

Compressed Nonparametric Language Modelling

Ehsan Shareghi

Monash University

Gholamreza Haffari

Monash University

Trevor Cohn

The University of Melbourne

Outline

- Infinite-Order Language Modelling and Challenges
- Compressed HPYP LM
- Inference and Sampling in Compressed HPYP LM
- Perplexity and Mixing
- Conclusion and Future Directions

Language Modelling (LM)

Donald trump is a p

donald trump is a **pokemon**

donald trump is a **potato**

donald trump is a **populist**

donald trump is a **politician**

donald trump is a **pragmatist**

donald trump is a **prophet**

donald trump is a **piece of garbage**

donald trump is a **pendejo**

Predictive typing/Auto completion

Président de la Chambre des représentants

President of the Bedroom of Representatives

President of the House of Representatives

$P(\text{House} \mid \text{President of the}) > P(\text{Bedroom} \mid \text{President of the})$

Machine Translation

Infinite order LM

$$P(w_1^N) = \prod_{i=1}^N P(w_i | w_1^{i-1})$$

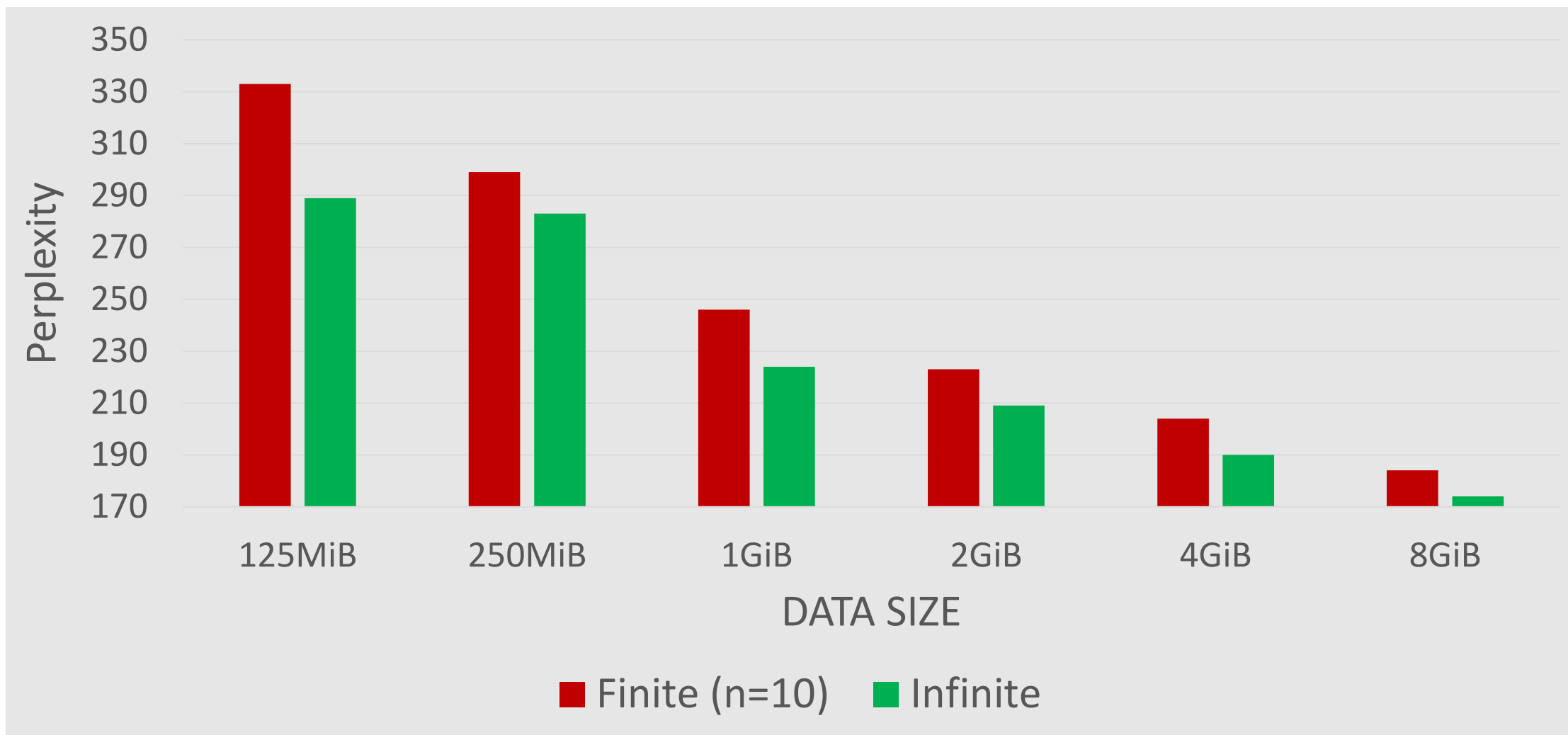
Statistical sparsity

Solution: smoothing (HPYP, etc)

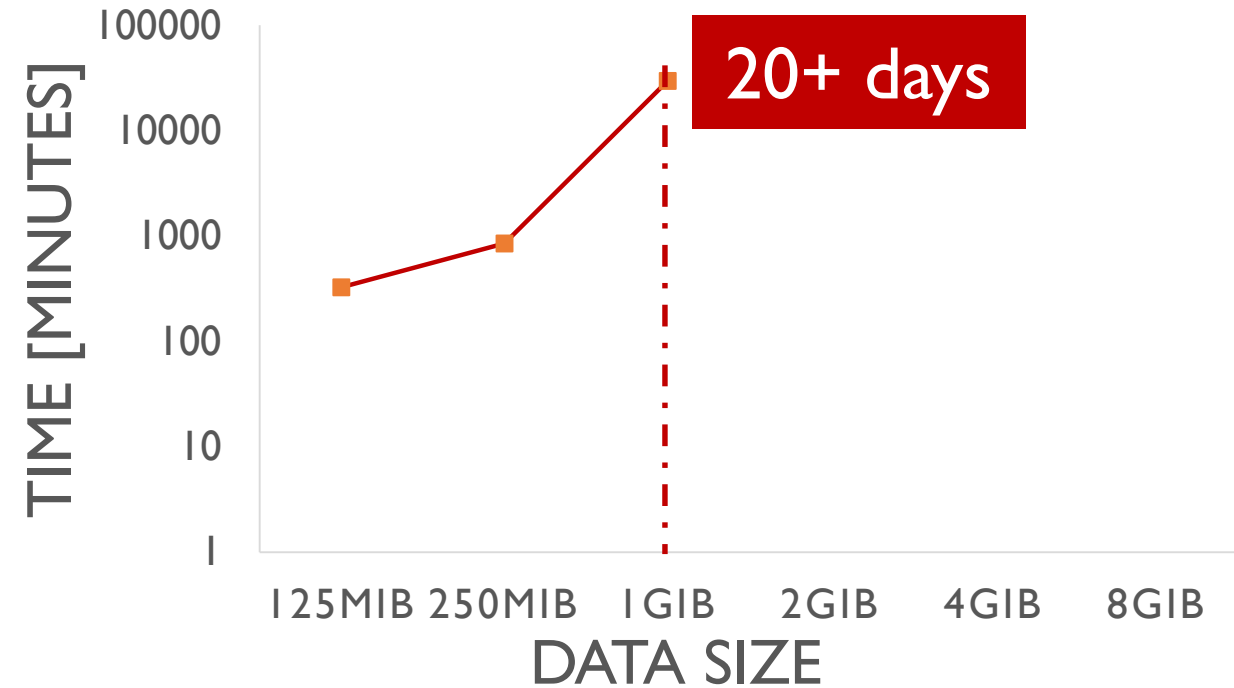
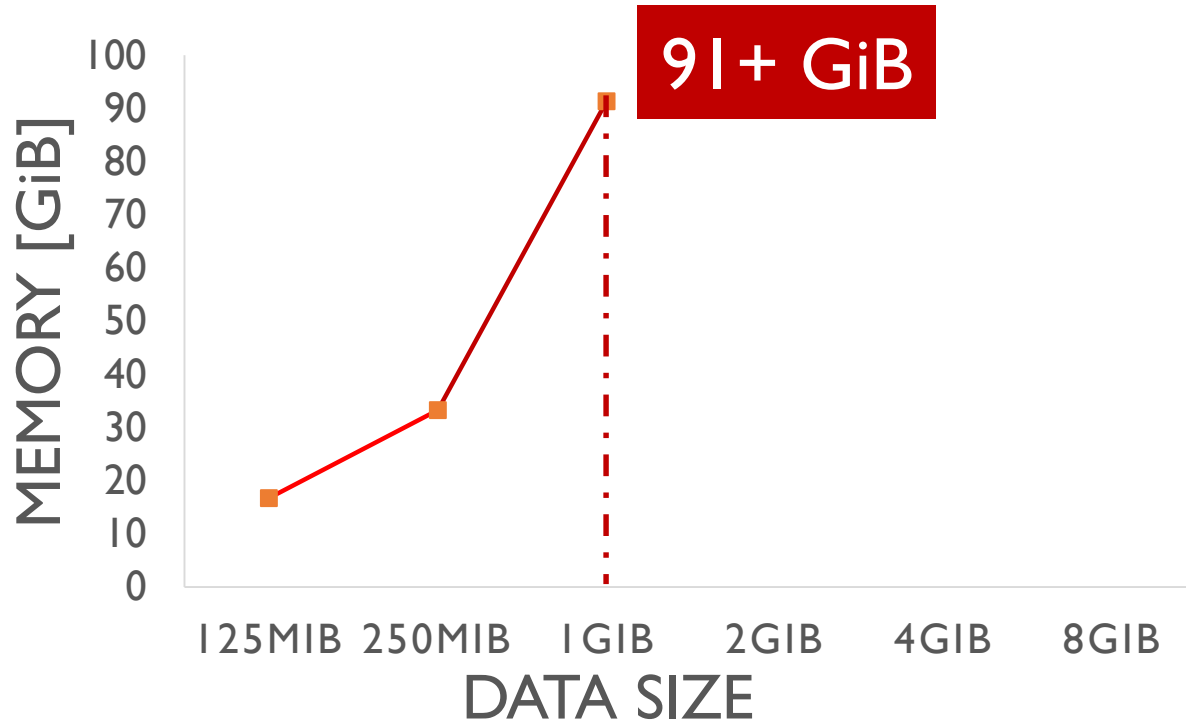
Computational cost of smoothing

No Scalable Solution

Why Infinite-order LM ?



Computational Cost of HPYP LM - Training



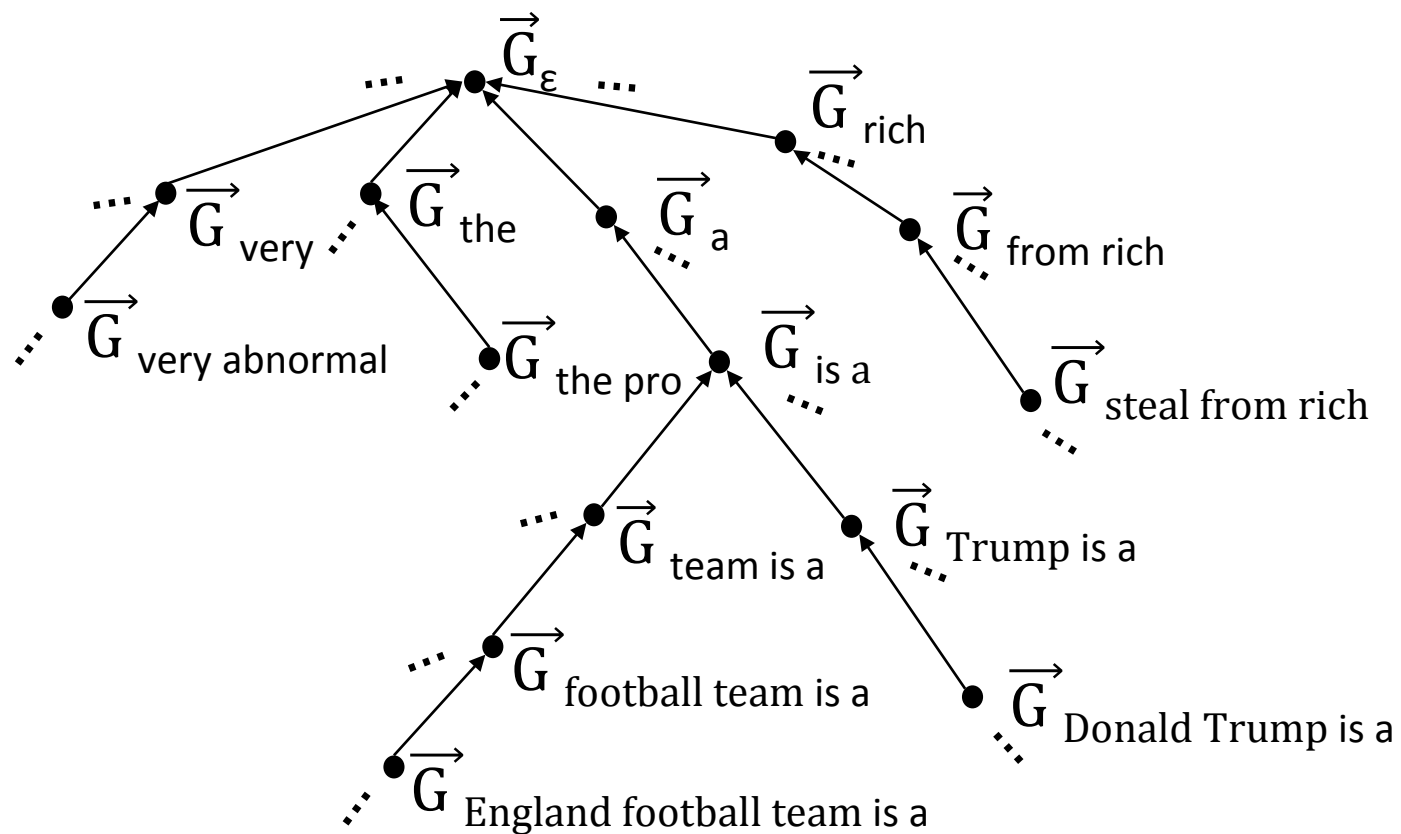
Involved Factors:

- Building Model (hierarchy)
- Parameters Sampling
- Storing the Model and Parameters

Outline

- Infinite-Order Language Modelling and Challenges
- **Compressed HPYP LM**
- Inference and Sampling in Compressed HPYP LM
- Perplexity and Mixing
- Conclusion and Future Directions

Hierarchical Pitman-Yor Process (HPYP)



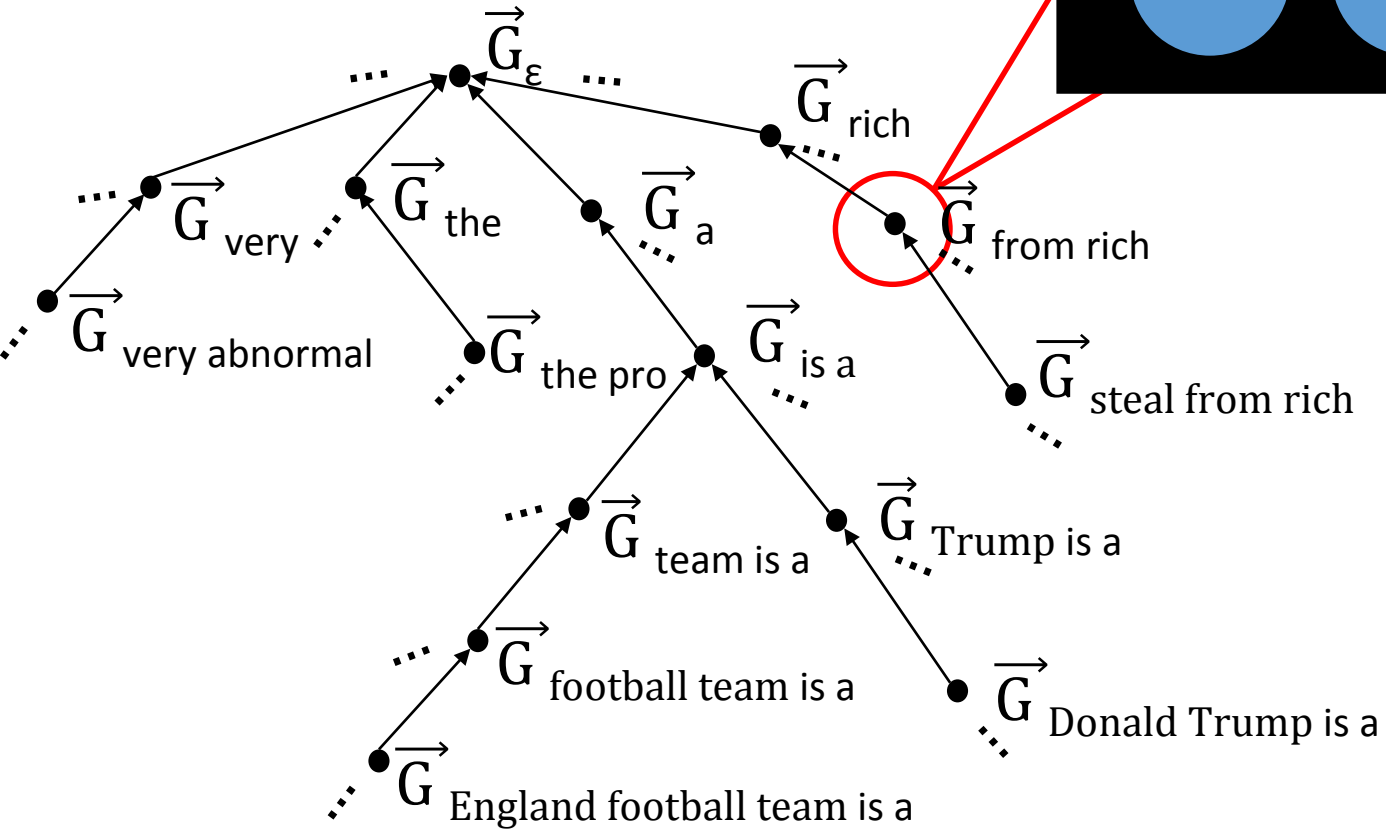
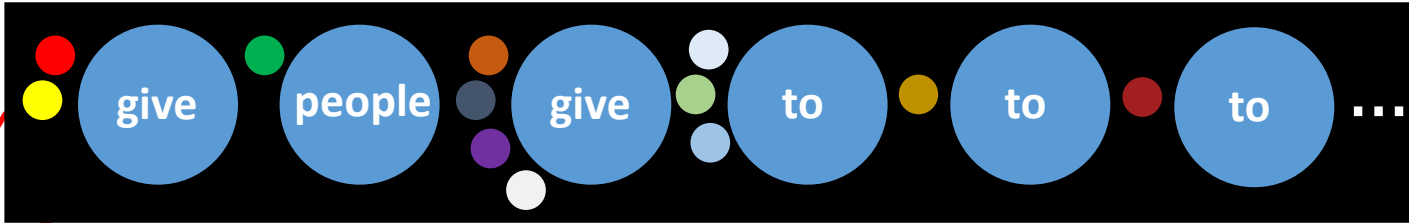
$$\vec{G}_u \sim \text{PYP}(\theta_u, d_u, \vec{G}_{\pi(u)})$$

$$\vec{G}_\epsilon \sim \text{PYP}(\theta_\epsilon, d_\epsilon, \frac{1}{|\text{vocab}|})$$

Same model as “Sequence Memoizer” Wood et al. (2011)

HPYP LM – Chinese Restaurant Process

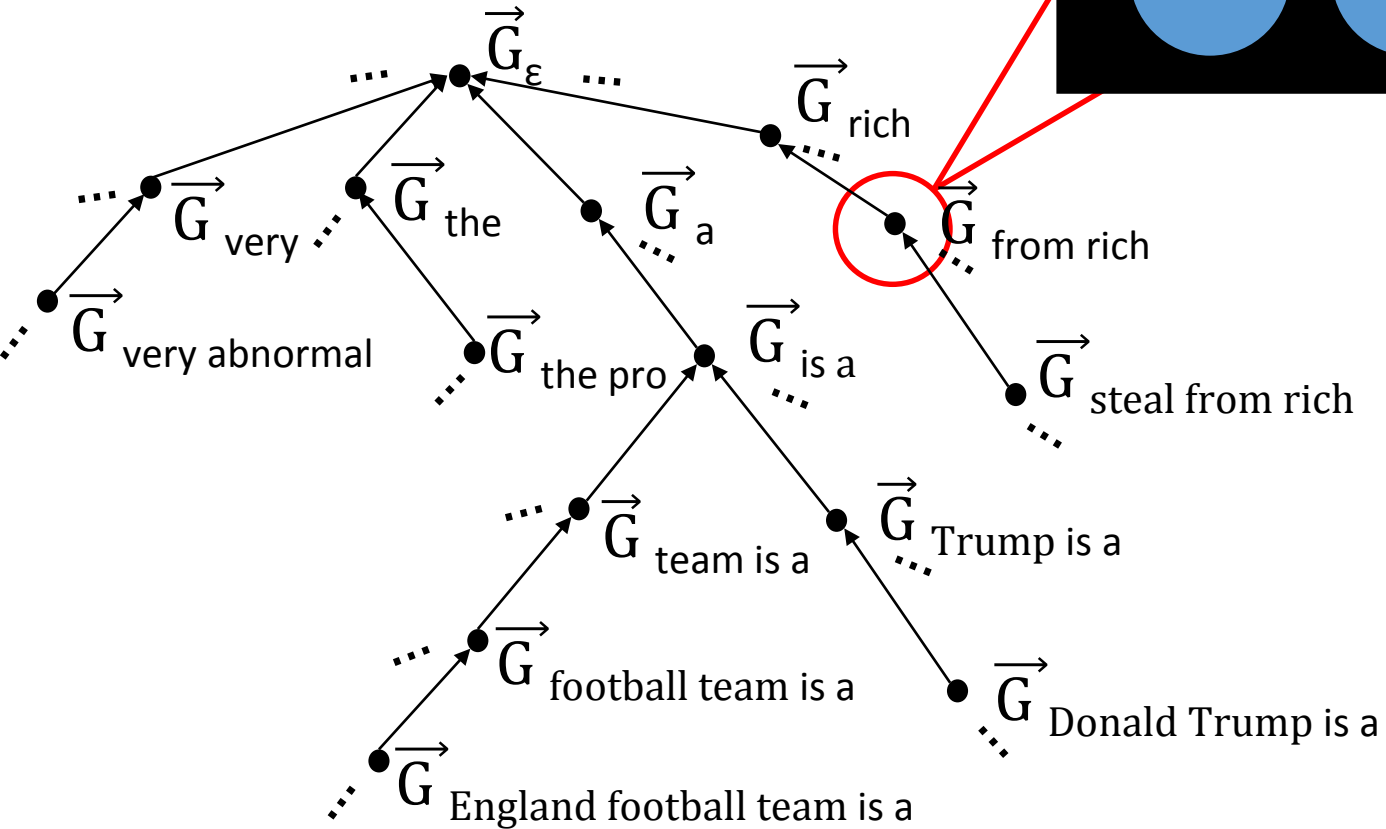
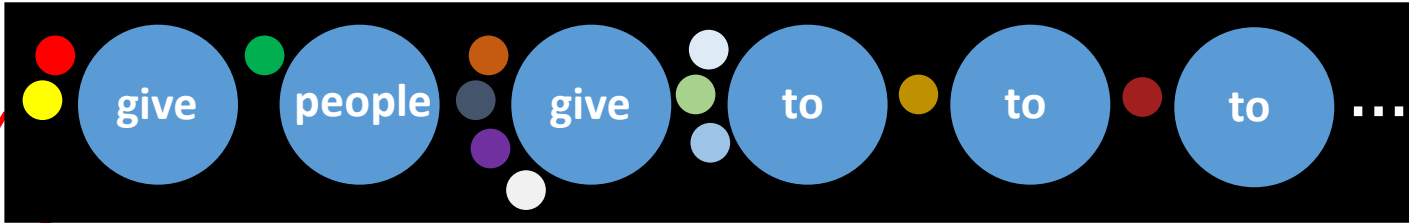
$u = \text{“from rich”}, \eta_u = \{d_u, \theta_u, \{n_u^w, t_u^w\}_{w \in u}\}$



$n_u^{give} = 6, t_u^{give} = 2$

HPYP LM – Chinese Restaurant Process

$u = \text{“from rich”}, \eta_u = \{d_u, \theta_u, \{n_u^w, t_u^w\}_{w \in u}\}$



$n_u^{give} = 6, t_u^{give} = 2$

$0 \leq t_u^w \leq n_u^w$

$n_{\pi(u)}^w = \sum_{v \in \text{children}(\pi(u))} t_v^w$

Compressed HPYP LM

- Hierarchy of KN and HPYP LMs are the same
- KN can serve as an approximate inference for HPYP ($\theta_u = 0$ and $t_u^W = I$)

Compressed HPYP LM

- Hierarchy of KN and HPYP LMs are the same
- KN can serve as an approximate inference for HPYP ($\theta_u = 0$ and $t_u^W = 1$)
- KN hierarchy can be recovered from a **compressed suffix tree** of data on-the-fly
- **Compressed Suffix Trees:**
 - Based on advanced data structures such as the **Wavelet Tree of the BWT of text**
 - Contain **all** the information about the **HPYP hierarchy and text itself** in a space matching the text size

Compressed HPYP LM

- Hierarchy of KN and HPYP LMs are the same
- KN can serve as an approximate inference for HPYP ($\theta_u = 0$ and $t_u^w = 1$)
- KN hierarchy can be recovered from a **compressed suffix tree** of data on-the-fly
- **Compressed Suffix Trees:**
 - Based on advanced data structures such as the **Wavelet Tree of the BWT of text**
 - Contain **all** the information about the **HPYP hierarchy and text itself** in a space matching the text size

Compressed HPYP

Constructs Compressed Suffix Tree of Data

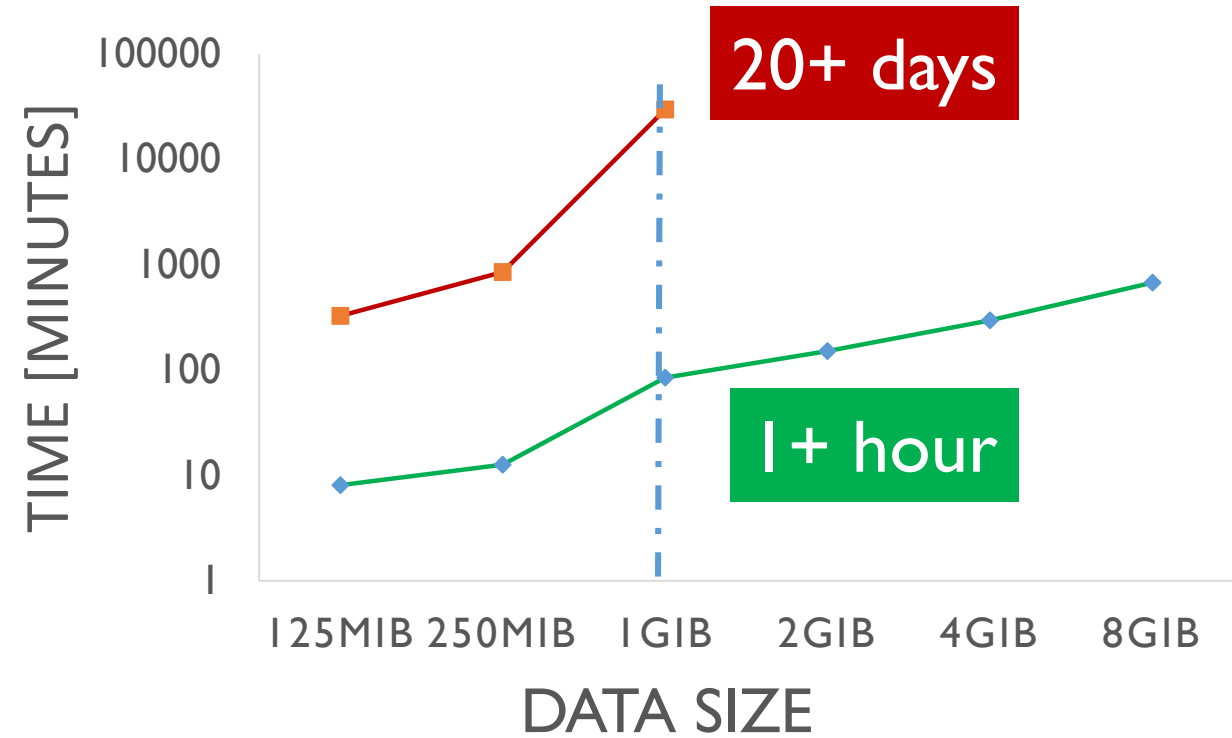
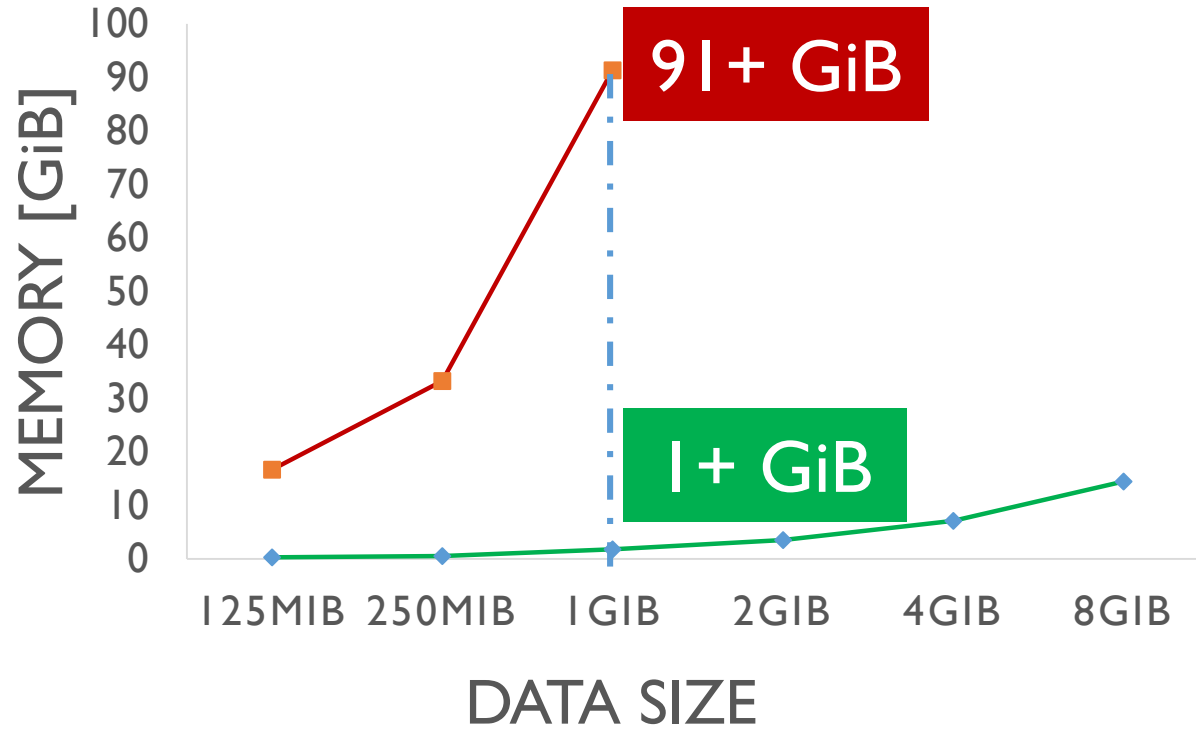
No Sampling

HPYP

Constructs Hierarchy of HPYP

Samples across all nodes and for all w

Training time comparison



Outline

- Infinite-Order Language Modelling and Challenges
- Compressed HPYP LM
- Inference and Sampling in Compressed HPYP LM
- Perplexity and Mixing
- Conclusion and Future Directions

Inference in Compressed HPYP

We need to compute the following **intractable** integral,

$$P(w|u) = \int P(w|u, \eta) P(\eta) d\eta$$

which we approximate using samples for $\eta_u = \{d_u, \theta_u, \{n_u^w, t_u^w\}_{w \in u}\}$.

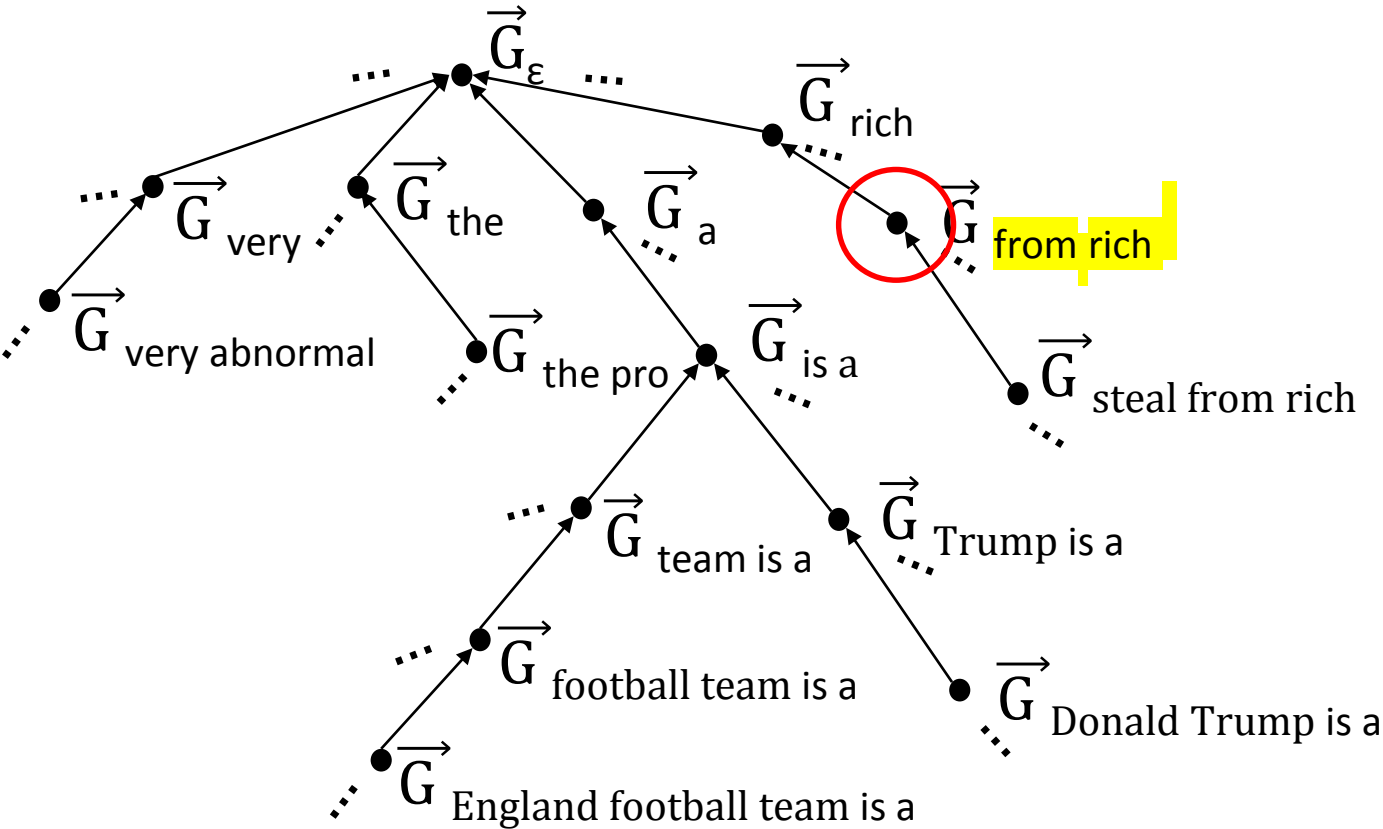
$$\begin{aligned} 0 &\leq t_u^w \leq n_u^w \\ n_{\pi(u)}^w &= \sum_{v \in \text{children}(\pi(u))} t_v^w \end{aligned}$$

Sampling $\{n_u^w, t_u^w\}_{w \in u}$

Given a query $P(\text{give} \mid \text{from rich})$

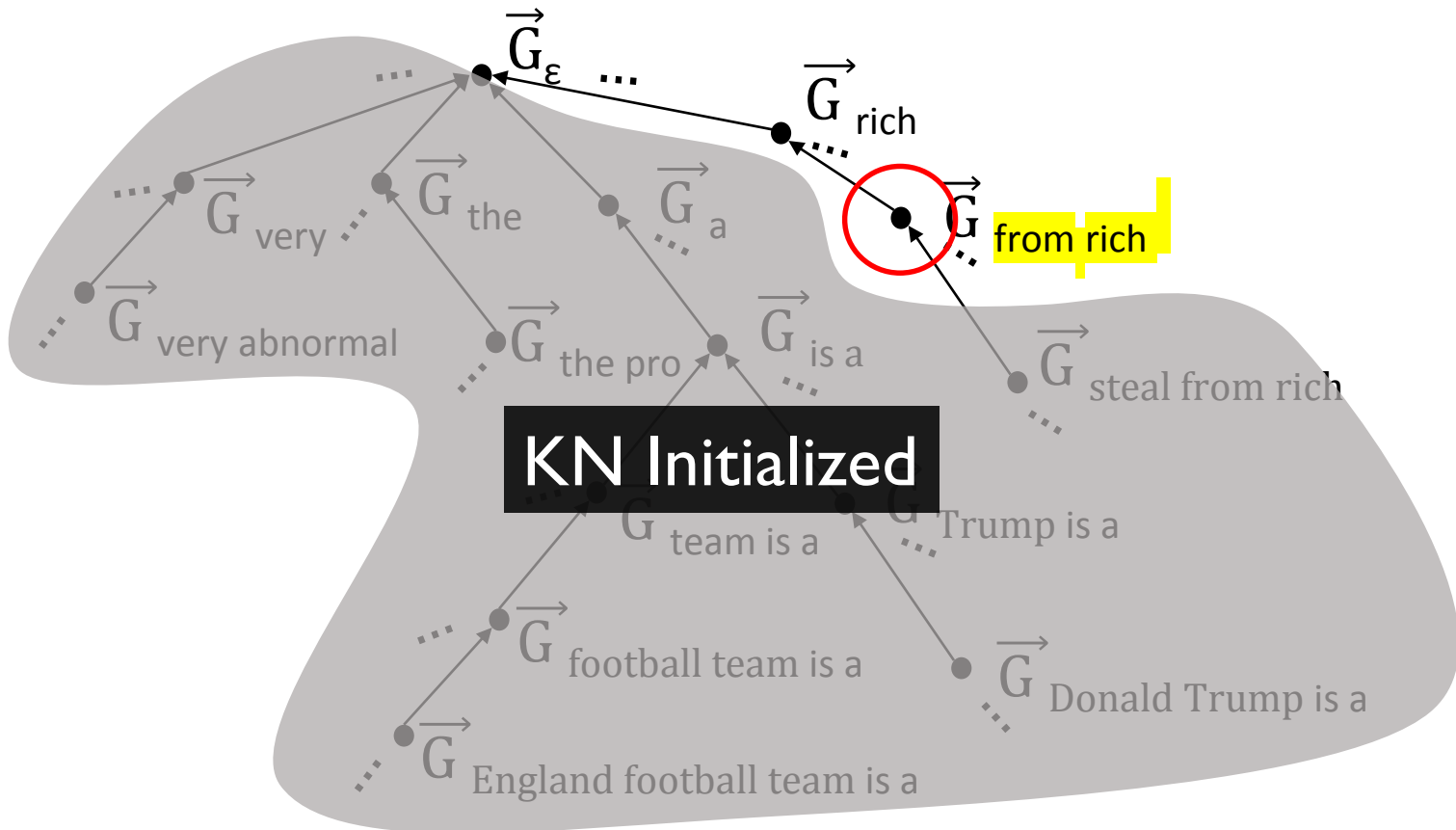
Sampling $\{n_u^w, t_u^w\}_{w \in \mathcal{U}}$

Given a query $P(\text{give} \mid \text{from rich})$



Sampling $\{n_u^w, t_u^w\}_{w \in \mathcal{U}}$

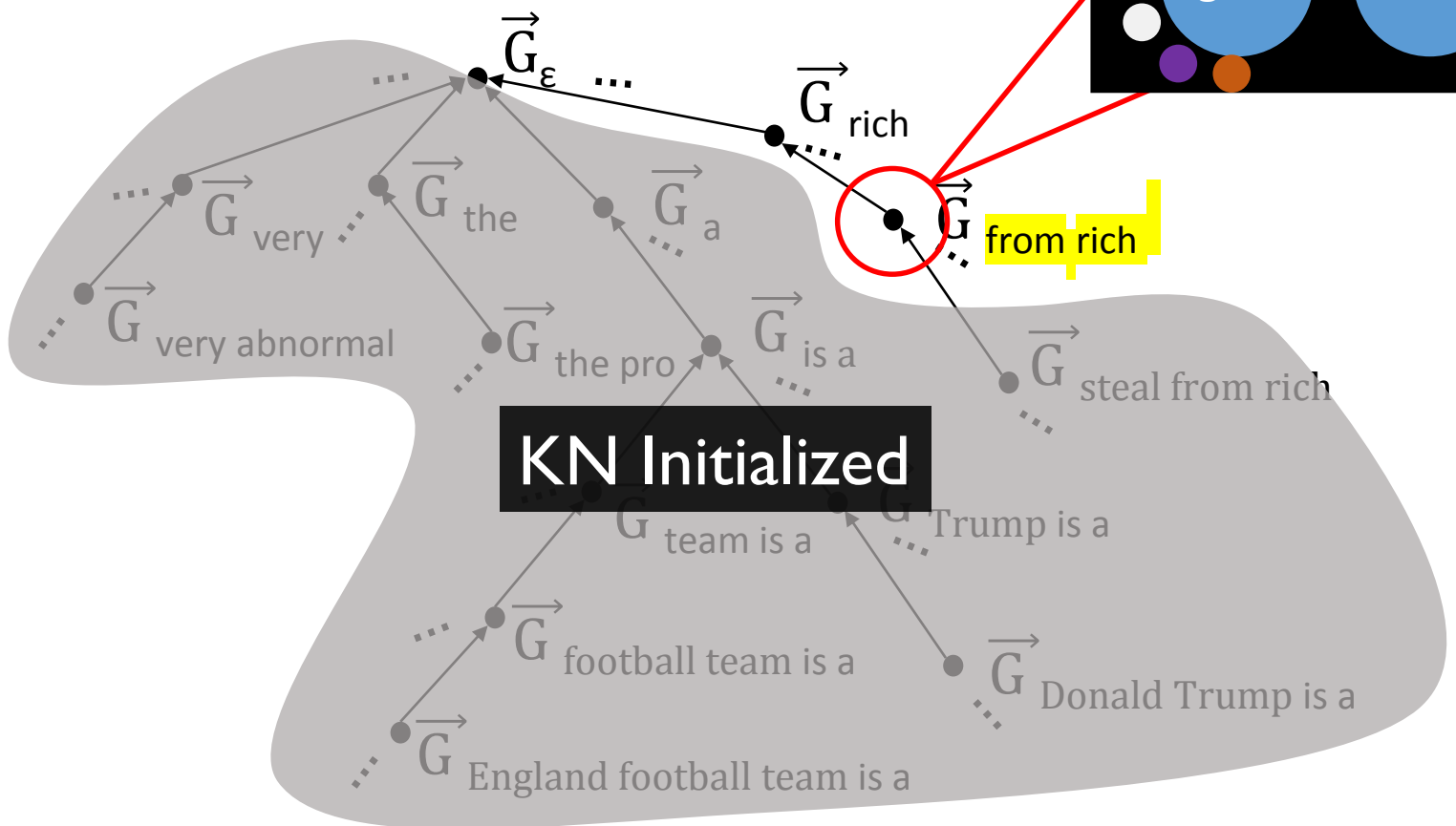
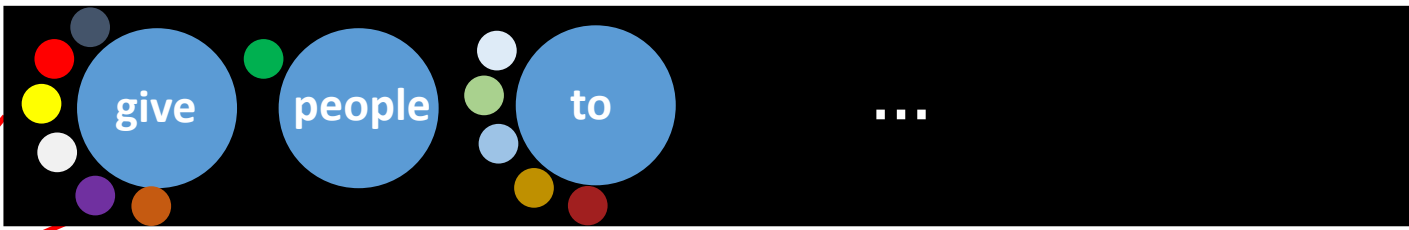
Given a query $P(\text{give} \mid \text{from rich})$



Sampling $\{n_u^w, t_u^w\}_{w \in u}$

Given a query $P(\text{give} \mid \text{from rich})$

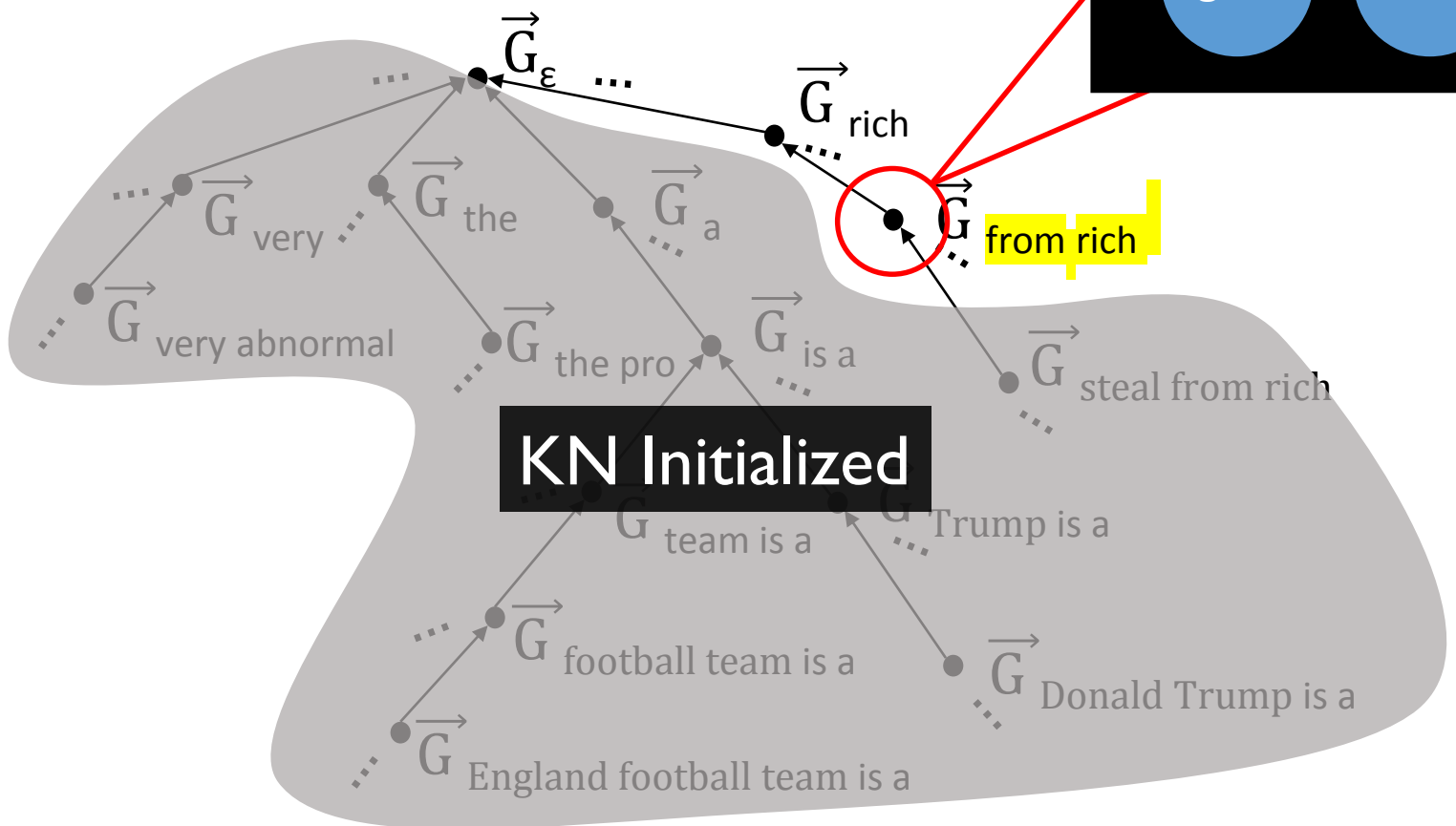
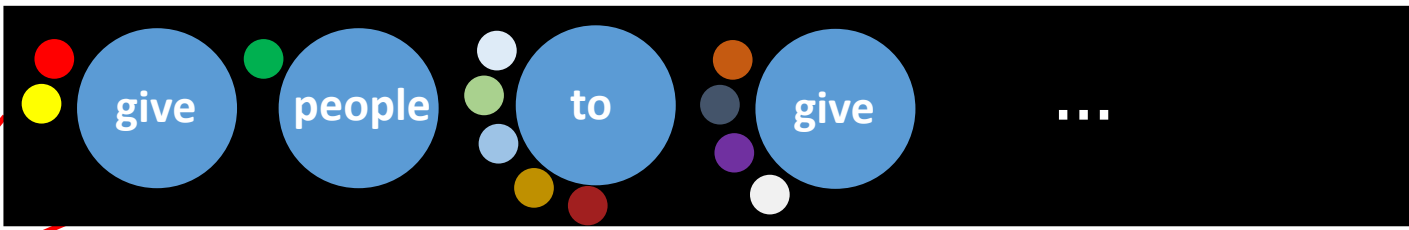
Read $n_{\text{from rich}}^{\text{give}}$ from Data



Sampling $\{n_u^w, t_u^w\}_{w \in \mathcal{U}}$

Given a query $P(\text{give} \mid \text{from rich})$

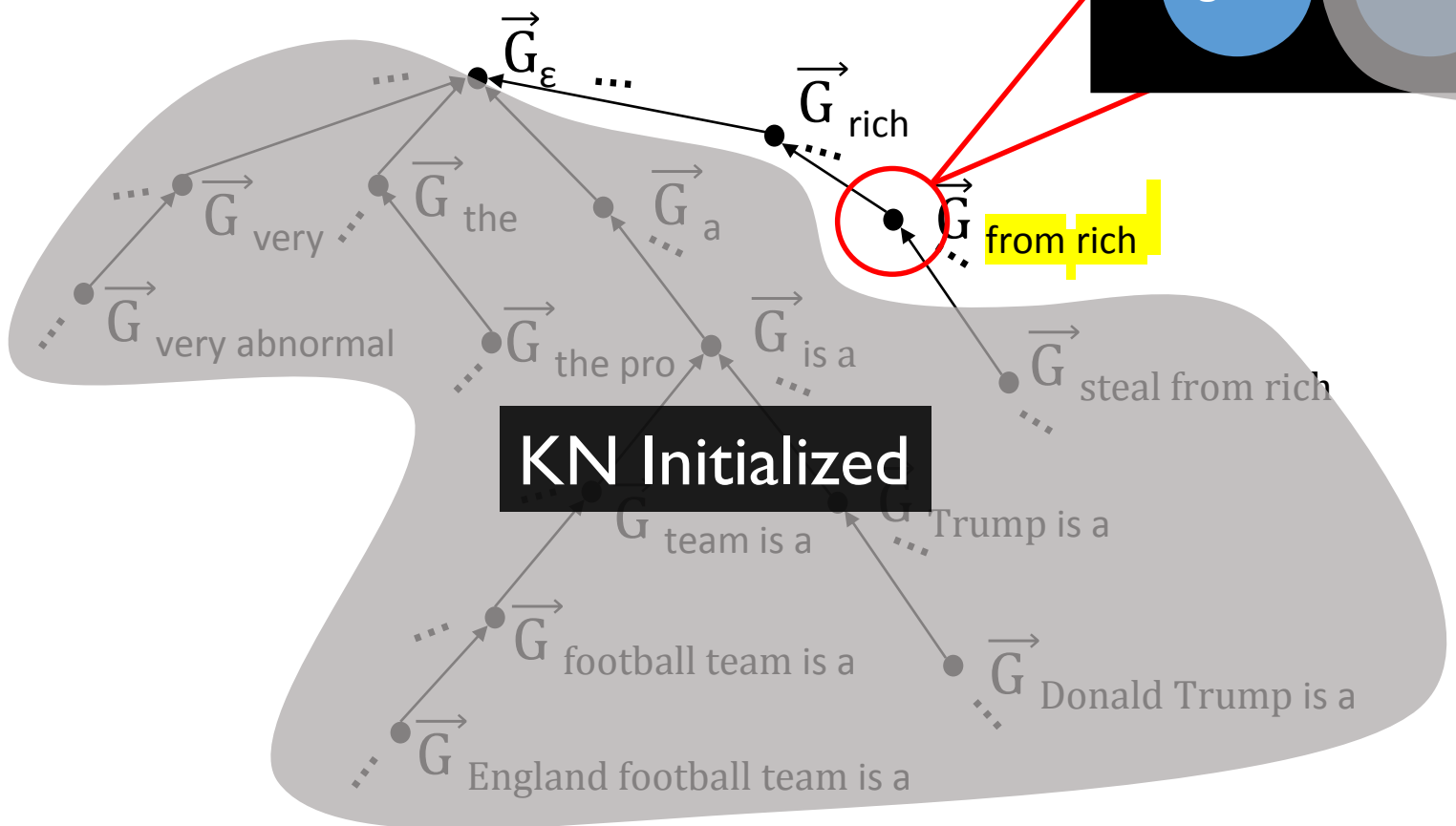
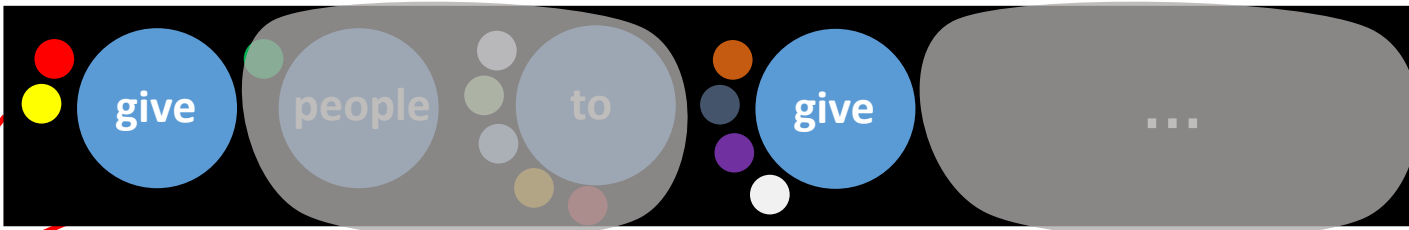
Sample $t_{\text{from rich}}^{\text{give}} = 2$



Sampling $\{n_u^w, t_u^w\}_{w \in U}$

Given a query $P(\text{give} \mid \text{from rich})$

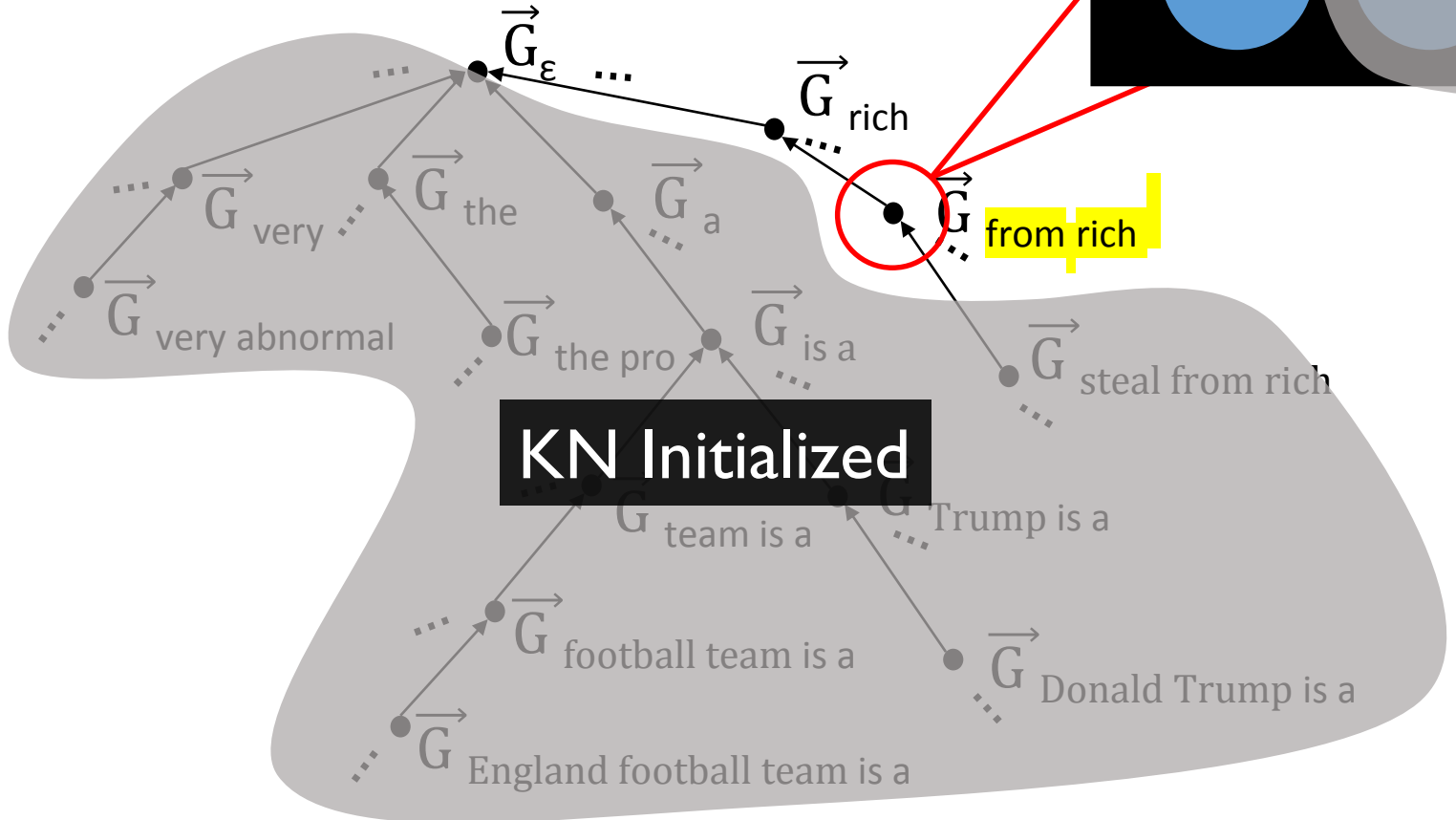
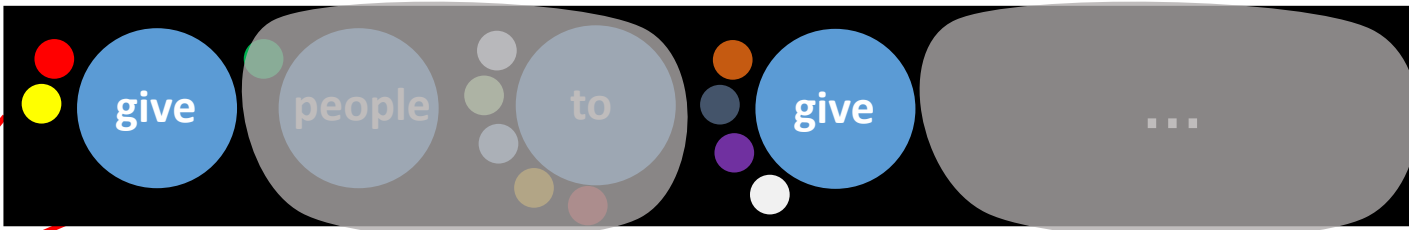
Sample $t_{\text{from rich}}^{\text{give}} = 2$



Sampling $\{n_u^w, t_u^w\}_{w \in \mathcal{U}}$

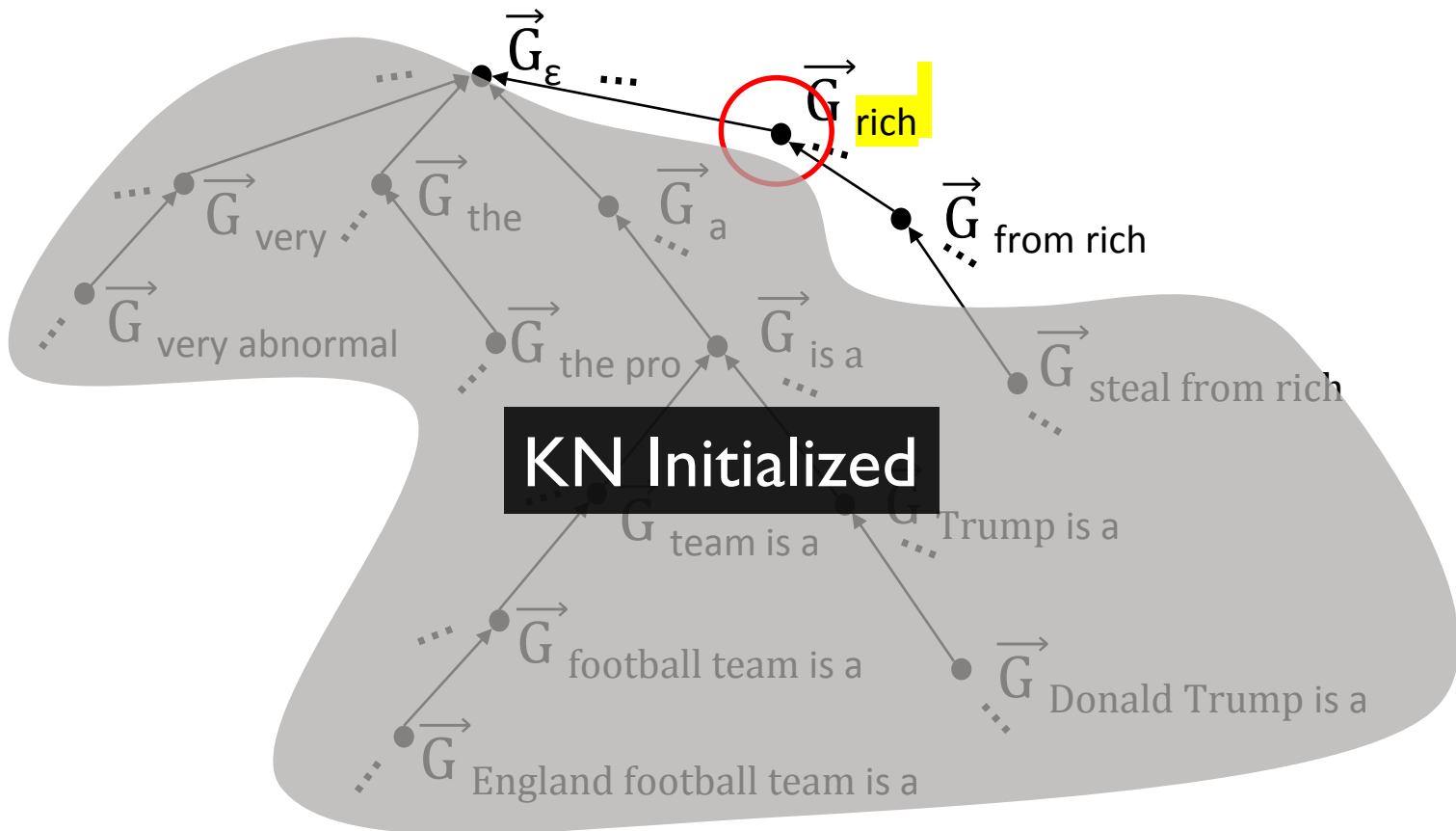
Given a query $P(\text{give} \mid \text{from rich})$

Update $n_{\text{rich}}^{\text{give}}$



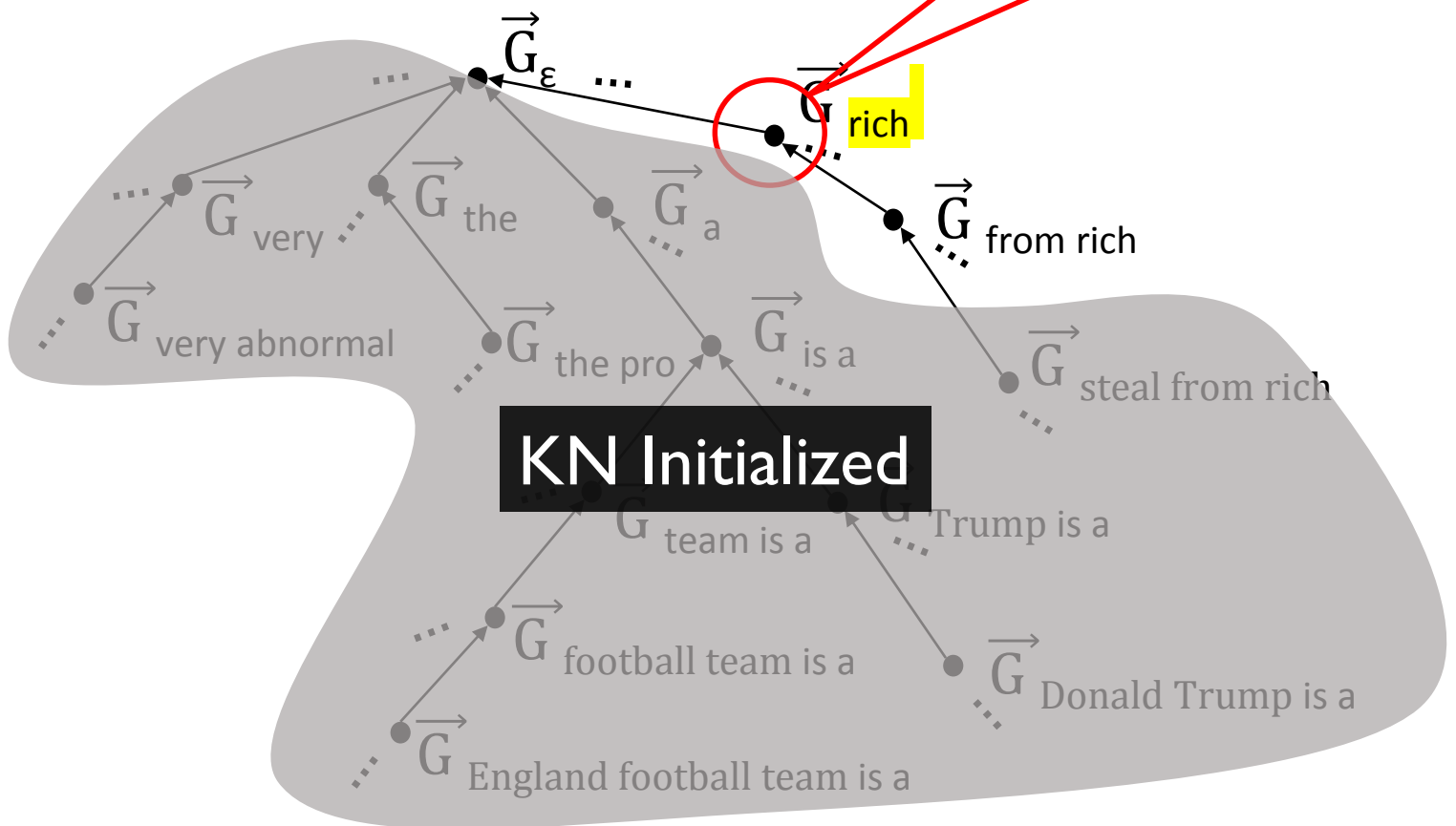
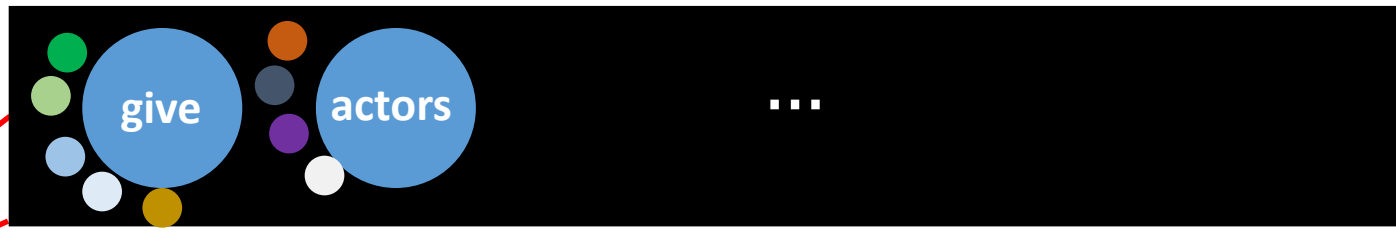
Sampling $\{n_u^w, t_u^w\}_{w \in U}$

Given a query $P(\text{give} \mid \text{from rich})$



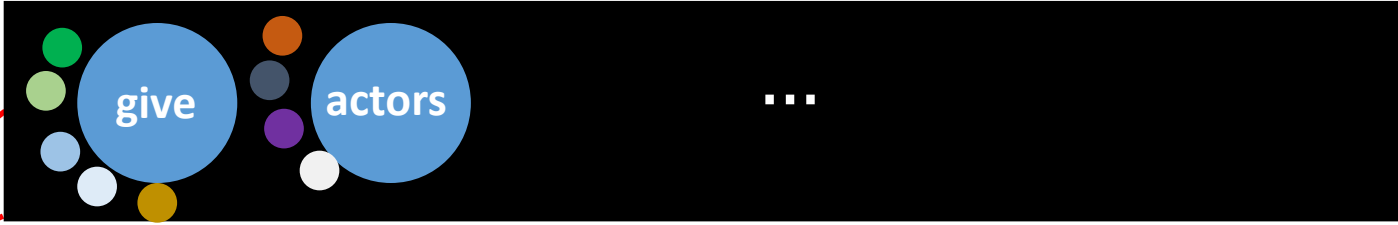
Sampling $\{n_u^w, t_u^w\}_{w \in U}$

Given a query $P(\text{give} \mid \text{from rich})$

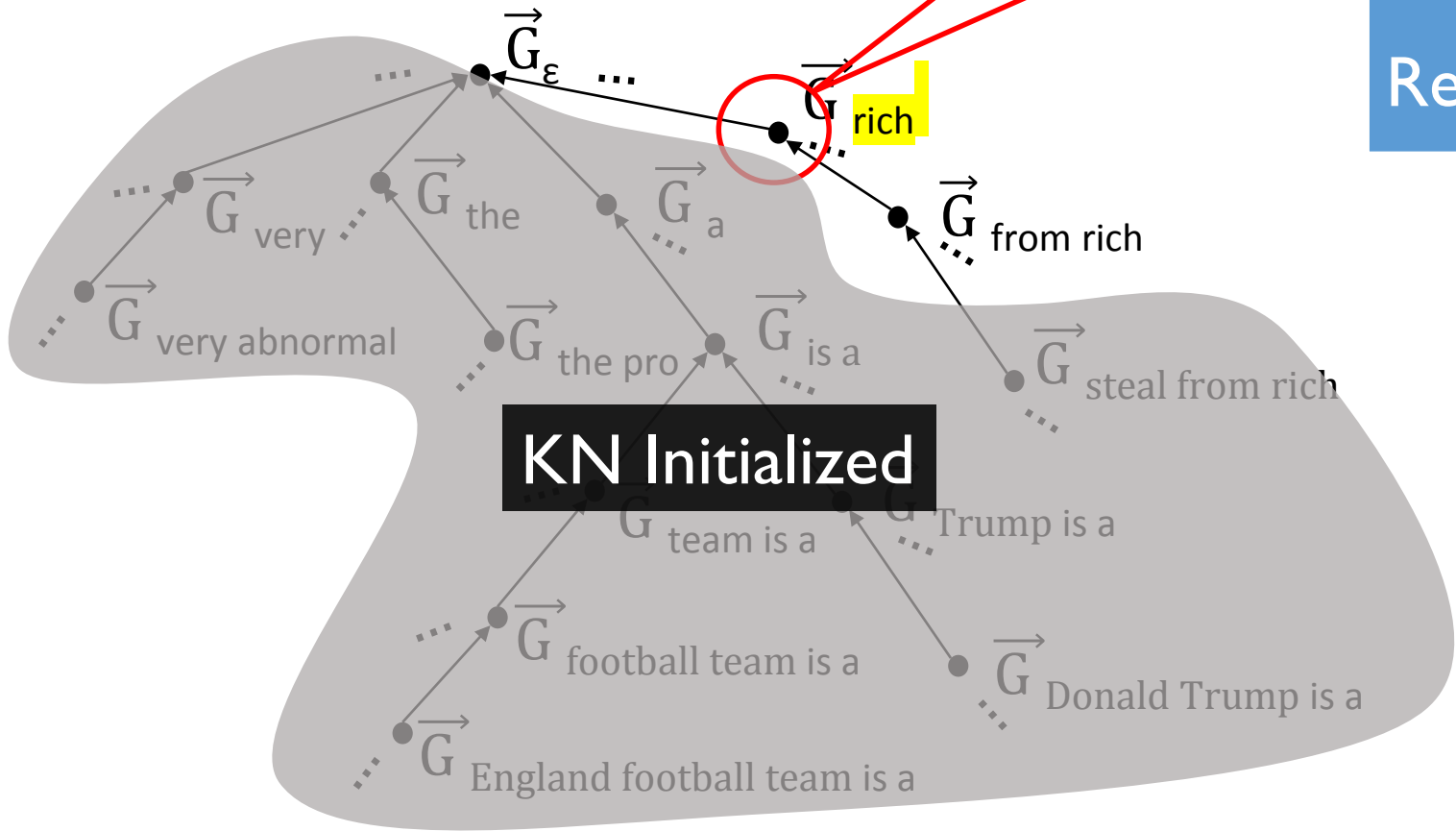


Sampling $\{n_u^w, t_u^w\}_{w \in \mathcal{U}}$

Given a query $P(\text{give} \mid \text{from rich})$

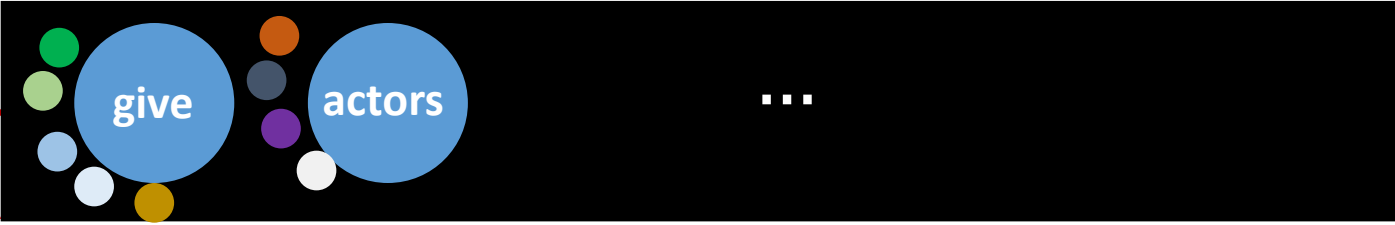


Read $n_{\text{rich}}^{\text{give}}$ from proxy counts



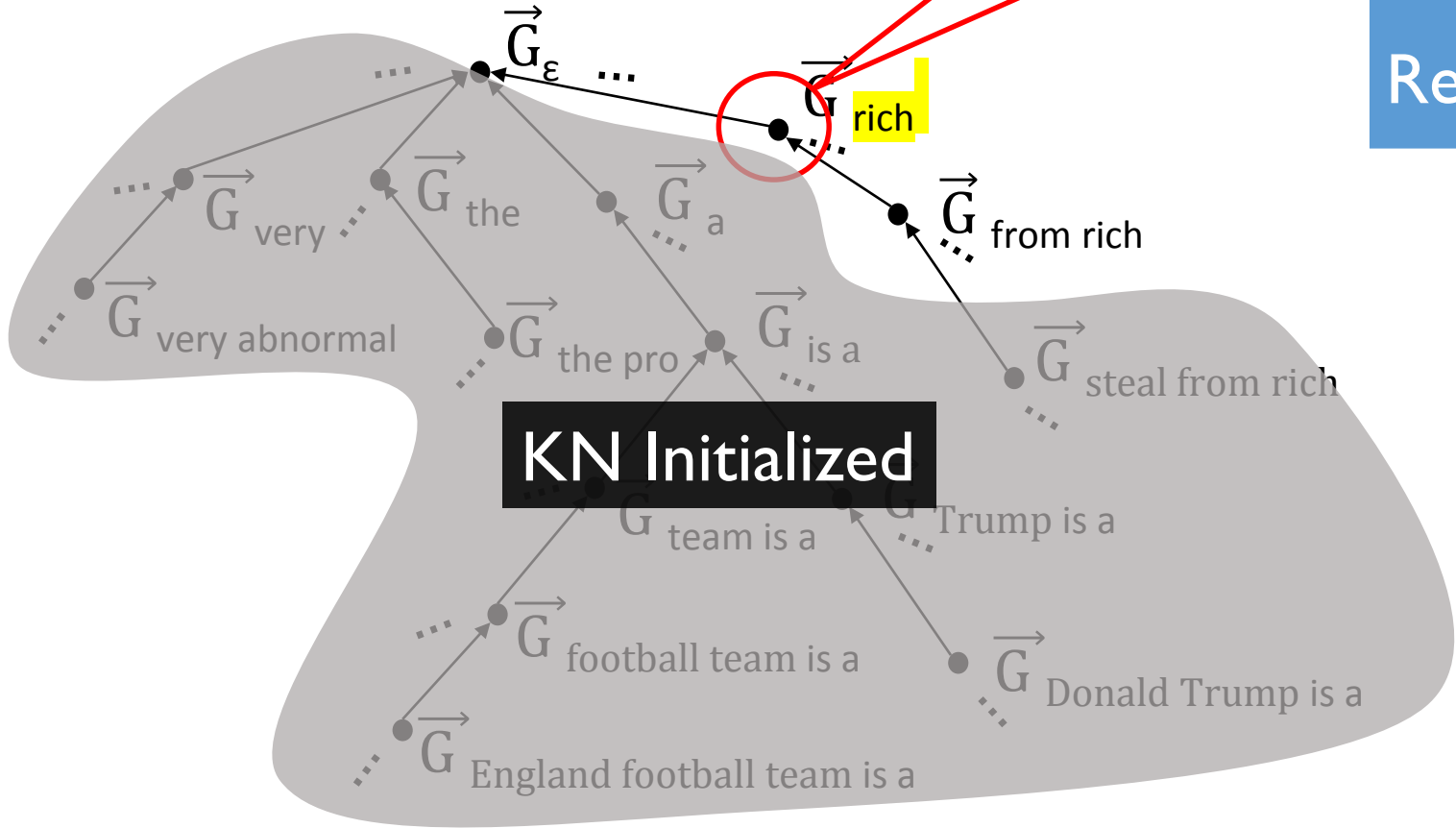
Sampling $\{n_u^w, t_u^w\}_{w \in \mathcal{U}}$

Given a query $P(\text{give} \mid \text{from rich})$



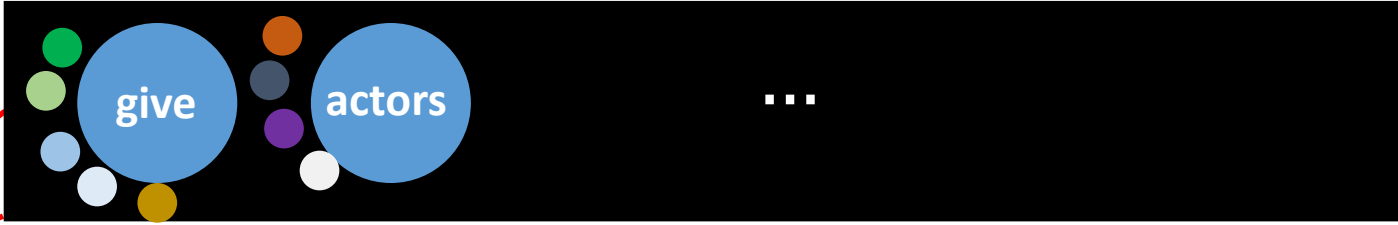
Read $n_{\text{rich}}^{\text{give}}$ from proxy counts

and so on



Sampling $\{n_u^w, t_u^w\}_{w \in \mathcal{U}}$

Given a query $P(\text{give} \mid \text{from rich})$

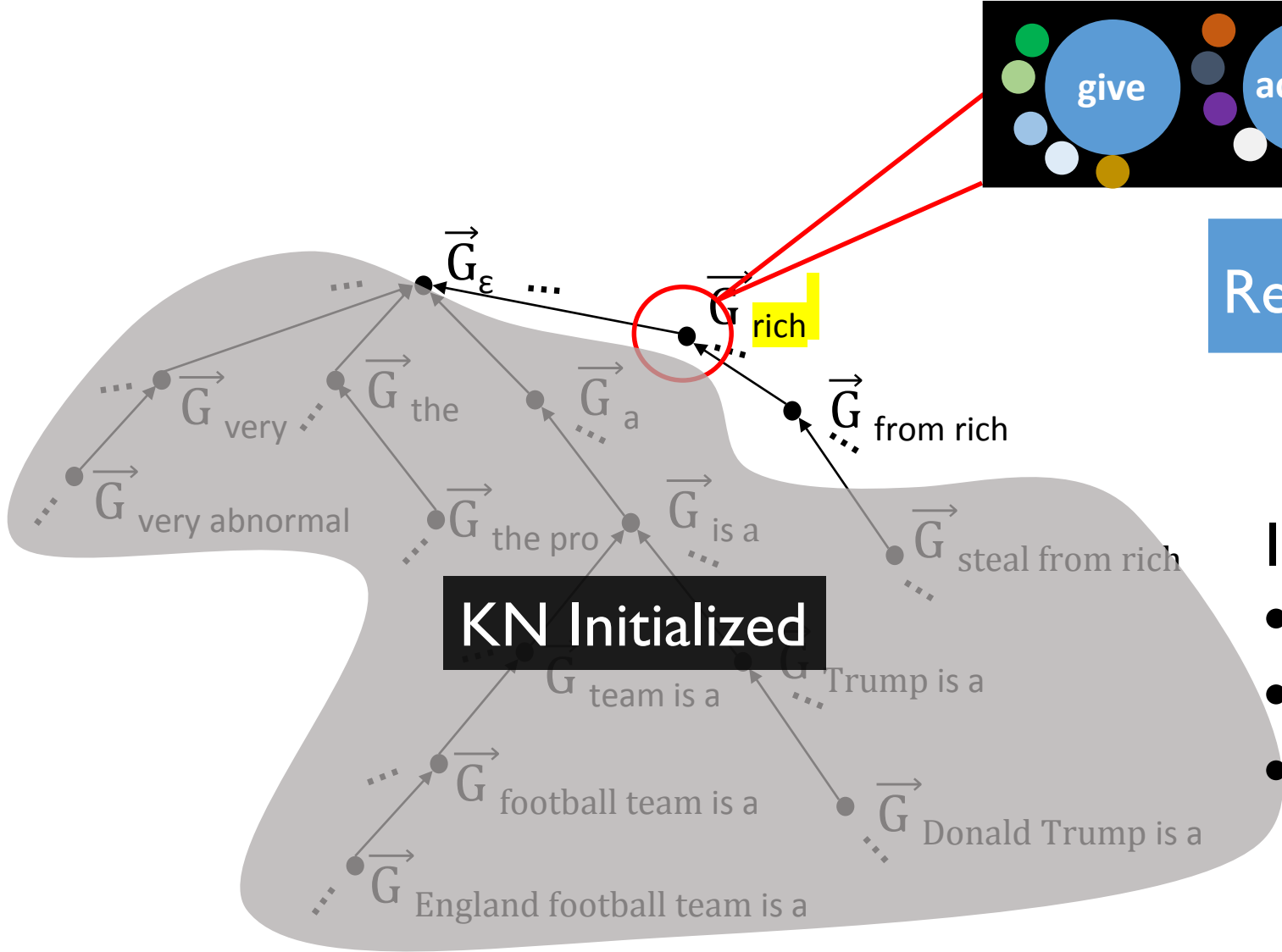


Read $n_{\text{rich}}^{\text{give}}$ from proxy counts

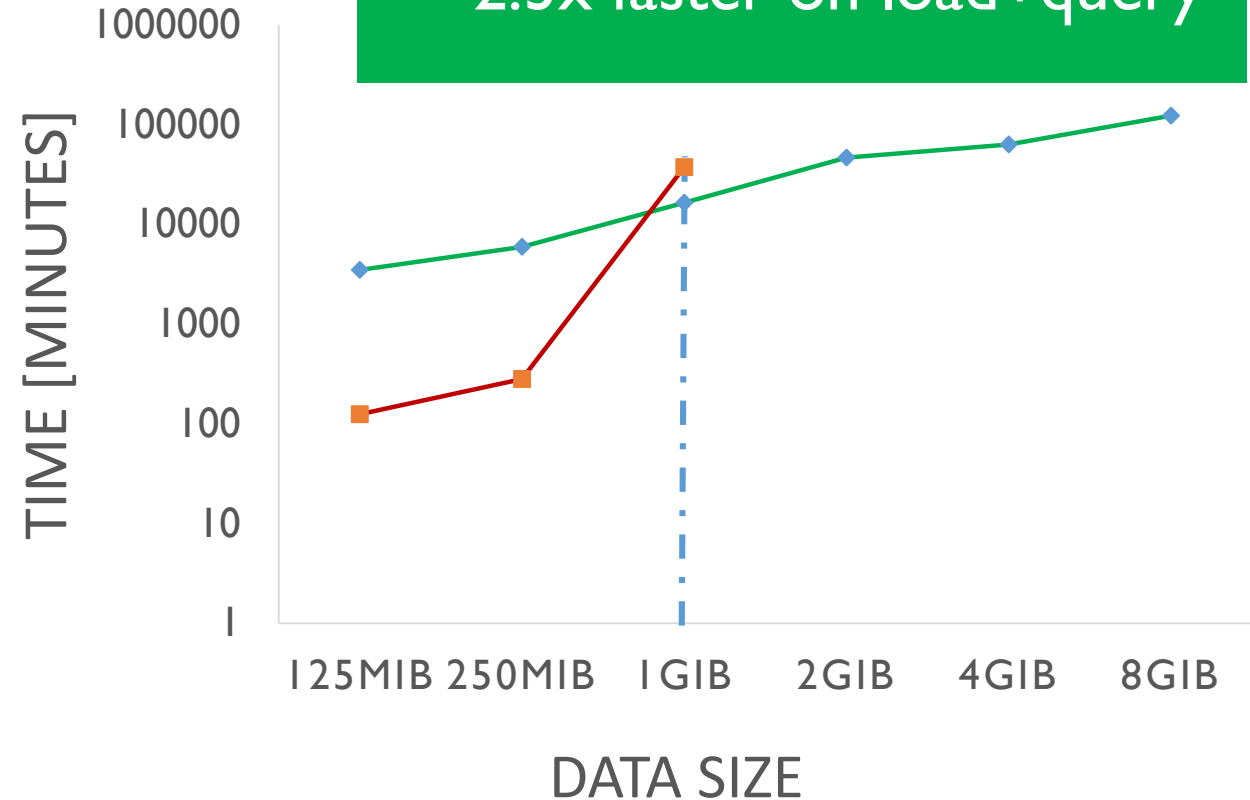
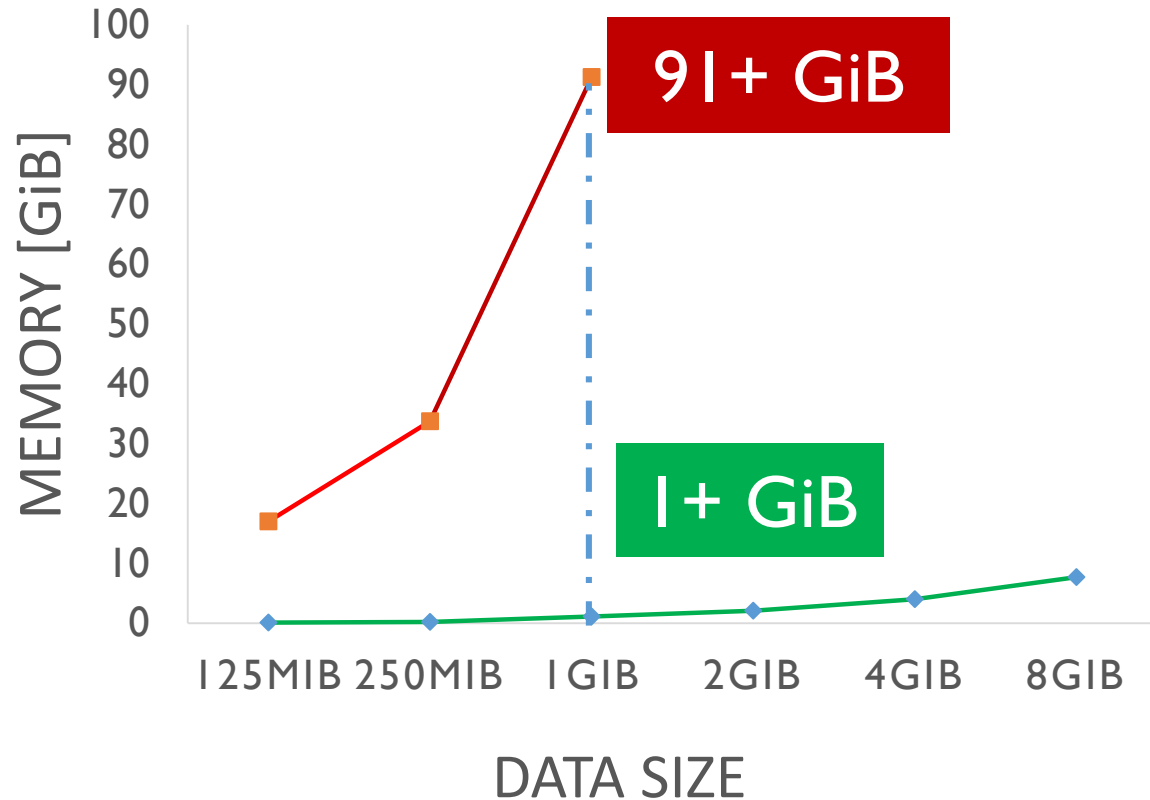
and so on

In practice:

- $0 \leq t_u^w \leq \text{Min}(n_u^w, M)$
- Generate 100 samples per node
- Forget samples



Test time comparison



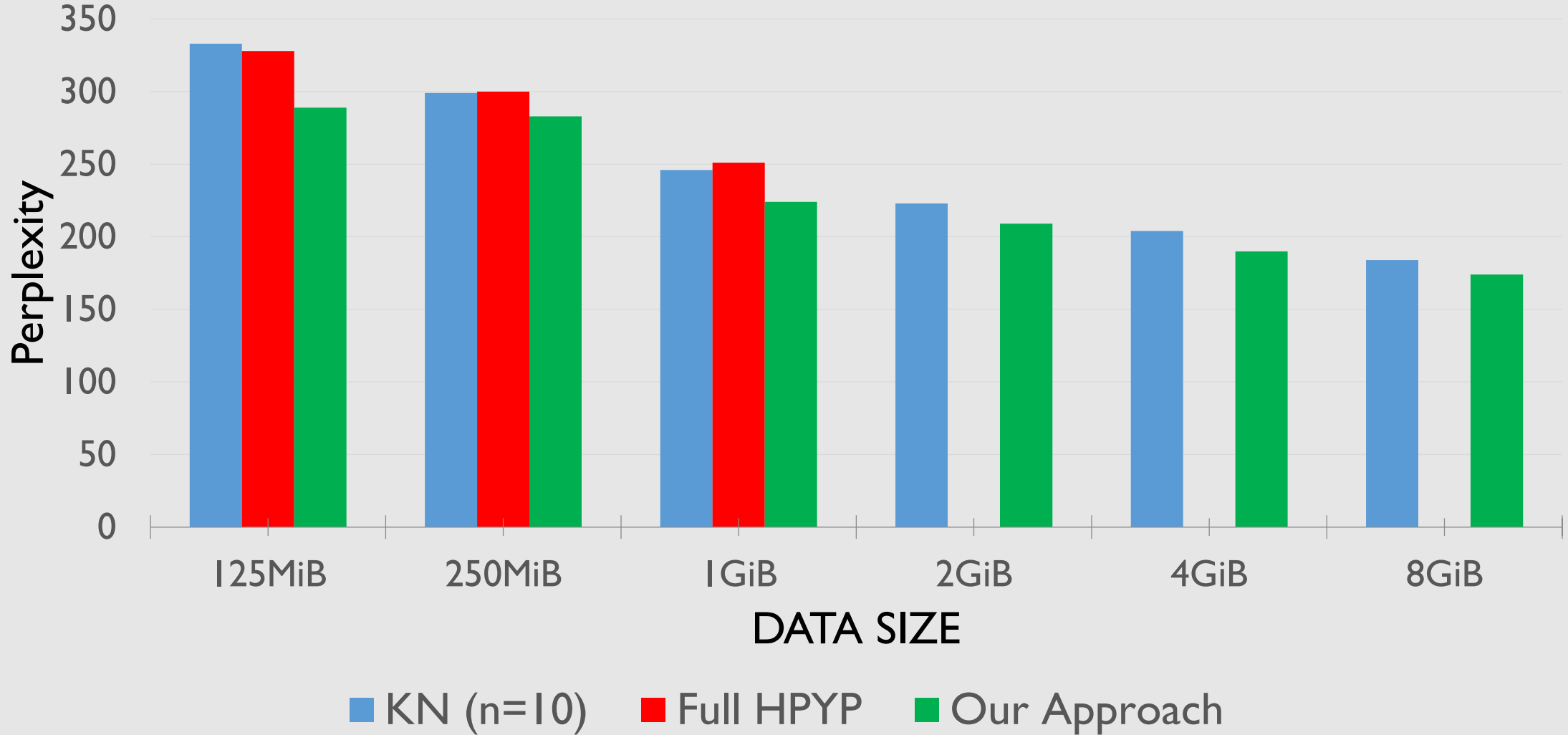
On 1 GiB:

- 1.8x slower on query
- 2.3x faster on load+query

Outline

- Infinite-Order Language Modelling and Challenges
- Compressed HPYP LM
- Inference and Sampling in Compressed HPYP LM
- Perplexity and Mixing
- Conclusion and Future Directions

Perplexity (and Mixing)



Outline

- Infinite-Order Language Modelling and Challenges
- Compressed HPYP LM
- Inference and Sampling in Compressed HPYP LM
- Perplexity and Mixing
- **Conclusion and Future Directions**

Conclusion

- Proposed a Compressed HPYP LM and a fast and memory-efficient approximate inference scheme.
- Proposed approach is several orders of magnitude smaller than the existing models.
- Avoided potential mixing issues, while consistently outperforming the state-of-the-art count-based language models by a significant margin.

Conclusion

- Proposed a Compressed HPYP LM and a fast and memory-efficient approximate inference scheme.
- Proposed approach is several orders of magnitude smaller than the existing models.
- Avoided potential mixing issues, while consistently outperforming the state-of-the-art count-based language models by a significant margin.

Future Directions

- Sampling speedup (i.e., learning an approximation for Stirling numbers)
- Exploring continuous space approximations of HPYP
- Exploring other applications

Conclusion

- Proposed a Compressed HPYP LM and a fast and memory-efficient approximate inference scheme.
- Proposed approach is several orders of magnitude smaller than the existing models.
- Avoided potential mixing issues, while consistently outperforming the state-of-the-art count-based language models by a significant margin.

Thanks!

Compressed Nonparametric Language Modelling

Slides, supplementary materials, more results available on : eehsan.github.io

Contact: Ehsan.Shareghi@gmail.com