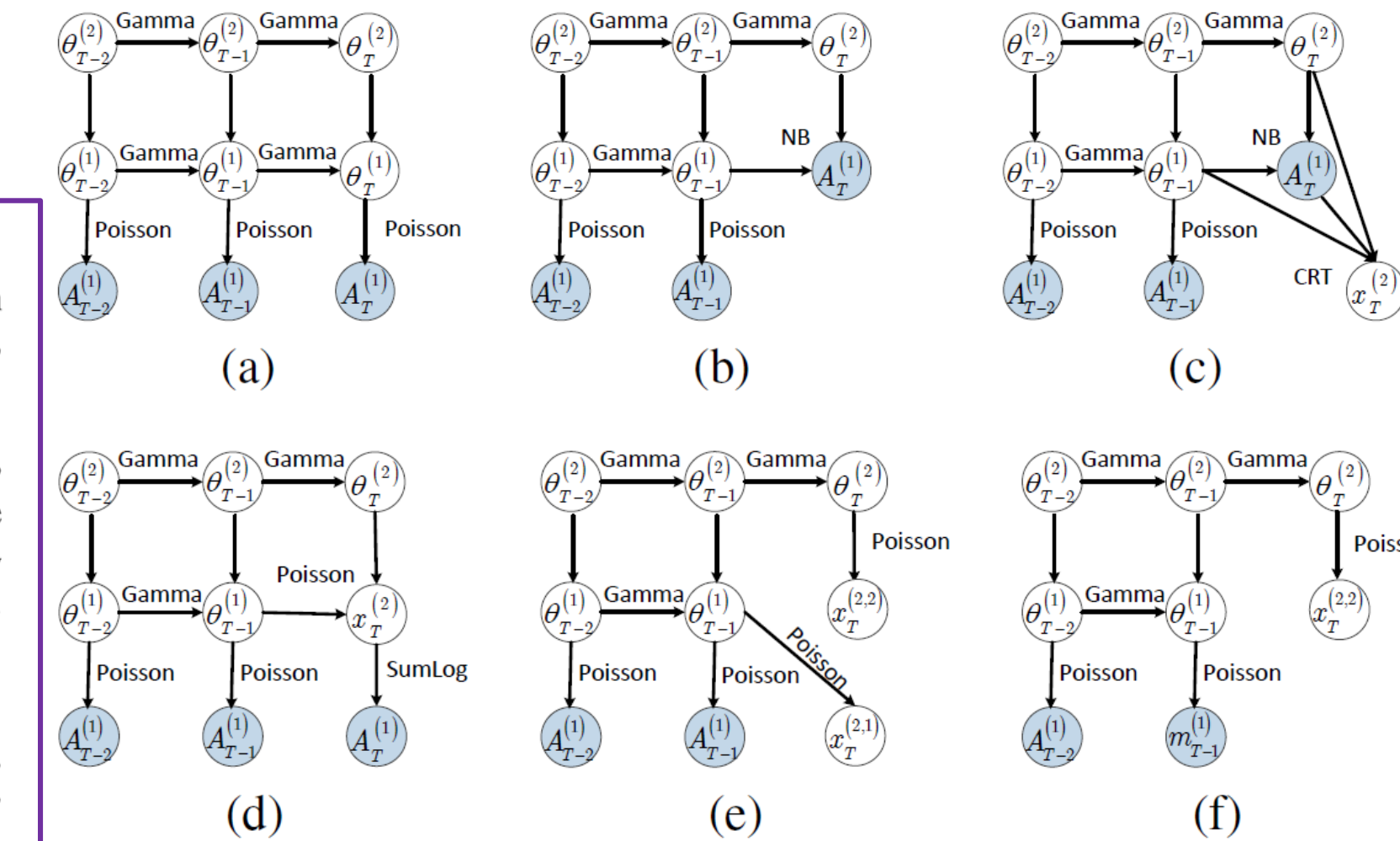
- Data augmentation
- Marginalization

Repeat the process

**Property 1 (P1):** if  $y_{-kt} = \sum_{n=1}^N y_n$ , where  $y_n \sim \text{Pois}(\theta_n)$  are independent Poisson-distributed random variables, then  $(y_1, \dots, y_n) \sim \text{Multi}\left(y, \frac{\theta_1}{\sum_{n=1}^N \theta_n}, \dots, \frac{\theta_N}{\sum_{n=1}^N \theta_n}\right)$  and  $y \sim \text{Pois}(\sum_{n=1}^N \theta_n)$ .

**Property 2 (P2):**  $y \sim \text{Pois}(c\theta)$ , where  $c$  is a constant, and  $\theta \sim \text{Gam}(a, b)$  then  $y \sim \text{NB}\left(a, \frac{c}{c+b}\right)$  is a negative binomial distributed random variable. We can equivalently parameterize it as  $y \sim \text{NB}(a, g(\zeta))$ , where  $g(\zeta) = 1 - \exp(-\zeta)$  is the Bernoulli-Poisson link and  $\zeta = \ln(1 + \frac{b}{a})$ .

**Property 3 (P3):** if  $y \sim \text{NB}(a, g(\zeta))$  and  $l \sim \text{CRT}(y, a)$  is a Chinese restaurant table-distributed random variable, then  $y$  and  $l$  are equivalently jointly distributed as  $y \sim \text{SumLog}(l, g(\zeta))$  and  $l \sim \text{Pois}(a\zeta)$ . Refer to [1,2]



➤ Working backward for  $t = T, \dots, 2$  and upward for  $l = 1, 2, \dots, L$ , we draw

$$(A_{k1t}^{(l)}, \dots, A_{kK_t}^{(l)}) \sim \text{Multi}\left(x_{kt}^{(l,l)}; \frac{\phi_{k1}^{(l)} \theta_{k1}^{(l)}}{\sum_{k_1=1}^{K_1} \phi_{kk_1}^{(l)} \theta_{k_1}^{(l)}}, \dots, \frac{\phi_{kK_t}^{(l)} \theta_{kK_t}^{(l)}}{\sum_{k_1=1}^{K_1} \phi_{kk_1}^{(l)} \theta_{k_1}^{(l)}}\right),$$

$$x_{kt}^{(l+1)} \sim \text{CRT}\left[A_{k,t}^{(l)} + Z_{k,t+1}^{(l)}, \tau_0 \left(\sum_{k_{l+1}=1}^{K_{l+1}} \phi_{kk_{l+1}}^{(l+1)} \theta_{k_{l+1}}^{(l+1)} + \sum_{k_1=1}^{K_l} \pi_{kk_1}^{(l)} \theta_{k_1}^{(l)}\right)\right]$$

$$(x_{kt}^{(l+1,l)}, x_{kt}^{(l+1,l+1)}) \sim \text{Multi}\left(x_{kt}^{(l+1)}, p_1/(p_1 + p_2), p_2/(p_1 + p_2)\right)$$

$$(Z_{k1t}^{(l)}, \dots, Z_{kK_t}^{(l)}) \sim \text{Multi}\left(x_{kt}^{(l+1,l)}; \frac{\pi_{k1}^{(l)} \theta_{k1}^{(l)}}{\sum_{k_1=1}^{K_l} \pi_{kk_1}^{(l)} \theta_{k_1}^{(l)}}, \dots, \frac{\pi_{kK_t}^{(l)} \theta_{kK_t}^{(l)}}{\sum_{k_1=1}^{K_l} \pi_{kk_1}^{(l)} \theta_{k_1}^{(l)}}\right)$$

➤ Working forward for  $t = 1, \dots, T$  and downward for  $l = L, \dots, 1$ , we sample

$$\theta_{kt}^{(l)} \sim \text{Gamma}\left[A_{k,t}^{(l)} + Z_{k,t+1}^{(l)} + \tau_0 \left(\sum_{k_{l+1}=1}^{K_{l+1}} \phi_{kk_{l+1}}^{(l+1)} \theta_{k_{l+1}}^{(l+1)} + \sum_{k_1=1}^{K_l} \pi_{kk_1}^{(l)} \theta_{k_1}^{(l)}\right), \tau_0 (1 + \zeta_t^{(l-1)} + \zeta_{t+1}^{(l)})\right]$$

### Stochastic gradient MCMC inference for simplex

$$\begin{aligned} (\pi_k^{(l)})_{n+1} = & \left[ (\pi_k^{(l)})_n + \frac{\varepsilon_n}{M_k^{(l)}} \left[ (\rho_{\tilde{z}_{k,n}}^{(l)} + \eta_{k,n}^{(l)}) - (\rho_{\tilde{z}_{k,n}}^{(l)} + \eta_{k,n}^{(l)}) (\pi_k^{(l)})_n \right] \right. \\ & \left. + \mathcal{N}\left(0, \frac{2\varepsilon_n}{M_k^{(l)}} \left[ \text{diag}(\pi_k^{(l)})_n - (\pi_k^{(l)})_n (\pi_k^{(l)})_n^T \right] \right) \right]_{\angle}, \text{ Refer to [4]} \end{aligned}$$

**Algorithm 2** Stochastic-gradient MCMC for DPGDS

Input: Data mini-batches; Output: Global parameters of DPGDS.

**for**  $i = 1, 2, \dots$  **do**  
  \ \* Collect local information  
  Backward-upward Gibbs sampling on the  $i$ th mini-batch for  $\{A_{vkt}^{(l)}\}_{v,k,t}; \{x_{kt}^{(l+1)}\}_{k,t}; \{x_{kt}^{(l+1,l)}\}_{k,t}; \{x_{kt}^{(l+1,l+1)}\}_{k,t}; \{Z_{k_1k_2,t}^{(l)}\}_{k_1,k_2,t}$  with (9)(10)(11)(12);  
  Backward-upward calculating for the  $\{\zeta_t^{(l)}\}_t$ ;  
  Forward-downward Gibbs sampling for the  $\{\theta_t^{(l)}\}_t$  with (13);  
  Sampling  $\delta^{(l)}$  with (17) or (18);  
  \ \* Update global parameters  
  **for**  $l = 1, 2, \dots, L$  **and**  $k = 1, 2, \dots, K_L$  **do**  
    Update  $M_k^{(l)}$  according to [28]; then  $\{\phi_k^{(l)}\}_k$  with (35); Update  $M_k^{(l)}$  according to [28]; then  $\{\pi_k^{(l)}\}_k$  with (34);  
  **end for**  
  Update  $\xi^{(l)}, \{\nu_k^{(l)}\}_k$ , and  $\beta^{(l)}$  with SGNHT [21]

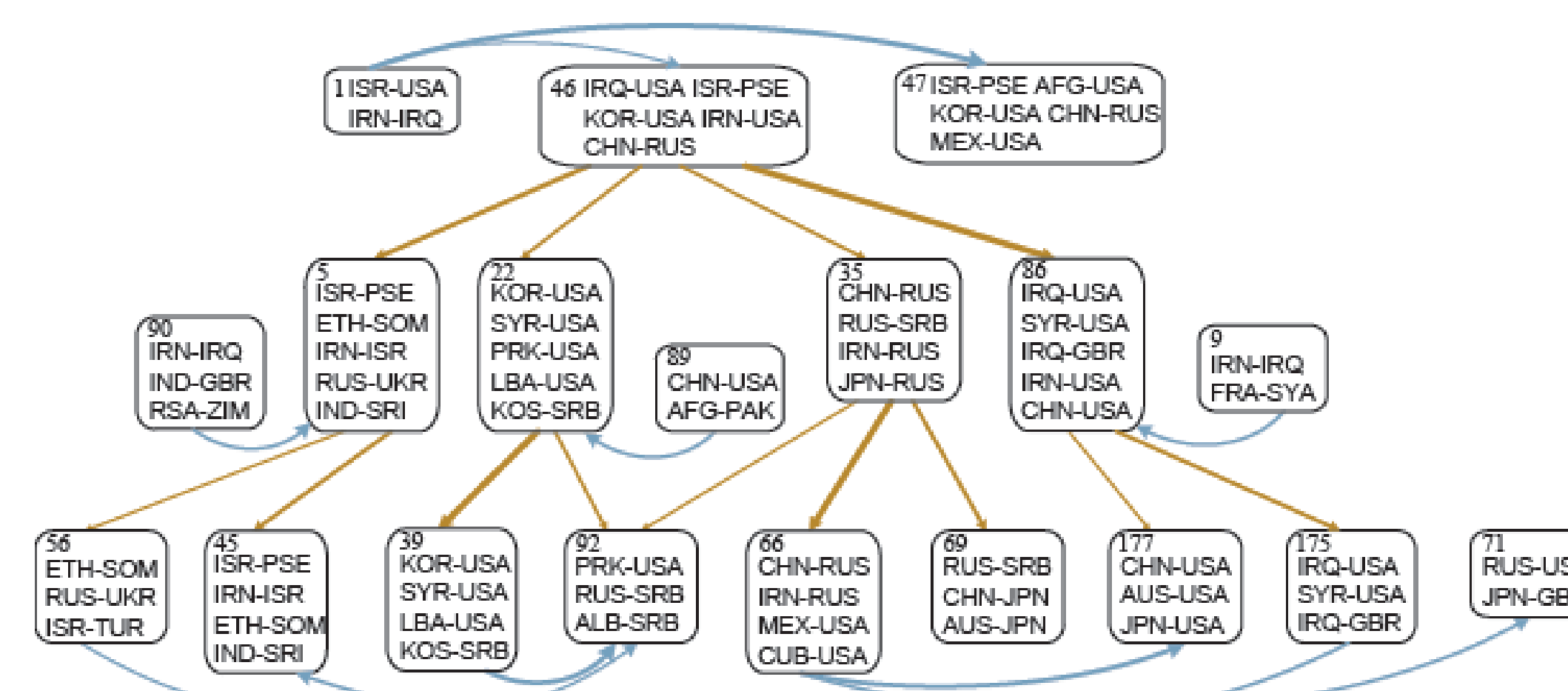
**end for**

## Experiments

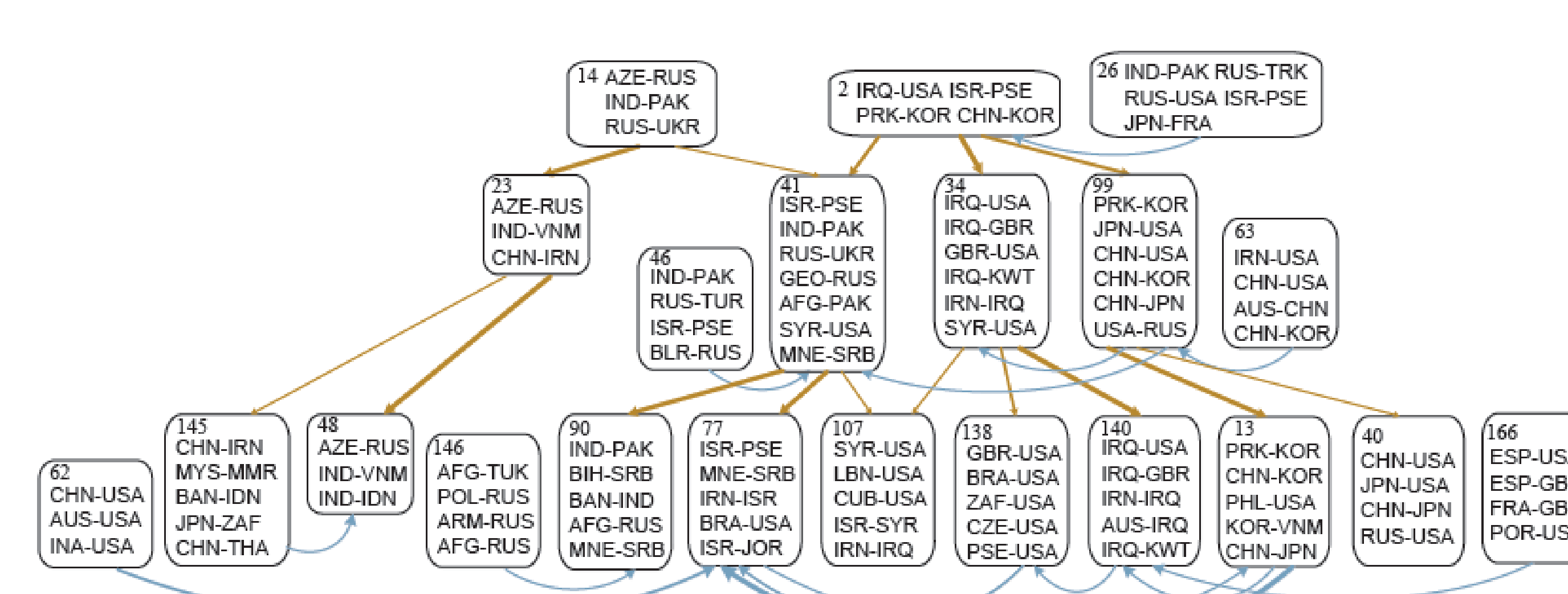
### Top@M results on real-world text data

Model	Top@M	Data				
		GDEL	ICEWS	SOTU	DBLP	NIPS
		$T = 365, V \approx 9000$	$T = 365, V \approx 3000$	$T = 225, V = 7518$	$T = 14, V = 1771$	$T = 17, V = 9836$
GDPFA	MP	0.611 $\pm$ 0.001	0.607 $\pm$ 0.002	0.379 $\pm$ 0.002	<b>0.435</b> $\pm$ 0.009	0.843 $\pm$ 0.005
	MR	0.145 $\pm$ 0.002	0.235 $\pm$ 0.005	0.369 $\pm$ 0.002	0.254 $\pm$ 0.005	0.050 $\pm$ 0.001
	PP	0.447 $\pm$ 0.014	0.465 $\pm$ 0.008	0.617 $\pm$ 0.013	0.581 $\pm$ 0.011	0.807 $\pm$ 0.006
PGDS	MP	0.679 $\pm$ 0.001	0.658 $\pm$ 0.001	0.375 $\pm$ 0.002	0.419 $\pm$ 0.004	0.864 $\pm$ 0.004
	MR	<b>0.150</b> $\pm$ 0.001	<b>0.245</b> $\pm$ 0.005	0.373 $\pm$ 0.002	0.252 $\pm$ 0.004	0.050 $\pm$ 0.001
	PP	0.420 $\pm$ 0.017	0.455 $\pm$ 0.008	0.612 $\pm$ 0.018	0.566 $\pm$ 0.008	0.802 $\pm$ 0.020
GPDM	MP	0.520 $\pm$ 0.001	0.530 $\pm$ 0.002	0.274 $\pm$ 0.001	0.388 $\pm$ 0.004	0.355 $\pm$ 0.008
	MR	0.141 $\pm$ 0.001	0.234 $\pm$ 0.001	0.261 $\pm$ 0.002	0.146 $\pm$ 0.005	0.050 $\pm$ 0.001
	PP	0.362 $\pm$ 0.021	0.185 $\pm$ 0.017	0.587 $\pm$ 0.016	0.509 $\pm$ 0.008	0.384 $\pm$ 0.028
TSBN	MP	0.594 $\pm$ 0.007	0.471 $\pm$ 0.001	0.360 $\pm$ 0.001	0.403 $\pm$ 0.012	0.788 $\pm$ 0.005
	MR	0.124 $\pm$ 0.001	0.158 $\pm$ 0.001	0.275 $\pm$ 0.001	0.194 $\pm$ 0.001	0.050 $\pm$ 0.001
	PP	0.418 $\pm$ 0.019	0.445 $\pm$ 0.031	0.611 $\pm$ 0.001	0.527 $\pm$ 0.003	0.692 $\pm$ 0.017
DTSBN-2	MP	0.439 $\pm$ 0.001	0.475 $\pm$ 0.002	0.370 $\pm$ 0.004	0.407 $\pm$ 0.003	0.756 $\pm$ 0.001
	MR	0.134 $\pm$ 0.001	0.208 $\pm$ 0.001	0.361 $\pm$ 0.001	0.248 $\pm$ 0.007	0.050 $\pm$ 0.001
	PP	0.391 $\pm$ 0.001	0.446 $\pm$ 0.001	0.587 $\pm$ 0.027	0.522 $\pm$ 0.005	0.737 $\pm$ 0.004
DTSBN-3	MP	0.411 $\pm$ 0.001	0.431 $\pm$ 0.001	0.390 $\pm$ 0.002	0.390 $\pm$ 0.002	0.774 $\pm$ 0.002
	MR	0.141 $\pm$ 0.001	0.189 $\pm$ 0.001	0.274 $\pm$ 0.001	0.252 $\pm$ 0.004	0.050 $\pm$ 0.001
	PP	0.367 $\pm$ 0.011	0.451 $\pm$ 0.026	0.548 $\pm$ 0.013	0.510 $\pm$ 0.006	0.715 $\pm$ 0.009
DPGDS-2	MP	0.688 $\pm$ 0.002	0.659 $\pm$ 0.001	0.379 $\pm$ 0.002	0.430 $\pm$ 0.009	0.867 $\pm$ 0.008
	MR	0.149 $\pm$ 0.001	0.242 $\pm$ 0.007	0.373 $\pm$ 0.001	0.254 $\pm$ 0.005	0.050 $\pm$ 0.001
	PP	0.443 $\pm$ 0.025	0.473 $\pm$ 0.012	0.622 $\pm$ 0.014	0.582 $\pm$ 0.007	0.814 $\pm$ 0.035
DPGDS-3	MP	<b>0.689</b> $\pm$ 0.002	0.660 $\pm$ 0.001	0.380 $\pm$ 0.001	0.431 $\pm$ 0.012	<b>0.887</b> $\pm$ 0.002
	MR	<b>0.150</b> $\pm$ 0.001	0.244 $\pm$ 0.003	<b>0.374</b> $\pm$ 0.002	<b>0.255</b> $\pm$ 0.004	0.050 $\pm$ 0.001
	PP	<b>0.456</b> $\pm$ 0.015	<b>0.478</b> $\pm$ 0.024	<b>0.628</b> $\pm$ 0.021	<b>0.600</b> $\pm$ 0.001	<b>0.839</b> $\pm$ 0.007

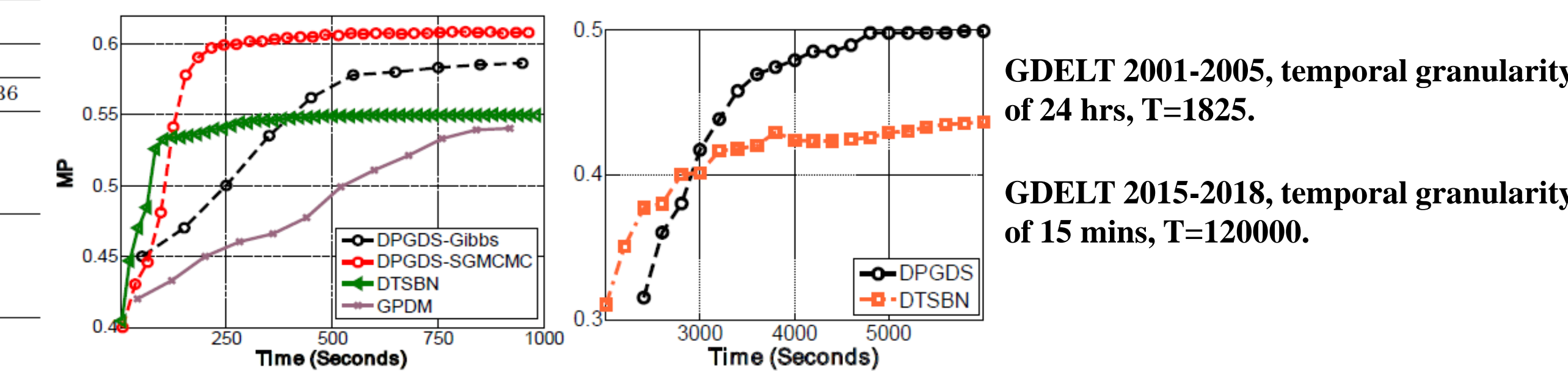
### Hierarchical topics learned from ICEWS 2007-2009



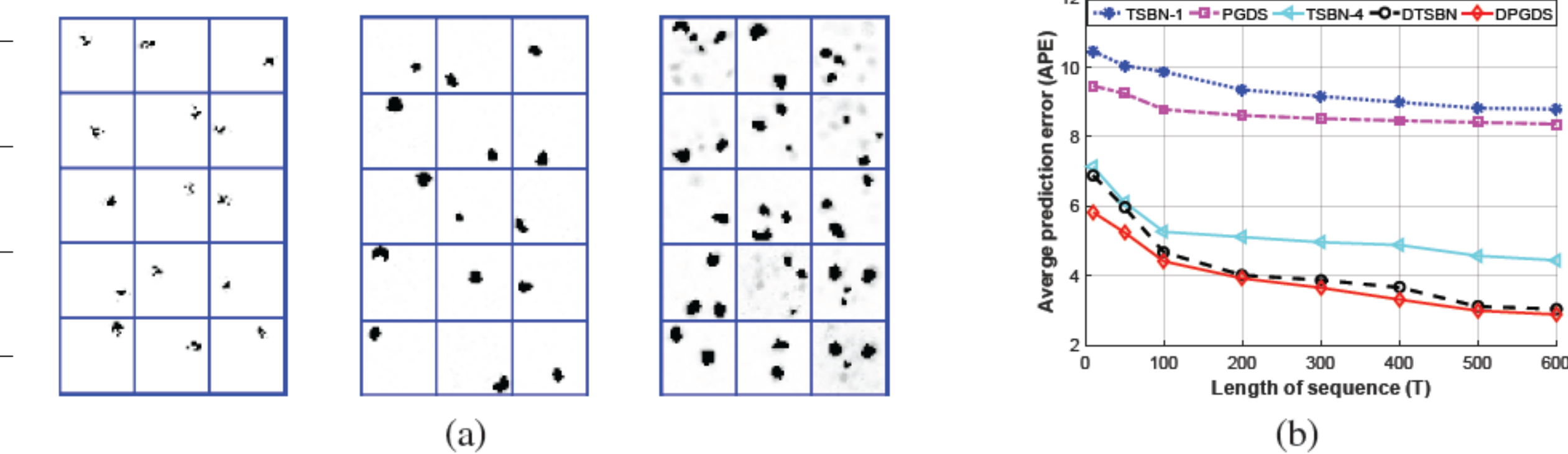
### Hierarchical topics learned from ICEWS 2001-2003



### Scalability

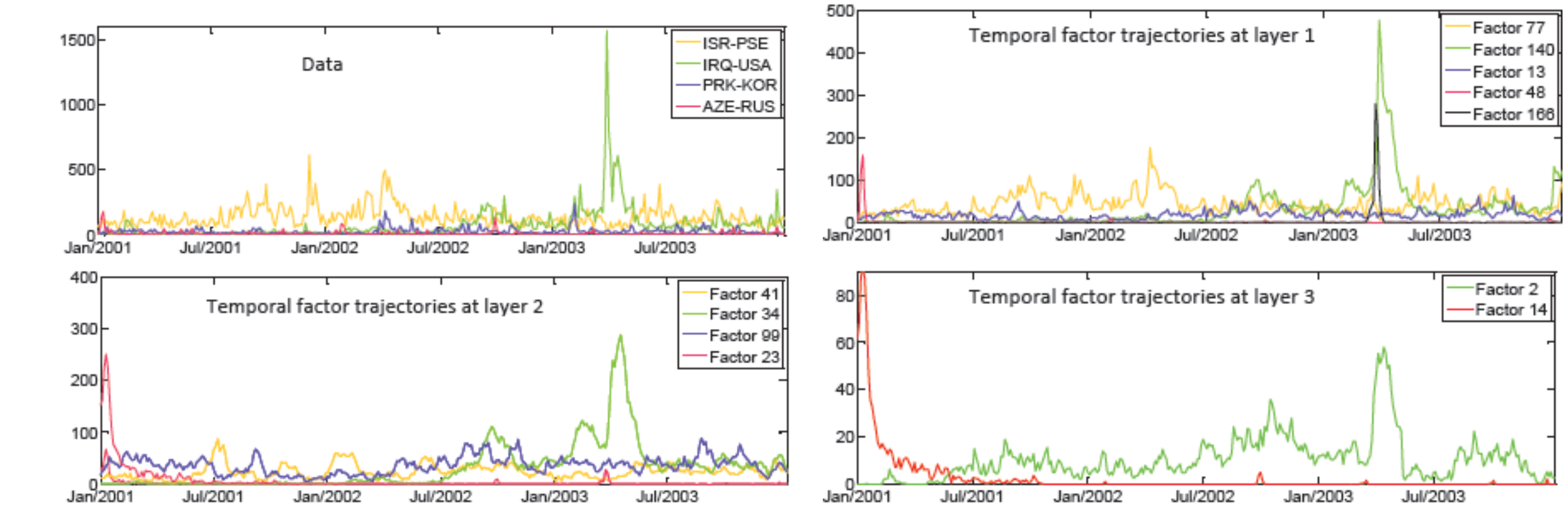


Shown in left and right are MP, as the function of time for GDEL 2001-2005, 2015-2018.

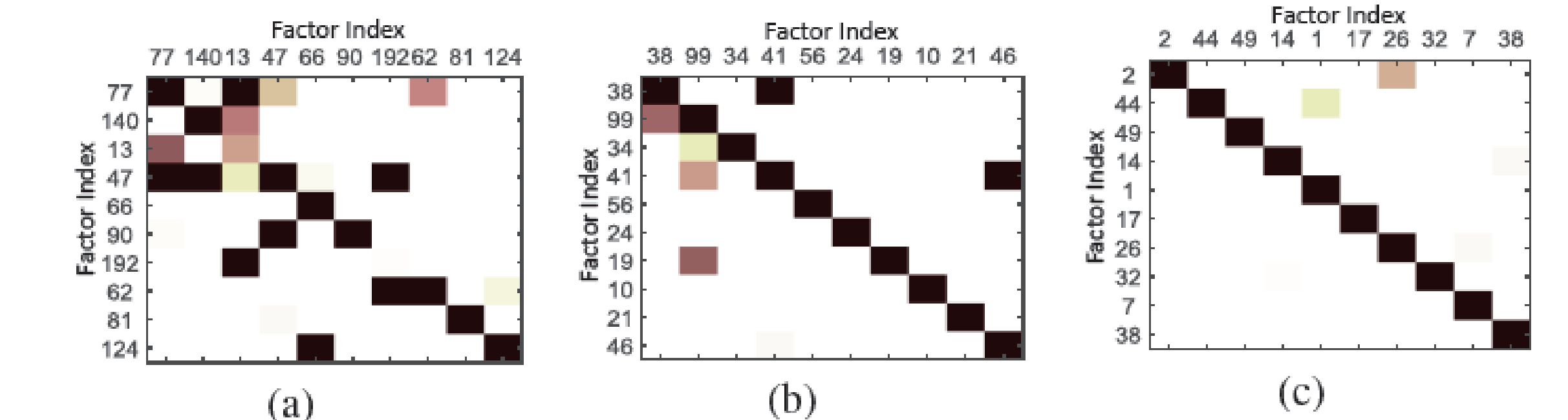


Results on the bouncing ball data set. (a) Top fifteen topics learned by three layer DPGDS at the layer 1, 2, 3. (b) APE as a function of the sequence length for various algorithms.

### Temporal trajectories at different layers for ICEWS 2001-2003



### Transition structure on ICEWS 2001-2003



Shown in (a)-(c) are transition matrices for layers 1, 2 and 3, respectively

[1] M. Zhou and L. Carin, "Negative binomial process count and mixture modeling," TPAMI, 2015.

[2] A. Schein, M. Zhou, and H. Wallach "Poisson-gamma dynamical systems," in NIPS, 2016.

[3] Z. Gan, C. Li, R. Henao, D. E. Carlson, and L. Carin, "Deep temporal sigmoid belief networks for sequence modeling," in NIPS, 2015.

[4] Y. Cong, B. Chen, H. Liu, and M. Zhou, "Deep latent Dirichlet allocation with topic-layer adaptive stochastic gradient Riemannian MCMC," in ICML, 2017.