

An Empirical Investigation of Discounting in Cross-Domain Language Models

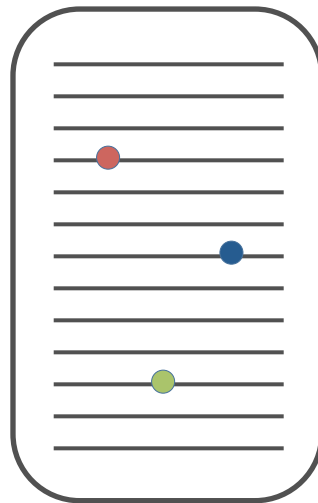


Greg Durrett and Dan Klein
UC Berkeley



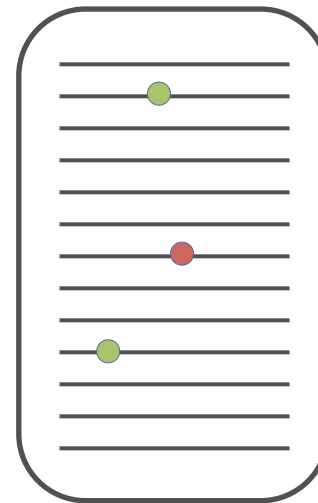
Discounting

Train

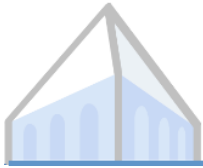


1

Test

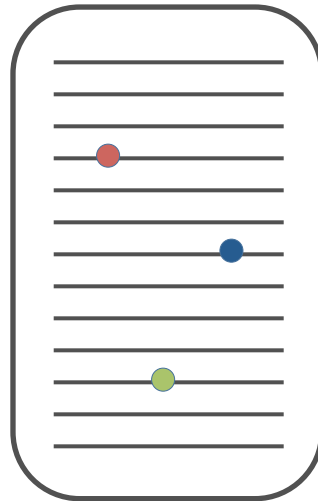


$c(\text{test}) = 2$

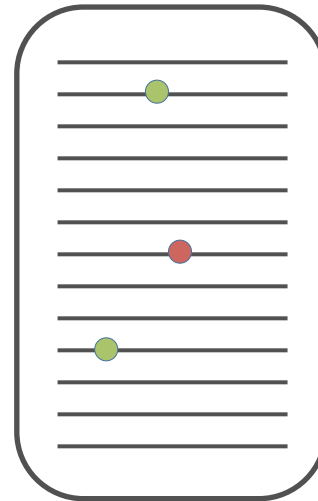


Discounting

Train



Test



1

$$c^*(1) = 0.28$$

2

$$c^*(2) = 0.97$$

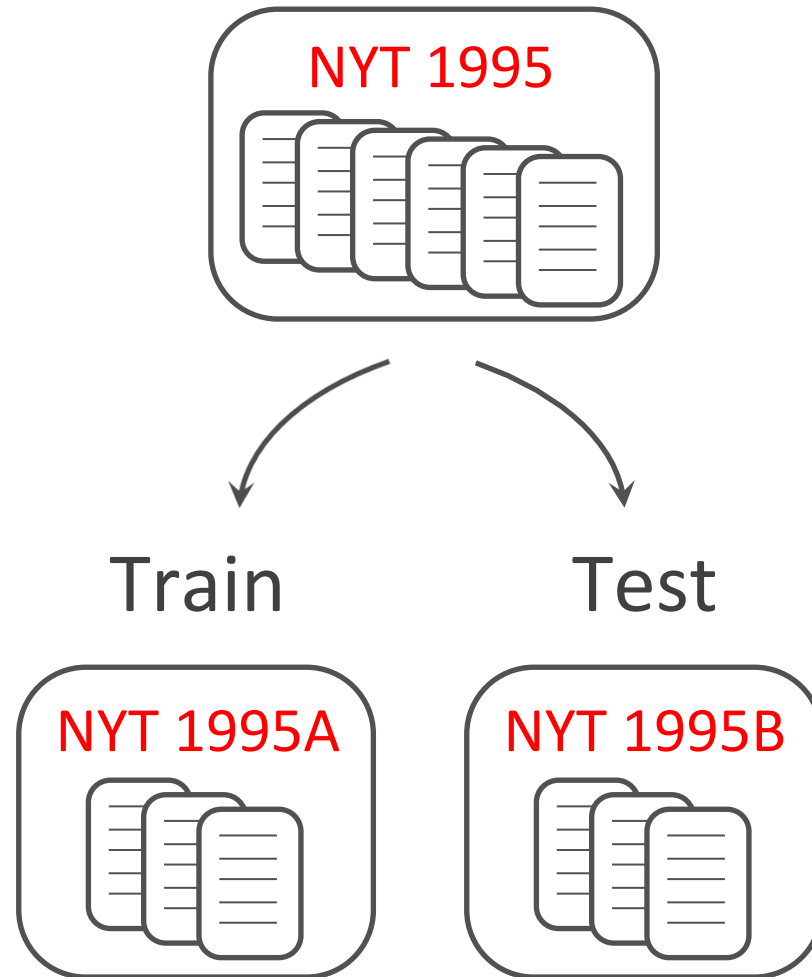
3

$$c^*(3) = 1.96$$

“future
count”



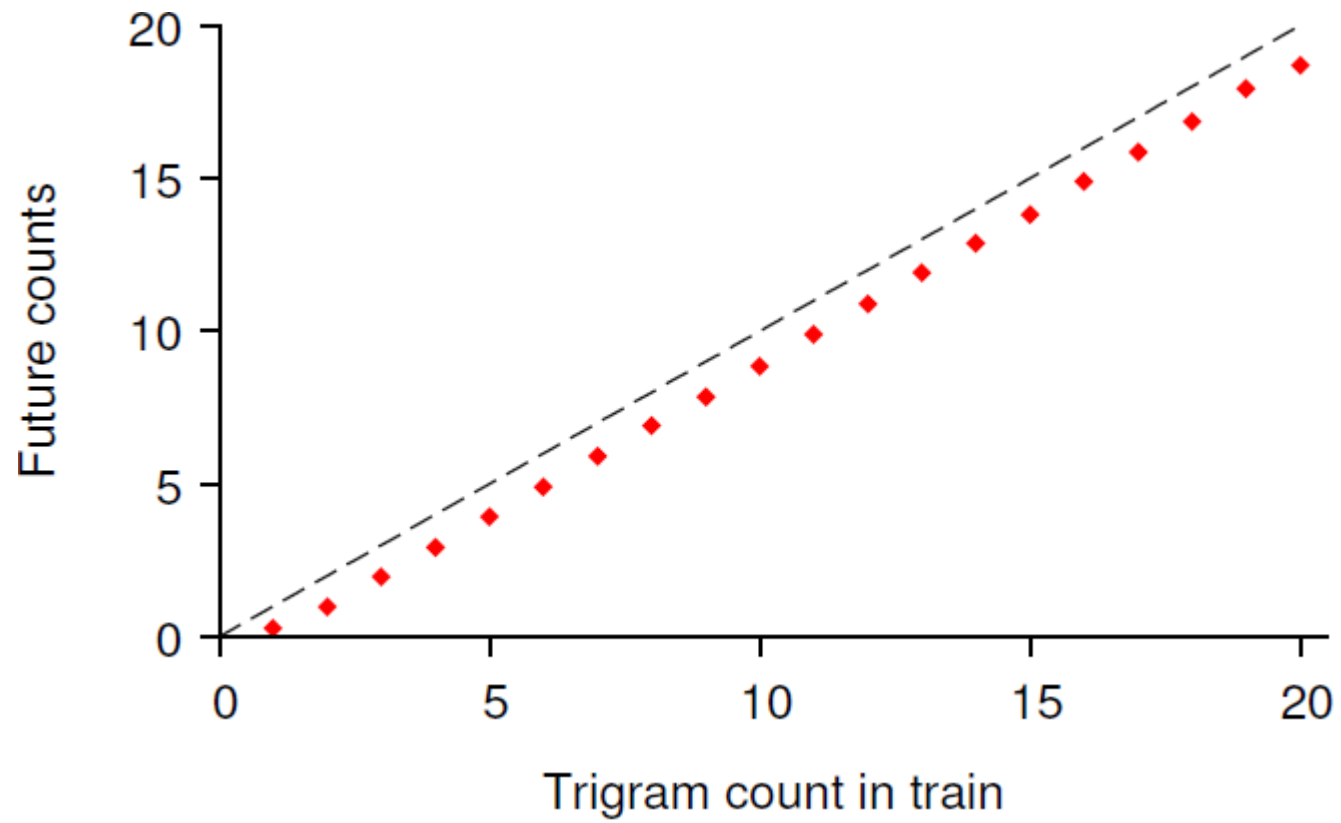
Experimental Setup



[Church and Gale, 1991]

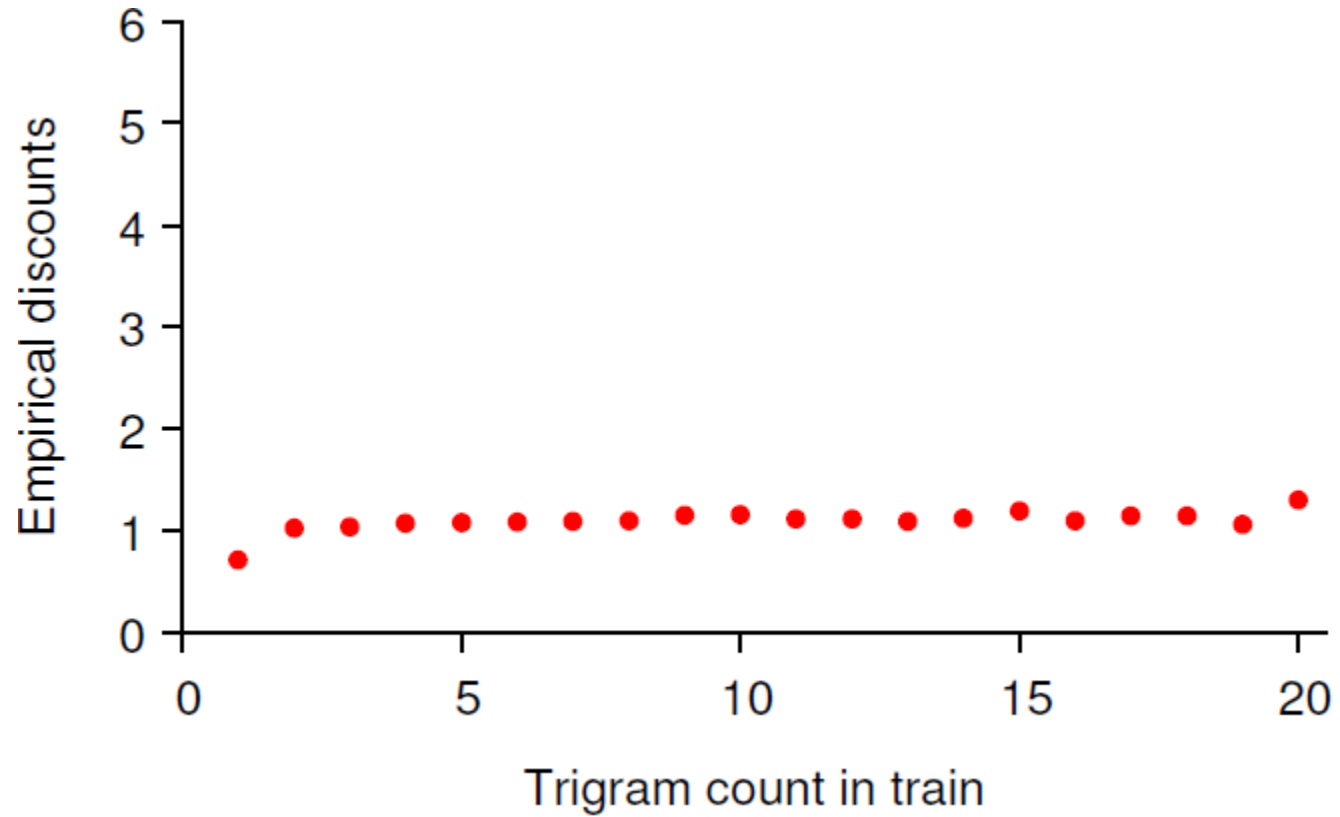


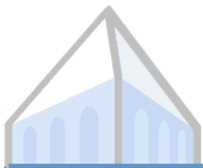
Future Counts



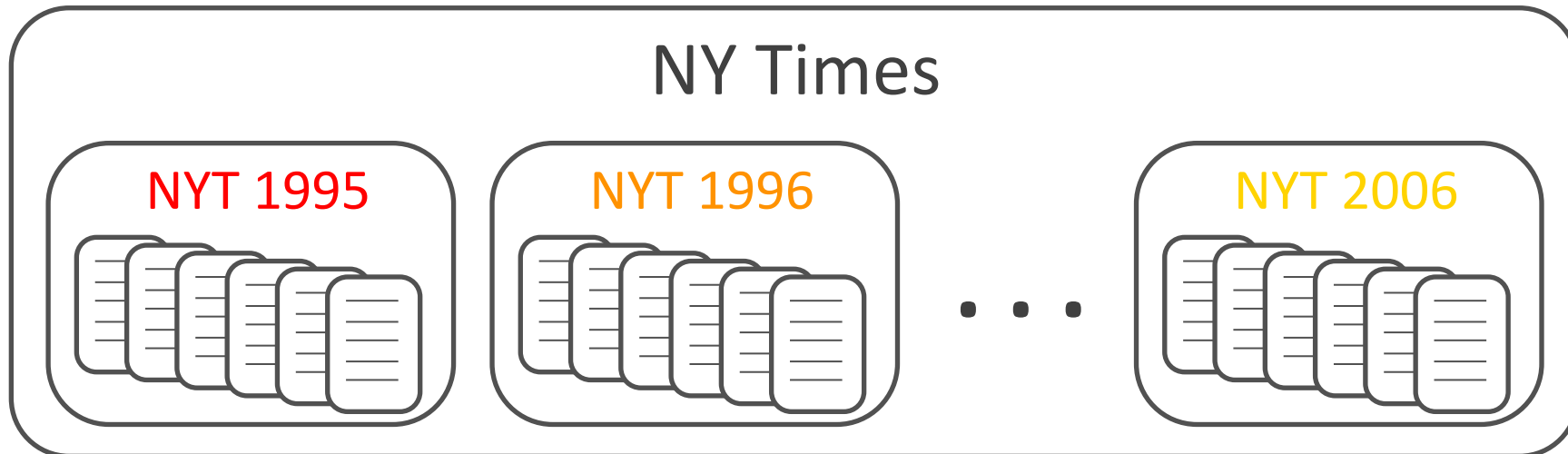


Empirical Discounts

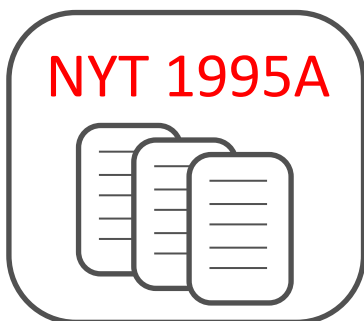




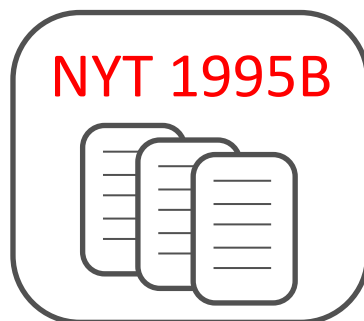
Experimental Setup



Train

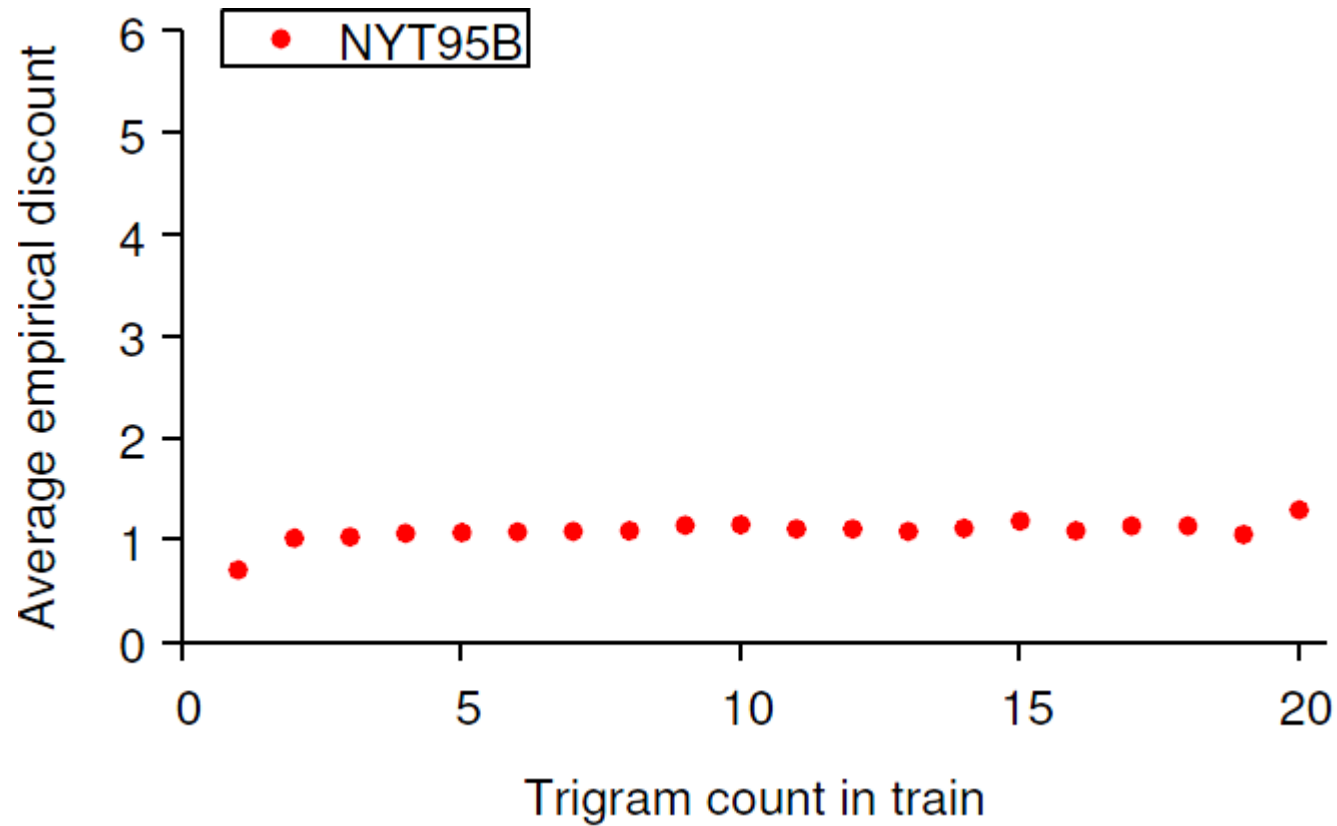


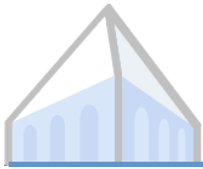
Test



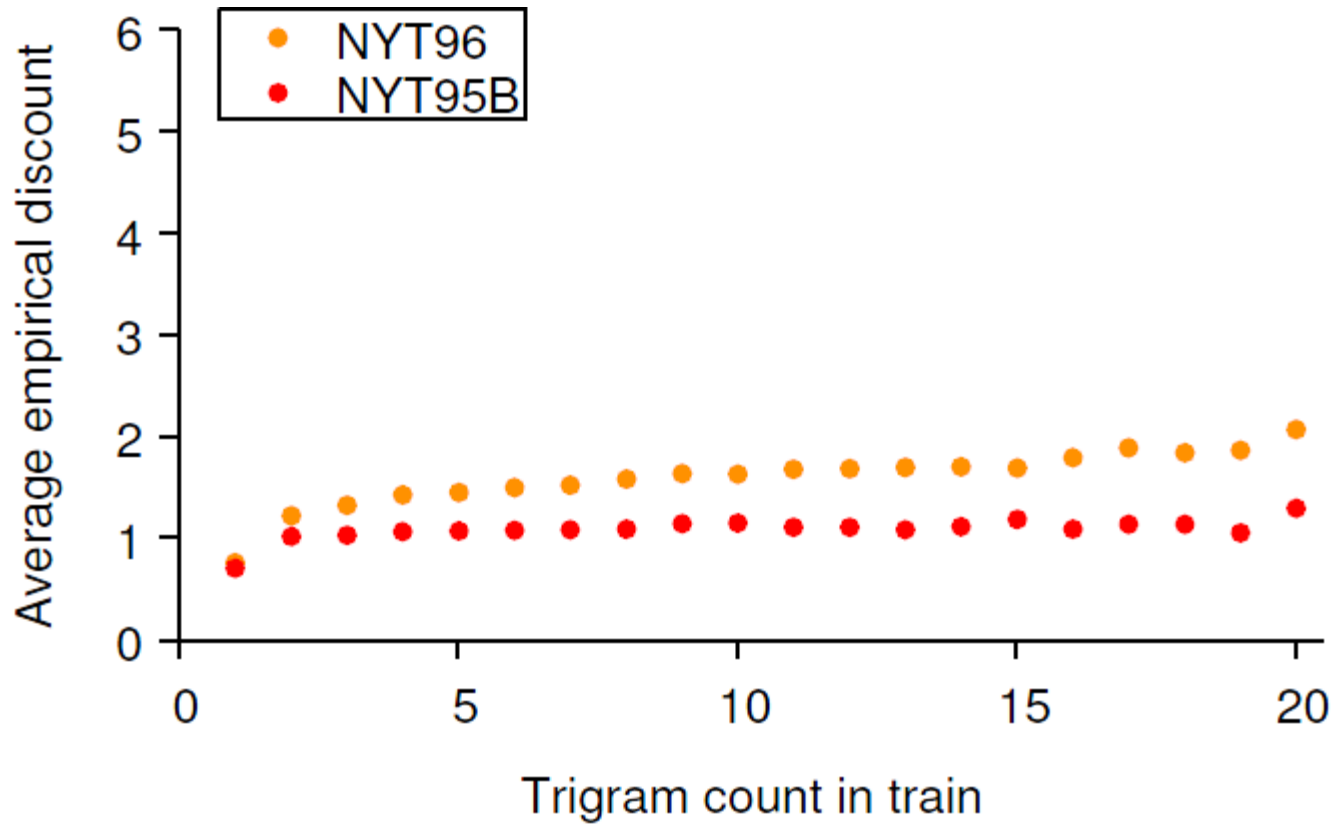


Empirical Discounts



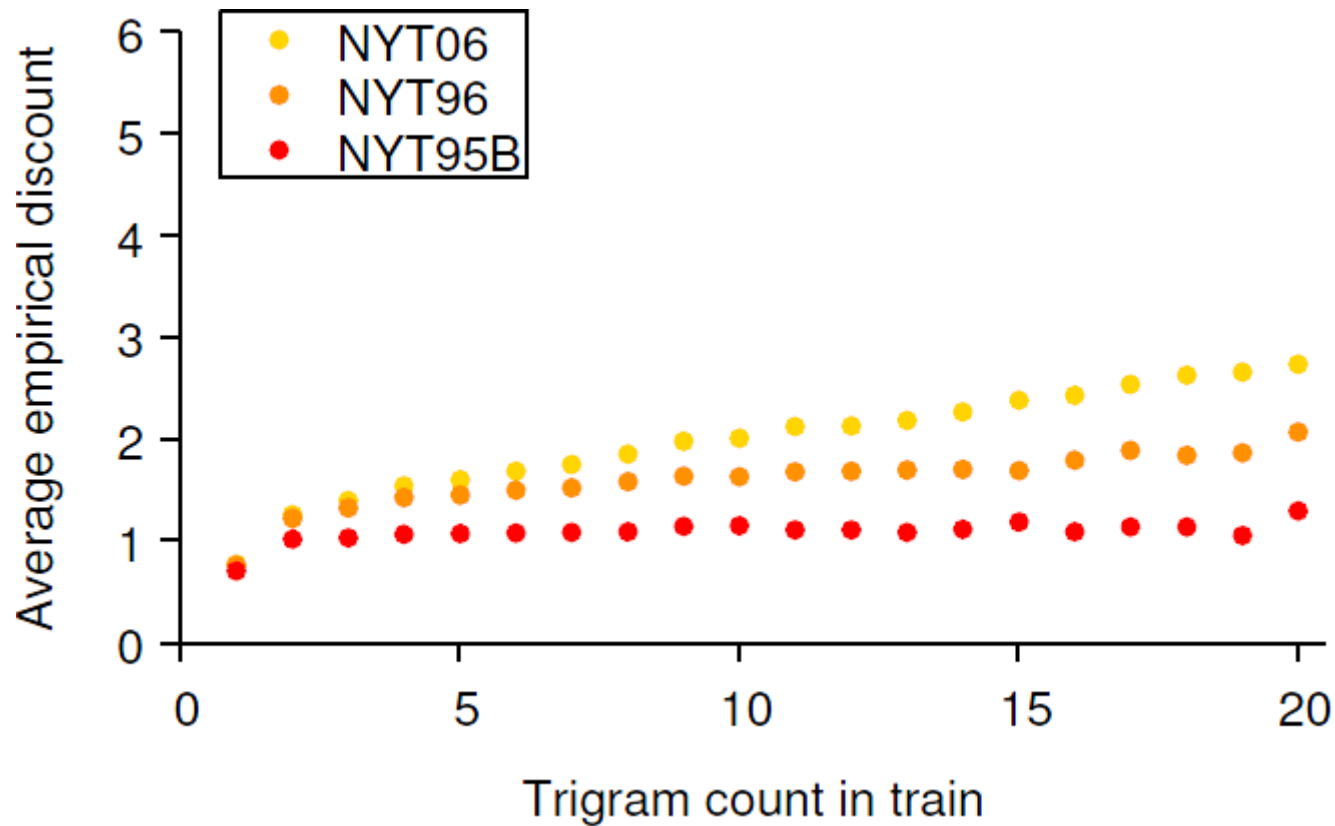


Empirical Discounts





Empirical Discounts



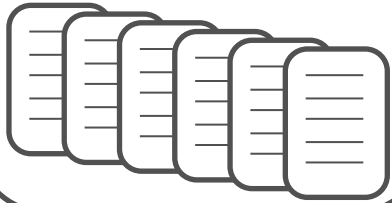
Changing the training year results in growing discounts



Experimental Setup

NY Times

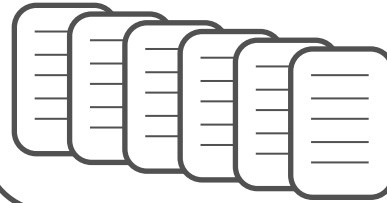
NYT 1995



...

Agence France-Presse

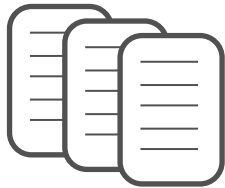
AFP 1995



...

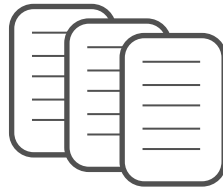
Train

NYT 1995A

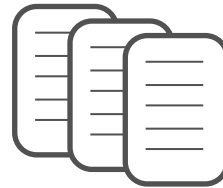


Test

NYT 1995B

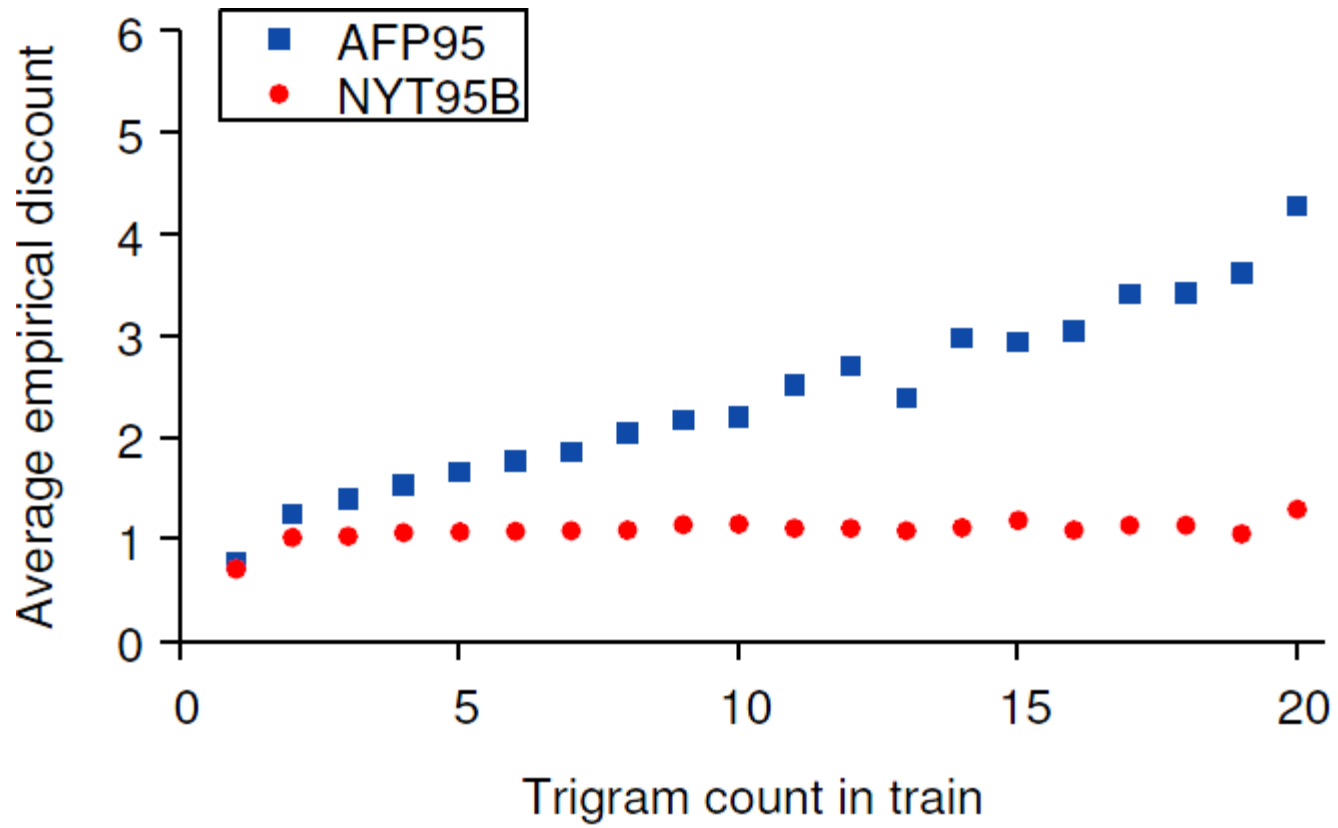


AFP 1995



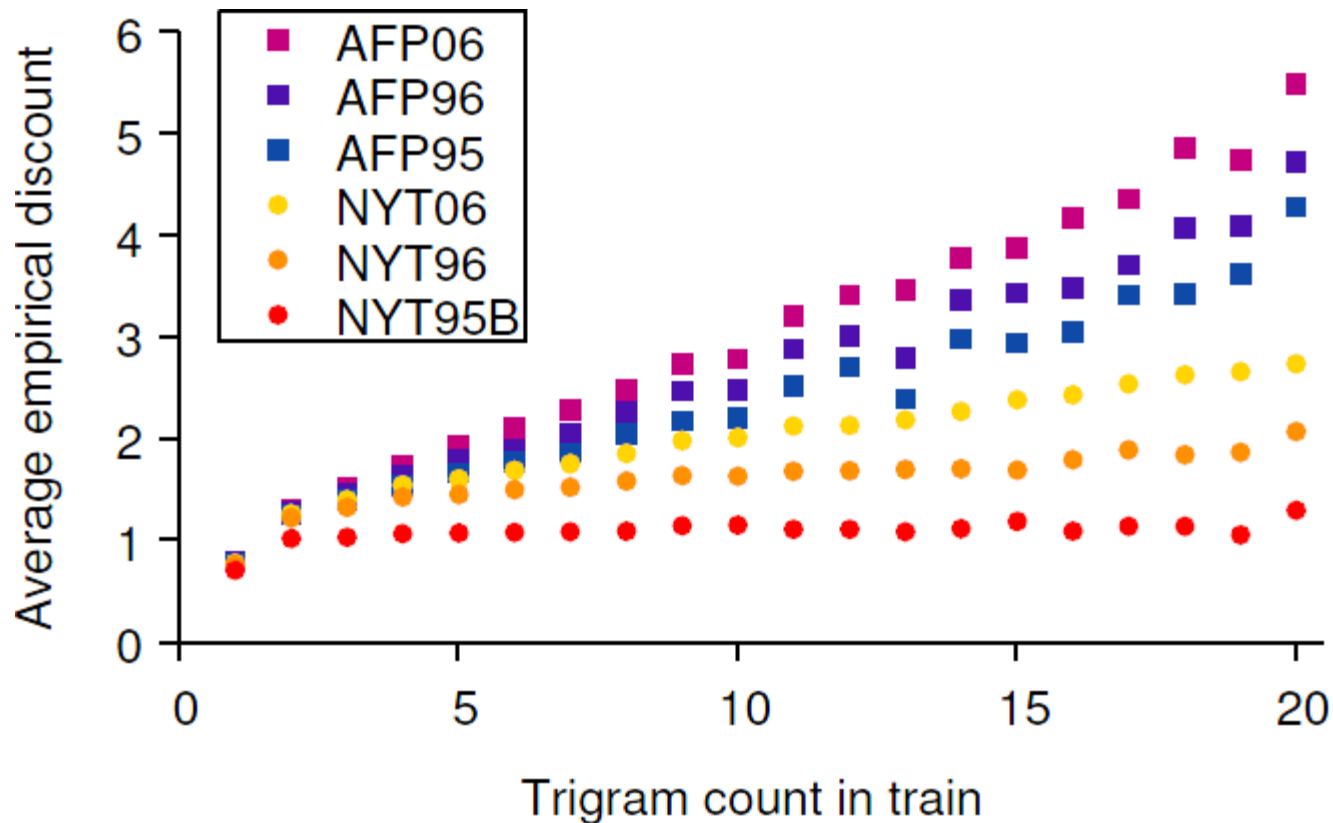


Empirical Discounts





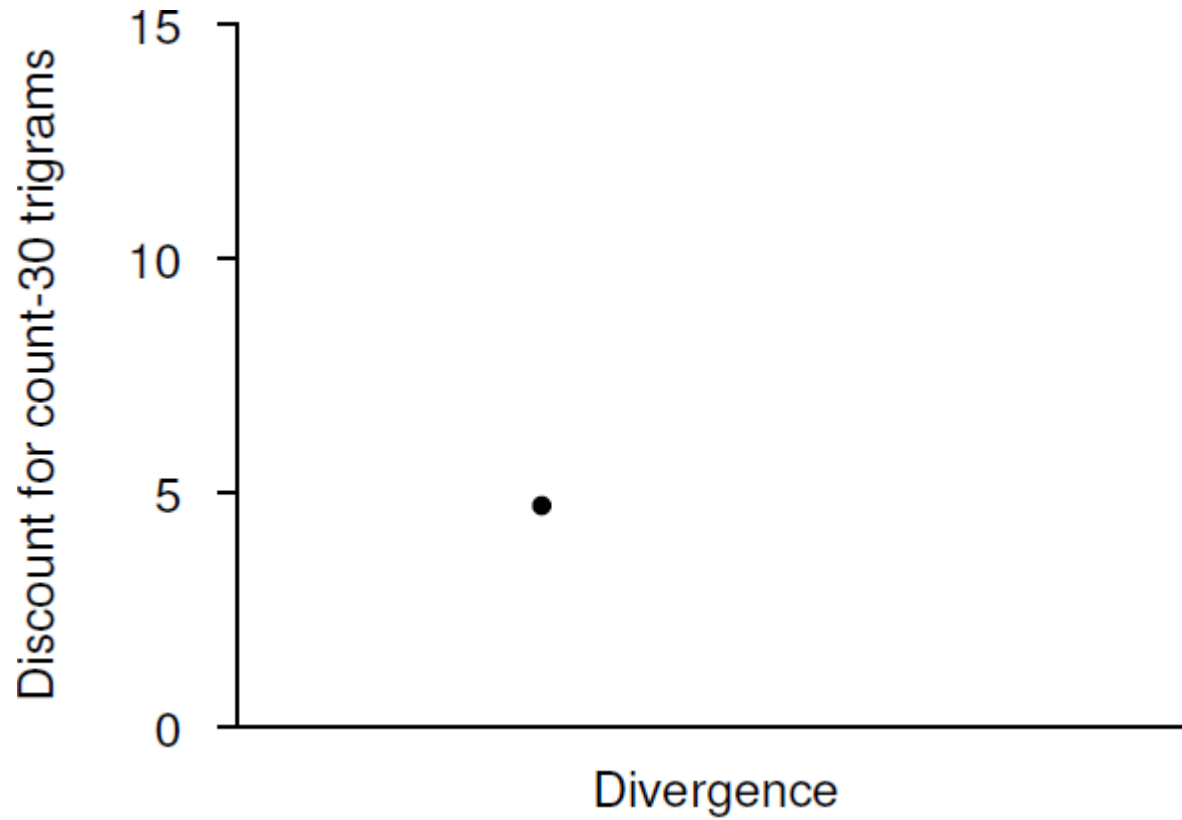
Empirical Discounts

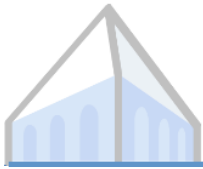


More discount growth as temporal divergence increases and source changes

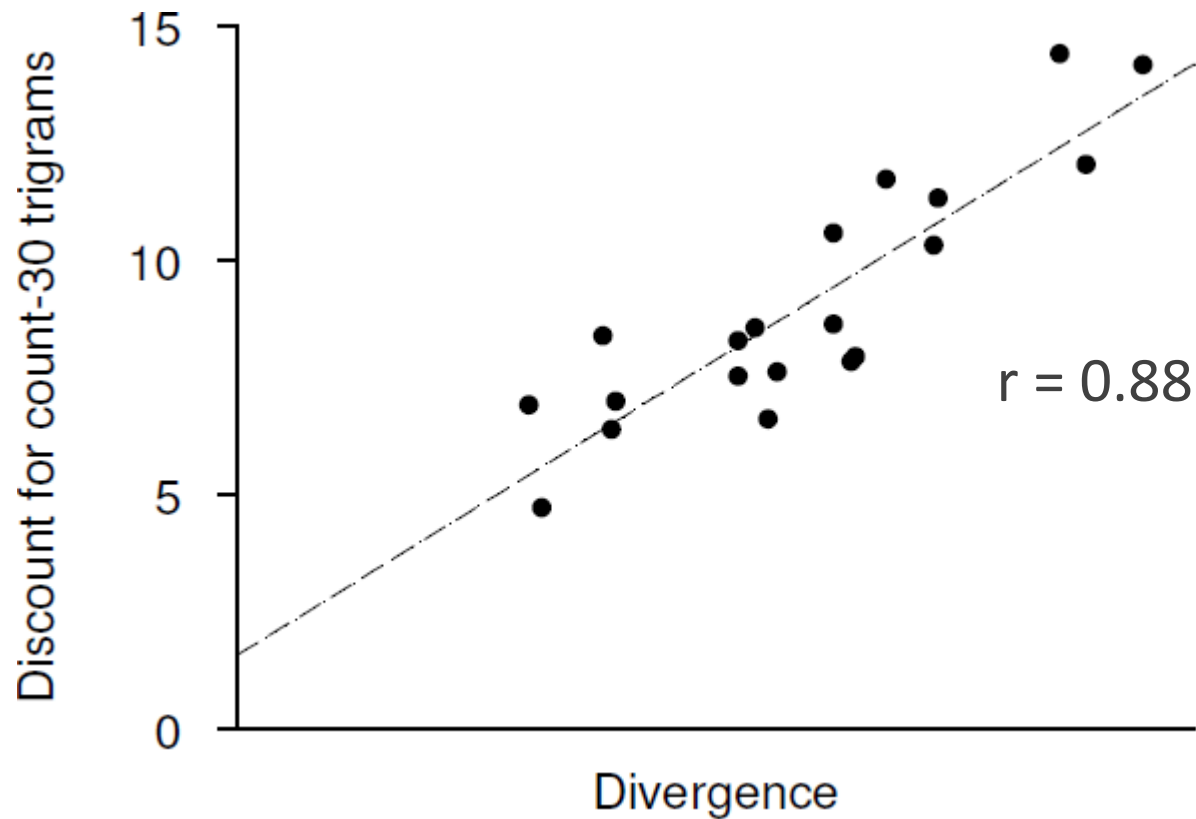


Predicting Discount Growth

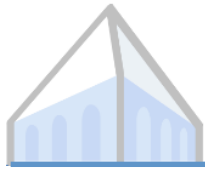




Predicting Discount Growth



More divergence yields more discounting



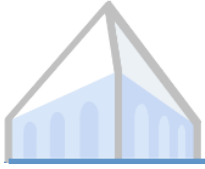
A Growing Discounts LM

- ▶ Interpolated Kneser-Ney

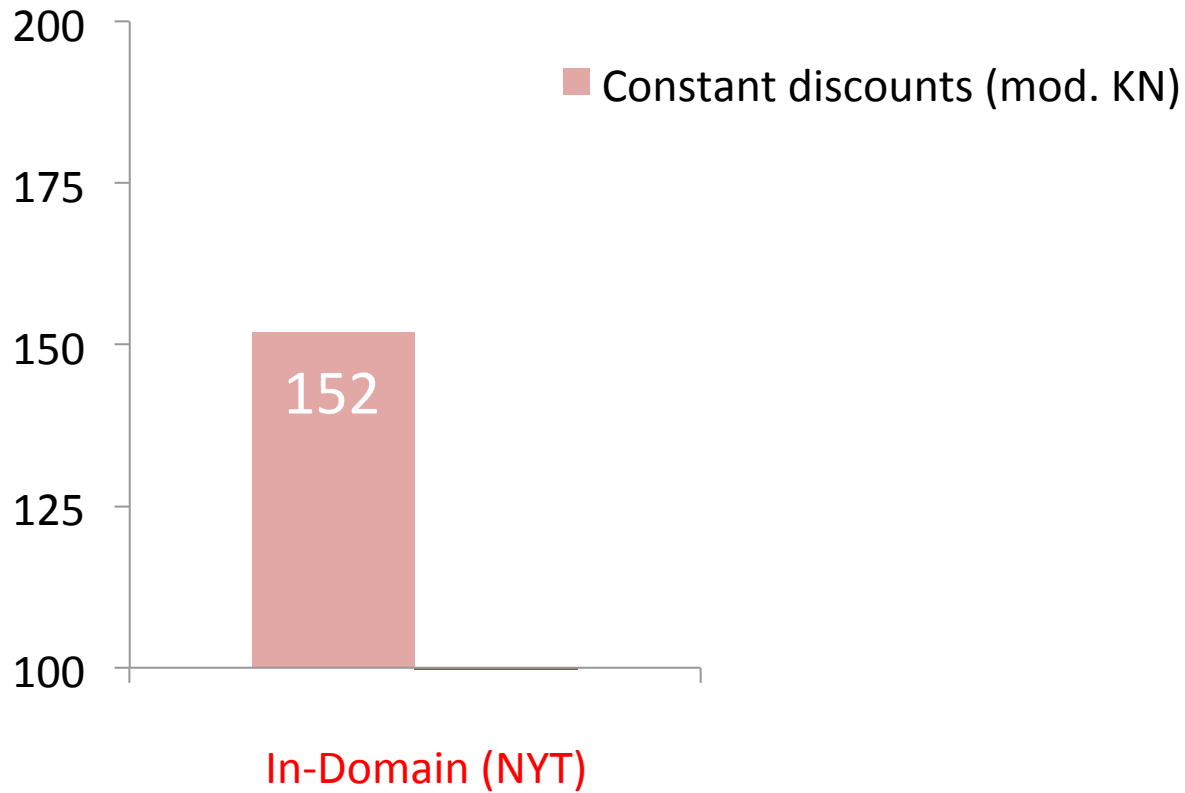
$$p(n\text{-gram}) = \frac{c(n\text{-gram}) - d}{c(\text{context})} + \dots$$

- ▶ Linearly Growing Discounts

$$p(n\text{-gram}) = \frac{c(n\text{-gram}) - (d_1 + d_2 c(n\text{-gram}))}{c(\text{context})} + \dots$$



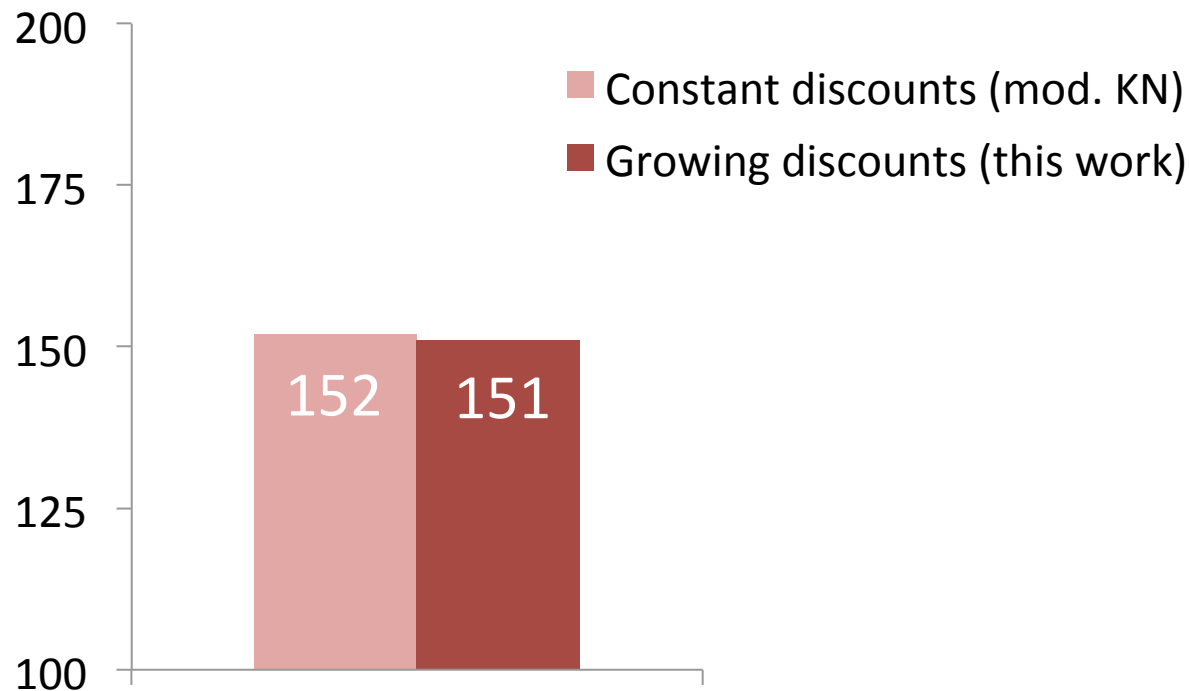
Perplexity Results



(Parameters tuned to maximize held-out NYT perplexity)



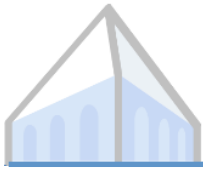
Perplexity Results



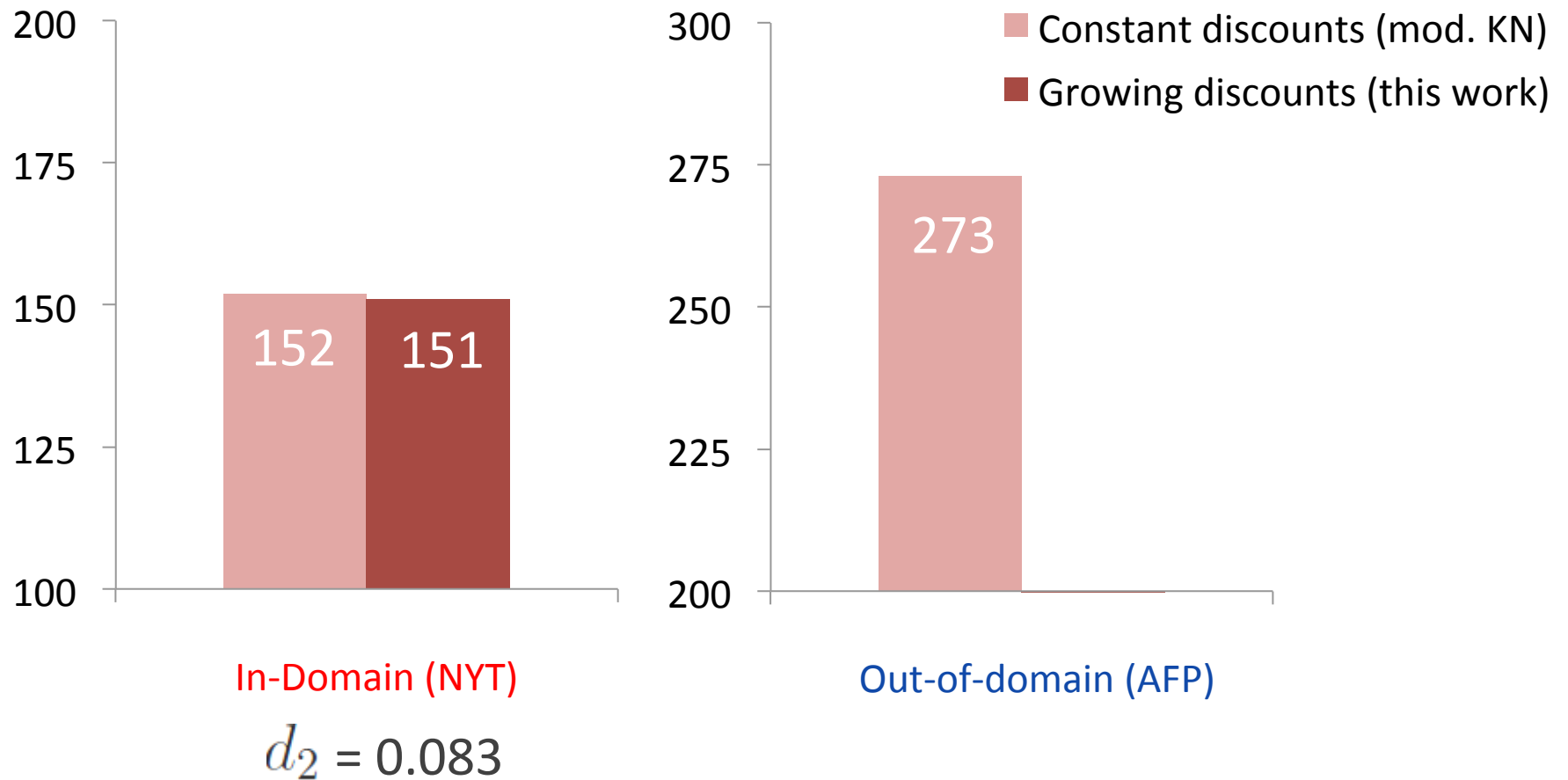
In-Domain (NYT)

$$d_2 = 0.083$$

(Parameters tuned to maximize held-out NYT perplexity)



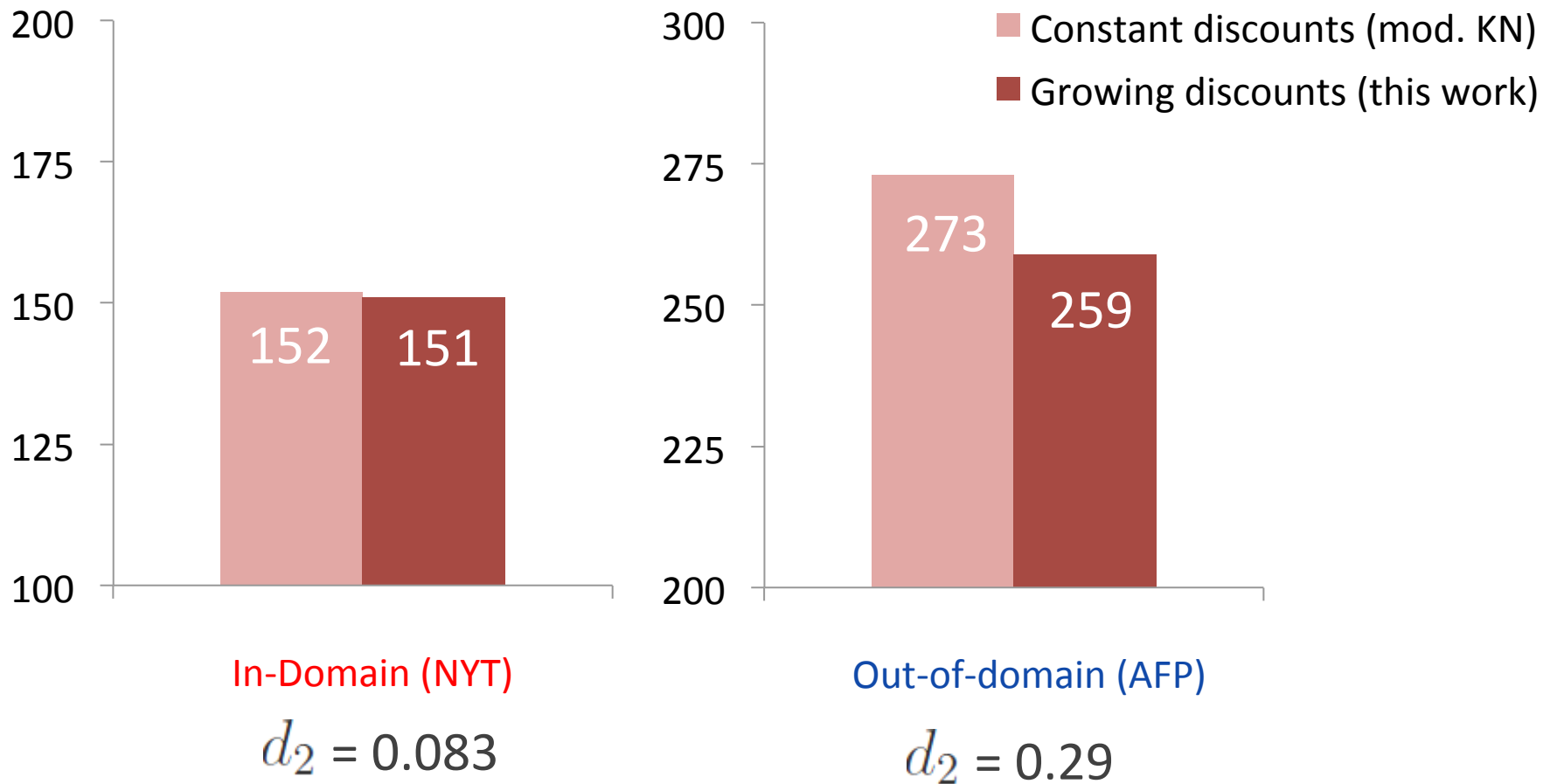
Perplexity Results



(Parameters tuned to maximize held-out NYT perplexity)



Perplexity Results



(Parameters tuned to maximize held-out NYT perplexity)



Conclusion

- ▶ Shape of discount must be changed (should grow with n-gram count) as corpora diverge
- ▶ Subtle cross-domain effects suggest using a qualitatively different model

Thank you!

